

THE ∞ -WASSERSTEIN DISTANCE: LOCAL SOLUTIONS AND EXISTENCE OF OPTIMAL TRANSPORT MAPS*

THIERRY CHAMPION[†], LUIGI DE PASCALE[‡], AND PETRI JUUTINEN[§]

Abstract. We consider the non-nonlineal optimal transportation problem of minimizing the cost functional $\mathcal{C}_\infty(\lambda) = \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|$ in the set of probability measures on Ω^2 having prescribed marginals. This corresponds to the question of characterizing the measures that realize the infinite Wasserstein distance. We establish the existence of “local” solutions and characterize this class with the aid of an adequate version of cyclical monotonicity. Moreover, under natural assumptions, we show that local solutions are induced by transport maps.

Key words. infinite Wasserstein distance, restrictable solutions, infinite cyclical monotonicity

AMS subject classifications. 49Q20, 49K30

DOI. 10.1137/07069938X

1. Introduction. In this paper, we consider the non-nonlineal optimal transportation problem that can be mathematically stated as the problem of minimizing the cost functional

$$(1.1) \quad \mathcal{C}_\infty(\lambda) := \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|$$

in the set of probability measures on Ω^2 having prescribed marginals. Here, and throughout the paper, we assume that Ω is a compact subset of \mathbb{R}^d , $d \geq 1$, $|\cdot|$ denotes the usual Euclidean norm in \mathbb{R}^d , μ, ν are the two (given) Borel probability measures on Ω , and $\Pi(\mu, \nu)$ denotes the set of admissible transport plans, i.e., the set of Borel probability measures λ on $\Omega^2 := \Omega \times \Omega$ with first marginal $\pi_1 \# \lambda = \mu$ and second marginal $\pi_2 \# \lambda = \nu$. Informally, if λ is induced by a transport map $T : \Omega \rightarrow \Omega$, i.e., $\lambda = (id \times T) \# \mu$, then $\mathcal{C}_\infty(\lambda)$ is simply the maximum of the transport distances $|T(x) - x|$.

The problem formulated above corresponds to the question of characterizing the measures that realize the infinite Wasserstein distance

$$(P_\infty) \quad W_\infty(\mu, \nu) = \inf \left\{ \mathcal{C}_\infty(\lambda) = \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x| : \lambda \in \Pi(\mu, \nu) \right\}$$

*Received by the editors August 6, 2007; accepted for publication (in revised form) November 16, 2007; published electronically March 26, 2008. Part of this research was conducted while the authors were visiting each other at their respective institutions.

<http://www.siam.org/journals/sima/40-1/69938.html>

[†]Laboratoire d’Analyse Non Linéaire Appliquée, U.F.R. des Sciences et Techniques, Université du Sud Toulon-Var, Avenue de l’Université, BP 20132, 83957 La Garde cedex, France (champion@univ-tln.fr). This author’s visit to Pisa in November, 2006, was supported by the I.N.D.A.M project “Traffic flows and optimization on complex networks” and by the G.N.A.M.P.A project “Fenomeni di evoluzione non lineari suggeriti dalla Fisica e dalla Biologia”, and his visit to Jyväskylä in June, 2007, was supported by the ESF program “Global and geometric properties of solutions of nonlinear partial differential equations”.

[‡]Dipartimento di Matematica Applicata, Università di Pisa, Via Buonarroti 1/c, 56127 Pisa, Italy (depascal@dm.unipi.it). This author’s research was partially supported by the Italian M.I.U.R. project “Metodi variazionali nella teoria del trasporto ottimo di massa e nella teoria geometrica della misura” and by the “Fondo di ateneo per la ricerca” of the University of Pisa.

[§]Department of Mathematics and Statistics, P.O. Box 35 (MaD), FI-40014 University of Jyväskylä, Finland (peanju@maths.jyu.fi). This author’s research was supported by Academy of Finland project 108374.

between μ and ν . Clearly, this is the limiting case, as $p \rightarrow \infty$, of the more familiar (see, e.g., [1, 2, 32]) p -Wasserstein distance problem

$$(P_p) \quad W_p(\mu, \nu) = \inf \left\{ \left(\int_{\Omega^2} |y - x|^p d\lambda(x, y) \right)^{\frac{1}{p}} : \lambda \in \Pi(\mu, \nu) \right\}, \quad 1 \leq p < \infty,$$

which is a model example of a Monge–Kantorovich-type optimal transport problem. Despite the close relationship, there are fundamental differences between these two problems. Most importantly, while (P_p) is linear in λ (removing the $1/p$ -power does not change the solution set), the mapping $\lambda \mapsto \mathcal{C}_\infty(\lambda)$ is not even convex. In particular, the problem (P_∞) is not additive, which implies that, unlike in the case of (P_p) , a restriction of an optimal transport plan need not be optimal for its own marginals.

In view of simple examples, it turns out that only imposing this “local optimality” property can lead to a satisfactory class of solutions. Hence we introduce in this paper the notion of *restrictable solutions*; this subclass of minimizers of (1.1) is characterized by the property that every portion of μ is transported onto its target in an optimal way; see Definition 4.1 below. The existence of a restrictable solution is obtained with the aid of approximating (P_∞) by the problems (P_p) . The same strategy also provides us with the notion of *infinite cyclical monotonicity*, which is derived from the standard c -cyclical monotonicity by applying it to the sequence of costs $c_p(x, y) = |x - y|^p$ and then taking the limit as $p \rightarrow \infty$. A reader familiar with the theory of infinity Laplacian and related problems [5] should recognize the analogy between restrictable solutions and absolute minimizers of supremum functionals.

It is one of the main results of this paper that restrictable and infinitely cyclically monotone solutions coincide; see Theorems 3.4 and 4.4 below. We would like to emphasize that although both of these notions are derived via an approximation argument, the proof for their equivalence is completely independent of the derivation. Moreover, this result holds without any further assumptions on the marginals μ and ν . The second principal question we address in this paper is existence and uniqueness of an optimal transport map. Our main result in this direction, Theorem 5.5, states that if $\mu \ll \mathcal{L}^d$, then any infinitely cyclically monotone solution γ to (P_∞) is induced by a map $T : \Omega \rightarrow \Omega$, i.e., $\gamma = (id \times T)_\# \mu$. Regarding the question of uniqueness, we are able to show that if, in addition to the previous assumptions, the second marginal ν is discrete, then the infinitely cyclically monotone solution to (P_∞) is unique.

A major technical difficulty that we are facing in the proofs is the absence of a useful duality theory, which is due to the nonconvexity of the objective functional (see section 5.4). As a consequence, we must rely on ad hoc arguments designed for the problem at hand. On the other hand, it is quite clear from the proofs that the machinery we are developing applies to more general problems than just (P_∞) . In fact, we could have just as well considered a functional $\lambda \mapsto \lambda\text{-ess sup}_{(x,y) \in \Omega^2} c(x, y)$, where $c(x, y)$ is, say, nonnegative and lower semicontinuous to begin with. In this work we concentrate on the model case $c(x, y) = |y - x|$ so as to identify the useful tools and notions without coping with the additional technical difficulties required by a more general cost c , which seems to be the natural next step in this study.

Let us finish this introduction by discussing some applications in which the infinite Wasserstein distance W_∞ appears. First on our list is the optimal design problem

$$(1.2) \quad \sup \left\{ \frac{W_\infty(\mu, \nu)^{p+d}}{W_p(\mu, \nu)^p \|(\frac{d\mu}{dx})^{-1}\|_{L^\infty(U)}} : (\mu, \nu) \in \mathcal{P}_{a.c.}(U) \times \mathcal{P}(\bar{U}) \right\}$$

that appears in [9] in connection with stability estimates for optimal transport maps; here \mathcal{P} and $\mathcal{P}_{a.c.}$ denote the spaces of probability measures and absolutely continuous probability measures, respectively, and $\frac{d\mu}{dx}$ is the Radon–Nikodym derivative of μ with respect to the Lebesgue measure. In [9], the authors prove that if $U \subset \mathbb{R}^d$ is a bounded Lipschitz domain, then the estimate

$$(1.3) \quad W_\infty(\mu, \nu)^{p+d} \leq C_{p,d}(U) \left\| \left(\frac{d\mu}{dx} \right)^{-1} \right\|_{L^\infty(U)} W_p(\mu, \nu)^p$$

holds for every $p > 1$. The inequality (1.3) is an intrinsic counterpart of a beautiful uniform estimate for the optimal transport maps proved in [9] under stronger regularity assumptions. The optimal constant in (1.3) is given by the supremum in (1.2), and it is conjectured (based on 1-dimensional examples and remarks on increasing transport maps) that it does not blow up when $p \rightarrow 1$.

Second, during the last few years, models of branching processes using in one way or another tools from the optimal transportation theory have been proposed by several authors; see, for example, [8, 10, 23, 34]. Roughly speaking, these models favor joint transportation, which in many real world situations, such as in the design of communication or irrigation networks, is more economical than individualized transportation. In particular, in [10] the authors propose minimizing a certain cost functional (which penalizes diffused measures) on the p -Wasserstein space of probability measures. It is remarked in [29] that this model is somewhat less realistic than the others cited above, but it has the advantage of being mathematically simpler. Moreover, as pointed out in Remark 6.2.7, Chapter 4, and section 0.2 of [29], the use of the infinite Wasserstein distance W_∞ in the model of [10] produces results which are closer to the ones derived from the other models.

Then there are applications to PDEs. The metric structure associated with the ∞ -Wasserstein distance is a crucial tool in proving the existence of stable solutions for a compressible fluid model of rotating binary stars in [24]. The same metric was used also to bound the growth of the wetted regions in the porous medium flow [13] and to study the long time asymptotics of nonlinear scalar conservation laws [12]. Moreover, the ∞ -Wasserstein distance is being used in some N -particle approximations of the Vlasov equations [21, 22].

Finally, in [11] the authors have considered a mathematical model of the optimal pricing policy for the use of a public transportation network. This model assumes that the price of a ticket (for the use of the network) is a function of the distance traveled. This seems reasonable in the case when each citizen is associated with a single journey, but it is not so realistic if we allow multiple journeys and an inexpensive season ticket is available. In the latter case, the price of a season ticket could be assumed to be a piecewise constant function of the maximal distance traveled, and hence it might be a good idea to insert a component similar to the functional we have considered into the model. It is also quite easy to imagine that in many other transportation problems a significant portion of the total cost is in one way or another connected with the maximal transportation distance. For example, if we assume that the physical transportation device (airplane, car, etc.) is the same for all distances, then it has to be chosen so that the longest transportation can be handled.

2. Existence of global solutions. As pointed out in the introduction, the objective functional

$$\begin{aligned} \lambda \mapsto \mathcal{C}_\infty(\lambda) &:= \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x| \\ &= \inf \left\{ t \geq 0 : \lambda(\{(x, y) \in \Omega^2 : |y - x| > t\}) = 0 \right\} \end{aligned}$$

is not linear (and not even convex) in λ , contrary to what is usually the case in classical optimal transport problems. However, it is, quite interestingly, level convex in the sense that if $\lambda_1, \lambda_2 \in \Pi(\mu, \nu)$, then

$$\mathcal{C}_\infty(t\lambda_1 + (1-t)\lambda_2) \leq \max\{\mathcal{C}_\infty(\lambda_1), \mathcal{C}_\infty(\lambda_2)\} \quad \text{for all } t \in (0, 1).$$

Note that this implies that the set of solutions to (P_∞) is convex. Moreover, it should be observed that $\mathcal{C}_\infty(\lambda)$ depends on the measure λ only via its support. More precisely, one has

$$(2.1) \quad \mathcal{C}_\infty(\lambda) = \sup\{|y - x| : (x, y) \in \text{supp}(\lambda)\}.$$

Thanks to this last property, we are in position to give the following existence result.

PROPOSITION 2.1. *Assume that Ω is a compact subset of \mathbb{R}^d and μ, ν are two probability measures on Ω . Then the problem*

$$(P_\infty) \quad W_\infty(\mu, \nu) = \inf \left\{ \mathcal{C}_\infty(\lambda) := \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x| : \lambda \in \Pi(\mu, \nu) \right\}$$

admits at least one solution $\lambda \in \Pi(\mu, \nu)$.

The optimal set of (P_∞) may be very large thanks to (2.1).

Example 2.2. Let $\mu := \frac{1}{2} \mathcal{L}^2|_{[0,1]^2 \cup [2,3]^2}$ and $\nu := \frac{1}{2}(\delta_{(2,1)} + \delta_{(1,2)})$. Then it is clear that the value of (P_∞) is $\sqrt{5}$ and that any admissible transport plan $\lambda \in \Pi(\mu, \nu)$ is a solution of (P_∞) .

In the proof of Proposition 2.1, we shall use the following lemma.

LEMMA 2.3. *If the sequence $(\lambda_n)_n$ converges weakly to λ in $\Pi(\mu, \nu)$, then for any $(x, y) \in \text{supp}(\lambda)$ there exists a sequence $((x_n, y_n))_{n \in \mathbb{N}}$ such that*

$$(2.2) \quad (x_n, y_n) \rightarrow (x, y) \text{ as } n \rightarrow \infty \quad \text{and} \quad (x_n, y_n) \in \text{supp}(\lambda_n) \quad \text{for all } n \in \mathbb{N}.$$

Proof. Suppose $(x, y) \in \Omega$ is such that (2.2) does not hold. Then we may assume without loss of generality that there exists $r > 0$ such that $B((x, y), r) \cap \text{supp}(\lambda_n) = \emptyset$ for any $n \in \mathbb{N}$. It then obviously follows from the weak convergence that $B((x, y), r) \cap \text{supp}(\lambda) = \emptyset$, which concludes the proof. \square

Proof of Proposition 2.1. Since the set Ω is a compact subset of \mathbb{R}^d and the measures μ and ν are probability measures on Ω , the nonempty set $\Pi(\mu, \nu)$ is compact for the weak convergence of measures (cf. [32, p. 49]). To apply the direct method of the calculus of variations, it remains to notice that $\lambda \mapsto \mathcal{C}_\infty(\lambda)$ is lower semicontinuous for this topology: this is a direct consequence of (2.1) and Lemma 2.3. \square

3. Infinitely cyclically monotone solutions. The proof of the existence of a solution to (P_∞) given in Proposition 2.1 is intrinsic, but one may obtain this result also via an approximation argument involving the family of problems $(P_p)_{p \geq 1}$ given by

$$(P_p) \quad W_p(\mu, \nu) = \inf \left\{ \mathcal{C}_p(\lambda) := \left(\int_{\Omega^2} |y - x|^p d\lambda(x, y) \right)^{\frac{1}{p}} : \lambda \in \Pi(\mu, \nu) \right\};$$

that is, the functional $\lambda \mapsto \mathcal{C}_p(\lambda)$ is being minimized over the set $\Pi(\mu, \nu)$.

Alternative Proof of Proposition 2.1. Under the assumptions made on Ω , μ , and ν , for any $p \geq 1$ the problem (P_p) admits at least one solution $\gamma_p \in \Pi(\mu, \nu)$; see, e.g.,

[32, Theorem 1.3]. Since $\Pi(\mu, \nu)$ is compact, we infer that $(\gamma_p)_{p \geq 1}$ converges weakly (up to a subsequence) to some $\gamma_\infty \in \Pi(\mu, \nu)$ as $p \rightarrow \infty$. Then, for any $\lambda \in \Pi(\mu, \nu)$, we have by the optimality of γ_p and Hölder's inequality that

$$\mathcal{C}_q(\gamma_p) = \left(\int_{\Omega^2} |y - x|^q d\gamma_p(x, y) \right)^{\frac{1}{q}} \leq \mathcal{C}_p(\gamma_p) \leq \mathcal{C}_p(\lambda)$$

for any $p \geq q \geq 1$. For a fixed $q \geq 1$, since the function $(x, y) \mapsto |y - x|^q$ is continuous and bounded on Ω^2 one has $\mathcal{C}_q(\gamma_p) \rightarrow \mathcal{C}_q(\gamma_\infty)$ as $p \rightarrow \infty$. Therefore, taking the limit in p and then in q in the above inequality we obtain $\mathcal{C}_\infty(\gamma_\infty) \leq \mathcal{C}_\infty(\lambda)$. Since this holds for any $\lambda \in \Pi(\mu, \nu)$, γ_∞ is a minimizer of (P_∞) . \square

The reason for considering the problems (P_p) in this context is not merely the fact that they provide an alternative route to the existence. Namely, it is known that an element $\gamma_p \in \Pi(\mu, \nu)$ is a solution of (P_p) if and only if its support is p -cyclically monotone; that is,

$$(3.1) \quad \sum_{i=1}^n |y_i - x_i|^p \leq \sum_{i=1}^n |y_{\sigma(i)} - x_i|^p$$

for every $n \geq 2$, $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\gamma_p)$, and for every permutation $\sigma \in \mathcal{S}_n$. We refer the reader, for example, to Theorem 3.2 in [2] (or to [20, 28, 33]).

By analogy with the p -cyclical monotonicity when $1 \leq p < \infty$, we introduce the corresponding notion for the case $p = \infty$ obtained by taking the limit in (3.1).

DEFINITION 3.1. *A transport plan $\gamma \in \Pi(\mu, \nu)$ is infinitely cyclically monotone if*

$$\max_{1 \leq i \leq n} |y_i - x_i| \leq \max_{1 \leq i \leq n} |y_{\sigma(i)} - x_i|$$

for every $n \geq 2$, $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\gamma)$, and $\sigma \in \mathcal{S}_n$.

Using again the approximation of (P_∞) by the problems (P_p) , we obtain the existence of an infinitely cyclically monotone solution to (P_∞) .

THEOREM 3.2. *For $1 \leq p < \infty$, let $\gamma_p \in \Pi(\mu, \nu)$ be a solution to (P_p) . Then any cluster point γ_∞ of $(\gamma_p)_{p \geq 1}$ in $\Pi(\mu, \nu)$ as $p \rightarrow \infty$ is an infinitely cyclically monotone solution to (P_∞) .*

Proof. For simplicity, let us assume that the entire family $(\gamma_p)_{p \geq 1}$ converges weakly to $\gamma_\infty \in \Pi(\mu, \nu)$. It suffices to show that γ_∞ is infinitely cyclically monotone. To this end, let $n \geq 2$, $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\gamma)$, and $\sigma \in \mathcal{S}_n$. We apply Lemma 2.3 to each pair (x_i, y_i) to obtain the existence of sequences $(x_1^p, y_1^p), \dots, (x_n^p, y_n^p)$ such that $(x_i^p, y_i^p) \rightarrow (x_i, y_i)$ for any i as $p \rightarrow \infty$, and $(x_i^p, y_i^p) \in \text{supp}(\gamma_p)$ for all $1 \leq p < \infty$ and $i = 1, \dots, n$. Since the support of γ_p is p -cyclically monotone, one has

$$\sum_{i=1}^n |y_i^p - x_i^p|^p \leq \sum_{i=1}^n |y_{\sigma(i)}^p - x_i^p|^p \quad \text{for all } 1 < p < \infty.$$

Taking the $1/p$ -power on both sides and letting p go to ∞ , one obtains the desired inequality. \square

Since for any $1 \leq p < \infty$ an admissible transport plan $\lambda \in \Pi(\mu, \nu)$ is a minimizer of (P_p) if and only if it is p -cyclically monotone, it is natural to ask whether this still holds for $p = \infty$. It is, however, quite clear that a generic minimizer of (P_∞) need

not be infinitely cyclically monotone; the following gives a counterexample for this implication.

Example 3.3. As an admissible transport plan for the measures μ and ν in Example 2.2 one may take

$$\lambda := \frac{1}{2} (\mathcal{L}^2|_{[0,1]^2} \times \delta_{(2,1)} + \mathcal{L}^2|_{[2,3]^2} \times \delta_{(1,2)}).$$

Then λ is a minimizer of (P_∞) , but it is not infinitely cyclically monotone: for example, $((0, 1), (2, 1))$ and $((3, 2), (1, 2))$ belong to $\text{supp}(\lambda)$, and

$$\max\{|(0, 1) - (1, 2)|, |(3, 2) - (2, 1)|\} < \max\{|(0, 1) - (2, 1)|, |(3, 2) - (1, 2)|\}.$$

Notice that in this case, for any $1 < p < \infty$, problem (P_p) admits a unique solution (up to a μ -negligible set) (see [20]), which in fact does not depend on p .

On the other hand, the following result shows that the reverse implication does hold: infinite cyclical monotonicity is indeed a sufficient condition for an admissible plan to be a minimizer of (P_∞) .

THEOREM 3.4. *Any infinitely cyclically monotone transport plan $\gamma \in \Pi(\mu, \nu)$ is a solution of the problem (P_∞) .*

Proof. We make a proof by contradiction. Let $\gamma \in \Pi(\mu, \nu)$ be infinitely cyclically monotone, and assume that

$$(3.2) \quad \gamma\text{-ess sup}_{(x,y) \in \Omega^2} |y - x| \geq 10\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|$$

for some $\tilde{\gamma} \in \Pi(\mu, \nu)$ and $\varepsilon > 0$.

Since Ω is compact, there exists a finite family $(c_i)_{1 \leq i \leq k}$ such that $\Omega \subset \bigcup_{i=1}^k B(c_i, \varepsilon)$. We shall denote $C := \{c_1, \dots, c_k\}$ and $V_1 := B(c_1, \varepsilon)$, and for any $i \in \{2, \dots, k\}$ we set $V_i := B(c_i, \varepsilon) \setminus \bigcup_{j=1}^{i-1} V_j$; without loss of generality, we assume that $V_i \neq \emptyset$ for all $i \in \{1, \dots, k\}$.

Next we define two discrete measures γ^ε and $\tilde{\gamma}^\varepsilon$ on Ω^2 by

$$\gamma^\varepsilon := \sum_{1 \leq i, j \leq k} \gamma(V_i \times V_j) \delta_{(c_i, c_j)}$$

and

$$\tilde{\gamma}^\varepsilon := \sum_{1 \leq i, j \leq k} \tilde{\gamma}(V_i \times V_j) \delta_{(c_i, c_j)}.$$

Notice that since γ and $\tilde{\gamma}$ have the same marginals, the same holds for γ^ε and $\tilde{\gamma}^\varepsilon$. In particular, one has

$$(3.3) \quad (x, y) \in \text{supp}(\gamma^\varepsilon) \Rightarrow \text{there exists } \tilde{x} \in C \text{ such that } (\tilde{x}, y) \in \text{supp}(\tilde{\gamma}^\varepsilon)$$

and

$$(3.4) \quad (\tilde{x}, \tilde{y}) \in \text{supp}(\tilde{\gamma}^\varepsilon) \Rightarrow \text{there exists } y \in C \text{ such that } (\tilde{x}, y) \in \text{supp}(\gamma^\varepsilon).$$

The following properties will also be useful in our argument.

Claim 1. There exists (x_0, y_0) in the support of γ^ε such that

$$|y_0 - x_0| \geq 5\varepsilon + \max\{|y - x| : (x, y) \in \text{supp}(\tilde{\gamma}^\varepsilon)\}.$$

Claim 2. For any $n \geq 1$, $(x_1, y_1), \dots, (x_n, y_n) \in \text{supp}(\gamma^\varepsilon)$, and $\sigma \in \mathcal{S}_n$,

$$\max_{1 \leq i \leq n} |y_i - x_i| \leq 4\varepsilon + \max_{1 \leq i \leq n} |y_{\sigma(i)} - x_i|.$$

Above, the first claim is simply a counterpart of the antithesis (3.2) for the discretized measures, while the second says that γ^ε is “almost” infinitely cyclically monotone. We postpone the verification of these two claims until the end of this proof.

Let $(x_0, y_0) \in \text{supp}(\gamma^\varepsilon)$ be given by Claim 1. Owing to (3.3) and (3.4), we can recursively define two sequences $(D_m)_{m \geq 1}$ and $(E_m)_{m \geq 0}$ of subsets of C by setting $E_0 := \{y_0\}$, and for $m \geq 1$,

$$D_m := \{\tilde{x} : \text{there exists } y \in E_{m-1} \text{ such that } (\tilde{x}, y) \in \text{supp}(\tilde{\gamma}^\varepsilon)\}$$

and

$$E_m := \{y : \text{there exists } \tilde{x} \in D_m \text{ such that } (\tilde{x}, y) \in \text{supp}(\gamma^\varepsilon)\}.$$

We then set $D := \bigcup_{m \geq 1} D_m$ and $E := \bigcup_{m \geq 0} E_m$.

There are now two alternatives: either x_0 belongs to D or not.

First case: $x_0 \in D$. In this case, there exists $m \geq 1$ such that $x_0 \in D_m$, and by going backwards from D_m to E_0 it is possible to define two finite families $(x_i)_{0 \leq i \leq m}$ and $(y_i)_{0 \leq i \leq m-1}$ such that

$$\text{for all } i \in \{0, \dots, m-1\}, \quad (x_i, y_i) \in \text{supp}(\gamma^\varepsilon) \quad \text{and} \quad (x_{i+1}, y_i) \in \text{supp}(\tilde{\gamma}^\varepsilon),$$

where we have set $x_m := x_0$. Claim 2 then yields

$$\max_{0 \leq i \leq m-1} |y_i - x_i| - 4\varepsilon \leq \max_{0 \leq i \leq m-1} |y_i - x_{i+1}|.$$

Since $\max_{0 \leq i \leq m-1} |y_i - x_i| \geq |y_0 - x_0|$, we infer from Claim 1 and the previous inequality that

$$\max \{|y - x| : (x, y) \in \text{supp}(\tilde{\gamma}^\varepsilon)\} + \varepsilon \leq \max_{0 \leq i \leq m-1} |y_i - x_{i+1}|.$$

Since $(x_{i+1}, y_i) \in \text{supp}(\tilde{\gamma}^\varepsilon)$ for any $i \in \{0, \dots, m-1\}$, this yields a contradiction.

Second case: $x_0 \notin D$. From the definitions of D and E , we notice the following two facts:

$$(3.5) \quad x \in D, (x, y) \in \text{supp}(\gamma^\varepsilon) \quad \Rightarrow \quad y \in E$$

and

$$(3.6) \quad \tilde{y} \in E, (\tilde{x}, \tilde{y}) \in \text{supp}(\tilde{\gamma}^\varepsilon) \quad \Rightarrow \quad \tilde{x} \in D.$$

As a consequence of (3.5) and since γ^ε and $\tilde{\gamma}^\varepsilon$ have the same marginals, one has

$$\gamma^\varepsilon(D \times E) = \gamma^\varepsilon(D \times C) = \tilde{\gamma}^\varepsilon(D \times C).$$

Similarly, one has

$$\tilde{\gamma}^\varepsilon(D \times E) = \tilde{\gamma}^\varepsilon(C \times E) = \gamma^\varepsilon(C \times E).$$

We then obtain

$$\gamma^\varepsilon(D \times E) = \tilde{\gamma}^\varepsilon(D \times C) \geq \tilde{\gamma}^\varepsilon(D \times E) = \gamma^\varepsilon(C \times E).$$

This implies that $\gamma^\varepsilon((C \setminus D) \times E) = 0$, whereas by hypothesis one has $(x_0, y_0) \in (C \setminus D) \times E$ and $\gamma^\varepsilon(\{(x_0, y_0)\}) > 0$ since (x_0, y_0) belongs to the support of the discrete measure γ^ε . This yields a contradiction.

To complete the proof of Theorem 3.4, it remains to prove Claims 1 and 2.

Proof of Claim 1. We infer from (3.2) that

$$\gamma\left(\left\{(x, y) : |y - x| \geq 9\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|\right\}\right) > 0.$$

As a consequence, there exist $i_1, i_2 \in \{1, \dots, k\}$ such that

$$\gamma\left(\left((V_{i_1} \times V_{i_2}) \cap \left\{(x, y) : |y - x| \geq 9\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|\right\}\right)\right) > 0.$$

Since $V_m \subset B(c_m, \varepsilon)$ for $m = i_1, i_2$, one then has

$$(3.7) \quad |c_{i_2} - c_{i_1}| \geq 7\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|,$$

and (c_{i_1}, c_{i_2}) belongs to the support of γ^ε . On the other hand, if (c_{j_1}, c_{j_2}) belongs to the support of $\tilde{\gamma}^\varepsilon$, then $\tilde{\gamma}(V_{j_1} \times V_{j_2}) > 0$, and thus

$$|c_{j_2} - c_{j_1}| \leq 2\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|.$$

Since this inequality holds whenever $(c_{j_1}, c_{j_2}) \in \text{supp}(\tilde{\gamma}^\varepsilon)$, one has

$$\max\{|y - x| : (x, y) \in \text{supp}(\tilde{\gamma}^\varepsilon)\} \leq 2\varepsilon + \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|.$$

This together with (3.7) shows that $(x_0, y_0) := (c_{i_1}, c_{i_2})$ has the desired property.

Proof of Claim 2. Let $n \geq 1$, $\sigma \in \mathcal{S}_n$, and (x_i, y_i) belong to the support of γ^ε for $i \in \{1, \dots, n\}$. For any $i \in \{1, \dots, n\}$, one has $(x_i, y_i) = (c_{j_1}, c_{j_2})$ for some $j_1, j_2 \in \{1, \dots, k\}$ with $\gamma(V_{j_1} \times V_{j_2}) > 0$, and thus there exists $(x'_i, y'_i) \in (V_{j_1} \times V_{j_2}) \cap \text{supp}(\gamma)$. As a consequence,

$$||y'_r - x'_s| - |y_r - x_s|| \leq 2\varepsilon \quad \text{for all } r, s \in \{1, \dots, n\}.$$

Since γ is infinitely cyclically monotone, we have

$$\max_{1 \leq i \leq n} |y'_{\sigma(i)} - x'_i| \geq \max_{1 \leq i \leq n} |y'_i - x'_i|.$$

It follows that

$$4\varepsilon + \max_{1 \leq i \leq n} |y_{\sigma(i)} - x_i| \geq \max_{1 \leq i \leq n} |y_i - x_i|,$$

which proves the claim. \square

Remark 3.5. Observe that in the proof above, we in fact always have $x_0 \in D$; that is, the *first case* always occurs. This is a consequence of the conservation of the masses: all the mass transported to E by $\tilde{\gamma}^\varepsilon$ originates from D , while all the mass in

D is transported to E by γ^ε . Since γ^ε and $\tilde{\gamma}^\varepsilon$ have the same marginals, this implies that both plans transport D exactly to E . The fact that the *second case* never occurs in the above proof also underlies the recursive construction of *Step II* in the proof of Theorem A in [25], even if the arguments are different. That paper deals with the sufficiency of cyclical monotonicity for optimality in the classical case (see also [31]).

Remark 3.6. There is a variation of the self-contained proof given above that relies directly on the fact that the total cost $\mathcal{C}_\infty(\lambda) = \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x|$ depends only on the support of λ and not on its density. In the discrete case this means that, as far as the total cost is concerned, the exact amount of mass transferred from any given point to another is irrelevant: it matters only whether the amount is positive or not. Hence in the proof we are allowed to change the transport plans γ_ε and $\tilde{\gamma}_\varepsilon$, along with their marginals, as long as we do not change their supports and make sure that the marginals of the new transport plans agree with each other. Now assuming that we can change the measures in such a way that all the point masses are of integer size, the problem can be interpreted as a pairing problem in which the infinite cyclical monotonicity is both a necessary and sufficient condition for optimality.

The existence of the required integer transport plans with given supports is non-trivial and follows from Dines' algorithm [17], which provides positive solutions for a system of linear equations.

4. Restrictable solutions. In the previous section, we derived the notion of infinitely cyclically monotone plans from the approximation of the problem (P_∞) by the family of problems (P_p) . Another interesting notion may be derived in the same way: let $1 \leq p < \infty$ and $\gamma_p \in \Pi(\mu, \nu)$ be a solution to (P_p) . Then it follows from the linearity of the functional

$$\gamma \mapsto \mathcal{C}_p(\gamma)^p = \int_{\Omega^2} |y - x|^p d\lambda(x, y)$$

that any nonzero measure γ' that is majorized by γ_p , i.e., $\gamma'(B) \leq \gamma_p(B)$ for all Borel sets $B \subset \Omega \times \Omega$, is an optimal transport plan for the problem

$$\mathcal{C}_p(\gamma') = \inf \{ \mathcal{C}_p(\gamma) : \gamma \in \Pi(\mu', \nu') \},$$

where $\mu' := \pi_{1\#}\gamma'$ and $\nu' := \pi_{2\#}\gamma'$. In other words, optimality is automatically inherited by restriction, and hence we may say that a solution $\gamma_p \in \Pi(\mu, \nu)$ of (P_p) is a *restrictable solution* of this problem. By analogy, we may define a similar notion of restrictable solutions for problem (P_∞) as follows.

DEFINITION 4.1. *A transport plan $\gamma \in \Pi(\mu, \nu)$ is a restrictable solution of (P_∞) if any nonzero Borel measure γ' in $\Omega \times \Omega$ that is majorized by γ is a solution to the problem*

$$(P'_\infty) \quad \inf \left\{ \lambda\text{-ess sup}_{(x,y) \in \Omega^2} |y - x| : \lambda \in \Pi(\mu', \nu') \right\},$$

where $\mu' := \pi_{1\#}\gamma'$ and $\nu' := \pi_{2\#}\gamma'$.

The reader should notice the obvious abuse of notation above, as the measures μ' and ν' in (P'_∞) are not, in general, probability measures. However, $\mu'(\Omega) = \nu'(\Omega) > 0$, and that is really all that is needed.

Example 4.2. It is quite clear that not every solution of (P_∞) is restrictable. Indeed, the optimal plan λ considered in Example 3.3 admits the following restriction:

$$\lambda' := \frac{1}{2} (\mathcal{L}^2 \llcorner_{S_1} \times \frac{1}{2} \delta_{(2,1)} + \mathcal{L}^2 \llcorner_{S_2} \times \frac{1}{2} \delta_{(1,2)}),$$

where $S_1 := [0, \frac{1}{2}] \times [\frac{1}{2}, 1]$ and $S_2 := [\frac{5}{2}, 3] \times [2, \frac{5}{2}]$. But λ' is not optimal for its own marginals: a better transport plan is the one that takes all the mass that lies in S_1 to $(1, 2)$ and the mass in S_2 to $(2, 1)$.

Remark 4.3. The notion of a restrictable solution bears a strong resemblance to that of an *absolute minimizer* used in connection with the L^∞ variational problems. We recall that a locally Lipschitz continuous function $u : D \rightarrow \mathbb{R}^m$, $m \geq 1$, is called an absolute minimizer of the functional $S(\varphi, D) := \text{ess sup}_{x \in D} H(x, \varphi(x), D\varphi(x))$ if

$$S(u, V) \leq S(v, V)$$

for every open $V \subset\subset D$ and $v \in W^{1,\infty}(V) \cap C(\bar{V})$ such that $v|_{\partial V} = u|_{\partial V}$. Absolute minimizers were introduced by Aronsson [3, 4]. It has turned out that this is the proper notion of a solution for this type of minimization problem in the sense that important properties such as uniqueness, regularity, and characterization via an Euler–Lagrange equation can be obtained for this class of functions. Absolute minimizers were introduced by Aronsson [3]; see also, e.g., [6, 5, 30, 14] for further details and background.

It is natural to ask whether any cluster point γ_∞ of $(\gamma_p)_{p \geq 1}$ in $\Pi(\mu, \nu)$ as $p \rightarrow \infty$ is a restrictable solution of (P_∞) . In view of Theorem 3.2, this can be established by showing that the restrictable solutions coincide with the class of infinitely cyclically monotone solutions.

THEOREM 4.4. *A transport plan $\gamma \in \Pi(\mu, \nu)$ is infinitely cyclically monotone if and only if it is a restrictable solution of the problem (P_∞) .*

Notice that Theorem 3.2 provides the existence of restrictable solutions to (P_∞) .

Proof. If $\gamma \in \Pi(\mu, \nu)$ is infinitely cyclically monotone, then the same holds for any restriction $\gamma' \leq \gamma$, and Theorem 3.4 then yields that such a restriction γ' is a solution of the corresponding problem (P'_∞) .

We now turn to the proof of the sufficiency. Let $\gamma \in \Pi(\mu, \nu)$ be a restrictable solution to (P_∞) , and let us fix points $(x_1, y_1), \dots, (x_m, y_m) \in \text{supp}(\gamma)$, $m \geq 2$, and a permutation σ of $\{1, \dots, m\}$. Without loss of generality, we may assume that $(x_i, y_i) \neq (x_j, y_j)$ whenever $i \neq j$. Then there is $\varepsilon_0 > 0$ such that for all $0 < \varepsilon \leq \varepsilon_0$, the sets $B_i := B(x_i, \varepsilon) \times B(y_i, \varepsilon)$ are pairwise disjoint and $\gamma(B_i) > 0$ for all $i = 1, \dots, m$.

We define two measures γ' and γ_σ by setting

$$\gamma' := \sum_{i=1}^m c_i \gamma|_{B_i} \quad \text{and} \quad \gamma_\sigma := \sum_{i=1}^m c_i T_{\#}^i \gamma|_{B_i}.$$

Here $T_{\#}^i \gamma|_{B_i}$ is the push-forward of $\gamma|_{B_i}$ by the mapping $T^i(x, y) := (x, y + y_{\sigma(i)} - y_i)$, and the positive numbers

$$c_i := \frac{\min_k \gamma(B_k)}{\gamma(B_i)}$$

are chosen so that $\gamma'(B_i) = \min_k \gamma(B_k) > 0$ is independent of i . Observe that the support of $T_{\#}^i \gamma|_{B_i}$ lies in $B_i^\sigma := B(x_i, \varepsilon) \times B(y_{\sigma(i)}, \varepsilon)$ and $\gamma_\sigma(B_i^\sigma) = \gamma'(B_i)$. Moreover, the first marginals $\mu' = \pi_{1\#} \gamma'$ and $\mu_\sigma = \pi_{1\#} \gamma_\sigma$ are equal.

Since γ' is majorized by the restrictable solution γ , we have

$$\gamma'\text{-ess sup}_{(x,y) \in \Omega^2} |x - y| = W_\infty(\mu', \nu') = \inf \left\{ \tilde{\gamma}\text{-ess sup}_{(x,y) \in \Omega^2} |x - y| : \tilde{\gamma} \in \Pi(\mu', \nu') \right\},$$

where $\nu' = \pi_2 \# \gamma'$. On the other hand, since the supports of both ν' and $\nu_\sigma = \pi_2 \# \gamma_\sigma$ are contained in the union of the balls $B(y_i, \varepsilon)$ and $\nu'(B(y_i, \varepsilon)) = \nu_\sigma(B(y_i, \varepsilon))$ for every i by the construction of γ' and γ_σ , we can rearrange ν' to ν_σ by transporting mass only within the balls $B(y_i, \varepsilon)$. Thereby we obtain that $W_\infty(\nu', \nu_\sigma) \leq 2\varepsilon$, and hence

$$\begin{aligned} \gamma'\text{-ess sup}_{(x,y) \in \Omega^2} |x - y| &= W_\infty(\mu', \nu') \leq W_\infty(\mu_\sigma, \nu_\sigma) + W_\infty(\nu_\sigma, \nu') \\ &\leq 2\varepsilon + \gamma_\sigma\text{-ess sup}_{(x,y) \in \Omega^2} |x - y|. \end{aligned}$$

Now clearly

$$\gamma'\text{-ess sup}_{(x,y) \in \Omega^2} |x - y| \geq \max_{1 \leq i \leq m} |x_i - y_i|$$

and

$$\gamma_\sigma\text{-ess sup}_{(x,y) \in \Omega^2} |x - y| \leq \max_{1 \leq i \leq m} |x_i - y_{\sigma(i)}| + 2\varepsilon,$$

and thus the preceding inequality yields

$$\max_{1 \leq i \leq m} |x_i - y_i| \leq \max_{1 \leq i \leq m} |x_i - y_{\sigma(i)}| + 4\varepsilon.$$

Since this holds for all $\varepsilon > 0$ small enough we are done. \square

Observe that for any $\gamma \in \Pi(\mu, \nu)$ and any Borel set $B \subset \Omega \times \Omega$ such that $\gamma(B) > 0$, the measure $\gamma|_B$ is majorized by γ . Thus, if γ is a restrictable solution of (P_∞) , then each such measure $\gamma|_B$ is an optimal transport plan for its own marginals. It turns out that in general the converse is false; that is, this family of measures alone does not suffice for characterizing the restrictable solutions, as shown in Example 4.5. Notice that this is another difference with the case of integral costs functionals (like the costs \mathcal{C}_p), since for those functionals the converse would be true: if $\gamma|_B$ is an optimal transport plan for its own marginals for any $B \subset \Omega \times \Omega$ with $\gamma(B) > 0$, then γ is a restrictable solution.

Example 4.5. Take $\Omega = [0, 1]$, and let

$$\mu := \frac{1}{3}\delta_0 + \frac{2}{3}\delta_1 \quad \text{and} \quad \nu := \frac{2}{3}\delta_0 + \frac{1}{3}\delta_1.$$

Then the plan $\gamma := \frac{1}{3}\delta_{(0,1)} + \frac{2}{3}\delta_{(1,0)}$ is not a restrictable solution since $\gamma' := \frac{1}{3}\delta_{(0,1)} + \frac{1}{3}\delta_{(1,0)}$ is majorized by γ , and it is clearly not an optimal transport plan for its own marginals $\mu' = \frac{1}{3}\delta_0 + \frac{1}{3}\delta_1$ and $\nu' = \frac{1}{3}\delta_0 + \frac{1}{3}\delta_1$. On the other hand, one can check that $\gamma|_B$ is an optimal transport plan, for its own marginals, for each Borel set $B \subset \Omega^2$. Notice that the only restrictable solution of (P_∞) in this case (which is also the unique solution to (P_p) when $p \geq 1$) is $\gamma_\infty = \frac{1}{3}\delta_{(0,0)} + \frac{1}{3}\delta_{(1,0)} + \frac{1}{3}\delta_{(1,1)}$.

In view of the above example, and under some regularity assumption on the measures μ and ν , it is possible to obtain the following refined version of Theorem 4.4.

PROPOSITION 4.6. *Let $\gamma \in \Pi(\mu, \nu)$ be an optimal transport plan, and assume that neither μ nor ν concentrates on sets of dimension $d - 1$. Then the following are equivalent:*

- (1) γ is infinitely cyclically monotone;
- (2) for each Borel set $B \subset \Omega \times \Omega$, $\gamma|_B$ is optimal between its projections.

Proof. We need only prove that under these assumptions (2) implies (1). Assume by contradiction that γ is not infinitely cyclically monotone. Then there exist a family $\{(x_i, y_i)\}_{i=1, \dots, n}$ in $\text{supp}(\gamma)$ and a permutation $\sigma \in \mathcal{S}_n$ such that

$$\max\{|x_1 - y_{\sigma(1)}|, \dots, |x_n - y_{\sigma(n)}|\} < \max\{|x_1 - y_1|, \dots, |x_n - y_n|\}.$$

By continuity, the same inequality holds true for any family $\{(x'_i, y'_i)\}_{i=1, \dots, n}$ for which $(x'_i, y'_i) \in B(x_i, \varepsilon) \times B(y_i, \varepsilon)$ for all $i = 1, \dots, n$ and for some small enough $\varepsilon > 0$. Notice that $\gamma(B(x_i, \varepsilon) \times B(y_i, \varepsilon))$ is positive for any i .

For each i we define $g_i : [0, \varepsilon] \rightarrow \mathbb{R}^+$ by $g_i(r) := \gamma(B(x_i, r) \times B(y_i, r))$. Since μ and ν do not concentrate on sets of dimension $d - 1$, the function g_i is continuous. Let $\alpha = \min_i \gamma(B(x_i, \varepsilon) \times B(y_i, \varepsilon))$, and choose $0 < \tilde{\varepsilon}_i \leq \varepsilon$ so that $g_i(\tilde{\varepsilon}_i) = \alpha$ for all i ; this is possible since $\lim_{r \rightarrow 0} g_i(r) = 0$.

Following now the proof of sufficiency of the previous theorem we obtain that $\gamma|_B$ for $B = \bigcup_{i=1}^n B(x_i, \tilde{\varepsilon}_i) \times B(y_i, \tilde{\varepsilon}_i)$ and for $\varepsilon > 0$ small enough is not an optimal transport between its marginals, which contradicts (2). \square

Remark 4.7. It is natural to ask what happens if $|x - y|$ is replaced by a more general (real-valued) cost function $c(x, y)$; that is, we consider the functional

$$\gamma \mapsto \gamma\text{-ess sup}_{(x, y) \in \Omega^2} c(x, y).$$

As expected, the basic existence results, Proposition 2.1 and Theorem 3.2, remain valid, provided that c is nonnegative and lower semicontinuous with all the relevant concepts appropriately redefined: in particular, in the definition of *infinite cyclical monotonicity* one should replace the support of γ with some appropriate set on which γ is concentrated. Moreover, the proof of the equivalence of restrictable and infinitely cyclically monotone solutions works in this generality if $c(x, y)$ is (uniformly) continuous.

5. Existence and uniqueness of an optimal transport map. In this section, we prove that under reasonably weak assumptions an infinitely cyclically monotone transport plan is induced by a transport map. Moreover, we start the analysis of the uniqueness of such transport maps and then comment on our method of proof in light of the duality issue for problem (P_∞) .

5.1. Properties of transport plans. We begin by considering some generic properties of transport plans. This subsection is largely independent of the cost, and the technique detailed below has applications also in the framework of classical transportation problems involving cost functionals in integral form; see [15].

DEFINITION 5.1. *Let $y \in \Omega$ and $r > 0$, and let $\gamma \in \Pi(\mu, \nu)$ be a transport plan. We define*

$$\gamma^{-1}(B(y, r)) := \pi^1((\Omega \times B(y, r)) \cap \text{supp } \gamma).$$

In other words, $\gamma^{-1}(B(y, r))$ is the set of points whose mass is partially or completely transported to $B(y, r)$ by γ . We recognize the slight abuse of notation, but if γ is thought of as a device that transports mass, then this seems justifiable. Notice also that $\gamma^{-1}(B(y, r))$ is a Borel set. In fact, it is a countable union of compact sets as shown by the equation

$$\begin{aligned}\pi^1((\Omega \times B(y, r)) \cap \text{supp } \gamma) &= \pi^1\left(\bigcup_n (\Omega \times \overline{B(y, r - 1/n)}) \cap \text{supp } \gamma\right) \\ &= \bigcup_n \pi^1((\Omega \times \overline{B(y, r - 1/n)}) \cap \text{supp } \gamma).\end{aligned}$$

Since this notion is important in what follows, we recall that when A is \mathcal{L}^d -measurable, one has

$$\lim_{r \rightarrow 0} \frac{\mathcal{L}^d(A \cap B(x, r))}{\mathcal{L}^d(B(x, r))} = 1$$

for almost every x in A : we shall call such a point x a Lebesgue point of A , this terminology deriving from the fact that such a point may also be considered as a Lebesgue point of χ_A .

The following lemma, although quite simple, is the cornerstone of the proof of Theorem 5.5 below.

LEMMA 5.2. *Let $\gamma \in \Pi(\mu, \nu)$, and assume that $\mu \ll \mathcal{L}^d$. Then γ is concentrated on a σ -compact set $R(\gamma)$ such that for all $(x, y) \in R(\gamma)$ the point x is a Lebesgue point of $\gamma^{-1}(B(y, r))$ for all $r > 0$.*

Proof. In the following, we shall denote by $\text{Leb}(E)$ the set of points $x \in E$ which are Lebesgue points of E . Let

$$A := \{(x, y) \in \text{supp}(\gamma) : x \notin \text{Leb}(\gamma^{-1}(B(y, r))) \text{ for some } r > 0\};$$

we first intend to show that $\gamma(A) = 0$. To this end, for each positive integer n we consider a finite covering $\Omega \subset \bigcup_{i \in I(n)} B(y_i^n, \frac{1}{2n})$ by balls of radius $\frac{1}{2n}$. We notice that if $(x, y) \in \text{supp}(\gamma)$ and x is not a Lebesgue point of $\gamma^{-1}(B(y, r))$ for some $r > 0$, then for any $n \geq \frac{1}{r}$ and y_i^n such that $|y_i^n - y| < \frac{1}{2n}$ the point x belongs to $\gamma^{-1}(B(y_i^n, \frac{1}{2n}))$ but is not a Lebesgue point of this set. Then

$$\pi^1(A) \subset \bigcup_{n \geq 1} \bigcup_{i \in I(n)} \left(\gamma^{-1} \left(B \left(y_i^n, \frac{1}{2n} \right) \right) \setminus \text{Leb} \left(\gamma^{-1} \left(B_\gamma \left(y_i^n, \frac{1}{2n} \right) \right) \right) \right).$$

Notice that the set on the right-hand side has Lebesgue measure 0 and thus μ -measure 0. It follows that $\gamma(A) \leq \gamma(\pi^1(A) \times \Omega) = \mu(\pi^1(A)) = 0$.

Finally, since $\mathcal{L}^d(\pi^1(A)) = 0$, there exists a sequence $(U_k)_{k \geq 0}$ of open sets such that

$$\text{for all } k \geq 0, \quad \pi^1(A) \subset U_k \quad \text{and} \quad \lim_{k \rightarrow \infty} \mathcal{L}^d(U_k) = 0.$$

Then the set $R(\gamma) := \text{supp}(\gamma) \cap (\bigcup_{k \geq 0} (\Omega \setminus U_k) \times \Omega)$ has the desired properties. \square

The above lemma yields the introduction of the following notion.

DEFINITION 5.3. *The couple $(x, y) \in \Omega \times \Omega$ is a γ -regular point if $x \in \gamma^{-1}(B(y, r))$ is a Lebesgue point of this set for any positive r .*

Notice that any element of the set $R(\gamma)$ of Lemma 5.2 is a γ -regular point.

For future use, we introduce a suitable notation to indicate a cone: let $x_0, \xi \in \mathbb{R}^d$ with $|\xi| = 1$, and let $\delta \in [0, 2]$. Then we define

$$C(x_0, \xi, \delta) := \left\{ x \in \mathbb{R}^d \setminus \{x_0\} : \frac{x - x_0}{|x - x_0|} \cdot \xi \geq 1 - \delta \right\}.$$

Notice that if $\delta = 0$, $C(x_0, \xi, 0)$ degenerates to a half-line, while in the case $\delta = 2$, $C(x_0, \xi, 2)$ is $\mathbb{R}^d \setminus \{x_0\}$.

We now remark the following property for the regular points of a transport plan.

PROPOSITION 5.4. *Let (x_0, y_0) be a γ -regular point, $r > 0$, $\alpha \in (0, 1)$, and $\delta > 0$. Then for $\varepsilon > 0$ sufficiently small the set of points $x \in \gamma^{-1}(B(y_0, r))$ such that $x \in C(x_0, \xi, \delta) \cap (B(x_0, \varepsilon) \setminus B(x_0, \alpha\varepsilon))$ has positive \mathcal{L}^d measure.*

Proof. It is enough to remark that $x_0 \in \text{Leb}(\gamma^{-1}(B(y_0, r)))$ and then

$$\lim_{\varepsilon \rightarrow 0} \frac{\mathcal{L}^d((B(x_0, \varepsilon) \setminus B(x_0, \alpha\varepsilon)) \cap C(x_0, \xi, \delta) \cap \gamma^{-1}(B(y_0, r)))}{\mathcal{L}^d(B(x_0, \varepsilon))} = c(\alpha, \delta),$$

where $c(\alpha, \delta) := \frac{\mathcal{L}^d((B(x_0, 1) \setminus B(x_0, \alpha)) \cap C(x_0, \xi, \delta))}{\mathcal{L}^d(B(x_0, 1))} > 0$. \square

5.2. Existence of an optimal transport map. Our main result in this subsection is the following theorem, which states that under the hypothesis that μ is absolutely continuous with respect to the Lebesgue measure \mathcal{L}^d , any optimal infinitely cyclically monotone transport plan for (P_∞) is induced by a transport map. This generalizes the corresponding result for the problem (P_p) when $p \in]1, \infty[$: if one assumes that

$$(5.1) \quad \mu(B) = 0 \quad \text{whenever} \quad \mathcal{H}^{d-1}(B) < \infty,$$

then any p -cyclically monotone transport plan is induced by a transport map (see Remark 4.7 in [20]).

It is in doubt whether (5.1) is sufficient to ensure that the conclusion of Theorem 5.5 holds. In the case $p = 1$, the hypothesis (5.1) is not sufficient to ensure that any 1-cyclically monotone transport plan is induced by a transport map. Even worse, for $\mu = \mathcal{L}^1|_{[0,1]}$ and $\nu = \frac{1}{2}\mathcal{L}^1|_{[0,2]}$ there exists an optimal (and hence 1-cyclically monotone) transport plan which is not induced by a transport map; see [2, p. 125].

We now state the main result of this section and refer the reader to section 5.4 for further comments.

THEOREM 5.5. *Assume that $\mu \ll \mathcal{L}^d$, and let $\gamma \in \Pi(\mu, \nu)$ be an infinitely cyclically monotone transport plan. Then there exists a Borel transport map $T : \Omega \rightarrow \Omega$ such that $\gamma = (id \times T)_\# \mu$.*

Proof. By Proposition 2.1 in [1], it is sufficient to prove that γ is concentrated on a γ -measurable graph. In view of Lemma 5.2, it is then sufficient to prove that $R(\gamma)$ is included in a graph, or more generally that if (x_0, y_0) and (x_0, y'_0) are both γ -regular points, then $y_0 = y'_0$.

We divide the proof into two parts and first show that $|x_0 - y_0| = |x_0 - y'_0|$. Arguing by contradiction, we assume that $|x_0 - y_0| < |x_0 - y'_0|$ and suppose for the time being that $x_0 \neq y_0$. Let $\xi' = \frac{y'_0 - x_0}{|y'_0 - x_0|}$, $0 < \varepsilon < |x_0 - y_0|$, and $0 < r < |y'_0 - x_0| - |y_0 - x_0|$. We claim that for $\delta := 1 - \frac{|x_0 - y_0|}{|x_0 - y'_0|}$ one has

$$(5.2) \quad \max\{|x - y'_0|, |x_0 - y|\} < \max\{|x - y|, |x_0 - y'_0|\}$$

for any (x, y) such that $x \in C(x_0, \xi', \delta) \cap B(x_0, \varepsilon) \setminus B(x_0, \frac{1}{2}\varepsilon)$ and $y \in B(y_0, r)$. Indeed, take (x, y) as above; it then follows from the choice of r that $|x_0 - y| < |x_0 - y'_0|$, while on the other hand

$$\begin{aligned} |x - y'_0|^2 &= |x - x_0| \left(|x - x_0| - 2 \frac{x - x_0}{|x - x_0|} \cdot (y'_0 - x_0) \right) + |x_0 - y'_0|^2 \\ &\leq |x - x_0| (|x_0 - y_0| - 2(1 - \delta)|x_0 - y'_0|) + |x_0 - y'_0|^2 < |x_0 - y'_0|^2. \end{aligned}$$

This proves the claim. We now infer from Proposition 5.4 that the set of points $x \in \gamma^{-1}(B(y_0, r))$ such that $x \in C(x_0, \xi', \delta) \cap B(x_0, \varepsilon) \setminus B(x_0, \frac{1}{2}\varepsilon)$ has positive measure when ε is small enough. In particular, this set is nonempty for small ε , and (5.2) then clearly contradicts the infinite cyclical monotonicity of γ . As a consequence, $|x_0 - y_0| = |x_0 - y'_0|$ in the case $x_0 \neq y_0$.

If $x_0 = y_0$, we repeat the argument above with the choices $0 < \varepsilon < \frac{1}{4}|x_0 - y'_0|$, $0 < r < \frac{1}{4}|y'_0 - x_0|$, and $\delta = \frac{1}{2}$. Then for any (x, y) such that $x \in C(x_0, \xi', \delta) \cap B(x_0, \varepsilon) \setminus B(x_0, \frac{1}{2}\varepsilon)$ and $y \in B(y_0, r)$, we clearly have $|x_0 - y| < |x_0 - y'_0|$, and the other inequality $|x - y'_0|^2 < |x_0 - y'_0|^2$ follows as above.

We now prove by contradiction that $y_0 = y'_0$. Note that since we already know that $|x_0 - y_0| = |x_0 - y'_0|$, we may assume that $|x_0 - y_0| = |x_0 - y'_0| > 0$. If $y_0 \neq y'_0$, we can find $\xi \in \mathbb{R}^d$ such that

$$\xi \cdot \frac{x_0 - y_0}{|x_0 - y_0|} < 0 \quad \text{and} \quad \xi \cdot \frac{x_0 - y'_0}{|x_0 - y'_0|} > 0.$$

Next we choose $r > 0$ such that

$$\sup \left\{ \xi \cdot \frac{x_0 - y}{|x_0 - y|} : y \in B(y_0, r) \right\} < 0$$

and $\delta > 0$ such that

$$(5.3) \quad \alpha := \inf \left\{ \frac{x_0 - x}{|x_0 - x|} \cdot \frac{x_0 - y'_0}{|x_0 - y'_0|} : x \in C(x_0, \xi, \delta) \right\} > 0$$

as well as

$$(5.4) \quad \sup \left\{ \frac{x_0 - x}{|x_0 - x|} \cdot \frac{x_0 - y}{|x_0 - y|} : x \in C(x_0, \xi, \delta), y \in B(y_0, r) \right\} < 0.$$

We now claim that for $\varepsilon > 0$ small enough, (5.2) holds for any (x, y) such that $x \in C(x_0, \xi, \delta) \cap B(x_0, \varepsilon) \setminus B(x_0, \frac{1}{2}\varepsilon)$ and $y \in B(y_0, r)$. Notice that this claim concludes the proof of $y_0 = y'_0$ modulo applying Proposition 5.4 as before. To verify that (5.2) holds, we first notice that (5.4) implies that

$$(5.5) \quad \text{for all } x \in C(x_0, \xi, \delta), y \in B(y_0, r), \quad |x_0 - y| < |x - y|$$

since $|x - y|^2 = |x - x_0|^2 - 2(x_0 - x) \cdot (x_0 - y) + |x_0 - y|^2$. We can also infer from (5.3) that

$$|x - y'_0|^2 \leq |x - x_0| (|x_0 - x| - 2\alpha|x_0 - y'_0|) + |x_0 - y'_0|^2$$

for any $x \in C(x_0, \xi, \delta)$. It follows that

$$(5.6) \quad \text{for all } x \in C(x_0, \xi, \delta) \cap B(x_0, \varepsilon), \quad |x - y'_0| < |x_0 - y'_0|$$

whenever $0 < \varepsilon < 2\alpha|x_0 - y'_0|$. We then get (5.2) from (5.5) and (5.6), which concludes the proof. \square

5.3. Uniqueness of the infinitely cyclically monotone transport map.

We now consider the question of the uniqueness of the infinitely cyclically monotone transport map obtained in the preceding section. We recall that when (5.1) holds and $p \in]1, \infty[$, problem (P_p) admits a unique (up to a μ negligible set) p -cyclically

monotone transport map (see, for example, [20] and section 5.4). Notice that in contrast this result does not hold for $p = 1$, not even under the stronger hypothesis that $\mu \ll \mathcal{L}^d$, as shown by Example 1.3 in [1]: when $\mu = \mathcal{L}^1|_{[0,1]}$ and $\nu = \mathcal{L}^1|_{[1,2]}$, both transport maps $t \mapsto t + 1$ and $t \mapsto 2 - t$ are optimal.

In the case of problem (P_∞) , the question of uniqueness is largely open. At the moment, we have only the following partial result stating the uniqueness of the infinitely cyclically monotone transport map under the hypothesis that ν is purely atomic with finite support.

THEOREM 5.6. *Suppose that $\mu \ll \mathcal{L}^d$ and $\nu = \sum_{i=0}^k a_i \delta_{y_i}$ for some $(y_i)_{0 \leq i \leq k} \subset \Omega$ and positive numbers a_0, \dots, a_k . Then there exists a unique (up to a μ -negligible set) infinitely cyclically monotone Borel transport map T from μ to ν .*

Proof. Suppose that there are two distinct infinitely cyclically monotone Borel transport maps T and \tilde{T} , and let us introduce the sets

$$U_j^i = T^{-1}(y_j) \cap \tilde{T}^{-1}(y_i).$$

We first claim that it is possible to define a sequence of integers $(i(p))_{p \geq 0}$ such that

$$\text{for all } p \geq 0, \quad i(p) \neq i(p+1) \quad \text{and} \quad \mu\left(U_{i(p)}^{i(p+1)}\right) > 0.$$

Indeed, the fact that the two transport maps are distinct means that it is possible to choose two indices $i(0) \neq i(1)$ such that $\mu\left(U_{i(0)}^{i(1)}\right) > 0$. Next we notice that since \tilde{T} maps μ to ν ,

$$\nu(\{y_{i(1)}\}) = \sum_{p=0}^k \mu\left(U_p^{i(1)}\right) \geq \mu\left(U_{i(0)}^{i(1)}\right) + \mu\left(U_{i(1)}^{i(1)}\right) > \mu\left(U_{i(1)}^{i(1)}\right).$$

Since T also maps μ to ν , we infer from the above inequality that there exists $p \neq i(1)$ such that $\mu\left(U_{i(1)}^p\right) > 0$: we then set $i(2) = p$ and start again from $U_{i(1)}^{i(2)}$. By repeating the above argument we can build recursively the sequence $(i(p))_{p \geq 0}$ with the desired properties.

Since the sequence $(i(p))_{p \geq 0}$ takes its values in the finite set $\{0, \dots, k\}$, we may assume that there exists some $m \geq 2$ such that $i(m) = i(0)$. For any $p \in \{0, \dots, m-1\}$, the set $U_{i(p)}^{i(p+1)}$ has nonzero Lebesgue measure, so we may choose a Lebesgue point x_p of $U_{i(p)}^{i(p+1)}$ for which $|y_{i(p+1)} - x_p| \neq |y_{i(p)} - x_p|$ and then set $x_m = x_0$. By definition,

$$T(x_p) = y_{i(p)} \quad \text{and} \quad \tilde{T}(x_p) = y_{i(p+1)} \quad \text{for all } p \in \{0, \dots, m-1\}.$$

Since T and \tilde{T} are infinitely cyclically monotone, we have

$$\max_{0 \leq p \leq m-1} |y_{i(p)} - x_p| \leq \max_{0 \leq p \leq m-1} |y_{i(p+1)} - x_p| \leq \max_{0 \leq p \leq m-1} |y_{i(p)} - x_p|,$$

so that

$$(5.7) \quad \max_{0 \leq p \leq m-1} |y_{i(p)} - x_p| = \max_{0 \leq p \leq m-1} |y_{i(p+1)} - x_p|.$$

Then let $I := \{q : |y_{i(q)} - x_q| = \max_{0 \leq p \leq m-1} |y_{i(p)} - x_p|\}$. We infer from (5.7) and the choice of the points x_p that for any $q \in I$ one has $|y_{i(q)} - x_q| > |y_{i(q+1)} - x_q|$. Since x_q is a Lebesgue point of $U_{i(q)}^{i(q+1)}$, there exists $\tilde{x} \in U_{i(q)}^{i(q+1)}$ arbitrarily close to x_q for which

$$|y_{i(q)} - \tilde{x}| > |y_{i(q)} - x_q|,$$

and thus we can choose $\tilde{x} \in U_{i(q)}^{i(q+1)}$ such that

$$(5.8) \quad |y_{i(q)} - \tilde{x}| > \max\{|y_{i(q+1)} - \tilde{x}|, |y_{i(q)} - x_q|\}.$$

We now set $\tilde{x}_q := \tilde{x}$ as well as $\tilde{x}_p := x_p$ for $p \neq q$ and notice that (5.7) should in fact also hold for this new choice of elements \tilde{x}_p in $U_{i(p)}^{i(p+1)}$. This leads to a contradiction since we can infer from the definition of I , the choice of q , and (5.8) that

$$\begin{aligned} \max_{0 \leq p \leq m-1} |y_{i(p)} - \tilde{x}_p| &= |y_{i(q)} - \tilde{x}| \\ &> \max\{|y_{i(q+1)} - \tilde{x}|, |y_{i(q)} - x_q|\} \geq \max_{0 \leq p \leq m-1} |y_{i(p+1)} - \tilde{x}_p|. \end{aligned}$$

This is in contradiction with (5.7) and concludes the proof of the theorem. \square

Remark 5.7. The construction of the sequence $(x_p)_{0 \leq p \leq m}$ is close to that proposed in the proof of Theorem 3.4, but it is easier since we do not need that it loops at x_0 . Indeed, in the above proof we do not really need that $x_m = x_0$, and we assume this only for convenience of notation, while in the course of the proof of Theorem 3.4 we intended to use Claim 1 and then had to start from the special x_0 found there.

At the moment we are not able to generalize this result to the case where ν is any probability measure on Ω . On the other hand, it is clear that the above uniqueness theorem requires that μ does not concentrate, as the following example shows.

Example 5.8. Assume that $\mu := \mathcal{H}^1 \llcorner_{[0,1] \times \{0\}}$, while $\nu := \frac{1}{2}(\delta_{(0,-1)} + \delta_{(0,1)})$. Then any transport map T (i.e., any μ -measurable function for which $T_{\#}\mu = \nu$) is an infinitely cyclically monotone optimal transport map from μ to ν .

5.4. Comments around duality. For $1 \leq p < \infty$, the mass transport problem (P_p) may be rewritten as

$$(P_p) \quad W_p^p(\mu, \nu) = \inf \left\{ \mathcal{C}_p^p(\lambda) = \int_{\Omega^2} |y - x|^p d\lambda(x, y) : \lambda \in \Pi(\mu, \nu) \right\}.$$

In this form, the objective functional $\lambda \mapsto \mathcal{C}_p^p(\lambda)$ is linear over the compact convex set $\Pi(\mu, \nu)$, and it is then quite natural to associate with (P_p) its dual problem

$$(D_p) \quad \sup \left\{ \int_{\Omega} \phi(x) d\mu(x) + \int_{\Omega} \psi(y) d\nu(y) : \phi(x) + \psi(y) \leq |y - x|^p \right\},$$

where $\phi \in L^1(d\mu)$, $\psi \in L^1(d\nu)$, and the constraint is required to hold for μ a.e. x and ν a.e. y . Due to the regularity of the integrand $c_p(x, y) := |y - x|^p$, the supremum of (D_p) is achieved for a couple (φ, φ^{c_p}) where the Kantorovich potential φ is continuous and c_p -concave. We refer the reader, for example, to section 3 of [2], Part I of [20], section 3.3 of [26], or section 2.4 of [32] for more on the related concepts and results.

The Kantorovich dual problem (D_p) appears to be a fundamental tool in understanding and solving the problems of the characterization, existence, and uniqueness

for an optimal transport map for (P_p) . For example, the notion of p -cyclical monotonicity (3.1) naturally appears via the equality $\sup(D_p) = \inf(P_p)$; see, e.g., the proof of Theorem 3.2 in [2] (an alternative and direct proof, for example, that of Theorem 2.3 in [20]). Moreover, a key point in the construction of an optimal transport map $T_p : \text{supp}(\mu) \rightarrow \text{supp}(\nu)$ for (P_p) is to identify the directions of transportation (known as *transport rays* for $p = 1$), that is, to associate with μ a.e. $x \in \text{supp}(\mu)$ the direction $\frac{T_p(x)-x}{|T_p(x)-x|}$ to which the mass present at x is transferred. It is now well understood that this direction may be obtained as the adequate c_p -gradient of an optimal Kantorovich potential φ (see, e.g., [20, 19, 27] or section 2.4 of [32]). In fact, the definition as well as the regularity properties of the transport rays are deeply linked with the fact that the support of an optimal transport plan γ_p for (P_p) is p -cyclically monotone and thus inherits good properties from being included in the subdifferential of a c_p -concave function (which, in turn, turns out to be a Kantorovich potential; see, e.g., section 2.4 of [32]).

In light of the preceding discussion, it is natural to try to develop a duality theory for the problem (P_∞) as well. We hereafter informally discuss this issue.

First, in view of the Example 3.3, it does not seem realistic that one could obtain useful information from an intrinsic approach. Indeed, the optimal transport plan λ proposed in Example 3.3 is not induced by any transport map, so we cannot expect that a dual problem directly associated with (P_∞) via some general construction gives information on the geometry of the optimal transport plans. On the other hand, notice that Theorems 5.5 and 5.6 do indicate that there exists a unique infinitely cyclically monotone optimal transport map for the particular problem of Example 3.3.

In view of the results of the two preceding sections, and since the notion of infinitely cyclically monotone plan was at first obtained via a limiting argument, one is led to study the asymptotic behavior of the family of dual problems (D_p) as $p \rightarrow \infty$. But as mentioned above, (D_p) is not directly related to (P_p) but to a reformulation of (P_p) which requires taking the p -power of the objective functional \mathcal{C}_p . As a consequence, one should in fact take the $\frac{1}{p}$ -power of the objective functional of (D_p) and then study the limiting problem as $p \rightarrow \infty$; unfortunately, our research in this direction has been unfruitful up to now. Finally, since the functional \mathcal{C}_p is not convex in λ , the convex duality theory does not apply directly to (P_p) . But one may wonder whether it is possible to overcome this difficulty and associate with (P_p) a dual problem with a structure similar to that of (D_p) : this is a quite involved question known as Dudley's problem (see, e.g., equation (1.1.10) and Remark 2.6.2 in [26]), and it is out of the scope of the present study.

Despite the above difficulties, we believe that developing a duality theory for the problem (P_∞) is an important issue since it would yield a deeper understanding of the problem of the existence and uniqueness for a particular optimal transport map. Further explorations in this direction could follow the methods of [7, 16, 18].

Acknowledgments. The authors would like to thank Guy Bouchitté for proposing the problem and for several fruitful discussions. The authors also thank Jouni Parkkonen for bringing [17] to their attention.

REFERENCES

- [1] L. AMBROSIO, *Lecture notes on optimal transportation problems*, in Mathematical Aspects of Evolving Interfaces (Funchal, 2000), Lecture Notes in Math. 1812, Springer, Berlin, 2003, pp. 1–52.

- [2] L. AMBROSIO AND A. PRATELLI, *Existence and stability results in the L^1 theory of optimal transportation*, in *Optimal Transportation and Applications* (Martina Franca, 2001), Lecture Notes in Math. 1813, Springer, Berlin, 2003, pp. 123–160.
- [3] G. ARONSSON, *Minimization problems for the functional $\sup_x F(x, f(x), f'(x))$* , Ark. Mat., 6 (1965), pp. 33–53.
- [4] G. ARONSSON, *Extension of functions satisfying Lipschitz conditions*, Ark. Mat., 6 (1967), pp. 551–561.
- [5] G. ARONSSON, M. G. CRANDALL, AND P. JUUTINEN, *A tour of the theory of absolutely minimizing functions*, Bull. Amer. Math. Soc. (N.S.), 41 (2004), pp. 439–505.
- [6] E. N. BARRON, R. R. JENSEN, AND C. Y. WANG, *The Euler equation and absolute minimizers of L^∞ functionals*, Arch. Ration. Mech. Anal., 157 (2001), pp. 255–283.
- [7] J. D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the Monge–Kantorovich mass transfer problem*, Numer. Math., 84 (2000), pp. 375–393.
- [8] M. BERNOT, V. CASELLES, AND J. M. MOREL, *Are there infinite irrigation trees?*, J. Math. Fluid Mech., 8 (2006), pp. 311–332.
- [9] G. BOUCHITTÉ, C. JIMENEZ, AND M. RAJESH, *A new L^∞ estimate in optimal mass transport*, Proc. Amer. Math. Soc., 135 (2007), pp. 3525–3535.
- [10] A. BRANCOLINI, G. BUTTAZZO, AND F. SANTAMBROGIO, *Path functionals over Wasserstein spaces*, J. Eur. Math. Soc. (JEMS), 8 (2006), pp. 415–434.
- [11] G. BUTTAZZO, A. PRATELLI, AND E. STEPANOV, *Optimal pricing policies for public transportation networks*, SIAM J. Optim., 16 (2006), pp. 826–853.
- [12] J. A. CARRILLO, M. DI FRANCESCO, AND C. LATTANZIO, *Contractivity and asymptotics in Wasserstein metrics for viscous nonlinear scalar conservation laws*, Boll. Unione Mat. Ital. Sez. B Artic. Ric. Mat. (8), 10 (2007), pp. 277–292.
- [13] J. A. CARRILLO, M. P. GUALDANI, AND G. TOSCANI, *Finite speed of propagation in porous media by mass transportation methods*, C. R. Math. Acad. Sci. Paris, 338 (2004), pp. 815–818.
- [14] T. CHAMPION AND L. DE PASCALE, *A principle of comparison with distance functions for absolute minimizers*, J. Convex Anal., 14 (2007), pp. 515–541.
- [15] T. CHAMPION AND L. DE PASCALE, in preparation.
- [16] L. DE PASCALE AND A. PRATELLI, *Regularity properties for Monge transport density and for solutions of some shape optimization problem*, Calc. Var. Partial Differential Equations, 14 (2002), pp. 249–274.
- [17] L. L. DINES, *On positive solutions of a system of linear equations*, Ann. of Math. (2), 28 (1926/27), pp. 386–392.
- [18] L. C. EVANS, *Partial differential equations and Monge–Kantorovich mass transfer*, in *Current Developments in Mathematics, 1997*, International Press, Boston, MA, 1999, pp. 65–126.
- [19] L. C. EVANS AND W. GANGBO, *Differential equations methods for the Monge–Kantorovich mass transfer problem*, Mem. Amer. Math. Soc., 137 (1999).
- [20] W. GANGBO AND R. J. MCCANN, *The geometry of optimal transportation*, Acta Math., 177 (1996), pp. 113–161.
- [21] M. HAURAY AND P.-E. JABIN, *N -particles approximation of the Vlasov equations with singular potential*, Arch. Ration. Mech. Anal., 183 (2007), pp. 489–524.
- [22] M. HAURAY, *private communication*.
- [23] F. MADDALENA, J. M. MOREL, AND S. SOLIMINI, *A variational model of irrigation patterns*, Interfaces Free Bound., 5 (2003), pp. 391–415.
- [24] R. J. MCCANN, *Stable rotating binary stars and fluid in a tube*, Houston J. Math., 32 (2006), pp. 603–631.
- [25] A. PRATELLI, *On the sufficiency of c -cyclical monotonicity for optimality of transport plans*, Math. Z., 258 (2008), pp. 677–690.
- [26] S. RACHEV AND L. RÜSCHENDORF, *Mass Transportation Problems. Vol. I. Theory*, Springer-Verlag, New York, 1998.
- [27] L. RÜSCHENDORF, *Optimal solutions of multivariate coupling problems*, Appl. Math. (Warsaw), 23 (1995), pp. 325–338.
- [28] L. RÜSCHENDORF, *On c -optimal random variables*, Statist. Probab. Lett., 27 (1996), pp. 267–270.
- [29] F. SANTAMBROGIO, *Variational Problems in Transport Theory with Mass Concentration*, Ph.D. thesis, Scuola Normale Superiore, Pisa, Italy, 2007 (available at <http://cvgmt.sns.it>).
- [30] O. SAVIN, *C^1 regularity for infinity harmonic functions in two dimensions*, Arch. Ration. Mech. Anal., 176 (2005), pp. 351–361.
- [31] W. SCHACHERMAYER AND J. TEICHMANN, *Characterization of optimal transport plans for the Monge–Kantorovich-problem*, Proc. Amer. Math. Soc., to appear.

- [32] C. VILLANI, *Topics in Optimal Transportation*, Grad. Stud. Math. 58, AMS, Providence, RI, 2003.
- [33] C. VILLANI, *Optimal Transport, Old and New*, <http://www.umpa.ens-lyon.fr/~cvillani/surveys.html#oldnew>.
- [34] Q. XIA, *Optimal paths related to transport problems*, Commun. Contemp. Math., 5 (2003), pp. 251–279.

REGULARITY UP TO THE BOUNDARY FOR NONLINEAR ELLIPTIC SYSTEMS ARISING IN TIME-INCREMENTAL INFINITESIMAL ELASTO-PLASTICITY*

PATRIZIO NEFF[†] AND DOROTHEE KNEES[‡]

Abstract. In this paper we investigate the question of higher regularity up to the boundary for quasi-linear elliptic systems which originate from the time discretization of models from infinitesimal elasto-plasticity. Our main focus lies on an elasto-plastic Cosserat model. More specifically we show that the time discretization renders H^2 -regularity of the displacement and H^1 -regularity for the symmetric plastic strain ε_p up to the boundary, provided that the plastic strain of the previous time step is in H^1 as well. This result contrasts with classical Hencky and Prandtl–Reuss formulations where it is known not to hold due to the occurrence of slip lines and shear bands. Similar regularity statements are obtained for other regularizations of ideal plasticity such as viscosity or isotropic hardening. In the first part we recall the time continuous Cosserat elasto-plasticity problem, provide the update functional for one time step, and show various preliminary results for the update functional (Legendre–Hadamard/monotonicity). Using nonstandard difference quotient techniques we are able to show the higher global regularity. Higher regularity is crucial for qualitative statements of finite element convergence. As a result we may obtain estimates linear in the mesh-width h in error estimates.

Key words. polar materials, perfect plasticity, higher global regularity, quasi-linear elliptic systems, error estimates, time increments

AMS subject classifications. 35B65, 49N60, 74A35, 74C05, 74G40

DOI. 10.1137/070695824

1. Introduction.

1.1. Plasticity and Cosserat models. This article addresses the regularity question for time-incremental formulations of *geometrically linear* elasto-plasticity. As a representative model problem we consider generalized continua of *Cosserat-micropolar* type.

The basic difference of a Cosserat model as compared with classical continuum models is the appearance of a nonsymmetric stress tensor which is augmented by a generalized balance of angular momentum equation allowing one to model interaction of particles not only by surface forces (classical Cauchy continuum) but also through surface couples (Cosserat continuum). General continuum models involving *independent rotations* as additional degrees of freedom were first introduced by the Cosserat brothers in [15]. For an introduction to the theory of Cosserat and micropolar models we refer the reader to the introduction in [49, 43, 45, 44, 48]; see also [22, 9].

There are a great many proposals for extensions of the elastic Cosserat framework to infinitesimal elasto-plasticity. We mention only [17, 19, 31, 55]. Recently the finite-strain formulation has been put into focus; see, e.g., [56, 62, 23] and the references therein.

*Received by the editors June 29, 2007; accepted for publication (in revised form) November 16, 2007; published electronically March 26, 2008.

<http://www.siam.org/journals/sima/40-1/69582.html>

[†]Fachbereich Mathematik, Technische Universität Darmstadt, Schlossgartenstrasse 7, 64289 Darmstadt, Germany (neff@mathematik.tu-darmstadt.de).

[‡]Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstr. 39, 10117 Berlin, Germany (knees@wias-berlin.de).

The first author has also proposed an elasto-plastic Cosserat model [45, 44] in a finite-strain framework. A geometrical linearization of this model has been investigated in [46, 48] and is shown to be well-posed also in the rate-independent limit for both quasi-static and dynamic processes.

When it comes to numerically solving problems in elasto-plasticity, then it is common practice to discretize the time evolution in the flow rule for the plastic variable with a backward Euler method and to consider a sequence of discrete-in-time problems [50]. Provided that the elasto-plastic model has certain variational features (hyperelasticity of the elastic response, associative flow rule) it is possible to recast the problem for one time step (called the update problem in the following) itself into a variational framework: the updated displacement is obtained as a minimizer of some update functional; see, e.g., [61, 60, 2, 66, 67]. This line of thought can be nicely extended to finite-strain multiplicative plasticity; see [52, 37, 36, 38] and the references therein. In the geometrically linear setting the resulting variational update problem usually has the form of a quasi-linear elliptic system whose corresponding energy has only linear growth (in case of perfect plasticity).

For qualitative statements on the rate of convergence of finite element methods it is necessary to know precisely the regularity of the function to be approximated. This then is the question on the regularity of the solution of the quasi-linear elliptic system constituting the update problem.

As far as classical rate-independent (perfect) elasto-plasticity is concerned we remark that global existence for the displacement has been shown only in a very weak, measure-valued sense, while the stresses could be shown to remain in $L^2(\Omega)$, provided that a safe load condition is assumed. For these results we refer the reader, for example, to [3, 13, 64]. If hardening or viscosity is added, then global H^1 -displacement solutions are found (see, e.g., [1, 12, 11]), already without safe load assumption. A complete theory for the classical rate-independent case remains, however, elusive; see also the remarks in [13].

Since classical perfect plasticity is, therefore, notoriously ill-posed (the updated displacements have derivatives only in a measure-valued sense) we focus in our investigation of higher regularity on certain modified update functionals which might allow for more regular updates. The Cosserat elasto-plastic model in [46] is our basic candidate. Based on this time-continuous model we investigate the time-incremental formulation and study the global regularity of minimizers of the corresponding update functional. In [49] this time-incremental formulation is the basis of a finite element approximation.

Our focus on Cosserat models is justified by the fact that the Cosserat-type models are today increasingly advocated as a means to regularize the pathological mesh size dependence of localization computations where shear failure mechanisms [14, 40, 4] play a dominant role; for applications in plasticity, see the nonexhaustive list [31, 19, 55, 17].

1.2. Outline of this contribution. Our contribution is organized as follows: first, we recall the time-continuous geometrically linear elasto-plastic Cosserat model as introduced in [45, 44] and investigated mathematically in [46, 48, 47].

Referring to the development in [49] we provide in section 2 the corresponding time-discretized formulation based on a fully implicit backward Euler discretization of the plastic flow rule in time. It is shown in [49] that at each time step t_n the updated displacement field u^n and the updated ‘‘Cosserat-microrotation-matrix’’ A^n can equivalently be obtained from a convex minimization problem which involves only

data from the previous time step. The plastic strain ε_p^n is then derived from A^n and u^n via a simple update formula. Furthermore, in [49] it has been shown that the update problem admits unique minimizers $u^n \in H^1(\Omega, \mathbb{R}^3)$, $A^n \in H^1(\Omega, \mathfrak{so}(3))$, and $\varepsilon_p^n \in L^2(\Omega, \text{Sym}(3))$, provided that the data coming from the previous time step are smooth enough. In order to quantify the rate of convergence of corresponding finite element methods for the update problem we investigate the regularity of the displacements u^n by studying the corresponding weak Euler–Lagrange equations. These equations form a quasi-linear elliptic system of partial differential equations. The main result of this paper is Theorem 5.2 in section 5, where we formulate a global regularity result for weak solutions of a rather general class of quasi-linear elliptic systems of second order. The time-incremental Cosserat plasticity formulation satisfies all the necessary assumptions of the regularity result, which allows us to show higher regularity to the extent that for all $n \in \mathbb{N}$: $u^n \in H^2(\Omega, \mathbb{R}^3)$, $A^n \in H^2(\Omega, \mathfrak{so}(3))$, and $\varepsilon_p^n \in H^1(\Omega, \text{Sym}(3))$ if pure Dirichlet data are assumed. Let us remark that it remains an open problem whether a similar regularity result is also valid for the time-continuous Cosserat model or other regularized time-continuous plasticity formulations.

The general quasi-linear elliptic systems, which we study in section 5, are of the following type: Find $u \in H_0^1(\Omega)$ such that for every $v \in H_0^1(\Omega)$

$$\int_{\Omega} \langle \mathcal{M}(x, \nabla u(x), z(x)), \nabla v(x) \rangle dx = \int_{\Omega} \langle f, v \rangle dx.$$

Here, $z \in L^2(\Omega, \mathbb{R}^N)$ and $f \in L^2(\Omega, \mathbb{R}^3)$ are the given data. For the Cosserat model, z is identified with (ε_p^n, A^n) ; the explicit structure of $\mathcal{M} = \mathcal{M}_C$ is given in section 2.4. It is shown that \mathcal{M}_C is rank-one monotone in ∇u and Lipschitz continuous but not differentiable. Consequently, we assume in the general case that the function $\mathcal{M} : \Omega \times \mathbb{M}^{m \times d} \times \mathbb{R}^N \rightarrow \mathbb{M}^{m \times d}$ is Lipschitz continuous, is rank-one monotone in ∇u , and induces a Gårding inequality. The precise conditions on \mathcal{M} are formulated as R1–R3 in section 5.1. Our main result is Theorem 5.2, where we prove for smooth domains that $u \in H^2(\Omega)$, provided that $z \in H^1(\Omega)$ and $f \in L^2(\Omega)$. We emphasize that we do not need the differentiability of \mathcal{M} and that we require \mathcal{M} to be rank-one monotone, only, instead of uniformly or strongly monotone. A further new aspect compared to systems studied in the literature is the presence of the function z in the definition of the differential operator.

Let us give a short overview on global regularity results for quasi-linear second order systems. Systems with quadratic growth or, more generally, with p growth are studied by several authors. We mention here the books [42, 39, 6] and the paper [53], where global regularity results for systems of the type

$$\text{Div } \mathcal{M}(x, \nabla u(x)) + f(x) = 0, \quad u|_{\partial\Omega} = g_D,$$

are shown for smooth domains assuming that \mathcal{M} is differentiable and strongly monotone. Further results for Lipschitz domains were obtained in [21, 20, 57], again assuming that \mathcal{M} is strongly monotone (or uniformly monotone if $p \neq 2$), that it is differentiable, and that there is a function W such that $\mathcal{M} = DW$. These results are proved with a difference quotient technique which relies on the standard finite differences $\delta_h u(x) := u(x+h) - u(x)$.

In [16] the authors study systems where $\mathcal{M}(x, u, \nabla u) = B(x)\nabla u + h(x, u, \nabla u)$. The main assumption in [16] is that B is uniformly positive definite, h is Hölder-continuous with respect to ∇u , and $h(x, u, \cdot)$ is uniformly monotone in zero. They prove that the gradient of solutions belongs locally to certain Campanato–Spanne

spaces. With our main result we can treat the case where h does not depend on u , where it is Lipschitz continuous and monotone but not necessarily uniformly monotone, and where B induces a rank-one positive quadratic form. We obtain $u \in H^2(\Omega)$ globally.

In [58] a nonlinear elliptic system is studied which is more related to our Cosserat model. There, \mathcal{M} is chosen as $\mathcal{M}(\nabla u) = \frac{h(|\varepsilon(u)|)}{|\varepsilon(u)|} \varepsilon(u)$, where $\varepsilon(u)$ is the linearized strain tensor, and it is assumed that h is differentiable except for a finite number of points and that h is strongly monotone. It is shown for smooth domains that $u \in H^2(\Omega)$ by investigating the regularity of functions u_δ with $\text{Div}(\delta\varepsilon(u_\delta) + \mathcal{M}(\varepsilon(u_\delta))) + f = 0$ for $\delta \searrow 0$. The results for u_δ are obtained with standard finite differences. Further results for related models were obtained in [54, 7]. Let us remark that the quasi-linear system we are interested in contains the above described systems as special cases (if $p = 2$) and that our main result is not covered by the above references. The local and global regularity of the stress fields of a class of degenerated quasi-linear elliptic systems is investigated in the papers [10, 33].

Note that higher regularity is not known to hold for the displacements of the classical limit of our formulation, which is the classical time-incremental Prandtl–Reuss model. In the first update step this model in turn is nothing else than the total deformation Hencky plasticity model. The Hencky model does not allow for regular displacements. Here, it is known that the displacement $u \in L^{\frac{3}{2}}(\Omega, \mathbb{R}^3)$ (see, e.g., [6, p. 423]), while the classical symmetric stresses satisfy $\sigma \in H_{\text{loc}}^1(\Omega, \text{Sym}(3)) \cap H^{\frac{1}{2}-\delta}(\Omega)$ for every $\delta > 0$ if the data are sufficiently regular and if Ω is a Lipschitz domain. See [59, 24, 5, 51, 18] for the local and [32] for the global result.

The proof of our own regularity result is split into the three classical steps. In the first step we investigate the tangential regularity of weak solutions in the case where Ω is a cube. Since we assumed rank-one monotonicity, only, we cannot apply the standard difference quotient technique in this step. Instead, we use finite differences which are based on inner variations: $\Delta_h u(x) = u(\tau_h(x)) - u(x)$, where $\tau_h(x) = x + \varphi^2(x)h$ for $h \in \mathbb{R}^d$ and a cut-off function φ . This will be explained in more detail in Remark 5.5. Let us note that these nonstandard differences were recently applied by Nesenenko [51] in order to obtain higher local regularity for models from elasto-plasticity with linear hardening. In the second step we prove higher regularity in directions normal to the boundary. Due to the lack of differentiability of \mathcal{M} we cannot apply the usual arguments (i.e., solving the equation for the normal derivatives) to obtain the differentiability of ∇u in the normal direction. Instead, we exploit the rank-one monotonicity of \mathcal{M} in order to get more information on the missing derivative. In the final step we prove the result for arbitrary bounded $\mathcal{C}^{1,1}$ -smooth domains by the usual localization procedure. The notation is found in the appendix.

2. The infinitesimal elasto-plastic Cosserat model. In this section we recall the specific isotropic infinitesimal elasto-plastic Cosserat model which has been proposed in a finite-strain setting in [44] and which was analyzed in [46]. Moreover, we derive a discrete formulation. This section does not contain new results; it serves to provide the clear definition of the problem and to introduce some of the notation.

2.1. Time-continuous infinitesimal elasto-plastic Cosserat model. The geometrically linear time continuous system in variational form with nondissipative Cosserat effects reads as follows: for given body forces $f(t) \in L^2(\Omega, \mathbb{R}^3)$ and given Dirichlet data find the *displacement* $u(t) \in H^1(\Omega, \mathbb{R}^3)$, the *skew-symmetric microrotation* $A(t) \in H^1(\Omega, \mathfrak{so}(3))$, and the *symmetric plastic strain* $\varepsilon_p(t) \in L^2(\Omega, \text{Sym}(3))$

with

$$\begin{aligned}
& \int_{\Omega} W(\nabla u, A, \varepsilon_p(t)) - \langle f(t), u \rangle \, dx \mapsto \min \quad \text{w.r.t. } (u, A) \text{ at fixed } \varepsilon_p(t), \\
& W(\nabla u, A, \varepsilon_p) = \mu \|\text{sym } \nabla u - \varepsilon_p\|^2 \\
& \quad + \mu_c \|\text{skew}(\nabla u - A)\|^2 + \frac{\lambda}{2} \text{tr} [\nabla u]^2 + 2\mu L_c^2 \|\nabla \text{axl}(A)\|^2, \\
& \varepsilon_p(t) \in \partial\chi(T_E(t)), \quad T_E = 2\mu(\varepsilon - \varepsilon_p), \quad \varepsilon_p \in \text{Sym}(3) \cap \mathfrak{sl}(3), \quad \varepsilon_p(0) = \varepsilon_p^0, \\
(2.1) \quad & u|_{\Gamma_D} = g_D(t, x) - x, \quad A|_{\Gamma_D} = \text{skew}(\nabla g_D(t, x))|_{\Gamma_D}.
\end{aligned}$$

Here, $\Omega \subset \mathbb{R}^3$ is a bounded smooth domain and $\Gamma_D \subset \partial\Omega$ is that part of the boundary where Dirichlet data are prescribed. The parameters $\mu, \lambda > 0$ are the Lamé constants of isotropic linear elasticity, $\mu_c > 0$ is the Cosserat couple modulus, and $L_c > 0$ is an internal length parameter.¹ The classical symmetric elastic strain $\text{sym } \nabla u$ is denoted by ε . The linear operator $\text{axl} : \mathfrak{so}(3) \rightarrow \mathbb{R}^3$ provides the canonical identification between the Lie algebra $\mathfrak{so}(3)$ of skew-symmetric matrices and vectors in \mathbb{R}^3 . The Lie algebra of trace-free matrices is denoted by $\mathfrak{sl}(3)$, and $\text{dev} : \mathbb{M}^{3 \times 3} \rightarrow \mathfrak{sl}(3)$, $\text{dev } X = X - \frac{1}{3}\mathbb{I}$ is the orthogonal projection onto $\mathfrak{sl}(3)$. As regards the plastic flow rule, $\partial\chi$ is the subdifferential of a convex flow potential $\chi : \mathbb{M}^{3 \times 3} \rightarrow \mathbb{R}^+$ acting on the generalized conjugate forces, i.e., the Eshelby stress tensor $T_E = -\partial_{\varepsilon_p} W(\nabla u, A, \varepsilon_p)$, where W is the free energy used in (2.1).²

The corresponding system of partial differential equations coupled with the flow rule is given by (note that $\|A\|_{\mathbb{M}^{3 \times 3}}^2 = 2\|\text{axl}(A)\|_{\mathbb{R}^3}^2$ for $A \in \mathfrak{so}(3, \mathbb{R})$)

$$\begin{aligned}
& \text{Div } \sigma = -f, \quad x \in \Omega, \quad \text{balance of forces,} \\
& \quad \sigma = 2\mu(\varepsilon - \varepsilon_p) + 2\mu_c(\text{skew}(\nabla u) - A) + \lambda \text{tr} [\varepsilon] \cdot \mathbb{I}, \\
& -\mu L_c^2 \Delta \text{axl}(A) = \mu_c \text{axl}(\text{skew}(\nabla u) - A), \quad \text{balance of angular momentum,} \\
& \quad \dot{\varepsilon}_p(t) \in \partial\chi(T_E), \quad T_E = 2\mu(\varepsilon - \varepsilon_p), \\
& \quad u|_{\Gamma_D}(t, x) = g_D(t, x) - x, \quad A|_{\Gamma_D} = \text{skew}(\nabla g_D(t, x))|_{\Gamma_D}, \\
& \quad \sigma \cdot \vec{n}|_{\partial\Omega \setminus \Gamma_D}(t, x) = 0, \quad \mu L_c^2 \nabla \text{axl}(A) \cdot \vec{n}|_{\partial\Omega \setminus \Gamma_D}(t, x) = 0, \\
& \quad \varepsilon_p(0) \in \text{Sym}(3) \cap \mathfrak{sl}(3).
\end{aligned}$$

Note that in this model the force stresses σ need not be symmetric and that the Cosserat effects, active through the microrotations A , appear only in the balance equations but not in the plastic flow rule since T_E does not depend on A . It is worth noting that this model is intrinsically thermodynamically correct. If $\Gamma_D = \partial\Omega$, then the model admits global weak solutions with the regularity [46]:

$$\begin{aligned}
& u \in L^\infty([0, T], H^1(\Omega, \mathbb{R}^3)), \quad A \in L^\infty([0, T], H^1(\Omega, \mathfrak{so}(3))), \\
& \varepsilon_p \in L^\infty([0, T], L^2(\Omega, \text{Sym}(3) \cap \mathfrak{sl}(3))).
\end{aligned}$$

2.2. Backward Euler time discretization of the flow rule. For a numerical treatment we consider the time discretization of the flow rule with the fully implicit backward Euler scheme. Let $0 = t_0 < t_1 < \dots < t_N = T$ be a subdivision of the

¹Observe that for $\mu_c = 0$ or $L_c = 0$ one recovers the classical Prandtl–Reuss formulation for the displacement u .

²The specification $\chi = I_K$ as indicator function of some elastic domain is not necessary at this point.

time interval $[0, T]$ with $t_j - t_{j-1} = \Delta t$. Let $f^n(x) = f(x, t_n)$, and assume that at time t_{n-1} a sufficiently regular plastic strain field $\varepsilon_p^{n-1} \in \text{Sym}(3) \cap \mathfrak{sl}(3)$ is given. We want to determine the *updated displacement* $u^n \in H^1(\Omega, \mathbb{R}^3)$, the *updated skew-symmetric microrotation* $A^n \in H^1(\Omega, \mathfrak{so}(3))$ and the *updated symmetric plastic strain* $\varepsilon_p^n \in L^2(\Omega, \text{Sym}(3) \cap \mathfrak{sl}(3))$ satisfying

$$\begin{aligned}
& \text{Div } \sigma^n = -f^n, \quad x \in \Omega, \\
& \sigma^n = 2\mu(\varepsilon^n - \varepsilon_p^n) + 2\mu_c(\text{skew}(\nabla u^n) - A^n) + \lambda \text{tr}[\varepsilon^n] \cdot \mathbb{1}, \\
(2.2) \quad & -\mu L_c^2 \Delta \text{axl}(A^n) = \mu_c \text{axl}(\text{skew}(\nabla u^n) - A^n), \\
& \frac{\varepsilon_p^n - \varepsilon_p^{n-1}}{\Delta t} \in \partial\chi(T_E^n), \quad T_E^n = 2\mu(\varepsilon^n - \varepsilon_p^n), \\
& u_{|\Gamma_D}^n(x) = g_D^n(x) - x, \quad A_{|\Gamma_D}^n = \text{skew}(\nabla g_D^n(x)), \\
& \sigma^n \cdot \vec{n}|_{\partial\Omega \setminus \Gamma_D}(x) = 0, \quad \mu L_c^2 \nabla \text{axl}(A^n) \cdot \vec{n}|_{\partial\Omega \setminus \Gamma_D}(x) = 0, \\
& \varepsilon_p^{n-1} \in L^2(\Omega, \text{Sym}(3) \cap \mathfrak{sl}(3)).
\end{aligned}$$

It is possible to explicitly solve the discretized flow rule (2.2)₄ for ε_p^n in terms of ε_p^{n-1} , ε^n , and Δt . To see this, consider

$$\begin{aligned}
(2.3) \quad & \frac{\varepsilon_p^n - \varepsilon_p^{n-1}}{\Delta t} \in \partial\chi(2\mu(\varepsilon^n - \varepsilon_p^n)) \Leftrightarrow 0 \in \partial\chi(2\mu(\varepsilon^n - \varepsilon_p^n)) - \frac{\varepsilon_p^n - \varepsilon_p^{n-1}}{\Delta t} \\
& \Leftrightarrow 0 \in \partial_{\varepsilon_p^n} \left(\mu \|\varepsilon_p^n - \varepsilon_p^{n-1}\|^2 + \Delta t \chi(2\mu(\varepsilon^n - \varepsilon_p^n)) \right).
\end{aligned}$$

Thus we can define the local potential for the local flow rule

$$(2.4) \quad V^{\text{time}}(\varepsilon^n, \varepsilon_p^n, \varepsilon_p^{n-1}, \Delta t) := \mu \|\varepsilon_p^n - \varepsilon_p^{n-1}\|^2 + \Delta t \chi(2\mu(\varepsilon^n - \varepsilon_p^n)).$$

It is easy to see that V^{time} is strictly convex in ε_p^n ; thus V^{time} admits a unique minimizer satisfying (2.3)₃. Moreover, we have

$$\begin{aligned}
V^{\text{time}}(\varepsilon^n, \varepsilon_p^n, \varepsilon_p^{n-1}, \Delta t) &= \mu \|\varepsilon_p^n - \varepsilon_p^{n-1}\|^2 + \Delta t \chi(2\mu(\varepsilon^n - \varepsilon_p^n)) \\
&= \frac{1}{4\mu} \|2\mu(\varepsilon_p^n - \varepsilon^n + \varepsilon^n - \varepsilon_p^{n-1})\|^2 + \Delta t \chi(2\mu(\varepsilon^n - \varepsilon_p^n)) \\
&= \frac{1}{4\mu} \|\Sigma^n - \Sigma_{\text{trial}}^n\|^2 + \Delta t \chi(\Sigma^n) = \tilde{V}(\Sigma^n, \Sigma_{\text{trial}}^n),
\end{aligned}$$

where $\Sigma^n = 2\mu(\varepsilon^n - \varepsilon_p^n)$ and the so-called trial stresses $\Sigma_{\text{trial}}^n = 2\mu(\varepsilon^n - \varepsilon_p^{n-1})$. Minimizing V^{time} with respect to ε_p^n is equivalent to minimizing \tilde{V} with respect to Σ^n . Proceeding further, we specialize χ . Let us define the elastic domain in stress space

$$K := \{\Sigma \in \mathbb{M}^{3 \times 3} \mid \|\text{dev } \Sigma\| \leq \sigma_y\},$$

with initial yield stress σ_y , $[\sigma_y] = [\text{MPa}]$, and corresponding indicator function

$$I_K(\Sigma) = \begin{cases} 0, & \|\text{dev } \Sigma\| \leq \sigma_y, \\ \infty, & \|\text{dev } \Sigma\| > \sigma_y, \end{cases}$$

and let $\chi = I_K$. We have therefore $\partial\chi = \partial I_K$ in the sense of the subdifferential. With this choice, the unique minimizer of \tilde{V} is simply characterized by

$$\inf_{\Sigma^n \in K} \|\Sigma^n - \Sigma_{\text{trial}}^n\|^2,$$

independent of Δt . The solution is the orthogonal projection of Σ_{trial}^n onto the convex set K , denoted by

$$\Sigma^n = P_K(\Sigma_{\text{trial}}^n) \Rightarrow 2\mu(\varepsilon^n - \varepsilon_p^n) = P_K(2\mu(\varepsilon^n - \varepsilon_p^{n-1})).$$

Reintroducing the last result into the balance of forces equation (2.2)₁ delivers

$$(2.5) \quad \begin{aligned} \text{Div } \sigma^n &= -f^n, \quad x \in \Omega, \\ \sigma^n &= P_K(2\mu(\varepsilon^n - \varepsilon_p^{n-1})) + 2\mu_c(\text{skew}(\nabla u^n) - A^n) + \lambda \text{tr}[\varepsilon^n] \cdot \mathbb{1}. \end{aligned}$$

This step is called *return mapping* [61, 60] in an engineering context of classical plasticity. At the given plastic strain of the previous time step ε_p^{n-1} this equation is the strong form of the update problem for the force-balance equation.

Gathering the previous development the formal problem for the update consists of determining $u^n \in H^1(\Omega, \mathbb{R}^3)$, $A^n \in H^1(\Omega, \mathfrak{so}(3))$, and $\varepsilon_p^n \in L^2(\Omega, \text{Sym}(3) \cap \mathfrak{sl}(3))$ satisfying

$$(2.6) \quad \begin{aligned} \text{Div } \sigma^n &= -f^n, \quad x \in \Omega, \\ \sigma^n &= P_K(2\mu(\varepsilon^n - \varepsilon_p^{n-1})) + 2\mu_c(\text{skew}(\nabla u^n) - A^n) + \lambda \text{tr}[\varepsilon^n] \cdot \mathbb{1}, \\ -\mu L_c^2 \Delta \text{axl}(A^n) &= \mu_c \text{axl}(\text{skew}(\nabla u^n) - A^n). \end{aligned}$$

The updated plastic strain field is then given by

$$(2.7) \quad \varepsilon_p^n = \varepsilon^n - \frac{1}{2\mu} P_K(2\mu(\varepsilon^n - \varepsilon_p^{n-1})).$$

For the precise formulation of this system we need the projection operator onto the yield surface which we recall in the following.

2.3. The projection onto the yield surface. Let K be a convex domain in stress space defined as

$$(2.8) \quad K := \{ \Sigma \in \mathbb{M}^{3 \times 3} \mid \|\text{dev } \Sigma\| \leq \sigma_y \}.$$

The orthogonal projection $P_K : \mathbb{M}^{3 \times 3} \rightarrow K$ onto this set is uniquely given by (see, e.g., [29, 30])

$$\begin{aligned} P_K(\Sigma) &= \begin{cases} \Sigma, & \Sigma \in K, \\ \Sigma - (\|\text{dev } \Sigma\| - \sigma_y) \frac{\text{dev } \Sigma}{\|\text{dev } \Sigma\|}, & \Sigma \notin K \end{cases} \\ &= \begin{cases} \Sigma, & \|\text{dev } \Sigma\| \leq \sigma_y, \\ \frac{1}{3} \text{tr}[\Sigma] \mathbb{1} + \frac{\sigma_y}{\|\text{dev } \Sigma\|} \text{dev } \Sigma, & \|\text{dev } \Sigma\| > \sigma_y. \end{cases} \end{aligned}$$

It is easy to see that P_K is Lipschitz continuous but not differentiable at Σ with $\|\text{dev } \Sigma\| = \sigma_y$.³ From convex analysis it is clear that P_K represents a monotone operator which is nonexpansive. Therefore, P_K has Lipschitz constant 1. Observe also that

$$(2.10) \quad P_K(\Sigma) = \frac{1}{3} \text{tr}[\Sigma] \mathbb{1} + P_K(\text{dev } \Sigma).$$

³Consider the simple example $p : \mathbb{R} \rightarrow \mathbb{R}$,

$$(2.9) \quad p(x) = \begin{cases} x, & |x| \leq \sigma_y, \\ \sigma_y \frac{x}{|x|}, & |x| > \sigma_y. \end{cases}$$

For future reference we calculate also

$$\begin{aligned}
(2.11) \quad \Sigma - P_K(\Sigma) &= \begin{cases} 0, & \|\operatorname{dev} \Sigma\| \leq \sigma_y, \\ \operatorname{dev} \Sigma \left(1 - \frac{\sigma_y}{\|\operatorname{dev} \Sigma\|}\right), & \|\operatorname{dev} \Sigma\| > \sigma_y \end{cases} \\
&= [\|\operatorname{dev} \Sigma\| - \sigma_y]_+ \frac{\operatorname{dev} \Sigma}{\|\operatorname{dev} \Sigma\|}, \\
\|\Sigma - P_K(\Sigma)\|^2 &= [\|\operatorname{dev} \Sigma\| - \sigma_y]_+^2,
\end{aligned}$$

where $[x]_+ := \max\{0, x\}$.

2.4. Weak form of the reduced update problem. From now onwards we take $\Gamma_D = \partial\Omega$ and assume $g_D = x$; i.e., the body is fixed everywhere on its boundary and subject only to body forces. This assumption allows us to confine attention to the simpler setting in $H_0^1(\Omega)$. We introduce the nonlinear mapping

$$\begin{aligned}
(2.12) \quad \mathcal{M}_C : \mathbb{M}^{3 \times 3} \times \operatorname{Sym}(3) \times \mathfrak{so}(3) &\rightarrow \mathbb{M}^{3 \times 3}, \\
\mathcal{M}_C(X, \varepsilon_p, A) &:= P_K(2\mu(\operatorname{sym} X - \varepsilon_p)) + \lambda \operatorname{tr}[X] \mathbb{1} + 2\mu_c(\operatorname{skew}(X) - A).
\end{aligned}$$

The weak form of the update problem (2.6) now reads as follows: for given $f^n \in L^2(\Omega, \mathbb{R}^3)$ and $\varepsilon_p^{n-1} \in L^2(\Omega, \operatorname{Sym}(3) \cap \mathfrak{sl}(3))$ find $(u^n, A^n) \in H_0^1(\Omega, \mathbb{R}^3) \times H_0^1(\Omega, \mathfrak{so}(3))$ satisfying for all $v \in H_0^1(\Omega, \mathbb{R}^3)$ and all $B \in H_0^1(\Omega, \mathfrak{so}(3))$

$$(2.13) \quad \int_{\Omega} \langle \mathcal{M}_C(\nabla u^n, \varepsilon_p^n, A^n), \nabla v \rangle \, dx = \int_{\Omega} \langle f^n, v \rangle \, dx,$$

$$(2.14) \quad \mu L_c^2 \int_{\Omega} \langle DA^n, DB \rangle \, dx = \mu_c \int_{\Omega} \langle \operatorname{skew} \nabla u^n - A^n, B \rangle \, dx.$$

The updated plastic strain field ε_p^n is then obtained by (2.7). It is shown in [49] that for every n the system (2.13)–(2.14) admits a unique weak solution $u^n \in H_0^1(\Omega, \mathbb{R}^3)$ and $A^n \in H_0^1(\Omega, \mathfrak{so}(3))$. Equation (2.13) represents the quasi-linear elliptic system for determining u^n , which will be discussed with respect to regularity. Together with $\varepsilon_p^{n-1}, \varepsilon^n \in H^1(\Omega, \operatorname{Sym}(3))$, which we will obtain from the regularity result to be proven below, using (2.7) we see that $\varepsilon_p^n \in H^1(\Omega, \operatorname{Sym}(3))$.

LEMMA 2.1 (strong Legendre–Hadamard ellipticity). *Let $\mu > 0$, $2\mu + 3\lambda > 0$, and $0 < \mu_c$. Then the matrix-valued function \mathcal{M}_C is strongly rank-one monotone; i.e., there exists a constant $c_{LH}^+ > 0$ such that for every $X \in \mathbb{M}^{3 \times 3}$, $\varepsilon_p \in \operatorname{Sym}(3)$, $A \in \mathfrak{so}(3)$ and for all $\xi, \eta \in \mathbb{R}^3$ we have*

$$(2.15) \quad \langle \mathcal{M}_C(X + \xi \otimes \eta, \varepsilon_p, A) - \mathcal{M}_C(X, \varepsilon_p, A), \xi \otimes \eta \rangle \geq c_{LH}^+ \|\xi\|^2 \|\eta\|^2.$$

Proof. The projection P_K itself is monotone, and for $\mu > 0$ there is no sign change. Thus the map $X \rightarrow P_K(2\mu(\operatorname{sym} X - \varepsilon_p))$ is also monotone in X . Since (2.10) holds we have

$$\langle P_K(2\mu(\operatorname{sym} X + \xi \otimes \eta - \varepsilon_p)) - P_K(2\mu(\operatorname{sym} X - \varepsilon_p)), \xi \otimes \eta \rangle \geq \frac{2\mu}{3} \operatorname{tr}[\xi \otimes \eta]^2.$$

For the remaining linear contribution we have

$$\begin{aligned}
&\langle \lambda \operatorname{tr}[X + \xi \otimes \eta] \mathbb{1} + 2\mu_c \operatorname{skew}(X + \xi \otimes \eta - A) - [\lambda \operatorname{tr}[X] \mathbb{1} + 2\mu_c \operatorname{skew}(X - A)], \xi \otimes \eta \rangle \\
&= \lambda \operatorname{tr}[\xi \otimes \eta]^2 + 2\mu_c \|\operatorname{skew}(\xi \otimes \eta)\|^2.
\end{aligned}$$

Thus

(2.16)

$$\begin{aligned}
& \langle \mathcal{M}_C(X + \xi \otimes \eta, \varepsilon_p, A) - \mathcal{M}_C(X, \varepsilon_p, A), \xi \otimes \eta \rangle \\
& \geq \frac{2\mu+3\lambda}{3} \operatorname{tr} [\xi \otimes \eta]^2 + 2\mu_c \|\operatorname{skew}(\xi \otimes \eta)\|^2 = \frac{2\mu+3\lambda}{3} \langle \xi, \eta \rangle^2 + \mu_c (\|\xi\|^2 \|\eta\|^2 - \langle \xi, \eta \rangle^2) \\
& \text{split } \mu_c^1 + \mu_c^2 = \mu_c \\
& = \left(\frac{2\mu+3\lambda}{3} - \mu_c^1 \right) \langle \xi, \eta \rangle^2 + \mu_c^1 \|\xi\|^2 \|\eta\|^2 + \underbrace{\mu_c^2 (\|\xi\|^2 \|\eta\|^2 - \langle \xi, \eta \rangle^2)}_{\geq 0} \\
& \geq \left(\frac{2\mu+3\lambda}{3} - \mu_c^1 \right) \langle \xi, \eta \rangle^2 + \mu_c^1 \|\xi\|^2 \|\eta\|^2 \geq \mu_c^1 \|\xi\|^2 \|\eta\|^2
\end{aligned}$$

if $0 < \mu_c^1 < \frac{3\lambda+2\mu}{3}$. Thus \mathcal{M}_C generates a strongly Legendre–Hadamard elliptic operator with ellipticity constant $c_{LH}^+ = \min(\mu_c, \frac{2\mu+3\lambda}{3})$. \square

Obviously, \mathcal{M} is Lipschitz continuous: for every $X_i \in \mathbb{M}^{3 \times 3}$, $P_i \in \operatorname{Sym}(3)$, $A_i \in \mathfrak{so}(3)$ we have

$$\|\mathcal{M}_C(X_1, P_1, A_1) - \mathcal{M}_C(X_2, P_2, A_2)\| \leq L_{\mathcal{M}_C} (\|X_1 - X_2\| + \|P_1 - P_2\| + \|A_1 - A_2\|).$$

LEMMA 2.2. *Let $\mu > 0$, $2\mu + 3\lambda > 0$, and $\mu_c > 0$. The operator \mathcal{M}_C generates a strongly monotone operator on $H_0^1(\Omega, \mathbb{R}^3)$; that is, there exists a constant $c_{\mathcal{M}_C} > 0$ such that for every $v_1, v_2 \in H_0^1(\Omega, \mathbb{R}^3)$ and for all $\varepsilon_p \in L^2(\Omega, \operatorname{Sym}(3))$ and $A \in L^2(\Omega, \mathfrak{so}(3))$ we have*

$$(2.17) \quad \int_{\Omega} \langle \mathcal{M}_C(\nabla v_1, \varepsilon_p, A) - \mathcal{M}_C(\nabla v_2, \varepsilon_p, A), \nabla v_1 - \nabla v_2 \rangle \, dx \geq c_{\mathcal{M}_C} \|v_1 - v_2\|_{H_0^1(\Omega, \mathbb{R}^3)}^2.$$

Proof. The same calculation as in the proof of Lemma 2.1 yields the estimate

$$\begin{aligned}
& \langle \mathcal{M}_C(\nabla v_1, \varepsilon_p, A) - \mathcal{M}_C(\nabla v_2, \varepsilon_p, A), \nabla v_1 - \nabla v_2 \rangle \\
& \geq \frac{2\mu + 3\lambda}{3} \operatorname{tr} [\nabla v_1 - \nabla v_2]^2 + 2\mu_c \|\operatorname{skew}(\nabla v_1 - \nabla v_2)\|^2.
\end{aligned}$$

Set $u = v_1 - v_2$ and consider

$$(2.18) \quad \frac{2\mu + 3\lambda}{3} \operatorname{tr} [\nabla u]^2 + 2\mu_c \|\operatorname{skew} \nabla u\|^2 = \frac{2\mu + 3\lambda}{3} |\operatorname{Div} u|^2 + \mu_c \|\operatorname{curl} u\|^2.$$

The Div/Curl inequality on the space $H_0^1(\Omega)$ guarantees that there exists $C^+ > 0$ such that

$$(2.19) \quad \forall u \in H_0^1(\Omega, \mathbb{R}^3) : \int_{\Omega} |\operatorname{Div} u|^2 + \|\operatorname{curl} u\|^2 \, dx \geq C^+ \|u\|_{H_0^1(\Omega, \mathbb{R}^3)}^2;$$

see, for example, [28]. Applying this inequality to (2.18) implies finally (2.17). \square

It is instructive to realize that although the quadratic form (2.18) is formally positive in the sense of Nečas [41] and strongly Legendre–Hadamard elliptic with constant coefficients it is impossible to extend the analysis to Dirichlet boundary conditions given only on a part of the boundary $\partial\Omega$. We observe that

$$(2.20) \quad \left\| \sqrt{\mu_c} \operatorname{skew} X + \sqrt{\frac{\lambda}{2 \cdot 3}} \operatorname{tr} [X] \mathbb{1} \right\|^2 = \frac{\lambda}{2} \operatorname{tr} [X]^2 + \mu_c \|\operatorname{skew} X\|^2.$$

Let $\widehat{\mathcal{A}}$ be the constant-coefficient first order differential operator

$$\widehat{\mathcal{A}}.\nabla u = \sqrt{\mu_c} \operatorname{skew}(\nabla u) + \sqrt{\frac{\lambda}{2\cdot 3}} \operatorname{tr}[\nabla u] \mathbb{1}.$$

The corresponding Fourier symbol is given as a linear operator $\mathcal{A}(\xi) : \mathbb{C}^3 \rightarrow \mathbb{C}^{3 \times 3}$ with

$$(2.21) \quad \mathcal{A}(\xi).\hat{u} := \sqrt{\mu_c} \operatorname{skew}(\xi \otimes \hat{u}) + \sqrt{\frac{\lambda}{2\cdot 3}} \operatorname{tr}[\xi \otimes \hat{u}] \mathbb{1}.$$

From (2.20) it follows that

$$\|\mathcal{A}(\xi).\hat{u}\|^2 = \frac{\lambda}{2} \operatorname{tr}[\xi \otimes \hat{u}]^2 + \mu_c \|\operatorname{skew} \xi \otimes \hat{u}\|^2.$$

By algebraic completeness of the symbol $\mathcal{A}(\xi) : \mathbb{C}^3 \rightarrow \mathbb{C}^{3 \times 3}$ it is meant

$$\forall \xi \in \mathbb{C}^3, \xi \neq 0 : \quad \mathcal{A}(\xi).\hat{u} = 0_{\mathbb{C}^{3 \times 3}} \Rightarrow \hat{u} = 0_{\mathbb{C}^3}.$$

Recall that the corresponding statement for real ξ , i.e.,

$$\forall \xi \in \mathbb{R}^3, \xi \neq 0 : \quad \mathcal{A}(\xi).\hat{u} = 0_{\mathbb{R}^{3 \times 3}} \Rightarrow \hat{u} = 0_{\mathbb{R}^3},$$

is a consequence of strict Legendre–Hadamard ellipticity of $\widehat{\mathcal{A}}$. If the symbol is algebraically complete, then, using the result in Nečas [41] the induced quadratic form

$$\int_{\Omega} \|\widehat{\mathcal{A}}.\nabla u\|^2 + \|u\|^2 \, dx$$

is an equivalent norm on $H^1(\Omega, \mathbb{R}^3)$. However, we proceed to show that \mathcal{A} as defined in (2.21) corresponding to our quadratic form (2.18) is not algebraically complete.

Proof. To this end we write

$$\mathcal{A}(\xi).\hat{u} = 0 \Rightarrow \operatorname{tr}[\xi \otimes \hat{u}] = 0, \quad \text{and} \quad \operatorname{skew}(\xi \otimes \hat{u}) = 0 \Rightarrow \xi = \hat{u}, \quad \operatorname{tr}[\xi \otimes \xi] = 0.$$

Consider for simplicity the two-dimensional case:

$$\begin{aligned} \xi &= \begin{pmatrix} \alpha_1 + i\beta_1 \\ \alpha_2 + i\beta_2 \end{pmatrix}, \quad \xi \otimes \xi = \begin{pmatrix} \xi_1 \xi_1 & \xi_1 \xi_2 \\ \xi_2 \xi_1 & \xi_2 \xi_2 \end{pmatrix}, \\ \operatorname{tr}[\xi \otimes \xi] &= \xi_1 \xi_1 + \xi_2 \xi_2 = \alpha_1^2 + \alpha_2^2 - (\beta_1^2 + \beta_2^2) + 2i(\alpha_1 \beta_1 + \alpha_2 \beta_2) = 0. \end{aligned}$$

Choosing $\xi = (i, 1)^T$ shows that $\operatorname{tr}[\xi \otimes \xi] = 0$, which proves the claim. \square

Thence, the quadratic form is not algebraically complete, and this excludes the treatment of mixed boundary conditions on u in the following: we are forced to assume $\Gamma_D = \partial\Omega$. However, inhomogeneous Dirichlet conditions may be prescribed as far as the use of the Div/Curl estimate is concerned.

2.5. Variational form of the update problem. Due to the underlying variational formulation, the weak form (2.13) of the time-incremental Cosserat problem still has a variational structure. In [49] it is shown that solving (2.13)–(2.14) is equivalent to the following minimization problem: find $(u^n, A^n) \in H_0^1(\Omega, \mathbb{R}^3) \times H_0^1(\Omega, \mathfrak{so}(3))$ which minimize the functional

$$(2.22) \quad I_{\text{incr}}^n(u, A) = \mathcal{E}_{\text{incr}}(u, A, \varepsilon_p^{n-1}) - \int_{\Omega} \langle f^n, u \rangle \, dx$$

in $H_0^1(\Omega, \mathbb{R}^3) \times H_0^1(\Omega, \mathfrak{so}(3))$. Here, $\mathcal{E}_{\text{incr}}$ denotes the free energy of the incremental problem defined by

$$(2.23) \quad \begin{aligned} \mathcal{E}_{\text{incr}}(u, A, \varepsilon_p) &= \frac{1}{2\mu} \int_{\Omega} \Psi(2\mu(\text{sym}(\nabla u) - \varepsilon_p)) \, dx + \frac{\lambda}{2} \int_{\Omega} \text{tr} [\nabla u]^2 \, dx \\ &\quad + \mu_c \int_{\Omega} \|\text{skew}(\nabla u) - A\|^2 \, dx + \mu L_c^2 \int_{\Omega} \|DA\|^2 \, dx, \end{aligned}$$

with a potential function $\Psi : \mathbb{M}^{3 \times 3} \rightarrow \mathbb{R}^+$ of the form

$$(2.24) \quad \begin{aligned} \Psi(X) &:= \begin{cases} \frac{1}{2} \|X\|^2, & \|\text{dev } X\| \leq \sigma_y, \\ \frac{1}{2} \left(\frac{1}{3} \text{tr} [X]^2 + 2\sigma_y \|\text{dev } X\| - \sigma_y^2 \right), & \|\text{dev } X\| > \sigma_y \end{cases} \\ &= \frac{1}{2} \|X\|^2 - \frac{1}{2} [\|\text{dev } X\| - \sigma_y]_+^2. \end{aligned}$$

Clearly, Ψ is convex but not strongly convex outside the yield surface. Moreover, it has only linear growth outside the yield surface. Note that for the first time step $n = 1$ and $\varepsilon_p^0 = 0$, $\mu_c = 0$, $L_c = 0$ the functional $I_{\text{incr}}^1(u, 0)$ reduces to the primal plastic functional of static perfect plasticity (Hencky plasticity) [35, 63, 24, 25, 6].

Calculating the subdifferential of the convex potential shows that

$$(2.25) \quad \begin{aligned} \partial\Psi(\Sigma).H &= \begin{cases} \langle \Sigma, H \rangle, & \|\text{dev } \Sigma\| \leq \sigma_y, \\ \frac{1}{3} \text{tr} [\Sigma] \text{tr} [H] + \frac{\sigma_y}{\|\text{dev } \Sigma\|} \langle \text{dev } \Sigma, \text{dev } H \rangle, & \|\text{dev } \Sigma\| > \sigma_y \end{cases} \\ &= \langle P_K(\Sigma), H \rangle. \end{aligned}$$

Hence $\partial\Psi(\Sigma) = P_K(\Sigma)$, motivating the variational structure. The following relationship between the potential Ψ and the projection P_K is also valid:

$$\Psi(X) = \frac{1}{2} \|X\|^2 - \frac{1}{2} \|X - P_K(X)\|^2.$$

For future reference the second differential of the potential Ψ can be calculated in those points where the potential is differentiable. It holds that

$$(2.26) \quad D_X^2 \Psi(X).(H, H) = \begin{cases} \|H\|^2, & \|\text{dev } X\| < \sigma_y, \\ \text{does not exist}, & \|\text{dev } X\| = \sigma_y, \\ \frac{1}{3} \text{tr} [H]^2 + \sigma_y \left(\frac{\|\text{dev } H\|^2}{\|\text{dev } X\|} - \frac{\langle \text{dev } X, H \rangle^2}{\|\text{dev } X\|^3} \right), & \|\text{dev } X\| > \sigma_y. \end{cases}$$

The potential Ψ is not strictly rank-one convex in X , since taking $H = \xi \otimes \eta$ with $\langle \xi, \eta \rangle = 0$ yields

$$D_X^2 \Psi(X).(\xi \otimes \eta, \xi \otimes \eta) = \begin{cases} \|\xi\|^2 \|\eta\|^2, & \|\text{dev } X\| \leq \sigma_y, \\ \sigma_y \left(\frac{\|\text{dev } \xi \otimes \eta\|^2}{\|\text{dev } X\|} - \frac{\langle \text{dev } X, \xi \otimes \eta \rangle^2}{\|\text{dev } X\|^3} \right), & \|\text{dev } X\| > \sigma_y. \end{cases}$$

Taking $X = \xi \otimes \eta$ shows finally

$$D_X^2 \Psi(X).(\xi \otimes \eta, \xi \otimes \eta) = \begin{cases} \|\xi\|^2 \|\eta\|^2, & \|\text{dev } X\| \leq \sigma_y, \\ 0, & \|\text{dev } X\| > \sigma_y. \end{cases}$$

3. Improved error estimates for Cosserat plasticity. Let $h > 0$ be the mesh size of a finite element method, and let $V_h \subset H_0^1(\Omega, \mathbb{R}^3)$ be a corresponding discrete finite element space. Let us concentrate on the displacement approximation only. In [49, Thm. 8] the following error estimate for the discrete solution $u_h^{\mu_c, n} \in V_h$ of the Galerkin approximation of (2.23) in V_h has been shown:

$$(3.1) \quad \|u^{\mu_c, n} - u_h^{\mu_c, n}\|_{H_0^1(\Omega)} \leq \frac{C_1}{\mu_c} \inf_{v_h \in V_h} \|u^{\mu_c, n} - v_h\|_{H_0^1(\Omega)},$$

with a constant $C_1 > 0$. Here, $u^{\mu_c, n} = u^n$ is the exact solution of (2.13).

Using our regularity result from section 5, i.e., $u^{\mu_c, n} \in H^2(\Omega, \mathbb{R}^3)$, the right-hand side can be estimated qualitatively. If V_h is chosen to be the space of piecewise linear finite elements, then it holds [8, p. 107] that

$$(3.2) \quad \|u^{\mu_c, n} - u_h^{\mu_c, n}\|_{H_0^1(\Omega)} \leq \frac{C_2}{\mu_c} h \|u^{\mu_c, n}\|_{H^2(\Omega)}.$$

In [49] it has also been shown that for $\mu_c \rightarrow 0$ the classical Prandtl–Reuss symmetric Cauchy stresses σ^0 are approximated by the sequence of nonsymmetric stresses σ^{μ_c} whenever a safe load condition is satisfied. The estimate (3.2) strongly suggests therefore to balance h against μ_c to obtain optimal rates of convergence to the classical solution as in [54], where hardening-type approximations have been considered.

4. Higher regularity for alternative regularized update potentials. Our regularity result can also be applied to many other problems arising in the context of infinitesimal plasticity. There exist several other possibilities to regularize the classical update problem for the Prandtl–Reuss model. We recall the classical update problem: find a minimizer $u^n \in BD(\Omega, \mathbb{R}^3)$ of the functional

$$(4.1) \quad I_{\text{incr}}^{\text{class}}(u) = \mathcal{E}_{\text{incr}}^{\text{class}}(u, \varepsilon_p^{n-1}) - \int_{\Omega} \langle f^n, u \rangle \, dx,$$

where $\mathcal{E}_{\text{incr}}^{\text{class}}$ denotes the free energy of the classical incremental problem defined by

$$(4.2) \quad \mathcal{E}_{\text{incr}}^{\text{class}}(u, \varepsilon_p) = \frac{1}{2\mu} \int_{\Omega} \Psi(2\mu(\text{sym}(\nabla u) - \varepsilon_p)) \, dx + \int_{\Omega} \frac{\lambda}{2} \text{tr} [\nabla u]^2 \, dx,$$

with the potential Ψ as in (2.24). There is a vast literature on this Prandtl–Reuss update problem, mostly for the first time step $n = 1$, in which case it is the classical Hencky problem of total deformation plasticity [63, 54, 24, 25]. In this case, the plastic strain field ε_p is a symmetric bounded measure [63, 6]. The classical symmetric Cauchy stresses $\sigma = 2\mu(\text{sym} \nabla u - \varepsilon_p) + \lambda \text{tr} [\nabla u] \mathbb{1}$ satisfy $\sigma \in L^2(\Omega, \text{Sym}(3))$; indeed higher regularity for the stresses can be shown in the sense that $\sigma \in H_{\text{loc}}^1(\Omega, \text{Sym}(3)) \cap H^{\frac{1}{2}-\delta}(\Omega)$.

For regularization purposes the following proposals are usually made:

$$(4.3) \quad \mathcal{E}_{\text{incr}}^{\text{reg}}(u, \varepsilon_p) = \frac{1}{2\mu} \int_{\Omega} \Psi(2\mu(\text{sym}(\nabla u) - \varepsilon_p)) \, dx + \int_{\Omega} \frac{\lambda}{2} \text{tr} [\nabla u]^2 + \text{Reg}(\nabla u, \varepsilon_p) \, dx,$$

with the function Reg in the form

$$(4.4) \quad \text{Reg}(\nabla u, \varepsilon_p) = \frac{\mu \delta}{2} \|\text{dev sym } \nabla u - \varepsilon_p\|^2, \quad \text{Fuchs and Seregin [24, p. 60],}$$

$$\text{Reg}(\nabla u, \varepsilon_p) = \frac{1}{2\mu(1 + \frac{\Delta t}{\eta})} [\|\mu(\text{dev sym } \nabla u - \varepsilon_p) - \sigma_y\|_+^2], \quad \text{linear viscosity } \eta,$$

$$\text{Reg}(\nabla u, \varepsilon_p) = \frac{\mu\delta}{2} \|\nabla u - \varepsilon_p\|^2, \quad \text{locally strictly convex in } \nabla u.$$

In each case, for $\delta > 0$ the density of the update problem is then uniformly convex in the symmetric strain $\varepsilon = \text{sym } \nabla u$. Moreover,

$$(4.5) \quad \text{Reg}(\nabla u, \varepsilon_p) + \frac{\lambda}{2} \text{tr} [\nabla u]^2 \geq c^+ \|\varepsilon - \varepsilon_p\|^2,$$

Korn's first inequality establishes quadratic growth, and we have uniform convexity for the regularized problem. Our main regularity result applies therefore also to these models.

In the case with linear hardening it is simpler to write the update potential directly. We consider as an example isotropic hardening with the hardening variable $\alpha \geq 0$ (a measure for the accumulated plastic strain in the previous time step). Here, the energy $\mathcal{E}_{\text{incr}}$ can be expressed as (cf. [60, p. 124])

$$(4.6) \quad \mathcal{E}_{\text{incr}}^{\text{hard}}(u, \varepsilon_p, \alpha) = \frac{1}{2\mu} \int_{\Omega} \Psi_{\text{hard}}(2\mu(\text{sym}(\nabla u) - \varepsilon_p), \alpha) \, dx + \int_{\Omega} \frac{\lambda}{2} \text{tr} [\nabla u]^2 \, dx,$$

with (cf. (2.24))

$$(4.7) \quad \Psi_{\text{hard}}(X, \alpha) = \begin{cases} \frac{1}{2} \|X\|^2, & \|\text{dev } X\| \leq \sigma_y + H\alpha, \\ \frac{1}{2(1 + \frac{H}{1[\text{MPa}]})} \left(\frac{H}{1[\text{MPa}]} \|X\|^2 + \frac{1}{3} \text{tr} [X]^2 \right. \\ \quad \left. + 2(\sigma_y + H\alpha) \|\text{dev } X\| - (\sigma_y + H\alpha)^2 \right), & \|\text{dev } X\| > \sigma_y + H\alpha \end{cases}$$

$$= \frac{1}{2} \|X\|^2 - \frac{1}{2(1 + \frac{H}{1[\text{MPa}]})} [\|\text{dev } X\| - (\sigma_y + H\alpha)]_+^2,$$

whose second derivative coincides with the consistent tangent method introduced already in [61]. The constant $H > 0$ is the hardening modulus with dimension [MPa]. In this form it is easy to see that for positive hardening modulus $H > 0$ the isotropic hardening update potential is uniformly convex in $\text{sym } \nabla u$ with quadratic growth and has a Lipschitz continuous derivative. Therefore, our main regularity result applies also to this functional.⁴ The relative merits of each individual regularization scheme depend on their ability to balance regularization and approximation. Linear viscosity and hardening can be justified on physical grounds, but the (small) viscosity parameter $\eta > 0$ is difficult to estimate, as is the linear hardening modulus $H > 0$. The physically motivated regularization terms have the property to control only the symmetric part of the displacement gradient. The regularization (4.4)₃, however, does not satisfy the linearized frame-indifference condition.

All alternative regularization procedures thus establish local coercivity in the strains. In contrast, the Cosserat regularization is weaker in the sense that only strong Legendre–Hadamard ellipticity is reestablished, which, provided that displacement boundary data are prescribed, suffices for existence, uniqueness, and higher regularity. Thus the Cosserat approach appears as the weakest regularization among the ones considered.

⁴Repin [54, eq. (2.3)] calls (4.4)₂ linear hardening and shows the regularity $u^\delta \in H_{\text{loc}}^2(\Omega, \mathbb{R}^3)$, while for the planar case $n = 2$ he obtains $u^\delta \in H^2(\Omega, \mathbb{R}^2)$ if $\Gamma = \partial\Omega$ is smooth.

5. The regularity theorem. We know already that problem (2.22) has solutions $u^n \in H^1(\Omega, \mathbb{R}^3)$. Looking at the system for the microrotations A^n at given $\nabla u^n \in L^2(\Omega, \mathbb{M}^{3 \times 3})$ we realize at once that the linearity in A^n together with the Laplacian structure allows us to use standard elliptic regularity results for linear systems, which yields higher regularity for the microrotations: $A^n \in H^2(\Omega, \mathfrak{so}(3))$. In this section we study the regularity of the displacement field u^n , which is determined through (2.13).

5.1. Higher regularity for a quasi-linear elliptic system. The quasi-linear elliptic system introduced in section 2.4 is a special case of the systems which we define here below. For $d, m, N \geq 1$ and $\Omega \subset \mathbb{R}^d$ let $\mathcal{M} : \Omega \times \mathbb{M}^{m \times d} \times \mathbb{R}^N \rightarrow \mathbb{M}^{m \times d}$ be a matrix-valued function with the following properties.

R1. The mapping $\mathcal{M} : \Omega \times \mathbb{M}^{m \times d} \times \mathbb{R}^N \rightarrow \mathbb{M}^{m \times d}$ is a Carathéodory function which is Lipschitz continuous in the following sense: there exist constants $L_1, L_2 > 0$ such that for every $x, x_i \in \Omega$, $a, a_i \in \mathbb{M}^{m \times d}$, and $z, z_i \in \mathbb{R}^N$ we have

$$\begin{aligned} \|\mathcal{M}(x_1, a, z) - \mathcal{M}(x_2, a, z)\| &\leq L_1(\|a\| + \|z\|) \|x_1 - x_2\|, \\ \|\mathcal{M}(x, a_1, z_1) - \mathcal{M}(x, a_2, z_2)\| &\leq L_2(\|a_1 - a_2\| + \|z_1 - z_2\|), \\ \mathcal{M}(x, 0, 0) &= 0. \end{aligned}$$

Assumption R1 implies the useful estimate

$$(5.1) \quad \begin{aligned} &\|\mathcal{M}(x_1, a_1, z_1) - \mathcal{M}(x_2, a_2, z_2)\| \\ &\leq L_1(\|a_1\| + \|z_1\|) \|x_1 - x_2\| + L_2(\|a_1 - a_2\| + \|z_1 - z_2\|). \end{aligned}$$

R2. The mapping \mathcal{M} is strongly rank-one monotone. That means that there exists a constant $c_{LH} > 0$ such that for every $x \in \overline{\Omega}$, $a \in \mathbb{M}^{m \times d}$, $z \in \mathbb{R}^N$, $\xi \in \mathbb{R}^m$, and $\eta \in \mathbb{R}^d$ we have

$$(5.2) \quad \langle \mathcal{M}(x, a + \xi \otimes \eta, z) - \mathcal{M}(x, a, z), \xi \otimes \eta \rangle \geq c_{LH} \|\xi\|^2 \|\eta\|^2.$$

R3. The Gårding inequality shall be satisfied: there exist constants $C_G > 0$, $c_G \in \mathbb{R}$ such that for every $u_1, u_2 \in H^1(\Omega)$ with $u_1 - u_2 \in H_0^1(\Omega)$ and for every $z \in L^2(\Omega)$ the following inequality is valid:

$$\begin{aligned} \int_{\Omega} \langle \mathcal{M}(x, \nabla u_1, z) - \mathcal{M}(x, \nabla u_2, z), \nabla(u_1 - u_2) \rangle dx \\ \geq C_G \|\nabla(u_1 - u_2)\|_{L^2(\Omega)}^2 - c_G \|u_1 - u_2\|_{L^2(\Omega)}^2. \end{aligned}$$

Remark 5.1. If \mathcal{M} is differentiable, then the Gårding inequality already implies that \mathcal{M} is rank-one monotone; see, for example, [65, Thm. 6.1].

We investigate the regularity properties of weak solutions to the following quasi-linear elliptic boundary value problem. For given $g \in H^{\frac{1}{2}}(\partial\Omega)$, $z \in L^2(\Omega, \mathbb{R}^N)$, and $f \in L^2(\Omega, \mathbb{R}^m)$ find $u \in H^1(\Omega, \mathbb{R}^m)$ with $u|_{\partial\Omega} = g$ such that for every $v \in H_0^1(\Omega, \mathbb{R}^m)$ we have

$$(5.3) \quad \int_{\Omega} \langle \mathcal{M}(x, \nabla u(x), z(x)), \nabla v(x) \rangle dx = \int_{\Omega} \langle f, v \rangle dx.$$

THEOREM 5.2. *Let $\Omega \subset \mathbb{R}^d$ be a bounded $\mathcal{C}^{1,1}$ -smooth domain, $m \geq 1$, and $N \geq 1$, and assume that $\mathcal{M} : \Omega \times \mathbb{M}^{m \times d} \times \mathbb{R}^N \rightarrow \mathbb{M}^{m \times d}$ satisfies R1–R3. Furthermore, let*

$g \in H^{\frac{3}{2}}(\partial\Omega)$, $z \in H^1(\Omega)$, and $f \in L^2(\Omega)$. Every weak solution $u \in H^1(\Omega)$ of (5.3) with $u|_{\partial\Omega} = g$ is an element of $H^2(\Omega)$ and satisfies

$$\|u\|_{H^2(\Omega)} \leq c (\|g\|_{H^{\frac{3}{2}}(\partial\Omega)} + \|z\|_{H^1(\Omega)} + \|f\|_{L^2(\Omega)} + \|u\|_{H^1(\Omega)}).$$

Before we prove Theorem 5.2, we apply it to the situation described in section 2.4. There, $m = d = 3$ and \mathbb{R}^N is identified with $\text{Sym}(3) \times \mathfrak{so}(3)$ so that $z = (\varepsilon_p, A)$. Moreover,

$$\begin{aligned} \mathcal{M}(x, \nabla u, z) &= \mathcal{M}_C(\nabla u, \varepsilon_p, A) \\ &= P_K(2\mu(\text{sym } \nabla u - \varepsilon_p)) + \lambda(\text{tr } [\nabla u])\mathbb{1} + 2\mu_c(\text{skew}(\nabla u) - A). \end{aligned}$$

Since P_K is a Lipschitz continuous mapping, we see immediately that assumption R1 is satisfied. R2 is proved in Lemma 2.1, and the Gårding inequality is satisfied since \mathcal{M}_C generates a strongly monotone operator on $H_0^1(\Omega)$; see Lemma 2.2. Therefore, we have the following result for the reduced update problem (2.6).

THEOREM 5.3. *Let Ω be $C^{1,1}$ -smooth, $f^n \in L^2(\Omega)$, and $\varepsilon_p^{n-1} \in H^1(\Omega)$. Then $u_n \in H^2(\Omega)$, $A^n \in H^2(\Omega)$, and $\varepsilon_p^n \in H^1(\Omega)$.*

The proof of Theorem 5.2 is carried out with a difference quotient technique. We cover the boundary of Ω with a finite number of domains and map each of these domains with a $C^{1,1}$ diffeomorphism onto the unit cube in such a way that the image of the boundary of Ω lies on the midplane of the unit cube. We first prove higher regularity in directions tangential to the midplane by estimating difference quotients. The regularity in the normal direction is then obtained on the basis of the tangential regularity and by using the differential equation together with the rank-one monotonicity of \mathcal{M} .

Since \mathcal{M} is nonlinear and since we assumed rank-one monotonicity instead of strong monotonicity, we cannot use as test functions the usual finite differences of the type $h^{-1}\varphi^2(x)(u(x+h) - u(x))$, where φ is a cut-off function. Instead, we use differences which are based on inner variations. We begin the proof of Theorem 5.2 by studying a model problem on a half cube.

5.2. A model problem on a half cube. Let $C_r = \{x \in \mathbb{R}^d; |x_i| < r, 1 \leq i \leq d\}$ be a cube with side length $2r$, C_r^\pm the upper and lower half cube, respectively, and $M_r = \{x \in C_r; x_d = 0\}$ the midplane.

LEMMA 5.4. *Let $\Omega = C_1^-$, $f \in L^2(C_1^-)$, and $z \in H^1(C_1^-)$, and assume that $u \in H^1(C_1^-)$ with $u|_{M_1} = 0$ satisfies (5.3). Then for every $r \in (0, 1)$ and for $1 \leq i \leq d-1$ we have $\partial_i u \in H^1(C_r^-)$. Moreover, there is a constant $c_r > 0$ such that*

$$(5.4) \quad \|\partial_i u\|_{H^1(C_r^-)} \leq c_r (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|f\|_{L^2(C_1^-)}).$$

Proof. Let $r \in (0, 1)$ and $\varphi \in C_0^\infty(C_1)$ with $\varphi(x) = 1$ on C_r . For $h \in \mathbb{R}^d$ we introduce the mapping

$$\tau_h : C_1 \rightarrow \mathbb{R}^d : x \rightarrow \tau_h(x) = x + \varphi(x)h.$$

Let $h_0 = \|\varphi\|_{W^{1,\infty}(C_1)}^{-1} \min\{1, \text{dist}(\text{supp } \varphi, \partial C_1)\}$. For every $h \in \mathbb{R}^d$ with $|h| < h_0$ and h parallel to the plane M_1 , the mapping τ_h is a diffeomorphism from C_1 onto itself with $\tau_h(C_1^\pm) = C_1^\pm$, $\tau_h(M_1) = M_1$, and $\tau_h(x) = x$ for every $x \in \partial C_1$; see, e.g., [26]. Moreover, τ_h has the following properties (if $|h| < h_0$):

$$\begin{aligned} \nabla \tau_h(x) &= (\mathbb{1} + h \otimes \nabla \varphi(x)), \quad \det[\nabla \tau_h(x)] = 1 + \langle h, \nabla \varphi(x) \rangle, \\ \nabla_y \tau_h^{-1}(y) &= (\mathbb{1} + h \otimes \nabla \varphi)^{-1}|_{\tau_h^{-1}(y)} = \mathbb{1} - ((1 + \langle h, \nabla \varphi \rangle)^{-1} h \otimes \nabla \varphi)|_{\tau_h^{-1}(y)}. \end{aligned}$$

For a function $v : C_1^- \rightarrow \mathbb{R}^s$ we introduce

$$\Delta_h v = v \circ \tau_h - v, \quad \Delta^h v = v - v \circ \tau_h^{-1}.$$

For $f, g \in L^2(C_1^-)$, $|h| < h_0$ and $h \parallel M_1$ the following product rule is valid:

$$(5.5) \quad \int_{C_1^-} f \Delta^h g \, dx = - \int_{C_1^-} g \Delta_h f \, dx - \int_{C_1^-} (f \circ \tau_h g) \langle h, \nabla \varphi \rangle \, dx.$$

This identity can be shown by a transformation of coordinates $y = \tau_h(x)$ in the term $(g \circ \tau_h^{-1})f$. Let $u \in H_0^1(C_1^-)$ be a solution of (5.3). For $h \in \mathbb{R}^d$ with $|h| < h_0$ and $h \parallel M_1$ we define the double difference $v_h(x) = \Delta^h(\Delta_h u(x))$. In view of the assumptions on φ , h_0 , and h it follows that $v_h \in H_0^1(C_1^-)$. Inserting v_h into (5.3) yields

$$(5.6) \quad \int_{C_1^-} \langle \mathcal{M}(x, \nabla u, z), \nabla v_h \rangle \, dx = \int_{C_1^-} \langle f, v_h \rangle \, dx.$$

Note that $\nabla v_h = \Delta^h \nabla(\Delta_h u) + [(\det[\nabla \tau_h])^{-1}(\nabla \Delta_h u) h \otimes \nabla \varphi] \circ \tau_h^{-1}$, and therefore (5.6) is equivalent to

$$\begin{aligned} & \int_{C_1^-} \langle \mathcal{M}(x, \nabla u, z), \Delta^h \nabla(\Delta_h u) \rangle \, dx \\ &= - \int_{C_1^-} \langle \mathcal{M}(x, \nabla u, z), (\det[\nabla \tau_h])^{-1}(\nabla \Delta_h u) h \otimes \nabla \varphi \rangle \circ \tau_h^{-1} \, dx + \int_{C_1^-} \langle f, \Delta^h \Delta_h u \rangle \, dx. \end{aligned}$$

Furthermore, the product rule (5.5) entails

$$\begin{aligned} & \int_{C_1^-} \langle \Delta_h \mathcal{M}(x, \nabla u, z), \nabla \Delta_h u \rangle \, dx = - \int_{C_1^-} \langle \mathcal{M}(x, \nabla u, z) \circ \tau_h, \nabla \Delta_h u \rangle \langle h, \nabla \varphi \rangle \, dx \\ & \quad + \int_{C_1^-} \langle \mathcal{M}(x, \nabla u, z), ((\det[\nabla \tau_h])^{-1}(\nabla \Delta_h u) h \otimes \nabla \varphi) \circ \tau_h^{-1} \rangle \, dx \\ & \quad + \int_{C_1^-} \langle f, \Delta^h \Delta_h u \rangle \, dx \\ (5.7) \quad & =: S_1 + S_2 + S_3. \end{aligned}$$

Finally we have

$$\begin{aligned} & \int_{C_1^-} \langle \mathcal{M}(x, \nabla(u \circ \tau_h), z) - \mathcal{M}(x, \nabla u, z), \nabla \Delta_h u \rangle \, dx \\ &= \int_{C_1^-} \langle \Delta_h \mathcal{M}(x, \nabla u, z), \nabla \Delta_h u \rangle \, dx \\ & \quad + \int_{C_1^-} \langle \mathcal{M}(x, \nabla(u \circ \tau_h), z) - \mathcal{M}(x, \nabla u, z) \circ \tau_h, \nabla \Delta_h u \rangle \, dx \\ & \stackrel{(5.7)}{=} S_1 + S_2 + S_3 + \int_{C_1^-} \langle \mathcal{M}(x, \nabla(u \circ \tau_h), z) - \mathcal{M}(x, \nabla u, z) \circ \tau_h, \nabla \Delta_h u \rangle \, dx \\ (5.8) \quad &= S_1 + \dots + S_4. \end{aligned}$$

The next task is to show that there is a constant $c > 0$, which does not depend on h , such that

$$(5.9) \quad |S_1 + \dots + S_4| \leq c |h| (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|f\|_{L^2(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)}.$$

Due to the Lipschitz assumptions on \mathcal{M} we have

$$\begin{aligned} |S_1| + |S_2| &\leq c|h| \|\mathcal{M}(\cdot, \nabla u, z)\|_{L^2(C_1^-)} \|\Delta_h u\|_{H^1(C_1^-)} \\ &\leq c|h| (\|u\|_{H^1(C_1^-)} + \|z\|_{L^2(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)}. \end{aligned}$$

Moreover, since $f \in L^2(C_1^-)$, the term S_3 can be estimated as

$$|S_3| \leq c|h| \|f\|_{L^2(C_1^-)} \|\Delta_h u\|_{H^1(C_1^-)}.$$

By inequality (5.1) we see that

$$\begin{aligned} |S_4| &\leq cL_1 |h| (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)} \\ &\quad + cL_2 (\|\nabla(u \circ \tau_h) - (\nabla u) \circ \tau_h\|_{L^2(C_1^-)} + c|h| \|z\|_{H^1(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)}. \end{aligned}$$

The identity $\nabla(u \circ \tau_h) - (\nabla u) \circ \tau_h = (\nabla u) \circ \tau_h (h \otimes \nabla \varphi)$ leads to

$$|S_4| \leq c|h| (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)}.$$

Collecting all the above estimates we finally arrive at inequality (5.9). Gårding's inequality (see R3) and Poincaré's inequality imply that

$$\begin{aligned} \int_{C_1^-} \langle \mathcal{M}(x, \nabla(u \circ \tau_h), z) - \mathcal{M}(x, \nabla u, z), \nabla \Delta_h u \rangle dx \\ \geq C_G \|\nabla \Delta_h u\|_{L^2(C_1^-)}^2 - c_G \|\Delta_h u\|_{L^2(C_1^-)}^2 \\ \geq c(\|\Delta_h u\|_{H^1(C_1^-)}^2 - |h|^2 \|u\|_{H^1(C_1^-)}^2). \end{aligned}$$

Combining the above estimates with (5.8) and (5.9) results finally in

$$\begin{aligned} \|\Delta_h u\|_{H^1(C_1^-)}^2 &\leq c|h| (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|f\|_{L^2(C_1^-)}) \|\Delta_h u\|_{H^1(C_1^-)} \\ &\quad + c|h|^2 \|u\|_{H^1(C_1^-)}^2, \end{aligned}$$

and the constant c is independent of h . From Young's inequality we obtain

$$(5.10) \quad |h|^{-1} \|\Delta_h u\|_{H^1(C_1^-)} \leq c(\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|f\|_{L^2(C_1^-)}).$$

It follows from this inequality that $\partial_i u \in H^1(C_r^-)$ for $1 \leq i \leq d-1$ and that $\|\partial_i u\|_{H^1(C_r^-)}$ is bounded by the right-hand side in (5.10); see, e.g., [34]. \square

Remark 5.5. If we choose the usual finite differences as test functions, i.e., $\tilde{v}_h(x) = \delta_{-h}(\varphi^2 \delta_h u)$, where $\delta_h u = u(x+h) - u(x)$, then similar calculations as those for v_h lead to the estimate

$$(5.11) \quad \begin{aligned} \int_{C_1^-} \varphi^2(x) \langle \mathcal{M}(x, \nabla u(x+h), z(x)) - \mathcal{M}(x, \nabla u(x), z(x)), \delta_h \nabla u \rangle dx \\ \leq c|h| \|\varphi^2 \delta_h u\|_{H^1(C_1^-)}; \end{aligned}$$

compare also (5.8) and (5.9). But now neither R2 nor R3 helps us to find a lower bound for the left-hand side of (5.11) in terms of $\|\varphi^2 \delta_h \nabla u\|_{L^2(C_1^-)}^2$, since in general $\delta_h \nabla u$ is not a rank-one matrix, and since we cannot interchange φ and \mathcal{M} due to the nonlinearity of \mathcal{M} .

LEMMA 5.6 (regularity in the normal direction). *With the same assumptions as in Lemma 5.4 it follows for every $r \in (0, 1)$ that $\partial_d u \in H^1(C_r^-)$. Furthermore, there exists a constant $c_r > 0$ such that*

$$(5.12) \quad \|u\|_{H^2(C_r^-)} \leq c_r (\|z\|_{H^1(C_1^-)} + \|f\|_{L^2(C_1^-)} + \|u\|_{H^1(C_1^-)}).$$

Proof. Let $r \in (0, 1)$. Equation (5.3) implies that

$$(5.13) \quad \operatorname{Div} \mathcal{M}(x, \nabla u(x), z(x)) + f(x) = 0$$

for almost every $x \in C_1^-$. Let \mathcal{M}_i denote the columns of the matrix-valued function \mathcal{M} , i.e., $\mathcal{M}_i(x, a, z) = (\mathcal{M}_i^\alpha(x, a, z))_{1 \leq \alpha \leq m} \in \mathbb{R}^m$ for $1 \leq i \leq d$. The Lipschitz continuity of \mathcal{M} and the tangential regularity proved in Lemma 5.4 guarantee that $\partial_i \mathcal{M}_i(\cdot, \nabla u, z) \in L^2(C_r^-)$ for $1 \leq i \leq d-1$ and is bounded by the right-hand side in (5.4). Together with (5.13) we obtain therefore

$$\partial_d \mathcal{M}_d(\cdot, \nabla u, z) = -f - \sum_{i=1}^{d-1} \partial_i \mathcal{M}_i(\cdot, \nabla u, z) \in L^2(C_r^-).$$

By Lemma 7.23 in [27] the derivative ∂_d can be replaced with a finite difference in the following way: For every $\Omega' \subset\subset C_r^-$ and every $h \in \mathbb{R}^d$ with $|h| < \operatorname{dist}(\Omega', \partial C_r^-)$ and $h \perp M_1$ we have

$$(5.14) \quad \begin{aligned} \|\delta_h \mathcal{M}_d(\cdot, \nabla u, z)\|_{L^2(\Omega')} &\leq \left(\|f\|_{L^2(C_r^-)} + \sum_{i=1}^{d-1} \|\partial_i \mathcal{M}_i(\cdot, \nabla u, z)\|_{L^2(C_r^-)} \right) |h| \\ &=: c_0 |h|. \end{aligned}$$

Here, $\delta_h v(x) := v(x+h) - v(x)$ for $h \in \mathbb{R}^d$. Thus, for every $h \perp M_1$ with $|h| < \operatorname{dist}(\Omega', \partial C_r^-)$ we have

$$(5.15) \quad \int_{\Omega'} \langle \delta_h \mathcal{M}_d(x, \nabla u, z), \delta_h \partial_d u \rangle dx \leq c_0 |h| \|\delta_h \partial_d u\|_{L^2(\Omega')},$$

where c_0 is the constant from (5.14). We now split the left-hand side into a term which can be estimated from below due to the rank-one monotonicity of \mathcal{M} and into terms which may be estimated from above using the Lipschitz continuity of \mathcal{M} and the regularity results from Lemma 5.4. For functions $v : C_1^- \rightarrow \mathbb{R}^m$ we define $\tilde{\nabla} v(x) = (\partial_1 v(x), \dots, \partial_{d-1} v(x), 0) \in \mathbb{M}^{m \times d}$. Furthermore, $v_h(x) := v(x+h)$ and $e_d = (0, \dots, 0, 1)^\top \in \mathbb{R}^d$. With these notations we have

$$(5.16) \quad \begin{aligned} &\int_{\Omega'} \langle \mathcal{M}_d(x, \tilde{\nabla} u + \partial_d u_h \otimes e_d, z) - \mathcal{M}_d(x, \nabla u, z), \delta_h \partial_d u \rangle dx \\ &= \int_{\Omega'} \langle \delta_h \mathcal{M}_d(x, \nabla u, z), \delta_h \partial_d u \rangle dx \\ &\quad + \int_{\Omega'} \langle \mathcal{M}_d(x, \tilde{\nabla} u + \partial_d u_h \otimes e_d, z) - \mathcal{M}_d(x+h, \nabla u_h, z_h), \delta_h \partial_d u \rangle dx \\ &= S_1 + S_2. \end{aligned}$$

The term S_1 is already estimated in (5.15). From the Lipschitz continuity of \mathcal{M} (see (5.1)) and the regularity results of Lemma 5.4 we obtain by straightforward

calculations

$$(5.17) \quad \begin{aligned} |S_2| &\leq c \|\delta_h \partial_d u\|_{L^2(\Omega')} \left((\|\tilde{\nabla} u + \partial_d u_h \otimes e_d\|_{L^2(\Omega')} + \|z\|_{H^1(C_1^-)}) |h| + \|\delta_h \tilde{\nabla} u\|_{L^2(\Omega')} \right) \\ &\leq c |h| (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|\partial_d \tilde{\nabla} u\|_{L^2(C_r^-)}) \|\delta_h \partial_d u\|_{L^2(\Omega')}, \end{aligned}$$

and the constant c is independent of Ω' and h . Moreover, choosing $\xi = \partial_d u_h$ and $\eta = e_d$ in (5.2), we obtain for the left-hand side in (5.16) from the rank-one monotonicity of \mathcal{M} that

$$(5.18) \quad \begin{aligned} \int_{\Omega'} \langle \mathcal{M}_d(x, \tilde{\nabla} u + \partial_d u_h \otimes e_d, z) - \mathcal{M}_d(x, \nabla u, z), \delta_h \partial_d u \rangle dx \\ \geq c_{LH} \|\delta_h \partial_d u\|_{L^2(\Omega')}^2. \end{aligned}$$

Estimates (5.15)–(5.18) together with Young's inequality finally imply that

$$(5.19) \quad |h|^{-1} \|\delta_h \partial_d u\|_{L^2(\Omega')} \leq c (\|u\|_{H^1(C_1^-)} + \|z\|_{H^1(C_1^-)} + \|\partial_d \tilde{\nabla} u\|_{L^2(C_r^-)})$$

for every $h \perp M_1$. The constant c is independent of h and $\Omega' \Subset C_r^-$. This implies that $\partial_d^2 u \in L^2(C_r^-)$ and that $\|\partial_d^2 u\|_{L^2(C_r^-)}$ is bounded by the right-hand side in (5.19). Estimate (5.12) is a combination of (5.19) and (5.4). \square

5.3. Proof of Theorem 5.2. Let the assumptions of Theorem 5.2 be valid, and assume that $g = 0$. Choose $x_0 \in \partial\Omega$, and let U_{x_0} be a neighborhood of x_0 such that there exists a $\mathcal{C}^{1,1}$ diffeomorphism $\Phi_{x_0} : U_{x_0} \rightarrow C_1$, where C_1 is the unit cube in \mathbb{R}^d , with the following properties (we omit the index x_0): $\Phi(U) = C_1$, $\Phi(U \cap \Omega) = C_1^-$, $\Phi(U \setminus \bar{\Omega}) = C_1^+$, $\Phi(U \cap \partial\Omega) = M_1$, and $\Phi(x_0) = 0$. Let $u \in H_0^1(\Omega)$ be a solution for (5.3) with the data $f \in L^2(\Omega)$ and $z \in H^1(\Omega)$. It follows that

$$\int_{U \cap \Omega} \langle \mathcal{M}(x, \nabla u, z), \nabla v \rangle dx = \int_{U \cap \Omega} \langle f, v \rangle dx$$

for every $v \in H_0^1(\Omega \cap U)$. After a transformation of coordinates with $y = \Phi(x)$ and $\Psi := \Phi^{-1}$, the previous equation can be written as follows: Let $\tilde{u}(y) = u(\Psi(y))$. For every $v \in H_0^1(C_1^-)$ we have

$$\int_{C_1^-} \langle \tilde{\mathcal{M}}(y, \nabla \tilde{u}, \tilde{z}), \nabla v \rangle dy = \int_{C_1^-} \langle \tilde{f}, v \rangle dy.$$

Here, we use the abbreviations

$$(5.20) \quad \tilde{\mathcal{M}}(y, a, \zeta) = |\det[\nabla \Psi(y)]| \mathcal{M}(\Psi(y), a(\nabla \Psi(y))^{-1}, \zeta)(\nabla \Psi(y))^{-\top},$$

$$(5.21) \quad \tilde{f}(y) = |\det[\nabla \Psi(y)]| f(\Psi(y)),$$

$$(5.22) \quad \tilde{z}(y) = z(\Psi(y))$$

for $y \in C_1^-$, $a \in \mathbb{M}^{m \times d}$, and $\zeta \in \mathbb{R}^N$. It follows immediately from the properties of the diffeomorphism Φ and from those of \mathcal{M} that $\tilde{\mathcal{M}}$ satisfies R1–R3 with respect to C_1^- . Furthermore, \tilde{f} and \tilde{z} have the smoothness required in Lemma 5.4. Thus, Lemmata 5.4 and 5.6 guarantee that $\tilde{u} \in H^2(C_r^-)$ for every $r < 1$ and that estimate (5.12) is valid. After applying the inverse transformation $\Psi : C_1^- \rightarrow U \cap \Omega$, we have finally shown the following: For every $x_0 \in \bar{\Omega}$ there exists an open neighborhood \tilde{U}_{x_0} such

that $u|_{\tilde{U}_{x_0} \cap \Omega} \in H^2(\tilde{U}_{x_0} \cap \Omega)$ and estimate (5.12) is valid with respect to $\tilde{U}_{x_0} \cap \Omega$. The constants may depend on x_0 . Since Ω is assumed to be bounded, we can cover $\bar{\Omega}$ by a finite number of the domains \tilde{U}_{x_0} and obtain finally that $u \in H^2(\Omega)$ with

$$(5.23) \quad \|u\|_{H^2(\Omega)} \leq c(\|z\|_{H^1(\Omega)} + \|f\|_{L^2(\Omega)} + \|u\|_{H^1(\Omega)}).$$

This proves Theorem 5.2 for the case of vanishing Dirichlet conditions. The general case can be seen as follows. There exists a linear and continuous extension operator $F : H^{\frac{3}{2}}(\partial\Omega) \rightarrow H^2(\Omega)$ with $(F(g))|_{\partial\Omega} = g$ for every $g \in H^{\frac{3}{2}}(\partial\Omega)$; see, for example, [68]. Then $u \in H^1(\Omega)$ with $u|_{\partial\Omega} = g$ for some $g \in H^{\frac{3}{2}}(\partial\Omega)$ is a solution to (5.3) if and only if there exists an element $\tilde{u} \in H_0^1(\Omega)$ with $u = \tilde{u} + F(g)$ and for every $v \in H_0^1(\Omega)$, \tilde{u} satisfies

$$\int_{\Omega} \langle \hat{\mathcal{M}}(x, \nabla \tilde{u}, \tilde{z}), \nabla v \rangle \, dx = \int_{\Omega} \langle f, v \rangle \, dx,$$

where $\tilde{z} = (F(g), z)$ and $\hat{\mathcal{M}}(x, a, \tilde{z}) = \mathcal{M}(x, a + F(g)(x), z)$. Clearly, $\hat{\mathcal{M}}$ satisfies R1–R3 as well, and by the first part of this proof it follows that $\tilde{u} \in H^2(\Omega)$. This finishes the proof of Theorem 5.2.

6. Discussion. We have shown that the time-incremental Cosserat elasto-plasticity problem admits $H^1(\Omega)$ -regular updates of the symmetric plastic strain ε_p^n , provided that the previous plastic strain ε_p^{n-1} is in $H^1(\Omega)$ and the domain and data are suitably regular. Altogether, the time-incremental problem allows the regularity for all $n \in \mathbb{N} : u^n \in H^2(\Omega, \mathbb{R}^3)$, $\varepsilon_p^n \in H^1(\Omega, \text{Sym}(3))$, and $A^n \in H^2(\Omega, \mathfrak{so}(3))$. Uniform bounds in time are missing, and it is an open question whether a similar result holds for the time-continuous problem.

The presented method of proof for higher regularity uses a difference quotient method which is based on inner variations and can be extended to more general problems. This will be the subject of further investigations.

Appendix. Notation. We denote by $\mathbb{M}^{3 \times 3}$ the set of real 3×3 second order tensors, written with capital letters. The standard Euclidean scalar product on $\mathbb{M}^{3 \times 3}$ is given by $\langle X, Y \rangle_{\mathbb{M}^{3 \times 3}} = \text{tr}[XY^T]$, and thus the Frobenius tensor norm is $\|X\|^2 = \langle X, X \rangle_{\mathbb{M}^{3 \times 3}}$ (we use these symbols indifferently for tensors and vectors). The identity tensor on $\mathbb{M}^{3 \times 3}$ will be denoted by $\mathbb{1}$, so that $\text{tr}[X] = \langle X, \mathbb{1} \rangle$. We let Sym and PSym denote the symmetric and positive definite symmetric tensors, respectively. We adopt the usual abbreviations of Lie algebra theory; i.e., $\mathfrak{so}(3) := \{X \in \mathbb{M}^{3 \times 3} \mid X^T = -X\}$ are skew symmetric second order tensors, and $\mathfrak{sl}(3) := \{X \in \mathbb{M}^{3 \times 3} \mid \text{tr}[X] = 0\}$ are traceless tensors. We set $\text{sym}(X) = \frac{1}{2}(X^T + X)$ and $\text{skew}(X) = \frac{1}{2}(X - X^T)$ such that $X = \text{sym}(X) + \text{skew}(X)$. For $X \in \mathbb{M}^{3 \times 3}$ we set for the deviatoric part $\text{dev } X = X - \frac{1}{3} \text{tr}[X] \mathbb{1} \in \mathfrak{sl}(3)$.

For a second order tensor X we let $X.e_i$ be the application of the tensor X to the column vector e_i . The first and second differential of a scalar-valued function $W(F)$ are written $D_F W(F).H$ and $D_F^2 W(F).(H, H)$, respectively. Sometimes we use also $\partial_X W(X)$ to denote the first derivative of W with respect to X . We employ the standard notation of Sobolev spaces, i.e., $L^2(\Omega)$, $H^{1,2}(\Omega)$, $H_0^{1,2}(\Omega)$, which we use indifferently for scalar-valued functions as well as for vector-valued and tensor-valued functions.

Acknowledgments. The idea for this work was conceived during a visit of P.N. at the Weierstrass Institute, Berlin. The kind hospitality of A. Mielke is gratefully

acknowledged. P.N. has also profited from the exchange with M. Costabel regarding Div/Curl systems. P.N. and D.K. both acknowledge the help of inspiring discussions on regularity of nonlinear elliptic problems with J. Frehse.

REFERENCES

- [1] H. D. ALBER, *Materials with Memory. Initial-Boundary Value Problems for Constitutive Equations with Internal Variables*, Lecture Notes in Math. 1682, Springer, Berlin, 1998.
- [2] J. ALBERTY, C. CARSTENSEN, AND D. ZARRABI, *Adaptive numerical analysis in primal elastoplasticity with hardening*, Comput. Methods Appl. Mech. Engrg., 171 (1999), pp. 175–204.
- [3] G. ANZELLOTTI AND S. LUCKHAUS, *Dynamical evolution of elasto-perfectly plastic bodies*, Appl. Math. Optim., 15 (1987), pp. 121–140.
- [4] J. P. BARDET, *Observations on the effects of particle rotations on the failure of idealized granular materials*, Mech. Mater., 18 (1994), pp. 159–182.
- [5] A. BENSOUSSAN AND J. FREHSE, *Asymptotic behaviour of Norton-Hoff’s law in plasticity theory and H^1 regularity*, in Boundary Value Problems for Partial Differential Equations and Applications, J.-L. Lions and C. Baiocchi, eds., RMA Res. Notes Appl. Math. 29, Masson, Paris, 1993, pp. 3–25.
- [6] A. BENSOUSSAN AND J. FREHSE, *Regularity Results for Nonlinear Elliptic Systems and Applications*, Appl. Math. Sci. 151, Springer, Berlin, 2002.
- [7] M. BILDHAUER AND M. FUCHS, *Smoothness of weak solutions of the Ramberg/Osgood equations on plane domains*, ZAMM Z. Angew. Math. Mech., 87 (2007), pp. 70–76.
- [8] S. C. BRENNER AND L. R. SCOTT, *The Mathematical Theory of Finite Element Methods*, Springer, Heidelberg, 1994.
- [9] G. CAPRIZ, *Continua with Microstructure*, Springer, Heidelberg, 1989.
- [10] C. CARSTENSEN AND S. MÜLLER, *Local stress regularity in scalar nonconvex variational problems*, SIAM J. Math. Anal., 34 (2002), pp. 495–509.
- [11] K. CHELMIŃSKI, *Coercive approximation of viscoplasticity and plasticity*, Asymptot. Anal., 26 (2001), pp. 105–133.
- [12] K. CHELMIŃSKI, *Perfect plasticity as a zero relaxation limit of plasticity with isotropic hardening*, Math. Methods Appl. Sci., 24 (2001), pp. 117–136.
- [13] K. CHELMIŃSKI, *Global existence of weak-type solutions for models of monotone type in the theory of inelastic deformations*, Math. Methods Appl. Sci., 25 (2002), pp. 1195–1230.
- [14] B. D. COLEMAN AND M. L. HODGDON, *On shear bands in ductile materials*, Arch. Rational Mech. Anal., 90 (1985), pp. 219–247.
- [15] E. COSSERAT AND F. COSSERAT, *Théorie des Corps Déformables*, Librairie Scientifique A. Hermann et Fils (Translation: *Theory of Deformable Bodies*, NASA TT F-11 561, Washington, DC, 1968), Paris, 1909.
- [16] J. DANĚČEK AND E. VISZUS, *$\mathcal{L}^{2,\Phi}$ regularity for nonlinear elliptic systems of second order*, Electron. J. Differential Equations, 2002, article 20.
- [17] R. DE BORST, *A generalization of J_2 -flow theory for polar continua*, Comput. Methods Appl. Mech. Engrg., 103 (1992), pp. 347–362.
- [18] A. DEMYANOV, *Regularity in Prandtl–Reuss perfect plasticity*, in Workshop: Analysis and Numerics for Rate-Independent Processes, G. Dal Maso, G. Francfort, A. Mielke, and T. Roubíček, eds., Oberwolfach Report, Vol. 4, Oberwolfach, Germany, 2007, pp. 591–666.
- [19] A. DIETSCHÉ, P. STEINMANN, AND K. WILLAM, *Micropolar elastoplasticity and its role in localization*, Int. J. Plasticity, 9 (1993), pp. 813–831.
- [20] C. EBMEYER, *Mixed boundary value problems for nonlinear elliptic systems with p -structure in polyhedral domains*, Math. Nachr., 236 (2002), pp. 91–108.
- [21] C. EBMEYER AND J. FREHSE, *Mixed boundary value problems for nonlinear elliptic equations in multidimensional non-smooth domains*, Math. Nachr., 203 (1999), pp. 47–74.
- [22] A. C. ERINGEN, *Microcontinuum Field Theories*, Springer, Heidelberg, 1999.
- [23] S. FOREST, G. CAILLETAUD, AND R. SIEVERT, *A Cosserat theory for elastoviscoplastic single crystals at finite deformation*, Arch. Mech., 49 (1997), pp. 705–736.
- [24] M. FUCHS AND G. SEREGIN, *Variational Methods for Problems from Plasticity Theory and for Generalized Newtonian Fluids*, Ann. Univ. Sarav. Ser. Math. 10, FB Mathematik, University of Saarbrücken, Saarbrücken, Germany, 1999.
- [25] M. FUCHS AND G. SEREGIN, *Variational Methods for Problems from Plasticity Theory and for Generalized Newtonian Fluids*, Lecture Notes in Math. 1749, Springer, Berlin, 2000.
- [26] M. GIAQUINTA AND S. HILDEBRANDT, *Calculus of Variations I*, Springer, Berlin, Heidelberg, 1996.

- [27] D. GILBARG AND N. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, Grundlehren Math. Wiss. 224, Springer, Berlin, 1977.
- [28] V. GIRAULT AND P. A. RAVIART, *Finite Element Approximation of the Navier-Stokes Equations*, Lecture Notes in Math. 749, Springer, Heidelberg, 1979.
- [29] W. HAN AND B. D. REDDY, *Plasticity. Mathematical Theory and Numerical Analysis*, Springer, Berlin, 1999.
- [30] I. R. IONESCU AND M. SOFONEA, *Functional and Numerical Methods in Viscoplasticity*, 1st ed., Oxford University Press, Oxford, UK, 1993.
- [31] M. M. IORDACHE AND K. WILLAM, *Localized failure analysis in elastoplastic Cosserat continua*, Comput. Methods Appl. Mech. Engrg., 151 (1998), pp. 559–586.
- [32] D. KNEES, *Global regularity of the elastic fields of a power-law model on Lipschitz domains*, Math. Methods Appl. Sci., 29 (2006), pp. 1363–1391.
- [33] D. KNEES, *Global stress regularity of convex and some nonconvex variational problems*, Ann. Mat. Pura Appl. (4), 187 (2008), pp. 157–184.
- [34] D. KNEES AND A. MIELKE, *Energy release rate for cracks in finite-strain elasticity*, Math. Methods Appl. Sci., 31 (2008), pp. 501–528.
- [35] R. V. KOHN AND R. TEMAM, *Dual spaces of stresses and strains, with applications to Hencky plasticity*, Appl. Math. Optim., 10 (1983), pp. 1–35.
- [36] A. MIELKE, *Deriving new evolution equations for microstructures via relaxation of variational incremental problems*, Comput. Methods Appl. Mech. Engrg., 193 (2004), pp. 5095–5127.
- [37] A. MIELKE, *Existence of minimizers in incremental elasto-plasticity with finite strains*, SIAM J. Math. Anal., 36 (2004), pp. 384–404.
- [38] A. MIELKE AND S. MÜLLER, *Lower semicontinuity and existence of minimizers in incremental finite-strain elastoplasticity*, ZAMM Z. Angew. Math. Mech., 86 (2006), pp. 233–250.
- [39] C. B. MORREY, *Multiple Integrals in the Calculus of Variations*, Springer, New York, 1966.
- [40] H. B. MÜHLHAUS, *Shear band analysis for granular materials within the framework of Cosserat theory*, Ing. Archiv., 56 (1989), pp. 389–399.
- [41] J. NEČAS, *Les Méthodes Directes en Théorie des Équations Elliptiques*, Masson et Cie, Éditeurs, Paris, 1967.
- [42] J. NEČAS, *Introduction to the Theory of Nonlinear Elliptic Equations*, Teubner Verlagsgesellschaft, Leipzig, 1983.
- [43] P. NEFF, *Existence of minimizers for a finite-strain micromorphic elastic solid*, Proc. Roy. Soc. Edinburgh Sect. A, 136 (2006), pp. 997–1012.
- [44] P. NEFF, *A finite-strain elastic-plastic Cosserat theory for polycrystals with grain rotations*, Internat. J. Engrg. Sci., 44 (2006), pp. 574–594.
- [45] P. NEFF, *Finite Multiplicative Elastic-Viscoplastic Cosserat Micropolar Theory for Polycrystals with Grain Rotations. Modelling and Mathematical Analysis*, Preprint 2297; <http://wwwbib.mathematik.tu-darmstadt.de/Math-Net/Preprints/Listen/pp03.html>.
- [46] P. NEFF AND K. CHELMIŃSKI, *Infinitesimal elastic-plastic Cosserat micropolar theory. Modelling and global existence in the rate independent case*, Proc. Roy. Soc. Edinburgh Sect. A, 135 (2005), pp. 1017–1039.
- [47] P. NEFF AND K. CHELMIŃSKI, *A note on approximation of Prandtl-Reuss plasticity through Cosserat plasticity*, Quart. Appl. Math., to appear.
- [48] P. NEFF AND K. CHELMIŃSKI, *Well-posedness of dynamic Cosserat plasticity*, Appl. Math. Optim., 56 (2007), pp. 19–35.
- [49] P. NEFF, K. CHELMIŃSKI, W. MÜLLER, AND C. WIENERS, *Numerical solution method for an infinitesimal elastic-plastic Cosserat model*, Math. Models Methods Appl. Sci., 17 (2007), pp. 1211–1239.
- [50] P. NEFF AND C. WIENERS, *Comparison of models for finite plasticity. A numerical study*, Comput. Vis. Sci., 6 (2003), pp. 23–35.
- [51] S. NESENEENKO, *Homogenization and Regularity in Viscoplasticity*. Ph.D. thesis, TU Darmstadt, Logos Verlag, Berlin, 2006.
- [52] M. ORTIZ AND L. STAINIER, *The variational formulation of viscoplastic constitutive updates*, Comput. Methods Appl. Mech. Engrg., 171 (1999), pp. 419–444.
- [53] J.-P. RAYMOND, *Régularité globale des solutions de systèmes elliptiques non linéaires*, Rev. Mat. Univ. Complut. Madrid, 2 (1989), pp. 241–270.
- [54] S. I. REPIN, *Errors of finite element method for perfectly elasto-plastic problems*, Math. Models Methods Appl. Sci., 6 (1996), pp. 587–604.
- [55] M. RISTINMAA AND M. VECCHI, *Use of couple-stress theory in elasto-plasticity*, Comput. Methods Appl. Mech. Engrg., 136 (1996), pp. 205–224.
- [56] C. SANSOUR, *A theory of the elastic-viscoplastic Cosserat continuum*, Arch. Mech., 50 (1998), pp. 577–597.

- [57] G. SAVARÉ, *Regularity results for elliptic equations in Lipschitz domains*, J. Funct. Anal., 152 (1998), pp. 176–201.
- [58] G. A. SEREGIN AND T. N. SHILKIN, *Regularity for minimizers of some variational problems in plasticity theory*, J. Math. Sci., 99 (2000), pp. 969–988.
- [59] G. A. SEREGIN, *Differentiability of extremals of variational problems in the mechanics of elastic perfectly plastic media*, Differential Equations, 23 (1987), pp. 1349–1358.
- [60] J. C. SIMO AND J. R. HUGHES, *Computational Inelasticity*, Interdiscip. Appl. Math. 7, Springer, Berlin, 1998.
- [61] J. C. SIMO AND R. L. TAYLOR, *Consistent tangent operators for rate-independent elasto-plasticity*, Comput. Methods Appl. Mech. Engrg., 48 (1985), pp. 101–118.
- [62] P. STEINMANN, *A micropolar theory of finite deformation and finite rotation multiplicative elastoplasticity*, Internat. J. Solids Structures, 31 (1994), pp. 1063–1084.
- [63] R. TEMAM, *Mathematical Problems in Plasticity*, Gauthier–Villars, New York, 1985.
- [64] R. TEMAM, *A generalized Norton-Hoff model and the Prandtl-Reuss law of plasticity*, Arch. Rational Mech. Anal., 95 (1986), pp. 137–183.
- [65] T. VALENT, *Boundary Value Problems of Finite Elasticity*, Springer, Berlin, 1988.
- [66] C. WIENERS, *Multigrid methods for Prandtl-Reuss plasticity*, Numer. Linear Algebra Appl., 6 (1999), pp. 457–478.
- [67] C. WIENERS, *Efficient elasto-plastic simulation*, in Multifield Problems. State of the Art, A. M. Sändig, W. Schiehlen, and W. L. Wendland, eds., Springer, Berlin, 2000, pp. 209–216.
- [68] J. WLOKA, *Partielle Differentialgleichungen*, Teubner Verlag, Stuttgart, 1982.

GLOBAL EXISTENCE RESULTS AND UNIQUENESS FOR DISLOCATION EQUATIONS*

GUY BARLES[†], PIERRE CARDALIAGUET[‡], OLIVIER LEY[†], AND RÉGIS MONNEAU[§]

Abstract. We are interested in nonlocal eikonal equations arising in the study of the dynamics of dislocation lines in crystals. For these nonlocal but also nonmonotone equations, only the existence and uniqueness of Lipschitz and local-in-time solutions were available in some particular cases. In this paper, we propose a definition of weak solutions for which we are able to prove the existence for all time. Then we discuss the uniqueness of such solutions in several situations, both in the monotone and the nonmonotone case.

Key words. nonlocal Hamilton–Jacobi equations, dislocation dynamics, nonlocal front propagation, level-set approach, geometrical properties, lower-bound gradient estimate, viscosity solutions, eikonal equation, L^1 -dependence in time

AMS subject classifications. 49L25, 35F25, 35A05, 35D05, 35B50, 45G10

DOI. 10.1137/070682083

1. Introduction. In this article we are interested in the dynamics of defects in crystals, called dislocations. The dynamics of these dislocations is the main microscopic explanation of the macroscopic behavior of metallic crystals (see, for instance, the physical monographs of Nabarro [24], Hirth and Lothe [19], or Lardner [21] for a mathematical presentation). A dislocation is a line moving in a crystallographic plane, called a slip plane. The typical length of such a dislocation line is of the order of 10^{-6} m. Its dynamics is given by a normal velocity proportional to the Peach–Koehler force acting on this line.

This Peach–Koehler force may have two possible contributions: the first one is the self-force created by the elastic field generated by the dislocation line itself (i.e., this self-force is a nonlocal function of the shape of the dislocation line); the second one is the force created by everything exterior to the dislocation line, such as the exterior stress applied on the material, or the force created by other defects. In this paper, we study a particular model introduced in Rodney, Le Bouar, and Finel [27].

More precisely, if, at time t , the dislocation line is the boundary of an open set $\Omega_t \subset \mathbb{R}^N$ with $N = 2$ for the physical application, the normal velocity to the set Ω_t is given by

$$(1) \quad V_n = c_0 \star \mathbb{1}_{\overline{\Omega}_t} + c_1,$$

where $\mathbb{1}_{\overline{\Omega}_t}(x)$ is the indicator function of the set $\overline{\Omega}_t$, which is equal to 1 if $x \in \overline{\Omega}_t$ and equal to 0 otherwise. The function $c_0(x, t)$ is a kernel which depends only on the physical properties of the crystal and on the choice of the dislocation line, whose

*Received by the editors February 6, 2007; accepted for publication (in revised form) November 16, 2007; published electronically March 26, 2008. This work was supported by contracts ACI JC 1025 (2003-2005) and ACI JC 1041 (2002-2004) of the French Ministry of Research.

<http://www.siam.org/journals/sima/40-1/68208.html>

[†]Laboratoire de Mathématiques et Physique Théorique, Fédération Denis Poisson, Université de Tours, Parc de Grandmont, 37200 Tours, France (barles@lmpt.univ-tours.fr, ley@lmpt.univ-tours.fr).

[‡]Université de Bretagne Occidentale, UFR des Sciences et Techniques, 6 Av. Le Gorgeu, BP 809, 29285 Brest, France (pierre.cardaliaguet@univ-brest.fr).

[§]CERMICS, Ecole Nationale des Ponts et Chaussées, 6 et 8 avenue Blaise Pascal, Cité Descartes, Champs-sur-Marne, 77455 Marne-la-Vallée Cedex 2, France (monneau@cermics.enpc.fr).

evolution we follow. In the special case of application to dislocations, the kernel c_0 does not depend on time, but to keep a general setting we allow here a dependence on the time variable. Here \star denotes the convolution in space, namely

$$(2) \quad (c_0(\cdot, t) \star \mathbf{1}_{\bar{\Omega}_t})(x) = \int_{\mathbb{R}^N} c_0(x - y, t) \mathbf{1}_{\bar{\Omega}_t}(y) dy,$$

and this term appears to be the Peach–Koehler self-force created by the dislocation itself, while $c_1(x, t)$ is an additional contribution to the velocity, created by everything exterior to the dislocation line. We refer the reader to Alvarez et al. [3] for a detailed presentation and a derivation of this model.

We proceed as in the level-set approach to derive an equation for the dislocation line. We replace the evolution of a set Ω_t (the strong solution) by the evolution of a function u such that $\Omega_t = \{u(\cdot, t) > 0\}$. Roughly speaking the dislocation line is represented by the zero level-set of the function u which solves the following equation:

$$(3) \quad \begin{cases} \frac{\partial u}{\partial t} = (c_0(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq 0\}})(x) + c_1(x, t) |Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where (2) now reads

$$(4) \quad c_0(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x) = \int_{\mathbb{R}^N} c_0(x - y, t) \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(y) dy.$$

Note that (3) is not really a level-set equation since it is not invariant under non-decreasing changes of functions $u \rightarrow \varphi(u)$, where φ is nondecreasing. As noticed by Slepčev [28], the natural level-set equation should be (11); see section 1.2.

Although (3) seems very simple, there are only a few known results. Under suitable assumptions on the initial data and on c_0, c_1 , the existence and uniqueness of the solution is known in two particular cases: either for short time (see [3]), or for all time under the additional assumption that $V_n \geq 0$, which is, for instance, always satisfied for c_1 satisfying $c_1(x, t) \geq |c_0(\cdot, t)|_{L^1(\mathbb{R}^N)}$ (see [2, 12, 5] for a level-set formulation).

In the general case, the existence for all time of solutions to (3) is not known and, in particular, in the case when the kernel c_0 has negative values; indeed, in this case, the front propagation problem (3) does not satisfy any monotonicity property (preservation of inclusions), and therefore, even if a level-set-type equation can be derived, viscosity solution theory cannot be readily used. At this point, it is worth pointing out that a key property in the level-set approach is the comparison principle for viscosity solutions which is almost equivalent to this monotonicity property (see, for instance, Giga's monograph [18]). On the other hand, one may try to partly use viscosity solution theory together with some other approximation and/or compactness arguments to prove at least the existence of weak solutions (in a suitable sense). But here also the bad sign of the kernel creates difficulties since one cannot readily use the classical half-relaxed limits techniques to pass to the limit in the approximate problems. Additional arguments are needed to obtain weak solutions.

The aim of this paper is to describe a general approach of these dislocations' dynamics, based on the level-set approach, which allows us to introduce a suitable notion of weak solutions, to prove the existence of these weak solutions for all time, and to analyze the uniqueness (or nonuniqueness) of these solutions.

1.1. Weak solutions of the dislocation equation. We introduce the following definition of weak solutions, which itself uses the definition of L^1 -viscosity solutions, recalled in Appendix A.

DEFINITION 1.1 (classical and weak solutions). *For any $T > 0$, we say that a function $u \in W^{1,\infty}(\mathbb{R}^N \times [0, T])$ is a weak solution of (3) on the time interval $[0, T]$ if there is some measurable map $\chi : \mathbb{R}^N \times (0, T) \rightarrow [0, 1]$ such that u is an L^1 -viscosity solution of*

$$(5) \quad \begin{cases} \frac{\partial u}{\partial t} = \bar{c}(x, t)|Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where

$$(6) \quad \bar{c}(x, t) = c_0(\cdot, t) \star \chi(\cdot, t)(x) + c_1(x, t)$$

and

$$(7) \quad \mathbf{1}_{\{u(\cdot, t) > 0\}}(x) \leq \chi(x, t) \leq \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x)$$

for almost all $(x, t) \in \mathbb{R}^N \times [0, T]$. We say that u is a classical solution of (3) if u is a weak solution to (5) and if

$$(8) \quad \mathbf{1}_{\{u(\cdot, t) > 0\}}(x) = \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x)$$

for almost all $(x, t) \in \mathbb{R}^N \times [0, T]$.

Note that we have $\chi(x, t) = \mathbf{1}_{\{u(\cdot, t) > 0\}}(x) = \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x)$ for almost all $(x, t) \in \mathbb{R}^N \times [0, T]$ for classical solutions.

To state our first existence result, we use the following assumptions.

(H0) $u_0 \in W^{1,\infty}(\mathbb{R}^N)$, $-1 \leq u_0 \leq 1$, and there exists $R_0 > 0$ such that $u_0(x) \equiv -1$ for $|x| \geq R_0$,

(H1) $c_0 \in C([0, T]; L^1(\mathbb{R}^N))$, $D_x c_0 \in L^\infty([0, T]; L^1(\mathbb{R}^N))$, $c_1 \in C(\mathbb{R}^N \times [0, T])$, and there exist constants M_1, L_1 such that, for any $x, y \in \mathbb{R}^N$ and $t \in [0, T]$,

$$|c_1(x, t)| \leq M_1 \quad \text{and} \quad |c_1(x, t) - c_1(y, t)| \leq L_1|x - y|.$$

In what follows, we denote by M_0, L_0 constants such that, for any (or almost every) $t \in [0, T]$, we have

$$|c_0(\cdot, t)|_{L^1(\mathbb{R}^N)} \leq M_0 \quad \text{and} \quad |D_x c_0(\cdot, t)|_{L^1(\mathbb{R}^N)} \leq L_0.$$

Our first main result is the following.

THEOREM 1.2 (existence of weak solutions). *Under assumptions (H0)–(H1), for any $T > 0$ and for any initial data u_0 , there exists a weak solution of (3) on the time interval $[0, T]$ in the sense of Definition 1.1.*

Our second main result states that a weak solution is a classical one if the evolving set is expanding and if the following additional condition is fulfilled

(H2) c_1 and c_0 satisfy (H1), and there exist constants m_0, N_1 and a positive function $N_0 \in L^1(\mathbb{R}^N)$ such that, for any $x, h \in \mathbb{R}^N$, $t \in [0, T]$, we have

$$\begin{aligned} |c_0(x, t)| &\leq m_0, \\ |c_1(x + h, t) + c_1(x - h, t) - 2c_1(x, t)| &\leq N_1|h|^2, \\ |c_0(x + h, t) + c_0(x - h, t) - 2c_0(x, t)| &\leq N_0(x)|h|^2. \end{aligned}$$

THEOREM 1.3 (some links between weak solutions and classical continuous viscosity solutions and uniqueness results). *Assume (H0)–(H1), and suppose that there is some $\delta \geq 0$ such that, for all measurable maps $\chi : \mathbb{R}^N \times (0, T) \rightarrow [0, 1]$,*

$$(9) \quad \text{for all } (x, t) \in \mathbb{R}^N \times [0, T], \quad c_0(\cdot, t) \star \chi(\cdot, t)(x) + c_1(x, t) \geq \delta,$$

and that the initial data u_0 satisfies (in the viscosity sense)

$$(10) \quad -|u_0| - |Du_0| \leq -\eta_0 \quad \text{in } \mathbb{R}^N$$

for some $\eta_0 > 0$. Then any weak solution u of (3) in the sense of Definition 1.1 is a classical continuous viscosity solution of (3). This solution is unique if (H2) holds and

- (i) *either $\delta > 0$,*
- (ii) *or $\delta = 0$ and u_0 is semiconvex, i.e., satisfies for some constant $C > 0$*

$$u_0(x+h) + u_0(x-h) - 2u_0(x) \geq -C|h|^2 \quad \text{for all } x, h \in \mathbb{R}^N.$$

Assumption (9) ensures that the velocity V_n in (1) is positive for positive δ . Of course, we can state similar results in the case of negative velocity. Assumption (10) means that u_0 is a viscosity subsolution of $-|v(x)| - |Dv(x)| + \eta_0 \leq 0$. When u_0 is C^1 , it follows that the gradient of u_0 does not vanish on the set $\{u_0 = 0\}$ (see [22] for details). Point (ii) of the theorem is the main result of [2, 5]. We also point out that, with adapted proofs, only a bound from below could be required in (H2) on $c_1(x+h, t) + c_1(x-h, t) - 2c_1(x, t)$ and $c_0(x+h, t) + c_0(x-h, t) - 2c_0(x, t)$.

Remark 1.1. In particular Theorem 1.3 implies uniqueness in the case $c_0 \geq 0$ and $c_1 \equiv 0$. The general study of nonnegative kernels is provided below (see Theorem 1.5 and Remark 1.3).

1.2. Nonnegative kernel $c_0 \geq 0$. In the special case where the kernel c_0 is nonnegative, an inclusion principle for the dislocation lines, or equivalently a comparison principle for the functions of the level-set formulations, is expected (cf. Cardaliaguet [11] and Slepčev [28]).

Moreover, in the classical level-set approach, all the level-sets of u should have the same type of normal velocity, and Slepčev [28] remarked that a formulation with a nonlocal term of the form $\{u(\cdot, t) \geq u(x, t)\}$ is more appropriate. Therefore it is natural to start studying the following equation (which replaces (3)):

$$(11) \quad \begin{cases} \frac{\partial u}{\partial t} = (c_0(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq u(x, t)\}})(x) + c_1(x, t)|Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where \star denotes the convolution in space as in (4).

The precise meaning of a viscosity solution of (11) is given in Definition 5.1.

In this context, assumption (H0) can be weakened into the following condition, which allows us to consider unbounded evolving sets:

$$(H0') \quad u_0 \in BUC(\mathbb{R}^N).$$

Our main result for this equation is the following.

THEOREM 1.4 (existence and uniqueness). *Assume that $c_0 \geq 0$ on $\mathbb{R}^N \times [0, T]$ and that (H0')–(H1) hold. Then there exists a unique viscosity solution u of (11).*

Remark 1.2. The comparison principle for this equation (see Theorem 5.2) is a generalization of [28, Theorem 2.3]: indeed, in [28], everything takes place in a fixed

bounded set, whereas here one has to deal with unbounded sets. See also [15] for related results.

Now we turn to the connections with weak solutions. To do so, if u is the unique continuous solution of (11) given by Theorem 1.4, we introduce the functions $\rho^+, \rho^- : \mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}$ defined by

$$\rho^+ := \mathbf{1}_{\{u \geq 0\}} \quad \text{and} \quad \rho^- := \mathbf{1}_{\{u > 0\}}.$$

Our result is the following.

THEOREM 1.5 (maximal and minimal weak solutions). *Under the assumptions of Theorem 1.4, the maximal and minimal weak solutions of (3) are the continuous functions v^+, v^- which are the unique L^1 -viscosity solutions of the equations*

$$(12) \quad \begin{cases} \frac{\partial v^\pm}{\partial t} = c[\rho^\pm](x, t) |Dv^\pm| & \text{in } \mathbb{R}^N \times (0, T), \\ v^\pm(x, 0) = u_0(x) & \text{in } \mathbb{R}^N, \end{cases}$$

where

$$c[\rho](x, t) := c_0(\cdot, t) \star \rho(\cdot, t)(x) + c_1(x, t) \quad \text{in } \mathbb{R}^N \times (0, T).$$

The functions v^\pm satisfy $\{v^+(\cdot, t) \geq 0\} = \{u(\cdot, t) \geq 0\}$ and $\{v^-(\cdot, t) > 0\} = \{u(\cdot, t) > 0\}$, where u is the solution of (11).

Moreover, if the set $\{u(\cdot, t) = 0\}$ has a zero-Lebesgue measure for almost all $t \in (0, T)$, then problem (3) has a unique weak solution which is also a classical one.

Remark 1.3. 1. Theorem 1.5 shows that, in the case when $c_0 \geq 0$, Slepčev's approach allows us to identify the maximal and minimal weak solutions as being associated with ρ^\pm .

2. Equalities $\{v^-(\cdot, t) \geq 0\} = \{u(\cdot, t) \geq 0\}$ and $\{v^+(\cdot, t) > 0\} = \{u(\cdot, t) > 0\}$ do not hold in general (see, for instance, Example 3.1 in section 3).

3. If the set $\{v^\pm(\cdot, t) = 0\}$ develops an interior, a dramatic loss of uniqueness for the weak solution of (3) may occur. This is illustrated by Example 3.1 below, where we are able to build infinitely many solutions after the onset of fattening.

4. We have uniqueness for (3) if $\{u(\cdot, t) = 0\}$ has a zero-Lebesgue measure for almost all $t \in (0, T)$. This condition is fulfilled when, for instance, $c[\rho] \geq 0$ holds for any indicator function ρ and (10) holds (see also Remark 1.1).

1.3. Organization of the paper. In section 2, we recall basic results for the classical eikonal equation which are used throughout the paper. In section 3, we prove the existence of weak solutions for (3), namely Theorem 1.2, and give a counterexample to the uniqueness in general. Let us mention that this first part of the paper, even if it requires rather deep results of viscosity solution theory, is of a general interest for a wide audience and can be read without having an expertise in this theory since one just needs to apply the results which, anyway, are rather natural. In section 4, we prove Theorem 1.3 in the case of expanding dislocations. The arguments we use here are far more involved from a technical point of view: in particular we need some fine estimates of the perimeter of the evolving sets. In section 5, we study the Slepčev formulation in the case of nonnegative kernels and prove Theorems 1.4 and 1.5. In spirit, this section is closely related to the classical level-set approach but is more technical. Finally, for sake of completeness, we recall in Appendix A the definition of L^1 -viscosity solutions and a new stability result proved by Barles in [4].

2. Some basic results for the classical (local) eikonal equation. We want to recall in this section some basic results on the level-set equation

$$(13) \quad \begin{cases} \frac{\partial v}{\partial t} = a(x, t)|Dv| & \text{in } \mathbb{R}^N \times (0, T), \\ v(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where $T > 0$ and $a : \mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}$ is, at least, a continuous function.

We provide some classical estimates on the solutions to (13) when a satisfies suitable assumptions. Our result is the following.

THEOREM 2.1. *If u_0 satisfies (H0) and a satisfies the assumptions of c_1 in (H1), then (13) has a unique continuous solution v which is Lipschitz continuous in $\mathbb{R}^N \times [0, T]$ and which satisfies*

- (i) $-1 \leq v \leq 1$ in $\mathbb{R}^N \times (0, T)$, $v(x, t) \equiv -1$ for $|x| \geq R_0 + M_1 t$,
- (ii) $|Dv(\cdot, t)|_\infty \leq |Du_0|_\infty e^{L_1 t}$,
- (iii) $|v_t(\cdot, t)|_\infty \leq M_1 |Du_0|_\infty e^{L_1 t}$.

We skip the very classical proof of Theorem 2.1; we just point out that the first point comes from the comparison result for (13) and the “finite speed of propagation property” (see Crandall and Lions [14]), while the second one is a basic gradient estimate (see, for example, Ley [22]), and the last one comes directly from the fact that the equation is satisfied almost everywhere.

The main consequence of this result is that the solution remains in a compact subset of the Banach space $(C(\mathbb{R}^N \times [0, T]), |\cdot|_\infty)$ as long as u_0 and a satisfy (H0)–(H1) with fixed constants.

Let us introduce the following.

DEFINITION 2.2 (interior ball property). *We say that a closed set $K \subset \mathbb{R}^N$ has an interior ball property of radius $r > 0$ if, for any $x \in K$, there exists $p \in \mathbb{R}^N \setminus \{0\}$ such that $B(x - r \frac{p}{|p|}, r) \subset K$.*

We will also use the following result, due to Cannarsa and Frankowska [10], the proof of which is given in Appendix B for the sake of completeness.

LEMMA 2.3 (interior ball regularization). *Suppose (H0) and that a satisfies the assumptions of c_1 in (H1)–(H2) and there exists a constant $\delta > 0$ such that*

$$c_1 \geq \delta > 0 \quad \text{on } \mathbb{R}^N \times [0, T].$$

Then there exists a constant γ (depending in particular on $\delta > 0$ and T and on the other constants of the problem) such that for the solution v of (13), the set $\{v(\cdot, t) \geq 0\}$ has an interior ball property of radius $r_t \geq \gamma t$ for $t \in (0, T)$.

3. Existence of weak solutions for (3). We aim to solve (3), i.e.,

$$\begin{cases} \frac{\partial u}{\partial t} = (c_0(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x) + c_1(x, t))|Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

proving Theorem 1.2, which states the existence of weak solutions as introduced in Definition 1.1.

A key difficulty in solving (3) comes from the fact that, in this kind of level-set equation, one may face the so-called nonempty interior difficulty, i.e., that the 0-level-set of the solution is “fat,” which may mean that it has either a nonempty interior or a nonzero Lebesgue measure. Clearly, in both cases, $\mathbf{1}_{\{u(\cdot, t) \geq 0\}}$ is different from

$\mathbf{1}_{\{u(\cdot,t)>0\}}$, and this leads to rather bad stability properties for (3) and therefore to difficulties in proving the existence of a solution (and even more for the uniqueness). The notion of weak solution (5)–(7) emphasizes this difficulty. On the contrary, if $\bar{c}(x,t) \geq 0$ in $\mathbb{R}^N \times [0, T]$, it is known that the “nonempty interior difficulty” cannot happen (see Barles, Soner, and Souganidis [6] and Ley [22]), and we recover a more classical formulation. We discuss this question in the next section as well as some uniqueness issues for our weak solutions. Let us finally note that weak solutions for (3) satisfy the following inequalities.

PROPOSITION 3.1. *Let u be a weak solution to (3). Then u also satisfies in the L^1 -sense*

$$(14) \quad \frac{\partial u}{\partial t} \leq (c_0^+(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) \geq 0\}}(x) - c_0^-(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) > 0\}}(x) + c_1(x, t)) |Du|,$$

$$(15) \quad \frac{\partial u}{\partial t} \geq (c_0^+(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) > 0\}}(x) - c_0^-(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) \geq 0\}}(x) + c_1(x, t)) |Du|$$

in $\mathbb{R}^N \times (0, T)$, where $c_0^+ = \max(0, c_0)$ and $c_0^- = \max(0, -c_0)$.

Proof of Proposition 3.1. Let \bar{c} be associated with u as in (5)–(7). Then we have

$$\bar{c}(x, t) \geq c_0^+(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) > 0\}}(x) - c_0^-(\cdot, t) \star \mathbf{1}_{\{u(\cdot,t) \geq 0\}}(x) + c_1(x, t)$$

for every $x \in \mathbb{R}^N$ and almost every $t \in (0, T)$. We note that the right-hand side of the inequality is lower semicontinuous. Following Lions and Perthame [23], u then solves (15) in the usual viscosity sense. The proof of (14) can be achieved in a similar way. \square

Proof of Theorem 1.2. 1. *Introduction of a perturbed equation.* First we are going to solve the equation

$$(16) \quad \frac{\partial u}{\partial t} = (c_0(\cdot, t) \star \psi_\varepsilon(u(\cdot, t)))(x) + c_1(x, t) |Du| \quad \text{in } \mathbb{R}^N \times (0, T),$$

where $\psi_\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ is a sequence of continuous functions such that $\psi_\varepsilon(t) \equiv 0$ for $t \leq -\varepsilon$, $\psi_\varepsilon(t) \equiv 1$ for $t \geq 0$, and ψ_ε is an affine function on $[-\varepsilon, 0]$.

We aim at applying Schauder’s fixed point theorem to a suitable map. We note that an alternative proof could be given by using techniques developed by Alibaud in [1].

2. *Definition of a map \mathcal{T} .* We introduce the convex and compact (by Ascoli’s theorem) subset

$$X = \{u \in C(\mathbb{R}^N \times [0, T]) : u \equiv -1 \text{ in } \mathbb{R}^N \setminus B(0, R_0 + MT), \\ |Du|, |u_t|/M \leq |Du_0|_\infty e^{LT}\}$$

of $(C(\mathbb{R}^N \times [0, T]), |\cdot|_\infty)$, for $M = M_0 + M_1$ and $L = L_0 + L_1$, and the map $\mathcal{T} : X \rightarrow X$ defined as follows: if $u \in C(\mathbb{R}^N \times [0, T])$, then $\mathcal{T}(u)$ is the unique solution v of (13) for

$$c_\varepsilon(x, t) = c_0(\cdot, t) \star \psi_\varepsilon(u(\cdot, t))(x) + c_1(x, t) \\ = \int_{\mathbb{R}^N} c_0(x - z, t) \psi_\varepsilon(u(z, t)) dz + c_1(x, t).$$

This definition is justified by the fact that, under assumption (H1) on c_1 and c_0 , c_ε satisfies (H1) with fixed constants $M = M_0 + M_1$ and $L = L_0 + L_1$; indeed M is a bound on $\sup_{[0,T]} |c_0(\cdot, t)|_{L^1} + M_1$, while L is estimated by the following calculation: for all $x, y \in \mathbb{R}^N$, $t \in [0, T]$, and $u \in X$, we have

$$\begin{aligned}
(17) \quad & c_\varepsilon(x, t) - c_\varepsilon(y, t) \\
&= \int_{\mathbb{R}^N} (c_0(x - z, t) - c_0(y - z, t)) \psi_\varepsilon(u(z, t)) dz + c_1(x, t) - c_1(y, t) \\
&\leq \int_{\mathbb{R}^N} |c_0(x - z, t) - c_0(y - z, t)| dz + |c_1(x, t) - c_1(y, t)| \\
&\leq (L_0 + L_1)|x - y|,
\end{aligned}$$

since $0 \leq \psi_\varepsilon \leq 1$.

Finally, under assumptions (H0)–(H1), for any $u \in X$, the results of Theorem 2.1 apply to (16), which imply that $\mathcal{T}(u) \in X$. It follows that \mathcal{T} is well defined.

3. *Application of Schauder's fixed point theorem to \mathcal{T} .* The map \mathcal{T} is continuous since ψ_ε is continuous by using the classical stability result for viscosity solutions (see, for instance, (30) in section 4). Therefore \mathcal{T} has a fixed point u_ε which is bounded in $W^{1,\infty}(\mathbb{R}^N \times [0, T])$ uniformly with respect to ε (since M and L are independent of ε).

4. *Convergence of the fixed point when $\varepsilon \rightarrow 0$.* From Ascoli's theorem, we extract a subsequence $(u_{\varepsilon'})_{\varepsilon'}$ which converges locally uniformly to a function denoted by u (in fact globally since the $u_{\varepsilon'}$ are equal to -1 outside a fixed compact subset).

The functions $\chi_{\varepsilon'} := \psi_{\varepsilon'}(u_{\varepsilon'})$ satisfy $0 \leq \chi_{\varepsilon'} \leq 1$. Therefore we can extract a subsequence—still denoted $(\chi_{\varepsilon'})$ —which converges weakly- $*$ in $L_{\text{loc}}^\infty(\mathbb{R}^N \times [0, T])$ to some function $\chi : \mathbb{R}^N \times (0, T) \rightarrow [0, 1]$. Therefore, for all $\varphi \in L_{\text{loc}}^1(\mathbb{R}^N \times [0, T])$,

$$(18) \quad \int_0^T \int_{\mathbb{R}^N} \varphi \chi_{\varepsilon'} dx dt \rightarrow \int_0^T \int_{\mathbb{R}^N} \varphi \chi dx dt.$$

From Fatou's lemma, if φ is nonnegative, it follows that

$$\begin{aligned}
\int_0^T \int_{\mathbb{R}^N} \varphi(x, t) \chi(x, t) dx dt &\leq \int_0^T \int_{\mathbb{R}^N} \varphi(x, t) \limsup_{\varepsilon' \rightarrow 0} \chi_{\varepsilon'}(x, t) dx dt \\
&\leq \int_0^T \int_{\mathbb{R}^N} \varphi(x, t) \limsup_{\varepsilon' \rightarrow 0, x' \rightarrow x, t' \rightarrow t} \chi_{\varepsilon'}(x', t') dx dt \\
&\leq \int_0^T \int_{\mathbb{R}^N} \varphi(x, t) \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x) dx dt.
\end{aligned}$$

Since the previous inequalities hold for any nonnegative $\varphi \in L_{\text{loc}}^1(\mathbb{R}^N \times [0, T])$, we obtain that, for almost every $(x, t) \in \mathbb{R}^N \times (0, T)$,

$$\chi(x, t) \leq \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x).$$

Similarly we get

$$\mathbf{1}_{\{u(\cdot, t) > 0\}}(x) \leq \chi(x, t).$$

Furthermore, setting $c_{\varepsilon'} = c_0 \star \chi_{\varepsilon'} + c_1$, from (18), we have, for all $(x, t) \in \mathbb{R}^N \times [0, T]$,

$$\begin{aligned} \int_0^t c_{\varepsilon'}(x, s) ds &= \int_0^t \int_{\mathbb{R}^N} c_0(x-y, s) \chi_{\varepsilon'}(y, s) dy ds + \int_0^t c_1(x, s) ds \\ &\rightarrow \int_0^t \bar{c}(x, s) ds, \end{aligned}$$

where $\bar{c}(x, t) = c_0(\cdot, t) \star \chi(\cdot, t)(x) + c_1(x, t)$. The above convergence is pointwise, but, noticing that $c_{\varepsilon'}$ satisfies (H3) (with $M := M_0 + M_1$ and $L := L_0 + L_1$) and using Remark A.1, we can apply the stability theorem (Theorem A.3) given in Appendix A. We obtain that u is an L^1 -viscosity solution to (5) with \bar{c} satisfying (6)–(7). \square

The following example is inspired from [6].

Example 3.1 (counterexample to the uniqueness of weak solutions). Let us consider, in dimension $N = 1$, the following equation of type (3),

$$(19) \quad \begin{cases} \frac{\partial U}{\partial t} = (1 \star \mathbf{1}_{\{U(\cdot, t) \geq 0\}}(x) + c_1(t)) |DU| & \text{in } \mathbb{R} \times (0, 2], \\ U(\cdot, 0) = u_0 & \text{in } \mathbb{R}, \end{cases}$$

where we set $c_0(x, t) := 1$, $c_1(x, t) := c_1(t) = 2(t-1)(2-t)$, and $u_0(x) = 1 - |x|$. Note that $1 \star \mathbf{1}_A = \mathcal{L}^1(A)$ for any measurable set $A \subset \mathbb{R}$, where $\mathcal{L}^1(A)$ is the Lebesgue measure on \mathbb{R} .

We start by solving auxiliary problems for time in $[0, 1]$ and $[1, 2]$ in order to produce a family of solutions for the original problem in $[0, 2]$.

1. *Construction of a solution for $0 \leq t \leq 1$.* The function $x_1(t) = (t-1)^2$ is the solution of the ODE

$$\dot{x}_1(t) = c_1(t) + 2x_1(t) \text{ for } 0 \leq t \leq 1, \quad \text{and } x_1(0) = 1$$

(note that $\dot{x}_1 \leq 0$ in $[0, 1]$). Consider

$$(20) \quad \begin{cases} \frac{\partial u}{\partial t} = \dot{x}_1(t) \left| \frac{\partial u}{\partial x} \right| & \text{in } \mathbb{R} \times (0, 1], \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}. \end{cases}$$

There exists a unique continuous viscosity solution u of (20). Looking for u under the form $u(x, t) = v(x, \Gamma(t))$ with $\Gamma(0) = 0$, we obtain that v satisfies

$$\frac{\partial v}{\partial t} \dot{\Gamma}(t) = \dot{x}_1(t) \left| \frac{\partial v}{\partial x} \right|.$$

Choosing $\Gamma(t) = -x_1(t) + 1$, we get that v is the solution of

$$\begin{cases} \frac{\partial v}{\partial t} = - \left| \frac{\partial v}{\partial x} \right| & \text{in } \mathbb{R} \times (0, 1], \\ v(\cdot, 0) = u_0 & \text{in } \mathbb{R}. \end{cases}$$

By the Oleinik–Lax formula, $v(x, t) = \inf_{|x-y| \leq t} u_0(y)$. Since u_0 is even, we have, for all $(x, t) \in \mathbb{R} \times [0, 1]$,

$$u(x, t) = \inf_{|x-y| \leq \Gamma(t)} u_0(y) = u_0(|x| + \Gamma(t)) = u_0(|x| - x_1(t) + 1).$$

Therefore, for $0 \leq t \leq 1$,

$$(21) \quad \{u(\cdot, t) > 0\} = (-x_1(t), x_1(t)) \quad \text{and} \quad \{u(\cdot, t) \geq 0\} = [-x_1(t), x_1(t)].$$

We will see in step 3 that u is a solution of (19) in $[0, 1]$.

2. *Construction of solutions for $1 \leq t \leq 2$.* Consider now, for any measurable function $0 \leq \gamma(t) \leq 1$, the unique solution y_γ of the ODE

$$(22) \quad \dot{y}_\gamma(t) = c_1(t) + 2\gamma(t)y_\gamma(t) \quad \text{for } 1 \leq t \leq 2, \quad \text{and } y_\gamma(1) = 0.$$

By comparison, we have $0 \leq y_0(t) \leq y_\gamma(t) \leq y_1(t)$ for $1 \leq t \leq 2$, where y_0, y_1 are the solutions of (22) obtained with $\gamma(t) \equiv 0, 1$. In particular, it follows that $\dot{y}_\gamma \geq 0$ in $[1, 2]$. Consider

$$\begin{cases} \frac{\partial u_\gamma}{\partial t} = \dot{y}_\gamma(t) \left| \frac{\partial u_\gamma}{\partial x} \right| & \text{in } \mathbb{R} \times (1, 2], \\ u_\gamma(\cdot, 1) = u(\cdot, 1) & \text{in } \mathbb{R}, \end{cases}$$

where u is the solution of (20). Again, this problem has a unique continuous viscosity solution u_γ , and setting $\Gamma_\gamma(t) = y_\gamma(t) \geq 0$ for $1 \leq t \leq 2$, we obtain that v_γ defined by $v_\gamma(x, \Gamma_\gamma(t)) = u_\gamma(x, t)$ is the unique continuous viscosity solution of

$$\begin{cases} \frac{\partial v_\gamma}{\partial t} = \left| \frac{\partial v_\gamma}{\partial x} \right| & \text{in } \mathbb{R} \times (0, \Gamma_\gamma(2)], \\ v_\gamma(\cdot, 0) = u(\cdot, 1) & \text{in } \mathbb{R}. \end{cases}$$

Therefore, for all $(x, t) \in \mathbb{R} \times [1, 2]$, we have

$$u_\gamma(x, t) = \sup_{|x-y| \leq y_\gamma(t)} u(y, 1) = \begin{cases} 0 & \text{if } |x| \leq y_\gamma(t), \\ u(|x| - y_\gamma(t), 1) & \text{otherwise.} \end{cases}$$

(Note that $u(-x, t) = u(x, t)$ since u_0 is even and, since $u(\cdot, 1) \leq 0$, by the maximum principle, we have $u_\gamma \leq 0$ in $\mathbb{R} \times [1, 2]$.) It follows that, for all $1 \leq t \leq 2$,

$$(23) \quad \{u_\gamma(\cdot, t) > 0\} = \emptyset \quad \text{and} \quad \{u_\gamma(\cdot, t) \geq 0\} = \{u_\gamma(\cdot, t) = 0\} = [-y_\gamma(t), y_\gamma(t)].$$

3. *There are several weak solutions of (19).* Set, for $0 \leq \gamma(t) \leq 1$,

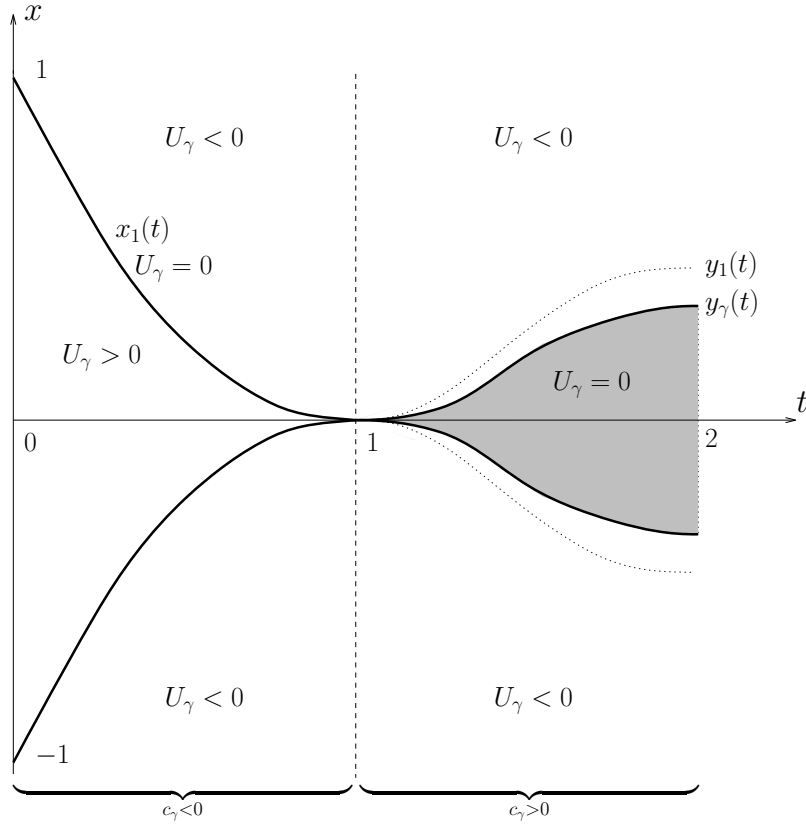
$$\begin{aligned} c_\gamma(t) &= c_1(t) + 2x_1(t), & U_\gamma(x, t) &= u(x, t) & \text{if } (x, t) \in \mathbb{R} \times [0, 1], \\ c_\gamma(t) &= c_1(t) + 2\gamma(t)y_\gamma(t), & U_\gamma(x, t) &= u_\gamma(x, t) & \text{if } (x, t) \in \mathbb{R} \times [1, 2]. \end{aligned}$$

Then, from steps 1 and 2, U_γ is the unique continuous viscosity solution of

$$(24) \quad \begin{cases} \frac{\partial U_\gamma}{\partial t} = c_\gamma(t) \left| \frac{\partial U_\gamma}{\partial x} \right| & \text{in } \mathbb{R} \times (0, 2], \\ U_\gamma(\cdot, 0) = u_0 & \text{in } \mathbb{R}. \end{cases}$$

Taking $\chi_\gamma(\cdot, t) = \gamma(t) \mathbf{1}_{[-y_\gamma(t), y_\gamma(t)]}$ for $1 \leq t \leq 2$, from (21) and (23), we have

$$\mathbf{1}_{\{U_\gamma(\cdot, t) > 0\}} \leq \chi_\gamma(\cdot, t) \leq \mathbf{1}_{\{U_\gamma(\cdot, t) \geq 0\}}$$

FIG. 1. Fattening phenomenon for the functions U_γ .

(see Figure 1). It follows that all the U_γ 's, for measurable $0 \leq \gamma(t) \leq 1$, are all weak solutions of (19), so we do not have uniqueness and the set of solutions is quite large.

Let us complete this counterexample by pointing out the following:

(i) As in [6], nonuniqueness comes from the fattening phenomenon for the front which is due to the fact that c_γ in (24) changes its sign at $t = 1$. It is even possible to build an autonomous counterexample up to start with a front with several connected components.

(ii) $c_0 \geq 0$, and therefore it also complements the results of section 5; indeed the unique solution u of (11) has the same 0-level-set as U_1 (obtained with $\gamma(t) \equiv 1$), and, with the notation of Theorem 1.5, $\rho^+ = \mathbf{1}_{\{U^1(\cdot, t) \geq 0\}}$ and $\rho^- = \mathbf{1}_{\{U^1(\cdot, t) > 0\}}$. In particular, for $t \geq 1$,

$$\{v^+(\cdot, t) \geq 0\} = \{U^1(\cdot, t) \geq 0\} = [-y_1(t), y_1(t)]$$

and

$$\{v^-(\cdot, t) > 0\} = \{U^1(\cdot, t) > 0\} = \emptyset.$$

Finally, we note that there are no strong solutions since (8) is obviously never satisfied.

(iii) $c_0 = 1$ does not satisfy (H1), but because of the finite speed of propagation property, it is possible to keep the same solution on a large ball in space and for $t \in (0, T)$ if we replace c_0 by a function with compact support in space such that

$c_0(x, t) = 1$ for $|x| \leq R$ with R large enough. In this way, it is possible for c_0 to satisfy (H1).

4. Uniqueness results for weak solutions of (3). Uniqueness of weak solutions of (3) is false in general, as shown in the counterexample of the previous section for sign changing velocities c_1 . This is in particular related to the “fattening phenomenon.” In [2] and [5] the authors proved that there is a unique “classical” viscosity solution for (3) under the assumptions that the initial set $\{u(\cdot, 0) \geq 0\}$ has the “interior sphere property” and that $c_1(\cdot, t) \geq |c_0(\cdot, t)|_{L^1}$ for any $t \geq 0$ —a condition which ensures that the velocity \bar{c} is nonnegative. By “classical” continuous viscosity solutions we mean that $t \rightarrow \mathbf{1}_{\{u(\cdot, t) \geq 0\}}$ is continuous in L^1 , which entails that $(x, t) \mapsto \bar{c}(x, t)$ is continuous, and that (3) holds in the usual viscosity sense.

Here we prove Theorem 1.3. If the condition $c_1(\cdot, t) \geq |c_0(\cdot, t)|_{L^1}$ is satisfied, then weak solutions are viscosity solutions. We also prove that the weak solution is unique if we suppose, moreover, either that the initial condition has the interior sphere property or that the strict inequality $c_1(\cdot, t) > |c_0(\cdot, t)|_{L^1}$ holds for any $t \geq 0$.

Proof of Theorem 1.3. 1. *Weak solutions are classical continuous viscosity solutions.* Let u be a weak solution, and let \bar{c} be associated with u as in Definition 1.1. Then, for any $x \in \mathbb{R}^N$ and for almost all $t \in [0, T]$, we have

$$\begin{aligned} \bar{c}(x, t) &\geq c_1(x, t) + c_0^+(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) > 0\}}(x) - c_0^-(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x) \\ &\geq c_1(x, t) - |c_0^-(\cdot, t)|_{L^1} \\ &\geq \delta \geq 0. \end{aligned}$$

From [22, Theorem 4.2], there exists a constant η which depends on T such that (10) implies

$$(25) \quad -|u| - |Du| \leq -\eta \quad \text{on } \mathbb{R}^N \times (0, T)$$

when we assume, moreover, that \bar{c} is continuous. In our case, where \bar{c} is not assumed continuous in time, (10) follows from the L^1 -stability result, Theorem A.3, where we approximate \bar{c} by a continuous function, and from the usual stability for L^1 -viscosity subsolutions.

Let us note that from the proof of [5, Corollary 2.5] we have in the viscosity sense

$$(26) \quad -|u(\cdot, t)| - |Du(\cdot, t)| \leq -\eta \quad \text{on } \mathbb{R}^N \text{ for every } t \in (0, T).$$

Following [5, Corollary 2.5], we get that, for every $t \in (0, T)$, the 0-level-set of $u(\cdot, t)$ has a zero-Lebesgue measure. Then we deduce that

$$\chi(x, t) = \mathbf{1}_{\{u(\cdot, t) \geq 0\}}(x) \quad \text{for a.e. } x \in \mathbb{R}^N \quad \text{and for all } t \in (0, T),$$

which (with (7)) entails that

$$\bar{c}(x, t) = c_1 + c_0 \star \mathbf{1}_{\{u(\cdot, t) \geq 0\}}$$

for any (x, t) . Moreover, $t \mapsto \mathbf{1}_{\{u(\cdot, t) \geq 0\}}$ is also continuous in L^1 , and then \bar{c} is continuous. Therefore u is a classical viscosity solution of (3).

2. *Uniqueness when u_0 is semiconvex (part (ii)).* If we assume that (9) and (10) hold and that u_0 is semiconvex, then weak solutions are viscosity solutions, and we can apply the uniqueness result for viscosity solutions given in [5], namely Theorem 4.2

(which remains true under our assumptions), which requires, in particular, semiconvexity of the velocity; see assumption (H2).

3. *A Gronwall-type inequality (part (i)).* From now on we assume that $\delta > 0$, and we aim to prove that the solution to (3) is unique. Let u_1, u_2 be two solutions. We set

$$\rho_i = \mathbf{1}_{\{u_i(\cdot, t) \geq 0\}} \quad \text{and} \quad \bar{c}_i(x, t) = c_0 \star \rho_i + c_1 \quad \text{for } i = 1, 2.$$

We want to prove in a first step the following Gronwall-type inequality for any t sufficiently small:

$$(27) \quad \begin{aligned} & |\rho_1(\cdot, t) - \rho_2(\cdot, t)|_{L^1} \\ & \leq C [\text{per}(\{u_1(\cdot, t) \geq 0\}) + \text{per}(\{u_2(\cdot, t) \geq 0\})] \int_0^t |\rho_1(\cdot, s) - \rho_2(\cdot, s)|_{L^1} ds, \end{aligned}$$

where C is a constant depending on the constants of the problem, where $\text{per}(\{u_i(\cdot, t) \geq 0\})$ is the \mathcal{H}^{N-1} measure of the set $\partial\{u_i(\cdot, t) \geq 0\}$ (for $i = 1, 2$).

We have

$$(28) \quad |\rho_1(\cdot, t) - \rho_2(\cdot, t)|_{L^1} \leq \mathcal{L}^N(\{-\alpha_t \leq u_1(\cdot, t) < 0\}) + \mathcal{L}^N(\{-\alpha_t \leq u_2(\cdot, t) < 0\}),$$

where \mathcal{L}^N is the Lebesgue measure in \mathbb{R}^N and

$$(29) \quad \alpha_t = \sup_{s \in [0, t]} |(u_1 - u_2)(\cdot, s)|_\infty \quad \text{for any } t \in (0, T).$$

In order to estimate the right-hand side of inequality (28), as in the proof of [5, Theorem 4.2], we need a lower-gradient bound as well as a semiconvexity property for u_1 and u_2 . We already know from step 1 that \bar{c}_i is continuous for $i = 1, 2$.

Let us start to estimate the right-hand side of (28). From the ‘‘stability estimates’’ on the solutions with respect to variations of the velocity (see [5, Lemma 2.2]), we have

$$(30) \quad \alpha_t \leq |Du_0|_\infty e^{Lt} \int_0^t |(\bar{c}_1 - \bar{c}_2)(\cdot, s)|_\infty ds,$$

where $L = L_0 + L_1$. Therefore

$$(31) \quad \alpha_t \leq m_0 |Du_0|_\infty e^{Lt} \int_0^t |(\rho_1 - \rho_2)(\cdot, s)|_{L^1} ds.$$

where the constant m_0 is given in (H2). In particular, since the $\rho_i(\cdot, t)$ are continuous in L^1 and equal at time $t = 0$ for $i = 1, 2$, we have $\alpha_t/t \rightarrow 0$ as $t \rightarrow 0^+$.

From now on, we mimic the proof of [5, Proposition 4.5]. Using the lower-gradient bound (25) for $u_i(\cdot, t)$ combined with the increase principle (see [5, Lemma 2.3]), we obtain for $\alpha_t < \eta/2$ that

$$\{-\alpha_t \leq u_i(\cdot, t) < 0\} \subset \{u_i(\cdot, t) \geq 0\} + (2\alpha_t/\eta)\bar{B}(0, 1)$$

for $i = 1, 2$. From the interior ball regularization lemma (Lemma 2.3), the set $\{u_i(\cdot, t) \geq 0\}$ satisfies for $t \in (0, T)$ the interior ball property of radius $r_t = \eta t/C_0$.

Applying [2, Lemmas 2.5 and 2.6], we obtain for $\sigma_t = 2\alpha_t/\eta$ that

$$\begin{aligned} \mathcal{L}^N(\{-\alpha_t \leq u_i(\cdot, t) < 0\}) &\leq \mathcal{L}^N((\{u_i(\cdot, t) \geq 0\} + \sigma_t \bar{B}(0, 1)) \setminus \{u_i(\cdot, t) \geq 0\}) \\ &\leq \frac{r_t}{N} \left[\left(1 + \frac{\sigma_t}{r_t}\right)^N - 1 \right] \text{per}(\{u_i(\cdot, t) \geq 0\}) \\ &\leq \frac{2^N \alpha_t}{\eta} \text{per}(\{u_i(\cdot, t) \geq 0\}) \end{aligned}$$

(using $(1+a)^N - 1 \leq aN(1+a)^{N-1}$ for $a \geq 0$) for $t \in [0, \tau]$, where $0 < \tau \leq T$ is defined by

$$(32) \quad \tau = \sup \left\{ t > 0 : \alpha_t < \frac{\eta}{2} \text{ and } 2 \frac{C_0}{\eta^2} \frac{\alpha_t}{t} \leq 1 \right\}.$$

Putting together (31), (28), and the previous inequality proves (27).

4. *Uniqueness when $\delta > 0$ (part (i)).* We now complete the uniqueness proof under the assumption $\delta > 0$. For this we first show that $\rho_1 = \rho_2$ in $[0, \tau]$. In order to apply the Gronwall lemma to the L^1 -estimate (27) obtained in step 3, it is enough to prove that the functions $t \mapsto \text{per}(\{u_i(\cdot, t) \geq 0\})$ belong to L^1 . For this, let us set

$$w_i(x) = \inf \{ t \geq 0 : u_i(x, t) \geq 0 \}.$$

Since u_i solves the eikonal equation $(u_i)_t = \bar{c}_i(x, t)|Du_i|$, from classical representation formulae, we have

$$\begin{aligned} \{u_i(\cdot, t) \geq 0\} &= \{x : \exists y(\cdot), |\dot{y}(s)| \leq c(y(s), s), 0 \leq s \leq t, \\ &\quad u_0(y(0)) \geq 0 \text{ and } y(t) = x\}. \end{aligned}$$

Therefore

$$\begin{aligned} w_i(x) &= \inf \{ t \geq 0 : \exists y(\cdot), |\dot{y}(s)| \leq c(y(s), s), 0 \leq s \leq t, \\ &\quad u_0(y(0)) \geq 0 \text{ and } y(t) = x\}. \end{aligned}$$

Applying the dynamic programming principle, since $\bar{c}_i \geq \delta > 0$, we obtain that w_i is Lipschitz continuous and is a viscosity solution of the autonomous equation $\bar{c}_i(x, w_i(x))|Dw_i(x)| = 1$. Note that $\{u_i(\cdot, t) \geq 0\} = \{w_i \leq t\}$. In particular, by Theorem 2.1(i), $\{w_i \leq t\} \subset B(0, R_0 + (M_0 + M_1)t)$ is bounded for any t . From the coarea formula, we have

$$\begin{aligned} \int_0^t \text{per}(\{u_i(\cdot, s) \geq 0\}) ds &= \int_0^t \text{per}(\{w_i \leq s\}) ds \\ &= \int_{\{w_i \leq t\}} |Dw_i(x)| dx, \end{aligned}$$

which is finite since w_i is Lipschitz continuous. Therefore we have proved that $t \mapsto \text{per}(\{u_i(\cdot, t) \geq 0\})$ belongs to $L^1([0, \tau])$, which entails from the Gronwall lemma that $\rho_1 = \rho_2$ in $[0, \tau]$ since $\rho_1(\cdot, 0) = \rho_2(\cdot, 0)$. Hence $\bar{c}_1 = \bar{c}_2$ and $u_1 = u_2$ in $[0, \tau]$. From the definition of α_t and τ (see (29) and (32)), necessarily $\tau = T$. It completes the proof. \square

5. Nonnegative kernel c_0 and Slepčev formulation for the nonlocal term. In this section, we deal with nonnegative kernels $c_0 \geq 0$. In this monotone framework, inclusion principles for evolving sets and comparison for solutions to the dislocation equation are expected (see Cardaliaguet [11] for related results). We start by studying the right level-set equation using a Slepčev formulation with the convolution term using all the level-sets $\{u(\cdot, t) \geq u(x, t)\}$ instead of only one level-set $\{u(\cdot, t) \geq 0\}$. This choice is motivated by the good stability properties of the Slepčev formulation.

The equation we are concerned with is

$$(33) \quad \begin{cases} \frac{\partial u}{\partial t} = c^+[u](x, t)|Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where the nonlocal velocity is

$$(34) \quad \begin{aligned} c^+[u](x, t) &= c_1(x, t) + c_0(\cdot, t) \star \mathbf{1}_{\{u(\cdot, t) \geq u(x, t)\}}(x) \\ &= c_1(x, t) + \int_{\mathbb{R}^N} c_0(x - z, t) \mathbf{1}_{\{u(\cdot, t) \geq u(x, t)\}}(z) dz \end{aligned}$$

and the additional velocity c_1 has no particular sign.

We denote

$$c^-[u](x, t) = c_1(x, t) + \int_{\mathbb{R}^N} c_0(x - z, t) \mathbf{1}_{\{u(\cdot, t) > u(x, t)\}}(z) dz.$$

We recall the notion of viscosity solutions for (33) as it appears in [28].

DEFINITION 5.1 (Slepčev viscosity solutions). *An upper semicontinuous function $u : \mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}$ is a viscosity subsolution of (33) if, for any $\varphi \in C^1(\mathbb{R}^N \times [0, T])$, for any maximum point (\bar{x}, \bar{t}) of $u - \varphi$, if $\bar{t} > 0$, then*

$$\frac{\partial \varphi}{\partial t}(\bar{x}, \bar{t}) \leq c^+[u](\bar{x}, \bar{t})|D\varphi(\bar{x}, \bar{t})|$$

and $u(\bar{x}, 0) \leq u_0(\bar{x})$ if $\bar{t} = 0$.

A lower semicontinuous function $u : \mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}$ is a viscosity supersolution of (33) if, for any $\varphi \in C^1(\mathbb{R}^N \times [0, T])$, for any minimum point (\bar{x}, \bar{t}) of $u - \varphi$, if $\bar{t} > 0$, then

$$\frac{\partial \varphi}{\partial t}(\bar{x}, \bar{t}) \geq c^-[u](\bar{x}, \bar{t})|D\varphi(\bar{x}, \bar{t})|$$

and $u(\bar{x}, 0) \geq u_0(\bar{x})$ if $\bar{t} = 0$.

A locally bounded function is a viscosity solution of (33) if its upper semicontinuous envelope is a subsolution and its lower semicontinuous envelope is a supersolution of (33).

Note that for the supersolution, we require the viscosity inequality with c^- instead of c^+ . It is the definition providing the expected stability results (see [28]).

THEOREM 5.2 (comparison principle). *Assume (H0') and that the kernels $c_0 \geq 0$ and c_1 satisfy (H1). Let u (respectively, v) be a bounded upper semicontinuous subsolution (respectively, a bounded lower semicontinuous supersolution) of (33). Then $u \leq v$ in $\mathbb{R}^N \times [0, T]$.*

Remark 5.1. We could deal with second-order terms in (33) (for instance we can add the mean curvature to the velocity (1)). See Forcadel [17] and Srour [29] for related results.

Before giving the proof of Theorem 5.2, let us note the following consequence.

Proof of Theorem 1.4. The uniqueness of a continuous viscosity solution to (33) is an immediate consequence of Theorem 5.2. Then existence is proved by Perron's method using classical arguments (see, for instance, [16, Theorem 1.2]), so we skip the details. \square

Proof of Theorem 5.2. 1. *The test function.* Since $u - v$ is a bounded upper semicontinuous function, for any $\varepsilon, \eta, \alpha > 0$ and $K = 2(L_0 + L_1) \geq 0$, the supremum

$$M_{\varepsilon, \eta, \alpha} = \sup_{(x, y, t) \in (\mathbb{R}^N)^2 \times [0, T]} \left\{ u(x, t) - v(y, t) - e^{Kt} \left(\frac{|x - y|^2}{\varepsilon^2} + \alpha|x|^2 + \alpha|y|^2 \right) - \eta t \right\}$$

is finite and achieved at a point $(\bar{x}, \bar{y}, \bar{t})$. Classical arguments show that

$$\liminf_{\varepsilon, \eta, \alpha \rightarrow 0} M_{\varepsilon, \eta, \alpha} = \sup_{\mathbb{R}^N \times [0, T]} \{u - v\}$$

and that

$$(35) \quad \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2}, \alpha|\bar{x}|^2, \alpha|\bar{y}|^2 \leq M_\infty,$$

where $M_\infty = |u|_\infty + |v|_\infty$.

2. *Viscosity inequalities when $\bar{t} > 0$.* Writing the viscosity inequalities for the subsolution u and the supersolution v , we obtain

$$(36) \quad Ke^{K\bar{t}} \left(\frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} + \alpha|\bar{x}|^2 + \alpha|\bar{y}|^2 \right) + \eta \leq c^+ [u(\bar{x}, \bar{t})|\bar{p} + \bar{q}_x| - c^- [v(\bar{y}, \bar{t})|\bar{p} - \bar{q}_y|],$$

where $\bar{p} = 2e^{K\bar{t}}(\bar{x} - \bar{y})/\varepsilon^2$, $\bar{q}_x = 2e^{K\bar{t}}\alpha\bar{x}$, and $\bar{q}_y = 2e^{K\bar{t}}\alpha\bar{y}$. We point out a difficulty in obtaining this inequality: in general, one gets it by doubling the time variable first and then by passing to the limit in the time penalization. This is not straightforward here because of the dependence with respect to time of the nonlocal terms. But the stability arguments of the Slepčev formulation take care of this difficulty.

3. *Difference between $\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}$ and $\{v(\cdot, \bar{t}) > v(\bar{y}, \bar{t})\}$.* We have

$$(37) \quad \{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\} \subset \{v(\cdot, \bar{t}) > v(\bar{y}, \bar{t})\} \cup \mathcal{E},$$

where $\mathcal{E} = \{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\} \cap \{v(\cdot, \bar{t}) \leq v(\bar{y}, \bar{t})\}$. If $x \in \mathcal{E}$, then $u(\bar{x}, \bar{t}) - v(\bar{y}, \bar{t}) \leq u(x, \bar{t}) - v(x, \bar{t})$. But from the definition of $M_{\varepsilon, \eta, \alpha}$,

$$\begin{aligned} & u(x, \bar{t}) - v(x, \bar{t}) - e^{K\bar{t}}2\alpha|x|^2 - \eta\bar{t} \\ & \leq u(\bar{x}, \bar{t}) - v(\bar{y}, \bar{t}) - e^{K\bar{t}} \left(\frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} + \alpha|\bar{x}|^2 + \alpha|\bar{y}|^2 \right) - \eta\bar{t}. \end{aligned}$$

It follows that

$$\mathcal{E} \subset \left\{ x \in \mathbb{R}^N : |x|^2 \geq \frac{1}{2}(|\bar{x}|^2 + |\bar{y}|^2) + \frac{|\bar{x} - \bar{y}|^2}{2\alpha\varepsilon^2} \right\}.$$

4. *Upper bound for $c^+[u](\bar{x}, \bar{t})$.* We have

$$(38) \quad \begin{aligned} c^+[u](\bar{x}, \bar{t}) &= \int_{\mathbb{R}^N} c_0(\bar{x} - z, \bar{t}) \mathbf{1}_{\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}}(z) dz + c_1(\bar{x}, \bar{t}) \\ &\leq \int_{\mathbb{R}^N} (c_0(\bar{x} - z, \bar{t}) - c_0(\bar{y} - z, \bar{t})) \mathbf{1}_{\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}}(z) dz \\ &\quad + \int_{\mathbb{R}^N} c_0(\bar{y} - z, \bar{t}) \mathbf{1}_{\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}}(z) dz + c_1(\bar{x}, \bar{t}). \end{aligned}$$

Using that $c_0 \geq 0$ and (37), we obtain

$$\int_{\mathbb{R}^N} c_0(\bar{y} - z, \bar{t}) \mathbf{1}_{\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}}(z) dz \leq \int_{\{v(\cdot, \bar{t}) > v(\bar{y}, \bar{t})\} \cup \mathcal{E}} c_0(\bar{y} - z, \bar{t}) dz.$$

From (38), we get

$$(39) \quad c^+[u](\bar{x}, \bar{t}) \leq c^-[v](\bar{y}, \bar{t}) + \mathcal{I}_1 + \mathcal{I}_2 + c_1(\bar{x}, \bar{t}) - c_1(\bar{y}, \bar{t}),$$

where

$$\mathcal{I}_1 = \int_{\mathbb{R}^N} (c_0(\bar{x} - z, \bar{t}) - c_0(\bar{y} - z, \bar{t})) \mathbf{1}_{\{u(\cdot, \bar{t}) \geq u(\bar{x}, \bar{t})\}}(z) dz$$

and

$$\mathcal{I}_2 = \int_{\mathcal{E}} c_0(\bar{y} - z, \bar{t}) dz.$$

5. *Estimate of \mathcal{I}_1 using (H1).* We have

$$c_0(\bar{x} - z, \bar{t}) - c_0(\bar{y} - z, \bar{t}) = \int_0^1 D_x c_0((1 - \lambda)(\bar{y} - z) + \lambda(\bar{x} - z), \bar{t})(\bar{x} - \bar{y}) d\lambda.$$

It follows that

$$(40) \quad \begin{aligned} |\mathcal{I}_1| &\leq \int_{\mathbb{R}^N} \int_0^1 |D_x c_0((1 - \lambda)(\bar{y} - z) + \lambda(\bar{x} - z), \bar{t})| |\bar{x} - \bar{y}| d\lambda dz \\ &\leq |D_x c_0(\cdot, \bar{t})|_{L^1} |\bar{x} - \bar{y}| \\ &\leq L_0 |\bar{x} - \bar{y}|. \end{aligned}$$

6. *Estimate of the right-hand side of inequality (36).* Noticing that $|\bar{p}||\bar{x} - \bar{y}| = 2e^{K\bar{t}}|\bar{x} - \bar{y}|^2/\varepsilon^2$ and using (40), (39), and (H1), we have

$$\begin{aligned} &c^+[u](\bar{x}, \bar{t})|\bar{p} + \bar{q}_x| - c^-[v](\bar{y}, \bar{t})|\bar{p} - \bar{q}_y| \\ &\leq (c^-[v](\bar{y}, \bar{t}) + \mathcal{I}_1 + \mathcal{I}_2 + c_1(\bar{x}, \bar{t}) - c_1(\bar{y}, \bar{t}))|\bar{p} + \bar{q}_x| - c^-[v](\bar{y}, \bar{t})|\bar{p} - \bar{q}_y| \\ &\leq |c^-[v](\bar{y}, \bar{t})| |\bar{q}_x + \bar{q}_y| + (L_0 + L_1)|\bar{x} - \bar{y}||\bar{p} + \bar{q}_x| + \mathcal{I}_2|\bar{p} + \bar{q}_x| \\ &\leq (M_0 + M_1)(|\bar{q}_x| + |\bar{q}_y|) + 2e^{K\bar{t}}(L_0 + L_1) \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} \\ &\quad + (L_0 + L_1)|\bar{x} - \bar{y}||\bar{q}_x| + \mathcal{I}_2 \left(|\bar{q}_x| + 2e^{K\bar{t}} \frac{|\bar{x} - \bar{y}|}{\varepsilon^2} \right). \end{aligned}$$

Since $|\bar{q}_x|, |\bar{q}_y| \rightarrow 0$ as $\alpha \rightarrow 0$ (see (35)) and \mathcal{I}_2 is bounded by $|c_0(\cdot, \bar{t})|_{L^1} \leq L_0$, there exists a modulus $m_\varepsilon(\alpha) \rightarrow 0$ as $\alpha \rightarrow 0$ such that (36) becomes

$$\begin{aligned} & Ke^{K\bar{t}} \left(\frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} + \alpha|\bar{x}|^2 + \alpha|\bar{y}|^2 \right) + \eta \\ & \leq m_\varepsilon(\alpha) + 2(L_0 + L_1)e^{K\bar{t}} \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} + 2\mathcal{I}_2 e^{K\bar{t}} \frac{|\bar{x} - \bar{y}|}{\varepsilon^2}. \end{aligned}$$

Recalling that we chose $K \geq 2(L_0 + L_1)$, we finally obtain

$$(41) \quad 0 < \eta \leq m_\varepsilon(\alpha) + 2\mathcal{I}_2 e^{K\bar{t}} \frac{|\bar{x} - \bar{y}|}{\varepsilon^2}.$$

7. *Limit when $\alpha \rightarrow 0$.* First, suppose that

$$(42) \quad \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} \rightarrow 0 \quad \text{as } \alpha \rightarrow 0.$$

It follows that $|\bar{x} - \bar{y}| \rightarrow 0$ as $\alpha \rightarrow 0$. Passing to the limit in (41), we obtain a contradiction. Therefore, (42) cannot hold, and, up to extracting a subsequence, there exists $\delta > 0$ such that

$$(43) \quad \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} \geq \delta > 0 \quad \text{for } \alpha > 0 \text{ small enough.}$$

From (41) and (35), we get

$$(44) \quad \eta \leq \limsup_{\alpha \rightarrow 0} 2\mathcal{I}_2 e^{K\bar{t}} \frac{|\bar{x} - \bar{y}|}{\varepsilon^2} \leq \frac{2e^{K\bar{t}} M_\infty^{1/2}}{\varepsilon} \limsup_{\alpha \rightarrow 0} \mathcal{I}_2.$$

To obtain a contradiction, it suffices to show that $\limsup_{\alpha \rightarrow 0} \mathcal{I}_2 = 0$.

8. *Convergence of \mathcal{I}_2 to 0 when $\alpha \rightarrow 0$.* By a change of variable, we have

$$\mathcal{I}_2 = \int_{\bar{\mathcal{E}}} c_0(\bar{y} - z, \bar{t}) dz \leq \int_{\bar{\mathcal{E}}} c_0(z, \bar{t}) dz,$$

where

$$\bar{\mathcal{E}} = \left\{ x \in \mathbb{R}^N : |x - \bar{y}|^2 \geq \frac{1}{2}(|\bar{x}|^2 + |\bar{y}|^2) + \frac{|\bar{x} - \bar{y}|^2}{2\alpha\varepsilon^2} \right\}.$$

Since $|c_0(\cdot, \bar{t})|_{L^1} \leq L_0$, to prove that $\mathcal{I}_2 \rightarrow 0$, it suffices to show that $\bar{\mathcal{E}} \subset \mathbb{R}^N \setminus B(0, R_\alpha)$ with $R_\alpha \rightarrow +\infty$. From (43), if $x \in \bar{\mathcal{E}}$, then

$$\begin{aligned} |x|^2 & \geq -2|x||\bar{y}| + \frac{1}{2}|\bar{x}|^2 - \frac{1}{2}|\bar{y}|^2 + \frac{\delta}{2\alpha} \\ & \geq -2|x||\bar{y}| - |\bar{y}||\bar{x} - \bar{y}| + \frac{\delta}{2\alpha} \\ & \geq \frac{1}{2\alpha}(\delta - 2(C^2 + 2C|x|)\sqrt{\alpha}), \end{aligned}$$

since by (35), there exists $C > 0$ such that $|\bar{x}|, |\bar{y}| \leq C/\sqrt{\alpha}$ and $|\bar{x} - \bar{y}| \leq C$. It follows that

$$|x| \geq \frac{1}{\sqrt{\alpha}} \left(-C + \sqrt{C^2 + \delta/2 - C^2\sqrt{\alpha}} \right) := R_\alpha \xrightarrow{\alpha \rightarrow 0} +\infty.$$

9. *End of the proof.* Finally, for every ε , if $\alpha = \alpha_\varepsilon$ is small enough, the supremum $M_{\varepsilon, \eta, \alpha}$ is necessarily achieved for $\bar{t} = 0$. It follows that

$$M_{\varepsilon, \eta, \alpha} \leq u(\bar{x}, 0) - v(\bar{y}, 0) - \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} \leq u_0(\bar{x}, 0) - u_0(\bar{y}, 0) - \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2}.$$

Since u_0 is uniformly continuous, for all $\rho > 0$, there exists $C_\rho > 0$ such that

$$M_{\varepsilon, \eta, \alpha} \leq \rho + C_\rho |\bar{x} - \bar{y}| - \frac{|\bar{x} - \bar{y}|^2}{\varepsilon^2} \leq \rho + \frac{C_\rho^2 \varepsilon^2}{4}.$$

Passing to the limits $\varepsilon \rightarrow 0$ and then $\rho, \alpha, \eta \rightarrow 0$, we obtain that $\sup\{u - v\} \leq 0$. \square

Now we turn to the connections with *discontinuous solutions* and weak solutions, which are closely connected. To do so, if u is the unique continuous solution of (33) given by Theorem 1.4, we recall that we use the notation

$$\rho^+ := \mathbf{1}_{\{u \geq 0\}}, \quad \rho^- := \mathbf{1}_{\{u > 0\}}, \quad \text{and} \quad c[\rho](x, t) = c_0(\cdot, t) \star \rho(\cdot, t)(x) + c_1(x, t).$$

Proof of Theorem 1.5. 1. *Claim:* Under the assumptions of Theorem 1.5, the functions ρ^+ and ρ^- are L^1 -viscosity solutions of the equation

$$(45) \quad \begin{cases} \rho_t = c[\rho]|D\rho| & \text{in } \mathbb{R}^N \times [0, T], \\ \rho(x, 0) = \mathbf{1}_{\{u_0 \geq 0\}} & \text{in } \mathbb{R}^N. \end{cases}$$

We consider two sequences of smooth nondecreasing functions $(\psi_\alpha)_\alpha, (\psi^\alpha)_\alpha$, taking values in $[0, 1]$, such that, for any $s \in \mathbb{R}$,

$$\psi_\alpha(s) \leq \mathbf{1}_{\{s > 0\}}(s) \leq \mathbf{1}_{\{s \geq 0\}}(s) \leq \psi^\alpha(s),$$

and such that, as $\alpha \rightarrow 0$, $\psi_\alpha \uparrow \mathbf{1}_{\{s > 0\}}$, $\psi^\alpha \downarrow \mathbf{1}_{\{s \geq 0\}}$.

We first remark that u satisfies, in the sense of Definition 5.1,

$$\begin{aligned} u_t &\leq c[\psi^\alpha(u(\cdot, t) - u(x, t))]|Du| & \text{in } \mathbb{R}^N \times (0, T), \\ u_t &\geq c[\psi_\alpha(u(\cdot, t) - u(x, t))]|Du| & \text{in } \mathbb{R}^N \times (0, T) \end{aligned}$$

since u is a continuous solution of (33), $c^+[u](x, t) \leq c[\psi^\alpha(u(\cdot, t) - u(x, t))]$, and $c^-[u](x, t) \geq c[\psi_\alpha(u(\cdot, t) - u(x, t))]$. The point for doing that is that the functions $c[\psi^\alpha(u(\cdot, t) - u(x, t))]$, $c[\psi_\alpha(u(\cdot, t) - u(x, t))]$ are now continuous in x and t .

Then we show that ρ^+, ρ^- satisfy the same inequalities, the functions $c[\psi^\alpha(u(\cdot, t) - u(x, t))]$, $c[\psi_\alpha(u(\cdot, t) - u(x, t))]$ being considered as fixed functions (in other words, we forget that they depend on u). In fact, we just provide the proof in detail for ρ^+ , the one for ρ^- being analogous. Following the proof of [6], we set

$$u_\varepsilon(x, t) := \frac{1}{2} \left(1 + \tanh \left(\varepsilon^{-1} (u(x, t) + \varepsilon^{1/2}) \right) \right).$$

Noticing that $u_\varepsilon = \phi_\varepsilon(u)$ for an increasing function ϕ_ε , we have that the function u_ε still satisfies the two above inequalities. It is easy to see that

$$\rho^+ = \limsup^* u_\varepsilon \quad \text{and} \quad (\rho^+)_* = \liminf_* u_\varepsilon,$$

and the half-relaxed limits method indeed shows that

$$\begin{aligned} (\rho^+)_t^* &\leq c[\psi^\alpha(u(\cdot, t) - u(x, t))] |D(\rho^+)^*| \quad \text{in } \mathbb{R}^N \times (0, T), \\ ((\rho^+)_*)_t &\geq c[\psi_\alpha(u(\cdot, t) - u(x, t))] |D(\rho^+)_*| \quad \text{in } \mathbb{R}^N \times (0, T). \end{aligned}$$

The next step consists of remarking that the viscosity sub- and supersolution inequalities for ρ^+ are obviously satisfied in the complementary of $\partial\{u \geq 0\}$ since ρ^+ is locally constant there, and therefore it is a classical solution of the problem. The only nontrivial viscosity sub- and supersolution inequalities we have to check are at points $(x, t) \in \partial\{u \geq 0\}$, i.e., such that $u(x, t) = 0$ since u is continuous. For such points, as $\alpha \rightarrow 0$,

$$c[\psi^\alpha(u(\cdot, t) - u(x, t))] \rightarrow c[\rho^+](x, t) = c[(\rho^+)^*](x, t)$$

since ρ^+ is upper semicontinuous, and

$$c[\psi_\alpha(u(\cdot, t) - u(x, t))] \rightarrow c[\rho^-](x, t).$$

The stability result for equations with an L^1 -dependence in time yields the inequalities

$$(46) \quad \begin{aligned} (\rho^+)_t^* &\leq c[(\rho^+)^*] |D(\rho^+)^*| \quad \text{in } \mathbb{R}^N \times (0, T), \\ ((\rho^+)_*)_t &\geq c[\rho^-] |D(\rho^+)_*| \quad \text{in } \mathbb{R}^N \times (0, T). \end{aligned}$$

The second inequality is weaker than the one we claim: to obtain $c[(\rho^+)_*]$ instead of $c[\rho^-]$, we have to play with the different level-sets of u : for $\beta > 0$ small, we set $\rho_\beta^+ = \mathbf{1}_{\{u \geq -\beta\}}$. Since u is a solution of (33) and $\psi_{\alpha, \beta} := \psi_\alpha(\cdot + \beta)$ is nondecreasing, then $\psi_{\alpha, \beta}(u)$ is a (continuous) supersolution of

$$(\psi_{\alpha, \beta}(u))_t \geq c^-[\psi_{\alpha, \beta}(u)] |D\psi_{\alpha, \beta}(u)| \quad \text{in } \mathbb{R}^N \times (0, T).$$

By stability we get, as $\alpha \rightarrow 0$,

$$((\rho_\beta^+)_*)_t \geq c^-[(\rho_\beta^+)_*] |D(\rho_\beta^+)_*| \quad \text{in } \mathbb{R}^N \times (0, T).$$

But, for $(x, t) \in \partial\{u \geq -\beta\}$,

$$(47) \quad c^-[(\rho_\beta^+)_*](x, t) = c[\mathbf{1}_{\{(\rho_\beta^+)_*(\cdot, t) > 0\}}](x, t) = c[(\rho_\beta^+)_*](x, t) \geq c[\rho^+](x, t).$$

It follows that $(\rho_\beta^+)_*$ is a supersolution of the eikonal equation with $c[\rho^+](x, t)$ (as before, the only nontrivial inequalities we have to check are on $\partial\{u \geq -\beta\}$, and they are true because of (47)). Letting β tend to 0 and using that $(\rho^+)_* = \liminf_* (\rho_\beta^+)_*$, we obtain the expected inequality (even something better since (46) holds actually with $c[\rho^+](x, t)$). In particular, we get that ρ^+ is a solution of

$$(48) \quad \begin{cases} \rho_t = c[\rho^+] |D\rho| & \text{in } \mathbb{R}^N \times [0, T), \\ \rho(x, 0) = \mathbf{1}_{\{u_0 \geq 0\}} & \text{in } \mathbb{R}^N. \end{cases}$$

The proof of the claim is complete.

2. *The functions v^\pm are weak solutions of (3).* Let us start with the “+” case. We first remark that the existence and uniqueness of v^+ follows from the standard theory for equations with an L^1 -dependence in time (see Appendix A).

To prove that v^+ is a weak solution of (3), it remains to prove that (7) holds. It is sufficient to show that

$$(49) \quad \{v^+(\cdot, t) > 0\} \subset \{u(\cdot, t) \geq 0\} \subset \{v^+(\cdot, t) \geq 0\}.$$

We use again the functions ψ_α, ψ^α introduced above. We remark that

$$\psi_\alpha(u_0) \leq \rho^+(x, 0) \leq \psi^\alpha(u_0) \quad \text{in } \mathbb{R}^N.$$

Moreover, v^+ and ρ^+ are solutions of the same equation, namely (48) with $c[\rho^+]$, which is considered as a fixed function, and so are $\psi_\alpha(v^+)$ and $\psi^\alpha(v^+)$ because the equation is geometric. Therefore, a standard comparison result implies

$$\psi_\alpha(v^+) \leq (\rho^+)_* \leq \rho^+ = (\rho^+)^* \leq \psi^\alpha(v^+) \quad \text{in } \mathbb{R}^N \times [0, T].$$

Letting α tend to 0, these inequalities imply (49).

We can prove the symmetric result for v^- , the only difference being that inclusion (49) has to be replaced by

$$(50) \quad \{v^-(\cdot, t) > 0\} \subset \{u(\cdot, t) > 0\} \subset \{v^-(\cdot, t) \geq 0\}.$$

3. *Claim: If v is a weak solution of (3), then $\mathbf{1}_{\{v(\cdot, t) \geq 0\}}$ is an L^1 -subsolution of (45).* From Proposition 3.1 and since $c_0 \geq 0$, v satisfies in the L^1 -sense

$$(51) \quad v_t \leq c[\mathbf{1}_{\{v(\cdot, t) \geq 0\}}] |Dv| \quad \text{in } \mathbb{R}^N \times [0, T].$$

By similar arguments as we used above, the function $\mathbf{1}_{\{v(\cdot, t) \geq 0\}}$ satisfies the same inequality which gives the result.

4. *The function ρ^+ is the maximal L^1 -subsolution of (45).* Let w be an L^1 - (upper semicontinuous) subsolution of (45). First we have $w \leq 1$ in $\mathbb{R}^N \times [0, T)$ by comparison with the constant supersolution 1 for the equation with $c[w]$ fixed. By considering $\max(w, 0)$ we may assume that $0 \leq w \leq 1$ in $\mathbb{R}^N \times [0, T)$. By similar arguments as we already used in step 1, we can show that $\mathbf{1}_{\{w(\cdot, t) > 0\}}$ is also an L^1 -subsolution of (45); thus we can assume that w is a characteristic function.

Then we remark that w is also a subsolution of (33): indeed, again, the only nontrivial viscosity inequalities are on the boundary of the set $\{w = 1\}$, and if (x, t) is such a point, we have $w(x, t) = 1$ because w is upper semicontinuous and $w = \mathbf{1}_{\{w(\cdot, t) \geq w(x, t)\}}$. Since u is a solution of the geometric equation (33), $\psi^\alpha(u)$ is still a solution which satisfies $\psi^\alpha(u)(x, 0) \geq \mathbf{1}_{\{u_0 \geq 0\}} \geq w(x, 0)$ in \mathbb{R}^N . By Theorem 5.2 we obtain

$$w \leq \psi^\alpha(u) \quad \text{in } \mathbb{R}^N \times [0, T).$$

Letting α tend to 0 provides $w \leq \rho^+$, which proves that ρ^+ is the maximal subsolution of (45).

5. *The function v^+ is the maximal weak solution of (3).* Let v be a weak solution of (3). From steps 3 and 4, we get $\mathbf{1}_{\{v(\cdot, t) \geq 0\}} \leq \rho^+(\cdot, t)$ in $\mathbb{R}^N \times [0, T)$, and (51) implies

$$v_t \leq c[\rho^+] |Dv| \quad \text{in } \mathbb{R}^N \times [0, T).$$

Therefore v is a subsolution of (12), and by a standard comparison result, this leads to $v \leq v^+$ in $\mathbb{R}^N \times [0, T)$, which proves the result.

6. We have $\{v^+(\cdot, t) \geq 0\} = \{u(\cdot, t) \geq 0\}$ and $\{v^-(\cdot, t) > 0\} = \{u(\cdot, t) > 0\}$. From step 5, we get $\mathbf{1}_{\{v^+(\cdot, t) \geq 0\}} \leq \rho^+(\cdot, t) = \mathbf{1}_{\{u(\cdot, t) \geq 0\}}$. The conclusion follows for v^+ using (49). The inclusion for v^- uses symmetric arguments.

7. *Uniqueness when $\{u(\cdot, t) = 0\}$ has Lebesgue measure 0.* If $\mathcal{L}^N(\{u(\cdot, t) = 0\}) = 0$, then $c[\rho^+] = c[\rho^-]$. Hence $v^+ = v^-$ is the unique weak solution of (3), and it is obviously a classical one. \square

Appendix A. A stability result for eikonal equations with L^1 -dependence in time. The aim of this appendix is to provide a self-contained presentation of a stability result for viscosity solutions of eikonal equations with L^1 -dependences in time which handles the case of weak convergence of the equations instead of the classical strong L^1 -convergence. This stability result is a particular case of a general stability result proved by Barles in [4].

For $T > 0$, we are interested in solutions of the following equation:

$$(52) \quad \begin{cases} \frac{\partial v}{\partial t} = \bar{c}(x, t)|Dv| & \text{in } \mathbb{R}^N \times (0, T), \\ v(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N, \end{cases}$$

where the velocity $\bar{c} : \mathbb{R}^N \times (0, T) \rightarrow \mathbb{R}$ is defined for almost every $t \in (0, T)$. We also assume that \bar{c} satisfies the following.

(H3) The function \bar{c} is continuous with respect to $x \in \mathbb{R}^N$ and measurable in t . For all $x, y \in \mathbb{R}^N$ and almost all $t \in [0, T]$,

$$|\bar{c}(x, t)| \leq M \quad \text{and} \quad |\bar{c}(x, t) - \bar{c}(y, t)| \leq L|x - y|.$$

Let us underline that we do not assume any continuity in time of \bar{c} . We recall the following (under assumption (H0)).

DEFINITION A.1 (L^1 -viscosity solutions). *An upper semicontinuous (respectively, lower semicontinuous) function v on $\mathbb{R}^N \times [0, T]$ is an L^1 -viscosity subsolution (respectively, supersolution) of (52) if*

$$v(0, \cdot) \leq u_0 \quad (\text{respectively, } v(0, \cdot) \geq u_0)$$

and if for every $(x_0, t_0) \in \mathbb{R}^N \times [0, T]$, $b \in L^1(0, T)$, $\varphi \in C^\infty(\mathbb{R}^N \times (0, T))$, and continuous function $G : \mathbb{R}^N \times (0, T) \times \mathbb{R}^N \rightarrow \mathbb{R}$ such that

(i) the function

$$(x, t) \mapsto v(x, t) - \int_0^t b(s)ds - \varphi(x, t)$$

has a local maximum (respectively, minimum) at (x_0, t_0) over $\mathbb{R}^N \times (0, T)$ and such that

(ii) for almost every $t \in (0, T)$ in some neighborhood of t_0 and for every (x, p) in some neighborhood of (x_0, p_0) with $p_0 = \nabla \varphi(x_0, t_0)$, we have

$$\bar{c}(x, t)|p| - b(t) \leq G(x, t, p) \quad (\text{respectively, } \bar{c}(x, t)|p| - b(t) \geq G(x, t, p)),$$

then

$$\frac{\partial \varphi}{\partial t}(x_0, t_0) \leq G(x_0, t_0, p_0) \quad \left(\text{respectively, } \frac{\partial \varphi}{\partial t}(x_0, t_0) \geq G(x_0, t_0, p_0) \right).$$

Finally, we say that a locally bounded function v defined on $\mathbb{R}^N \times [0, T]$ is an L^1 -viscosity solution of (52) if its upper semicontinuous (respectively, lower semicontinuous) envelope is an L^1 -viscosity subsolution (respectively, supersolution).

Let us recall that viscosity solutions in the L^1 -sense were introduced in Ishii's paper [20]. We refer the reader to Nunziante [25, 26] and Bourgoing [7, 8] for a complete presentation of the theory.

Then we have the following result.

THEOREM A.2 (existence and uniqueness). *For any $T > 0$, under assumptions (H0) and (H3), there exists a unique L^1 -viscosity solution to (52).*

Finally, let us consider the solutions v^ε to the following equation:

$$(53) \quad \begin{cases} \frac{\partial v^\varepsilon}{\partial t} = \bar{c}^\varepsilon(x, t) |Dv^\varepsilon| & \text{in } \mathbb{R}^N \times (0, T), \\ v^\varepsilon(\cdot, 0) = u_0 & \text{in } \mathbb{R}^N. \end{cases}$$

Then we have the following.

THEOREM A.3 (L^1 -stability [4]). *Under assumption (H0), let us assume that the velocity \bar{c}^ε satisfies (H3) (with some constants M, L independent of ε). Let us consider the L^1 -viscosity solution v^ε to (53). Assume that v^ε converges locally uniformly to a function v and, for all $x \in \mathbb{R}^N$,*

$$(54) \quad \int_0^t \bar{c}^\varepsilon(x, s) ds \rightarrow \int_0^t \bar{c}(x, s) ds \quad \text{locally uniformly in } (0, T).$$

Then v is a L^1 -viscosity solution of (52).

Remark A.1. Theorem A.3 is stated as in [4], but note that, under (H3), assumption (54) is automatically satisfied as soon as the convergence is merely pointwise. Indeed, since

$$\left| \int_0^t \bar{c}^\varepsilon(x, s) ds \right| \leq MT \quad \text{and} \quad \left| \int_0^t \bar{c}^\varepsilon(x, s) ds - \int_0^{t'} \bar{c}^\varepsilon(x, s) ds \right| \leq M|t - t'|,$$

from Ascoli's theorem, the convergence is uniform.

Appendix B. Interior ball regularization (proof of Lemma 2.3). The proof of this result can be adapted from those of Cannarsa and Frankowska [10] or [2] (see also [9] for related perimeter estimates for general equations). For the sake of completeness, we give a proof close to the one of [10] (this latter holds for much more general, but time-independent, dynamics). The unique (and small) contribution of this part amounts to explaining how this proof can be simplified in the particular case of dynamics of the form (55) and to point out that the time dependence is not an issue for the results of [10] to hold.

We first prove that the reachable set for controlled dynamics of the form

$$(55) \quad \dot{x}(t) = c(x(t), t)u(t), \quad u \in L^\infty([0, T], \bar{B}(0, 1)),$$

enjoys the interior ball property for positive time. We assume that $c : [0, T] \times \mathbb{R}^N \rightarrow \mathbb{R}$ satisfies, for any $x, y \in \mathbb{R}^N$ and $t \in [0, T]$,

$$\begin{cases} \text{(i)} & c \text{ is Borel measurable,} \\ & \text{differentiable with respect to the space variable for a.e. time,} \\ \text{(ii)} & |c(x, t) - c(y, t)| \leq L_1|x - y|, \\ \text{(iii)} & |D_x c(x, t) - D_x c(y, t)| \leq N_1|x - y|, \\ \text{(iv)} & M_1 \geq c(x, t) \geq \delta > 0, \end{cases}$$

where $L_1, N_1 \geq 0$ and $M_1, \delta > 0$ are given constants. Let $K_0 \subset \mathbb{R}^N$ be the initial set. We define the reachable set $\mathcal{R}(t)$ from K_0 for (55) at time t by

$$\mathcal{R}(t) = \{x(t), x(\cdot) \text{ solution to (55) with } x(0) \in K_0\}.$$

It is known that $\mathcal{R}(t)$ is a closed subset of \mathbb{R}^N . Let y_0 be an extremal solution on the time interval $[0, T]$, i.e., a solution of (55) such that

$$y_0(0) \in K_0 \quad \text{and} \quad y_0(T) \in \partial\mathcal{R}(T).$$

From the Pontryagin maximum principle for extremal trajectories (see, for instance, [13]), there is some adjoint function $p_0 : [0, T] \rightarrow \mathbb{R}^N \setminus \{0\}$ such that (y_0, p_0) is a solution to

$$\begin{cases} \dot{y}_0(t) = c(y_0(t), t) \frac{p_0(t)}{|p_0(t)|}, \\ -\dot{p}_0(t) = D_x c(y_0(t), t) |p_0(t)|. \end{cases}$$

Since the system is positively homogeneous with respect to p , we can assume, without loss of generality, that $|p_0(T)| = 1$, and we set $\theta_0 := p_0(T)$.

Let P be the matrix-valued solution to

$$\begin{cases} \dot{P}(t) = \frac{p_0(t)}{|p_0(t)|} [D_x c(y_0(t), t)]^* P(t), \\ P(T) = Id. \end{cases}$$

A straightforward computation shows that $P^*(t)p_0(t) = \theta_0$ for any $t \in [0, T]$.

Let us fix some parameter $\gamma > 0$ to be chosen later, $\theta \in B(0, 1)$ and let us set, for all $t \in [0, T]$,

$$y_\theta(t) = y_0(t) - \gamma t P(t)(\theta_0 - \theta).$$

Our aim is to show that y_θ is a solution to (55). Indeed we have

$$\begin{aligned} |\dot{y}_\theta|^2 &= \left| c(y_0, t) \frac{p_0}{|p_0|} - \gamma P(\theta_0 - \theta) - \gamma t \frac{p_0}{|p_0|} D_x c^* P(\theta_0 - \theta) \right|^2 \\ &= |c(y_\theta, t)|^2 - 2\gamma c(y_0, t) \left\langle \frac{p_0}{|p_0|}, P(\theta_0 - \theta) \right\rangle \\ &\quad + |c(y_0, t)|^2 - |c(y_\theta, t)|^2 - 2\gamma t c(y_0, t) \left\langle \frac{p_0}{|p_0|}, \frac{p_0}{|p_0|} D_x c^* P(\theta_0 - \theta) \right\rangle \\ &\quad + \gamma^2 \left| P(\theta_0 - \theta) + t \frac{p_0}{|p_0|} D_x c^* P(\theta_0 - \theta) \right|^2 \\ &\leq |c(y_\theta, t)|^2 - 2\gamma \frac{c(y_0, t)}{|p_0|} \langle \theta_0, \theta_0 - \theta \rangle \quad (\text{because } P^* p_0 = \theta_0) \\ &\quad + c^2(y_0, t) - c^2(y_\theta, t) - \langle D_x(c^2)(y_0, t), y_0 - y_\theta \rangle + \gamma^2 M |\theta_0 - \theta|^2 \\ &\leq |c(y_\theta, t)|^2 - \gamma \frac{\delta}{|p_0|} |\theta_0 - \theta|^2 + M'_1 |y_0 - y_\theta|^2 + \gamma^2 M |\theta_0 - \theta|^2 \\ &\leq |c(y_\theta, t)|^2 - \gamma \frac{\delta}{|p_0|} |\theta_0 - \theta|^2 + \gamma^2 M' |\theta_0 - \theta|^2, \end{aligned}$$

with $M'_1 = L_1^2 + M_1 N_1$, and where M and M' depend only on T, L_1, N_1, M_1 because $|p_0(t)|$ is bounded from below by a constant depending only on T, L_1 . Hence, for γ sufficiently small, y_θ is a solution of (55) starting from $y_0(0) \in K_0$ and therefore $y_\theta(T) \in \mathcal{R}(T)$.

Finally, $\mathcal{R}(T)$ contains all the $y_\theta(T)$ for $\theta \in B(0, 1)$, i.e., the ball centered at $y_0(T) - \gamma T \theta_0$ and of radius γT (since $P(T) = Id$).

We apply the previous result with $c = c_1$ and $K_0 = \{v(\cdot, 0) \geq 0\}$. Then $\{v(\cdot, t) \geq 0\} = \mathcal{R}(t)$ for all $t > 0$.

We end with a remark: in the statement of Lemma 2.3, c_1 is assumed to be continuous in time. As we have seen, it is not necessary; c_1 can be merely measurable in time up to considering the L^1 -solution v of (13), as recalled in Appendix A. \square

REFERENCES

- [1] N. ALIBAUD, *Existence, uniqueness and regularity for nonlinear degenerate parabolic equations with nonlocal terms*, NoDEA Nonlinear Differential Equations Appl., to appear.
- [2] O. ALVAREZ, P. CARDALIAGUET, AND R. MONNEAU, *Existence and uniqueness for dislocation dynamics with nonnegative velocity*, Interfaces Free Bound., 7 (2005), pp. 415–434.
- [3] O. ALVAREZ, P. HOCH, Y. LE BOUAR, AND R. MONNEAU, *Dislocation dynamics: Short-time existence and uniqueness of the solution*, Arch. Ration. Mech. Anal., 181 (2006), pp. 449–504.
- [4] G. BARLES, *A new stability result for viscosity solutions of nonlinear parabolic equations with weak convergence in time*, C. R. Math. Acad. Sci. Paris, 343 (2006), pp. 173–178.
- [5] G. BARLES AND O. LEY, *Nonlocal first-order Hamilton-Jacobi equations modelling dislocations dynamics*, Comm. Partial Differential Equations, 31 (2006), pp. 1191–1208.
- [6] G. BARLES, H. M. SONER, AND P. E. SOUGANIDIS, *Front propagation and phase field theory*, SIAM J. Control Optim., 31 (1993), pp. 439–469.
- [7] M. BOURGOING, *Viscosity solutions of fully nonlinear second order parabolic equations with L^1 -time dependence and Neumann boundary conditions*, Discrete Contin. Dyn. Syst., to appear.
- [8] M. BOURGOING, *Viscosity solutions of fully nonlinear second order parabolic equations with L^1 -time dependence and Neumann boundary conditions. Existence and applications to the level-set approach*, Discrete Contin. Dyn. Syst., to appear.
- [9] P. CANNARSA AND P. CARDALIAGUET, *Perimeter estimates for reachable sets of control systems*, J. Convex Anal., 13 (2006), pp. 253–267.
- [10] P. CANNARSA AND H. FRANKOWSKA, *Interior sphere property of attainable sets and time optimal control problems*, ESAIM Control Optim. Calc. Var., 12 (2006), pp. 350–370.
- [11] P. CARDALIAGUET, *On front propagation problems with nonlocal terms*, Adv. Differential Equations, 5 (2000), pp. 213–268.
- [12] P. CARDALIAGUET AND C. MARCHI, *Regularity of the eikonal equation with Neumann boundary conditions in the plane: Application to fronts with nonlocal terms*, SIAM J. Control Optim., 45 (2006), pp. 1017–1038.
- [13] F. H. CLARKE, *Optimization and Nonsmooth Analysis*, John Wiley and Sons, New York, 1983.
- [14] M. G. CRANDALL AND P.-L. LIONS, *Viscosity solutions of Hamilton-Jacobi equations*, Trans. Amer. Math. Soc., 277 (1983), pp. 1–42.
- [15] F. DA LIO, N. FORCADEL, AND R. MONNEAU, *Convergence of a non-local eikonal equation to anisotropic mean curvature motion. Application to dislocations dynamics*, J. Eur. Math. Soc. (JEMS), to appear.
- [16] F. DA LIO, C. I. KIM, AND D. SLEPČEV, *Nonlocal front propagation problems in bounded domains with Neumann-type boundary conditions and applications*, Asymptot. Anal., 37 (2004), pp. 257–292.
- [17] N. FORCADEL, *Dislocation Dynamics with a Mean Curvature Term: Short Time Existence and Uniqueness*, preprint, 2005.
- [18] Y. GIGA, *Surface Evolution Equations: A Level-Set Method*, Monogr. Math. 99, Birkhäuser Verlag, Basel, 2006.
- [19] J. R. HIRTH AND L. LOTHE, *Theory of Dislocations*, 2nd ed., Krieger, Malabar, FL, 1992.
- [20] H. ISHII, *Hamilton-Jacobi equations with discontinuous Hamiltonians on arbitrary open sets*, Bull. Fac. Sci. Engrg. Chuo Univ., 28 (1985), pp. 33–77.

- [21] R.W. LARDNER, *Mathematical Theory of Dislocations and Fracture*, Mathematical Expositions 17, University of Toronto Press, Toronto, ON, Canada, 1974.
- [22] O. LEY, *Lower-bound gradient estimates for first-order Hamilton-Jacobi equations and applications to the regularity of propagating fronts*, Adv. Differential Equations, 6 (2001), pp. 547–576.
- [23] P.-L. LIONS AND B. PERTHAME, *Remarks on Hamilton-Jacobi equations with measurable time-dependent Hamiltonians*, Nonlinear Anal., 11 (1987), pp. 613–621.
- [24] F. R. N. NABARRO, *Theory of Crystal Dislocations*, Clarendon Press, Oxford, UK, 1969.
- [25] D. NUNZIANTE, *Uniqueness of viscosity solutions of fully nonlinear second order parabolic equations with discontinuous time-dependence*, Differential Integral Equations, 3 (1990), pp. 77–91.
- [26] D. NUNZIANTE, *Existence and uniqueness of unbounded viscosity solutions of parabolic equations with discontinuous time-dependence*, Nonlinear Anal., 18 (1992), pp. 1033–1062.
- [27] D. RODNEY, Y. LE BOUAR, AND A. FINEL, *Phase field methods and dislocations*, Acta Mater., 51 (2003), pp. 17–30.
- [28] D. SLEPČEV, *Approximation schemes for propagation of fronts with nonlocal velocities and Neumann boundary conditions*, Nonlinear Anal., 52 (2003), pp. 79–115.
- [29] A. SROUR, *Nonlocal Second-Order Hamilton-Jacobi Equations Arising in Tomographic Reconstruction*, preprint, 2007.

REFINABLE FUNCTIONS AND CASCADE ALGORITHMS IN WEIGHTED SPACES WITH HÖLDER CONTINUOUS MASKS*

BIN HAN†

Abstract. Refinable functions and cascade algorithms play a fundamental role in wavelet analysis, which is useful in many applications. In this paper we shall study several properties of refinable functions, cascade algorithms, and wavelets, associated with Hölder continuous masks, in the weighted subspaces $L_{2,p,\gamma}(\mathbb{R})$ of $L_2(\mathbb{R})$, where $1 \leq p \leq \infty$, $\gamma \geq 0$ and $f \in L_{2,p,\gamma}(\mathbb{R})$ means $\|f\|_{L_{2,p,\gamma}(\mathbb{R})} := \|\sum_{k \in \mathbb{Z}} |e^{\gamma|\cdot|} f(\cdot + 2\pi k)|^2\|_{L_p(\mathbb{T})}^{1/2} < \infty$. In particular, $\|f\|_{L_{2,1,\gamma}(\mathbb{R})} = \|f e^{\gamma|\cdot|}\|_{L_2(\mathbb{R})}$ and $\|f\|_{L_{2,\infty,0}(\mathbb{R})} = \|\sum_{k \in \mathbb{Z}} |\hat{f}(\cdot + 2\pi k)|^2\|_{L_\infty(\mathbb{T})}^{1/2}$. For a mask $\hat{a} \in C^\beta(\mathbb{T})$ with $\beta > 0$ and $\hat{a}(0) = 1$ (that is, \hat{a} is a Hölder continuous mask with Hölder exponent β), we prove that the cascade algorithm associated with the mask \hat{a} converges in the space $L_{2,\infty,0}(\mathbb{R})$ if and only if $\nu_2(\hat{a}) > 0$, where the quantity $\nu_2(\hat{a})$ will be defined in this paper and plays an important role in our study of refinable functions and cascade algorithms with Hölder continuous masks. In particular, if the shifts of a refinable function ϕ , satisfying $\hat{\phi}(2\cdot) = \hat{a}\hat{\phi}$, are stable in $L_2(\mathbb{R})$, then we must have $\nu_2(\hat{a}) > 0$, and therefore the cascade algorithm associated with mask \hat{a} converges in the space $L_{2,\infty,0}(\mathbb{R})$. Based on this result, we are able to settle several problems on refinable functions, cascade algorithms, and wavelets associated with masks having infinitely many nonzero Fourier coefficients. As an application of the characterization of the convergence of a cascade algorithm in the space $L_{2,\infty,0}(\mathbb{R})$, we are able to show that for a mask \hat{a} having exponential decay of order $r > 0$, the cascade algorithm associated with mask \hat{a} converges in the weighted space $L_{2,1,\gamma}(\mathbb{R})$ for $0 < \gamma < 2r$ if and only if $\nu_2(\hat{a}) > 0$. Consequently, if a mask \hat{a} has exponential decay of order $r > 0$ and $\nu_2(\hat{a}) > 0$, then its standard refinable function ϕ , defined by $\hat{\phi}(\xi) := \prod_{j=1}^{\infty} \hat{a}(2^{-j}\xi)$, must have exponential decay of order $2r$ in $L_2(\mathbb{R})$; that is, $\|\phi\|_{L_{2,1,\gamma}(\mathbb{R})}^2 = \int_{\mathbb{R}} |\phi(x)|^2 e^{2\gamma|x|} dx < \infty$ for all $0 < \gamma < 2r$. As another application of the characterization of the convergence of a cascade algorithm in the space $L_{2,\infty,0}(\mathbb{R})$, we completely characterize biorthogonal wavelets and Riesz wavelets in $L_2(\mathbb{R})$, which are derived from refinable functions and whose involved wavelet filters in the frequency domain are Hölder continuous. We shall also investigate some basic properties of the quantity $\nu_2(\hat{a})$ and discuss how to calculate and estimate $\nu_2(\hat{a})$. Examples using fractional splines and the Butterworth filters will be given to illustrate the results in this paper.

Key words. refinable functions, cascade algorithms, Hölder continuous masks, masks having infinitely many nonzero Fourier coefficients, transition operator, weighted L_2 spaces, biorthogonal wavelets, Riesz wavelets, exponential decay

AMS subject classifications. 42C40, 41A15, 46B15, 47B37

DOI. 10.1137/060661016

1. Introduction and motivation. A wavelet system is generally derived from a refinable function via a multiresolution analysis (MRA). We say that ϕ is a *refinable function* if it satisfies the refinement equation:

$$(1.1) \quad \hat{\phi}(2\xi) = \hat{a}(\xi)\hat{\phi}(\xi) \quad \text{a.e. } \xi \in \mathbb{R},$$

where \hat{a} is a 2π -periodic function, called the *mask* (or filter) for ϕ , and the Fourier transform is defined to be $\hat{f}(\xi) := \int_{\mathbb{R}} f(x)e^{-ix\xi} dx$ for $f \in L_1(\mathbb{R})$.

*Received by the editors May 26, 2006; accepted for publication (in revised form) December 4, 2007; published electronically March 26, 2008. This work was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC) under grant RGP 228051.

<http://www.siam.org/journals/sima/40-1/66101.html>

†Department of Mathematical and Statistical Sciences, University of Alberta, Edmonton T6G 2G1, AB, Canada (bhan@math.ualberta.ca).

For a 2π -periodic function \hat{a} , we say that \hat{a} has *exponential decay of order r* if \hat{a} is the restriction of a 2π -periodic function $\hat{a}(z)$ on the real line $\text{Im}(z) = 0$ such that $\hat{a}(z)$ is holomorphic in the strip $\{z \in \mathbb{C} : |\text{Im}(z)| < r\}$, where $\text{Im}(z)$ denotes the imaginary part of the complex number z . Write $\hat{a}(\xi) = \sum_{k \in \mathbb{Z}} a_k e^{-ik\xi}$ in terms of its Fourier series. It is easy to check that \hat{a} has exponential decay of order r if and only if for every $0 \leq \gamma < r$, there is a positive constant C_γ such that $|a_k| \leq C_\gamma e^{-\gamma|k|}$ for all $k \in \mathbb{Z}$. A particular family of masks with exponential decay consists of rational masks that can be written in the form of $\hat{b}(\xi)/\hat{c}(\xi)$ for some 2π -periodic trigonometric polynomials \hat{b} and \hat{c} such that $\hat{c}(\xi) \neq 0$ for all $\xi \in \mathbb{R}$. For example, the well-known Butterworth filters are rational masks [5, 11, 18, 25].

Masks with infinitely many nonzero Fourier coefficients, or equivalently, filters with infinite support, are called infinite impulse response (IIR) filters in electrical engineering. Due to some desirable properties, IIR filters, including masks with exponential decay and masks for bandlimited wavelets [7] and fractional splines [27], are of interest in some applications and have been extensively designed for various purposes in the area of digital signal processing in electrical engineering [3, 5, 18, 25, 27]. In contrast to the case of trigonometric polynomial masks whose various mathematical properties have been extensively studied and more or less well understood in the literature (see [1, 4, 6, 7, 8, 10, 15, 19, 26, 28] and the references therein), there are still several unsolved questions related to refinable functions and wavelets with masks having infinitely many nonzero Fourier coefficients. For example, Daubechies and Huang in [9] showed that if the standard refinable function ϕ , defined by

$$(1.2) \quad \hat{\phi}(\xi) := \prod_{j=1}^{\infty} \hat{a}(2^{-j}\xi), \quad \xi \in \mathbb{R},$$

lies in $L_1(\mathbb{R})$ with a mask \hat{a} having an absolutely summable sequence of Fourier coefficients, and if ϕ has exponential decay, then the mask \hat{a} must have exponential decay. But to the best of our knowledge, there are very few results on the converse direction, and it is widely believed that for a mask \hat{a} with exponential decay, its standard refinable function ϕ in (1.2) should also have exponential decay in certain spaces in some sense. This is one of our motivations to study refinable functions and cascade algorithms, associated with masks having infinitely many nonzero Fourier coefficients, in some weighted subspaces of $L_2(\mathbb{R})$.

Before proceeding further, let us introduce some definitions and notation. Let $\mathbb{T} := \mathbb{R}/[2\pi\mathbb{Z}]$ and $L_p(\mathbb{T})$ denote the linear space of all 2π -periodic measurable functions $f : \mathbb{R} \mapsto \mathbb{C}$ such that $2\pi\|f\|_{L_p(\mathbb{T})}^p := \int_{-\pi}^{\pi} |f(x)|^p dx < \infty$ for $1 \leq p < \infty$ and $\|f\|_{L_\infty(\mathbb{T})}$ denotes its essential upper bound. For $\gamma \geq 0$ and $1 \leq p \leq \infty$, throughout the paper, the space $L_{2,p,\gamma}(\mathbb{R})$ denotes the subspace of all $f \in L_2(\mathbb{R})$ such that

$$(1.3) \quad \|f\|_{L_{2,p,\gamma}(\mathbb{R})} := \left\| \left[\widehat{e^{\gamma|\cdot|} f}, \widehat{e^{\gamma|\cdot|} f} \right] \right\|_{L_p(\mathbb{T})}^{1/2} < \infty,$$

where the *bracket product* [20] is defined to be

$$(1.4) \quad [f, g](\xi) := \sum_{k \in \mathbb{Z}} f(\xi + 2\pi k) \overline{g(\xi + 2\pi k)}, \quad \xi \in \mathbb{R}, f, g \in L_2(\mathbb{R}).$$

It is easy to verify that $[f, g] \in L_2(\mathbb{T})$ for $f, g \in L_2(\mathbb{R})$ and $L_{2,p,\gamma}(\mathbb{R})$ is a Banach space. In particular, by Plancherel's theorem, we have

$$\|f\|_{L_{2,1,\gamma}(\mathbb{R})}^2 = \|e^{\gamma|\cdot|} f\|_{L_2(\mathbb{R})}^2 = \int_{\mathbb{R}} |f(x)|^2 e^{2\gamma|x|} dx.$$

Therefore, $L_{2,1,\gamma}(\mathbb{R})$ is a weighted subspace of $L_2(\mathbb{R})$, and it is a natural candidate of subspaces to measure the exponential decay of a function in $L_2(\mathbb{R})$. Note that $L_{2,1,0}(\mathbb{R}) = L_2(\mathbb{R})$.

We say that the shifts of a function f are *stable* in $L_2(\mathbb{R})$ if there exists a positive constant C such that $C^{-1} \leq [f, \hat{f}](\xi) \leq C$ for almost every $\xi \in \mathbb{R}$. It is well known that stability plays an important role in wavelet analysis [3, 4, 7]. If the shifts of a function f are stable in $L_2(\mathbb{R})$, then it is obvious that $f \in L_{2,\infty,0}(\mathbb{R})$, since $\|f\|_{L_{2,\infty,0}(\mathbb{R})} := \|[\hat{f}, \hat{f}]\|_{L_\infty(\mathbb{T})}^{1/2} < \infty$. On the other hand, we shall see in Proposition 6.1 that $L_{2,1,\gamma}(\mathbb{R}) \subseteq L_{2,\infty,0}(\mathbb{R})$ for any $\gamma > 0$. In particular, all compactly supported functions in $L_2(\mathbb{R})$ are included in $L_{2,\infty,0}(\mathbb{R})$. Therefore, the space $L_{2,\infty,0}(\mathbb{R})$ is a large enough subspace of $L_2(\mathbb{R})$ and includes most interesting functions in wavelet analysis. In this paper, we are particularly interested in the subspaces $L_{2,\infty,0}(\mathbb{R})$ and $L_{2,1,\gamma}(\mathbb{R})$ for $\gamma > 0$.

In the following, let us introduce a basic quantity $\nu_2(\hat{a})$, which plays a critical role in our study of refinable functions, cascade algorithms, and wavelets with masks having infinitely many nonzero Fourier coefficients. For 2π -periodic functions \hat{a} and f , the *transition operator* $T_{\hat{a}}$ is defined to be

$$(1.5) \quad [T_{\hat{a}}f](\xi) := |\hat{a}(\xi/2)|^2 f(\xi/2) + |\hat{a}(\xi/2 + \pi)|^2 f(\xi/2 + \pi), \quad \xi \in \mathbb{R}.$$

For $\tau \in \mathbb{R}$ and $1 \leq p \leq \infty$, we define a quantity

$$(1.6) \quad \rho_\tau(\hat{a}, p) := \limsup_{n \rightarrow \infty} \|T_{\hat{a}}^n (|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})}^{1/n}.$$

Now we define the quantity $\nu_2(\hat{a})$ in this paper as follows:

$$(1.7) \quad \nu_2(\hat{a}) := -[\log_2 \rho(\hat{a})]/2,$$

where

$$(1.8) \quad \rho(\hat{a}) := \inf\{\rho_\tau(\hat{a}, \infty) : |\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T}) \text{ and } \tau \geq 0\}.$$

When \hat{a} is a 2π -periodic trigonometric polynomial, the quantity $\nu_2(\hat{a})$ in (1.7) agrees with the one in [13] and can be calculated by finding the spectral radius of an associated finite matrix generated by \hat{a} . See section 4 for details on calculating and estimating the quantity $\nu_2(\hat{a})$. The quantity $\nu_2(\hat{a})$ in (1.7), whose definition appears to be a little bit technical in its nature, plays a very important role in investigating many problems in wavelet analysis. See [13] for applications and the importance of $\nu_2(\hat{a})$ in wavelet analysis with 2π -periodic trigonometric polynomial masks \hat{a} .

We say that f belongs to the *Hölder class* $C^\beta(\mathbb{T})$ with $\beta > 0$ if f is a 2π -periodic continuous function such that $f \in C^n(\mathbb{T})$ and there exists a positive constant C satisfying $|f^{(n)}(x) - f^{(n)}(y)| \leq C|x - y|^{\beta-n}$ for all $x, y \in \mathbb{T}$, where n is the largest integer such that $n \leq \beta$ and $f^{(n)}$ denotes the n th derivative of f . Throughout the paper, we say that \hat{a} is a *Hölder continuous mask* if $\hat{a} \in C^\beta(\mathbb{T})$ for some $\beta > 0$ and $\hat{a}(0) = 1$. That \hat{a} is a Hölder continuous mask is a very natural and weak condition to guarantee that as the Fourier transform of the standard refinable function ϕ with mask \hat{a} , the function $\hat{\phi}$, which is defined through the infinite product in (1.2), is well defined.

In section 2, we shall present a necessary and sufficient condition in Theorem 2.1 for the convergence of a cascade algorithm in the space $L_{2,\infty,0}(\mathbb{R})$. Theorem 2.1 plays a central role in our study of refinable functions with exponential decay in $L_2(\mathbb{R})$ and

of MRA Riesz wavelet bases in $L_2(\mathbb{R})$. More precisely, we prove in Theorem 2.1 that for a mask $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$ (that is, \hat{a} is a Hölder continuous mask), the cascade algorithm associated with mask \hat{a} converges in the space $L_{2,\infty,0}(\mathbb{R})$ if and only if $\nu_2(\hat{a}) > 0$. As a direct consequence of Theorem 2.1, we show in Corollary 2.2 that if the refinement equation $\hat{\phi}(2\cdot) = \hat{a}\hat{\phi}$ has a solution ϕ such that the shifts of ϕ are stable in $L_2(\mathbb{R})$, then we must have $\nu_2(\hat{a}) > 0$, and therefore the cascade algorithm associated with mask \hat{a} converges in the space $L_{2,\infty,0}(\mathbb{R})$.

Cascade algorithms in $l_2(\mathbb{R})$ and other spaces with 2π -periodic trigonometric polynomial masks have been extensively studied in the literature. To cite only a few references here, we refer the reader to [1, 4, 6, 7, 8, 10, 15, 19, 22, 26, 28] and the references therein, where the property of the masks being 2π -periodic trigonometric polynomials plays a critical role. The study of cascade algorithms and refinable functions with Hölder continuous masks in this paper is not a trivial generalization of the known results in the literature, as we shall see in sections 4–6.

As an application of Theorem 2.1, we are able to prove in Theorem 2.3 that for a mask \hat{a} having exponential decay of order $r > 0$, the cascade algorithm associated with the mask \hat{a} converges in the space $L_{2,1,\gamma}(\mathbb{R})$ for $0 < \gamma < 2r$ if and only if $\nu_2(\hat{a}) > 0$. Consequently, the standard refinable function ϕ associated with mask \hat{a} in (1.2) must have exponential decay of order $2r$, namely,

$$(1.9) \quad \|\phi\|_{L_{2,1,\gamma}(\mathbb{R})}^2 = \int_{\mathbb{R}} |\phi(x)|^2 e^{2\gamma|x|} dx < \infty \quad \forall 0 \leq \gamma < 2r.$$

An MRA wavelet function ψ is obtained from a refinable function ϕ with mask \hat{a} via

$$(1.10) \quad \hat{\psi}(2\xi) := \hat{b}(\xi)\hat{\phi}(\xi), \quad \xi \in \mathbb{R},$$

for some 2π -periodic measurable function \hat{b} . We say that ψ generates a *Riesz wavelet basis* in $L_2(\mathbb{R})$ if $\text{span}\{\psi_{j,k} := 2^{j/2}\psi(2^j \cdot -k) : j, k \in \mathbb{Z}\}$ is dense in $L_2(\mathbb{R})$ and there exists a positive constant C such that

$$(1.11) \quad C^{-1} \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} |c_{j,k}|^2 \leq \left\| \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} c_{j,k} \psi_{j,k} \right\|_{L_2(\mathbb{R})}^2 \leq C \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} |c_{j,k}|^2$$

for all finitely supported sequences $\{c_{j,k}\}_{j,k \in \mathbb{Z}}$. MRA Riesz wavelet bases in $L_2(\mathbb{R})$ are of interest in some applications [2, 5, 17, 21, 24]. A natural and important question here is when ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$. MRA Riesz wavelet bases with compact support have been investigated in [5, 14, 16, 17, 21, 24], where some necessary and sufficient conditions have been obtained for trigonometric polynomial masks. Most approaches in these papers rely largely on an interesting result of Cohen and Daubechies in [5] saying that for a mask \hat{a} with exponential decay, the transition operator $T_{\hat{a}}$ acting on some weighted subspaces of $\ell_2(\mathbb{Z})$ is a compact operator. Built on this interesting result of [5], a characterization of Riesz wavelets with trigonometric polynomial masks is obtained in [16] (also cf. [5]) in terms of the spectrum of $T_{\hat{a}}$. However, the approach in [5, 16, 24] seems difficult, if not impossible, to be generalized to masks without exponential decay, since the compactness of the operator $T_{\hat{a}}$ may be lost.

In order to study biorthogonal wavelets and Riesz wavelets with Hölder continuous masks, using a quite different approach in this paper, as another application

of Theorem 2.1, we shall prove in Theorem 3.2 that for $\hat{a}, \hat{b} \in C^\beta(\mathbb{T})$ with $\beta > 0$, assuming that the shifts of the standard refinable function ϕ with mask \hat{a} are stable in $L_2(\mathbb{R})$, then ψ in (1.10) generates a Riesz wavelet basis in $L_2(\mathbb{R})$ if and only if (i) $\hat{b}(0) = 0$ and $d(\xi) := \hat{a}(\xi)\hat{b}(\xi + \pi) - \hat{a}(\xi + \pi)\hat{b}(\xi) \neq 0$ for all $\xi \in \mathbb{R}$; (ii) $\nu_2(\hat{a}) > 0$ and $\nu_2(\hat{a}) > 0$, where $\hat{a}(\xi) := \hat{b}(\xi + \pi)/d(\xi)$. Moreover, in the case that \hat{a} is a 2π -periodic trigonometric polynomial, we prove that the shifts of $\phi \in L_2(\mathbb{R})$ must be stable in $L_2(\mathbb{R})$ if ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$. Even for the case that both \hat{a} and \hat{b} are 2π -periodic trigonometric polynomials, the mask \hat{a} is generally not a 2π -periodic trigonometric polynomial, and consequently refinable functions with masks being nontrigonometric polynomials will naturally appear in our study of MRA Riesz wavelet bases in $L_2(\mathbb{R})$. This is another motivation for us to study refinable functions, cascade algorithms, and wavelets with masks having infinitely many nonzero Fourier coefficients.

To illustrate the results in this paper, we shall apply the results in sections 2 and 3 to a family of Hölder continuous masks $\widehat{a_{\beta_1, \beta_2, \beta_3}}$ and $|\widehat{a_{\beta_1, \beta_2, \beta_3}}|$, where

$$(1.12) \quad \widehat{a_{\beta_1, \beta_2, \beta_3}}(\xi) := \frac{2^{-2\beta_1}(1 + e^{-i\xi})^{2\beta_1}}{(|\cos(\xi/2)|^{2\beta_2} + |\sin(\xi/2)|^{2\beta_2})^{\beta_3}}, \quad \beta_1, \beta_2, \beta_3 > 0.$$

In fact, the masks for the B -splines correspond to the case $\widehat{a_{\beta_1, \beta_2, \beta_3}}$ with $\beta_2 = 1$ and $2\beta_1 \in \mathbb{N}$. The masks for various types of fractional splines in [27] correspond to the case $\widehat{a_{\beta_1, \beta_2, \beta_3}}$ or $|\widehat{a_{\beta_1, \beta_2, \beta_3}}|$ with $\beta_2 = 1$. The classical Butterworth filters in [5, 11, 18, 25] correspond to the case $|\widehat{a_{\beta_1, \beta_2, \beta_3}}|$ with $\beta_3 = 1$ and $\beta_1 = \beta_2 \in \mathbb{N}$. To illustrate the main results in sections 2 and 3, we shall study the convergence of cascade algorithms and MRA Riesz wavelet bases associated with the masks given in (1.12).

Since the quantity $\nu_2(\hat{a})$ plays a very important role in our study of refinable functions, cascade algorithms, and wavelets with Hölder continuous masks, we shall investigate in section 4 some basic properties of the quantity $\nu_2(\hat{a})$ and discuss how to calculate and estimate the quantity $\nu_2(\hat{a})$ for a mask \hat{a} being a general Lebesgue measurable 2π -periodic function. In section 4, we shall generalize a well-known result on $\nu_2(\hat{a})$, whose proof in the general case of Hölder continuous masks is nontrivial and will be presented in the last section of this paper. The general discussion on the quantity $\nu_2(\hat{a})$ in section 4 is of interest in its own right, and the results in section 4 may be useful elsewhere.

For simplicity of presentation and readability of this paper, the proofs of Theorems 2.1 and 2.3 in section 2, which are a little bit technical in their nature, will be postponed to sections 5 and 6, respectively. The results in this paper can be nontrivially generalized to high dimensions and multiwavelets, which we shall discuss elsewhere.

2. Convergence of cascade algorithms in subspaces of $L_2(\mathbb{R})$. In this section, we shall present the main results on the convergence of cascade algorithms in the subspaces $L_{2, \infty, 0}(\mathbb{R})$ and $L_{2, 1, \gamma}(\mathbb{R})$. For simplicity of presentation, the proofs of the main results in this section will be postponed to sections 5 and 6. To illustrate the results in this section, we shall apply these results to the masks in (1.12), which include the Butterworth filters in [25] and the masks for fractional splines in [27] as special cases.

For a quotient function f/g , throughout the paper, we use the convention that $(f/g)(\xi)$ is equal to $f(\xi)/g(\xi)$ if $g(\xi) \neq 0$, 1 if $f(\xi) = g(\xi) = 0$, or $+\infty$ if $g(\xi) = 0$ but

$f(\xi) \neq 0$. We say that a function $f \in L_2(\mathbb{R})$ is *admissible with respect to* \hat{a} if there exists a positive number $\tau > 0$ such that

$$(2.1) \quad [\hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}, \hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}]/|\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T}).$$

Note that all compactly supported functions in $L_2(\mathbb{R})$ belong to $L_{2,1,\gamma}(\mathbb{R})$ for all $\gamma > 0$. For every $f \in L_{2,1,\gamma}(\mathbb{R})$ with $\gamma > 0$ such that $\hat{f}(2\pi k) = 0$ for all $k \in \mathbb{Z} \setminus \{0\}$, we shall show in Proposition 6.2 that f is admissible with respect to \hat{a} for any $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$, $\hat{a}(\pi) = 0$, and $\beta > 0$.

For $0 < \beta \leq 1$, we say that f belongs to the Lipschitz class $\Lambda^\beta(\mathbb{T})$ if there is a positive constant C such that $|f(x) - f(y)| \leq C|x - y|^\beta$ for all $x, y \in \mathbb{T}$.

Now we have the following result on the convergence of cascade algorithms in the space $L_{2,\infty,0}(\mathbb{R})$, which plays a central role in this paper and whose proof will be given in section 5.

THEOREM 2.1. *Let $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$ (that is, \hat{a} is a Hölder continuous mask). Then the following are equivalent:*

- (i) $\hat{a}(\pi) = 0$, and for every admissible function $f \in L_{2,\infty,0}(\mathbb{R})$ with respect to \hat{a} , $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$, where the functions f_n are defined by

$$(2.2) \quad \widehat{f_n}(\xi) := \hat{a}(\xi/2)\widehat{f_{n-1}}(\xi/2) = \hat{f}(2^{-n}\xi) \prod_{j=1}^n \hat{a}(2^{-j}\xi), \quad \xi \in \mathbb{R}, n \in \mathbb{N}.$$

- (ii) For one admissible function $f \in L_{2,\infty,0}(\mathbb{R})$ with respect to \hat{a} such that the shifts of f are stable, $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$.
- (iii) For every $\tau > 0$, $\rho_\tau(\hat{a}, \infty) < 1$, where $\rho_\tau(\hat{a}, \infty)$ is defined in (1.6).
- (iv) $\nu_2(\hat{a}) > 0$; that is, $\rho(\hat{a}) < 1$, where $\nu_2(\hat{a})$ and $\rho(\hat{a})$ are defined in (1.7) and (1.8), respectively.
- (v) For at least one $\tau > 0$, $\rho_\tau(\hat{a}, \infty) < 1$ and $|\hat{a}(\cdot + \pi)|^2/|\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$.

Let ϕ denote the standard refinable function associated with the mask \hat{a} in (1.2). If $\nu_2(\hat{a}) > 0$, then $\hat{\phi}$ belongs to the Lipschitz class $\Lambda^{\min(1,\beta)}(\mathbb{R})$, $[\hat{\phi}, \hat{\phi}] \in \Lambda^{\min(1,\beta)}(\mathbb{T})$, and for any $0 \leq \nu < \nu_2(\hat{a})$,

$$(2.3) \quad [\hat{\phi}, \hat{\phi}]_\nu := \sum_{k \in \mathbb{Z}} (1 + |\cdot + 2\pi k|^2)^\nu |\hat{\phi}(\cdot + 2\pi k)|^2 \in C(\mathbb{T}).$$

In fact, by a more complicated argument, we could show in Theorem 2.1 that $\hat{\phi} \in C^\beta(\mathbb{R})$ and $[\hat{\phi}, \hat{\phi}] \in C^\beta(\mathbb{T})$, which we shall address elsewhere. According to the proof of Theorem 2.1 in section 5, (ii) implies $\nu_2(\hat{a}) > 0$ without the admissibility condition on f ; that is, if $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$ for a function $f \in L_{2,\infty,0}(\mathbb{R})$ with stability, then $\nu_2(\hat{a}) > 0$. We say that the cascade algorithm associated with a mask \hat{a} converges in a given function space if for every admissible function f in that space with respect to \hat{a} , the sequence $\{f_n\}_{n=1}^\infty$ defined in (2.2) is a Cauchy sequence in that space.

As a direct consequence of Theorem 2.1, we have the following corollary.

COROLLARY 2.2. *Let $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$. Suppose that ϕ is a (not necessarily the standard) refinable function such that $\hat{\phi}(2\xi) = \hat{a}(\xi)\hat{\phi}(\xi)$ a.e. $\xi \in \mathbb{R}$ and the shifts of ϕ are stable in $L_2(\mathbb{R})$. Then $\nu_2(\hat{a}) > 0$ and the cascade algorithm associated with mask \hat{a} must converge in the space $L_{2,\infty,0}(\mathbb{R})$.*

Proof. Since the shifts of ϕ are stable in $L_2(\mathbb{R})$, we have $[\hat{\phi}, \hat{\phi}] \in L_\infty(\mathbb{T})$, and so $\phi \in L_{2,\infty,0}(\mathbb{R})$. Since $\hat{a}(\cdot/2)\hat{\phi}(\cdot/2) - \hat{\phi} = 0$, $\hat{\phi}$ is an admissible function in $L_{2,\infty,0}(\mathbb{R})$

with respect to \hat{a} . So, for $f = \phi$, (ii) of Theorem 2.1 holds since $f_n = \phi$ for all $n \in \mathbb{N}$. Now by Theorem 2.1, $\nu_2(\hat{a}) > 0$. \square

For a mask with exponential decay, using Theorem 2.1, we characterize the convergence of a cascade algorithm with an exponentially decaying mask in the spaces $L_{2,p,\gamma}(\mathbb{R})$ in the following result, whose proof will be given in section 6.

THEOREM 2.3. *Let \hat{a} be a mask such that $\hat{a}(0) = 1$ and \hat{a} has exponential decay of order r for some $r > 0$. Then the following are equivalent:*

- (i) $\hat{a}(\pi) = 0$, and for every $0 < \gamma < 2r$ and every admissible function $f \in L_{2,1,\gamma}(\mathbb{R})$ with respect to \hat{a} , $\{f_n\}_{n=1}^{\infty}$ is a Cauchy sequence in $L_{2,1,\gamma}(\mathbb{R})$, where f_n are defined in (2.2).
- (ii) $\hat{a}(\pi) = 0$, and for every $0 < \gamma < 2r$, every $1 \leq p \leq \infty$, and every admissible function $f \in L_{2,p,\gamma}(\mathbb{R})$ with respect to \hat{a} , $\{f_n\}_{n=1}^{\infty}$ is a Cauchy sequence in $L_{2,p,\gamma}(\mathbb{R})$.
- (iii) $\hat{a}(\pi) = 0$, and for some $0 < \gamma < 2r$, some $1 \leq p \leq \infty$, and every admissible function $f \in L_{2,p,\gamma}(\mathbb{R})$ with respect to \hat{a} , $\{f_n\}_{n=1}^{\infty}$ is a Cauchy sequence in $L_{2,p,\gamma}(\mathbb{R})$.
- (iv) For some $0 < \gamma < 2r$, some $1 \leq p \leq \infty$, and one admissible function $f \in L_{2,p,\gamma}(\mathbb{R})$ with respect to \hat{a} such that the shifts of f are stable in $L_2(\mathbb{R})$, the sequence $\{f_n\}_{n=1}^{\infty}$ is a Cauchy sequence in $L_{2,p,\gamma}(\mathbb{R})$.
- (v) $\nu_2(\hat{a}) > 0$.

In particular, if $\nu_2(\hat{a}) > 0$ and \hat{a} has exponential decay of order r , then the standard refinable function ϕ with mask \hat{a} in (1.2) must have exponential decay of order $2r$; that is, (1.9) holds.

In the following, we shall apply the above results to the masks in (1.12), which include both the classical Butterworth filters in [25] and the masks for the fractional splines in [27] as special cases.

Example 2.4. Let $\hat{a} = \widehat{a_{\beta_1, \beta_2, \beta_3}}$ or $\hat{a} = |\widehat{a_{\beta_1, \beta_2, \beta_3}}|$, where the masks $\widehat{a_{\beta_1, \beta_2, \beta_3}}$ are defined in (1.12). By Proposition 4.2, it is easy to check that $\hat{a} \in C^\beta(\mathbb{T})$ for some $\beta > 0$. Denote

$$\hat{B}(\xi) = (|\cos(\xi/2)|^{2\beta_2} + |\sin(\xi/2)|^{2\beta_2})^{\beta_3}.$$

Then $|\hat{a}(\xi)| = 2^{-2\beta_1} |1 + e^{-i\xi}|^{2\beta_1} / \hat{B}(\xi)$. By calculation, it is easy to deduce that

$$(2.4) \quad \min(1, 2^{1-\beta_2}) \leq |\cos(\xi/2)|^{2\beta_2} + |\sin(\xi/2)|^{2\beta_2} \leq \max(1, 2^{1-\beta_2}) \quad \forall \xi \in \mathbb{R}.$$

Consequently, we have $\hat{B}(\xi) \geq \min(1, 2^{(1-\beta_2)\beta_3})$ for all $\xi \in \mathbb{R}$ and $\beta_2, \beta_3 > 0$.

For $0 < \beta_2 \leq 1$, by Theorem 4.1 and Lemma 4.3, it follows from $\hat{B}(\xi) \geq \min(1, 2^{(1-\beta_2)\beta_3}) = 1$ that

$$\rho(\hat{a}) = \rho_{4\beta_1}(\hat{a}, \infty) = \rho_0(2^{-2\beta_1} / \hat{B}, \infty) \leq \rho_0(2^{-2\beta_1}, \infty) = 2^{1-4\beta_1},$$

since for a constant c , $T_c^n 1 = 2^n |c|^{2n}$, and therefore

$$\rho_0(c, \infty) = \limsup_{n \rightarrow \infty} \|T_c^n 1\|_{L^\infty(\mathbb{T})}^{1/n} = 2|c|^2.$$

For $\beta_2 > 1$, by Theorem 4.1 and Lemma 4.3, it follows from $\hat{B}(\xi) \geq 2^{(1-\beta_2)\beta_3}$ that

$$\rho(\hat{a}) = \rho_{4\beta_1}(\hat{a}, \infty) = \rho_0(2^{-2\beta_1} / \hat{B}, \infty) \leq \rho_0(2^{-2\beta_1 - (1-\beta_2)\beta_3}, \infty) = 2^{1-4\beta_1 - 2(1-\beta_2)\beta_3}.$$

Therefore, for

$$(2.5) \quad \beta_1, \beta_2, \beta_3 > 0 \quad \text{satisfying} \quad \beta_1 > 1/4 + \max(0, (\beta_2 - 1)\beta_3/2),$$

we have $\rho(\hat{a}) < 1$; that is, $\nu_2(\hat{a}) > 0$. By Theorem 2.1, the cascade algorithm associated with mask \hat{a} converges in the space $L_{2,\infty,0}(\mathbb{R})$.

The classical Butterworth filters in [25] correspond to the case $\beta_3 = 1$ and $\beta_1 = \beta_2 \in \mathbb{N}$. The condition in (2.5) holds for all $\beta_3 = 1$ and $\beta_1 = \beta_2 \in \mathbb{N}$, that is, holds for all Butterworth filters.

For the fractional splines in [27], since $\beta_2 = 1$, the condition in (2.5) becomes $\beta_1 > 1/4$. Note that when $\beta_2 = 1$, the standard refinable function ϕ associated with mask \hat{a} in (1.2) satisfies $|\hat{\phi}(2\xi)| = |(\sin \xi)/\xi|^{2\beta_1}$. Clearly, if $0 < \beta_1 \leq 1/4$, then $\hat{\phi} \notin L_2(\mathbb{R})$. So, the condition in (2.5), that is, $\beta_1 > 1/4$, is sharp for the case of fractional splines in [27].

Now we consider the case that \hat{a} has exponential decay; that is, $\hat{a}(\xi) = a_{\widehat{\beta_1, \beta_2, \beta_3}}(\xi)$ with $2\beta_1 \in \mathbb{N}$ and $\beta_2 \in \mathbb{N}$. In this case, by the definition of the masks $a_{\beta_1, \beta_2, \beta_3}$ in (1.12), it is easy to see that \hat{a} can be extended into a holomorphic function on some strip $\{z \in \mathbb{C} : |\text{Im}(z)| < r\}$ for some $r > 0$ depending only on β_2 . So, \hat{a} has exponential decay of order r . Now by Theorem 2.3, if (2.5) holds with $2\beta_1 \in \mathbb{N}$ and $\beta_2 \in \mathbb{N}$, then the cascade algorithm associated with mask \hat{a} converges in the spaces $L_{2,p,\gamma}(\mathbb{R})$ for all $0 \leq \gamma < 2r$ and the standard refinable function ϕ associated with mask \hat{a} must have exponential decay of order $2r$ in $L_2(\mathbb{R})$.

3. Characterization of MRA biorthogonal wavelets and Riesz wavelets.

In this section, using Theorem 2.1, we shall study MRA biorthogonal wavelets and Riesz wavelets with Hölder continuous masks.

Let us first recall the definition of biorthogonal wavelets in [6]. For two functions $\psi, \tilde{\psi} \in L_2(\mathbb{R})$, we say that $(\psi, \tilde{\psi})$ generates a pair of biorthogonal wavelet bases in $L_2(\mathbb{R})$ if each of ψ and $\tilde{\psi}$ generates a Riesz wavelet basis in $L_2(\mathbb{R})$ and the following biorthogonality relation holds:

$$(3.1) \quad \langle \psi_{j,k}, \tilde{\psi}_{j',k'} \rangle := \int_{\mathbb{R}} \psi_{j,k}(x) \overline{\tilde{\psi}_{j',k'}(x)} dx = \delta_{j-j'} \delta_{k-k'} \quad \forall j, j', k, k' \in \mathbb{Z},$$

where $\psi_{j,k} := 2^{j/2} \psi(2^j \cdot -k)$ and δ denotes the Dirac sequence such that $\delta_0 = 1$ and $\delta_k = 0$ for all $k \neq 0$. Compactly supported biorthogonal wavelets have been investigated in [4, 6, 7, 12] and other papers.

As an application of Theorem 2.1, we have the following result on biorthogonal wavelets with Hölder continuous masks.

THEOREM 3.1. *Let $\hat{a}, \hat{\hat{a}} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = \hat{\hat{a}}(0) = 1$ and $\beta > 0$. Define two refinable functions $\hat{\phi}$ and $\hat{\hat{\phi}}$ associated with masks \hat{a} and $\hat{\hat{a}}$ by*

$$(3.2) \quad \hat{\phi}(\xi) := \prod_{j=1}^{\infty} \hat{a}(2^{-j}\xi) \quad \text{and} \quad \hat{\hat{\phi}}(\xi) := \prod_{j=1}^{\infty} \hat{\hat{a}}(2^{-j}\xi), \quad \xi \in \mathbb{R}.$$

Then the following are equivalent:

- (i) $[\hat{\phi}, \hat{\hat{\phi}}], [\hat{\hat{\phi}}, \hat{\phi}] \in L_\infty(\mathbb{T})$ and $[\hat{\phi}, \hat{\hat{\phi}}] = 1$; or equivalently, the shifts of both ϕ and $\hat{\phi}$ are stable in $L_2(\mathbb{R})$ and the biorthogonality relation holds:

$$(3.3) \quad \langle \phi, \tilde{\phi}(\cdot - k) \rangle = \int_{\mathbb{R}} \phi(x) \overline{\tilde{\phi}(x - k)} dx = \delta_k \quad \forall k \in \mathbb{Z}.$$

- (ii) $\nu_2(\hat{a}) > 0$, $\nu_2(\hat{\hat{a}}) > 0$, and $\hat{\hat{a}}$ is a dual mask of \hat{a} , where we say that $\hat{\hat{a}}$ is a dual mask of \hat{a} if

$$(3.4) \quad \overline{\hat{a}(\xi)\hat{\hat{a}}(\xi)} + \overline{\hat{a}(\xi+\pi)\hat{\hat{a}}(\xi+\pi)} = 1.$$

Let $\hat{b}, \hat{\hat{b}} \in C^\beta(\mathbb{T})$ such that $\hat{b}(0) = \hat{\hat{b}}(0) = 0$. Define two wavelet functions ψ and $\tilde{\psi}$ by

$$(3.5) \quad \hat{\psi}(\xi) := \hat{b}(\xi/2)\hat{\phi}(\xi/2) \quad \text{and} \quad \hat{\tilde{\psi}}(\xi) := \hat{\hat{b}}(\xi/2)\hat{\hat{\phi}}(\xi/2).$$

If $\nu_2(\hat{a}) > 0$, $\nu_2(\hat{\hat{a}}) > 0$, and

$$(3.6) \quad \begin{bmatrix} \hat{a}(\xi) & \hat{a}(\xi+\pi) \\ \hat{b}(\xi) & \hat{b}(\xi+\pi) \end{bmatrix} \overline{\begin{bmatrix} \hat{\hat{a}}(\xi) & \hat{\hat{a}}(\xi+\pi) \\ \hat{\hat{b}}(\xi) & \hat{\hat{b}}(\xi+\pi) \end{bmatrix}}^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix},$$

then $(\psi, \tilde{\psi})$ generates a pair of biorthogonal wavelet bases in $L_2(\mathbb{R})$.

Proof. Suppose that (i) holds. By $1 = |[\hat{\phi}, \hat{\phi}]|^2 \leq [\hat{\phi}, \hat{\phi}][\hat{\hat{\phi}}, \hat{\hat{\phi}}]$, we have $[\hat{\phi}, \hat{\phi}] \geq \|[\hat{\phi}, \hat{\phi}]\|_{L^\infty(\mathbb{T})}^{-1}$ and $[\hat{\phi}, \hat{\phi}] \geq \|[\hat{\phi}, \hat{\phi}]\|_{L^\infty(\mathbb{T})}^{-1}$. Therefore, the shifts of ϕ and $\tilde{\phi}$ are stable in $L_2(\mathbb{R})$. By Corollary 2.2, we have $\nu_2(\hat{a}) > 0$ and $\nu_2(\hat{\hat{a}}) > 0$. Now by $[\hat{\phi}, \hat{\phi}] = 1$, it follows directly from the refinement equations $\hat{\phi}(2\xi) = \hat{a}(\xi)\hat{\phi}(\xi)$ and $\hat{\hat{\phi}}(2\xi) = \hat{\hat{a}}(\xi)\hat{\hat{\phi}}(\xi)$ that

$$\begin{aligned} 1 &= [\hat{\phi}, \hat{\phi}](2\xi) = \overline{\hat{a}(\xi)\hat{\hat{a}}(\xi)}[\hat{\phi}, \hat{\phi}](\xi) + \overline{\hat{a}(\xi+\pi)\hat{\hat{a}}(\xi+\pi)}[\hat{\phi}, \hat{\phi}](\xi+\pi) \\ &= \overline{\hat{a}(\xi)\hat{\hat{a}}(\xi)} + \overline{\hat{a}(\xi+\pi)\hat{\hat{a}}(\xi+\pi)}. \end{aligned}$$

So, $\hat{\hat{a}}$ is a dual mask of \hat{a} . Therefore, (i) \Rightarrow (ii).

To prove (ii) \Rightarrow (i), since $\nu_2(\hat{a}) > 0$ and $\nu_2(\hat{\hat{a}}) > 0$, by Theorem 2.1, we see that (i) of Theorem 2.1 holds for both \hat{a} and $\hat{\hat{a}}$. Moreover, by (2.3), we have $[\hat{\phi}, \hat{\phi}], [\hat{\hat{\phi}}, \hat{\hat{\phi}}] \in L^\infty(\mathbb{T})$.

Take $\hat{f} := \chi_{[-\pi, \pi]}$, the characteristic function of the interval $[-\pi, \pi]$. By (i) of Theorem 2.1, we have $\hat{a}(\pi) = \hat{\hat{a}}(\pi) = 0$. Since $\hat{a}, \hat{\hat{a}} \in C^\beta(\mathbb{T})$ and $\hat{a}(0) = \hat{\hat{a}}(0) = 1$, it is easy to directly verify that f is an admissible function in $L_{2,\infty,0}(\mathbb{R})$ with respect to both \hat{a} and $\hat{\hat{a}}$. Define

$$\widehat{f}_n(\xi) := \hat{f}(2^{-n}\xi) \prod_{j=1}^n \hat{a}(2^{-j}\xi) \quad \text{and} \quad \widehat{\tilde{f}}_n(\xi) := \hat{f}(2^{-n}\xi) \prod_{j=1}^n \hat{\hat{a}}(2^{-j}\xi), \quad n \in \mathbb{N}.$$

Then by (i) of Theorem 2.1, both $\{f_n\}_{n=1}^\infty$ and $\{\tilde{f}_n\}_{n=1}^\infty$ are Cauchy sequences in $L_{2,\infty,0}(\mathbb{R})$. Note that $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \hat{\phi}(\xi)$ and $\lim_{n \rightarrow \infty} \widehat{\tilde{f}}_n(\xi) = \hat{\hat{\phi}}(\xi)$. So, we must have

$$\lim_{n \rightarrow \infty} \|f_n - \phi\|_{L_{2,\infty,0}(\mathbb{R})} = \lim_{n \rightarrow \infty} \|\tilde{f}_n - \tilde{\phi}\|_{L_{2,\infty,0}(\mathbb{R})} = 0.$$

Since $[\hat{f}, \hat{f}] = 1$ and the discrete biorthogonality relation in (3.4) holds, it is easy to show by induction that $[\widehat{f}_n, \widehat{f}_n] = 1$ for all $n \in \mathbb{N}$. Consequently, we must have $[\hat{\phi}, \hat{\phi}] = 1$. Therefore, (ii) \Rightarrow (i).

If $\nu_2(\hat{a}) > 0$, $\nu_2(\hat{\hat{a}}) > 0$, and (3.6) holds, then all the conditions in (ii) hold. So, $[\hat{\phi}, \hat{\phi}] = 1$. Now it follows from (3.6) that $[\hat{\phi}, \hat{\psi}] = 0$, $[\hat{\psi}, \hat{\phi}] = 0$, and $[\hat{\psi}, \hat{\psi}] = 1$. By a

standard argument on MRA [4, 6], we deduce that (3.1) holds. Since $\hat{\phi}(0) = \hat{\tilde{\phi}}(0) = 1$ and both $\hat{\phi}$ and $\hat{\tilde{\phi}}$ are continuous, by the standard argument on MRA, we see that both $\{\psi_{j,k} : j, k \in \mathbb{Z}\}$ and $\{\tilde{\psi}_{j,k} : j, k \in \mathbb{Z}\}$ are dense in $L_2(\mathbb{R})$. To show that both ψ and $\tilde{\psi}$ generate Riesz wavelet bases in $L_2(\mathbb{R})$, we need to show that ψ and $\tilde{\psi}$ satisfy (1.11). By the biorthogonality relation in (3.1), it suffices to show that the right-hand inequality in (1.11) holds for both ψ and $\tilde{\psi}$ [4], which is equivalent to showing that there exists a positive constant C such that

$$(3.7) \quad \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} [|\langle f, \psi_{j,k} \rangle|^2 + |\langle f, \tilde{\psi}_{j,k} \rangle|^2] \leq C \|f\|^2 \quad \forall f \in L_2(\mathbb{R}).$$

Since $\nu_2(\hat{a}) > 0$ and $\nu_2(\hat{\tilde{a}}) > 0$, by Theorem 2.1, we have $[\hat{\phi}, \hat{\phi}]_\nu, [\hat{\tilde{\phi}}, \hat{\tilde{\phi}}]_\nu \in L_\infty(\mathbb{T})$ for $0 < \nu < \min(\nu_2(\hat{a}), \nu_2(\hat{\tilde{a}}))$. Using Fourier transform and the Parseval's identity, by $\hat{\psi}(\xi) = \hat{b}(\xi/2)\hat{\phi}(\xi/2)$, we have

$$\begin{aligned} 2\pi \sum_{k \in \mathbb{Z}} |\langle f, \psi_{j,k} \rangle|^2 &= 2^j \int_{-\pi}^{\pi} |[\hat{f}(2^j \cdot), \hat{\psi}](\xi)|^2 d\xi \\ &\leq 2^{j+1} \int_{-\pi}^{\pi} |\hat{b}(\xi)|^2 |[\hat{f}(2^{j+1} \cdot), \hat{\phi}](\xi)|^2 d\xi \\ &\leq 2^{j+1} \int_{-\pi}^{\pi} |\hat{b}(\xi)|^2 |[\hat{f}(2^{j+1} \cdot), \hat{f}(2^{j+1} \cdot)]_{-\nu}(\xi)|^2 |[\hat{\phi}, \hat{\phi}]_\nu(\xi)|^2 d\xi \\ &\leq \|[\hat{\phi}, \hat{\phi}]_\nu\|_{L_\infty(\mathbb{T})} 2^{j+1} \int_{-\pi}^{\pi} |\hat{b}(\xi)|^2 |[\hat{f}(2^{j+1} \cdot), \hat{f}(2^{j+1} \cdot)]_{-\nu}(\xi)|^2 d\xi \\ &= \|[\hat{\phi}, \hat{\phi}]_\nu\|_{L_\infty(\mathbb{T})} \int_{\mathbb{R}} \frac{|\hat{b}(2^{-j-1}\xi)|^2}{(1 + |2^{-j-1}\xi|^2)^\nu} |\hat{f}(\xi)|^2 d\xi. \end{aligned}$$

Note that $\hat{b}(0) = \hat{\tilde{b}}(0) = 0$ and $\hat{b}, \hat{\tilde{b}} \in C^\beta(\mathbb{T})$ with $\beta > 0$. Consequently, we see that (3.7) holds with

$$\begin{aligned} C &= \|[\hat{\phi}, \hat{\phi}]_\nu\|_{L_\infty(\mathbb{T})} \left\| \sum_{j \in \mathbb{Z}} \frac{|\hat{b}(2^j \cdot)|^2}{(1 + |2^j \cdot|^2)^\nu} \right\|_{L_\infty(\mathbb{R})} \\ &\quad + \|[\hat{\tilde{\phi}}, \hat{\tilde{\phi}}]_\nu\|_{L_\infty(\mathbb{T})} \left\| \sum_{j \in \mathbb{Z}} \frac{|\hat{\tilde{b}}(2^j \cdot)|^2}{(1 + |2^j \cdot|^2)^\nu} \right\|_{L_\infty(\mathbb{R})} < \infty. \end{aligned}$$

Therefore, $(\psi, \tilde{\psi})$ generates a pair of biorthogonal wavelet bases in $L_2(\mathbb{R})$. \square

As an application of Theorems 2.1 and 3.1, we characterize MRA Riesz wavelet bases in $L_2(\mathbb{R})$ in the following result, which improves and generalizes [14, Theorem 6] and [16, Theorem 1.1] by taking a different approach.

THEOREM 3.2. *Let $\hat{a}, \hat{b} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$. Define ϕ and ψ by*

$$(3.8) \quad \hat{\phi}(\xi) := \prod_{j=1}^{\infty} \hat{a}(2^{-j}\xi) \quad \text{and} \quad \hat{\psi}(\xi) := \hat{b}(\xi/2)\hat{\phi}(\xi/2), \quad \xi \in \mathbb{R}.$$

Then the shifts of ϕ are stable in $L_2(\mathbb{R})$ and ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$ if and only if

- (i) $\hat{b}(0) = 0$ and $d(\xi) := \hat{a}(\xi)\hat{b}(\xi + \pi) - \hat{a}(\xi + \pi)\hat{b}(\xi) \neq 0$ for all $\xi \in \mathbb{R}$,
- (ii) $\nu_2(\hat{a}) > 0$ and $\nu_2(\hat{a}) > 0$, where $\hat{a}(\xi) := \overline{\hat{b}(\xi + \pi)/d(\xi)}$.

In the case that \hat{a} is a 2π -periodic trigonometric polynomial, then the shifts of $\phi \in L_2(\mathbb{R})$ must be stable in $L_2(\mathbb{R})$ if ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$.

Proof. Suppose that (i) and (ii) hold. Define $\hat{b}(\xi) := -\overline{\hat{a}(\xi + \pi)/d(\xi)}$. Since $\hat{a}, \hat{b} \in C^\beta(\mathbb{T})$ and $d(\xi) \neq 0$ for all $\xi \in \mathbb{R}$, it is evident that $\hat{a}, \hat{b} \in C^\beta(\mathbb{T})$. By $\nu_2(\hat{a}) > 0$ and Theorem 2.1, $\hat{a}(\pi) = 0$. Since $\hat{a}(0) = 1$ and $\hat{b}(0) = 0$, we must have $\hat{a}(0) = 1$ and $\hat{a}(\pi) = 0$. Moreover, it is easy to check that (3.6) holds. Define $\tilde{\phi}$ and $\tilde{\psi}$ as in (3.2) and (3.5). Now it follows from Theorem 3.1 that $(\psi, \tilde{\psi})$ generates a pair of biorthogonal wavelet bases in $L_2(\mathbb{R})$. In particular, we conclude that ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$.

Conversely, suppose that the shifts of ϕ are stable in $L_2(\mathbb{R})$ and ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$. Since the shifts of ϕ are stable in $L_2(\mathbb{R})$, by Corollary 2.2, we have $\nu_2(\hat{a}) > 0$ and $\hat{a}(\pi) = 0$. Since $\hat{\psi}$ is continuous and ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$, we must have $\hat{\psi}(0) = 0$, and therefore $\hat{b}(0) = 0$ by $\hat{\phi}(0) = 1$. So, $\hat{a}(0) = \overline{\hat{b}(\pi)/d(0)} = 1$. By [14, Lemma 1], we see that (i) must hold and there exists a function $\tilde{\phi} \in L_2(\mathbb{R})$ such that the shifts of $\tilde{\phi}$ are stable in $L_2(\mathbb{R})$ and $\hat{\phi}(2\xi) = \hat{a}(\xi)\hat{\phi}(\xi)$. Since $\hat{a} \in C^\beta(\mathbb{T})$ and $\hat{a}(0) = 1$, by Corollary 2.2, we have $\nu_2(\hat{a}) > 0$. Therefore, both (i) and (ii) hold.

Suppose that \hat{a} is a 2π -periodic trigonometric polynomial and the shifts of $\phi \in L_2(\mathbb{R})$ are not stable in $L_2(\mathbb{R})$. Since ϕ is compactly supported, by [19, Theorem 5.3], there exists a compactly supported refinable function $\eta \in L_2(\mathbb{R})$ with stable shifts in $L_2(\mathbb{R})$ such that $\hat{\eta}(2\xi) = \hat{c}(\xi)\hat{\eta}(\xi)$ for some 2π -periodic trigonometric polynomial \hat{c} , and $\hat{\phi}(\xi) = \theta(\xi)\hat{\eta}(\xi)$ for some 2π -periodic trigonometric polynomial θ . Note that $[\hat{\phi}, \hat{\phi}](\xi) = |\theta(\xi)|^2[\hat{\eta}, \hat{\eta}](\xi)$ and $\hat{a}(\xi) = \theta(2\xi)\hat{c}(\xi)/\theta(\xi)$. Since the shifts of the compactly supported function ϕ are not stable in $L_2(\mathbb{R})$, there is $\xi_0 \in \mathbb{R} \setminus [2\pi\mathbb{Z}]$ such that $\theta(\xi_0) = 0$.

Since $\hat{\psi}(2\xi) = \hat{b}(\xi)\hat{\phi}(\xi) = \hat{b}(\xi)\theta(\xi)\hat{\eta}(\xi)$ and the shifts of η are stable in $L_2(\mathbb{R})$, if ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$, by what has been proved, then we must have

$$(3.9) \quad 0 \neq \hat{d}(\xi) := \hat{c}(\xi)\hat{b}(\xi + \pi)\theta(\xi + \pi) - \hat{c}(\xi + \pi)\hat{b}(\xi)\theta(\xi) = \frac{\theta(\xi)\theta(\xi + \pi)}{\theta(2\xi)}d(\xi) \\ \forall \xi \in \mathbb{R}$$

and $\nu_2(\hat{c}) > 0$, $\nu_2(\hat{a}) > 0$, where $\hat{a}(\xi) := \overline{\hat{b}(\xi + \pi)\theta(\xi + \pi)/\hat{d}(\xi)}$.

By (3.9), we conclude that if $\xi \in \mathbb{R}$ is a zero of θ , then 2ξ must also be a zero of θ . Since $\theta(\xi_0) = 0$, we now see that $\theta(2^j\xi_0) = 0$ for all $j \in \mathbb{N} \cup \{0\}$. By the definition of $\hat{d}(\xi)$ in (3.9), for all $j \in \mathbb{N} \cup \{0\}$, we have

$$\hat{d}(2^j\xi_0) := \hat{c}(2^j\xi_0)\hat{b}(2^j\xi_0 + \pi)\theta(2^j\xi_0 + \pi) - \hat{c}(2^j\xi_0 + \pi)\hat{b}(2^j\xi_0)\theta(2^j\xi_0) \\ = \hat{c}(2^j\xi_0)\hat{b}(2^j\xi_0 + \pi)\theta(2^j\xi_0 + \pi).$$

Now by the definition of \hat{a} , we deduce that

$$\hat{a}(2^j \xi_0) = \frac{\overline{\hat{b}(2^j \xi_0 + \pi) \theta(2^j \xi_0 + \pi)}}{\hat{c}(2^j \xi_0) \overline{\hat{b}(2^j \xi_0 + \pi) \theta(2^j \xi_0 + \pi)}} = \frac{1}{\hat{c}(2^j \xi_0)} \quad \forall j \in \mathbb{N} \cup \{0\},$$

from which we see that

$$(3.10) \quad \prod_{j=0}^n \hat{a}(2^j \xi_0) = \frac{1}{\prod_{j=0}^n \hat{c}(2^j \xi_0)} \quad \forall n \in \mathbb{N}.$$

Since the shifts of η are stable, there exists $k_0 \in \mathbb{Z}$ such that $\hat{\eta}(\xi_0 + 2\pi k_0) \neq 0$. Since $\xi_0 \notin 2\pi\mathbb{Z}$, we have $\xi_0 + 2\pi k_0 \neq 0$, and therefore $\lim_{n \rightarrow \infty} \hat{\eta}(2^n(\xi_0 + 2\pi k_0)) = 0$ by $\eta \in L_1(\mathbb{R}) \cap L_2(\mathbb{R})$. Now by (3.10), we have

$$\begin{aligned} \left| \prod_{j=0}^n \hat{a}(2^j \xi_0) \right| &= \frac{1}{\left| \prod_{j=0}^n \hat{c}(2^j \xi_0) \right|} = \frac{1}{\left| \prod_{j=0}^n \hat{c}(2^j(\xi_0 + 2\pi k_0)) \right|} \\ &= \frac{|\hat{\eta}(\xi_0 + 2\pi k_0)|}{|\hat{\eta}(2^{n+1}(\xi_0 + 2\pi k_0))|} \rightarrow \infty \end{aligned}$$

as $n \rightarrow \infty$, which is a contradiction to [14, Lemma 1] (also see (5.19)), since $\hat{a} \in C^\beta(\mathbb{T})$. Therefore, the shifts of $\phi \in L_2(\mathbb{R})$ must be stable in $L_2(\mathbb{R})$. \square

To illustrate the results in this section, we consider MRA Riesz wavelet bases in $L_2(\mathbb{R})$ using the masks in (1.12). The following result generalizes [17, Theorem 2.2] for the case of B -splines.

THEOREM 3.3. *Let $\hat{a} = \widehat{a_{\beta_1, \beta_2, \beta_3}}$ or $\hat{a} = |\widehat{a_{\beta_1, \beta_2, \beta_3}}|$, where the masks $\widehat{a_{\beta_1, \beta_2, \beta_3}}$ are defined in (1.12). Let ϕ denote the standard refinable function associated with mask \hat{a} in (1.2) and define a wavelet function ψ by*

$$(3.11) \quad \hat{\psi}(2\xi) := e^{-i\xi \overline{\hat{a}(\xi + \pi)}} \hat{\phi}(\xi), \quad \xi \in \mathbb{R}.$$

Then the shifts of ϕ are stable in $L_2(\mathbb{R})$ and the wavelet function ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$, provided that $\beta_1, \beta_2, \beta_3 > 0$ satisfy

$$(3.12) \quad \beta_1 > 1/4 + |\beta_2 - 1|/2, \quad (\beta_2 - 1)\beta_3 > -1/2 \quad \text{or} \quad \beta_1 \geq 1, \quad (\beta_2 - 1)\beta_3 \leq -1/2.$$

Proof. Let $\hat{b}(\xi) := e^{-i\xi \overline{\hat{a}(\xi + \pi)}}$. Then $\hat{\psi}(2\xi) = \hat{b}(\xi) \hat{\phi}(\xi)$. Clearly, $\hat{a}, \hat{b} \in C^\beta(\mathbb{T})$ for some $\beta > 0$ and $\hat{b}(0) = 0$. By calculation, we have

$$d(\xi) = e^{-i(\xi + \pi)} (|\hat{a}(\xi)|^2 + |\hat{a}(\xi + \pi)|^2) = e^{-i(\xi + \pi)} \frac{|\cos(\xi/2)|^{4\beta_1} + |\sin(\xi/2)|^{4\beta_1}}{(|\cos(\xi/2)|^{2\beta_2} + |\sin(\xi/2)|^{2\beta_2})^{2\beta_3}}.$$

So, $d(\xi) \neq 0$ for all $\xi \in \mathbb{R}$, and (i) of Theorem 3.2 holds.

By calculation, we have

$$|\hat{a}(\xi)| = |\overline{\hat{b}(\xi + \pi)/d(\xi)}| = 2^{-2\beta_1} |1 + e^{-i\xi}|^{2\beta_1} \hat{c}(\xi),$$

where

$$\hat{c}(\xi) := \frac{(|\cos(\xi/2)|^{2\beta_2} + |\sin(\xi/2)|^{2\beta_2})^{\beta_3}}{|\cos(\xi/2)|^{4\beta_1} + |\sin(\xi/2)|^{4\beta_1}}.$$

Now it follows from the inequalities in (2.4) that

$$0 < \hat{c}(\xi) \leq c_{\beta_1, \beta_2, \beta_3} := \frac{\max(1, 2^{(1-\beta_2)\beta_3})}{\min(1, 2^{1-2\beta_1})} \quad \forall \xi \in \mathbb{R}.$$

If (3.12) holds, then (2.5) must be true, and therefore, by Example 2.4, we have $\nu_2(\hat{a}) > 0$. By Theorem 4.1 and Lemma 4.3, if (3.12) holds, then we have

$$\rho(\hat{a}) = \rho_{4\beta_1}(\hat{a}, \infty) = \rho_0(2^{-2\beta_1}\hat{c}, \infty) \leq \rho_0(2^{-2\beta_1}c_{\beta_1, \beta_2, \beta_3}, \infty) = 2^{1-4\beta_1}|c_{\beta_1, \beta_2, \beta_3}|^2 < 1,$$

since the last inequality combined with (2.5) is equivalent to (3.12). Now by Theorem 3.2, the shifts of ϕ are stable in $L_2(\mathbb{R})$, and the function ψ generates a Riesz wavelet basis in $L_2(\mathbb{R})$. \square

For the fractional splines in [27], we have $\beta_2 = 1$, and the condition in (3.12) becomes $\beta_1 > 1/4$. As discussed in Example 2.4, the standard refinable function associated with mask \hat{a} and $\beta_2 = 1$ does not belong to $L_2(\mathbb{R})$ for $0 < \beta_1 \leq 1/4$. So, Theorem 3.3 is sharp for all the fractional splines in [27].

For the classical Butterworth filters, we have $\beta_3 = 1$ and $\beta_1 = \beta_2 \in \mathbb{N}$. Now the condition in (3.12) becomes $\beta_1 > 1/2$. Therefore, Theorem 3.3 holds for all the classical Butterworth filters in [25].

4. Some properties and estimate of the quantity $\nu_2(\hat{a})$. In this section, we shall investigate some properties of the quantity $\nu_2(\hat{a})$ in (1.7) and discuss how to estimate the quantity $\nu_2(\hat{a})$. Some results in this section will be needed in our study of refinable functions, cascade algorithms, and wavelets with Hölder continuous masks. The results in this section for the general case of Lebesgue measurable masks are also of interest in their own right and may be useful elsewhere.

The following result generalizes a well-known result for a univariate 2π -periodic trigonometric polynomial \hat{a} and a positive integer τ in the wavelet literature. The proof of the following result for the general case of Hölder continuous masks is non-trivial and will be presented in section 7.

THEOREM 4.1. *Let \hat{a} be a 2π -periodic measurable function such that $|\hat{a}|^2 \in C^\beta(\mathbb{T})$ with $|\hat{a}(0)|^2 \neq 0$ and $\beta > 0$. If $|\hat{a}(\xi)|^2 = |1 + e^{-i\xi}|^{2\tau} |\hat{A}(\xi)|^2$ a.e. $\xi \in \mathbb{R}$ for some $\tau \geq 0$ such that $\hat{A} \in L_\infty(\mathbb{T})$, then*

$$\begin{aligned} (4.1) \quad \rho_{2\tau}(\hat{a}, \infty) &= \inf_{n \rightarrow \infty} \left\| \frac{T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})}{|\sin(\cdot/2)|^{2\tau}} \right\|_{L_\infty(\mathbb{T})}^{1/n} \\ &= \lim_{n \rightarrow \infty} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} = \inf_{n \in \mathbb{N}} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n}. \end{aligned}$$

As in (1.7), we define a similar quantity as follows:

$$(4.2) \quad \nu_2(\hat{a}, p) := -[\log_2 \rho(\hat{a}, p)]/2, \quad 1 \leq p \leq \infty,$$

where the quantity $\rho(\hat{a}, p)$, using $\rho_\tau(\hat{a}, p)$ in (1.6), is defined to be

$$(4.3) \quad \rho(\hat{a}, p) := \inf\{\rho_\tau(\hat{a}, p) : |\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T}) \text{ and } \tau \geq 0\}.$$

Clearly, $\nu_2(\hat{a}) = \nu_2(\hat{a}, \infty)$ and $\rho(\hat{a}) = \rho(\hat{a}, \infty)$. For a 2π -periodic trigonometric polynomial \hat{a} , we can write $\hat{a}(\xi) = (1 + e^{-i\xi})^\tau \hat{A}(\xi)$ for some nonnegative integer τ and some 2π -periodic trigonometric polynomial \hat{A} with $\hat{A}(\pi) \neq 0$. Write $|\hat{A}(\xi)|^2 = \sum_{k=-K}^K c_k e^{-ik\xi}$. It is known [7, 13, 15] that $\nu_2(\hat{a}, p) = \nu_2(\hat{a}) = -1/2 - \log_2 \sqrt{\rho}$, where ρ is the spectral radius of the square matrix $(c_{2j-k})_{-K \leq j, k \leq K}$.

In the following, we shall investigate the mutual relations among the quantities $\nu_2(\hat{a}, p)$.

PROPOSITION 4.2. *The following statements hold:*

- (1) For $0 < \tau < 1$, $|\sin(\cdot/2)|^\tau \in C^\tau(\mathbb{T})$ and $|\cos(\cdot/2)|^\tau \in C^\tau(\mathbb{T})$.
- (2) For a 2π -periodic measurable function \hat{a} and $0 \leq \tau_1 \leq \tau_2$,

$$(4.4) \quad \rho_{\tau_2}(\hat{a}, p) \leq \rho_{\tau_1}(\hat{a}, p) \leq \rho_{\tau_1}(\hat{a}, q) \quad \text{and} \quad \nu_2(\hat{a}, q) \leq \nu_2(\hat{a}, p) \\ \forall 1 \leq p \leq q \leq \infty.$$

- (3) If $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$, then the condition

$$(4.5) \quad \liminf_{n \rightarrow \infty} \|T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)\|_{L_1(\mathbb{T})} = 0 \quad \text{for some } \tau \geq 0$$

implies $\hat{a}(\pi) = 0$. In particular, (4.5) holds if $\nu_2(\hat{a}, p) > 0$ for some $1 \leq p \leq \infty$.

Proof. It is easy to prove that $1 - x^\tau \leq (1 - x)^\tau$ for all $0 < \tau \leq 1$ and $0 \leq x \leq 1$. Consequently, we have

$$\left| |\sin(x/2)|^\tau - |\sin(y/2)|^\tau \right| \leq |\sin(x/2) - \sin(y/2)|^\tau \leq |x - y|^\tau.$$

So, $|\sin(\cdot/2)|^\tau \in C^\tau(\mathbb{T})$, and it follows that $|\cos(\cdot/2)|^\tau \in C^\tau(\mathbb{T})$. So, (1) holds.

Since $0 \leq \tau_1 \leq \tau_2$, it is evident that $|\sin(\xi/2)|^{\tau_2} \leq |\sin(\xi/2)|^{\tau_1}$ for all $\xi \in \mathbb{R}$. Therefore,

$$T_{\hat{a}}^n(|\sin(\cdot/2)|^{\tau_2}) \leq T_{\hat{a}}^n(|\sin(\cdot/2)|^{\tau_1}).$$

Now the claim in (4.4) follows directly from the above inequality and the fact that $\|\cdot\|_{L_p(\mathbb{T})} \leq \|\cdot\|_{L_q(\mathbb{T})}$ for all $1 \leq p \leq q \leq \infty$. So, (2) holds.

To prove (3), we denote $\Phi(\xi) := \prod_{j=1}^{\infty} |\hat{a}(2^{-j}\xi)|^2$. Since $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\beta > 0$, Φ is well defined and is continuous with $\Phi(0) = 1$. Suppose that $\hat{a}(\pi) \neq 0$. Then there exist $0 < \varepsilon < \pi/2$ and a positive constant C such that $|\hat{a}(\xi + \pi)|^2 |\sin(\xi/2 + \pi/2)|^\tau \geq C$ and $C \leq \Phi(\xi) \leq 1/C$ for all $\xi \in (-2\varepsilon, 2\varepsilon)$. By the

definition of $T_{\hat{a}}$, we deduce that

$$\begin{aligned}
\int_0^{2\pi} [T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)](\xi) d\xi &= 2^n \int_0^{2\pi} |\sin(\xi/2)|^\tau |\hat{a}(\xi)|^2 |\hat{a}(2\xi)|^2 \cdots |\hat{a}(2^{n-1}\xi)|^2 d\xi \\
&\geq 2^n \int_{\pi-\varepsilon}^{\pi+\varepsilon} |\sin(\xi/2)|^\tau |\hat{a}(\xi)|^2 |\hat{a}(2\xi)|^2 \cdots |\hat{a}(2^{n-1}\xi)|^2 d\xi \\
&\geq 2^n C \int_{\pi-\varepsilon}^{\pi+\varepsilon} |\hat{a}(2\xi)|^2 \cdots |\hat{a}(2^{n-1}\xi)|^2 d\xi \\
&= 2^{n-1} C \int_{-2\varepsilon}^{2\varepsilon} |\hat{a}(\xi)|^2 \cdots |\hat{a}(2^{n-2}\xi)|^2 d\xi \\
&= 2^{n-1} C \int_{-2\varepsilon}^{2\varepsilon} \frac{\Phi(2^{n-1}\xi)}{\Phi(\xi)} d\xi \\
&\geq 2^{n-1} C^2 \int_{-2\varepsilon}^{2\varepsilon} \Phi(2^{n-1}\xi) d\xi \\
&= C^2 \int_{-2^n\varepsilon}^{2^n\varepsilon} \Phi(\xi) d\xi.
\end{aligned}$$

Hence,

$$\liminf_{n \rightarrow \infty} \|T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)\|_{L_1(\mathbb{T})} \geq (2\pi)^{-1} C^2 \int_{\mathbb{R}} \Phi(\xi) d\xi > 0,$$

since $\Phi \geq 0$ is continuous and $\Phi(0) = 1$. This is a contradiction to our assumption in (4.5). So, we must have $\hat{a}(\pi) = 0$. If $\nu_2(\hat{a}, p) > 0$, then $\nu_2(\hat{a}, 1) \geq \nu_2(\hat{a}, p) > 0$, and therefore (4.5) holds. \square

The following result will be needed later in this section.

LEMMA 4.3. *Let \hat{a} and \hat{c} be 2π -periodic measurable functions such that $|\hat{a}(\xi)| \leq |\hat{c}(\xi)|$ for almost every $\xi \in \mathbb{R}$. Then*

$$(4.6) \quad \rho_\tau(\hat{a}, p) \leq \rho_\tau(\hat{c}, p) \quad \text{and} \quad \nu_2(\hat{c}, p) \leq \nu_2(\hat{a}, p) \quad \forall 1 \leq p \leq \infty, \tau \in \mathbb{R}.$$

Proof. To prove (4.6), since $|\hat{a}(\xi)|^2 \leq |\hat{c}(\xi)|^2$, it is obvious that

$$0 \leq [T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)](\xi) \leq [T_{\hat{c}}^n(|\sin(\cdot/2)|^\tau)](\xi) \quad \text{a.e. } \xi \in \mathbb{R}.$$

Therefore, $\|T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})} \leq \|T_{\hat{c}}^n(|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})}$, which implies the first part of (4.6).

If $|\hat{c}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$ for some τ , then we also have

$$|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$$

since $|\hat{a}(\xi)| \leq |\hat{c}(\xi)|$. Now by the definition of $\nu_2(\hat{a}, p)$ in (4.2) and the first part of (4.6), it is easy to see that $\nu_2(\hat{c}, p) \leq \nu_2(\hat{a}, p)$. \square

For a particular family of masks, the following result reveals the mutual relations among the quantities $\nu_2(\hat{a}, p)$ for different $1 \leq p \leq \infty$.

LEMMA 4.4. *Let \hat{a} be a 2π -periodic measurable function such that $|\hat{a}(\xi)| = |1 + e^{-i\xi}|^\tau |\hat{A}(\xi)|$ for some $\tau \geq 0$ and some 2π -periodic trigonometric polynomial \hat{A} with $\hat{A}(0) \neq 0$. Then*

$$(4.7) \quad \rho_{2\tau}(\hat{a}, p) = \lim_{n \rightarrow \infty} \|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L_p(\mathbb{T})}^{1/n} = \rho_0(\hat{A}, \infty) = \inf_{n \in \mathbb{N}} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} \\ \forall 1 \leq p \leq \infty.$$

In particular, $\nu_2(\hat{a}, p) = \nu_2(\hat{a}) = \nu_2(\hat{A})$ for all $1 \leq p \leq \infty$. That is, $\nu_2(\hat{a}, p)$ is independent of p .

Proof. Let N be an integer such that $N \geq \tau$. Define $\hat{c}(\xi) := (1 + e^{-i\xi})^N \hat{A}(\xi)$. By calculation, it follows from $|\hat{a}(\xi)|^2 = 2^{2\tau} |\cos(\xi/2)|^{2\tau} |\hat{A}(\xi)|^2$ that

$$\begin{aligned} [T_{\hat{a}}(f(\cdot)|\sin(\cdot/2)|^{2\tau})](\xi) &= |\hat{a}(\xi/2)|^2 |\sin(\xi/4)|^{2\tau} f(\xi/2) \\ &\quad + |\hat{a}(\xi/2 + \pi)|^2 |\sin(\xi/4 + \pi/2)|^{2\tau} f(\xi/2 + \pi) \\ &= |\sin(\xi/2)|^{2\tau} [|\hat{A}(\xi/2)|^2 f(\xi/2) + |\hat{A}(\xi/2 + \pi)|^2 f(\xi/2 + \pi)] \\ &= |\sin(\xi/2)|^{2\tau} [T_{\hat{A}} f](\xi). \end{aligned}$$

That is, when $|\hat{a}(\xi)|^2 = 2^{2\tau} |\cos(\xi/2)|^{2\tau} |\hat{A}(\xi)|^2$, for any 2π -periodic function f , by induction, we have

$$(4.8) \quad [T_{\hat{a}}^n(f(\cdot)|\sin(\cdot/2)|^{2\tau})](\xi) = |\sin(\xi/2)|^{2\tau} [T_{\hat{A}}^n f](\xi), \quad n \in \mathbb{N}.$$

Setting $f = 1$ in (4.8), we have

$$(4.9) \quad [T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})](\xi) = |\sin(\xi/2)|^{2\tau} [T_{\hat{A}}^n 1](\xi) \leq [T_{\hat{A}}^n 1](\xi).$$

Note that $2N - 2\tau \geq 0$. Similarly, by $|\hat{c}(\xi)|^2 = 2^{2N-2\tau} |\cos(\xi/2)|^{2N-2\tau} |\hat{a}(\xi)|^2$, setting $f(\xi) = |\sin(\xi/2)|^{2\tau}$ and replacing τ by $N - \tau$ in (4.8), we have

$$(4.10) \quad [T_{\hat{c}}^n(|\sin(\cdot/2)|^{2N})](\xi) = |\sin(\xi/2)|^{2N-2\tau} [T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})](\xi) \\ \leq [T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})](\xi).$$

Thus, it follows from (4.9) and (4.10) that

$$(4.11) \quad \|T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\|_{L_p(\mathbb{T})} \leq \|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L_p(\mathbb{T})} \leq \|T_{\hat{A}}^n 1\|_{L_p(\mathbb{T})} \\ \leq \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}.$$

Since both $|\hat{c}|^2$ and $\sin^{2N}(\cdot/2)$ are 2π -periodic trigonometric polynomials, by induction, it is known [4, 7, 15] that $\{T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\}_{n=1}^\infty$ spans a finite dimensional space and in fact the degrees of all trigonometric polynomials $T_{\hat{c}}^n(\sin^{2N}(\cdot/2))$ are uniformly bounded. Therefore, there exists a positive constant C , independent of all n , such that

$$C \|T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\|_{L_\infty(\mathbb{T})} \leq \|T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\|_{L_p(\mathbb{T})} \quad \forall 1 \leq p \leq \infty, n \in \mathbb{N}.$$

Hence, we conclude from the above inequality and (4.11) that for $1 \leq p \leq \infty$,

$$(4.12) \quad C^{1/n} \|T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\|_{L_\infty(\mathbb{T})}^{1/n} \leq \|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L_p(\mathbb{T})}^{1/n} \leq \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} \\ \forall n \in \mathbb{N}.$$

Since $\hat{c}(\xi) = (1 + e^{-i\xi})^N \hat{A}(\xi)$ and $\hat{A}(0) \neq 0$, by Theorem 4.1, we have

$$\rho_{2N}(\hat{c}, \infty) = \lim_{n \rightarrow \infty} \|T_{\hat{c}}^n(\sin^{2N}(\cdot/2))\|_{L_\infty(\mathbb{T})}^{1/n} = \lim_{n \rightarrow \infty} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} = \inf_{n \in \mathbb{N}} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n}.$$

Now (4.7) follows directly from (4.12).

Since \hat{A} is a 2π -periodic trigonometric polynomial, we can write $\hat{A}(\xi) = (1 + e^{-i\xi})^k \hat{B}(\xi)$ for some nonnegative integer k and some 2π -periodic trigonometric polynomial \hat{B} such that $\hat{B}(\pi) \neq 0$. So, $|\hat{a}(\xi)| = |1 + e^{-i\xi}|^{k+\tau} |\hat{B}(\xi)|$. Now by the definition of $\nu_2(\hat{a}, p)$ and Theorem 4.1, it follows from (4.7) that

$$\begin{aligned} \nu_2(\hat{a}, p) &:= -[\log_2 \rho_{2k+2\tau}(\hat{a}, p)]/2 = -[\log_2 \rho_0(\hat{B}, \infty)]/2 \\ &= -[\log_2 \rho_{2k}(\hat{A}, \infty)]/2 = \nu_2(\hat{A}). \end{aligned}$$

Therefore, $\nu_2(\hat{a}, p) = \nu_2(\hat{A})$ for all $1 \leq p \leq \infty$. \square

In the following, we shall discuss how to approximate the quantity $\rho_0(\hat{A}, \infty)$.

PROPOSITION 4.5. *Let $\hat{A}, \widehat{A}_j \in L_\infty(\mathbb{T})$, $j \in \mathbb{N}$ such that $\lim_{j \rightarrow \infty} \|\widehat{A}_j - \hat{A}\|_{L_\infty(\mathbb{T})} = 0$. Then*

$$(4.13) \quad \limsup_{j \rightarrow \infty} \rho_0(\widehat{A}_j, \infty) \leq \rho_0(\hat{A}, \infty).$$

If in addition $|\hat{A}(\xi)| \leq |\widehat{A}_j(\xi)|$ for almost every $\xi \in \mathbb{R}$ and for all $j \in \mathbb{N}$, then

$$(4.14) \quad \lim_{j \rightarrow \infty} \rho_0(\widehat{A}_j, \infty) = \rho_0(\hat{A}, \infty).$$

Proof. By Proposition 7.1, $\rho_0(\hat{A}, \infty) = \inf_{n \in \mathbb{N}} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}$. Therefore, for any $\varepsilon > 0$, there exists a positive integer N such that $\|T_{\hat{A}}^N 1\|_{L_\infty(\mathbb{T})}^{1/N} < \rho_0(\hat{A}, \infty) + \varepsilon$. Now by $\lim_{j \rightarrow \infty} \|\widehat{A}_j - \hat{A}\|_{L_\infty(\mathbb{T})} = 0$, there exists a positive integer J such that

$$\|T_{\widehat{A}_j}^N 1\|_{L_\infty(\mathbb{T})}^{1/N} < \rho_0(\hat{A}, \infty) + \varepsilon \quad \forall j \geq J.$$

Consequently, by Proposition 7.1, we deduce that

$$\rho_0(\widehat{A}_j, \infty) = \inf_{n \in \mathbb{N}} \|T_{\widehat{A}_j}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} \leq \|T_{\widehat{A}_j}^N 1\|_{L_\infty(\mathbb{T})}^{1/N} < \rho_0(\hat{A}, \infty) + \varepsilon \quad \forall j \geq J.$$

Hence, we have $\limsup_{j \rightarrow \infty} \rho_0(\widehat{A}_j, \infty) \leq \rho_0(\hat{A}, \infty) + \varepsilon$. Taking $\varepsilon \rightarrow 0$, we conclude that (4.13) holds. If $|\hat{A}| \leq |\widehat{A}_j|$, then by Lemma 4.3, $\rho_0(\hat{A}, \infty) \leq \rho_0(\widehat{A}_j, \infty)$ for all $j \in \mathbb{N}$. Therefore, $\rho_0(\hat{A}, \infty) \leq \liminf_{j \rightarrow \infty} \rho_0(\widehat{A}_j, \infty)$. Now it follows from (4.13) that (4.14) holds. \square

As a consequence of Proposition 4.5, we have the following corollary.

COROLLARY 4.6. *Let $\hat{A} \in C^\beta(\mathbb{T})$ with $\hat{A}(0) \neq 0$, $\hat{A}(\pi) \neq 0$, and $\beta > 0$. Suppose that there is a sequence $\{\widehat{A}_j\}_{j \in \mathbb{N}}$ in $C^\beta(\mathbb{T})$ such that $\lim_{j \rightarrow \infty} \|\widehat{A}_j - \hat{A}\|_{L_\infty(\mathbb{T})} = 0$, and*

$|\hat{A}(\xi)| \leq |\widehat{A}_j(\xi)|$ for all $\xi \in \mathbb{R}$ and $j \in \mathbb{N}$. For $\tau \geq 0$ and a 2π -periodic trigonometric polynomial \hat{c} with $\hat{c}(0) \neq 0$, let $\hat{a}(\xi) := (1 + e^{-i\xi})^\tau \hat{c}(\xi) \hat{A}(\xi)$ and $\hat{a}_j(\xi) := (1 + e^{-i\xi})^\tau \hat{c}(\xi) \widehat{A}_j(\xi)$. Then $\lim_{j \rightarrow \infty} \nu_2(\hat{a}_j) = \nu_2(\hat{a})$ and $\nu_2(\hat{a}_j) \leq \nu_2(\hat{a})$ for all $j \in \mathbb{N}$.

Proof. Write $\hat{c}(\xi) = (1 + e^{-i\xi})^k \hat{B}(\xi)$ for some nonnegative integer k and some 2π -periodic trigonometric polynomial \hat{B} with $\hat{B}(\pi) \neq 0$. Since $\hat{A}(\pi) \neq 0$ and $\lim_{j \rightarrow \infty} \|\widehat{A}_j - \hat{A}\|_{L^\infty(\mathbb{T})} = 0$, by $\widehat{A}_j \in C^\beta(\mathbb{T})$, we have $\widehat{A}_j(\pi) \neq 0$ for sufficiently large j . Now by Theorem 4.1 and Proposition 4.5, we have

$$\begin{aligned} \rho(\hat{a}, \infty) &= \rho_{2k+2\tau}(\hat{a}, \infty) = \rho_0(\hat{A}\hat{B}, \infty) = \lim_{j \rightarrow \infty} \rho_0(\widehat{A}_j \hat{B}, \infty) \\ &= \lim_{j \rightarrow \infty} \rho_{2k+2\tau}(\hat{a}_j, \infty) = \lim_{j \rightarrow \infty} \rho(\hat{a}_j, \infty). \end{aligned}$$

It follows from the definition of $\nu_2(\hat{a})$ that $\nu_2(\hat{a}) = \lim_{j \rightarrow \infty} \nu_2(\hat{a}_j)$. Now by $|\hat{a}(\xi)| \leq |\widehat{a}_j(\xi)|$, it follows from Lemma 4.3 that $\nu_2(\hat{a}_j) \leq \nu_2(\hat{a})$. \square

PROPOSITION 4.7. *Let \hat{a} and \hat{a}_j , $j \in \mathbb{N}$, be 2π -periodic measurable functions such that*

$$(4.15) \quad \lim_{j \rightarrow \infty} \|\widehat{a}_j/\hat{a}\|_{L^\infty(\mathbb{T})} = \lim_{j \rightarrow \infty} \|\hat{a}/\widehat{a}_j\|_{L^\infty(\mathbb{T})} = 1,$$

where by convention $(\widehat{a}_j/\hat{a})(\xi)$ is equal to $\widehat{a}_j(\xi)/\hat{a}(\xi)$ if $\hat{a}(\xi) \neq 0$, 1 if $\hat{a}(\xi) = \widehat{a}_j(\xi) = 0$, or $+\infty$ if $\hat{a}(\xi) = 0$ but $\widehat{a}_j(\xi) \neq 0$. Then

$$(4.16) \quad \lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j, p) = \nu_2(\hat{a}, p) \quad \forall 1 \leq p \leq \infty.$$

Moreover, if $\nu_2(\widehat{a}_j, p) = \nu_2(\widehat{a}_j)$ for all $j \in \mathbb{N}$, then $\nu_2(\hat{a}, p) = \nu_2(\hat{a})$. In particular, if a mask \hat{a} has exponential decay, then $\nu_2(\hat{a}, p) = \nu_2(\hat{a})$ for all $1 \leq p \leq \infty$.

Proof. Denote $\widehat{c}_j = \widehat{a}_j/\hat{a}$. By convention and (4.15), $0 < |\widehat{c}_j(\xi)| < \infty$ for almost every $\xi \in \mathbb{R}$, and it is easy to check that $\widehat{a}_j(\xi) = \widehat{c}_j(\xi)\hat{a}(\xi)$ and $\hat{a}(\xi) = \widehat{a}_j(\xi)/\widehat{c}_j(\xi)$ for almost every $\xi \in \mathbb{R}$. For $\tau \geq 0$, we now have

$$\begin{aligned} \|1/\widehat{c}_j\|_{L^\infty(\mathbb{T})}^{-2n} [T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)](\xi) &\leq [T_{\widehat{a}_j}^n(|\sin(\cdot/2)|^\tau)](\xi) \\ &\leq \|\widehat{c}_j\|_{L^\infty(\mathbb{T})}^{2n} [T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)](\xi). \end{aligned}$$

Consequently, we have

$$\begin{aligned} \|1/\widehat{c}_j\|_{L^\infty(\mathbb{T})}^{-2} \|T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})}^{1/n} &\leq \|T_{\widehat{a}_j}^n(|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})}^{1/n} \\ &\leq \|\widehat{c}_j\|_{L^\infty(\mathbb{T})}^2 \|T_{\hat{a}}^n(|\sin(\cdot/2)|^\tau)\|_{L_p(\mathbb{T})}^{1/n}. \end{aligned}$$

Hence, we deduce that

$$(4.17) \quad \|1/\widehat{c}_j\|_{L^\infty(\mathbb{T})}^{-2} \rho_\tau(\hat{a}, p) \leq \rho_\tau(\widehat{a}_j, p) \leq \|\widehat{c}_j\|_{L^\infty(\mathbb{T})}^2 \rho_\tau(\hat{a}, p) \quad \forall j \in \mathbb{N}, 1 \leq p \leq \infty.$$

By our assumption in (4.15), we have $\lim_{j \rightarrow \infty} \|1/\widehat{c}_j\|_{L^\infty(\mathbb{T})} = \lim_{j \rightarrow \infty} \|\widehat{c}_j\|_{L^\infty(\mathbb{T})} = 1$. Now by the definition of $\nu_2(\hat{a})$ and (4.17), we must have $\lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j, p) = \nu_2(\hat{a}, p)$ for all $1 \leq p \leq \infty$.

If $\nu_2(\widehat{a}_j, p) = \nu_2(\widehat{a}_j)$, then $\nu_2(\hat{a}, p) = \lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j, p) = \lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j) = \nu_2(\hat{a})$.

If \hat{a} has exponential decay, then we can write $\hat{a}(\xi) = \hat{c}(\xi)\hat{A}(\xi)$, where \hat{c} is a 2π -periodic trigonometric polynomial and \hat{A} has exponential decay satisfying $\hat{A}(\xi) \neq 0$

for all $\xi \in \mathbb{R}$. Now it is easy to see that there is a sequence $\{\widehat{A}_j\}_{j=1}^\infty$ of 2π -periodic trigonometric polynomials such that $\lim_{j \rightarrow \infty} \|\widehat{A}_j/\widehat{A}\|_{L_\infty(\mathbb{T})} = \lim_{j \rightarrow \infty} \|\widehat{A}/\widehat{A}_j\|_{L_\infty(\mathbb{T})} = 1$. Taking $\widehat{a}_j(\xi) := \widehat{a}(\xi)\widehat{A}_j(\xi)$, then (4.15) holds, and by Lemma 4.4, $\nu_2(\widehat{a}_j, p) = \nu_2(\widehat{a}_j)$ for all $1 \leq p \leq \infty$ and $j \in \mathbb{N}$. By what has been proved, we have $\nu_2(\widehat{a}, p) = \lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j, p) = \lim_{j \rightarrow \infty} \nu_2(\widehat{a}_j) = \nu_2(\widehat{a})$. \square

In passing, we mention that the quantity defined in [14, equation (2.16)] corresponds to $\nu_2(\widehat{a}, 1)$ in (4.2) of this paper. For a mask with exponential decay, by Proposition 4.7, the quantity $\nu_2(\widehat{a})$ in [14, equation (2.16)] agrees with the one in this paper. However, it is not clear whether $\nu_2(\widehat{a}, p) = \nu_2(\widehat{a})$ for all $1 \leq p \leq \infty$ if $\widehat{a} \in C^\beta(\mathbb{T})$ for some $\beta > 0$.

5. Proof of Theorem 2.1. Since $\widehat{a}(0) = 1$ and $\widehat{a} \in C^\beta(\mathbb{T})$ with $\beta > 0$, we see that $\widehat{\phi}$ in (1.2) is well defined and $\widehat{\phi}$ is continuous. If $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \widehat{f}(2^{-n}\xi) = 1$ for almost every $\xi \in \mathbb{R}$, by $\widehat{\phi}(\xi) = \lim_{n \rightarrow \infty} \widehat{f}_n(\xi)$ for almost every $\xi \in \mathbb{R}$, then we must have $\phi \in L_{2,\infty,0}(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \|f_n - \phi\|_{L_{2,\infty,0}(\mathbb{R})} = 0$.

We shall prove in the order that (i) \Rightarrow (ii) \Rightarrow (iii) \Rightarrow (iv) \Rightarrow (v) \Rightarrow (ii) and (iii) \Rightarrow (i). Note that $\widehat{a} \in C^\beta(\mathbb{T})$ implies $\widehat{a} \in C^{\min(1,\beta)}(\mathbb{T})$. So, without loss of generality, we replace β by $\min(1,\beta)$. That is, we assume $0 < \beta \leq 1$, and the following proof depends only on the fact that $\widehat{a} \in C^\alpha(\mathbb{T})$ for a small number $\alpha > 0$.

We show (i) \Rightarrow (ii) by constructing an admissible function η in $L_{2,\infty,0}(\mathbb{R})$ such that the shifts of η are stable in $L_2(\mathbb{R})$ and $\widehat{\eta}$ is compactly supported; such an admissible initial function η will be used in several places in this proof. Since $\widehat{\phi}$ is a continuous function with $\widehat{\phi}(0) = 1$, there exists $0 < \varepsilon < \pi$ such that $1/2 \leq |\widehat{\phi}(\xi)| \leq 3/2$ for all $|\xi| \leq \varepsilon$. Define a function $\widehat{\eta}$ by

$$(5.1) \quad \widehat{\eta}(\xi) := \begin{cases} \widehat{\phi}(\xi) & \text{if } |\xi| \leq \varepsilon, \\ \widehat{\phi}(-\varepsilon)(3\pi + 2\xi)/(3\pi - 2\varepsilon) & \text{if } -3\pi/2 \leq \xi < -\varepsilon, \\ \widehat{\phi}(\varepsilon)(3\pi - 2\xi)/(3\pi - 2\varepsilon) & \text{if } \varepsilon < \xi \leq 3\pi/2, \\ 0 & \text{otherwise.} \end{cases}$$

Note that $\widehat{\eta}(\xi) - \widehat{a}(\xi/2)\widehat{\eta}(\xi/2) = 0$ for all $|\xi| \leq \varepsilon$. By the assumption $\widehat{a}(\pi) = 0$ in (i) and $\widehat{a} \in C^\beta(\mathbb{T})$, we see that $|\widehat{a}(\cdot + \pi)|^2/|\sin(\cdot/2)|^{2\beta} \in L_\infty(\mathbb{T})$. Since $\widehat{\eta}$ is supported inside $[-3\pi/2, 3\pi/2]$, now we can easily verify that η is an admissible function in $L_{2,\infty,0}(\mathbb{R})$ with respect to \widehat{a} , since the condition in (2.1) holds with $\tau = 2\beta > 0$. By the definition of $\widehat{\eta}$, we have $1/32 \leq [\widehat{\eta}, \widehat{\eta}] \leq 9/2$. Therefore, the shifts of η are stable in $L_2(\mathbb{R})$. Taking $f = \eta$, it follows directly from (i) that (ii) holds.

Now we prove (ii) \Rightarrow (iii) without the condition that the initial function f is admissible. By the definition of f_n in (2.2) and induction, we deduce that

$$(5.2) \quad [\widehat{f}_n, \widehat{f}_n](\xi) = \sum_{k=0}^{2^n-1} \prod_{j=1}^n |\widehat{a}(2^{-j}(\xi + 2\pi k))|^2 [\widehat{f}, \widehat{f}](2^{-n}(\xi + 2\pi k)) = (T_{\widehat{a}}^n[\widehat{f}, \widehat{f}])(\xi).$$

Since the shifts of f are stable in $L_2(\mathbb{R})$ and $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$, there exists a positive constant C_1 such that

$$\|f_n\|_{L_{2,\infty,0}(\mathbb{R})}^2 = \|[\widehat{f}_n, \widehat{f}_n]\|_{L_\infty(\mathbb{T})} \leq C_1$$

for all $n \in \mathbb{N}$ and $1/C_1 \leq [\hat{f}, \hat{f}](\xi) \leq C_1$ for almost every $\xi \in \mathbb{R}$. Now it follows from (5.2) that

$$0 \leq [T_{\hat{a}}^n 1](\xi) \leq C_1 (T_{\hat{a}}^n C_1^{-1})(\xi) \leq C_1 (T_{\hat{a}}^n [\hat{f}, \hat{f}])(\xi) = C_1 [\widehat{f_n}, \widehat{f_n}](\xi) \leq C_1^2.$$

That is, we have

$$(5.3) \quad \|T_{\hat{a}}^n 1\|_{L_\infty(\mathbb{T})} \leq C_1^2 \quad \forall n \in \mathbb{N} \cup \{0\}.$$

Since $\hat{a} \in C^\beta(\mathbb{T})$, we have $|\hat{a}|^2 \in C^\beta(\mathbb{T})$, and there exists a positive constant C_2 such that

$$(5.4) \quad \||\hat{a}|^2 - |\hat{a}(\cdot - h)|^2\|_{L_\infty(\mathbb{T})} \leq C_2 h^\beta \quad \forall h > 0.$$

Now we are going to extend an interesting idea in [23] to show that $T_{\hat{a}}^n g$, $n \in \mathbb{N}$, are equicontinuous if $g \in C^\tau(\mathbb{T})$ for some $\tau > 0$. For $g \in L_\infty(\mathbb{T})$, we denote

$$(5.5) \quad \omega_n(\xi, h) := |[T_{\hat{a}}^n g](\xi) - [T_{\hat{a}}^n g](\xi - h)|, \quad \xi \in \mathbb{R}, h > 0, n \in \mathbb{N} \cup \{0\}.$$

Setting

$$C_3 := \max(C_1^2, 2^{1-\beta} C_1^4 C_2 / (1 - 2^{-\beta})) < \infty,$$

we show that (5.3) and (5.4) imply that

$$(5.6) \quad \|h^{-\beta} \omega_{n+k}(\cdot, h)\|_{L_\infty(\mathbb{T})} \leq C_3 \left(2^{-n\beta} \|(2^{-n}h)^{-\beta} \omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} + \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} \right) \\ \forall h > 0, k, n \in \mathbb{N} \cup \{0\}, g \in L_\infty(\mathbb{T}).$$

By the definition of $T_{\hat{a}}$, we have

$$\begin{aligned} & [T_{\hat{a}}^n g](\xi) - [T_{\hat{a}}^n g](\xi - h) \\ &= |\hat{a}(\xi/2)|^2 ([T_{\hat{a}}^{n-1} g](\xi/2) - [T_{\hat{a}}^{n-1} g](\xi/2 - h/2)) \\ & \quad + |\hat{a}(\xi/2 + \pi)|^2 ([T_{\hat{a}}^{n-1} g](\xi/2 + \pi) - [T_{\hat{a}}^{n-1} g](\xi/2 + \pi - h/2)) \\ & \quad + (|\hat{a}(\xi/2)|^2 - |\hat{a}(\xi/2 - h/2)|^2) [T_{\hat{a}}^{n-1} g](\xi/2 - h/2) \\ & \quad + (|\hat{a}(\xi/2 + \pi)|^2 - |\hat{a}(\xi/2 + \pi - h/2)|^2) [T_{\hat{a}}^{n-1} g](\xi/2 + \pi - h/2). \end{aligned}$$

It follows from (5.4) that

$$\begin{aligned} \omega_n(\xi, h) &\leq |\hat{a}(\xi/2)|^2 \omega_{n-1}(\xi/2, h/2) + |\hat{a}(\xi/2 + \pi)|^2 \omega_{n-1}(\xi/2 + \pi, h/2) \\ & \quad + 2C_2 (h/2)^\beta \|T_{\hat{a}}^{n-1} g\|_{L_\infty(\mathbb{T})} \\ &= [T_{\hat{a}} \omega_{n-1}(\cdot, h/2)](\xi) + 2^{1-\beta} C_2 h^\beta \|T_{\hat{a}}^{n-1} g\|_{L_\infty(\mathbb{T})}. \end{aligned}$$

Consequently, by induction on n , we deduce from the above inequality that for all $k, n \in \mathbb{N} \cup \{0\}$,

$$(5.7) \quad \omega_{n+k}(\xi, h) \leq [T_{\hat{a}}^n \omega_k(\cdot, 2^{-n}h)](\xi) \\ + 2C_2 h^\beta \sum_{j=1}^n 2^{-j\beta} \|T_{\hat{a}}^{n+k-j} g\|_{L_\infty(\mathbb{T})} \|T_{\hat{a}}^{j-1} 1\|_{L_\infty(\mathbb{T})}.$$

In fact, (5.7) clearly holds for $n = 0$ and all $k \in \mathbb{N} \cup \{0\}$. Suppose that (5.7) holds for n and all $k \in \mathbb{N} \cup \{0\}$. Then by induction hypothesis we have

$$\begin{aligned} \omega_{n+1+k}(\xi, h) &\leq [T_{\hat{a}}^n \omega_{k+1}(\cdot, 2^{-n}h)](\xi) \\ &\quad + 2C_2 h^\beta \sum_{j=1}^n 2^{-j\beta} \|T_{\hat{a}}^{n+k+1-j} g\|_{L_\infty(\mathbb{T})} \|T_{\hat{a}}^{j-1} 1\|_{L_\infty(\mathbb{T})}. \end{aligned}$$

Note that we proved in the inequality above (5.7) that

$$\omega_{k+1}(\xi, 2^{-n}h) \leq [T_{\hat{a}} \omega_k(\cdot, 2^{-1-n}h)](\xi) + 2^{1-\beta} C_2 (2^{-n}h)^\beta \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})}.$$

Applying the operator $T_{\hat{a}}^n$ on both sides of the above inequality, we deduce that

$$T_{\hat{a}}^n \omega_{k+1}(\xi, 2^{-n}h) \leq [T_{\hat{a}}^{n+1} \omega_k(\cdot, 2^{-1-n}h)](\xi) + 2^{1-\beta} C_2 (2^{-n}h)^\beta \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} [T_{\hat{a}}^n 1](\xi).$$

Combining all the above inequalities together, we see that (5.7) holds for $n+1$ and all $k \in \mathbb{N} \cup \{0\}$. So, by induction, (5.7) holds for all $k, n \in \mathbb{N} \cup \{0\}$.

By (5.3), we have

$$\|T_{\hat{a}}^{n+k-j} g\|_{L_\infty(\mathbb{T})} \leq \|T_{\hat{a}}^{n-j} 1\|_{L_\infty(\mathbb{T})} \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} \leq C_1^2 \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})}$$

and

$$\|T_{\hat{a}}^n \omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} \leq \|T_{\hat{a}}^n 1\|_{L_\infty(\mathbb{T})} \|\omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} \leq C_1^2 \|\omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})}.$$

Therefore, we deduce from (5.7) that

$$\begin{aligned} \|\omega_{n+k}(\cdot, h)\|_{L_\infty(\mathbb{T})} &\leq C_1^2 \|\omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} + 2C_1^4 C_2 h^\beta \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} \sum_{j=1}^{\infty} 2^{-j\beta} \\ &\leq C_3 \left(\|\omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} + h^\beta \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} \right). \end{aligned}$$

That is, (5.6) has been proved. In particular, for any $\tau > 0$, we take $g(\xi) = |\sin(\xi/2)|^\tau$ and $\nu := \min(\beta, \tau) > 0$. By Proposition 4.2, $g \in C^\nu(\mathbb{T})$. It is evident that $\|T_{\hat{a}}^0 g\|_{L_\infty(\mathbb{T})} = \|g\|_{L_\infty(\mathbb{T})} = 1$. By $g \in C^\nu(\mathbb{T})$, there exists a positive constant C_4 such that

$$(5.8) \quad (2^{-n}h)^{-\nu} \omega_0(\xi, 2^{-n}h) = (2^{-n}h)^{-\nu} |g(\xi) - g(\xi - 2^{-n}h)| \leq C_4.$$

Since $0 < \nu \leq \beta$, we see that (5.4) holds with β being replaced by ν (now the constant C_2 in (5.4) may be different). That is, for all $k, n \in \mathbb{N} \cup \{0\}$ and $h > 0$, (5.6) becomes

$$(5.9) \quad \begin{aligned} &\|h^{-\nu} \omega_{n+k}(\cdot, h)\|_{L_\infty(\mathbb{T})} \\ &\leq C_3 \left(2^{-n\nu} \|(2^{-n}h)^{-\nu} \omega_k(\cdot, 2^{-n}h)\|_{L_\infty(\mathbb{T})} + \|T_{\hat{a}}^k g\|_{L_\infty(\mathbb{T})} \right). \end{aligned}$$

Setting $k = 0$ in (5.9), we deduce that

$$(5.10) \quad \|h^{-\nu} \omega_n(\cdot, h)\|_{L_\infty(\mathbb{T})} \leq C_3 (2^{-n\nu} C_4 + 1) \leq C_3 (C_4 + 1) < \infty \quad \forall h > 0, n \in \mathbb{N}.$$

By the definition of $\omega_n(\xi, h)$ in (5.5), this is equivalent to saying that

$$|[T_{\hat{a}}^n g](\xi_1) - [T_{\hat{a}}^n g](\xi_2)| \leq C_3 (C_4 + 1) |\xi_1 - \xi_2|^\nu \quad \forall n \in \mathbb{N}, \xi_1, \xi_2 \in \mathbb{R}.$$

Note that \hat{a} is continuous and $\|T_{\hat{a}}^n g\|_{L_\infty(\mathbb{T})} \leq \|T_{\hat{a}}^n 1\|_{L_\infty(\mathbb{T})} \|g\|_{L_\infty(\mathbb{T})} \leq C_1^2$. So, the sequence $\{T_{\hat{a}}^n g\}_{n=1}^\infty$ is bounded and equicontinuous in $C(\mathbb{T})$. By the Arzela–Ascoli theorem, there is a subsequence $\{T_{\hat{a}}^{n_k} g\}_{k=1}^\infty$ converging to $g_\infty \in C(\mathbb{T})$ as $k \rightarrow \infty$; that is,

$$(5.11) \quad \lim_{k \rightarrow \infty} \|T_{\hat{a}}^{n_k} g - g_\infty\|_{C(\mathbb{T})} = 0.$$

Since all $T_{\hat{a}}^{n_k} g \geq 0$, we must have $g_\infty \geq 0$. Now we show that if (ii) holds, then we must have $g_\infty \equiv 0$. Define

$$\begin{aligned} g_n(\xi) &:= g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 = |\sin(2^{-1-n}\xi)|^\tau |\widehat{f_n}(\xi)|^2 \\ &= |\sin(2^{-1-n}\xi)|^\tau |\hat{f}(2^{-n}\xi)|^2 \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2. \end{aligned}$$

By induction, we observe that $[g_n, g_n] = T_{\hat{a}}^n([\hat{f}, \hat{f}]g)$ for all $n \in \mathbb{N}$. Since $[\hat{f}, \hat{f}] \geq 1/C_1$, we conclude that

$$[T_{\hat{a}}^n g](\xi) \leq C_1 [T_{\hat{a}}^n([\hat{f}, \hat{f}]g)](\xi) = C_1 [g_n, g_n](\xi).$$

Therefore,

$$(5.12) \quad \int_0^{2\pi} [T_{\hat{a}}^n g](\xi) d\xi \leq C_1 \int_0^{2\pi} [g_n, g_n](\xi) d\xi = C_1 \int_{\mathbb{R}} g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 d\xi.$$

By $g(\xi) = |\sin(\xi/2)|^\tau$ and $\tau > 0$, we have $\lim_{n \rightarrow \infty} g(2^{-n}\xi) = g(0) = 0$. Since $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$, we have $\lim_{n \rightarrow \infty} \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2 = |\hat{\phi}(\xi)|^2$. Observing that

$$\begin{aligned} 0 &\leq g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 = g(2^{-n}\xi) |\hat{f}(2^{-n}\xi)|^2 \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2 \\ &\leq \|f\|_{L_{2,\infty,0}(\mathbb{R})}^2 g(2^{-n}\xi) \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2, \end{aligned}$$

we see that $\lim_{n \rightarrow \infty} g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 = 0$ for almost every $\xi \in \mathbb{R}$.

Since $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$, the sequence $\{f_n\}_{n=1}^\infty$ must also be a Cauchy sequence in $L_2(\mathbb{R})$ by $\|f\|_{L_2(\mathbb{R})} \leq \|f\|_{L_{2,\infty,0}(\mathbb{R})}$. Consequently,

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} |\widehat{f_n}(\xi)|^2 d\xi = \lim_{n \rightarrow \infty} \int_0^{2\pi} [\widehat{f_n}, \widehat{f_n}](\xi) d\xi$$

exists and is finite. Now by $0 \leq g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 \leq |\widehat{f_n}(\xi)|^2$ and the generalized Lebesgue dominated convergence theorem, we conclude that

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 d\xi = \int_{\mathbb{R}} \lim_{n \rightarrow \infty} g(2^{-n}\xi) |\widehat{f_n}(\xi)|^2 d\xi = \int_{\mathbb{R}} 0 d\xi = 0.$$

Hence, by (5.12), we get

$$(5.13) \quad \lim_{n \rightarrow \infty} \int_0^{2\pi} [T_{\hat{a}}^n g](\xi) d\xi = 0.$$

Now it follows from (5.11) and (5.13) that

$$\int_0^{2\pi} g_\infty(\xi) d\xi = \lim_{k \rightarrow \infty} \int_0^{2\pi} [T_{\hat{a}}^{n_k} g](\xi) d\xi = 0.$$

Since $g_\infty \geq 0$, we must have $g_\infty = 0$. That is, (5.11) becomes

$$(5.14) \quad \lim_{k \rightarrow \infty} \|T_{\hat{a}}^{n_k} g\|_{L_\infty(\mathbb{T})} = 0.$$

Setting $k = n_k$ in (5.9), by (5.10), we deduce that

$$\|h^{-\nu} \omega_{n+n_k}(\cdot, h)\|_{L_\infty(\mathbb{T})} \leq C_3(2^{-n\nu} C_3(C_4 + 1) + \|T_{\hat{a}}^{n_k} g\|_{L_\infty(\mathbb{T})}).$$

So, when n and k are large enough, by (5.14) and $\nu > 0$, setting $N = n + n_k$, we must have

$$h^{-\nu} \omega_N(\xi, h) \leq \pi^{-\nu} \quad \forall \xi \in \mathbb{R}, h > 0.$$

In particular, setting $\xi = 0$ or $-h$ in the above inequality, we have

$$|[T_{\hat{a}}^N g](h) - [T_{\hat{a}}^N g](0)| \leq \pi^{-\nu} |h|^\nu \quad \forall h \in \mathbb{R}.$$

By (5.13) and Proposition 4.2, we see that $\hat{a}(\pi) = 0$. Since $g(0) = \hat{a}(\pi) = 0$, by induction, one can verify that $[T_{\hat{a}}^n g](0) = 0$ for all $n \in \mathbb{N}$. Therefore, it follows from the above inequality that

$$\frac{[T_{\hat{a}}^N g](h)}{|\sin(h/2)|^\nu} \leq (\pi/2)^\nu \frac{[T_{\hat{a}}^N g](h)}{|h|^\nu} \leq 2^{-\nu} < 1 \quad \forall h \in [-\pi, \pi] \setminus \{0\}.$$

In order to show that (ii) \Rightarrow (iii), by Proposition 4.2, it suffices to show $\rho_\tau(\hat{a}, \infty) < 1$ for sufficiently small $\tau > 0$. Since $\hat{a} \in C^\beta(\mathbb{T})$ with $\beta > 0$, we must have $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^{2\beta} \in L_\infty(\mathbb{T})$. So, for any $0 < \tau < \beta$, by $\nu = \min(\beta, \tau)$, we have $\nu = \tau$, and by Theorem 4.1, we must have

$$\rho_\tau(\hat{a}, \infty) = \inf_{n \in \mathbb{N}} \|[T_{\hat{a}}^n g]/g\|_{L_\infty(\mathbb{T})}^{1/n} \leq \|[T_{\hat{a}}^N g]/g\|_{L_\infty(\mathbb{T})}^{1/N} \leq 2^{-\nu/N} < 1.$$

So, (iii) holds for all $0 < \tau < \beta$, and therefore (iii) holds for all $\tau > 0$. Hence, (ii) \Rightarrow (iii).

By (iii) and Proposition 4.2, we have $\hat{a}(\pi) = 0$. So, it follows from $\hat{a} \in C^\beta(\mathbb{T})$ that $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^{2\beta} \in L_\infty(\mathbb{T})$. Now it is straightforward to see that (iii) \Rightarrow (iv), since $\rho(\hat{a}) \leq \rho_{2\beta}(\hat{a}, \infty) < 1$.

By the definition of $\rho(\hat{a})$, (iv) implies that $\rho_\tau(\hat{a}, \infty) < 1$ and $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$ for some $\tau \geq 0$. On the other hand, since $\hat{a}(0) = 1$ and \hat{a} is continuous, by $[T_{\hat{a}}^n 1](0) \geq |\hat{a}(0)|^n = 1$, we must have $\rho_0(\hat{a}, \infty) \geq 1$. Therefore, $\tau > 0$, and so (iv) \Rightarrow (v).

Now we show that (v) \Rightarrow (ii). By (v) and Proposition 4.2, $\hat{a}(\pi) = 0$ and $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$. Let $\hat{\eta}$ be defined in (5.1). Then $\eta \in L_{2,\infty,0}(\mathbb{R})$ and the shifts of η are stable in $L_2(\mathbb{R})$. Consequently, it is easy to directly verify that

$$(5.15) \quad H := [\hat{a}(\cdot/2)\hat{\eta}(\cdot/2) - \hat{\eta}, \hat{a}(\cdot/2)\hat{\eta}(\cdot/2) - \hat{\eta}] / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T}).$$

So, η is an admissible function in $L_{2,\infty,0}(\mathbb{R})$ with respect to \hat{a} such that the shifts of η are stable in $L_2(\mathbb{R})$. Taking $f = \eta$ and defining f_n as in (2.2), we show that $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$.

Note that

$$(5.16) \quad |\widehat{f_{n+1}}(\xi) - \widehat{f_n}(\xi)| = |\hat{a}(2^{-1-n}\xi)\hat{f}(2^{-1-n}\xi) - \hat{f}(2^{-n}\xi)| \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|.$$

Consequently, we have

$$[\widehat{f_{n+1}} - \widehat{f_n}, \widehat{f_{n+1}} - \widehat{f_n}](\xi) = \left[T_{\hat{a}}^n([\hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}, \hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}]) \right](\xi) \quad \forall n \in \mathbb{N}.$$

By (5.15), we have

$$[\hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}, \hat{a}(\cdot/2)\hat{f}(\cdot/2) - \hat{f}](\xi) = H(\xi)g(\xi)$$

and $H \in L_\infty(\mathbb{T})$, where $g(\xi) := |\sin(\xi/2)|^\tau$. Hence, we have

$$(5.17) \quad \|f_{n+1} - f_n\|_{L_{2,\infty,0}(\mathbb{R})}^2 = \|[\widehat{f_{n+1}} - \widehat{f_n}, \widehat{f_{n+1}} - \widehat{f_n}]\|_{L_\infty(\mathbb{T})} \leq \|H\|_{L_\infty(\mathbb{T})} \|T_{\hat{a}}^n g\|_{L_\infty(\mathbb{T})}.$$

By our assumption in (v), we have $\rho_\tau(\hat{a}, \infty) < 1$, and therefore, for any ρ such that $\rho_\tau(\hat{a}, \infty) < \rho < 1$, there exists a positive constant C such that $\|T_{\hat{a}}^n g\|_{L_\infty(\mathbb{T})} \leq C\rho^n$ for all $n \in \mathbb{N}$. Now it follows from (5.17) that

$$(5.18) \quad \|f_{n+1} - f_n\|_{L_{2,\infty,0}(\mathbb{R})} \leq \|H\|_{L_\infty(\mathbb{T})}^{1/2} C^{1/2} \rho^{n/2} \quad \forall n \in \mathbb{N},$$

which yields, by $0 < \rho < 1$, that $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$. So, (v) \Rightarrow (ii).

Finally, we show that (iii) \Rightarrow (i). Let f be an admissible function in $L_{2,\infty,0}(\mathbb{R})$ with respect to \hat{a} such that (2.1) holds for some $\tau > 0$. By Proposition 4.2 and (iii), we have $\hat{a}(\pi) = 0$. Now by $\hat{a} \in C^\beta(\mathbb{T})$, we must have $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^{2\beta} \in L_\infty(\mathbb{T})$. Replacing τ by $\min(\tau, 2\beta)$, we see that (2.1) still holds and $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$. By our assumption in (iii) and $\tau > 0$, we have $\rho_\tau(\hat{a}, \infty) < 1$. Now by the same proof for showing (v) \Rightarrow (ii), we see that $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$. Therefore, (iii) \Rightarrow (i).

Now we prove the rest of Theorem 2.1. Since $\nu_2(\hat{a}) > 0$, we must have $\hat{a}(\pi) = 0$, and (5.3) holds. Consequently,

$$(5.19) \quad \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2 \leq \|T_{\hat{a}}^n 1\|_{L_\infty(\mathbb{T})} \leq C_1^2 \quad \forall n \in \mathbb{N}, \xi \in \mathbb{R}.$$

Since $\hat{a} \in C^\beta(\mathbb{T})$ and since we assumed $0 < \beta \leq 1$, there exists a positive constant C such that $|\hat{a}(\xi_1) - \hat{a}(\xi_2)| \leq C|\xi_1 - \xi_2|^\beta$ for all $\xi_1, \xi_2 \in \mathbb{R}$. We deduce that

$$\begin{aligned} |\hat{\phi}(\xi_1) - \hat{\phi}(\xi_2)| &= \left| \sum_{j=1}^{\infty} \left[\prod_{k=1}^{j-1} \hat{a}(2^{-k}\xi_1) \right] [\hat{a}(2^{-j}\xi_1) - \hat{a}(2^{-j}\xi_2)] \left[\prod_{\ell=j+1}^{\infty} \hat{a}(2^{-\ell}\xi_2) \right] \right| \\ &\leq \sum_{j=1}^{\infty} \left| \prod_{k=1}^{j-1} \hat{a}(2^{-k}\xi_1) \right| \times |\hat{a}(2^{-j}\xi_1) - \hat{a}(2^{-j}\xi_2)| \times \left| \prod_{\ell=j+1}^{\infty} \hat{a}(2^{-\ell}\xi_2) \right| \\ &\leq C_1^2 C \sum_{j=1}^{\infty} 2^{-j\beta} |\xi_1 - \xi_2|^\beta \leq |\xi_1 - \xi_2|^\beta C_1^2 C / (1 - 2^{-\beta}). \end{aligned}$$

So, $\hat{\phi} \in \Lambda^\beta(\mathbb{R})$. Let $\hat{\eta}$ be defined in (5.1). Take $f = \eta$ and define f_n as in (2.2). Since $\nu_2(\hat{a}) > 0$, $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$. In particular, by $\lim_{n \rightarrow \infty} \widehat{f_n}(\xi) = \hat{\phi}(\xi)$, we have $\phi \in L_{2,\infty,0}(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \|f_n - \phi\|_{L_{2,\infty,0}(\mathbb{R})} = 0$. By induction, we have $[\widehat{f_n}, \widehat{f_n}] = T_{\hat{a}}^n[\hat{f}, \hat{f}]$. Note that $[\widehat{f_n}, \widehat{f_n}]$ is continuous. Now define $g := [\hat{\eta}, \hat{\eta}]$. By the definition of $\hat{\eta}$ in (5.1) and $\hat{\phi} \in C^\beta(\mathbb{R})$, it is easy to check that $g \in C^\beta(\mathbb{T})$. Consequently, taking $k = 0$ in (5.6), we see that a similar result as in (5.10) holds; that is, there exists a positive constant C_5 such that

$$(5.20) \quad h^{-\beta} \left\| [\widehat{f_n}, \widehat{f_n}] - [\widehat{f_n}, \widehat{f_n}](\cdot - h) \right\|_{L_\infty(\mathbb{T})} = h^{-\beta} \|\omega_n(\cdot, h)\|_{L_\infty(\mathbb{T})} \leq C_5 \quad \forall h > 0,$$

where ω_n is defined in (5.5) and we used the fact $[\widehat{f_n}, \widehat{f_n}] = T_{\hat{a}}^n[\hat{f}, \hat{f}] = T_{\hat{a}}^n g$. Since $\lim_{n \rightarrow \infty} \|f_n - \phi\|_{L_{2,\infty,0}(\mathbb{R})} = 0$, we deduce that $\lim_{n \rightarrow \infty} \|[\widehat{f_n}, \widehat{f_n}] - [\hat{\phi}, \hat{\phi}]\|_{C(\mathbb{T})} = 0$. Now it follows from (5.20) that $\|[\hat{\phi}, \hat{\phi}] - [\hat{\phi}, \hat{\phi}](\cdot - h)\|_{C(\mathbb{T})} \leq C_5 h^\beta$ for all $h > 0$. So, $[\hat{\phi}, \hat{\phi}] \in \Lambda^\beta(\mathbb{T})$.

Now we prove (2.3). Since $0 \leq \nu < \nu_2(\hat{a})$, by the definition of $\nu_2(\hat{a})$, there exists $\tau > 0$ such that $\rho_\tau(\hat{a}, \infty) < 2^{-2\nu}$ and $|\hat{a}(\cdot + \pi)|^2 / |\sin(\cdot/2)|^\tau \in L_\infty(\mathbb{T})$. Let $\hat{\eta}$ be defined in (5.1). Take $f = \eta$ and define f_n as in (2.2). As in the proof of (v) \Rightarrow (ii), we see that (5.18) holds for $\rho_\tau(\hat{a}, \infty) < \rho < 2^{-2\nu}$. Denote $g_n := \widehat{f_{n+1}} - \widehat{f_n}$ and

$$[g_n, g_n]_\nu(\xi) := \sum_{k \in \mathbb{Z}} |g_n(\xi + 2\pi k)|^2 (1 + |\xi + 2\pi k|^2)^\nu, \quad \xi \in [-\pi, \pi].$$

Since $\hat{f} = \hat{\eta}$ is supported inside $[-3\pi/2, -\varepsilon] \cup [\varepsilon, 3\pi/2]$, it follows from (5.16) that g_n is supported inside $[-3\pi 2^n, -\varepsilon 2^n] \cup [\varepsilon 2^n, 3\pi 2^n]$. Therefore, for $\xi \in [-\pi, \pi]$, it follows from (5.18) that

$$\begin{aligned} [g_n, g_n]_\nu(\xi) &= \sum_{|\xi + 2\pi k| \leq 3\pi 2^n} |g_n(\xi + 2\pi k)|^2 (1 + |\xi + 2\pi k|^2)^\nu \\ &\leq (1 + (3\pi 2^n)^2)^\nu \sum_{k \in \mathbb{Z}} |g_n(\xi + 2\pi k)|^2 \\ &= (1 + (3\pi 2^n)^2)^\nu \|f_{n+1} - f_n\|_{L_{2,\infty,0}(\mathbb{R})}^2 \leq C_6 (2^{2\nu} \rho)^n, \end{aligned}$$

where $C_6 := 18^\nu \pi^{2\nu} C \| [h, h] \|_{L_\infty(\mathbb{T})} < \infty$. Since $[g_n, g_n]$ is continuous, we conclude that

$$(5.21) \quad [g_n, g_n]_\nu(\xi) \leq C_6 (2^{2\nu} \rho)^n \quad \forall n \in \mathbb{N}, \xi \in [-\pi, \pi].$$

Since $\nu_2(\hat{a}) > 0$, we have $\phi = \eta + \sum_{n=0}^\infty (f_{n+1} - f_n)$ in $L_{2,\infty,0}(\mathbb{R})$ with $f_0 := \eta$. Since all $\hat{\phi}, \hat{\eta}$, and $g_n, n \in \mathbb{N}$, are continuous and g_n is supported inside $[-3\pi 2^n, -\varepsilon 2^n] \cup [\varepsilon 2^n, 3\pi 2^n]$, we must have $\hat{\phi}(\xi) = \hat{\eta}(\xi) + \sum_{n=0}^\infty g_n(\xi)$ for all $\xi \in \mathbb{R}$, where the series is in fact a finite sum for any ξ in any bounded set. Therefore, for all $\xi \in [-2\pi, 2\pi]$ and $N \geq 3$, we have

$$\begin{aligned} \left(\sum_{|k|=N+1}^\infty (1 + |\xi + 2\pi k|^2)^\nu |\hat{\phi}(\xi + 2\pi k)|^2 \right)^{1/2} &\leq \sum_{n=\log_2(N/3)}^\infty [g_n, g_n]_\nu^{1/2}(\xi) \\ &\leq C_6^{1/2} \sum_{n=\log_2(N/3)}^\infty (2^\nu \rho^{1/2})^n \leq C_6^{1/2} (N/3)^{\log_2(2^\nu \rho^{1/2})} / (1 - 2^\nu \rho^{1/2}). \end{aligned}$$

Since $2^\nu \rho^{1/2} < 1$, we have $\log_2(2^\nu \rho^{1/2}) < 0$, and therefore $\lim_{N \rightarrow \infty} (N/3)^{\log_2(2^\nu \rho^{1/2})} = 0$. Hence, the series $\sum_{k=-N}^N (1 + |\xi + 2\pi k|^2)^\nu |\hat{\phi}(\xi + 2\pi k)|^2$ is uniformly convergent as $N \rightarrow \infty$ for all $\xi \in [-2\pi, 2\pi]$. Since $\hat{\phi}$ is continuous, we conclude that $[\hat{\phi}, \hat{\phi}]_\nu \in C(\mathbb{T})$. That is, (2.3) holds. \square

For a mask $\hat{a} \in C^\beta(\mathbb{T})$ with $\beta > 0$, a more complicated argument can be used to show a stronger statement that $\hat{\phi} \in C^\beta(\mathbb{R})$ and $[\hat{\phi}, \hat{\phi}] \in C^\beta(\mathbb{T})$, instead of the result $\hat{\phi} \in \Lambda^{\min(1, \beta)}(\mathbb{R})$ and $[\hat{\phi}, \hat{\phi}] \in \Lambda^{\min(1, \beta)}(\mathbb{T})$ stated in Theorem 2.1. We shall address this technical issue elsewhere.

6. Proof of Theorem 2.3. Before we present the proof of Theorem 2.3, we need the two following auxiliary results.

PROPOSITION 6.1. *Let $f \in L_{2,1,\gamma_2}(\mathbb{R})$ with $\gamma_2 > 0$. Then for $0 \leq \gamma_1 < \gamma_2$,*

$$(6.1) \quad [\hat{f}, \hat{f}](\xi + i\zeta) := \sum_{k \in \mathbb{Z}} |\hat{f}(\xi + i\zeta)|^2 \leq C_{\gamma_2 - \gamma_1} \|f\|_{L_{2,1,\gamma_2}(\mathbb{R})}^2$$

$$\forall \xi \in \mathbb{R}, \zeta \in [-\gamma_1, \gamma_1],$$

where

$$C_\gamma := \left\| \sum_{k \in \mathbb{Z}} e^{-2\gamma|\cdot - k|} \right\|_{L_\infty(\mathbb{R})} < \infty$$

for $\gamma > 0$. If in addition $\hat{f}(2\pi k) = 0$ for all $k \in \mathbb{Z}$, then $f = h - h(\cdot - 1)$ and $h \in L_{2,1,\gamma_1}(\mathbb{R})$ for all $0 \leq \gamma_1 < \gamma_2$, where $h := \sum_{k=0}^{\infty} f(\cdot - k)$.

Proof. Since $f \in L_{2,1,\gamma_2}(\mathbb{R})$, by the Cauchy-Schwarz inequality, for $0 \leq \gamma_1 < \gamma_2$, we have

$$(6.2) \quad \int_0^1 \left(\sum_{k \in \mathbb{Z}} |f(x - k)| e^{\gamma_1|x - k|} \right)^2 dx \leq C_\gamma \int_0^1 \sum_{k \in \mathbb{Z}} |f(x - k)|^2 e^{2\gamma_2|x - k|} dx$$

$$= C_\gamma \|f\|_{L_{2,1,\gamma_2}(\mathbb{R})}^2,$$

where $\gamma := \gamma_2 - \gamma_1 > 0$. Using the Fourier series of $[\hat{f}, \hat{f}]$ and $\hat{f}(\xi + i\zeta) = \widehat{e^{\zeta \cdot} f}(\xi)$, for any fixed $\zeta \in [-\gamma_1, \gamma_1]$ and almost every $\xi \in \mathbb{R}$, we deduce from (6.2) that

$$[\hat{f}, \hat{f}](\xi + i\zeta) = \left| \sum_{k \in \mathbb{Z}} e^{ik\xi} \int_{\mathbb{R}} e^{\zeta x} f(x) \overline{e^{\zeta(x+k)} f(x+k)} dx \right|$$

$$\leq \sum_{k \in \mathbb{Z}} \int_{\mathbb{R}} e^{|\zeta x|} |f(x)| e^{|\zeta(x+k)|} |f(x+k)| dx$$

$$= \int_0^1 \left(\sum_{k \in \mathbb{Z}} |f(x+k)| e^{|\zeta| \times |x+k|} \right)^2 dx \leq C_\gamma \|f\|_{L_{2,1,\gamma_2}(\mathbb{R})}^2.$$

Since \hat{f} is continuous, it follows from the above inequality that (6.1) holds for all $\xi \in \mathbb{R}$.

If $\hat{f}(2\pi k) = 0$ for all $k \in \mathbb{Z}$, then we have $\sum_{k \in \mathbb{Z}} f(\cdot - k) = 0$. From the definition of h , we see that $f = h - h(\cdot - 1)$ and $h = -\sum_{k=-\infty}^{-1} f(\cdot - k)$. We now verify

that $h \in L_{2,1,\gamma_3}$ for all $0 < \gamma_3 < \gamma_2$. Take γ_1 such that $\gamma_3 < \gamma_1 < \gamma_2$. Then by $h = \sum_{k=0}^{\infty} f(\cdot - k)$, we have

$$\begin{aligned} \sum_{j=0}^{\infty} |h(x-j)|e^{\gamma_3|x-j|} &\leq \sum_{j=0}^{\infty} \sum_{k=0}^{\infty} |f(x-j-k)|e^{\gamma_3|x-j|} \\ &= \sum_{k=0}^{\infty} |f(x-k)|e^{\gamma_1|x-k|} e^{-\gamma_1|x-k|} \sum_{j=0}^k e^{\gamma_3|x-j|}. \end{aligned}$$

For $x \in [0, 1]$, by $\gamma_3 < \gamma_1 < \gamma_2$ and $k \geq 0$, we have

$$e^{-\gamma_1|x-k|} \sum_{j=0}^k e^{\gamma_3|x-j|} \leq e^{-\gamma_1(k-1)} \sum_{j=0}^k e^{\gamma_3(j+1)} \leq \frac{e^{\gamma_2+2\gamma_3}}{e^{\gamma_3}-1} := C.$$

Therefore, we deduce that

$$\sum_{j=0}^{\infty} |h(x-j)|e^{\gamma_3|x-j|} \leq C \sum_{k=0}^{\infty} |f(x-k)|e^{\gamma_1|x-k|}, \quad x \in [0, 1].$$

Similarly, using $h = -\sum_{k=-\infty}^{-1} f(\cdot - k)$, we have

$$\sum_{j=-\infty}^{-1} |h(x-j)|e^{\gamma_3|x-j|} \leq C \sum_{k=-\infty}^{-1} |f(x-k)|e^{\gamma_1|x-k|}, \quad x \in [0, 1].$$

Hence,

$$\sum_{j \in \mathbb{Z}} |h(x-j)|e^{\gamma_3|x-j|} \leq C \sum_{k \in \mathbb{Z}} |f(x-k)|e^{\gamma_1|x-k|}.$$

By (6.2), we conclude that

$$\begin{aligned} \|h\|_{L_{2,1,\gamma_3}(\mathbb{R})}^2 &= \int_0^1 \sum_{j \in \mathbb{Z}} |h(x-j)e^{\gamma_3|x-k|}|^2 dx \leq \int_0^1 \left(\sum_{j \in \mathbb{Z}} |h(x-j)|e^{\gamma_3|x-j|} \right)^2 dx \\ &\leq C \int_0^1 \left(\sum_{k \in \mathbb{Z}} |f(x-k)|e^{\gamma_1|x-k|} \right)^2 dx \leq CC_{\gamma_2-\gamma_1} \|f\|_{L_{2,1,\gamma_2}(\mathbb{R})}^2. \end{aligned}$$

So, $h \in L_{2,1,\gamma_3}(\mathbb{R})$ for all $0 < \gamma_3 < \gamma_2$. This completes the proof. \square

It follows from (6.1) and $\widehat{f}(\xi + i\zeta) = e^{\widehat{C}\zeta} \widehat{f}(\xi)$ that for any $1 \leq p, q \leq \infty$,

$$(6.3) \quad \|f\|_{L_{2,p,\gamma_1}(\mathbb{R})}^2 \leq \|f\|_{L_{2,\infty,\gamma_1}(\mathbb{R})}^2 \leq C_{\gamma_2-\gamma_1} \|f\|_{L_{2,1,\gamma_2}(\mathbb{R})}^2 \leq C_{\gamma_2-\gamma_1} \|f\|_{L_{2,q,\gamma_2}(\mathbb{R})}^2$$

$$\forall 0 \leq \gamma_1 < \gamma_2.$$

Consequently, $L_{2,q,\gamma_2}(\mathbb{R}) \subseteq L_{2,p,\gamma_1}(\mathbb{R})$ for all $1 \leq p, q \leq \infty$ and $0 \leq \gamma_1 < \gamma_2$.

Now we have the following result on admissible functions in a cascade algorithm.

PROPOSITION 6.2. *Let $f \in L_{2,1,\gamma}(\mathbb{R})$ for some $\gamma > 0$ such that f satisfies*

$$(6.4) \quad \widehat{f}(2\pi k) = 0 \quad \forall k \in \mathbb{Z} \setminus \{0\}.$$

Then f is admissible with respect to \hat{a} for any $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$, $\hat{a}(\pi) = 0$, and $\beta > 0$.

Proof. Let $\hat{f}_1 := \hat{f}(2 \cdot + 2\pi)$ and $\hat{f}_2 := \hat{f} - \hat{f}(2 \cdot)$. Then $f_1, f_2 \in L_{2,1,\gamma/2}(\mathbb{R})$ and $\hat{f}_1(2\pi k) = \hat{f}_2(2\pi k) = 0$ for all $k \in \mathbb{Z}$. By Proposition 6.1, $f_1 = h_1 - h_1(\cdot - 1)$ and $f_2 = h_2 - h_2(\cdot - 1)$ for some $h_1, h_2 \in L_{2,1,\gamma_1}(\mathbb{R})$ with $0 < \gamma_1 < \gamma/2$. Therefore, we have

$$(6.5) \quad [\hat{f}_1, \hat{f}_1](\xi) = |1 - e^{-i\xi}|^2 [\widehat{h}_1, \widehat{h}_1](\xi) \quad \text{and} \quad [\hat{f}_2, \hat{f}_2](\xi) = |1 - e^{-i\xi}|^2 [\widehat{h}_2, \widehat{h}_2](\xi).$$

Since $\hat{a} \in C^\beta(\mathbb{T})$ with $\hat{a}(0) = 1$ and $\hat{a}(\pi) = 0$, there exists a positive constant C such that

$$(6.6) \quad |\hat{a}(\xi) - 1| \leq C|1 - e^{-i\xi}|^\beta \quad \text{and} \quad |\hat{a}(\xi + \pi)| \leq C|1 - e^{-i\xi}|^\beta \quad \forall \xi \in \mathbb{R}.$$

Now it follows from (6.5) and (6.6) that

$$\begin{aligned} \sum_{k \in \mathbb{Z}} |\hat{a}(\xi/2 + 2\pi k) \hat{f}(\xi/2 + 2\pi k) - \hat{f}(\xi + 4\pi k)|^2 \\ &= \sum_{k \in \mathbb{Z}} |[\hat{a}(\xi/2) - 1] \hat{f}(\xi/2 + 2\pi k) + \hat{f}_2(\xi/2 + 2\pi k)|^2 \\ &\leq 2|\hat{a}(\xi/2) - 1|^2 [\hat{f}, \hat{f}](\xi/2) + 2[\hat{f}_2, \hat{f}_2](\xi/2) \\ &\leq 2C^2 |1 - e^{-i\xi/2}|^{2\beta} (\|[\hat{f}, \hat{f}]\|_{L_\infty(\mathbb{T})} + 2^{2-2\beta} \|[\widehat{h}_2, \widehat{h}_2]\|_{L_\infty(\mathbb{T})}). \end{aligned}$$

Similarly, by $\hat{f}(\xi + 2\pi + 4\pi k) = \hat{f}_1(\xi/2 + 2\pi k)$, we have

$$\begin{aligned} \sum_{k \in \mathbb{Z}} |\hat{a}(\xi/2 + \pi + 2\pi k) \hat{f}(\xi/2 + \pi + 2\pi k) - \hat{f}(\xi + 2\pi + 4\pi k)|^2 \\ &\leq 2|\hat{a}(\xi/2 + \pi)|^2 [\hat{f}, \hat{f}](\xi/2 + \pi) + 2[\hat{f}_1, \hat{f}_1](\xi/2) \\ &\leq 2C^2 |1 - e^{-i\xi/2}|^{2\beta} (\|[\hat{f}, \hat{f}]\|_{L_\infty(\mathbb{T})} + 2^{2-2\beta} \|[\widehat{h}_1, \widehat{h}_1]\|_{L_\infty(\mathbb{T})}). \end{aligned}$$

Letting $C_1 := 8C^2 (\|[\hat{f}, \hat{f}]\|_{L_\infty(\mathbb{T})} + \|[\widehat{h}_1, \widehat{h}_1]\|_{L_\infty(\mathbb{T})} + \|[\widehat{h}_2, \widehat{h}_2]\|_{L_\infty(\mathbb{T})}) < \infty$, combining the above two inequalities, for almost every $\xi \in [-\pi, \pi]$, we have

$$|\hat{a}(\cdot/2) \hat{f}(\cdot/2) - \hat{f}, \hat{a}(\cdot/2) \hat{f}(\cdot/2) - \hat{f}](\xi) \leq C_1 |1 - e^{-i\xi/2}|^{2\beta} \leq 2^\beta C_1 |\sin(\xi/2)|^{2\beta}.$$

Therefore, (2.1) holds with $\tau = 2\beta > 0$. \square

Proof of Theorem 2.3. By (6.3), it is easy to see that (i) \Rightarrow (ii). (ii) \Rightarrow (iii) is obvious. To show that (iii) \Rightarrow (iv), we take $f = \max(0, 1 - |\cdot|)$. It is easy to check that f satisfies the condition in (6.4) and the shifts of f are stable. By $\hat{a}(\pi) = 0$ and Proposition 6.2, $f \in L_{2,p,\gamma}(\mathbb{R})$ is admissible with respect to \hat{a} . Therefore, (iii) \Rightarrow (iv).

If (iv) holds, by Proposition 6.1 or (6.3), $f \in L_{2,\infty,0}(\mathbb{R})$ and f is admissible with respect to \hat{a} . Now it follows from (iv) and (6.3) that $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,\infty,0}(\mathbb{R})$. Therefore, (ii) of Theorem 2.1 holds, and consequently $\nu_2(\hat{a}) > 0$. So, (iv) \Rightarrow (v).

To complete the proof, we have to show that (v) \Rightarrow (i), which is the major part of this proof. By Proposition 4.2, $\nu_2(\hat{a}) > 0$ implies $\hat{a}(\pi) = 0$. So, we can write $\hat{a}(\xi) = (1 + e^{-i\xi}) \hat{A}(\xi)$, where \hat{A} also has exponential decay of order r . By Theorems 2.1 and 4.1, it follows from $\nu_2(\hat{a}) > 0$ and $\hat{a}(\xi) = (1 + e^{-i\xi}) \hat{A}(\xi)$ that $\inf_{n \in \mathbb{N}} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} =$

$\rho_0(\hat{A}, \infty) = \rho_2(\hat{a}, \infty) < 1$. Therefore, there exist $0 < \rho < 1$ and $N \in \mathbb{N}$ such that $\|T_{\hat{A}}^N 1\|_{L_\infty(\mathbb{T})}^{1/N} < \rho < 1$. Since \hat{A} has exponential decay of order r , the operator $[T_{\hat{A}}^n 1](\xi)$ in (1.5) is well defined for $\xi \in \Gamma_{2r} := \{z \in \mathbb{C} : |\operatorname{Im}(z)| < 2r\}$. Since \hat{A} is a 2π -periodic continuous function on the strip Γ_r , there exists $\gamma_1 > 0$ such that

$$(6.7) \quad \|[T_{\hat{A}}^N 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} \leq \rho^N < 1 \quad \forall \zeta \in [-\gamma_1, \gamma_1].$$

Take m to be the smallest nonnegative integer such that $2^{1-m}r < \gamma_1$. For each $n \geq N + m$, we can uniquely write $n = Nk + j$, where $j \in \{m, m+1, \dots, m+N-1\}$ and $k \in \mathbb{N}$. So, for every $\zeta \in [-\gamma, \gamma]$ and $0 < \gamma < 2r$, by (6.7) and the extended definition of $T_{\hat{A}}$ in (1.5), we have $|2^{-j}\zeta| \leq 2^{-m}\gamma < 2^{1-m}r < \gamma_1$ and

$$\begin{aligned} \|[T_{\hat{A}}^n 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} &= \|[T_{\hat{A}}^j T_{\hat{A}}^{Nk} 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} \\ &\leq \|[T_{\hat{A}}^j 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} \|[T_{\hat{A}}^{Nk} 1](\cdot + i2^{-j}\zeta)\|_{L_\infty(\mathbb{T})} \\ &\leq \|[T_{\hat{A}}^j 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} \rho^{Nk} \leq C_1 \rho^n, \end{aligned}$$

where

$$C_1 := \rho^{1-m-N} \sup \{ \|[T_{\hat{A}}^j 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} : j = m, \dots, m+N-1, \zeta \in [-\gamma, \gamma] \} < \infty.$$

That is, we have

$$(6.8) \quad \|[T_{\hat{A}}^n 1](\cdot + i\zeta)\|_{L_\infty(\mathbb{T})} \leq C_1 \rho^n \quad \forall n \geq N + m, \zeta \in [-\gamma, \gamma].$$

Denote $\hat{F}(\xi) := \hat{a}(\xi/2)\hat{f}(\xi/2) - \hat{f}(\xi)$. Since f is admissible with respect to \hat{a} and $F \in L_{2,1,\gamma}(\mathbb{R})$, we must have $\hat{F}(2\pi k) = 0$ for all $k \in \mathbb{Z}$. By Proposition 6.1, $\hat{F}(\xi) = (1 - e^{-i\xi})\hat{h}(\xi)$ for some $h \in L_{2,1,\gamma/2}(\mathbb{R})$. Denote $g_n := \widehat{f_{n+1}} - \widehat{f_n}$ with $g := g_0 := \widehat{f_1} - \widehat{f}$. Then $g(\xi) = \hat{F}(\xi)$. Therefore, by Proposition 6.1, we conclude that

$$(6.9) \quad \begin{aligned} [g, g](\xi) &= [\hat{F}, \hat{F}](\xi) = |1 - e^{-i\xi}|^2 |\hat{h}, \hat{h}](\xi) \leq C_2 |1 - e^{-i\xi}|^2 \\ &\forall \xi \in \Gamma_{\gamma_2} := \{z \in \mathbb{C} : |\operatorname{Im}(z)| < \gamma_2\}, \end{aligned}$$

where $\gamma_2 := \gamma/4 > 0$, $C_2 := C_{\gamma/4} \|h\|_{L_{2,1,\gamma/2}(\mathbb{R})} < \infty$, and $C_{\gamma/4}$ is defined in Proposition 6.1.

By the definition of f_n and by induction on n , for $n \geq n_0 := 1 - \log_2(\gamma_2/r)$, we have that $g_n(\xi) = g(2^{-n}\xi) \prod_{j=1}^n \hat{a}(2^{-j}\xi)$ for $\xi \in \Gamma_{2r}$ and g_n is holomorphic on Γ_{2r} , since $2^{-n}\Gamma_{2r} \subseteq \Gamma_{\gamma_2}$ for all $n \geq n_0$. Now by induction, it follows from (6.9) that

$$(6.10) \quad [g_n, g_n](\xi) = (T_{\hat{a}}^n [g, g])(\xi) \leq C_2 (T_{\hat{a}}^n (|1 - e^{-i\cdot}|^2))(\xi), \quad \xi \in \Gamma_{2r}.$$

By $\hat{a}(\xi) = (1 + e^{-i\xi})\hat{A}(\xi)$, we have $|\hat{a}(\xi)|^2 = 2^2 |\cos(\xi/2)|^2 |\hat{A}(\xi)|^2$. So, (4.8) holds with $\tau = 1$. So, we deduce that

$$[T_{\hat{a}}^n (|1 - e^{-i\cdot}|^2)](\xi) = |1 - e^{-i\xi}|^2 [T_{\hat{A}}^n 1](\xi) \leq (1 + e^{2r})^2 [T_{\hat{A}}^n 1](\xi), \quad \xi \in \Gamma_{2r}.$$

Consequently, since $0 < \gamma < 2r$, it follows from (6.8) and (6.10) that

$$(6.11) \quad \begin{aligned} [g_n, g_n](\xi) &\leq C_2 (1 + e^{2r})^2 |[T_{\hat{A}}^n 1](\xi)| \leq C_1 C_2 (1 + e^{2r})^2 \rho^n \leq C_3 \rho^n \\ &\forall n \geq n_1, |\operatorname{Im}(\xi)| \leq \gamma, \end{aligned}$$

where $C_3 := C_1 C_2 (1 + e^{2r})^2 < \infty$ and $n_1 = \max(n_0, N + m)$. Note that $g_n(\cdot + i\zeta)$ is the Fourier transform of $e^{\zeta}(f_{n+1} - f_n)$. For $\zeta \in [-\gamma, \gamma]$, we have

$$\begin{aligned} 2\pi \|(f_{n+1} - f_n)e^{\zeta}\|_{L_2(\mathbb{R})}^2 &= \|g_n(\cdot + i\zeta)\|_{L_2(\mathbb{R})}^2 = 2\pi \| [g_n, g_n](\cdot + i\zeta) \|_{L_1(\mathbb{T})} \\ &\leq 2\pi \| [g_n, g_n](\cdot + i\zeta) \|_{L_\infty(\mathbb{T})}. \end{aligned}$$

Now by the definition of the space $L_{2,1,\gamma}(\mathbb{R})$, it follows from the above inequality and (6.11) that

$$\|f_{n+1} - f_n\|_{L_{2,1,\gamma}(\mathbb{R})}^2 = \|(f_{n+1} - f_n)e^{\gamma|\cdot|}\|_{L_2(\mathbb{R})}^2 \leq 2C_3\rho^n \quad \forall n \geq n_0.$$

Consequently, $\{f_n\}_{n=1}^\infty$ is a Cauchy sequence in $L_{2,1,\gamma}(\mathbb{R})$. So, (v) \Rightarrow (i).

Now we show that ϕ has exponential decay of order $2r$. Since $\hat{\phi}(\xi) := \prod_{j=1}^\infty \hat{a}(2^{-j}\xi)$ and \hat{a} has exponential decay of order r with $\hat{a}(0) = 1$, we see that $\hat{\phi}$ can be extended into a holomorphic function in the strip Γ_{2r} . If $\nu_2(\hat{a}) > 0$, for every $0 < \gamma < 2r$, by (i) and $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \hat{\phi}(\xi)$ for all $\xi \in \mathbb{R}$ (here we additionally assumed that $\lim_{\xi \rightarrow 0} \hat{f}(\xi) = 1$, which is satisfied by many initial admissible functions), we see that $\lim_{n \rightarrow \infty} \|f_n - \phi\|_{L_{2,1,\gamma}(\mathbb{R})} = 0$ and hence $\phi \in L_{2,1,\gamma}(\mathbb{R})$. Therefore, $\phi \in L_{2,1,\gamma}(\mathbb{R})$ for all $0 \leq \gamma < 2r$. That is, (1.9) holds. \square

7. Proof of Theorem 4.1. Before we present a proof of Theorem 4.1, we need the following result.

PROPOSITION 7.1. *Let \hat{a} and f be 2π -periodic measurable functions. Assume that $f(\xi) > 0$ for almost every $\xi \in \mathbb{R}$ and $[T_{\hat{a}}^n f]/f \in L_\infty(\mathbb{T})$ for all $n \geq n_0$, where $n_0 \in \mathbb{N}$. Then*

$$(7.1) \quad \limsup_{n \rightarrow \infty} \|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n} = \lim_{n \rightarrow \infty} \|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n} = \inf_{n \geq n_0} \|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n}.$$

Proof. Denote $\rho := \inf_{n \geq n_0} \|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n}$. Then $\rho < \infty$ by $[T_{\hat{a}}^{n_0} f]/f \in L_\infty(\mathbb{T})$. Since for all $n \geq n_0$, $\|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n} \geq \rho$, we have

$$(7.2) \quad \liminf_{n \rightarrow \infty} \|[T_{\hat{a}}^n f]/f\|_{L_\infty(\mathbb{T})}^{1/n} \geq \rho.$$

For any $\varepsilon > 0$, there exists $m \in \mathbb{N}$ such that $m \geq n_0$ and $\|[T_{\hat{a}}^m f]/f\|_{L_\infty(\mathbb{T})}^{1/m} \leq \rho + \varepsilon$. Since $f(\xi) > 0$ for almost every $\xi \in \mathbb{R}$, we deduce that

$$(7.3) \quad [T_{\hat{a}}^m f](\xi) \leq \|[T_{\hat{a}}^m f]/f\|_{L_\infty(\mathbb{T})} f(\xi) \leq f(\xi)(\rho + \varepsilon)^m \quad \text{a.e. } \xi \in \mathbb{R}.$$

For each $n \geq 2m$, we can uniquely write $n = mN + j$, where $N \in \mathbb{N}$ and $j \in \{m, m+1, \dots, 2m-1\}$. Therefore, by (7.3),

$$[T_{\hat{a}}^n f](\xi) = [T_{\hat{a}}^j (T_{\hat{a}}^{mN} f)](\xi) \leq (\rho + \varepsilon)^{mN} [T_{\hat{a}}^j f](\xi).$$

Since $f(\xi) > 0$ for almost every $\xi \in \mathbb{R}$, the above inequality yields that

$$\frac{[T_{\hat{a}}^n f](\xi)}{f(\xi)} \leq (\rho + \varepsilon)^{mN} \frac{[T_{\hat{a}}^j f](\xi)}{f(\xi)} \leq C(\rho + \varepsilon)^n \quad \text{a.e. } \xi \in \mathbb{R},$$

where

$$C := \max\{(\rho + \varepsilon)^{-j} \|[T_{\hat{a}}^j f]/f\|_{L_\infty(\mathbb{T})} : j = m, \dots, 2m-1\} < \infty.$$

Thus, $\|[T_{\hat{a}}^n f]/f\|_{L^\infty(\mathbb{T})}^{1/n} \leq C^{1/n}(\rho + \varepsilon)$ for all $n \geq 2m$. Consequently,

$$\limsup_{n \rightarrow \infty} \|[T_{\hat{a}}^n f]/f\|_{L^\infty(\mathbb{T})}^{1/n} \leq (\rho + \varepsilon).$$

Taking $\varepsilon \rightarrow 0$, we have $\limsup_{n \rightarrow \infty} \|[T_{\hat{a}}^n f]/f\|_{L^\infty(\mathbb{T})}^{1/n} \leq \rho$. By (7.2), (7.1) holds. \square

Proof of Theorem 4.1. Without loss of generality, we can assume that $|\hat{a}(0)|^2 = 1$; otherwise, we consider $\hat{a}/|\hat{a}(0)|$. By $|\hat{a}(\xi)|^2 = 2^{2\tau} |\cos(\xi/2)|^{2\tau} |\hat{A}(\xi)|^2$, (4.8) holds. Setting $f = 1$ in (4.8), we have

$$(7.4) \quad [T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})](\xi) = |\sin(\xi/2)|^{2\tau} [T_{\hat{A}}^n 1](\xi) \quad \forall n \in \mathbb{N} \text{ a.e. } \xi \in \mathbb{R}.$$

Since $|\sin(\xi/2)|^{2\tau} \leq 1$, it follows from (7.4) that

$$\|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L^\infty(\mathbb{T})} \leq \|T_{\hat{A}}^n 1\|_{L^\infty(\mathbb{T})}$$

for all $n \in \mathbb{N}$. Consequently, by $\hat{A} \in L^\infty(\mathbb{T})$ and Proposition 7.1, we conclude that

$$(7.5) \quad \rho := \rho_{2\tau}(\hat{a}, \infty) := \limsup_{n \rightarrow \infty} \|[T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})]\|_{L^\infty(\mathbb{T})}^{1/n} \leq \limsup_{n \rightarrow \infty} \|T_{\hat{A}}^n 1\|_{L^\infty(\mathbb{T})}^{1/n}.$$

On the other hand, by (7.4), we have

$$(7.6) \quad [T_{\hat{A}}^n 1](\xi) = \frac{[T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})](\xi)}{|\sin(\xi/2)|^{2\tau}} =: g_n(\xi) \quad \forall n \in \mathbb{N}.$$

By the definition of $T_{\hat{a}}$ and g_n , we have

$$g_n(\xi) = |\hat{A}(\xi/2)|^2 g_{n-1}(\xi/2) + \frac{|\hat{a}(\xi/2 + \pi)|^2}{|\sin(\xi/2)|^{2\tau}} [T_{\hat{a}}^{n-1}(|\sin(\cdot/2)|^{2\tau})](\xi/2 + \pi).$$

Since $\tau \geq 0$, for almost every $\xi \in [-\pi, \pi]$, we have

$$|\hat{a}(\xi/2 + \pi)|^2 / |\sin(\xi/2)|^{2\tau} = |\hat{A}(\xi/2 + \pi)|^2 / |\cos(\xi/4)|^{2\tau} \leq 2^\tau \|\hat{A}\|_{L^\infty(\mathbb{T})},$$

from which we see that

$$g_n(\xi) \leq |\hat{A}(\xi/2)|^2 g_{n-1}(\xi/2) + 2^\tau \|\hat{A}\|_{L^\infty(\mathbb{T})} \|T_{\hat{a}}^{n-1}(|\sin(\cdot/2)|^{2\tau})\|_{L^\infty(\mathbb{T})} \\ \text{a.e. } \xi \in [-\pi, \pi].$$

By induction on n and $g_0 = 1$, we deduce from the above inequality that for almost every $\xi \in [-\pi, \pi]$ and all $n \in \mathbb{N}$,

$$(7.7) \quad g_n(\xi) \leq \prod_{j=1}^n |\hat{A}(2^{-j}\xi)|^2 \\ + 2^\tau \|\hat{A}\|_{L^\infty(\mathbb{T})} \sum_{j=1}^{n-1} \|T_{\hat{a}}^j(|\sin(\cdot/2)|^{2\tau})\|_{L^\infty(\mathbb{T})} \prod_{k=1}^{n-j-1} |\hat{A}(2^{-k}\xi)|^2.$$

Since $|\hat{a}|^2 \in C^\beta(\mathbb{T})$ and $|\hat{a}(0)|^2 = 1$, we define

$$\Phi(\xi) := \prod_{j=1}^{\infty} |\hat{a}(2^{-j}\xi)|^2, \quad \xi \in \mathbb{R}.$$

Then Φ is continuous and $\Phi(0) = 1$. Thus, there exists $0 < \varepsilon_0 < \pi$ such that $1/2 \leq \Phi(\xi) \leq 3/2$ for all $\xi \in [-\varepsilon_0, \varepsilon_0]$. By

$$|\hat{a}(\xi)|^2 = 2^{2\tau} |\cos(\xi/2)|^{2\tau} |\hat{A}(\xi)|^2,$$

for $n \geq \log_2(\pi/\varepsilon_0)$ and $\xi \in [-\pi, \pi]$, we have

$$\begin{aligned} \prod_{j=1}^n |\hat{A}(2^{-j}\xi)|^2 &= \left| \frac{\sin(2^{-1-n}\xi)}{\sin(\xi/2)} \right|^{2\tau} \prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2 \\ &= \left| \frac{\sin(2^{-1-n}\xi)}{\sin(\xi/2)} \right|^{2\tau} \frac{\Phi(\xi)}{\Phi(2^{-n}\xi)} \leq \left| \frac{2^{-1-n}\xi}{\xi/\pi} \right|^{2\tau} \frac{\Phi(\xi)}{1/2} \\ &\leq 2^{1-2\tau} \pi^{2\tau} \|\Phi\|_{L_\infty([-\pi, \pi])} 2^{-2\tau n}. \end{aligned}$$

Therefore, we have

$$(7.8) \quad \prod_{j=1}^n |\hat{A}(2^{-j}\xi)|^2 \leq C_1 2^{-2\tau n} \quad \forall \xi \in [-\pi, \pi]$$

with $C_1 := 2^{1-2\tau} \pi^{2\tau} \|\Phi\|_{L_\infty([-\pi, \pi])} < \infty$. Since $\prod_{j=1}^n |\hat{a}(2^{-j}\xi)|^2 = \Phi(\xi)/\Phi(2^{-n}\xi)$ and $1/2 \leq \Phi(\xi) \leq 3/2$ for $\xi \in [-\varepsilon_0, \varepsilon_0]$, we have

$$\begin{aligned} \|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L_\infty(\mathbb{T})} &\geq \left\| \left| \sin(2^{-1-n}\cdot) \right|^{2\tau} \prod_{j=1}^n |\hat{a}(2^{-j}\cdot)|^2 \right\|_{L_\infty([-\pi, \pi])} \\ &= \left\| \left| \sin(2^{-1-n}\cdot) \right|^{2\tau} \Phi(\cdot)/\Phi(2^{-n}\cdot) \right\|_{L_\infty([-\pi, \pi])} \\ &\geq |\sin(2^{-1-n}\varepsilon_0)|^{2\tau} \Phi(\varepsilon_0)/\Phi(2^{-n}\varepsilon_0) \\ &\geq 3^{-1} \pi^{-2\tau} \varepsilon_0^{2\tau} 2^{-2\tau n}. \end{aligned}$$

By the definition of ρ , it follows from the above inequality that $\rho \geq 2^{-2\tau}$. By (7.8) and the definition of ρ in (7.5), for any $\varepsilon > 0$, there exists a positive constant C with $C \geq C_1$ such that

$$\left\| \prod_{j=1}^n |\hat{A}(2^{-j}\cdot)|^2 \right\|_{L_\infty([-\pi, \pi])} \leq C 2^{-2\tau n} \leq C(\rho + \varepsilon)^n$$

and

$$\|T_{\hat{a}}^n(|\sin(\cdot/2)|^{2\tau})\|_{L_\infty(\mathbb{T})} \leq C(\rho + \varepsilon)^n \quad \forall n \in \mathbb{N}.$$

Now it follows from (7.7) and the above inequalities that

$$\|g_n\|_{L_\infty(\mathbb{T})} = \|g_n\|_{L_\infty([-\pi, \pi])} \leq C(\rho + \varepsilon)^n + 2^\tau \|\hat{A}\|_{L_\infty(\mathbb{T})} C^2 n (\rho + \varepsilon)^{n-1},$$

from which we deduce that $\limsup_{n \rightarrow \infty} \|g_n\|_{L_\infty(\mathbb{T})}^{1/n} \leq \rho + \varepsilon$. Taking $\varepsilon \rightarrow 0$, by the definition of g_n in (7.6), we conclude that

$$\limsup_{n \rightarrow \infty} \|T_{\hat{A}}^n 1\|_{L_\infty(\mathbb{T})}^{1/n} = \limsup_{n \rightarrow \infty} \|g_n\|_{L_\infty(\mathbb{T})}^{1/n} \leq \rho.$$

So, by Proposition 7.1, the proof is completed by the above inequality and (7.5). \square

REFERENCES

- [1] A. S. CAVARETTA, W. DAHMEN, AND C. A. MICCHELLI, *Stationary Subdivision*, Mem. Amer. Math. Soc. 453, AMS, Providence, RI, 1991.
- [2] C. K. CHUI AND J. Z. WANG, *On compactly supported spline wavelets and a duality principle*, Trans. Amer. Math. Soc., 330 (1992), pp. 903–915.
- [3] A. COHEN, *Wavelets and Multiscale Signal Processing*, Chapman and Hall, New York, 1995.
- [4] A. COHEN AND I. DAUBECHIES, *A stability criterion for biorthogonal wavelet bases and their related subband coding scheme*, Duke Math. J., 68 (1992), pp. 313–335.
- [5] A. COHEN AND I. DAUBECHIES, *A new technique to estimate the regularity of refinable functions*, Rev. Mat. Iberoamericana, 12 (1996), pp. 527–591.
- [6] A. COHEN, I. DAUBECHIES, AND J. C. FEAUVEAU, *Biorthogonal bases of compactly supported wavelets*, Comm. Pure Appl. Math., 45 (1992), pp. 485–560.
- [7] I. DAUBECHIES, *Ten Lectures on Wavelets*, CBMS-NSF Regional Conf. Ser. in Appl. Math. 61, SIAM, Philadelphia, 1992.
- [8] I. DAUBECHIES AND J. C. LAGARIAS, *Two-scale difference equations I. Existence and global regularity of solutions*, SIAM J. Math. Anal., 22 (1991), pp. 1388–1410.
- [9] I. DAUBECHIES AND Y. HUANG, *A decay theorem for refinable functions*, Appl. Math. Lett., 7 (1994), pp. 1–4.
- [10] N. DYN AND D. LEVIN, *Subdivision schemes in geometric modeling*, Acta Numer., 11 (2002), pp. 73–144.
- [11] A.-H. FAN AND Q. Y. SUN, *Regularity of Butterworth refinable functions*, Asia J. Math., 5 (2001), pp. 433–440.
- [12] B. HAN, *Analysis and construction of optimal multivariate biorthogonal wavelets with compact support*, SIAM J. Math. Anal., 31 (1999), pp. 274–304.
- [13] B. HAN, *Vector cascade algorithms and refinable function vectors in Sobolev spaces*, J. Approx. Theory, 124 (2003), pp. 44–88.
- [14] B. HAN, *On a conjecture about MRA Riesz wavelet bases*, Proc. Amer. Math. Soc., 134 (2006), pp. 1973–1983.
- [15] B. HAN AND R.-Q. JIA, *Multivariate refinement equations and convergence of subdivision schemes*, SIAM J. Math. Anal., 29 (1998), pp. 1177–1199.
- [16] B. HAN AND R. Q. JIA, *Characterization of Riesz bases of wavelets generated from multiresolution analysis*, Appl. Comput. Harmon. Anal., 23 (2007), pp. 321–345.
- [17] B. HAN AND Z. SHEN, *Wavelets with short support*, SIAM J. Math. Anal., 38 (2006), pp. 530–556.
- [18] C. HERLEY AND M. VETTERLI, *Wavelets and recursive filter banks*, IEEE Trans. Signal Process., 41 (1993), pp. 2536–2556.
- [19] R. Q. JIA, *Subdivision schemes in L_p spaces*, Adv. Comput. Math., 3 (1995), pp. 309–341.
- [20] R. Q. JIA AND C. A. MICCHELLI, *Using the refinement equation for the construction of pre-wavelets II: Power of two*, in Curves and Surfaces, P. J. Laurent, A. Le Méhauté, and L. L. Schumaker, eds., Academic Press, New York, 1991, pp. 209–246.
- [21] R. Q. JIA, J. Z. WANG, AND D. X. ZHOU, *Compactly supported wavelet bases for Sobolev spaces*, Appl. Comput. Harmon. Anal., 15 (2003), pp. 224–241.
- [22] W. LAWTON, S. L. LEE, AND Z. W. SHEN, *Convergence of multidimensional cascade algorithms*, Numer. Math., 78 (1998), pp. 427–438.
- [23] P. G. LEMARIÉ-RIEUSSET, *Fonctions d'échelle interpolantes, polynômes de Bernstein et ondelettes non stationnaires*, Rev. Mat. Iberoamericana, 13 (1997), pp. 91–188.
- [24] R. LORENTZ AND P. OSWALD, *Criteria for hierarchical bases in Sobolev spaces*, Appl. Comput. Harmon. Anal., 8 (2000), pp. 32–85.
- [25] A. OPPENHEIM AND R. SCHAFER, *Digital Signal Processing*, Prentice-Hall, New York, 1975.
- [26] Q. SUN, *Convergence of cascade algorithms and smoothness of refinable distributions*, Chinese Ann. Math. Ser. B, 24 (2003), pp. 367–386.
- [27] M. UNSER AND T. BLU, *Fractional splines and wavelets*, SIAM Rev., 42 (2000), pp. 43–67.
- [28] D. X. ZHOU, *Norms concerning subdivision sequences and their applications in wavelets*, Appl. Comput. Harmon. Anal., 11 (2001), pp. 329–346.

PROPAGATION THROUGH GENERIC LEVEL CROSSINGS: A SURFACE HOPPING SEMIGROUP*

CLOTILDE FERMANIAN KAMMERER[†] AND CAROLINE LASSER[‡]

Abstract. We construct a surface hopping semigroup, which asymptotically describes nuclear propagation through crossings of electron energy levels. The underlying time-dependent Schrödinger equation has a matrix-valued potential, whose eigenvalue surfaces have a generic intersection of codimension two, three, or five in Hagedorn’s classification. Using microlocal normal forms reminiscent of the Landau–Zener problem, we prove convergence to the true solution with an error of the order $\varepsilon^{1/8}$, where ε is the semiclassical parameter. We present numerical experiments for an algorithmic realization of the semigroup illustrating the convergence of the algorithm.

Key words. time-dependent Schrödinger system, eigenvalue crossing, microlocal normal form, surface hopping

AMS subject classifications. 35Q40, 41A60, 81V55

DOI. 10.1137/070686810

1. Introduction. In the framework of time-dependent Born–Oppenheimer approximation, the dynamics of molecules can approximately be reduced to matrix-valued Schrödinger equations on the nucleonic configuration space,

$$(1.1) \quad \begin{cases} i\varepsilon\partial_t\psi^\varepsilon(q, t) = \left(-\frac{\varepsilon^2}{2}\Delta_q + V(q)\right)\psi^\varepsilon(q, t), & (q, t) \in \mathbb{R}^d \times \mathbb{R}, \\ \psi^\varepsilon(q, 0) = \psi_0^\varepsilon(q); \end{cases}$$

see, for example, [12, 19]. The linear Schrödinger equation (1.1) has a unique global solution $\psi^\varepsilon \in C(\mathbb{R}, L^2(\mathbb{R}^d, \mathbb{C}^N))$ for all square-integrable initial data ψ_0^ε . The parameter $\varepsilon > 0$ is small and causes a highly oscillatory behavior of the solution in space and time. It can be thought of as the square root of the ratio of electronic mass and the average mass of the nuclei. Moreover, the solution itself does not have any direct physical interpretation. It is the position density $|\psi^\varepsilon(q, t)|^2$ which gives the probability of finding the nuclei in the configuration $q \in \mathbb{R}^d$ at time t . We are interested in an asymptotic description for the time evolution of quadratic quantities like the position density with the following properties. First, it shall be effective in the sense that it unfolds characteristic dynamical properties. Second, it shall be explicit enough, such that it allows an algorithmic realization. Third, the resulting algorithm shall be applicable on high-dimensional nucleonic configuration spaces \mathbb{R}^d , $d \gg 1$.

Hagedorn rigorously derived and classified Schrödinger systems for molecular propagation through electron energy level crossings of minimal multiplicity [13]. He obtained potentials of the form

$$V(q) = v(q) \text{Id} + V_\ell(\phi(q)), \quad \ell \in \{2, 3, 3', 5\},$$

*Received by the editors March 30, 2007; accepted for publication (in revised form) December 3, 2007; published electronically April 2, 2008.

<http://www.siam.org/journals/sima/40-1/68681.html>

[†]Université de Paris 12 Val de Marne, Laboratoire d’analyse et de mathématiques appliquées, 61 avenue du général de Gaulle, 94 010 Créteil cedex, France (clotilde.fermanian@univ-paris12.fr).

[‡]Freie Universität Berlin, Fachbereich Mathematik und Informatik, Arnimallee 6, 14195 Berlin, Germany (lasser@math.fu-berlin.de).

where $v(q) \in \mathcal{C}^\infty(\mathbb{R}^d, \mathbb{R})$ is a smooth real-valued function, Id is the identity matrix in $\mathbb{C}^{2 \times 2}$ or $\mathbb{C}^{4 \times 4}$, and $q \mapsto \phi(q)$ is a smooth vector-valued function with $\phi(q) \in \mathbb{R}^2, \mathbb{R}^3$, or \mathbb{R}^5 . The matrices V_ℓ are given by

$$(1.2) \quad V_2(\phi) = \begin{pmatrix} \phi_1 & \phi_2 \\ \phi_2 & -\phi_1 \end{pmatrix}, \quad V_3(\phi) = \begin{pmatrix} \phi_1 & \phi_2 + i\phi_3 \\ \phi_2 - i\phi_3 & -\phi_1 \end{pmatrix},$$

$$(1.3) \quad V_{3'}(\phi) = \begin{pmatrix} \begin{pmatrix} \phi_1 & \phi_2 + i\phi_3 \\ \phi_2 - i\phi_3 & -\phi_1 \end{pmatrix} & \mathbf{0} \\ \mathbf{0} & \begin{pmatrix} \phi_1 & \phi_2 - i\phi_3 \\ \phi_2 + i\phi_3 & -\phi_1 \end{pmatrix} \end{pmatrix},$$

$$(1.4) \quad V_5(\phi) = \begin{pmatrix} \phi_0 \text{Id} & \begin{pmatrix} \phi_1 + i\phi_2 & \phi_3 + i\phi_4 \\ -\phi_3 + i\phi_4 & \phi_1 - i\phi_2 \end{pmatrix} \\ \begin{pmatrix} \phi_1 - i\phi_2 & -\phi_3 - i\phi_4 \\ \phi_3 - i\phi_4 & \phi_1 + i\phi_2 \end{pmatrix} & -\phi_0 \text{Id} \end{pmatrix}.$$

For these four matrices, the eigenvalues are $\pm|\phi|$. Therefore, the eigenvalues of $V(q)$ are $v(q) \pm |\phi(q)|$, and $q^* \in \mathbb{R}^d$ is a point of crossing eigenvalues if and only if $\phi(q^*) = 0$. We shall say that this crossing is generic if

$$d\phi \text{ is of maximal rank on } \{\phi = 0\},$$

i.e., of rank 2 for $\ell = 2$, of rank 3 for $\ell = 3, 3'$, and of rank 5 for $\ell = 5$. This explains why these crossings are usually referred to as codimension two, three, and five crossings and enlightens the choice of the index ℓ we have made. Hagedorn's codimension one crossings are not considered here, since they violate the above rank condition and also show a different dynamical behavior than systems with crossings of higher codimension. We set

$$N(2) = N(3) = 2, \quad N(3') = N(5) = 4,$$

so that the potential $V(q)$, the wave function $\psi^\varepsilon(q, t)$, and the differential $d\phi(q)$ belong to $\mathbb{C}^{N(\ell) \times N(\ell)}$, $\mathbb{C}^{N(\ell)}$, and $\mathbb{R}^{\ell \times d}$, respectively. For $\ell = 3'$ we set $\mathbb{R}^{3'} = \mathbb{R}^3$. The orthogonal eigenprojectors

$$\Pi^\pm(q) = \frac{1}{2} (\text{Id} \pm |\phi(q)|^{-1} V_\ell(\phi(q)))$$

of the matrix $V(q)$ have a conical singularity at points of crossing eigenvalues q^* ; that is, $\nabla \Pi^\pm(q) = O(|q - q^*|^{-1})$ as $q \rightarrow q^*$. This motivates the notion of conical intersections, by which especially codimension two crossings are frequently referred to.

Eigenvalue crossings are ubiquitous in the quantum mechanical description of polyatomic molecules, that is, molecules with more than two nuclei. The collection [4] provides an exposition of this active area of research in theoretical chemistry. As for a prominent example of an ultrafast isomerization on the femtosecond time scale, a codimension two crossing of energy levels explains the effectiveness of the first step of vision, the cis-trans isomerization of retinal in rhodopsin; see also [14] and section 3 below for related numerical experiments.

The analysis of scalar Schrödinger equations teaches us that the direct study of the time evolution of quadratic quantities like the position density $|\psi^\varepsilon(q, t)|^2$ is

impossible. The oscillations of $\psi^\varepsilon(q, t)$ have to be taken into account, and one has to work in the space of positions and momenta (q, p) , the phase space $\mathbb{R}_q^d \times \mathbb{R}_p^d$. Therefore, one studies the Wigner function of $\psi^\varepsilon(q, t)$ in a suitable ε -dependent scaling, which resolves the highly oscillatory features of the solution,

$$W^\varepsilon(\psi^\varepsilon(t))(q, p) = (2\pi)^{-d} \int_{\mathbb{R}^d} \psi^\varepsilon\left(q - \frac{\varepsilon}{2}v, t\right) \otimes \overline{\psi^\varepsilon}\left(q + \frac{\varepsilon}{2}v, t\right) e^{i v \cdot p} dv.$$

It plays the role of a generalized probability density on phase space. For square-integrable wave functions ψ , the Wigner function $W^\varepsilon(\psi)$ is a square-integrable function on phase space with values in the space of hermitian matrices. One recovers the position density by

$$|\psi(q)|^2 = \text{tr} \int_{\mathbb{R}^d} W^\varepsilon(\psi)(q, p) dp.$$

Besides, the action of the Wigner function against compactly supported smooth test functions $a \in C_c^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N(\ell) \times N(\ell)})$ is simply expressed in terms of the semiclassical pseudodifferential operator with symbol a , which is defined by

$$\text{op}_\varepsilon(a)\psi(q) = (2\pi\varepsilon)^{-d} \int_{\mathbb{R}^{2d}} a\left(\frac{1}{2}(q+v), p\right) e^{\frac{i}{\varepsilon}p \cdot (q-v)} \psi(v) dv dp$$

for $\psi \in L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})$. Indeed, we have

$$\text{tr} \int_{\mathbb{R}^{2d}} W^\varepsilon(\psi)(q, p) a(q, p) dq dp = (\text{op}_\varepsilon(a)\psi, \psi)_{L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})}.$$

It is our aim to construct an asymptotic semigroup, which approximately propagates the initial data's Wigner function for all generic level crossings of Hagedorn's classification. Our approximation relies on a microlocal normal form for operators with eigenvalue crossings, which has been derived in [1, 2, 6]. Roughly speaking, near a crossing point the Schrödinger operator

$$-i\varepsilon\partial_t - \frac{\varepsilon^2}{2}\Delta_q + V(q)$$

is equivalent to the normal form

$$-i\varepsilon\partial_t + V_\ell\left(t, \text{op}_\varepsilon(|d\phi(q)p|^{-\frac{1}{2}}\pi_\ell(q, p)\phi(q))\right),$$

where $\pi_\ell(q, p)$ denotes the orthogonal projection onto the hyperplane normal to the vector $d\phi(q)p \in \mathbb{R}^\ell$. In the case $\ell = 2$, this resembles the Landau–Zener system

$$i\varepsilon\frac{d}{dt}\psi(t) = \begin{pmatrix} t & \gamma \\ \gamma & -t \end{pmatrix} \psi(t), \quad \gamma > 0,$$

for which the probability that a solution starting at time $t = -\infty$ in the one eigenspace will have passed over to the other eigenspace at time $t = +\infty$, which has explicitly been computed by Landau [16] and Zener [21] in the 1930s. This famous Landau–Zener transition rate reads as

$$\exp\left(-\frac{\pi}{\varepsilon}\gamma^2\right),$$

and in [8] it is proven that the rate still gives the correct asymptotics if γ is replaced by a bounded operator. Our semigroup combines effective transitions between eigenspaces close to points of crossing eigenvalues on the one hand with classical transport in the adiabatic regime on the other hand. More precisely, the V -diagonal components $\Pi^\pm W^\varepsilon(\psi_0^\varepsilon)\Pi^\pm$ of the initial Wigner function are transported along the Hamiltonian curves of the eigenvalues of the Schrödinger operator's symbol

$$\frac{1}{2}|p|^2 + v(q) \pm |\phi(q)|.$$

Whenever one of the trajectories $(q^\pm(t), p^\pm(t))$ attains a local minimal gap between the eigenvalues, there is an effective nonadiabatic transfer of weight according to the ε -dependent transition rate

$$\exp\left(-\frac{\pi}{\varepsilon} \frac{|\pi_\ell(q, p)\phi(q)|^2}{|d\phi(q)p|}\right).$$

Since the rate is negligibly small, when the eigenvalue gap is larger than $\sqrt{\varepsilon}$, the nonadiabatic transfer of weight is effectively performed at times t^* with

$$t \mapsto |\phi(q^\pm(t))| \text{ has a local minimum in } t = t^* \text{ and } |\phi(q^\pm(t^*))| \leq R\sqrt{\varepsilon}$$

for some fixed $R > 0$. Our main result is that this dynamical description yields approximate solutions with an error of order $\varepsilon^{1/8}$ when choosing $R = \varepsilon^{-1/8}$. Moreover, it is explicit enough for an algorithmic realization, whose performance on a model for retinal in rhodopsin is studied here as well. The algorithm is a mathematical counterpart to the popular surface hopping algorithms of chemical physics introduced by Tully and Preston in [20].

Quantum dynamical descriptions in terms of classical transport as described above, that is, in the spirit of an Egorov theorem, are well established and have been given for Wigner functions, for example, in [10, 11]: for Schrödinger systems they hold to leading order in ε , until classical trajectories come close to a point of crossing eigenvalues. Then, as already mentioned, the adiabatic approximation is no longer valid, and there are leading order nonadiabatic transitions between the levels (the energy propagated on one level to the crossing may pass partially or completely on the other level). This phenomenon has been precisely analyzed in the case of Gaussian wave packet propagation by Hagedorn [13] for all generic electron level crossings. For initial data, which are less specific than Gaussian wave packets, the evolution of appropriate two-scale Wigner measures has been studied. These measures are weak limits of the Wigner function and incorporate information on concentration effects close to trajectories, which touch points of crossing eigenvalues, with respect to the second scale $\sqrt{\varepsilon}$. These Wigner measures have been analyzed for a linear codimension two crossing in [7], for general two-level systems in [8], and for all of Hagedorn's models in [5]. In [18], the results of [7] have been lifted to a leading order approximation of the Wigner function. Here, we aim at approximating the Wigner function for all generic crossings, while additionally proving a convergence rate.

We will proceed as follows. Section 2 constructs the surface hopping semigroup, states the main result, that is, the validity of our approximation with an error of order $\varepsilon^{1/8}$, and discusses the strategy of the proof. In section 3, numerical results are presented for an algorithmic realization of the semigroup applied to a retinal model. In section 4, the proof for propagation away from the crossing is carried out, while the microlocal normal form yields the correct nonadiabatic transition rates, as proven

in section 5. In section 6, the main result is extended to observables, which are more pertinent for the crossings with degenerate eigenvalues ($\ell = 3', 5$). Finally, the appendix summarizes basic facts of Weyl calculus.

2. Main result. Propagation through level crossings can be approximated by a proper combination of classical transport and nonadiabatic transitions. For this, we study the underlying classical flows and combine them with effective nonadiabatic transitions to an asymptotic semigroup.

2.1. Transport and transitions. We consider the classical flows

$$\Phi_{\pm}^{-t} : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}, \quad \Phi_{\pm}^{-t}(q_0, p_0) = (q^{\pm}(t), p^{\pm}(t))$$

associated with the Hamiltonian curves of $\frac{1}{2}|p|^2 + v(q) \pm |\phi(q)|$. These curves are solutions to the Hamiltonian systems

$$\begin{cases} \dot{q}^{\pm}(t) = p^{\pm}(t), & \dot{p}^{\pm}(t) = -\nabla v(q^{\pm}(t)) \mp {}^t d\phi(q^{\pm}(t)) \frac{\phi(q^{\pm}(t))}{|\phi(q^{\pm}(t))|}, \\ q^{\pm}(0) = q_0, & p^{\pm}(0) = p_0. \end{cases}$$

We consider only initial phase space points $(q_0, p_0) \in \mathbb{R}^{2d}$ such that for $t > 0$

$$(2.1) \quad \phi(q^{\pm}(t)) = 0 \Rightarrow d\phi(q^{\pm}(t))p^{\pm}(t) \neq 0.$$

This condition guarantees that classical trajectories arrive transversally at the crossing set and have a unique smooth continuation through this singularity; see Proposition 1 in [5].

For a large class of test functions and under suitable restrictions on the time interval, the classical flows are enough for approximating the dynamics up to an error of order ε . Indeed, one considers observables $a \in \mathcal{C}_c^{\infty}(\mathbb{R}^{2d}, \mathbb{C}^{N(\ell) \times N(\ell)})$ such that

$$(2.2) \quad a = a^+ \Pi^+ + a^- \Pi^-, \quad a^{\pm} \in \mathcal{C}_c^{\infty}(\mathbb{R}^{2d} \setminus \{\phi = 0\}, \mathbb{C}).$$

For $\ell = 2, 3$, the eigenspaces are one-dimensional, and these observables focus on the V -diagonal elements of the Wigner matrix, where V -diagonal means diagonal with respect to the decomposition of $\mathbb{C}^{N(\ell)}$ by the eigenprojectors $\Pi^+(q)$ and $\Pi^-(q)$. For $\ell = 3', 5$, however, the eigenspaces are two-dimensional, and observables of the form (2.2) are not enough to completely resolve all dynamical features within the eigenspaces. We will address this issue in section 6.

For all times $t \in [0, T]$, such that the classical trajectories Φ_{\pm}^t arriving on the support of a have not passed the crossing set $\{\phi = 0\}$, the action of the Wigner function on $a = a^{\pm} \Pi^{\pm}$ obeys

$$\begin{aligned} \text{tr} \int_{\mathbb{R}^{2d}} W^{\varepsilon}(\psi^{\varepsilon}(t))(q, p) \Pi^{\pm}(q) a^{\pm}(q, p) dq dp \\ - \text{tr} \int_{\mathbb{R}^{2d}} (\Pi^{\pm} W^{\varepsilon}(\psi_0^{\varepsilon}) \Pi^{\pm} \circ \Phi_{\pm}^{-t})(q, p) a^{\pm}(q, p) dq dp = O(\varepsilon) \end{aligned}$$

as $\varepsilon \rightarrow 0$. Such Egorov-type descriptions hold, until classical trajectories come close to a crossing point $(q, p) \in \{\phi = 0\}$ and leading order nonadiabatic transitions occur. These transitions depend on how the solution $\psi^{\varepsilon}(t)$ concentrates on the ingoing trajectories with respect to the scale $\sqrt{\varepsilon}$. For the linear codimension two crossing with $\phi(q) = q$, $q \in \mathbb{R}^2$, the two-scale Wigner measure's description of [7] is lifted

to an approximation of the Wigner function in [18]. This linear model has specific features; see also [9]. In particular, all classical trajectories which meet the crossing are included in the set $\{q \wedge p = 0\}$, where $q \wedge p := q_2 p_1 - q_1 p_2$ for $q, p \in \mathbb{R}^2$. The idea of [18] is to propagate the V -diagonal parts of the initial Wigner function along the classical trajectories and to apply the ε -dependent transition coefficient

$$T_{lin}^\varepsilon(q^*, p^*) = \exp\left(-\frac{\pi}{\varepsilon} \frac{|q^* \wedge p^*|^2}{|p^*|^3}\right),$$

as soon as the trajectories reach their minimal distance from the crossing set, which is easy to check since $q \cdot p = 0$ at such a point. Theorem 3.2 in [18] proves that under suitable conditions on the initial data this ε -dependent propagation is correct in the limit $\varepsilon \rightarrow 0$. We construct here an extension, which covers the general situation described above, and give a convergence proof including a convergence rate. Our approach draws from the understanding of the nonadiabatic mechanism as developed in [5].

2.2. A surface hopping semigroup. Let $R > 0$. In the general case, the crucial points in phase space are those where the classical trajectories attain a *local minimal gap between the two eigenvalues*. These points fulfill the condition

$$\frac{d}{dt} \left(|\phi(q^\pm(t))|^2 \right) = 2 \, d\phi(q^\pm(t)) p^\pm(t) \cdot \phi(q^\pm(t)) = 0,$$

and one performs an effective nonadiabatic transfer of weight, whenever a trajectory passes the set

$$S_{\varepsilon,R} = \{(q, p) \in \mathbb{R}^{2d} \mid |\phi(q)| \leq R\sqrt{\varepsilon}, \, d\phi(q)p \cdot \phi(q) = 0\}.$$

The microlocal normal form, which will be given later in Theorem 5.2, suggests the transition rate

$$T^\varepsilon(q^*, p^*) = \exp\left(-\frac{\pi}{\varepsilon} \frac{|\pi_\ell(q^*, p^*)\phi(q^*)|^2}{|d\phi(q^*)p^*|}\right),$$

where $\pi_\ell(q^*, p^*)$ is the orthogonal projection from the Euclidean space \mathbb{R}^ℓ into the hyperplane normal to the vector $d\phi(q^*)p^* \in \mathbb{R}^\ell$. Since $d\phi(q^\pm(t))p^\pm(t)$ does not vanish when the considered trajectories arrive at the crossing set $\{\phi = 0\}$, it is also nonzero when arriving at the jump manifold $S_{\varepsilon,R}$ if $R\sqrt{\varepsilon}$ is small enough. Besides, for $\ell = 2$, one has

$$|\pi_\ell(q, p)\phi(q)| = |\phi(q) \wedge \frac{d\phi(q)p}{|d\phi(q)p|}|,$$

and we recover the transition coefficient $T_{lin}^\varepsilon(q^*, p^*)$ for $\phi(q) = q$, $q \in \mathbb{R}^2$.

We attach the labels -1 and $+1$ to phase space. For points $(q, p, j) \in \mathbb{R}_\pm^{2d} := \mathbb{R}^{2d} \times \{-1, +1\}$, we consider trajectories

$$\mathcal{T}_{\varepsilon,R}^{(q,p,j)} : [0, +\infty) \rightarrow \mathbb{R}_\pm^{2d},$$

which combine deterministic classical transport and random jumps between the levels at the manifold $S_{\varepsilon,R}$. More precisely, we set $\mathcal{T}_{\varepsilon,R}^{(q,p,j)}(t) = (\Phi_j^t(q, p), j)$ as long as $\Phi_j^t(q, p) \notin S_{\varepsilon,R}$. Whenever the deterministic flow $\Phi_j^t(q, p)$ hits the manifold $S_{\varepsilon,R}$ at

a point (q^*, p^*) , a random jump from j to $-j$ occurs with probability $T^\varepsilon(q^*, p^*)$. For points (q, p, j) generating classical trajectories, which either violate the non-degeneracy condition (2.1) or do not leave the set $S_{\varepsilon, R}$, there is either no transport or no jump at all. Since the trajectories which touch the crossing set arrive there transversally, each path

$$(q, p, j) \rightarrow \mathcal{T}_{\varepsilon, R}^{(q, p, j)}(t)$$

has a finite number of jumps and remains in a bounded region of \mathbb{R}_{\pm}^{2d} within a bounded time interval $[0, T]$. Away from the jump manifold $S_{\varepsilon, R} \times \{-1, +1\}$ each path is smooth. Hence, the random trajectories define a time-dependent Markov process with state space \mathbb{R}_{\pm}^{2d} . The associated transition function $P_{\varepsilon, R}(p, q, j; t, \Gamma)$ describes the probability of being at time t in the measurable set $\Gamma \subset \mathbb{R}_{\pm}^{2d}$ having started in (q, p, j) . Its action on bounded measurable scalar functions $f : \mathbb{R}_{\pm}^{2d} \rightarrow \mathbb{C}$ defines a semigroup $(\mathcal{L}_{\varepsilon, R}^t)_{t \geq 0}$ by

$$(\mathcal{L}_{\varepsilon, R}^t f)(q, p, j) := \int_{\mathbb{R}^{2d} \times \{-1, +1\}} f(x, \xi, k) P_{\varepsilon, R}(q, p, j; t, d(x, \xi, k)).$$

For introducing the semigroup's action on Wigner functions, we use the following space of continuous V -diagonal test functions satisfying T^ε -dependent boundary conditions at the jump manifold.

DEFINITION 2.1. *A continuous function $a \in C_c(\mathbb{R}^{2d} \setminus S_{\varepsilon, R}, \mathbb{C}^{N(\ell)} \times N(\ell))$ belongs to the space $\mathcal{C}_{\varepsilon, R}$ if it has the following properties:*

- i. $a = a^+ \Pi^+ + a^- \Pi^-$ with $a^\pm \in C_c(\mathbb{R}^{2d} \setminus S_{\varepsilon, R}, \mathbb{C})$.
- ii. The function $f_a : (\mathbb{R}^{2d} \setminus S_{\varepsilon, R}) \times \{-1, +1\} \rightarrow \mathbb{C}$,

$$f_a(q, p, +) = a^+(q, p), \quad f_a(q, p, -) = a^-(q, p),$$

satisfies for all $(q, p, j) \in S_{\varepsilon, R} \times \{-1, +1\}$

$$\begin{aligned} & \lim_{\delta \rightarrow -0} f_a(q + \delta p, p + \delta(-\nabla v(q) - j^t d\phi(q)\phi(q)/|\phi(q)|), j) \\ &= \lim_{\delta \rightarrow +0} (T^\varepsilon f_a)(q + \delta p, p + \delta(-\nabla v(q) + j^t d\phi(q)\phi(q)/|\phi(q)|), -j) \\ &= \lim_{\delta \rightarrow +0} ((1 - T^\varepsilon) f_a)(q + \delta p, p + \delta(-\nabla v(q) - j^t d\phi(q)\phi(q)/|\phi(q)|), j). \end{aligned}$$

For test functions $a \in \mathcal{C}_{\varepsilon, R}$, the action of $(\mathcal{L}_{\varepsilon, R}^t)_{t \geq 0}$ is naturally given by

$$(\mathcal{L}_{\varepsilon, R}^t a)(q, p) := (\mathcal{L}_{\varepsilon, R}^t f_a)(q, p, +1) \Pi^+(q) + (\mathcal{L}_{\varepsilon, R}^t f_a)(q, p, -1) \Pi^-(q).$$

By construction, the semigroup leaves $\mathcal{C}_{\varepsilon, R}$ invariant, and duality allows us to define its action on Wigner functions. More precisely, let $W^\varepsilon(\psi)$ be the Wigner function of some wave function $\psi \in L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})$. Then, $\mathcal{L}_{\varepsilon, R}^t W^\varepsilon(\psi)$ acts on $a \in \mathcal{C}_{\varepsilon, R}$ by

$$(\mathcal{L}_{\varepsilon, R}^t W^\varepsilon(\psi), a) = \text{tr} \int_{\mathbb{R}^{2d}} W^\varepsilon(\psi)(q, p) (\mathcal{L}_{\varepsilon, R}^t a)(q, p) dq dp,$$

defining a locally integrable function on phase space. We finally choose an ε -dependent hopping range $R(\varepsilon) = \varepsilon^{-1/8}$ and set

$$(\mathcal{L}_\varepsilon^t)_{t \geq 0} := (\mathcal{L}_{\varepsilon, R(\varepsilon)}^t)_{t \geq 0}.$$

2.3. Assumptions and main result. Let $\psi^\varepsilon(t)$ be the solution of the Schrödinger equation (1.1) with initial datum ψ_0^ε . We now state the precise assumptions, under which the action of the semigroup $(\mathcal{L}_\varepsilon^t)_{t \geq 0}$ on the initial Wigner function $W^\varepsilon(\psi_0^\varepsilon)$ approximates the V -diagonal components of $W^\varepsilon(\psi^\varepsilon(t))$.

(A0) $V \in \mathcal{C}^\infty(\mathbb{R}^d, \mathbb{C}^{N(\ell) \times N(\ell)})$ is of subquadratic growth and of the form

$$V(q) = v(q) \text{Id} + V_\ell(\phi(q)), \quad \ell \in \{2, 3, 3', 5\},$$

where the matrices $V_\ell(\phi(q))$ have been defined in (1.2), (1.3), and (1.4). We assume the eigenvalue crossings to be generic in the sense that $d\phi$ is of maximal rank on the crossing set $\{\phi = 0\}$.

(A1) $(\psi_0^\varepsilon)_{\varepsilon > 0}$ is a bounded family in $L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})$ associated with $\text{Ran} \Pi^+$,

$$\|\Pi^- \psi_0^\varepsilon\|_{L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})} = O(\varepsilon^{\beta_1}), \quad \beta_1 \geq 1/8.$$

We suppose that the initial data are localized away from the crossing $\{\phi = 0\}$; i.e., for all $b \in \mathcal{C}_c^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N(\ell) \times N(\ell)})$ with $\text{supp}(b) \subset \{|\phi| \leq R\sqrt{\varepsilon}\}$, $R = \varepsilon^{-1/8}$,

$$\int_{\mathbb{R}^{2d}} W^\varepsilon(\psi_0^\varepsilon)(q, p) b(q, p) dq dp = O(\varepsilon^{\beta_2}), \quad \beta_2 \geq 1/8.$$

We also assume localization away from the set

$$\{(q_0, p_0) \in \mathbb{R}^{2d} \mid \exists t > 0 : \phi(q^\pm(t)) = 0, d\phi(q^\pm(t))p^\pm(t) = 0\},$$

which contains the points issuing classical trajectories, which arrive at the crossing without a unique continuation through it.

(A2) The test function $a \in \mathcal{C}_c^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N(\ell) \times N(\ell)})$ has its support at a distance larger than $R\sqrt{\varepsilon}$ with $R = \varepsilon^{-1/8}$ from the crossing; that is,

$$\text{supp}(a) \cap \{(q, p) \in \mathbb{R}^{2d} \mid |\phi(q)| \leq R\sqrt{\varepsilon}\} = \emptyset, \quad R = \varepsilon^{-1/8},$$

and

$$a(q, p) = a^+(q, p)\Pi^+(q) + a^-(q, p)\Pi^-(q), \quad (q, p) \in \mathbb{R}^{2d},$$

with scalar-valued $a^\pm \in \mathcal{C}_c^\infty(\mathbb{R}^{2d}, \mathbb{C})$.

(A3) Within the time interval $[0, T]$, each of the plus-trajectories arriving at the support of a^+ at time T has performed at most one jump, generating minus-trajectories arriving at the support of a^- , which have not jumped at all.

Alternatively, assumptions (A1) and (A3) could also require that the initial data are associated with $\text{Ran} \Pi^-$ and that each of the minus-trajectories arriving at the support of a^- at time T has performed at most one jump, generating plus-trajectories arriving at the support of a^+ , which have not jumped at all.

THEOREM 2.2. *Let the potential V , the initial data $(\psi_0^\varepsilon)_{\varepsilon > 0}$, the observable a , and the time interval $[0, T]$ fulfill assumptions (A0), (A1), (A2), and (A3). Let $\chi \in \mathcal{C}_c^\infty([0, T], \mathbb{R})$. Then, there exist positive constants $C, \varepsilon_0 > 0$ such that for all $0 < \varepsilon < \varepsilon_0$ the solution $\psi^\varepsilon(t)$ of the Schrödinger equation (1.1) satisfies*

$$(2.3) \quad \left| \text{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) (W^\varepsilon(\psi^\varepsilon(t)) - \mathcal{L}_\varepsilon^t W^\varepsilon(\psi_0))(q, p) a(q, p) dq dp dt \right| \leq C \varepsilon^{1/8}.$$

Before entering the proof, we add some remarks. First, if one allows initial data in assumption (A1) with $\beta_1, \beta_2 > 0$, then the result holds with convergence

rate $\varepsilon^{\min(\beta_1, \beta_2, 1/8)}$. Second, pointwise convergence holds on time intervals without nonadiabatic jumps, that is, when the solution has passed by the jump manifold; see also [18]. However, for pointwise convergence only the limit behavior without convergence rate can be deduced, since the constants C and ε_0 depend on the cut-off function χ and its derivatives in a way that possible oscillations in time are not controlled. Finally, in section 6 an extension of the approximation for the cases $\ell = 3', 5$ with degenerate eigenvalues is given. There, assumption (A2) is generalized to observables a which commute with V , that is, $a = \Pi^+ a \Pi^+ + \Pi^- a \Pi^-$.

2.4. Strategy of the proof. For notational convenience, we suppose $\beta_1 = \beta_2 = 1/2$. Otherwise, one has to add $O(\varepsilon^{1/8})$ in all the estimates. We work with the semigroup $(\mathcal{L}_{\varepsilon, R}^t)_{t \geq 0}$ and prove convergence with an error of order

$$O(1/(R^5 \sqrt{\varepsilon})) + O(R^3 \sqrt{\varepsilon}) + O(1/R^2) + O(\sqrt{\varepsilon} |\ln \varepsilon|)$$

as $\varepsilon \rightarrow 0$ and $R \rightarrow +\infty$, which gives the claimed rate when choosing $R = \varepsilon^{-1/8}$. We distinguish between regions of large and small eigenvalue gap, that is, between sets $\{|\phi| > C R \sqrt{\varepsilon}\}$ with $C = \frac{1}{2}, 1$ on the one hand and $\{|\phi| \leq R \sqrt{\varepsilon}\}$ on the other hand. For a large gap we prove classical transport with an error of size $O(1/(R^5 \sqrt{\varepsilon})) + O(1/R^2) + O(\sqrt{\varepsilon})$. Close to the crossing set, proving the relevance of nonadiabatic transitions, we use a microlocal normal form, which reduces the Schrödinger equation to a Landau–Zener-type problem with explicitly computable transition rates. There, we introduce an error of order $O(R^3 \sqrt{\varepsilon}) + O(1/R^2) + O(\sqrt{\varepsilon} |\ln \varepsilon|) + O(1/(R^5 \sqrt{\varepsilon}))$. The combination of both errors will then yield the final estimate of Theorem 2.2.

PROPOSITION 2.3. *Let $c \in \mathcal{C}_c^\infty(\mathbb{R}^{2d}, \mathbb{C})$, and let $b \in \mathcal{C}^\infty(\mathbb{R}^\ell, \mathbb{C})$ with ∇b compactly supported. If there exist $C > 0$ and $s_0 > 0$ such that*

$$\forall r \in [-s_0, s_0] : \quad \Phi_{\pm}^r(\text{supp}(c)) \subset \{|\phi| > C R \sqrt{\varepsilon}\},$$

then for all $\chi \in \mathcal{C}_c^\infty(\mathbb{R}, \mathbb{R})$ and for all $s \in [t - s_0, t + s_0]$

$$\begin{aligned} & \text{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) c(q, p) b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right) \Pi^\pm(q) W^\varepsilon(\psi^\varepsilon(t))(q, p) dq dp dt \\ &= \text{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) c(q, p) b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right) (\Pi^\pm W^\varepsilon(\psi^\varepsilon(s)) \Pi^\pm \circ \Phi_{\pm}^{-t+s})(q, p) dq dp dt \\ & \quad + O\left(\frac{1}{R^5 \sqrt{\varepsilon}}\right) + O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}). \end{aligned}$$

Proposition 2.3 will be proved in section 4. To use it for the main proof, we need to specify which points of $\text{supp}(a^\pm)$ arrive close to the crossing. We denote the sets of trajectories arriving at (respectively, arising from) the crossing set $\{\phi = 0\}$ by

$$M^{\pm, in} = \{\Phi_{\pm}^t(q, p) \in \mathbb{R}^{2d} \mid \Phi_{\pm}^t(q, p) \notin \{\phi = 0\}, \exists t_0 < t : \Phi_{\pm}^{t_0}(q, p) \in \{\phi = 0\}\},$$

$$M^{\pm, out} = \{\Phi_{\pm}^t(q, p) \in \mathbb{R}^{2d} \mid \Phi_{\pm}^t(q, p) \notin \{\phi = 0\}, \exists t_0 > t : \Phi_{\pm}^{t_0}(q, p) \in \{\phi = 0\}\}.$$

The sets $M^{\pm, in/out}$ are smooth submanifolds of \mathbb{R}^{2d} . By construction of the semigroup, all phase space points generating backward trajectories passing through the zone of small gap $\{|\phi| \leq R \sqrt{\varepsilon}\}$ and performing a jump are contained in a neighborhood Ω^\pm of the intersection of the support of a^\pm with $M^{\pm, out}$.

Some of the random trajectories reaching Ω^\pm touch the crossing set. We consider one of them, which arrives at time t_0 at (q_0, p_0) with $\phi(q_0) = 0$, $d\phi(q_0)p_0 \neq 0$,

and choose the associated point (q_0, t_0, p_0, τ_0) in the phase space of space-time, where $\tau_0 = -\frac{1}{2}|p_0|^2 - v(q_0)$ is the energy coordinate. The normal form theorems [1, 2, 6] give neighborhoods of these points, on which the Schrödinger equation (1.1) microlocally reduces to a Landau–Zener-type problem

$$-i\varepsilon\partial_s v^\varepsilon = V_\ell(s, \tilde{z} + \gamma_\varepsilon(z, \zeta))v^\varepsilon + O(\varepsilon^\infty).$$

These model problems have explicitly computable transition rates in the scattering regime; see [7, 8]. The compact subset of $\{|\phi| < R\sqrt{\varepsilon}\}$, which is touched by the backward trajectories coming from Ω^\pm , can be covered by finitely many of these neighborhoods projected to \mathbb{R}^{2d} , and without loss of generality we assume that one of them suffices.

Moreover, for each point in Ω^\pm being reached by a random trajectory at time t there are positive numbers $0 < t_f^\pm < t_i^\pm$, such that at time $t - t_i^\pm$ and $t - t_f^\pm$ the trajectories are contained in an annulus $\{C_1\sqrt{\varepsilon} < |\phi| < C_2\sqrt{\varepsilon}\}$ with $C_1, C_2 > 0$ and have performed their only jump within the interval $]t - t_i^\pm, t - t_f^\pm[$, whose length is denoted by $\delta_t^\pm = t_i^\pm - t_f^\pm$. These quantities are well defined, since the trajectories are transverse to the crossing set. Choosing $C_1 = \frac{R}{2}$ and $C_2 = R$, Ω^\pm can be covered by finitely many open sets, such that each of these sets can be associated with such points of time t_i^\pm and t_f^\pm . Without loss of generality, we assume that one of them is enough. Then, we have by Proposition 2.3 with $C = \frac{1}{2}$

$$\begin{aligned} & \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) a^\pm(q, p) \Pi^\pm(q) W^\varepsilon(\psi^\varepsilon(t))(q, p) dq dp dt \\ &= \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) a^\pm(q, p) \left(\Pi^\pm W^\varepsilon(\psi^\varepsilon(t - t_f^\pm)) \Pi^\pm \circ \Phi_\pm^{-t_f^\pm} \right) (q, p) dq dp dt \\ & \quad + O(1/(R^5\sqrt{\varepsilon})) + O(1/R^2) + O(\sqrt{\varepsilon}). \end{aligned}$$

Then, we perform a cut-off of the symbol $a^\pm \circ \Phi_\pm^{t_f^\pm}$. We consider a smooth compactly supported function $\chi_0 \in C_c^\infty(\mathbb{R}^\ell, \mathbb{R})$ such that $\chi_0(u) = 1$ on $\{|u| < 1\}$ and $\chi_0(u) = 0$ for $\{|u| > 2\}$. We write

$$\begin{aligned} \left(a^\pm \circ \Phi_\pm^{t_f^\pm} \right) (q, p) &= a_{BO}^\pm(q, p) + a_{LZ}^\pm(q, p), \\ a_{BO}^\pm(q, p) &= \left(a^\pm \circ \Phi_\pm^{t_f^\pm} \right) (q, p) \left(1 - \chi_0\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right) \right), \\ a_{LZ}^\pm(q, p) &= \left(a^\pm \circ \Phi_\pm^{t_f^\pm} \right) (q, p) \chi_0\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right). \end{aligned}$$

Since the trajectories, which pass within the time interval $]t - t_i^\pm, t - t_f^\pm[$ through the support of a_{BO}^\pm , do not jump, Proposition 2.3 with $C = 1$ is enough to deal with the Born–Oppenheimer part. The analysis of the Landau–Zener part, however, involves nonadiabatic transitions. For points $(q, p) \in \operatorname{supp}(a_{LZ}^\pm)$ we have $|\phi(q)| \leq 2R\sqrt{\varepsilon}$. Hence, not all of the trajectories passing through the support of a_{LZ}^\pm jump. Nevertheless, we argue as if all of them did. Indeed, the transition coefficients generated by these added jumps are exponentially small with respect to ε , since they occur for points (q, p) with $|\phi(q)| > R\sqrt{\varepsilon}$.

PROPOSITION 2.4. *Let $0 < t_f^\pm < t_i^\pm$ be such that at $t - t_i^\pm$ and $t - t_f^\pm$ all random trajectories arriving at time t in Ω^\pm are contained in $\{\frac{R}{2}\sqrt{\varepsilon} \leq |\phi| \leq R\sqrt{\varepsilon}\}$ and have*

performed their only jump within the interval $]t - t_i^\pm, t - t_f^\pm[$ of length $\delta_t^\pm = t_i^\pm - t_f^\pm$. Then, for all $\chi \in \mathcal{C}_c^\infty([0, T], \mathbb{R})$

$$\begin{aligned} & \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) W^\varepsilon(\psi^\varepsilon(t))(q, p) a_{LZ}^\pm(q, p) \Pi^\pm(q) \, dq \, dp \, dt \\ &= \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) W^\varepsilon(\psi^\varepsilon(t - \delta_t^\pm))(q, p) (\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} a_{LZ}^\pm)(q, p) \Pi^\pm(q) \, dq \, dp \, dt \\ (2.4) \quad & + O(1/R^2) + O(R^3\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O(1/(R^5\sqrt{\varepsilon})). \end{aligned}$$

One observes that in the right-hand side of (2.4) only the plus-projector Π^+ appears. This comes from the fact that at time $t - t_i^\pm$ the contribution on the minus-mode is of order $O(1/(R^5\sqrt{\varepsilon})) + O(1/R^2) + O(\sqrt{\varepsilon})$, which is due to Proposition 2.3 and the initial data being associated only with $\operatorname{Ran} \Pi^+$. Finally, again the classical transport result of Proposition 2.3 relates the right-hand side of (2.4) with the initial data, and the proof of our main result, Theorem 2.2, is complete.

3. Numerical experiments. Before giving the detailed proof, we present numerical experiments illustrating the theoretical convergence result and the effectiveness of the algorithm. We consider a two-level Schrödinger equation with codimension two crossing in two space dimensions, which models the photoisomerization of retinal in rhodopsin. This conformational change is considered as the first step of vision. In [14], computations with the model Hamiltonian

$$-\frac{\omega}{2}\partial_x^2 - \frac{1}{2m}\partial_\varphi^2 + \frac{1}{2}\omega x^2 + \begin{pmatrix} \frac{1}{2}W_0(1 - \cos \varphi) & \lambda x \\ \lambda x & E_1 - \frac{1}{2}W_1(1 - \cos \varphi) + \kappa x \end{pmatrix}$$

have qualitatively reproduced spectroscopic information of the molecule. The two effective coordinates $(\varphi, x) \in]-\frac{\pi}{2}, \frac{3\pi}{2}] \times \mathbb{R}$ consist of the reaction coordinate ϕ and a collective coordinate x . The parameters are $m^{-1} = 4.84 \cdot 10^{-4}$, $E_1 = 2.48$, $W_0 = 3.6$, $W_1 = 1.09$, $\omega = \lambda = 0.19$, and $\kappa = 0.1$ (all in eV, $\hbar = 1$); see note 18 in [14]. Setting

$$\varepsilon = m^{-1/2} = 0.022, \quad q_1 = \varphi, \quad q_2 = \frac{\varepsilon}{\sqrt{\omega}} x,$$

one obtains a rescaled Hamiltonian $-\frac{\varepsilon^2}{2}\Delta_q + V(q)$ with potential

$$V(q) = \frac{1}{2}(\beta q_2)^2 + \begin{pmatrix} \frac{1}{2}W_0(1 - \cos q_1) & \alpha_1 q_2 \\ \alpha_1 q_2 & E_1 - \frac{1}{2}W_1(1 - \cos q_2) + \alpha_2 q_2 \end{pmatrix},$$

whose parameters $(\alpha_1, \alpha_2) = \frac{\sqrt{\omega}}{\varepsilon}(\kappa, \lambda) \approx (2, 3.8)$ and $\beta = \omega/\varepsilon \approx 8.6$ are of order one with respect to ε . Fixing these values of (α_1, α_2) and β , we run a series of experiments for varying values of the semiclassical parameter ε and hopping ranges R ,

$$\varepsilon \in \{0.0005, 0.001, 0.005, 0.01, 0.015, 0.02, 0.022, 0.03\}, \quad R \in \{1, 2, 3\},$$

in the following set-up. We consider normalized Gaussian initial data associated with the plus-level, that is,

$$\psi_0^\varepsilon(q) = (\varepsilon\pi)^{-1/2} \exp\left(-\frac{1}{2\varepsilon}|q - q_0^\varepsilon|^2 + \frac{i}{\varepsilon} p_0 \cdot (q - q_0^\varepsilon)\right) v^+(q),$$

TABLE 1

The table shows final level populations and particle numbers as well as the accuracy of the reference solver. The population of the upper level $\|\Pi^+\psi^\varepsilon(t)\|^2$ at the final time $T = 7\sqrt{\varepsilon}$ is computed by the reference solver, a pseudospectral Strang splitting scheme, and illustrates leading order nonadiabatic transitions for all test cases. Depending on the hopping range R , the final particle numbers of the surface hopping algorithm vary between 3000 and 30000. The reference accuracy is the difference in L^2 -norm of the final wave function computed with full and half resolution. For all computations it is less than 10^{-4} .

ε	0.0005	0.001	0.005	0.01	0.015	0.02	0.022	0.03
$\ \Pi^+\psi^\varepsilon(t)\ ^2$	0.499	0.576	0.792	0.378	0.263	0.276	0.276	0.239
# particles, $R = 1$	5396	5229	3730	3548	3316	3274	3482	3292
# particles, $R = 2$	6650	6353	4732	4766	4801	5569	6106	8281
# particles, $R = 3$	7392	6881	5809	6155	7119	10 536	12 581	34 485
Ref. accuracy $\cdot 10^5$	3.44	2.89	7.57	2.56	2.47	2.45	2.47	2.63

where $v^+(q)$ is a normalized eigenvector of $V(q)$ for the eigenvalue $v(q) + |\phi(q)|$, which depends smoothly on q . The initial center in position space

$$q_0^\varepsilon = (1.63 - 4\sqrt{\varepsilon}, 0.5\sqrt{\varepsilon})$$

is chosen left of the two crossing points $(\gamma_l, 0)$, $(\gamma_r, 0)$, where $\gamma_{l,r}$ are the two solutions of $\cos \varphi = 1 - 2E_1/(W_0 + W_1)$ for $\varphi \in]-\frac{\pi}{2}, \frac{3\pi}{2}]$, that is, $\gamma_l \approx 1.63$ and $\gamma_r \approx 4.65$. The initial momentum center and the time interval,

$$p_0 = (1, 0), \quad [0, T] = [0, 7\sqrt{\varepsilon}],$$

are chosen such that the wave function passes only the left crossing point $(\gamma_l, 0)$ once without returning to it again. The upper level populations $\|\Pi^+\psi^\varepsilon(t)\|^2$ at the final time $T = 7\sqrt{\varepsilon}$, which are given in Table 1, confirm that the experimental set-up produces leading order nonadiabatic transitions.

The numerical realization of the surface hopping semigroup $(\mathcal{L}_{\varepsilon,R}^t)_{t \geq 0}$ is the same as for the simulations of models with a linear isotropic potential matrix presented in [17], up to adding the $R\sqrt{\varepsilon}$ -dependent jump criterion

$$t \mapsto |\phi(q^\pm(t))| \text{ has a local minimum in } t = t^* \text{ and } |\phi(q^\pm(t^*))| \leq R\sqrt{\varepsilon}.$$

The initial sampling is performed on 16×16 grids in position and momentum space, and the classical transport is discretized by the explicit Runge–Kutta method of Dormand and Prince DOPRI45. For comparison, we have also solved the Schrödinger equation (1.1) by a numerically converged pseudospectral Strang splitting scheme. The solutions obtained with a 1024×512 space grid on the computational domain $[1.63 - 8\sqrt{\varepsilon}, 1.63 + 16\sqrt{\varepsilon}] \times [-6\sqrt{\varepsilon}, 6\sqrt{\varepsilon}]$ and with 10^4 time steps are regarded as a reference, since they differ in L^2 -norm from the corresponding solution with a fourth of the grid points and half the time steps by less than 10^{-4} ; see Table 1. We have computed the following quadratic quantities of the wave function at final time $T = 7\sqrt{\varepsilon}$: the level populations $\|\Pi^\pm\psi^\varepsilon(t)\|^2$ and the expectation values of position and momentum on each level,

$$\langle \Pi^\pm(q)\psi^\varepsilon(q, t), q\Pi^\pm(q)\psi^\varepsilon(q, t) \rangle, \quad \langle \Pi^\pm(q)\psi^\varepsilon(q, t), -i\varepsilon\nabla_q\Pi^\pm(q)\psi^\varepsilon(q, t) \rangle.$$

We note that the reference solver restricts the length of the time interval to $7\sqrt{\varepsilon}$, since the dynamics can be resolved only as long as the solution stays well localized in the computational domain, while for the fixed number of 1024×512 grid points the size of the computational domain affects the numerical accuracy.

Comparing the outcome of the two algorithms, we find all errors within and below the corridor $[\frac{1}{5}\sqrt{\varepsilon}, 5\sqrt{\varepsilon}]$ (see Figure 1), which is better than the proven convergence rate $\varepsilon^{1/8}$. Moreover, the errors increase monotonically when increasing the semiclassical parameter; however, when entering the range $\varepsilon \geq 0.01$ we observe different dependencies. The level populations' error starts decreasing, while the error of the momentum expectation oscillates. We have no mathematical explanation for these observations and can make only an educated guess. The good convergence rate might be caused by the localization properties of the initial Gaussian wave packet combined with the short length of the time interval. Indeed, the error terms $O(1/R^2)$ and $O(1/(R^5\sqrt{\varepsilon}))$ in the approximation by classical transport in Proposition 2.3 might be negligible in this case as well as the contribution $O(1/R^2)$ in Proposition 2.4, which is due to localization in energy. The tendencies for a large semiclassical parameter, however, are clearly beyond the reach of our asymptotic analysis.

Table 1 shows that an increase of the hopping range R increases the number of final particles, that is, the number of jumps within the overall time interval. We obtain particle numbers around 3000 and 30000 for $R = 1$ and $R = 3$, respectively, resulting in half a minute and 5 minutes computing time for our implementation of the algorithm in MATLAB 7.0 on a 3GHz Pentium 4 computer. However, an enlarged hopping range need not improve the accuracy of the approximation. The plots in Figure 1 mostly display smaller errors for larger R , but the level populations in the physical relevant range of $\varepsilon = 0.022$ have the most accurate computation for $R = 2$.

Summarizing, the numerical experiments are consistent with the theoretical result, but also present a better convergence rate and tendencies in the range of a larger semiclassical parameter, which seem to be unexplainable by our asymptotic analysis. A systematic comparison with the well-established surface hopping algorithms of chemical physics is in progress.

4. Propagation outside the crossing zone. We now begin proving our main result. The first step is to establish the validity of the classical transport approximation in the zone of large eigenvalue gap $\{|\phi| \geq C R\sqrt{\varepsilon}\}$.

Proof of Proposition 2.3. Our aim is to prove

$$\begin{aligned} \text{tr} \int \left\{ \chi(t)c(q,p)b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right), \tau + \frac{1}{2}|p|^2 + v(q) \pm |\phi(q)| \right\} \Pi^\pm(q)W^\varepsilon(\psi^\varepsilon(t))(q,p) \, dq \, dp \, dt \\ = O\left(\frac{1}{R^2}\right) + O\left(\frac{1}{R^5\sqrt{\varepsilon}}\right) + O(\sqrt{\varepsilon}), \end{aligned}$$

since then classical transport follows immediately. The key argument is the estimation of the action of the commutator

$$K = \frac{1}{\varepsilon} \left[\chi(t)\text{op}_\varepsilon \left(c(q,p)b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right) \Pi^\pm(q) \right), -i\varepsilon\partial_t - \frac{\varepsilon^2}{2}\Delta_q + V(q) \right]$$

on the solution of the Schrödinger equation (1.1). Indeed, observing that

$$(K\psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})} = 0,$$

we are going to prove

$$\begin{aligned} (K\psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})} &= O\left(\frac{1}{R^2}\right) + O\left(\frac{1}{R^5\sqrt{\varepsilon}}\right) + O(\sqrt{\varepsilon}) \\ &+ \left(\text{op}_\varepsilon \left(\left\{ \chi(t)c(q,p)b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right), \tau + \frac{1}{2}|p|^2 + v(q) \pm |\phi(q)| \right\} \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})}, \end{aligned}$$

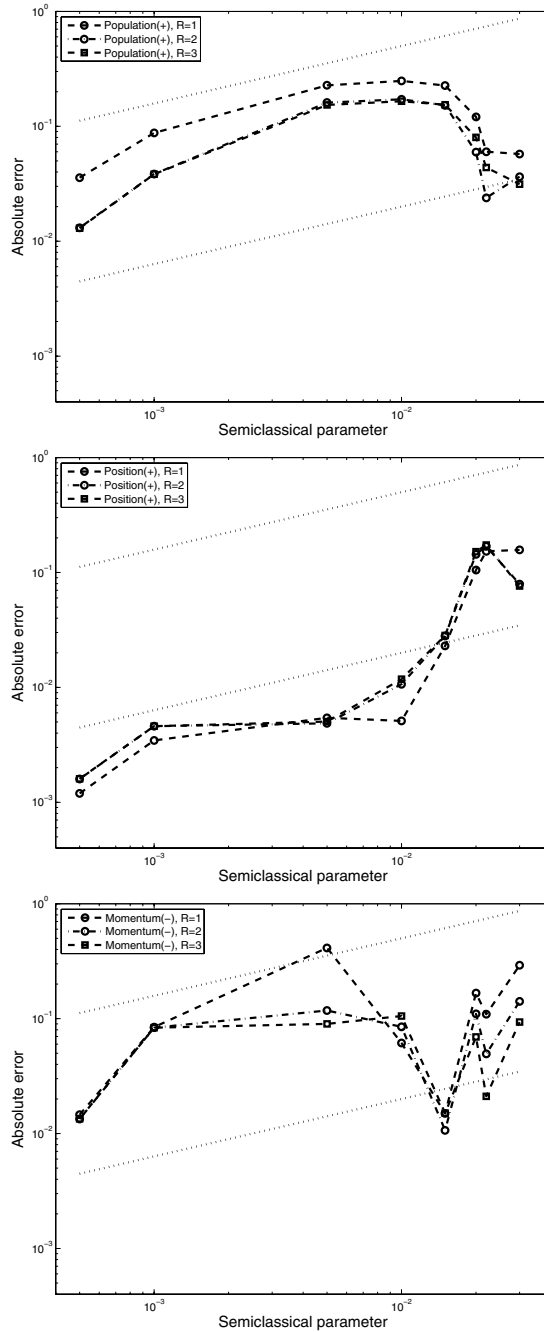


FIG. 1. Double logarithmic plots of the absolute error, when comparing the outcome of the surface hopping algorithm with a numerically converged pseudospectral splitting scheme. The semiclassical parameter ε varies in the set $\{0.0005, 0.001, 0.005, 0.01, 0.015, 0.02, 0.022, 0.03\}$. The dashed, dashed-dotted, and solid lines refer to a hopping range $R = 1, 2, 3$, respectively. The three plots show the error of the level population $\|\Pi^+(q)\psi^\varepsilon(q, t)\|^2$, of the position expectation value $\langle \Pi^+(q)\psi^\varepsilon(q, t), q\Pi^+(q)\psi^\varepsilon(q, t) \rangle$, and of the momentum expectation $\langle \Pi^-(q)\psi^\varepsilon(q, t), -i\varepsilon\nabla_q\Pi^-(q)\psi^\varepsilon(q, t) \rangle$ at the final time $T = 7\sqrt{\varepsilon}$. All errors lie in and below the corridor defined by the two functions $\varepsilon \mapsto 5\sqrt{\varepsilon}$ and $\varepsilon \mapsto \frac{1}{5}\sqrt{\varepsilon}$, which are represented by dotted lines.

where op_ε also denotes Weyl quantized operators acting on space-time variables. We use the scaling operator

$$(4.1) \quad T : L^2_{\text{loc}}(\mathbb{R}^{d+1}) \rightarrow L^2_{\text{loc}}(\mathbb{R}^{d+1}), \quad (T\psi)(t, q) = \varepsilon^{d/4} \psi(t, \sqrt{\varepsilon}q)$$

and write

$$T^*KT = \frac{1}{\varepsilon} \left[\text{op}_1 \left(\chi(t)c(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) b \left(\frac{\phi(\sqrt{\varepsilon}q)}{R\sqrt{\varepsilon}} \right) \Pi^\pm(\sqrt{\varepsilon}q) \right), -i\varepsilon\partial_t - \frac{\varepsilon}{2}\Delta_q + V(\sqrt{\varepsilon}q) \right].$$

We have to deal carefully with the ε, R dependence of the symbol in the left-hand side of the commutator. For all multi-indices $\alpha \in \mathbb{N}^d$,

$$D^\alpha(\Pi^\pm(q)) = O(|\phi(q)|^{-|\alpha|}).$$

Since $|\phi(q)| > C R\sqrt{\varepsilon}$ on the support of $b_{\varepsilon, R}(q, p) := c(q, p)b\left(\frac{\phi(q)}{R\sqrt{\varepsilon}}\right)$, we have

$$(4.2) \quad D_q^\alpha((b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(1), \quad D_p^\alpha((b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(\varepsilon^{|\alpha|/2})$$

for $\alpha \in \mathbb{N}^d$. By the symbolic calculus of Lemma A.2,

$$\begin{aligned} \frac{1}{\varepsilon} \left[\text{op}_1((b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)), \text{op}_1\left(\frac{\varepsilon}{2}|p|^2\right) \right] &= \frac{1}{i} \text{op}_1(\{(b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \frac{1}{2}|p|^2\}) \\ &= \frac{1}{i} \text{op}_1(\{b_{\varepsilon, R}(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \frac{1}{2}|p|^2\}\Pi^\pm(\sqrt{\varepsilon}q)) - \frac{1}{i} \text{op}_1(r_0(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)), \end{aligned}$$

where

$$(4.3) \quad r_0(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) = b_{\varepsilon, R}(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) \sum_{j=1}^d \sqrt{\varepsilon} p_j (\partial_{q_j} \Pi^\pm)(\sqrt{\varepsilon}q).$$

Moreover, in view of $[\Pi^\pm, V] = 0$ and (4.2),

$$\begin{aligned} \frac{1}{\varepsilon} \left[\text{op}_1((b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)), V(\sqrt{\varepsilon}q) \right] &= \frac{1}{2i\varepsilon} \text{op}_1(\{(b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), V(\sqrt{\varepsilon}q)\}) \\ &\quad - \frac{1}{2i\varepsilon} \text{op}_1(\{V(\sqrt{\varepsilon}q), (b_{\varepsilon, R}\Pi^\pm)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)\}) + O(\varepsilon). \end{aligned}$$

Working on the Poisson brackets involving $V = v + V_\ell(\phi)$, we first observe that

$$\{b_{\varepsilon, R}\Pi^\pm, V_\ell(\phi)\} - \{V_\ell(\phi), b_{\varepsilon, R}\Pi^\pm\} = \Pi^\pm \{b_{\varepsilon, R}, V_\ell(\phi)\} - \{V_\ell(\phi), b_{\varepsilon, R}\} \Pi^\pm.$$

Using that $V_\ell(\phi) = |\phi|(\Pi^+ - \Pi^-)$ and $\partial_{q_j}\Pi^\pm = \Pi^\pm(\partial_{q_j}\Pi^\pm) + (\partial_{q_j}\Pi^\pm)\Pi^\pm$, we get

$$\{b_{\varepsilon, R}\Pi^\pm, V_\ell(\phi)\} - \{V_\ell(\phi), b_{\varepsilon, R}\Pi^\pm\} = \pm 2\{b_{\varepsilon, R}, |\phi|\}\Pi^\pm \pm 2|\phi| \sum_{j=1}^d (\partial_{p_j} b_{\varepsilon, R}) \partial_{q_j} \Pi^\pm$$

and set

$$(4.4) \quad \begin{aligned} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) &= \frac{1}{\varepsilon} |\phi(\sqrt{\varepsilon}q)| \sum_{j=1}^d \partial_{p_j} (b_{\varepsilon, R}(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) \partial_{q_j} (\Pi^\pm(\sqrt{\varepsilon}q)) \\ &= |\phi(\sqrt{\varepsilon}q)| \sum_{j=1}^d (\partial_{p_j} c)(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) b \left(\frac{\phi(\sqrt{\varepsilon}q)}{R\sqrt{\varepsilon}} \right) (\partial_{q_j} \Pi^\pm)(\sqrt{\varepsilon}q). \end{aligned}$$

Now, collecting all the different pieces, we have

$$\begin{aligned} K &= \text{op}_\varepsilon(\{\chi(t)b_{\varepsilon,R}(q,p), \tau + \frac{1}{2}|p|^2 + v(q) \pm |\phi(q)|\} \Pi^\pm(q)) \\ &\quad + \text{op}_\varepsilon(\chi(t)r_0(q,p)) + \text{op}_\varepsilon(\chi(t)r_1(q,p)) + O(\varepsilon). \end{aligned}$$

Hence, our claim follows from the analysis of r_0 and r_1 , which is carried out in Lemma 4.1. \square

LEMMA 4.1. *Let ψ^ε solve the Schrödinger equation (1.1). For the matrix-valued functions r_0 and r_1 defined in (4.3) and (4.4), respectively, one has*

$$\begin{aligned} (\text{op}_\varepsilon(\chi(t)r_0(q,p)) \psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})} &= O(\sqrt{\varepsilon}) + O(1/R^2) + O(1/(R^5 \sqrt{\varepsilon})), \\ (\text{op}_\varepsilon(\chi(t)r_1(q,p)) \psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})} &= O(\sqrt{\varepsilon}) + O(1/R^2). \end{aligned}$$

Proof. Both functions have off-diagonal matrix structure; that is, $r_0(q,p)$ and $r_1(q,p)$ do not commute with $V(q)$. However, since r_1 contains an additional factor $|\phi|$, it is less singular than r_0 , in the sense that for $\alpha \in \mathbb{N}^d$ and (q,p) with $|\phi(\sqrt{\varepsilon}q)| > CR\sqrt{\varepsilon}$

$$D_q^\alpha(r_0(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(R^{-|\alpha|-1}\varepsilon^{-1/2}), \quad D_q^\alpha(r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(R^{-|\alpha|}).$$

We begin with r_1 . We write $r_1 = \Pi^- r_1 \Pi^+ + \Pi^+ r_1 \Pi^-$ and work successively with each part. Thus, without loss of generality, we suppose that $\Pi^- r_1 \Pi^+ = r_1$. The strategy is to reuse the Schrödinger equation, since

$$r_1 = \Pi^- r_1 \Pi^+ = \frac{1}{2|\phi|} [r_1, V_\ell(\phi)] = \frac{1}{2|\phi|} [r_1, \tau + \frac{1}{2}|p|^2 + V].$$

With the scaling operator T defined in (4.1), we obtain

$$\text{op}_\varepsilon(\chi(t)r_1(q,p)) = T^* \text{op}_1\left(\left[\frac{\chi(t)}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q)\right]\right) T.$$

Then, by the symbolic calculus of Lemma A.2,

$$\begin{aligned} \text{op}_1\left(\left[\frac{\chi(t)}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q)\right]\right) \\ = \left[\text{op}_1\left(\frac{\chi(t)}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)\right), \text{op}_1\left(\varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q)\right)\right] + \text{op}_1(r_2(t, \sqrt{\varepsilon}q, \sqrt{\varepsilon}p)), \end{aligned}$$

where

$$r_2(t, \sqrt{\varepsilon}q, \sqrt{\varepsilon}p) = \chi(t)\tilde{r}_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) + \varepsilon \frac{\chi'(t)}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) + \varepsilon \chi(t)\tilde{r}_2(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$$

with

$$\begin{aligned} \tilde{r}_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) &= \frac{1}{2i} \left\{ \frac{1}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q) \right\} \\ &\quad - \frac{1}{2i} \left\{ \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q), \frac{1}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) \right\}. \end{aligned}$$

The term $\tilde{r}_2(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$ contains the commutator of second derivatives in p of the symbol $|\phi(\sqrt{\varepsilon}q)|^{-1} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$ with second derivatives of $V(\sqrt{\varepsilon}q)$ and a linear combination of derivatives in (q,p) of order greater than or equal to three. Hence, by Lemma A.2

$$\text{op}_\varepsilon(\chi(t)\tilde{r}_2(q,p)) = O(\varepsilon\sqrt{\varepsilon}) + O(\varepsilon/R^4) \text{ in } \mathcal{L}(L^2(\mathbb{R}^{d+1})).$$

It remains to study \tilde{r}_1 . Since derivatives in p of $r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$ generate powers of $\sqrt{\varepsilon}$, the bracket with $V(\sqrt{\varepsilon}q)$ gives

$$\text{op}_1\left(\chi(t) \left\{ \frac{1}{2|\phi(q\sqrt{\varepsilon})|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), V(\sqrt{\varepsilon}q) \right\}\right) = O\left(\frac{\sqrt{\varepsilon}}{R}\right) \text{ in } \mathcal{L}(L^2(\mathbb{R}^{d+1})).$$

For the bracket with $\frac{\varepsilon}{2}|p|^2$ one has

$$\begin{aligned} & \left\{ \frac{1}{2|\phi(\sqrt{\varepsilon}q)|} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \frac{\varepsilon}{2}|p|^2 \right\} \\ &= \frac{\sqrt{\varepsilon}}{2|\phi(\sqrt{\varepsilon}q)|} \sqrt{\varepsilon}p \cdot \nabla_q (r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) + \frac{\varepsilon}{2} \frac{d\phi(\sqrt{\varepsilon}q)\sqrt{\varepsilon}p \cdot \phi(\sqrt{\varepsilon}q)}{|\phi(\sqrt{\varepsilon}q)|^3} r_1(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p). \end{aligned}$$

Both terms give a contribution of order $O(1/R^2)$, but since they are not purely off-diagonal, the preceding commutator argument cannot be reiterated. Hence,

$$\text{op}_\varepsilon(\chi(t)\tilde{r}_1(q, p)) = O\left(\frac{1}{R^2}\right) + O\left(\frac{\sqrt{\varepsilon}}{R}\right), \quad \text{op}_\varepsilon\left(\varepsilon \frac{\chi'(t)}{|\phi(q)|} r_1(q, p)\right) = O\left(\frac{\sqrt{\varepsilon}}{R}\right)$$

in $\mathcal{L}(L^2(\mathbb{R}^{d+1}))$, and we have proven one part of the lemma.

In the case of the more singular symbol r_0 , the previous strategy results in an error of size

$$\begin{aligned} & 1/R\sqrt{\varepsilon} (O(\varepsilon\sqrt{\varepsilon}) + O(\varepsilon/R^4) + O(\sqrt{\varepsilon}/R) + O(1/R^2)) \\ &= O(\sqrt{\varepsilon}) + O(1/R^2) + O(1/(R^3\sqrt{\varepsilon})). \end{aligned}$$

However, the special form of r_0 allows us to ameliorate the term of order $O(1/(R^3\sqrt{\varepsilon}))$, which stems from the Poisson bracket with $\frac{\varepsilon}{2}|p|^2$. Indeed, we observe that

$$p \cdot \nabla \Pi^+(q) = \frac{1}{4|\phi(q)|^3} [V_\ell(\phi(q)), [V_\ell(\phi(q)), V_\ell(d\phi(q)p)]] .$$

Therefore, the bracket with $\frac{\varepsilon}{2}|p|^2$ in the \tilde{r}_1 term writes as

$$\begin{aligned} & \left\{ \frac{1}{4|\phi(\sqrt{\varepsilon}q)|^3} b_{\varepsilon, R}(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) [V_\ell(\phi(\sqrt{\varepsilon}q)), V_\ell(d\phi(\sqrt{\varepsilon}q)\sqrt{\varepsilon}p)], \frac{\varepsilon}{2}|p|^2 \right\} \\ &= \frac{\varepsilon}{|\phi(\sqrt{\varepsilon}q)|^4} [V_\ell(\phi(\sqrt{\varepsilon}q)), G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)] \end{aligned}$$

for some matrix-valued function G_ε with

$$D_q^\alpha(G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(1), \quad D_p^\alpha(G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)) = O(\varepsilon^{|\alpha|/2})$$

for all $\alpha \in \mathbb{N}^d$. We then set

$$\begin{aligned} \tilde{r}_3(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) &:= \frac{\varepsilon}{|\phi(\sqrt{\varepsilon}q)|^4} [V_\ell(\phi(\sqrt{\varepsilon}q)), G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)] \\ &= - \left[\frac{\varepsilon}{|\phi(\sqrt{\varepsilon}q)|^4} G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q) \right] \end{aligned}$$

and obtain

$$\begin{aligned} & (\text{op}_\varepsilon(\chi(t)\tilde{r}_3(q, p)) \psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})} \\ &= \left(T^* \text{op}_1 \left(\left[\frac{\varepsilon\chi(t)}{|\phi(\sqrt{\varepsilon}q)|^4} G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q) \right] \right) T \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2i} \left(T^* \text{op}_1 \left(\left\{ \frac{\varepsilon \chi(t)}{|\phi(\sqrt{\varepsilon}q)|^4} G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p), \varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q) \right\} \right) T\psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\
&- \frac{1}{2i} \left(T^* \text{op}_1 \left(\left\{ \varepsilon\tau + \frac{\varepsilon}{2}|p|^2 + V(\sqrt{\varepsilon}q), \frac{\varepsilon \chi(t)}{|\phi(\sqrt{\varepsilon}q)|^4} G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p) \right\} \right) T\psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\
&+ (T^* \text{op}_1(\tilde{r}_4(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p))T\psi^\varepsilon, \psi^\varepsilon)_{L^2(\mathbb{R}^{d+1})},
\end{aligned}$$

where $\tilde{r}_4(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$ contains the commutator of second order derivatives in p of $\varepsilon \chi(t)|\phi(\sqrt{\varepsilon}q)|^{-4} G_\varepsilon(\sqrt{\varepsilon}q, \sqrt{\varepsilon}p)$ with second derivatives of $V(\sqrt{\varepsilon}q)$ and a linear combination of higher order derivatives in (q, p) . Hence, $\text{op}_\varepsilon(\tilde{r}_4(q\sqrt{\varepsilon}, p\sqrt{\varepsilon})) = O(\varepsilon/R^4)$ in $\mathcal{L}(L^2(\mathbb{R}^{d+1}))$. Since the Poisson brackets give a contribution of order $O(1/(R^5\sqrt{\varepsilon})) + O(1/R^4)$, the other part of the lemma is proven, too. \square

5. Transitions near the crossing. The microlocal normal form used for proving Proposition 2.4 holds locally near some point $(q_0, t_0, p_0, \tau_0) \in \mathbb{R}^{2d+2}$ of the phase space of space-time, which is a crossing point in the sense that $\phi(q_0) = 0$ and $\tau_0 + v(q_0) + \frac{1}{2}|p_0|^2 = 0$.

5.1. Localization in energy. For localization in energy, we consider a cut-off function $\theta \in C_c^\infty(\mathbb{R})$, $0 \leq \theta \leq 1$, with $\theta(x) = 1$ for $|x| \leq \frac{1}{2}$ and $\theta(x) = 0$ for $|x| > 1$. We set

$$\lambda^\pm(q, p, \tau) = \tau + v(q) + \frac{1}{2}|p|^2 \pm |\phi(q)|$$

and crucially use the following lemma for suitably reformulating Proposition 2.4.

LEMMA 5.1. *Let $c \in C_c^\infty(\mathbb{R}^{2d+1+\ell}, \mathbb{C})$. If $c_{\varepsilon, R}(t, q, p) = c(t, q, p, \phi(q)/(R\sqrt{\varepsilon}))$ is supported in $\{|\phi| \geq \frac{R}{2}\sqrt{\varepsilon}\}$, then*

$$\begin{aligned}
\text{tr} \int_{\mathbb{R}^{2d+1}} W^\varepsilon(\psi^\varepsilon(t))(q, p) c_{\varepsilon, R}(t, q, p) \Pi^\pm(q) dq dp dt &= O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}) \\
&+ \left(\text{op}_\varepsilon \left(c_{\varepsilon, R}(t, q, p) \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})}.
\end{aligned}$$

Proof. Writing $1 - \theta(u) = uG(u)$ with $G \in C^\infty(\mathbb{R})$, we have

$$1 - \theta\left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}}\right) = \frac{1}{R\sqrt{\varepsilon}} \lambda^\pm(q, p, \tau) G\left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}}\right).$$

We now argue as in section 4, using the estimates on Π^\pm and $\lambda^\pm(q, p, \tau)\Pi^\pm(q) = \Pi^\pm(q)(\tau + \frac{1}{2}|p|^2 + V(q))$. The symbolic calculus of Lemma A.2 yields in $\mathcal{L}(L^2(\mathbb{R}^{d+1}))$

$$\begin{aligned}
\text{op}_\varepsilon \left(c_{\varepsilon, R}(t, q, p) \left(1 - \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \right) \Pi^\pm(q) \right) &= O(\sqrt{\varepsilon}) + O\left(\frac{1}{R^2}\right) \\
&+ \frac{1}{R\sqrt{\varepsilon}} \text{op}_\varepsilon \left(c_{\varepsilon, R}(t, q, p) G \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \Pi^\pm(q) \right) \text{op}_\varepsilon \left(\tau + \frac{1}{2}|p|^2 + V(q) \right).
\end{aligned}$$

Indeed, the derivatives of the projectors are less harmful than in section 4, since they are divided only by $\sqrt{\varepsilon}$ and not by ε . Since ψ^ε solves the equation, we obtain

$$\left(\text{op}_\varepsilon \left(c_{\varepsilon, R}(t, q, p) \left(1 - \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \right) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} = O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}). \quad \square$$

By Lemma 5.1, we introduce the energy cut-off on both sides of equality (2.4), adding an error of order $O(1/R^2) + O(\sqrt{\varepsilon})$. Then, the left-hand side reads

$$\begin{aligned} \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t) W^\varepsilon(\psi^\varepsilon(t))(q, p) a_{LZ}^\pm(q, p) \Pi^\pm(q) dq dp dt &= O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}) \\ &+ \left(\operatorname{op}_\varepsilon \left(a_{LZ}^\pm(q, p) \chi(t) \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})}. \end{aligned}$$

Using the notation

$$f_a^\pm(q, p, j) = \mathbf{1}_{\{j=\pm 1\}}(j) a^j(q, p), \quad (q, p, j) \in \mathbb{R}_{\pm}^{2d},$$

for a V -diagonal matrix-valued symbol a , the action of the semigroup on a_{LZ}^\pm can be written as

$$(\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} a_{LZ}^\pm)(q, p) \Pi^\pm(q) = (\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} f_{a_{LZ}}^\pm)(q, p, +) \Pi^\pm(q).$$

Then, the right-hand side of (2.4) is

$$\begin{aligned} \operatorname{tr} \int_{\mathbb{R}^{2d+1}} \chi(t + \delta_t^\pm) W^\varepsilon(\psi^\varepsilon(t))(q, p) (\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} a_{LZ}^\pm)(q, p) \Pi^\pm(q) dq dp dt \\ = \left(\operatorname{op}_\varepsilon \left(\chi(t + \delta_t) \theta \left(\frac{\lambda^+(q, p, \tau)}{R\sqrt{\varepsilon}} \right) (\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} f_{a_{LZ}}^\pm)(q, p, +) \Pi^+(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\ + O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}), \end{aligned}$$

and the proof of Proposition 2.4 reduces to showing that

$$\begin{aligned} &\left(\operatorname{op}_\varepsilon \left(a_{LZ}^\pm(q, p) \chi(t) \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\ &= \left(\operatorname{op}_\varepsilon \left(\chi(t + \delta_t^\pm) \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right) (\mathcal{L}_{\varepsilon, R}^{\delta_t^\pm} f_{a_{LZ}}^\pm)(q, p, +) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\ (5.1) \quad &+ O\left(\frac{1}{R^2}\right) + O(R^3\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O\left(\frac{1}{R^5\sqrt{\varepsilon}}\right). \end{aligned}$$

For points $(q, t, p, \tau, j) \in \mathbb{R}^{2d+2} \times \{\pm 1\}$, we consider random trajectories

$$\mathcal{T}_{\varepsilon, R}^{(q, t, p, \tau, j)} : [0, +\infty) \rightarrow (\mathbb{R}^{2d+2} \times \{\pm 1\})$$

with $\mathcal{T}_{\varepsilon, R}^{(q, t, p, \tau, j)}(r) = (q^j(r), r + t, p^j(r), \tau, j)$ as long as $(q^j(r), p^j(r)) \notin S_{\varepsilon, R}$ and a jump from j to $-j$ with probability $T^\varepsilon(q^*, p^*)$, whenever $(q^j(r), p^j(r))$ hits $S_{\varepsilon, R}$ at a point (q^*, p^*) . We keep the notation $(\mathcal{L}_{\varepsilon, R}^r)_{r \geq 0}$ for the associated semigroup and set

$$(5.2) \quad c_{\varepsilon, R}^\pm(q, t, p, \tau) = a_{LZ}^\pm(q, p) \chi(t) \theta \left(\frac{\lambda^\pm(q, p, \tau)}{R\sqrt{\varepsilon}} \right).$$

Since $r \mapsto \lambda^\pm(q^\pm(r), p^\pm(r), \tau)$ is a constant function, and since within $[t - t_i, t - t_f]$ all involved random trajectories perform a jump, we have

$$(\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_{\varepsilon, R}}^+)(q, t, p, \tau, +) = \chi(t + \delta_t) \theta \left(\frac{\lambda^+(q, p, \tau)}{R\sqrt{\varepsilon}} \right) (\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{a_{LZ}}^+)(q, p, +),$$

$$(\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_{\varepsilon, R}}^-)(q, t, p, \tau, +) = \chi(t + \delta_t) \theta \left(\frac{\lambda^+(q, p, \tau)}{R\sqrt{\varepsilon}} \right) (\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{a_{LZ}}^-)(q, p, +).$$

With this notation, (5.1) and, consequently, Proposition 2.4 are equivalent to

$$\begin{aligned}
& \left(\text{op}_\varepsilon \left(c_{\varepsilon,R}^\pm(q, t, p, \tau) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\
&= \left(\text{op}_\varepsilon \left((\mathcal{L}_{\varepsilon,R}^{\delta_t^\pm} f_{c_{\varepsilon,R}^\pm}^\pm)(q, t, p, \tau, +) \Pi^\pm(q) \right) \psi^\varepsilon, \psi^\varepsilon \right)_{L^2(\mathbb{R}^{d+1})} \\
(5.3) \quad & + O(1/R^2) + O(R^3\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O(1/(R^5\sqrt{\varepsilon})).
\end{aligned}$$

We emphasize that the symbols $c_{\varepsilon,R}^\pm$ and $\mathcal{L}_{\varepsilon,R}^{\delta_t^\pm} f_{c_{\varepsilon,R}^\pm}^\pm$ are compactly supported inside the annulus $\{\frac{R}{2}\sqrt{\varepsilon} < |\phi| < R\sqrt{\varepsilon}\}$ at a distance of order $R\sqrt{\varepsilon}$ of $J^{\pm, \text{out}}$, where

$$J^{\pm, \text{in/out}} = \left\{ (q, t, p, \tau) \in \mathbb{R}^{2d+2} \mid (q, p) \in M^{\pm, \text{in/out}}, \tau + v(q) + \frac{1}{2}|p|^2 \pm |\phi(q)| = 0 \right\}$$

denote the submanifolds, which consist of all Hamiltonian trajectories entering (respectively, leaving) the crossing set.

5.2. The normal form. Let us first recall some basic facts about canonical transforms and Fourier integral operators. The phase space $T^*\mathbb{R}^{d+1} = \mathbb{R}^{d+1} \times \mathbb{R}^{d+1}$ is a symplectic space, once endowed with the symplectic form $\omega = d\tau \wedge dt + dp \wedge dq$. A canonical transform $\kappa : T^*\mathbb{R}^{d+1} \rightarrow T^*\mathbb{R}^{d+1}$ is a change of coordinates, which preserves the symplectic form. With a canonical transform κ , one associates a unitary operator $U : L^2(\mathbb{R}^{d+1}) \rightarrow L^2(\mathbb{R}^{d+1})$ such that for all $a \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+2}, \mathbb{C}^{N(\ell) \times N(\ell)})$

$$U^* \text{op}_\varepsilon(a) U = \text{op}_\varepsilon(a \circ \kappa) + O(\varepsilon^2)$$

as bounded operators on $L^2(\mathbb{R}^{d+1})$; see, for example, section 2.2 in [7]. The operator U is a Fourier integral operator. The last equality extends to symbols of the form

$$b_{\varepsilon,R}(q, t, p, \tau) = b\left(q, t, p, \tau, \frac{f(q,p)}{R\sqrt{\varepsilon}}\right)$$

with $b \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+3}, \mathbb{C}^{N(\ell) \times N(\ell)})$ and $f \in \mathcal{C}^\infty(\mathbb{R}^{2d}, \mathbb{R})$ according to

$$(5.4) \quad U^* \text{op}_\varepsilon(b_{\varepsilon,R}) U = \text{op}_\varepsilon(b_{\varepsilon,R} \circ \kappa) + O(\sqrt{\varepsilon}).$$

The proof of this statement follows the proof of Lemma 2 in [7]: one uses symbolic calculus for the commutator of a usual semiclassical pseudodifferential operator and a two-scale one of the form $\text{op}_\varepsilon(b_{\varepsilon,R})$, hence the gain of $\sqrt{\varepsilon}$.

We shall crucially use the following microlocal normal form result, which for codimension two and three crossings is proven in [1, 2], however, without the explicit equations (5.6)–(5.9). These equations, including the normal form for codimension five crossings, are provided in Theorem 1 and Proposition 4 of [6].

THEOREM 5.2 (see [6]). *We consider $\rho_0 = (q_0, t_0, p_0, \tau_0 = -v(q_0) - \frac{1}{2}|p_0|^2)$ such that $\phi(q_0) = 0$, $d\phi(q_0)p_0 \neq 0$, and $d\phi$ is of maximal rank near q_0 . Then, there exists a local canonical transform κ from a neighborhood of ρ_0 into some neighborhood Ω of 0,*

$$\kappa : (q, t, p, \tau) \mapsto (z, s, \zeta, \sigma), \quad \kappa(\rho_0) = 0.$$

There exist a Fourier integral operator U associated with κ^{-1} and an invertible matrix-valued symbol $A_\varepsilon = A_0 + \varepsilon A_1 + \varepsilon^2 A_2 + \dots$ such that $v^\varepsilon = U^ \text{op}_\varepsilon(A_\varepsilon)^{-1} \psi^\varepsilon$ satisfies for all $\varphi \in \mathcal{C}_c^\infty(\Omega, \mathbb{R})$*

$$(5.5) \quad \text{op}_\varepsilon(\varphi) \text{op}_\varepsilon(-\sigma + V_\ell(s, \tilde{z} + \gamma_\varepsilon(z, \zeta))) v^\varepsilon = O(\varepsilon^\infty)$$

in $L^2(\mathbb{R}^{d+1})$, where $z = (\tilde{z}, z') \in \mathbb{R}^d$ with $\tilde{z} \in \mathbb{R}^{\ell-1}$ and $\gamma_\varepsilon = \gamma_\varepsilon(z, \zeta)$ is a vector-valued symbol $\gamma_\varepsilon = \gamma_0 + \varepsilon\gamma_1 + \varepsilon^2\gamma_2 + \dots \in \mathbb{R}^{\ell-1}$ with

$$\gamma_\varepsilon = 0 \text{ for } \ell = 2, \quad \gamma_\varepsilon = O(|\tilde{z}|^2) \text{ for } \ell > 2.$$

$\tilde{z} \in \mathbb{R}^{\ell-1}$ contains the coordinates of the vector $|\mathrm{d}\phi(q)p|^{-1/2}\pi_\ell(q,p)\phi(q)$ in an orthonormal basis of the hyperplane normal to $\mathrm{d}\phi(q)p$ up to $O(|\phi(q)|^2)$, while

$$s = -|\mathrm{d}\phi(q)p|^{-1/2} \frac{\mathrm{d}\phi(q)p}{|\mathrm{d}\phi(q)p|} \cdot \phi(q) + O(s^2 + \sigma^2 + |\tilde{z}|^2), \quad (5.6)$$

$$\sigma = |\mathrm{d}\phi(q)p|^{-1/2} \left(\tau + \frac{1}{2}|p|^2 + v(q) \right) + O(s^2 + \sigma^2 + |\tilde{z}|^2).$$

Moreover,

$$(5.7) \quad J^{\pm, in} = \{\sigma \mp s = 0, \tilde{z} = 0, s \leq 0\}, \quad J^{\pm, out} = \{\sigma \pm s = 0, \tilde{z} = 0, s \geq 0\},$$

and there exists $\gamma \in \{-1, +1\}$ such that for all $\rho = (q, t, p, \tau)$ with $\kappa(\rho) = (z, s, \zeta, \sigma)$

$$(5.8) \quad A_0^*(\rho) \left(\tau + \frac{1}{2}|p|^2 + V(q) \right) A_0(\rho) = \gamma \left(-\sigma + V_\ell(s, \tilde{z} + \gamma_0(z, \zeta)) \right),$$

$$(5.9) \quad A_0^*(\rho) V_\ell \left(\frac{\pi_\ell(q,p)\phi(q)}{|\pi_\ell(q,p)\phi(q)|} \right) A_0(\rho) = \gamma V_\ell \left(0, \frac{\tilde{z}}{|\tilde{z}|} \right) + O(\sqrt{\sigma^2 + s^2 + |\tilde{z}|^2}).$$

We denote by

$$\begin{aligned} \tilde{\Pi}^\pm(z, s, \zeta) &= \frac{1}{2} \left(\mathrm{Id} \mp \frac{1}{\sqrt{s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2}} V_\ell(s, \tilde{z} + \gamma_0(z, \zeta)) \right), \\ \tilde{\lambda}^\pm(z, s, \zeta, \sigma) &= -\sigma \mp \sqrt{s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2} \end{aligned}$$

the spectral projectors and the eigenvalues of $-\sigma + V_\ell(s, \tilde{z} + \gamma_0(z, \zeta))$. Due to the relation $J^{\pm, in/out} \subseteq \{-\sigma \mp \sqrt{s^2 + |\tilde{z}|^2} = 0, \tilde{z} = 0\}$, the labeling of $\tilde{\Pi}^\pm$ coincides on $J^{\pm, in/out}$ with the one for Π^\pm .

PROPOSITION 5.3. *There exist functions k^\pm such that if $\kappa(q, t, p, \tau) = (z, s, \zeta, \sigma)$, the projectors Π^\pm and $\tilde{\Pi}^\pm$ are related by*

$$(5.10) \quad \tilde{\Pi}^\pm(z, s, \zeta, \sigma) = (k^\pm A_0^* \Pi^\pm A_0)(q, t, p, \tau) \text{ on } \Sigma^\mp = \{\lambda^\mp = 0\}.$$

If $S = \{\phi(q) = 0, \tau + \frac{1}{2}|p|^2 + v(q) = 0\}$, then $k_{|S}^+ = k_{|S}^- = e \neq 0$ and $(e A_0^* A_0)|_S = \mathrm{Id}|_S$. Moreover, on $\Sigma^\pm \cap \{0 < |\phi| \leq R\sqrt{\varepsilon}\}$,

$$(5.11) \quad \tilde{\Pi}^\pm(z, s, \zeta, \sigma) = e(A_0^* \Pi^\pm A_0)(q, t, p, \tau) + O(R\sqrt{\varepsilon}).$$

Proof. For convenience, we set

$$P = \tau + \frac{1}{2}|p|^2 + V(q), \quad \tilde{P} = -\sigma + V_\ell(s, \tilde{z} + \gamma_0(z, \zeta)).$$

By (5.8), the use of determinants gives $\Sigma^+ \cup \Sigma^- = \kappa^{-1}(\{\tilde{\lambda}^+ = 0\} \cup \{\tilde{\lambda}^- = 0\})$. Considering the equations of $J^{\pm, in/out}$, the only possibility is

$$\Sigma^\pm = \kappa^{-1}(\{\tilde{\lambda}^\pm = 0\}).$$

Therefore, for $\rho \in \Sigma^+$ we have

$$\gamma \tilde{P}(\kappa(\rho)) = \gamma(\tilde{\lambda}^- \tilde{\Pi}^-)(\kappa(\rho)) = (A_0^* P A_0)(\rho) = (\lambda^- A_0^* \Pi^- A_0)(\rho).$$

The same argument for Σ^- gives (5.10).

The fact $k_{|S}^+ = k_{|S}^-$ comes from the precise analysis of the Hamiltonian vector fields associated with the eigenvalues λ^\pm . Let $\rho \in S$. We find in [6, section 5] that if

$$H(\rho) = \lim_{\alpha \rightarrow 0^+} H_{\lambda^\pm}(\Phi_\pm^\alpha(q, p)), \quad H'(\rho) = \lim_{\alpha \rightarrow 0^\pm} H_{\lambda^\pm}(\Phi_\pm^\alpha(q, p)),$$

then there exists a nonzero function e such that

$$H = e(\partial_s + \partial_\sigma), \quad H' = e(\partial_s - \partial_\sigma) \quad \text{on } S.$$

Since $\lambda^\pm = 0$ and $\tilde{\lambda}^\pm = 0$ on S , we have $k_{|S}^+ = k_{|S}^- = e$. Next, we consider the limit of the projectors Π^\pm along outgoing trajectories,

$$\Pi_S^\mp(\rho) = \lim_{\alpha \rightarrow 0^-} \Pi^\mp(\Phi_\pm^\alpha(q, p)) = \frac{1}{2} \left(\text{Id} \mp V_\ell \left(\frac{d\phi(q)p}{|d\phi(q)p|} \right) \right).$$

Then, (5.10) gives on S

$$e A_0^* A_0 = e A_0^* (\Pi_S^+ + \Pi_S^-) A_0 = (\tilde{\Pi}^+ + \tilde{\Pi}^-) \circ \kappa = \text{Id}.$$

Finally, let $\rho \in \Sigma^+$. By relation (5.8), we have for any vector $w \in \mathbb{C}^{N(\ell)}$ that $w \in \text{Ker } \tilde{P}(\kappa(\rho))$ if and only if $A_0(\rho)w \in \text{Ker } P(\rho)$. Moreover, $\text{Ker } P(\rho) = \text{Ran } \Pi^+(\rho)$ and $\text{Ker } \tilde{P}(\kappa(\rho)) = \text{Ran } \tilde{\Pi}^+(\kappa(\rho))$. We therefore obtain

$$\text{Ran } \tilde{\Pi}^+(\kappa(\rho)) = \text{Ran } (A_0^{-1} \Pi^+ A_0)(\rho).$$

Since $\sqrt{e} A_0$ is unitary on S , we have

$$(A_0^{-1} \Pi^+ A_0)^*(\rho) = (A_0^{-1} \Pi^+ A_0)(\rho) + O(R\sqrt{\varepsilon})$$

for $\rho \in \Sigma^+ \cap \{|\phi| \leq R\sqrt{\varepsilon}\}$. Since the two projectors $\tilde{\Pi}^+(\kappa(\rho))$ and $(A_0^{-1} \Pi^+ A_0)(\rho)$ have the same range, while one of them is orthogonal and the other orthogonal up to $O(R\sqrt{\varepsilon})$, they coincide up to $O(R\sqrt{\varepsilon})$. The same argument holds for $\rho \in \Sigma^-$, and we have proven relation (5.11). \square

As the next step towards proving the claimed identity (5.3), we perform the canonical change of coordinates for arriving at the microlocal normal form.

PROPOSITION 5.4. *Let $c_{\varepsilon, R}^\pm \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+2}, \mathbb{C})$ be the functions defined in (5.2) and $v^\varepsilon = U^* \text{op}_\varepsilon(A_\varepsilon)^{-1} \psi^\varepsilon$. Denote*

$$(5.12) \quad \tilde{T}^\varepsilon(\tilde{z}) = \exp\left(-\frac{\pi}{\varepsilon} |\tilde{z}|^2\right).$$

Then, there exist functions $b^\pm \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+2}, \mathbb{C})$ and $s_1^\pm \in \mathbb{R}$, such that $b^\pm(z, s, \zeta, \eta)$ and $b^\pm(z, s_1^\pm + s, \zeta, \eta)$ are compactly supported in $\{s > 0\}$ and $\{s < 0\}$, respectively, and satisfy

$$\left(\text{op}_\varepsilon \left(c_{\varepsilon, R}^+(q, t, p, \tau) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2}$$

$$\begin{aligned}
&= \left(\text{op}_\varepsilon \left(b^+ \left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) v_2^\varepsilon(z, s), v_2^\varepsilon(z, s) \right)_{L^2} + O(R\sqrt{\varepsilon}), \\
&\left(\text{op}_\varepsilon \left((\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_{\varepsilon, R}}^+) (q, t, p, \tau, +) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\
&= \left(\text{op}_\varepsilon \left((1 - \tilde{T}^\varepsilon(\tilde{z})) b^+ \left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) v_1^\varepsilon(z, s), v_1^\varepsilon(z, s) \right)_{L^2} + O(R^3\sqrt{\varepsilon}), \\
&\left(\text{op}_\varepsilon \left(c_{\varepsilon, R}^- (q, t, p, \tau) \Pi^-(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\
&= \left(\text{op}_\varepsilon \left(b^- \left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) v_1^\varepsilon(z, s), v_1^\varepsilon(z, s) \right)_{L^2} + O(R\sqrt{\varepsilon}), \\
&\left(\text{op}_\varepsilon \left((\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_{\varepsilon, R}}^-) (q, t, p, \tau, +) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\
&= \left(\text{op}_\varepsilon \left(\tilde{T}^\varepsilon(\tilde{z}) b^- \left(z, s + s_1^-, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) v_1^\varepsilon(z, s), v_1^\varepsilon(z, s) \right)_{L^2} + O(R^3\sqrt{\varepsilon}).
\end{aligned}$$

Proof. We prove only the first two equalities, since one deals similarly with the other two. By symbolic calculus and the transformation property (5.4), the canonical transform κ^{-1} of Theorem 5.2 acts as

$$\begin{aligned}
&\left(\text{op}_\varepsilon \left(c_{\varepsilon, R}^+ (q, t, p, \tau) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\
&= \left(\text{op}_\varepsilon \left(((c_{\varepsilon, R}^+ A_0^* \Pi^+ A_0) \circ \kappa^{-1})(z, s, \zeta, \sigma) \right) v^\varepsilon(z, s), v^\varepsilon(z, s) \right)_{L^2} + O(\sqrt{\varepsilon}).
\end{aligned}$$

The compactly supported function $c_{\varepsilon, R}^+ \circ \kappa^{-1}$ is localized near $J^{+, out}$, that is, near $\{\sigma + s = 0, \tilde{z} = 0, s > 0\}$. The relation (5.11) between the projectors gives a function $b \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+2+\ell}, \mathbb{C})$ compactly supported in $\{s > 0\}$ such that

$$\begin{aligned}
&((c_{\varepsilon, R}^+ A_0^* \Pi^+ A_0) \circ \kappa^{-1})(z, s, \zeta, \sigma) \\
&= b \left(z, s, \zeta, \sigma, \frac{\tilde{z}}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z, s, \zeta, \sigma)}{R\sqrt{\varepsilon}} \right) \tilde{\Pi}^+(z, s, \zeta) + O(R\sqrt{\varepsilon}) \\
&=: b_{\varepsilon, R}(z, s, \zeta, \sigma) \tilde{\Pi}^+(z, s, \zeta) + O(R\sqrt{\varepsilon})
\end{aligned}$$

as functions in $\mathcal{C}_c^\infty(\mathbb{R}^{2d+2}, \mathbb{C})$. Hence, we obtain

$$\begin{aligned}
&\left(\text{op}_\varepsilon \left(c_{\varepsilon, R}^+ (q, t, p, \tau) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\
&= \left(\text{op}_\varepsilon \left(b_{\varepsilon, R}(z, s, \zeta, \sigma) \tilde{\Pi}^+(z, s, \zeta) \right) v^\varepsilon(z, s), v^\varepsilon(z, s) \right)_{L^2} + O(R\sqrt{\varepsilon}).
\end{aligned}$$

For $|\tilde{z}| = O(R\sqrt{\varepsilon})$ we have

$$\begin{aligned}
&\tilde{\Pi}^+(z, s, \zeta) = \begin{pmatrix} 0 & 0 \\ 0 & \text{Id} \end{pmatrix} + O(R\sqrt{\varepsilon}) \text{ in } \{s > 0\}, \\
(5.13) \quad &\tilde{\Pi}^+(z, s, \zeta) = \begin{pmatrix} \text{Id} & 0 \\ 0 & 0 \end{pmatrix} + O(R\sqrt{\varepsilon}) \text{ in } \{s < 0\},
\end{aligned}$$

and therefore

$$\begin{aligned} & \left(\text{op}_\varepsilon \left(b_{\varepsilon,R}(z, s, \zeta, \sigma) \tilde{\Pi}^+(z, s, \zeta) \right) v^\varepsilon(z, s), v^\varepsilon(z, s) \right)_{L^2} \\ &= \left(\text{op}_\varepsilon \left(b_{\varepsilon,R}(z, s, \zeta, \sigma) \right) v_2^\varepsilon(z, s), v_2^\varepsilon(z, s) \right)_{L^2} + O(R\sqrt{\varepsilon}). \end{aligned}$$

We now remove the σ -dependence of the symbol. Taylor expanding around $\sigma = -\sqrt{s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2}$, we write

$$\begin{aligned} b_{\varepsilon,R}(z, s, \zeta, \sigma) &= b \left(z, s, \zeta, -\sqrt{s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2}, \frac{\tilde{z}}{R\sqrt{\varepsilon}}, 0 \right) \\ &\quad + \frac{1}{R\sqrt{\varepsilon}} \tilde{\lambda}^+(s, z, \sigma, \zeta) G \left(z, s, \zeta, \sigma, \frac{\tilde{z}}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z, s, \sigma, \zeta)}{R\sqrt{\varepsilon}} \right) \end{aligned}$$

with $G \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+2+\ell}, \mathbb{C})$. Since

$$\tilde{\lambda}^+(z, s, \zeta, \sigma) \tilde{\Pi}^+(z, s, \zeta) = \tilde{\Pi}^+(z, s, \zeta) (\sigma - V_\ell(s, \tilde{z} + \gamma_0(z, \zeta))),$$

and since v^ε solves the Landau-Zener-type problem (5.5), an argument analogous to the proof of Lemma 5.1 yields

$$\begin{aligned} & \left(\text{op}_\varepsilon \left(b_{\varepsilon,R}(z, s, \zeta, \sigma) \right) v_2^\varepsilon(z, s), v_2^\varepsilon(z, s) \right)_{L^2} = O\left(\frac{1}{R^2}\right) + O(\sqrt{\varepsilon}) \\ & \quad + \left(\text{op}_\varepsilon \left(b \left(z, s, \zeta, -\sqrt{s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2}, \frac{\tilde{z}}{R\sqrt{\varepsilon}}, 0 \right) \right) v_2^\varepsilon(z, s), v_2^\varepsilon(z, s) \right)_{L^2}. \end{aligned}$$

Setting $b^+(z, s, \zeta, \eta) = b(z, s, \zeta, -s, \eta, 0)$, we obtain

$$\begin{aligned} & \left(\text{op}_\varepsilon \left(c_{\varepsilon,R}^+(q, t, p, \tau) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\ &= \left(\text{op}_\varepsilon \left(b^+ \left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) v_2^\varepsilon(z, s), v_2^\varepsilon(z, s) \right)_{L^2} + O(R\sqrt{\varepsilon}) + O\left(\frac{1}{R^2}\right). \end{aligned}$$

Next, we focus on the second claimed identity, which contains nonadiabatic transitions. We have

$$\begin{aligned} & \left(\text{op}_\varepsilon \left((\mathcal{L}_{\varepsilon,R}^{\delta_t} f_{c_{\varepsilon,R}}^+)(q, t, p, \tau, +) \Pi^+(q) \right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t) \right)_{L^2} \\ &= \left(\text{op}_\varepsilon \left(((\mathcal{L}_{\varepsilon,R}^{\delta_t} f_{c_{\varepsilon,R}}^+) A_0^* \Pi^+ A_0)(\kappa^{-1}(z, s, \zeta, \sigma), +) \right) v^\varepsilon(z, s), v^\varepsilon(z, s) \right)_{L^2} + O(\sqrt{\varepsilon}). \end{aligned}$$

Let $r \mapsto (z^+(r), s^+(r), \zeta^+(r), \sigma^+(r))$ be the Hamiltonian trajectory of

$$\begin{aligned} \dot{z} &= \partial_\zeta \tilde{\lambda}^+ = {}^t d_\zeta \gamma_0(z, \zeta) (\tilde{z} + \gamma_0(z, \zeta)) (s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2)^{-1/2}, \\ \dot{s} &= \partial_\sigma \tilde{\lambda}^+ = 1, \\ \dot{\zeta} &= -\partial_z \tilde{\lambda}^+ = -{}^t d_z (\tilde{z} + \gamma_0(z, \zeta)) (\tilde{z} + \gamma_0(z, \zeta)) (s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2)^{-1/2}, \\ \dot{\sigma} &= -\partial_s \tilde{\lambda}^+ = -s (s^2 + |\tilde{z} + \gamma_0(z, \zeta)|^2)^{-1/2} \end{aligned}$$

with $(z(0), s(0), \zeta(0), \sigma(0)) = (z, s, \zeta, \sigma) = \kappa(q, t, p, \tau)$. A trajectory jumps for $r = r_*$ if $s^+(r_*) = O(|\phi(q)|^2) = O(R^2\varepsilon)$. Then, $\sigma^+(r_*) = O(R^2\varepsilon)$ as well. Since

$$\frac{d}{dr} (s^+(r) + \sigma^+(r)) = O(|\tilde{z}|^2) \quad \text{on} \quad J^{+,out} = \{\sigma + s = 0, \tilde{z} = 0, s > 0\},$$

and since $|\tilde{z}| = O(R\sqrt{\varepsilon})$ on the support of our symbol, we have

$$\begin{aligned} z^+(r) &= z + O(R^3\varepsilon^{3/2}), & s^+(r) &= s + r, \\ \zeta^+(r) &= \zeta + O(R\sqrt{\varepsilon}), & \sigma^+(r) &= -s^+(r) + O(R^2\varepsilon). \end{aligned}$$

These asymptotics together with conservation of energy along Hamiltonian flows yield

$$\begin{aligned} &b\left(z^+(r), s^+(r), \zeta^+(r), \sigma^+(r), \frac{\tilde{z}^+(r)}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z^+(r), s^+(r), \zeta^+(r), \sigma^+(r))}{R\sqrt{\varepsilon}}\right) \\ &= b\left(z, s + r, \zeta, -(s + r), \frac{\tilde{z}}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z, s, \zeta, \sigma)}{R\sqrt{\varepsilon}}\right) + O(R\sqrt{\varepsilon}). \end{aligned}$$

Jumps occur for $|\phi(q)| \leq R\sqrt{\varepsilon}$, that is, for

$$\begin{aligned} |\tilde{z}|^2 &= |\mathrm{d}\phi(q)p|^{-1}|\pi_\ell(q, p)\phi(q)|^2 + O(|\phi(q)|^3) \\ &= |\mathrm{d}\phi(q)p|^{-1}|\pi_\ell(q, p)\phi(q)|^2 + O(R^3\varepsilon^{3/2}). \end{aligned}$$

Therefore, the transition rate $T^\varepsilon(q, p)$ reads in the new coordinates as

$$\begin{aligned} T^\varepsilon(q, p) &= \exp\left(-\frac{\pi}{\varepsilon}|\mathrm{d}\phi(q)p|^{-1}|\pi_\ell(q, p)\phi(q)|^2\right) \\ &= \exp\left(-\frac{\pi}{\varepsilon}|\tilde{z}|^2\right) + O(R^3\sqrt{\varepsilon}) = \tilde{T}^\varepsilon(\tilde{z}) + O(R^3\sqrt{\varepsilon}), \end{aligned}$$

and there exists $s_1^+ \in \mathbb{R}$ such that

$$\begin{aligned} &((\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_\varepsilon, R}^+) A_0^* \Pi^+ A_0)(\kappa^{-1}(z, s, \zeta, \sigma), +) \\ &= \left(1 - \tilde{T}^\varepsilon(\tilde{z})\right) b\left(z, s + s_1^+, \zeta, -(s + s_1^+), \frac{\tilde{z}}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z, s, \zeta, \sigma)}{R\sqrt{\varepsilon}}\right) \tilde{\Pi}^+(z, s, \zeta) + O(R^3\sqrt{\varepsilon}). \end{aligned}$$

By the asymptotics (5.13) of $\tilde{\Pi}^+$ above $\{s < 0\}$, we then get

$$\begin{aligned} &\left(\mathrm{op}_\varepsilon\left((\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_\varepsilon, R}^+) A_0^* \Pi^+ A_0)(\kappa^{-1}(z, s, \zeta, \sigma), +)\right) v^\varepsilon(z, s), v^\varepsilon(z, s)\right)_{L^2} = O(R^3\sqrt{\varepsilon}) \\ &+ \left(\mathrm{op}_\varepsilon\left(\left(1 - \tilde{T}^\varepsilon(\tilde{z})\right) b\left(z, s + s_1^+, \zeta, -(s + s_1^+), \frac{\tilde{z}}{R\sqrt{\varepsilon}}, \frac{\tilde{\lambda}^+(z, s, \zeta, \sigma)}{R\sqrt{\varepsilon}}\right)\right) v_1^\varepsilon(z, s), v_1^\varepsilon(z, s)\right)_{L^2}. \end{aligned}$$

As before, we remove the σ -dependence of the symbol and obtain

$$\begin{aligned} &\left(\mathrm{op}_\varepsilon\left((\mathcal{L}_{\varepsilon, R}^{\delta_t} f_{c_\varepsilon, R}^+)(q, t, p, \tau, +)\Pi^+(q)\right) \psi^\varepsilon(q, t), \psi^\varepsilon(q, t)\right)_{L^2} \\ &= \left(\mathrm{op}_\varepsilon\left(\left(1 - \tilde{T}^\varepsilon(\tilde{z})\right) b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) v_1^\varepsilon(z, s), v_1^\varepsilon(z, s)\right)_{L^2} + O(R^3\sqrt{\varepsilon}). \quad \square \end{aligned}$$

5.3. The transitions. It remains to analyze $v^\varepsilon = U^* \mathrm{op}_\varepsilon(A_\varepsilon^{-1})\psi^\varepsilon$ for proving the equality of each pair in Proposition 5.4 up to an error of $O(1/R^2) + O(R^3\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O(1/(R^5\sqrt{\varepsilon}))$, concluding the proof of our main result. We consider the solution u^ε of

$$(5.14) \quad -i\varepsilon\partial_s u^\varepsilon = \begin{pmatrix} s \mathrm{Id} & \sqrt{\varepsilon} G \\ \sqrt{\varepsilon} G^* & -s \mathrm{Id} \end{pmatrix} u^\varepsilon, \quad u^\varepsilon|_{s=0} = v^\varepsilon|_{s=0},$$

where G is one of the operators

$$G_2 = \frac{1}{\sqrt{\varepsilon}} \varphi \left(\frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \tilde{z}, \quad G_3 = \frac{1}{\sqrt{\varepsilon}} \varphi \left(\frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) (Z_1 + iZ_2),$$

$$G_5 = \frac{1}{\sqrt{\varepsilon}} \varphi \left(\frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \begin{pmatrix} Z_1 + iZ_2 & Z_3 + iZ_4 \\ -Z_3 + iZ_4 & Z_1 - iZ_2 \end{pmatrix}$$

with $Z = \text{op}_\varepsilon(\tilde{z} + \gamma_\varepsilon(z, \zeta))$ and $\varphi \in \mathcal{C}_c^\infty(\mathbb{R}^{\ell-1}, \mathbb{R})$. In all three cases, G is a bounded operator on $L^2(\mathbb{R}^d)$ with $\|G\| = O(R)$, and we have

$$\|u^\varepsilon - v^\varepsilon\|_{L^2_{loc}(\mathbb{R}^{d+1})} = O(\varepsilon^\infty).$$

The following Landau–Zener-type formula is given in Proposition 7 of [8], up to the explicit error terms.

PROPOSITION 5.5. *Let u^ε be the solution of (5.14). There exist vector-valued functions $\alpha^\varepsilon = (\alpha_1^\varepsilon, \alpha_2^\varepsilon)$, $\omega^\varepsilon = (\omega_1^\varepsilon, \omega_2^\varepsilon) \in L^2(\mathbb{R}^d, \mathbb{C}^{N(\ell)})$ such that for any function $\chi \in \mathcal{C}_c^\infty(\{x \in \mathbb{R} \mid |x| \leq R^2\}, \mathbb{R})$ the families $(\chi(GG^*)\alpha_1^\varepsilon)_{\varepsilon>0}$, $(\chi(G^*G)\alpha_2^\varepsilon)_{\varepsilon>0}$, $(\chi(GG^*)\omega_1^\varepsilon)_{\varepsilon>0}$, $(\chi(G^*G)\omega_2^\varepsilon)_{\varepsilon>0}$ are bounded in $L^2(\mathbb{R}^d, \mathbb{C})$ and satisfy*

$$\chi(GG^*)u_1^\varepsilon(z, s) = \chi(GG^*)e^{is^2/(2\varepsilon)} \left| \frac{s}{\sqrt{\varepsilon}} \right|^{i\frac{GG^*}{2}} k_1^\varepsilon(z) + O(R^2\sqrt{\varepsilon}),$$

$$\chi(G^*G)u_2^\varepsilon(z, s) = \chi(G^*G)e^{-is^2/(2\varepsilon)} \left| \frac{s}{\sqrt{\varepsilon}} \right|^{-i\frac{G^*G}{2}} k_2^\varepsilon(z) + O(R^2\sqrt{\varepsilon})$$

in $L^2(\mathbb{R}^d, \mathbb{C})$, where $k_j^\varepsilon = \alpha_j^\varepsilon$ and $k_j^\varepsilon = \omega_j^\varepsilon$ for $s < 0$ and $s > 0$, respectively, $j \in \{1, 2\}$. Moreover,

$$(5.15) \quad \begin{pmatrix} \omega_1^\varepsilon \\ \omega_2^\varepsilon \end{pmatrix} = \begin{pmatrix} a(GG^*) & -\bar{b}(GG^*)G \\ b(G^*G)G^* & a(G^*G) \end{pmatrix} \begin{pmatrix} \alpha_1^\varepsilon \\ \alpha_2^\varepsilon \end{pmatrix}$$

with

$$a(\lambda) = e^{-\pi\lambda/2}, \quad b(\lambda) = \frac{2ie^{i\pi/4}}{\lambda\sqrt{\pi}} 2^{-i\lambda/2} e^{-\pi\lambda/4} \Gamma(1 + i\frac{\lambda}{2}) \sinh(\frac{\pi\lambda}{2}).$$

Proof. Lemma 7 in [8] is the crucial step in the proof of Proposition 7 for which we have to check that the leading order error estimate is indeed $O(R^2\sqrt{\varepsilon})$. For this, we turn to the explicit calculations in the proof of Lemma 11 in [7] and study the two integrals

$$A_0 = s^{-1+i\eta^2/2} e^{-is^2/2} \int_{\mathbb{R}} \tilde{\chi} \left(\sqrt{2} - \sqrt{2 - \frac{2z}{s^2}} \right) \left| \sqrt{2} - \sqrt{2 - \frac{2z}{s^2}} \right|^{i\eta^2/2} \frac{e^{iz}}{\sqrt{2 - \frac{2z}{s^2}}} dz,$$

$$B_0 = s^{1+i\eta^2/2} \int_{\mathbb{R}} (1 - \tilde{\chi}(y)) e^{-\frac{i}{2}s^2(1+y^2-2\sqrt{2}y)} |y|^{i\eta^2/2} dy,$$

where $\tilde{\chi} \in \mathcal{C}_c^\infty(\mathbb{R}, \mathbb{R})$ is a function with $0 \leq \tilde{\chi} \leq 1$, $\tilde{\chi}(y) = 0$ for $|y| \geq \sqrt{2}/2$, and $\tilde{\chi}(y) = 1$ for $|y| \leq \sqrt{2}/4$. The phase function $y \mapsto -\frac{i}{2}(1+y^2-2\sqrt{2}y)$ of B_0 has the stationary point $y = \sqrt{2}$, and Taylor expansion of $y \mapsto \ln|y|$ around $y = \sqrt{2}$ yields

$$B_0 = \sqrt{2\pi} e^{-i\pi/4} 2^{i\eta^2/4} s^{i\eta^2/2} e^{is^2/2} + O(\eta^2 s^{-2}),$$

while integration by parts gives

$$A_0 = O(\eta^2 s^{-1})$$

as $\eta, s \rightarrow \infty$. Since the asymptotics of the other relevant integrals can be obtained analogously, the claimed error estimates follow by setting $s = O(\varepsilon^{-1/2})$ and $\eta = O(R)$. \square

For implementing these Landau–Zener asymptotics, we need the following additional relations, which are literally contained in the proofs of Lemmas 8 and 9 in [8].

LEMMA 5.6. *For any $\chi \in \mathcal{C}_c^\infty(\mathbb{R}, \mathbb{R})$ and $b \in \mathcal{C}_c^\infty(\mathbb{R}^{2d+N(\ell)-1}, \mathbb{C})$, we have*

$$\chi\left(\frac{|\xi|^2}{\varepsilon}\right) = \chi(GG^*) + O(\sqrt{\varepsilon}) = \chi(G^*G) + O(\sqrt{\varepsilon}),$$

$$\left|\frac{s}{\sqrt{\varepsilon}}\right|^{\pm i\frac{G^*G}{2}} \text{op}_\varepsilon\left(b\left(z, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \left|\frac{s}{\sqrt{\varepsilon}}\right|^{\mp i\frac{G^*G}{2}} = \text{op}_\varepsilon\left(b\left(z, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) + O(\sqrt{\varepsilon}|\ln \varepsilon|).$$

Now, we are ready to conclude the proof of Theorem 2.2. We discuss only the first pair of terms in Proposition 5.4, since the other pair can be dealt with analogously. We set

$$\begin{aligned} I_{\varepsilon,R}^1 &= \left(\text{op}_\varepsilon\left(b^+\left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) u_2^\varepsilon(z, s), u_2^\varepsilon(z, s)\right)_{L^2}, \\ I_{\varepsilon,R}^2 &= \left(\text{op}_\varepsilon\left((1 - \tilde{T}^\varepsilon(\tilde{z}))b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) u_1^\varepsilon(z, s), u_1^\varepsilon(z, s)\right)_{L^2}. \end{aligned}$$

By Lemma 5.6 and Proposition 5.5, we have

$$\begin{aligned} I_{\varepsilon,R}^1 &= \left(\text{op}_\varepsilon\left(b^+\left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \chi(G^*G)u_2^\varepsilon(z, s), \chi(G^*G)u_2^\varepsilon(z, s)\right)_{L^2} + O(\sqrt{\varepsilon}) \\ &= \left(\text{op}_\varepsilon\left(b^+\left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) e^{-i\frac{s^2}{2\varepsilon}} \left|\frac{s}{\sqrt{\varepsilon}}\right|^{-i\frac{G^*G}{2}} \omega_2^\varepsilon(z), e^{-i\frac{s^2}{2\varepsilon}} \left|\frac{s}{\sqrt{\varepsilon}}\right|^{-i\frac{G^*G}{2}} \omega_2^\varepsilon(z)\right)_{L^2} \\ &\quad + O(R^2\sqrt{\varepsilon}) \\ &= \left(\text{op}_\varepsilon\left(b^+\left(z, s, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \omega_2^\varepsilon(z), \omega_2^\varepsilon(z)\right)_{L^2} + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) \\ &= \left(\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \omega_2^\varepsilon(z), \omega_2^\varepsilon(z)\right)_{L^2} + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) \end{aligned}$$

because $s_1^\pm = O(R\sqrt{\varepsilon})$ and analogously

$$\begin{aligned} I_{\varepsilon,R}^2 &= \left(\text{op}_\varepsilon\left((1 - \tilde{T}^\varepsilon(\tilde{z}))b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \alpha_1^\varepsilon(z), \alpha_1^\varepsilon(z)\right)_{L^2} \\ &\quad + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|). \end{aligned}$$

By the scattering identity (5.15), we have $\omega_2^\varepsilon = b(G^*G)G^*\alpha_1^\varepsilon + a(G^*G)\alpha_2^\varepsilon$ and

$$\begin{aligned} I_{\varepsilon,R}^1 &= \left(\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) (b(G^*G)G^*\alpha_1^\varepsilon(z) + a(G^*G)\alpha_2^\varepsilon(z)), \right. \\ &\quad \left. b(G^*G)G^*\alpha_1^\varepsilon(z) + a(G^*G)\alpha_2^\varepsilon(z)\right)_{L^2} + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|). \end{aligned}$$

Since the wave function $\psi^\varepsilon(q, t)$ is of order $O(1/R^2) + O(\sqrt{\varepsilon}) + O(1/(R^5\sqrt{\varepsilon}))$ near the set $J^{-,in} = \{\sigma + s = 0, \tilde{z} = 0, s < 0\}$, we have

$$\begin{aligned} \begin{pmatrix} 0 \\ v_{\tilde{z}}^\varepsilon(z, s) \end{pmatrix} &= \text{op}_\varepsilon\left(\tilde{\Pi}^-(z, s, \zeta)\right) v^\varepsilon(z, s) + O(R\sqrt{\varepsilon}) \\ &= O(1/R^2) + O(R\sqrt{\varepsilon}) + O(1/(R^5\sqrt{\varepsilon})) \end{aligned}$$

as functions in $L_{loc}^2(\mathbb{R}^{d+1})$ localized near $J^{-,in}$. The preceding arguments expressing $I_{\varepsilon,R}^1$ and $I_{\varepsilon,R}^2$ in terms of α^ε and ω^ε then yield

$$a(G^*G)\alpha_2^\varepsilon(z) = O(1/R^2) + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O(1/(R^5\sqrt{\varepsilon}))$$

near $J^{-,in}$, and hence

$$\begin{aligned} I_{\varepsilon,R}^1 &= \left(\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) b(G^*G)G^*\alpha_1^\varepsilon(z), b(G^*G)G^*\alpha_1^\varepsilon(z) \right)_{L^2} \\ &\quad + O\left(\frac{1}{R^2}\right) + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O\left(\frac{1}{R^5\sqrt{\varepsilon}}\right). \end{aligned}$$

Lemma 5.6 together with the relations $G^*b(GG^*) = b(G^*G)G^*$ and

$$\lambda|b(\lambda)|^2 = 1 - e^{-\pi\lambda}$$

implies

$$\begin{aligned} &G\bar{b}(G^*G)\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) b(G^*G)G^*\alpha_1^\varepsilon(z) \\ &= G\bar{b}(G^*G)b(G^*G)G^*\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \alpha_1^\varepsilon(z) + O(\sqrt{\varepsilon}) \\ &= GG^*\bar{b}(GG^*)b(GG^*)\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \alpha_1^\varepsilon(z) + O(\sqrt{\varepsilon}) \\ (5.16) \quad &= (1 - \tilde{T}^\varepsilon(\tilde{z}))\text{op}_\varepsilon\left(b^+\left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}}\right)\right) \alpha_1^\varepsilon(z) + O(\sqrt{\varepsilon}) \end{aligned}$$

and finally

$$I_{\varepsilon,R}^1 = I_{\varepsilon,R}^2 + O(1/R^2) + O(R^2\sqrt{\varepsilon}) + O(\sqrt{\varepsilon}|\ln \varepsilon|) + O(1/(R^5\sqrt{\varepsilon})).$$

6. Eigenvalues of multiplicity two. We have assumed that the matrix-valued observable a is V -diagonal in the sense that $a = a^+\Pi^+ + a^-\Pi^-$ with scalar-valued functions a^\pm . The more natural assumption, that a commutes with V ,

$$a = \Pi^+ a \Pi^+ + \Pi^- a \Pi^-,$$

does not change the situation in the case $\ell = 2, 3$ but enlarges the class of observables for $\ell = 3', 5$. For observables of this form, one has to modify the Markov process to account for a polarization effect. The state space requires an additional component $w \in \mathbb{C}^4$, and when the deterministic flow $\Phi_j^t(q, p)$ hits the jump manifold $S_{\varepsilon,R}$ in a point (q^*, p^*) a more general branching occurs. The state (q^*, p^*, j, w) changes with probability $T^\varepsilon(q^*, p^*)$ to $(q^*, p^*, -j, w)$ and with probability $1 - T^\varepsilon(q^*, p^*)$ to $(q^*, p^*, j, \mathcal{R}(q^*, p^*)w)$, where

$$\mathcal{R}(q, p) = V_\ell \left(\frac{\pi_\ell(q, p)\phi(q)}{|\pi_\ell(q, p)\phi(q)|} \right).$$

This phenomenon is also described in Theorem 1 of [5] for two-scale Wigner measures. Our main result, Theorem 2.2, for the propagation of Wigner functions still applies for the semigroup, which incorporates polarization.

Proof of Theorem 2.2 for $a = \Pi^+ a \Pi^+ + \Pi^- a \Pi^-$ if $\ell = 3', 5$. Let us first prove classical transport. We set $A^+ = \Pi^+ a \Pi^+$ and focus on the $+$ mode. We extensively use $\Pi^+ A^+ = A^+ \Pi^+ = A^+$. The strategy is similar to the one of section 4, and we have to focus on the Poisson brackets

$$\frac{1}{2}\{A^+(q, p), \tau + \frac{1}{2}|p|^2 + v(q) + V_\ell(\phi(q))\} - \frac{1}{2}\{\tau + \frac{1}{2}|p|^2 + v(q) + V_\ell(\phi(q)), A^+(q, p)\}.$$

We set $\mu(q, p, \tau) = \tau + \frac{1}{2}|p|^2 + v(q)$ and write

$$\{A^+, \mu\} = \Pi^+ \{A^+, \mu\} \Pi^+ + A^+ \{\Pi^+, \mu\} + \{\Pi^+, \mu\} A^+.$$

We observe that

$$r_0 = A^+ \{\Pi^+, \mu\} + \{\Pi^+, \mu\} A^+ = -A^+ (\nabla_q \Pi^+ \cdot p) - (\nabla_q \Pi^+ \cdot p) A^+$$

can be treated as in section 4. Indeed, since A^+ commutes with $V_\ell(\phi)$, one has for any matrix G that $A^+[V_\ell(\phi), G] = [V_\ell(\phi), A^+G]$, $[V_\ell(\phi), G]A^+ = [V_\ell(\phi), GA^+]$, and consequently

$$\begin{aligned} r_0 &= -\frac{1}{4|\phi|^3} [V_\ell(\phi), [V_\ell(\phi), A^+V_\ell(d\phi p)]] - \frac{1}{4|\phi|^3} [V_\ell(\phi), [V_\ell(\phi), V_\ell(d\phi p)A^+]] \\ &= -\left[\mu + V_\ell(\phi), \frac{1}{4|\phi|^3} [V_\ell(\phi), A^+V_\ell(d\phi p)]\right] - \left[\mu + V_\ell(\phi), \frac{1}{4|\phi|^3} [V_\ell(\phi), V_\ell(d\phi p)A^+]\right]. \end{aligned}$$

The most harmful of the arising terms contains the brackets with $\frac{1}{2}|p|^2$, that is,

$$\tilde{r}_0 = -\left\{\frac{1}{2}|p|^2, \frac{1}{4|\phi|^3} [V_\ell(\phi), A^+V_\ell(d\phi p)]\right\} - \left\{\frac{1}{2}|p|^2, \frac{1}{4|\phi|^3} [V_\ell(\phi), V_\ell(d\phi p)A^+]\right\}.$$

Since the term containing the derivatives of $V_\ell(\phi)$ vanishes,

$$-\frac{1}{4|\phi|^3} [V_\ell(d\phi p), A^+V_\ell(d\phi p)] - \frac{1}{4|\phi|^3} [V_\ell(d\phi p), V_\ell(d\phi p)A^+] = 0,$$

there is a matrix-valued function G_ε with suitable bounds on its derivatives, such that $\tilde{r}_0 = |\phi|^{-4} [V_\ell(\phi), G_\varepsilon]$, and hence the other arguments of Lemma 4.1 apply for the analysis of r_0 .

For the brackets with the matrix part, we write

$$\frac{1}{2}\{A^+, V_\ell(\phi)\} - \frac{1}{2}\{V_\ell(\phi), A^+\} = \Pi^+ \{A^+, |\phi|\} \Pi^+ + |\phi| (\{A^+, \Pi^+\} - \{\Pi^+, A^+\}).$$

The second part,

$$r_1 = |\phi| (\{A^+, \Pi^+\} - \{\Pi^+, A^+\}) = |\phi| (\nabla_p A^+ \cdot \nabla_q \Pi^+ + \nabla_q \Pi^+ \cdot \nabla_p A^+),$$

is off-diagonal with respect to V , since $\nabla_p A^+ = \Pi^+ \nabla_p A^+ = \nabla_p A^+ \Pi^+$, $\Pi^\pm \Pi^\mp = 0$, and $\Pi^\pm \nabla_q \Pi^\pm = 0$ imply

$$\Pi^\pm (\nabla_p A^+ \cdot \nabla_q \Pi^+ + \nabla_q \Pi^+ \cdot \nabla_p A^+) \Pi^\pm = 0.$$

Hence, Lemma 4.1 applies.

The importance of $\mathcal{R}(q, p)$ for the nonadiabatic transitions becomes clear, when recasting (5.16) in the previous section as

$$\begin{aligned} & \begin{pmatrix} 0 & G^* \\ G & 0 \end{pmatrix} \bar{b}(G^*G) \operatorname{op}_\varepsilon \left(b^+ \left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) b(G^*G) \begin{pmatrix} 0 & G \\ G^* & 0 \end{pmatrix} \begin{pmatrix} \alpha_1^\varepsilon(z) \\ 0 \end{pmatrix} \\ &= V_\ell^* \left(0, \frac{\tilde{z}}{|\tilde{z}|} \right) (G^*G)^{\frac{1}{2}} \bar{b}(G^*G) \operatorname{op}_\varepsilon \left(b^+ \left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) \\ & \quad \times b(G^*G) (G^*G)^{\frac{1}{2}} V_\ell^* \left(0, \frac{\tilde{z}}{|\tilde{z}|} \right) \begin{pmatrix} \alpha_1^\varepsilon(z) \\ 0 \end{pmatrix} + O(\sqrt{\varepsilon}) \\ &= (1 - \tilde{T}^\varepsilon(\tilde{z})) V_\ell^* \left(0, \frac{\tilde{z}}{|\tilde{z}|} \right) \operatorname{op}_\varepsilon \left(b^+ \left(z, s + s_1^+, \zeta, \frac{\tilde{z}}{R\sqrt{\varepsilon}} \right) \right) V_\ell \left(0, \frac{\tilde{z}}{|\tilde{z}|} \right) \begin{pmatrix} \alpha_1^\varepsilon(z) \\ 0 \end{pmatrix} + O(\sqrt{\varepsilon}) \end{aligned}$$

and observing that the normal form transformation relates $V_\ell(0, \tilde{z}/|\tilde{z}|)$ and $\mathcal{R}(q, p)$ by identity (5.9) of Theorem 5.2. \square

Appendix A. Weyl calculus. For the convenience of the reader, we formulate the key technical lemma of the calculus of Weyl quantized pseudodifferential operators.

DEFINITION A.1. *A smooth matrix-valued function $a \in \mathcal{C}^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N \times N})$ is of subquadratic growth if for all $|\alpha| + |\beta| \geq 2$ there exists $C_{\alpha, \beta} > 0$ such that*

$$\|\partial_q^\alpha \partial_p^\beta a\|_\infty \leq C_\alpha.$$

LEMMA A.2. *Let $a \in \mathcal{C}_c^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N \times N})$, and let $b \in \mathcal{C}^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N \times N})$ be of subquadratic growth. Then, for all $\psi \in \mathcal{C}_c^\infty(\mathbb{R}^d, \mathbb{C}^N)$*

$$\operatorname{op}_1(a) \operatorname{op}_1(b) \psi = \left(\operatorname{op}_1(ab) + \frac{1}{2i} \operatorname{op}_1(\{a, b\}) + \operatorname{op}_1(c) + \operatorname{op}_1(r) \right) \psi,$$

with $\{a, b\} = \partial_p a \partial_q b - \partial_q a \partial_p b$ the Poisson bracket, c a linear combination of

$$\partial_{q_j \cdot p_j} a \partial_{q_j \cdot p_j} b, \partial_{q_j}^2 a \partial_{p_j}^2 b, \partial_{p_j}^2 a \partial_{q_j}^2 b,$$

and $r \in \mathcal{C}^\infty(\mathbb{R}^{2d}, \mathbb{C}^{N \times N})$ such that

$$N_k(r) := \sup_{|\alpha| + |\beta| \leq k} \|\partial_q^\alpha \partial_p^\beta r\|_\infty \leq C_k \sum_{m+m'=k} (N_m(D^3 a) N_{m'}(D^3 b))$$

for all $k \in \mathbb{N}$.

For a proof of this classical lemma, the reader can refer to [15] or to [3]. The theorem of Calderon and Vaillancourt implies that $\operatorname{op}_1(r)$ is a bounded operator on $L^2(\mathbb{R}^d, \mathbb{C}^N)$.

REFERENCES

- [1] Y. COLIN DE VERDIÈRE, *The level crossing problem in semi-classical analysis I. The symmetric case*, Ann. Inst. Fourier (Grenoble), 53 (2003), pp. 1023–1054.
- [2] Y. COLIN DE VERDIÈRE, *The level crossing problem in semi-classical analysis II. The Hermitian case*, Ann. Inst. Fourier (Grenoble), 54 (2004), pp. 1423–1441.
- [3] M. DIMASSI AND J. SJÖSTRAND, *Spectral Asymptotics in the Semi-classical Limit*, London Math. Soc. Lecture Note Ser. 268, Cambridge University Press, Cambridge, UK, 1999.
- [4] W. DOMCKE, D. YARKONY, AND H. KÖPPEL, EDS., *Conical Intersections*, World Scientific, Singapore, 2004.

- [5] C. FERMANIAN KAMMERER, *Wigner measures and molecular propagation through generic energy level crossings*, Rev. Math. Phys., 15 (2003), pp. 1285–1317.
- [6] C. FERMANIAN KAMMERER, *Normal forms for conical intersections in quantum chemistry*, Math. Phys. Electron. J., 13 (2007), paper 4.
- [7] C. FERMANIAN KAMMERER AND P. GÉRARD, *Mesures semi-classiques et croisements de modes*, Bull. Soc. Math. France, 130 (2002), pp. 123–168.
- [8] C. FERMANIAN KAMMERER AND P. GÉRARD, *A Landau-Zener formula for non-degenerated involutive codimension 3 crossings*, Ann. Henri Poincaré, 4 (2003), pp. 513–552.
- [9] C. FERMANIAN KAMMERER AND C. LASSER, *Wigner measures and codimension two crossings*, J. Math. Phys., 44 (2003), pp. 507–527.
- [10] P. GÉRARD, P. MARKOWICH, N. MAUSER, AND F. POUPAUD, *Homogenization limits and Wigner transforms*, Comm. Pure Appl. Math., 50 (1997), pp. 323–379.
- [11] P. GÉRARD, P. MARKOWICH, N. MAUSER, AND F. POUPAUD, *Erratum: Homogenization limits and Wigner transforms*, Comm. Pure Appl. Math., 53 (2000), pp. 280–281.
- [12] G. HAGEDORN, *A time dependent Born-Oppenheimer approximation*, Comm. Math. Phys., 77 (1980), pp. 1–19.
- [13] G. HAGEDORN, *Molecular Propagation through Electron Energy Level Crossings*, Mem. Amer. Math. Soc. 111, AMS, Providence, RI, 1994.
- [14] S. HAHN AND G. STOCK, *Quantum-mechanical modeling of the femtosecond isomerization in rhodopsin*, J. Phys. Chem. B, 104 (2000), pp. 1146–1149.
- [15] L. HÖRMANDER, *The Analysis of Linear Partial Differential Operators III. Pseudo-differential Operators*, Classics Math., Springer, Berlin, 1985.
- [16] L. LANDAU, *Collected Papers of L. Landau*, Pergamon Press, Oxford, UK, 1965.
- [17] C. LASSER, T. SWART, AND S. TEUFEL, *Construction and validation of a rigorous surface hopping algorithm for conical crossings*, Commun. Math. Sci., 5 (2007), pp. 789–814.
- [18] C. LASSER AND S. TEUFEL, *Propagation through conical crossings: An asymptotic semigroup*, Comm. Pure Appl. Math., 58 (2005), pp. 1188–1230.
- [19] H. SPOHN AND S. TEUFEL, *Adiabatic decoupling and time-dependent Born-Oppenheimer theory*, Comm. Math. Phys., 224 (2001), pp. 113–132.
- [20] J. TULLY AND R. PRESTON, *Trajectory surface hopping approach to nonadiabatic molecular collisions: The reaction of H^+ with D_2* , J. Chem. Phys., 55 (1971), pp. 562–572.
- [21] C. ZENER, *Non-adiabatic crossing of energy levels*, Proc. Roy. Soc. Lond., 137 (1932), pp. 696–702.

ON A MODEL OF MULTIPHASE FLOW*

DEBORA AMADORI[†] AND ANDREA CORLI[‡]

Abstract. We consider a hyperbolic system of three conservation laws in one space variable. The system is a model for fluid flow allowing phase transitions; in this case the state variables are the specific volume, the velocity, and the mass density fraction of the vapor in the fluid. For a class of initial data having large total variation we prove the global existence of solutions to the Cauchy problem.

Key words. hyperbolic systems of conservation laws, phase transitions

AMS subject classifications. 35L65, 35L60, 35L67, 76T30

DOI. 10.1137/07069211X

1. Introduction. We consider a model for the one-dimensional flow of an inviscid fluid capable of undergoing phase transitions. Both liquid and vapor phases are possible, as well as mixtures of them. In Lagrangian coordinates the model is

$$(1.1) \quad \begin{cases} v_t - u_x & = 0, \\ u_t + p(v, \lambda)_x & = 0, \\ \lambda_t & = 0. \end{cases}$$

Here $t > 0$ and $x \in \mathbb{R}$; moreover, $v > 0$ is the specific volume, u the velocity, and λ the mass density fraction of vapor in the fluid. Then $\lambda \in [0, 1]$, with $\lambda = 0$ characterizing the liquid and $\lambda = 1$ the vapor phase; the intermediate values of λ model the mixtures of the two pure phases. The pressure is denoted by $p = p(v, \lambda)$; under natural assumptions the system is strictly hyperbolic.

This model is a simplified version of a model proposed by Fan [13], where also viscous and relaxation terms were taken into account. The model is isothermal (see (2.1) below); in the presence of phase transitions this physical assumption is meaningful for retrograde fluids. A study of the Riemann problem for a 2×2 relaxation approximation of (1.1) has been done in [10]. We focus here on the global existence of solutions to the Cauchy problem for (1.1), namely, for initial data

$$(v, u, \lambda)(0, x) = (v_o(x), u_o(x), \lambda_o(x))$$

having finite total variation. This problem is motivated by the study of more complete models, where (1.1) is supplemented by source terms.

The problem of the global existence of solutions to strictly hyperbolic system of conservation laws has been studied for a long time; see [9, 11, 24, 25, 26] for general information. If the initial data have *small* total variation, then the Glimm theorem [14] applies; we refer the reader again to [9] for the analogous results obtained by a

*Received by the editors May 16, 2007; accepted for publication (in revised form) December 17, 2007; published electronically April 2, 2008. This work was supported by PROGETTO GNAMPA 2005 *Analisi Asintotica Per Sistemi Iperbolici Non Lineari*.

<http://www.siam.org/journals/sima/40-1/69211.html>

[†]Dipartimento di Matematica Pura e Applicata, Università degli Studi dell'Aquila, Via Vetoio, 67010 Coppito (AQ), Italy (amadori@univaq.it).

[‡]Dipartimento di Matematica, Università di Ferrara, Via Machiavelli 35, 44100 Ferrara, Italy (andrea.corli@unife.it).

wave-front tracking algorithm as well as for uniqueness and continuous dependence of the solutions on the initial data.

Some special systems allow, however, initial data with *large* total variation. For the system of isothermal gasdynamics Nishida [19] proved that it is sufficient that the variation $\text{TV}(v_o, u_o)$ of the initial data is finite in order to have globally defined solutions. This result was extended by Nishida and Smoller [20] to any pressure law $p = k/v^\gamma$, $\gamma > 1$, provided that $(\gamma - 1)\text{TV}(v_o, u_o)$ is small; related results are in [12]. For the full nonisentropic system of 3×3 gasdynamics, where $p = k \exp(\frac{\gamma-1}{R}s)/v^\gamma$ and s denotes the entropy, Liu [17, 16] proved the global existence of solutions if $(\gamma - 1)\text{TV}(v_o, p_o)$ is small and $\text{TV}(s_o)$ bounded. Temple [27] and Peng [21] obtained similar results. All these papers use the Glimm scheme. Analogous results making use of a wave-front tracking scheme have been given recently by Asakura [4, 5]; we point out that the use of wave-front tracking schemes in case of data with large variation is far from being trivial, and a deep analysis of the wave interactions is required. Very general results can be proved for systems with coinciding shock and rarefaction curves [8]; however, system (1.1) is not of this type.

In comparison with the above systems of gasdynamics, in (1.1) we keep a γ -law for the pressure with $\gamma = 1$ but add a dependence of p on λ : we then take $p = a(\lambda)/v$ for a suitable function a . System (1.1) has close connections to a system introduced by Benzoni-Gavage [6] and studied by Peng [22]; it seems, however, that the proof in [22] is not complete. A comparison of these models is done in subsection 3.1. We mention that the method of compensated compactness has also been applied to (1.1) (see [15, 7] and [18, sections 12.3 and 16]) but for different pressure laws.

In this paper we prove by a wave-front tracking scheme the global existence of solutions to (1.1) for a wide class of initial data with large total variation. We first introduce a *weighted total variation* (WTV) of $a(\lambda_o)$; this quantity arises in a natural way in the problem and also has an analytical meaning, being the logarithmic variation in the case of continuous functions. We prescribe a bound, on $\text{WTV}(a(\lambda_o))$; for the variation $\text{TV}(v_o, u_o)$ there is not such a bound, but, roughly speaking, the larger $\text{TV}(v_o, u_o)$ is, the smaller $\text{WTV}(a(\lambda_o))$ must be. An important point is that we give explicit expressions for these bounds; then our results are qualitatively different from some of those quoted above, where a generic smallness is required.

The plan of the paper is the following. The main result is stated in section 2, Theorem 2.2. The Riemann problem is reviewed in section 3 together with related results; proofs have been given in [2]. The definition of the algorithm is in section 4. The core of the proof is section 5—where interactions are studied in detail—and section 6—where we prove the convergence and consistence of the scheme. A careful analysis is needed due to the presence of large waves.

The paper is completed by two appendices. In the first one we prove the main result on the WTV. In the second we study the interaction of two shock waves to the light of section 5; namely, we look for precise bounds of the damping coefficient that controls the reflected wave produced in the interaction; we think that this analysis is interesting on its own. Good reading!

2. Main results. We consider the system of conservation laws (1.1). The pressure is given by

$$(2.1) \quad p(v, \lambda) = \frac{a^2(\lambda)}{v},$$

where a is a smooth (\mathbf{C}^1) function defined on $[0, 1]$ satisfying for every $\lambda \in [0, 1]$

$$(2.2) \quad a(\lambda) > 0, \quad a'(\lambda) > 0;$$

see Figure 2.1. For instance $a^2(\lambda) = k_0 + \lambda(k_1 - k_0)$ for $0 < k_0 < k_1$. As a consequence of (2.1) and (2.2) we have, for every $(v, \lambda) \in (0, +\infty) \times [0, 1]$,

$$(2.3) \quad p > 0, \quad p_v < 0, \quad p_{vv} > 0,$$

$$(2.4) \quad p_\lambda > 0, \quad p_{v\lambda} < 0.$$

Remark that assumptions (2.3) and (2.4) are analogous to those usually made on the pressure in the full nonisentropic case [17], the entropy replacing λ .

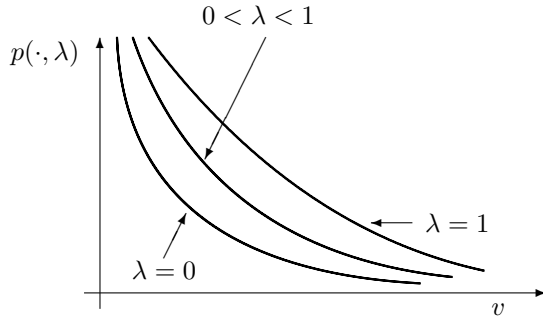


FIG. 2.1. Pressure curves as functions of v .

We denote $U = (v, u, \lambda) \in \Omega = (0, +\infty) \times \mathbb{R} \times [0, 1]$. Under assumptions (2.1) and (2.2) the system (1.1) is strictly hyperbolic in the whole Ω with eigenvalues $e_1 = -\sqrt{-p_v(v, \lambda)}$, $e_2 = 0$, $e_3 = \sqrt{-p_v(v, \lambda)}$. We write $c = \sqrt{-p_v} = a(\lambda)/v$. The eigenvectors associated with the eigenvalues e_i , $i = 1, 2, 3$, are $r_1 = (1, c, 0)$, $r_2 = (-p_\lambda, 0, p_v)$, $r_3 = (-1, c, 0)$. Because of the third inequality in (2.3) the eigenvalues e_1, e_3 are genuinely nonlinear with $\nabla e_i \cdot r_i = p_{vv}/(2c) > 0$, $i = 1, 3$, while e_2 is linearly degenerate. Pairs of Riemann invariants are $R_1 = \{u - a(\lambda) \log v, \lambda\}$, $R_2 = \{u, p\}$, $R_3 = \{u + a(\lambda) \log v, \lambda\}$.

We denote by $\text{TV}(f)$ the total variation of a function f . In the case $f : \mathbb{R} \rightarrow (0, +\infty)$ we define the *weighted total variation* of f by

$$\text{WTV}(f) = 2 \sup \sum_{j=1}^n \frac{|f(x_j) - f(x_{j-1})|}{f(x_j) + f(x_{j-1})},$$

where the supremum is taken over all $n \geq 1$ and $(n + 1)$ -tuples of points x_j with $x_0 < x_1 < \dots < x_n$. This variation is motivated by the definition (3.6) of strength for the waves of the second family. If f is bounded and bounded away from zero, then

$$\frac{1}{\sup f} \text{TV}(f) \leq \text{WTV}(f) \leq \frac{1}{\inf f} \text{TV}(f).$$

PROPOSITION 2.1. Consider $f : \mathbb{R} \rightarrow (0, +\infty)$; then

$$(2.5) \quad \frac{\inf f}{\sup f} \text{TV}(\log(f)) \leq \text{WTV}(f) \leq \text{TV}(\log(f)).$$

Moreover, if $f \in C(\mathbb{R})$, then $\text{WTV}(f) = \text{TV}(\log(f))$.

The proof is deferred to Appendix A. In (2.5), in the inequality on the right, the strict sign may occur if f is discontinuous; see Remark A.1.

We provide system (1.1) with initial data

$$(2.6) \quad U(x, 0) = U_o(x) = (v_o(x), u_o(x), \lambda_o(x))$$

for $x \in \mathbb{R}$. Denote $a_o(x) \doteq a(\lambda_o(x))$, $p_o(x) \doteq p(v_o(x), \lambda_o(x))$; remark that $\inf a_o(x) \geq a(0) > 0$. The main result of this paper now follows.

THEOREM 2.2. *Assume (2.1), (2.2). Consider initial data (2.6) with $v_o(x) \geq \underline{v} > 0$ for some constant \underline{v} and $0 \leq \lambda_o(x) \leq 1$. For every $m > 0$ and a suitable function $k(m) \in (0, 1/2)$ the following holds. If*

$$(2.7) \quad \text{TV}(\log(p_o)) + \frac{1}{\inf a_o} \text{TV}(u_o) < 2\left(1 - 2\text{WTV}(a_o)\right)m,$$

$$(2.8) \quad \text{WTV}(a_o) < k(m),$$

then the Cauchy problem (1.1), (2.6) has a weak entropic solution (v, u, λ) defined for $t \in [0, +\infty)$. Moreover, the solution is valued in a compact set of Ω , and there is a constant $C(m)$ such that for every $t \in [0, +\infty)$

$$(2.9) \quad \text{TV}(v(t, \cdot), u(t, \cdot)) \leq C(m).$$

The function $k(m)$, whose expression is given in (6.22), deserves some comments. The interaction of two waves α, α' of the same family $i = 1, 3$ produces a wave β of the same family i and a “reflected” wave δ of the other family j ($j = 1, 3, j \neq i$). For a suitable definition of the strengths of the waves we prove that $|\delta| \leq d \cdot \min\{|\alpha|, |\alpha'|\}$ for a damping coefficient $d < 1$ depending on α and α' ; see Lemma 5.6. The function k above depends essentially on the supremum of such coefficients d ; we prove that $k(0) = 1/2$ and that $k(m)$ decreases to 0 as $m \rightarrow +\infty$. In particular then $\text{WTV}(a_o) < 1/2$. The assumptions (2.7), (2.8) read as analogous to those in [20]: the larger m is, the smaller $k(m)$ is, and vice versa. The occurrence of a possible blow-up when the bound on $\text{WTV}(a_o)$ does not hold is an interesting open problem.

The variation of λ_o appears both in condition (2.7), because of p_o , and in (2.8). Using the definition of the pressure, we can replace (2.7) by the slightly stronger condition $\text{TV} \log(v_o) + 2\text{TV} \log(a_o) + \frac{1}{\inf a_o} \text{TV}(u_o) \leq 2(1 - 2\text{WTV}(a_o))m$ or even

$$\text{TV} \log(v_o) + \frac{1}{a(0)} \text{TV}(u_o) \leq 2m \left(1 - \frac{2m + 1}{m} \text{TV}(\log(a_o))\right)$$

by making use of (2.5). In particular if λ_o is constant, we recover the famous result by Nishida [19].

Clearly $\lambda(t, x) = \lambda_o(x)$ for any t because of the third equation in (1.1); this is why only v and u appear in the estimate (2.9). In other words system (1.1) can be rewritten as a p -system of two conservation laws with flux depending on x , namely, for the pressure law $p = p(v, \lambda_o(x)) = a^2(\lambda_o(x))/v$.

The proof of Theorem 2.2 makes use of a wave-front tracking scheme where we exploit the special structure of system (1.1) by differentiating the treatment of 1- and 3-waves from that of 2-waves. Our algorithm is a natural extension of that in [3], where the system for λ_o constant is studied, in the presence of a relaxation term.

Here we consider a linear functional as in [3] that accounts for the strengths of all 1- and 3-waves, with a weight $\xi > 1$ assigned to shock waves; a crucial point in

the proof is the choice of ξ as a function of m . This functional differs from that in [19, 4], where ξ is missing and only the variation of shocks is taken into account. Moreover, motivated again by [3], we do not introduce a simplified Riemann solver for interactions between 1- and 3-waves but only for interactions involving the 2-contact discontinuities. The interaction potential then uniquely considers interactions of 2-waves with 1- or 3-waves approaching it.

System (1.1) can be written in Eulerian coordinates. Denoting $\rho = 1/v$ the density, the pressure law becomes $p = a^2(\lambda)\rho$, and (1.1) turns into

$$(2.10) \quad \begin{cases} \rho_t + (\rho u)_x & = 0, \\ (\rho u)_t + (\rho u^2 + p(\rho, \lambda))_x & = 0, \\ (\rho \lambda)_t + (\rho \lambda u)_x & = 0. \end{cases}$$

A global existence result of weak solutions for (2.10) holds by Theorem 2.2 because of [28].

3. Preliminaries.

3.1. Comparison with other models. In [6] many models for diphasic flows are proposed and studied. In a simple case (no source terms, the fluid in either a dispersed or separated configuration) and keeping notation as in [6, page 35], they can be written as

$$(3.1) \quad \begin{cases} (\rho_l R_l)_t + (\rho_l R_l u_l)_x & = 0, \\ (\rho_g R_g)_t + (\rho_g R_g u_g)_x & = 0, \\ (\rho_l R_l u_l + \rho_g R_g u_g)_t + (\rho_l R_l u_l^2 + \rho_g R_g u_g^2 + p)_x & = 0. \end{cases}$$

Here the indexes l and g stand for *liquid* and *gas*. Therefore ρ_l , R_l , u_l are the liquid density, phase fraction, and velocity, and analogously for the gas; clearly $R_l + R_g = 1$. The pressure law is $p = a^2 \rho_g$ for $a > 0$ a constant. Equations (3.1) state the conservation of mass of either phase and the total momentum.

A case studied in [6, page 44] is when $u_l = u_g$ and ρ_l is constant, say, equal to 1. The unknown variables are then R_l , u , ρ_g , and it is assumed $0 < R_l < 1$ and $\rho_g > 0$; as a consequence $0 < R_g < 1$. Under these conditions, and still writing ρ_l instead of 1 for clarity, we define the concentration $c = \frac{\rho_g R_g}{\rho_l R_l} > 0$ and deduce the pressure law $p = a^2 c \frac{R_l}{1-R_l}$. We obtain exactly the model of [22]:

$$(3.2) \quad \begin{cases} (R_l)_t + (R_l u)_x & = 0, \\ (R_l c)_t + (R_l c u)_x & = 0, \\ (R_l(1+c)u)_t + (R_l(1+c)u^2 + p)_x & = 0. \end{cases}$$

This system is strictly hyperbolic for $c > 0$. Remark that the three eigenvalues of (3.2) coincide with u at $c = 0$, and if c vanishes identically, then (3.2) reduces to the pressureless gasdynamics system. System (3.2) is analogous to (2.10), but the pressure laws are different. In fact the variables ρ and λ of (2.10) write $\rho = \rho_l R_l + \rho_g R_g$ and $\lambda = \frac{\rho_g R_g}{\rho_l R_l + \rho_g R_g} = \frac{c}{1+c}$, and then $R_l = (1-\lambda)\rho$, $\rho_g = \frac{\lambda}{\rho^{-1} - (1-\lambda)}$. If we sum up the first two equations in (3.2), we find the first equation in (2.10); the third (resp., second) equation in (3.2) becomes the second (resp., third) equation in (2.10). The choice $p = a^2 \rho_g$ for the pressure in (3.1) gives $p = a^2 \frac{\lambda}{\rho^{-1} - (1-\lambda)}$.

Notice that the pressure vanishes in the presence of a pure liquid phase, and this is the main difference with (2.1), (2.2).

We now compare (2.10) and (3.2) in Lagrangian coordinates. Consider for (3.2) the change of coordinates $y = R_l dx - R_l u dt$ based on the streamlines of the liquid particles (because $R_l = \rho_l R_l$) [22]. Denote $w = \frac{1}{R_l} - 1 = \frac{R_g}{R_l} = \frac{c}{\rho_g}$. Then for $p = \frac{a^2 c}{w}$ system (3.2) turns into

$$(3.3) \quad \begin{cases} w_t - u_x & = 0, \\ ((1+c)u)_t + p_x & = 0, \\ c_t & = 0. \end{cases}$$

It is more interesting, however, to consider for system (3.2) the change $y = (1+c)R_l dx - (1+c)R_l u dt = \rho dx - \rho u dt$ into Lagrangian coordinates based on the streamlines of the full density ρ . Let w be as above and $v = \frac{w}{1+c} = \frac{R_g}{\rho}$. Then system (3.2) becomes system (1.1) with $a^2(\lambda) = a^2 \frac{c}{1+c} = a^2 \lambda$. As a consequence the pressure law $p(v, \lambda) = a^2(\lambda)/v$ does not satisfy (2.2). This difficulty can be overcome as follows. Fix any $0 < a_1 < a_2 < a$ and consider for $c_i = \frac{a_i^2}{a^2 - a_i^2}$ the invariant domain $\{(R_l, u, c): 0 < c_1 \leq c \leq c_2\}$ [22]. In this domain $0 < b_1 \leq \lambda \leq b_2 < 1$ for $b_i = a_i^2/a^2$. If we denote $\mu = \frac{\lambda - b_1}{b_2 - b_1}$, then the function $b(\mu) = a(\lambda) = a(b_1 + (b_2 - b_1)\mu)$ makes the pressure law $p(v, \lambda) = b^2(\mu)/v$, with $\mu \in [0, 1]$, satisfy both conditions in (2.2).

3.2. Wave curves and the Riemann problem. In this section we recall some results about the wave curves for system (1.1) and the solution to the Riemann problem; see [2] for more details.

The shock-rarefaction curves through the point $U_o = (v_o, u_o, \lambda_o)$ for (1.1) are

$$(3.4) \quad \Phi_i(v, U_o) = (v, \phi_i(v, U_o), \lambda_o), \quad i = 1, 3,$$

$$\phi_1(v, U_o) = \begin{cases} u_o + a(\lambda_o) \cdot (v - v_o) / \sqrt{v v_o}, & v < v_o, \quad \text{shock,} \\ u_o + a(\lambda_o) \log(v/v_o), & v > v_o, \quad \text{rarefaction,} \end{cases}$$

$$\phi_3(v, U_o) = \begin{cases} u_o - a(\lambda_o) \log(v/v_o), & v < v_o, \quad \text{rarefaction,} \\ u_o - a(\lambda_o) \cdot (v - v_o) / \sqrt{v v_o}, & v > v_o, \quad \text{shock,} \end{cases}$$

$$(3.5) \quad \Phi_2(\lambda, U_o) = \left(v_o \frac{a^2(\lambda)}{a^2(\lambda_o)}, u_o, \lambda \right), \quad \lambda \in [0, 1], \quad \text{contact discontinuity.}$$

The curves Φ_1, Φ_2 , and Φ_3 are plane curves: Φ_1 and Φ_3 lie on the plane $\lambda = \lambda_o$, while Φ_2 on $u = u_o$.

DEFINITION 3.1 (wave strengths). Under the notation (3.4), (3.5) we define the strength ε_i of an i -wave as

$$(3.6) \quad \varepsilon_1 = \frac{1}{2} \log \left(\frac{v}{v_o} \right), \quad \varepsilon_2 = 2 \frac{a(\lambda) - a(\lambda_o)}{a(\lambda) + a(\lambda_o)}, \quad \varepsilon_3 = \frac{1}{2} \log \left(\frac{v_o}{v} \right).$$

According to this definition, rarefaction waves have positive strengths and shock waves have negative strengths. Given the initial datum $\lambda_o = \lambda_o(x)$, denote

$$(3.7) \quad a^* \doteq \sup_{x \in \mathbb{R}} a(\lambda_o(x)), \quad a_* \doteq \inf_{x \in \mathbb{R}} a(\lambda_o(x)), \quad [a]_* \doteq \frac{a^* - a_*}{a^* + a_*}.$$

Then $[a]_* \leq \frac{a(1)-a(0)}{a(1)+a(0)} < 1$ and $|\varepsilon_2| \leq 2[a]_* < 2$. It is useful to also define the function (see [22])

$$(3.8) \quad h(\varepsilon) = \begin{cases} \varepsilon & \text{if } \varepsilon \geq 0, \\ \sinh \varepsilon & \text{if } \varepsilon < 0. \end{cases}$$

Then we have for $i = 1, 3$

$$(3.9) \quad \phi_i(v, U_o) = u_o + a(\lambda_o) \cdot 2h(\varepsilon_i).$$

At last we consider the Riemann problem. This is the initial value problem for (1.1) under the piecewise constant initial condition

$$(3.10) \quad (v, u, \lambda)(0, x) = \begin{cases} (v_\ell, u_\ell, \lambda_\ell) = U_\ell & \text{if } x < 0, \\ (v_r, u_r, \lambda_r) = U_r & \text{if } x > 0 \end{cases}$$

for U_ℓ and U_r in Ω . We denote $a_r = a(\lambda_r)$, $p_r = a_r^2/v_r$, and similarly for a_ℓ , p_ℓ .

PROPOSITION 3.2. *Fix any pair of states U_ℓ, U_r in Ω ; then the Riemann problem (1.1), (3.10) has a unique Ω -valued solution in the class of solutions consisting of simple Lax waves. If ε_i is the strength of the i -wave, $i = 1, 2, 3$, then*

$$\varepsilon_3 - \varepsilon_1 = \frac{1}{2} \log \left(\frac{p_r}{p_\ell} \right), \quad 2(a_\ell h(\varepsilon_1) + a_r h(\varepsilon_3)) = u_r - u_\ell.$$

Moreover, let $\underline{v} > 0$ be a fixed number. There exists a constant $C_1 > 0$ depending on \underline{v} and $a(\lambda)$ such that if $v_l, v_r \geq \underline{v}$, then

$$(3.11) \quad |\varepsilon_1| + |\varepsilon_2| + |\varepsilon_3| \leq C_1 |U_\ell - U_r|.$$

For the proof, see [2]. One can easily find that

$$(3.12) \quad |\varepsilon_1| + |\varepsilon_3| \leq \frac{1}{2} |\log(p_r) - \log(p_\ell)| + \frac{1}{2 \min\{a_\ell, a_r\}} |u_r - u_\ell| \\ \leq \frac{1}{2} |\log(v_r) - \log(v_\ell)| + |\log(a_r) - \log(a_\ell)| + \frac{1}{2 \min\{a_\ell, a_r\}} |u_r - u_\ell|.$$

We remark that for any Riemann data $(v_\ell, u_\ell, \lambda_\ell)$, (v_r, u_r, λ_r) , the λ component of the solution takes value λ_ℓ for $x < 0$ and λ_r for $x > 0$. The fact that the interfaces between different phases are connected by a stationary wave can then be interpreted as a “kinetic condition” [1], analogous to Maxwell’s rule.

4. The approximate solution. In this section we define a wave-front tracking scheme [9] to build up piecewise constant approximate solutions to (1.1). More precisely, we follow the algorithm introduced in [3].

First, we approximate the initial data. For any $\nu \in \mathbb{N}$ we take a sequence $(v_o^\nu, u_o^\nu, \lambda_o^\nu)$ of piecewise constant functions with a finite number of jumps such that

- (i) $\text{TV} p_o^\nu \leq \text{TV} p_o$, $\text{TV} u_o^\nu \leq \text{TV} u_o$, $\text{WTV} a(\lambda_o^\nu) \leq \text{WTV} a(\lambda_o)$, $\inf a_o^\nu \geq \inf a_o$;
- (ii) $\lim_{x \rightarrow -\infty} (v_o^\nu, u_o^\nu, \lambda_o^\nu)(x) = \lim_{x \rightarrow -\infty} (v_o, u_o, \lambda_o)(x)$;
- (iii) $\|(v_o^\nu, u_o^\nu, \lambda_o^\nu) - (v_o, u_o, \lambda_o)\|_{\mathbf{L}^1} \leq \frac{1}{\nu}$,

where $p_o^\nu = a^2(\lambda_o^\nu)/v_o^\nu$. Second, we define the approximate Riemann solver. We introduce positive parameters $\eta = \eta_\nu$, $\rho = \rho_\nu$; they control, respectively, the size of rarefactions and the threshold when a simplified Riemann solver is used. Define also a parameter $\hat{s} > 0$ strictly larger than all possible speeds of wave fronts of both families 1 and 3. These parameters will be determined at the end of section 6.

- At time $t = 0$ we solve the Riemann problems at each point of jump of $(v_o^\nu, u_o^\nu, \lambda_o^\nu)(0+, \cdot)$ as follows: shocks are not modified while rarefactions are approximated by fans of waves, each of them having size less than η . More precisely, a rarefaction of size ε is approximated by $N = \lceil \varepsilon/\eta \rceil + 1$ waves whose size is $\varepsilon/N < \eta$; we set their speeds to be equal to the characteristic speed of the state at the right. Then $(v, u, \lambda)(t, \cdot)$ is defined until some wave fronts interact; by slightly changing the speed of some waves [9] we can assume that only *two* fronts interact at a time.
- When two wave fronts of families either 1 or 3 interact we solve the Riemann problem at the interaction point. If one of the incoming waves is a rarefaction, after the interaction it is prolonged (if it still exists) as a single discontinuity with speed equal to the characteristic speed of the state at the right. If a new rarefaction is generated, we employ the Riemann solver described before and divide it into a fan of waves having size less than η .
- When a wave front of family either 1 or 3 interacts with a 2-wave we proceed as follows. Let δ_2 be the size of the 2-wave and δ the size of the other wave.
 - If $|\delta_2\delta| \geq \rho$, we solve the Riemann problem as above, that is, with the *accurate Riemann solver*.
 - If $|\delta_2\delta| < \rho$, we prolong the 1- or 3- wave with a wave of the same family and size. Since the two waves do not commute, a *nonphysical* front is introduced [9], with fixed speed $\hat{s} > 0$. The size of a nonphysical wave is set to be $|u_r - u_\ell|$, where u_ℓ, u_r are the u components of the left and right states of the wave. We call this solver the *simplified Riemann solver*.
- When a nonphysical front interacts with a front of family 1, 2, or 3 (“physical”), we prolong the solution with a physical wave of the same size and a nonphysical one, consequently computing the intermediate value.

We refer for the last two items to Proposition 5.12 below. Remark that two nonphysical fronts cannot interact since they have the same constant speed \hat{s} . We denote by \mathcal{NP} the set of nonphysical waves.

5. Interactions. Fix the index ν introduced in the previous section. We shall prove in subsection 6.1 that the algorithm described above is defined for any $t > 0$ and provides for any initial data $(v_o^\nu, u_o^\nu, \lambda_o^\nu)$ a piecewise constant approximate solution $(v^\nu, u^\nu, \lambda^\nu) = (v, u, \lambda)$, where we dropped for simplicity the index ν . Here we study the interaction of waves.

For $K_{np} > 0$ and $t > 0$ we define the functional L and the interaction potential Q , both referred to $(v, u, \lambda)(t, \cdot)$, by

$$\begin{aligned}
 L(t) &= \sum_{i=1,3} |\gamma_i| + K_{np} L_{np}, & L_{np} &= \sum_{\gamma \in \mathcal{NP}} |\gamma|, \\
 (5.1) \quad Q(t) &= \sum_{\gamma_3 \text{ at the left of } \delta_2} |\gamma_3| |\delta_2| + \sum_{\gamma_1 \text{ at the right of } \delta_2} |\delta_2| |\gamma_1|.
 \end{aligned}$$

Remark that L takes into account only the strengths of both 1- and 3-waves and that of nonphysical waves. For contact discontinuities we define

$$L_{cd} = \sum |\gamma_2| = \text{WTV}a(\lambda_o').$$

Finally, for $\xi \geq 1$ and $K \geq 0$ we introduce

$$(5.2) \quad L_\xi = \sum_{\substack{i=1,3 \\ \gamma_i > 0}} |\gamma_i| + \xi \sum_{\substack{i=1,3 \\ \gamma_i < 0}} |\gamma_i| + K_{np}L_{np},$$

$$(5.3) \quad F = L_\xi + KQ.$$

For simplicity we omitted noting the dependence on K_{np} in the functional L_ξ and on K_{np} , ξ , K in F ; the choice of K_{np} shall depend on that of K ; see Proposition 5.12.

Observe that if λ_o is constant, then $Q = 0$ and $F = L_\xi$, whose variation was analyzed in Lemma 3.2 of [3]. Hence we will assume from now on that

$$(5.4) \quad A_o \doteq \text{WTV}(a_o) > 0.$$

By assumption (i) in section 4, one has $L_{cd} \leq A_o$.

In the following sections we analyze in detail the different types of interactions. Recalling the definition of h , (3.8), and with the notation of Figure 5.1, we introduce the following identities (see (3.1), (3.2) in [22]):

$$(5.5) \quad \varepsilon_3 - \varepsilon_1 = \alpha_3 + \beta_3 - \alpha_1 - \beta_1,$$

$$(5.6) \quad a_\ell h(\varepsilon_1) + a_r h(\varepsilon_3) = a_\ell h(\alpha_1) + a_m h(\alpha_3) + a_m h(\beta_1) + a_r h(\beta_3).$$

Formula (5.5) does not depend on λ and follows easily by equating the specific volumes v before and after the interaction time. By equating the velocities u we obtain (5.6). These properties are a consequence of the definition (3.6) of the strengths for 1- and 3-waves and of (3.9).

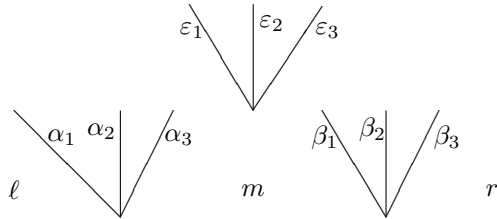


FIG. 5.1. A general interaction pattern.

5.1. Interactions with a 2-wave. We first consider the interactions of 1- or 3-waves with a 2-wave; see Figure 5.2.

PROPOSITION 5.1 (see [2]). Denote by λ_ℓ , λ_r the side states of a 2-wave. The interactions of 1- or 3-waves with the 2-wave give rise to the following pattern of solutions:

Interaction	Outcome	
	$\lambda_\ell < \lambda_r$	$\lambda_\ell > \lambda_r$
$2 \times 1R$	$1R + 2 + 3R$	$1R + 2 + 3S$
$2 \times 1S$	$1S + 2 + 3S$	$1S + 2 + 3R$
$3R \times 2$	$1S + 2 + 3R$	$1R + 2 + 3R$
$3S \times 2$	$1R + 2 + 3S$	$1S + 2 + 3S$

The next lemma is concerned instead with the *strengths* of waves involved in the interaction above. The inequalities (5.8) improve the inequality (3.3) in [22] in the special case of two interacting wave fronts, one of them being of the second family. More precisely, under the notation of [22] we find a term $1/(a_r + a_\ell)$ instead of $1/\min\{a_r, a_\ell\}$. The proof differs from Peng's. Our estimates are sharp: in some cases (5.8) reduces to an identity.

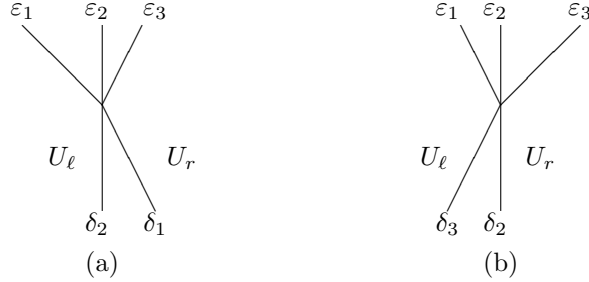


FIG. 5.2. Interactions. (a) from the right; (b) from the left.

LEMMA 5.2 (see [2]). Assume that a 1-wave of strength δ_1 or a 3-wave of strength δ_3 interacts with a 2-wave of strength $\delta_2 = 2(a_r - a_\ell)/(a_r + a_\ell)$. Then the strengths ε_i of the outgoing waves satisfy $\varepsilon_2 = \delta_2$ and

$$(5.7) \quad |\varepsilon_i - \delta_i| = |\varepsilon_j| \leq \frac{1}{2} |\delta_2| \cdot |\delta_i| \leq [a]_* |\delta_i|$$

for $i, j = 1, 3, i \neq j$. Moreover,

$$(5.8) \quad |\varepsilon_1| + |\varepsilon_3| \leq \begin{cases} |\delta_1| + |\delta_1|[\delta_2]_+ & \text{if 1 interacts,} \\ |\delta_3| + |\delta_3|[\delta_2]_- & \text{if 3 interacts.} \end{cases}$$

Here $[x]_+ = \max\{x, 0\}$, $[x]_- = \max\{-x, 0\}$, $x \in \mathbb{R}$. Remark that the colliding 1- or 3-wave does not change sign across the interaction. Moreover, the functional L increases iff the incoming and the reflected waves are of the same type; this happens when the colliding wave is moving toward a more liquid phase.

Now we prove that F is decreasing for suitable K when an interaction with a 2-wave occurs. The potential Q is needed to balance the possible increase of L_ξ .

PROPOSITION 5.3. Assume $A_o < 2$ and consider an interaction of a 1- or 3-wave with a 2-wave, with the notation of Lemma 5.2. Then $\Delta Q < 0$. If, moreover,

$$(5.9) \quad \xi \geq 1 \quad \text{and} \quad K > \frac{2\xi}{2 - A_o},$$

then

$$(5.10) \quad \xi |\varepsilon_j| = \xi (|\varepsilon_i| - |\delta_i|) < \frac{K}{2} |\Delta Q|,$$

and hence $\Delta F < 0$.

Proof. We consider the interaction of a 1-wave with a 2-wave; the symmetric case follows in an analogous way. We use the notation as in Figure 5.2(a). We define $L_{cd}^* = L_{cd}^- + L_{cd}^+$, L_{cd}^\pm , meaning *right* or *left* of the 2-wave under consideration.

By assumption, one has

$$(5.11) \quad L_{cd} = L_{cd}^- + L_{cd}^+ + |\delta_2| = L_{cd}^* + |\delta_2| \leq A_o < 2.$$

Recall that $\varepsilon_1 - \delta_1 = \varepsilon_3$ and by Proposition 5.1

$$(5.12) \quad |\varepsilon_1| - |\delta_1| = |\varepsilon_3| \quad \text{if } \delta_2 > 0,$$

$$(5.13) \quad |\varepsilon_1| - |\delta_1| = -|\varepsilon_3| \quad \text{if } \delta_2 < 0,$$

so that in particular $|\varepsilon_3| = ||\varepsilon_1| - |\delta_1||$. An estimate for ΔQ follows at once because of (5.11):

$$(5.14) \quad \begin{aligned} \Delta Q &= -|\delta_2 \delta_1| + (|\varepsilon_1| - |\delta_1|) L_{cd}^- + |\varepsilon_3| L_{cd}^+ \leq \frac{1}{2} |\delta_2 \delta_1| (L_{cd}^* - 2) \\ &\leq \frac{1}{2} |\delta_2 \delta_1| (A_o - 2) < 0. \end{aligned}$$

Hence, using (5.7), we get

$$(5.15) \quad \xi |\varepsilon_3| + \frac{K}{2} \Delta Q \leq \frac{1}{2} |\delta_2 \delta_1| \left\{ \xi + \frac{K}{2} (A_o - 2) \right\} < 0$$

because of (5.9); this proves (5.10). Finally, by using (5.15) we get

$$(5.16) \quad \Delta F = \Delta L_\xi + K \Delta Q \leq \xi |\varepsilon_3| + \xi |\varepsilon_1 - \delta_1| + K \Delta Q < 0. \quad \square$$

5.2. Interactions between 1- and 3-waves. Here we analyze the possible interactions between 1- and 3-waves. Two situations may occur; see Figure 5.3: either the waves belong to different families or they both belong to the same family. In this last case, at least one of the waves must be a shock.



FIG. 5.3. Interactions of 1- and 3-waves.

LEMMA 5.4 (different families interacting). *If a wave of the third family interacts with a wave of the first family, they cross each other without changing their strength.*

Proof. See also Lemma 3.1 in [3]. Using notation as in Figure 5.3(a) we have $\varepsilon_3 - \varepsilon_1 = \delta_3 - \delta_1$ and $h(\varepsilon_1) + h(\varepsilon_3) = h(\delta_1) + h(\delta_3)$. The uniqueness of solutions to the Riemann problem implies $\varepsilon_1 = \delta_1$, $\varepsilon_3 = \delta_3$.

Remark that here $\Delta L_\xi = 0 = \Delta Q$ and then $\Delta F = 0$ for all $\xi \geq 1$ and K . \square

LEMMA 5.5 (same family interacting: outcome). *Assume that a wave α_3 of the third family interacts with a wave β_3 of the third family, giving rise to waves ε_1 , ε_3 . Then*

$$(i) \quad \alpha_3 < 0, \beta_3 < 0 \Rightarrow \varepsilon_1 > 0, \varepsilon_3 < 0,$$

(ii) $\alpha_3\beta_3 < 0 \Rightarrow \varepsilon_1 < 0$.

An analogous result holds for interacting waves of the first family.

Proof. The proof can be done in a geometric way by observing the mutual positions of the curves [19, 26]. A simple alternative proof by analytical arguments now follows. We have

$$(5.17) \quad \varepsilon_3 - \varepsilon_1 = \alpha_3 + \beta_3,$$

$$(5.18) \quad h(\varepsilon_1) + h(\varepsilon_3) = h(\alpha_3) + h(\beta_3).$$

In case (i) these formulas read $\varepsilon_1 - \varepsilon_3 = |\alpha_3| + |\beta_3| > 0$ and $-h(\varepsilon_1) - h(\varepsilon_3) = \sinh(|\alpha_3|) + \sinh(|\beta_3|) > 0$. If it were $\varepsilon_3 > 0$, then $\varepsilon_1 > 0$ from the first equality and $\varepsilon_1 < 0$ from the second, a contradiction. Therefore $\varepsilon_3 < 0$ so that $\varepsilon_1 + |\varepsilon_3| = |\alpha_3| + |\beta_3|$ and $-h(\varepsilon_1) + \sinh(|\varepsilon_3|) = \sinh(|\alpha_3|) + \sinh(|\beta_3|)$. Analogously, if it were $\varepsilon_1 < 0$, using elementary inequalities we get $0 = \sinh(|\varepsilon_1|) + \sinh(|\alpha_3| + |\beta_3| + |\varepsilon_1|) - \sinh(|\alpha_3|) - \sinh(|\beta_3|) \geq 2\sinh(|\varepsilon_1|)$, a contradiction again. Hence $\varepsilon_1 > 0$.

In case (ii) assume $\alpha_3 < 0, \beta_3 > 0$; the other case is dealt with analogously since (5.17), (5.18) are symmetric in α_3, β_3 . We have $\varepsilon_3 - \varepsilon_1 = -|\alpha_3| + |\beta_3|$ and $h(\varepsilon_1) + h(\varepsilon_3) = -\sinh(|\alpha_3|) + \sinh(|\beta_3|)$; then $[h(\varepsilon_1) + \varepsilon_1] + [h(\varepsilon_3) - \varepsilon_3] = |\alpha_3| - \sinh(|\alpha_3|) < 0$. If $\varepsilon_3 > 0$, this last equality becomes $h(\varepsilon_1) + \varepsilon_1 = |\alpha_3| - \sinh(|\alpha_3|) < 0$, which implies $\varepsilon_1 < 0$. If $\varepsilon_3 < 0$, then $h(\varepsilon_1) + \varepsilon_1 = [|\alpha_3| - \sinh(|\alpha_3|)] - [|\varepsilon_3| - \sinh(|\varepsilon_3|)]$. If it were $\varepsilon_1 > 0$, it would be $|\alpha_3| < |\varepsilon_3|$, since the map $x \mapsto x - \sinh x$ is decreasing; but from $|\varepsilon_3| + |\varepsilon_1| = |\alpha_3| - |\beta_3|$ we would get that $|\varepsilon_3| < |\alpha_3|$, a contradiction. Hence in all cases one has $\varepsilon_1 < 0$. \square

Now we give sharper estimates for the interaction of waves of the same family: we prove that the strength of the reflected wave is bounded by the size of each incoming wave, multiplied by a damping factor smaller than 1. This property will be crucial in the next section, and it holds also for interactions with a 2-wave, with damping factor $[a]_*$; see (5.7). In the case below, however, the coefficient depends on the strengths of the incoming waves; this happens also when nonphysical waves are generated; see Proposition 5.12. We assume that

$$(5.19) \quad \begin{aligned} & \text{the strength of any interacting } i\text{-wave is less than } m \\ & \text{for some } m > 0 \text{ and } i = 1, 3. \end{aligned}$$

In the special case of interaction of waves of the same family producing two outgoing shocks we give a more precise result in Appendix B.

LEMMA 5.6 (same family interacting). *Consider the interaction of two waves of the same family, of sizes α_i and β_i , $i = 1, 3$, producing two outgoing waves $\varepsilon_1, \varepsilon_3$; assume (5.19). Then the following hold.*

(i) *There exists a damping coefficient $d = d(m)$, with $0 < d < 1$, such that*

$$(5.20) \quad |\varepsilon_j| \leq d(m) \cdot \min\{|\alpha_i|, |\beta_i|\}, \quad j \neq i.$$

(ii) *If the incoming waves are both shocks, the resulting shock satisfies $|\varepsilon_i| > \max\{|\alpha_i|, |\beta_i|\}$. If the incoming waves have different signs, both the amount of shocks and the amount of rarefactions of the i th family decrease across the interaction.*

In any case

$$(5.21) \quad |\varepsilon_i| \leq |\alpha_i| + |\beta_i|.$$

Proof. (i) To fix the ideas, assume $i = 3$. We have

$$(5.22) \quad \varepsilon_3 - \varepsilon_1 = \alpha_3 + \beta_3,$$

$$(5.23) \quad h(\varepsilon_1) + h(\varepsilon_3) = h(\alpha_3) + h(\beta_3),$$

and then

$$(5.24) \quad h(\varepsilon_1) + h(\varepsilon_1 + \alpha_3 + \beta_3) = h(\alpha_3) + h(\beta_3).$$

Remark that (5.24) is symmetric in α_3, β_3 ; from the implicit function theorem we find that $\varepsilon_1 = \varepsilon_1(\alpha_3, \beta_3)$ is \mathbf{C}^1 . Using the notation $\varepsilon_1 = \tau$, $\alpha_3 = a$, $\beta_3 = b$, the identity (5.24) rewrites as

$$(5.25) \quad h(\tau) + h(\tau + a + b) - h(a) - h(b) = 0,$$

with $\tau = \tau(a, b)$. One verifies that $\tau(a, 0) = \tau(0, b) = 0$ and that

$$\tau_a = \frac{h'(a) - h'(\tau + a + b)}{h'(\tau) + h'(\tau + a + b)}, \quad \tau_b = \frac{h'(b) - h'(\tau + a + b)}{h'(\tau) + h'(\tau + a + b)}.$$

As $a \rightarrow 0$, one has $\tau_a \rightarrow (1 - h'(b))/(1 + h'(b))$; then $|\tau_a(0, b)| < 1$, and it can be bounded by a positive constant less than 1 that depends on m . The same argument works for τ_b .

To complete the proof, we show that $|\tau| < \min\{|a|, |b|\}$ in the nontrivial case $a \neq 0 \neq b$. We argue by contradiction, using an argument of [23]. Suppose that $|\tau| \geq |a|$; we can assume $\tau > 0$, since the case $\tau < 0$ can be proved by using the equality (5.25) written in terms of $G(t) = -h(-t)$. Since the function h is increasing we have, for $\tau \geq |a|$,

$$h(\tau) \geq h(a), \quad h(\tau + a + b) \geq h(b).$$

Moreover, one of the two inequalities is strict: if $a < 0$, the first; if $a > 0$, the second. Hence we contradict (5.25).

(ii) From [3] we already know that $\Delta L = |\varepsilon_1| + |\varepsilon_3| - |\alpha_3| - |\beta_3| \leq 0$, and hence (5.21).

If the incoming waves are both shocks, then (5.22) becomes

$$(5.26) \quad |\varepsilon_3| + |\varepsilon_1| = |\alpha_3| + |\beta_3|.$$

From (i) we have $|\varepsilon_1| < |\alpha_3|, |\beta_3|$ and hence the first part of (ii). On the other hand, if $\alpha_3\beta_3 < 0$, we have $\varepsilon_1 < 0$. We can assume $\alpha_3 < 0 < \beta_3$; hence (5.22) becomes $\varepsilon_3 = |\beta_3| - |\alpha_3| - |\varepsilon_1|$. If $\varepsilon_3 > 0$, then $|\varepsilon_3| < |\beta_3|$; if $\varepsilon_3 < 0$, using (i) again one finds that $|\varepsilon_3| < |\alpha_3|$. \square

Remark 5.7. The damping coefficient $d(m)$ (see Figure 5.4) is given by

$$d(m) = \max_{\substack{|a| \leq m \\ |b| \leq m}} \frac{|\varepsilon(a, b)|}{\min\{|a|, |b|\}},$$

where the function $\varepsilon(a, b)$ satisfies $h(\varepsilon) + h(\varepsilon + a + b) - h(a) - h(b) = 0$; see (5.24). Hence $d(m)$ increases with m and vanishes as $m \rightarrow 0$ because quadratic interaction estimates hold for m small.

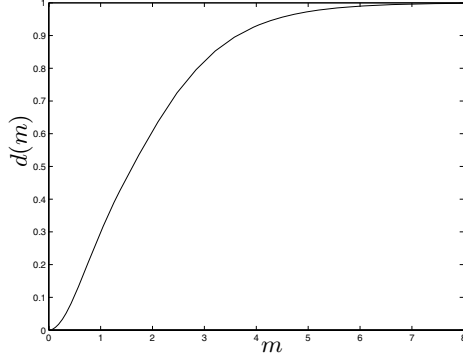


FIG. 5.4. The coefficient $d(m)$.

Moreover, it is asymptotic to 1 for m large. Indeed, from the proof of Lemma 5.6 we have $\tau_a(0, b) = \frac{1-h'(b)}{1+h'(b)}$; then $\tau_a(0, b) = 0$ if $b > 0$ and $|\tau_a(0, b)| = \frac{\cosh(b)-1}{\cosh(b)+1} \leq \frac{\cosh m-1}{\cosh m+1}$ if $b < 0$. Therefore $|\tau_a(0, b)| \leq \frac{\cosh m-1}{\cosh m+1}$ for every b , and an analogous estimate holds for $|\tau_b(0, a)|$. Hence $d(m) \geq \frac{\cosh m-1}{\cosh m+1} \doteq c(m)$; we refer to Lemma B.1 for the role of this quantity.

Remark 5.8. When a rarefaction interacts with a 1- or 3-wave, its size does not increase. Indeed, the size does not change upon interactions with waves of the other family by Lemma 5.4; if the rarefaction interacts with a shock of the same family, we apply Lemma 5.6(ii). Remark, moreover, that by Lemma 5.5, when a rarefaction and a shock of the same family interact, the reflected wave is never a rarefaction.

Remark 5.9. If two waves of the same family interact, the wave belonging to that family can be missing, while the “reflected” wave is always present. This follows easily from (5.22), (5.23).

PROPOSITION 5.10 (variation of F). *Consider the interactions of any two wave fronts of the same family, 1 or 3, and assume (5.4), (5.19). If*

$$(5.27) \quad 1 < \xi < \frac{1}{d} \quad \text{and} \quad K < \frac{\xi - 1}{A_o},$$

then $\Delta L_\xi < 0$ and $\Delta F < 0$.

Proof. Let two waves α_i, β_i interact, $i = 1, 3$, giving rise to waves $\varepsilon_1, \varepsilon_3$. We consider $i = 3$, the other case being analogous. Using (5.21), we get

$$(5.28) \quad \Delta Q = (|\varepsilon_3| - |\alpha_3| - |\beta_3|) L_{cd}^+ + |\varepsilon_1| L_{cd}^- \leq |\varepsilon_1| L_{cd}^- \leq |\varepsilon_1| A_o.$$

Now we claim that

$$(5.29) \quad \Delta L_\xi + |\varepsilon_1|(\xi - 1) \leq 0.$$

From this estimate it follows that $\Delta F = \Delta L_\xi + K\Delta Q \leq |\varepsilon_1|(1 - \xi + KA_o) < 0$ because of (5.27). To prove our claim we consider the possible cases; we make use of (5.22).

$\boxed{SS \rightarrow RS}$. Since $\Delta L = 0$, then

$$(5.30) \quad \Delta L_\xi + (\xi - 1)|\varepsilon_1| = \xi(|\varepsilon_1| + |\varepsilon_3| - |\alpha_3| - |\beta_3|) \leq 0.$$

$\boxed{SR, RS \rightarrow SR}$. Assume $\alpha_3 < 0 < \beta_3$; then (5.29) reads $(2\xi - 1)|\varepsilon_1| + |\varepsilon_3| - |\beta_3| - \xi|\alpha_3| \leq 0$. For later use we prove the stronger inequality

$$(5.31) \quad \xi^2|\varepsilon_1| + |\varepsilon_3| - |\beta_3| - \xi|\alpha_3| \leq 0.$$

Indeed, from Lemma 5.6(ii) we have $|\varepsilon_3| < |\beta_3|$, while $\xi|\varepsilon_1| \leq |\alpha_3|$ from (5.20), (5.27)₁.

$\boxed{SR, RS \rightarrow SS}$. Assume $\alpha_3 < 0 < \beta_3$; then (5.29) is $(2\xi - 1)|\varepsilon_1| + \xi(|\varepsilon_3| - |\alpha_3|) - |\beta_3| \leq 0$. We prove also in this case the stronger inequality

$$(5.32) \quad \xi^2|\varepsilon_1| + \xi(|\varepsilon_3| - |\alpha_3|) - |\beta_3| \leq 0.$$

Indeed, by (5.17) and again because of (5.20), (5.27)₁, one has

$$\begin{aligned} \xi^2|\varepsilon_1| + \xi(|\varepsilon_3| - |\alpha_3|) - |\beta_3| &= \xi^2|\varepsilon_1| + \xi(|\varepsilon_1| - |\beta_3|) - |\beta_3| \\ &= (\xi + 1)(\xi|\varepsilon_1| - |\beta_3|) \leq 0. \end{aligned}$$

This proves the claim and concludes the proof. \square

Remark 5.11. From the above proof we see that $\Delta L_\xi \leq 0$ for $\xi = 1$. This was a key point in [19], where, however, a different choice of strengths was made. In [3] the inequality $\Delta L_\xi \leq 0$ was proved to hold also for $1 < \xi \leq \xi_o$ for some $\xi_o > 1$; the condition (5.27)₁ gives an estimate of such a threshold.

More precisely, in the first two cases of Proposition 5.10 we have $\Delta L_\xi \leq 0$ for every $\xi \geq 1$. The third case is analyzed in detail in Lemma B.1; we prove there that $\Delta L_\xi \leq 0$ for any $\xi > 1$ if $c(m) \leq 1/2$, while we need $1 < \xi \leq \frac{1}{2c(m)-1}$ if $c(m) > 1/2$.

5.3. Nonphysical waves. In this subsection we compute the strength of a non-physical wave generated by an interaction and prove that it does not change in subsequent interactions. We introduce the following notation: given $U_\ell = (v_\ell, u_\ell, \lambda_\ell)$ and λ_r we define by

$$U_{\ell r}^* = \Phi_2(\lambda_r, U_\ell) = (A_{r\ell}v_\ell, u_\ell, \lambda_r)$$

the state on the right of a 2-wave with left state $U_\ell = (v_\ell, u_\ell, \lambda_\ell)$ and $\lambda = \lambda_r$ on the right, where $A_{r\ell} = a^2(\lambda_r)/a^2(\lambda_\ell)$. See (3.5) and [2].

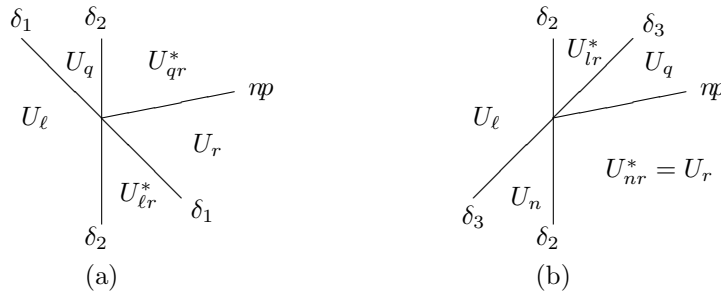


FIG. 5.5. *Simplified Riemann solver.*

PROPOSITION 5.12 (nonphysical waves). Consider $U_\ell = (v_\ell, u_\ell, \lambda_\ell)$. Let $U_r = (v_r, u_r, \lambda_r)$ be connected to $U_{\ell r}^*$ by a 1-wave of size δ_1 and $U_q = (v_q, u_q, \lambda_\ell)$ be connected to U_ℓ by a 1-wave of size δ_1 ; see Figure 5.5(a). Assume (5.19).

Then U_{qr}^* and U_r differ only in the u component; if δ_2 denotes the size of the 2-wave, there exists a constant $C_o = C_o(m)$ such that

$$(5.33) \quad \|U_{qr}^* - U_r\| = |u_q - u_r| \leq C_o |\delta_2 \delta_1|.$$

A similar result holds for the interaction of a 3-wave (see Figure 5.5(b)), again under (5.19).

Moreover, the size of a nonphysical wave does not change in subsequent interactions. For any $K > 0$ and $K_{np} < K/C_o$ at any interaction involving a nonphysical wave we have

$$(5.34) \quad \Delta F \leq 0,$$

with $\Delta F < 0$ when a nonphysical wave is generated.

Proof. Recalling [2, Lemma 2], only the u component will be different after commutation of the 1- and the 2-wave. We find that

$$u_q - u_\ell = 2a_\ell h(\delta_1), \quad u_r - u_\ell = 2a_r h(\delta_1),$$

and hence

$$|u_q - u_r| = 2|a_\ell - a_r| \cdot |h(\delta_1)| \leq |\delta_2 \delta_1| \cdot 2a(1) \max_{0 < \eta \leq m} \frac{\sinh \eta}{\eta}.$$

Then (5.33) follows with $C_o(m) \doteq 2a(1) \cdot \frac{\sinh m}{m}$.

Next, assume that a nonphysical wave interacts with a 2-wave. Since the values of u do not change across a 2-wave, the left and right values of u of the nonphysical wave do not change across the interaction; hence the size does not change.

Assume then that a nonphysical wave interacts with a 1- or 3-wave of size δ . Since λ is constant, we refer only to the components v, u . Let (v_ℓ, u_ℓ) and (v_ℓ, u_q) be the side states of the nonphysical wave before the interaction and $(v_\ell, u_q), (v_r, u_r)$ be the side states of the physical wave. After the interaction, let $(\tilde{v}_\ell, \tilde{u}_\ell)$ be the intermediate state. One has

$$u_r - u_q = 2a(\lambda)h(\delta) = \tilde{u}_\ell - u_\ell,$$

and hence $|u_q - u_\ell| = |u_r - \tilde{u}_\ell|$.

At last, we consider the functional F . The potential Q is unaltered when nonphysical waves interact with other waves. The only cases in which L_{np} changes are when a nonphysical wave arises. Assume that a 1- or a 3-wave of size δ interacts with a 2-wave of size δ_2 , producing a wave of the same size and a nonphysical wave. Then $\Delta Q = -|\delta\delta_2|$ and $\Delta L_\xi = K_{np}\Delta L_{np} \leq K_{np}C_o|\delta\delta_2|$; hence $\Delta F = \Delta L_\xi + K\Delta Q \leq |\delta\delta_2|(K_{np}C_o - K)$, and then (5.34). \square

5.4. Decreasing of the functional F and control of the variations. We first collect the previous results into a single proposition.

PROPOSITION 5.13 (local decreasing). *Consider the interaction of any two waves either of families 1, 2, 3 or nonphysical. Assume (5.19) for some $m > 0$; let $C_o = C_o(m)$ as in Proposition 5.12. Finally, let A_o satisfy*

$$(5.35) \quad 0 < A_o < 2 \frac{1-d}{3-d}.$$

If ξ , K , K_{np} satisfy

$$(5.36) \quad \frac{2 - A_o}{2 - 3A_o} < \xi < \frac{1}{d}, \quad \frac{2\xi}{2 - A_o} < K < \frac{\xi - 1}{A_o}, \quad K_{np} < \frac{K}{C_o},$$

then

$$(5.37) \quad \Delta F \leq 0.$$

Proof. The condition on K comes from (5.9) and (5.27)₂. The interval where K lies is not empty if $A_o < \frac{2}{3}$ and $\xi > \frac{2 - A_o}{2 - 3A_o}$; together with (5.27)₁ this gives the assumption required on ξ . In turn, it is possible to choose ξ in such an interval if (5.35) holds. Remark that $2\frac{1-d}{3-d} \leq \frac{2}{3}$, so the previous condition on A_o holds. Therefore the assumptions of Propositions 5.3, 5.10, and 5.12 hold, and then (5.37) follows. \square

PROPOSITION 5.14 (global decreasing). *Let $m > 0$; assume (5.35), (5.36),*

$$(5.38) \quad L(0+) < \frac{m}{2\xi - 1}$$

and that the approximate solution U is defined in $[0, T]$. Then $L(t) < m$ for any $t \in (0, T]$; as a consequence, condition (5.19) holds for any 1- or 3-wave in U . Finally, $\Delta F(t) \leq 0$ for all times $t \in [0, T]$.

Proof. Since L may change value only at the times of interaction, we use an induction argument based on Proposition 5.13.

First, we have $L(0+) < m$ because $\xi \geq 1$ and (5.38). Now assume that $L(\tau) < m$ for all $0 < \tau < t$. By Proposition 5.13 one has both $\Delta F(\tau) \leq 0$ and $\Delta F(t) \leq 0$, so that

$$F(t) \leq F(0+) \leq \xi L(0+) + KQ(0+).$$

Since $Q(0+) \leq L(0+)L_{cd} \leq L(0+)A_o$, and recalling (5.36)₂, we get

$$(5.39) \quad F(t) \leq L(0+) \cdot (\xi + KA_o) \leq L(0+) \cdot (2\xi - 1) < m,$$

again because of (5.38). Finally, since $L(t) \leq L_\xi(t) \leq F(t)$ we arrive at the conclusion. \square

Remark 5.15. From (5.38) we see that, in order to have $L(t) < m$, the smaller ξ is, the larger the $L(0+)$ can be chosen.

Remark 5.16. From Proposition 5.14 we deduce that v remains bounded away from zero. Indeed, recalling (3.6), (3.5) and that v is constant across nonphysical waves, we have

$$\frac{1}{2} \text{TV}(\log v(t, \cdot)) \leq L(t) + \text{TV}(\log a_o) \leq m + \text{TV}(\log a_o).$$

6. The convergence and the consistence of the algorithm. In this section we prove Theorem 2.2. We first show that for fixed ν the algorithm introduced in section 4 gives an approximate solution defined for every $t > 0$; more precisely, we prove that at every time the number of interactions is bounded. Then we prove that the total amount of nonphysical waves in each approximate solution is very small. The convergence of a suitable subsequence is assured by Helly's theorem; then consistence follows.

6.1. Control of the number of interactions. We prove first that the size of the rarefactions in the scheme is small.

LEMMA 6.1. *Consider a rarefaction of size ε ; then*

$$(6.1) \quad |\varepsilon| < \eta e^{\frac{A_o}{2}}.$$

Proof. We analyze all possible situations. When the rarefaction is generated, one has $0 < \varepsilon < \eta$. When it interacts with a 1- or 3-wave, the size does not increase; see Remark 5.8. By Proposition 5.12 the size does not change when interactions with nonphysical waves occur.

The last case to be considered is when a rarefaction interacts with a 2-wave. In this case the size may increase; however, a rarefaction can meet a fixed 2-wave only once. Consider the case of a 1-rarefaction of size δ_1 , as in Proposition 5.3, the other being analogous. If $\delta_2 < 0$, then the size decreases; see (5.13). If $\delta_2 > 0$, by (5.12) we have

$$|\varepsilon_1| = |\delta_1| + |\varepsilon_3| \leq |\delta_1| \left(1 + \frac{1}{2} |\delta_2| \right) < |\delta_1| e^{\frac{|\delta_2|}{2}}.$$

Summarizing the three cases above, we get $|\varepsilon| < \eta e^{L_{cd}^+/2}$ (or $|\varepsilon| < \eta e^{L_{cd}^-/2}$) for a 1-rarefaction (resp., 3-rarefaction), where L_{cd}^\pm is the sum of the 2-waves at the right or left of the rarefaction. Then (6.1) follows. \square

Next, we prove that the number of interactions remains bounded in finite time, so that the approximate solution is well defined for all $t > 0$. We first give a lemma.

LEMMA 6.2. *Consider the wave-front tracking algorithm described in section 4, under the assumptions of Proposition 5.14. Then*

- (i) *the number of interactions involving a 2-wave and solved by the accurate Riemann solver is finite;*
- (ii) *the number of interactions where a new rarefaction of size $\varepsilon \geq \eta$ arises is finite.*

Proof. First consider (i) and refer to Proposition 5.3. Then, using (5.16), we have $\Delta F \leq \rho(\xi + K(A_o - 2)/2) < 0$, and hence F decreases by a uniform positive quantity; since it is nonincreasing, this can happen only a finite number of times.

Then consider (ii). After (i), it remains only to consider the case of two shocks of the same family interacting.

Under the notation of the corresponding case in the proof of Proposition 5.10 we have $\varepsilon = \varepsilon_1 \geq \eta$ and

$$\Delta F \leq |\varepsilon_1| (1 - \xi + KA_o) \leq \eta (1 - \xi + KA_o) < 0$$

because of (5.36). Arguing similarly as in (i), this can happen only a finite number of times. \square

Regarding (ii) in Lemma 6.2, recall that if $\varepsilon \geq \eta$, then the new rarefaction must be split into more than one wave. Therefore Lemma 6.2 can be rephrased by saying that, except for finite interactions, the number of waves emitted in an interaction is at most three, and this case occurs precisely when a nonphysical wave is generated; moreover, in every interaction at most one wave per family is emitted.

In a schematic way, apart from a finite number of interactions, in our algorithm the following hold (we will consider the set of nonphysical waves as a fourth family of waves):

- (a) the interaction of an i -wave, $i = 1, 3$, with a 2-wave is solved by a single i -wave, a 2-wave, and a 4-wave;
- (b) in the interaction of just 1- and/or 3-waves, there is at most one outgoing wave of each family 1 and 3;
- (c) the interaction of a i -wave, $i = 1, 2, 3$, with a 4-wave is solved by an i -wave and a 4-wave.

The next proposition extends the result of Lemma 2.5 in [3].

PROPOSITION 6.3. *Consider the wave-front algorithm described in section 4 and assume in the strip $[0, T) \times \mathbb{R}$ the following:*

for some $a_1 < a_2 < 0 < b_1 < b_2$ the waves of the first (resp., third) family have speeds in the interval $[a_1, a_2]$ (resp., $[b_1, b_2]$).

Then the number of interactions in the region $[0, T) \times \mathbb{R}$ is finite.

Proof. Assume by contradiction that in the region $[0, T) \times \mathbb{R}$ there exists an infinite number of interactions. Unless we take a smaller T we can assume that the number of interactions is finite in every strip $[0, t) \times \mathbb{R}$, $0 < t < T$, and that T is an accumulation point for the times of interaction. Because of the finite propagation speed, the interaction points are also bounded in space.

Then there exists a bounded sequence (t_j, x_j) , $j = 1, 2, \dots$, of interaction points such that

$$0 < t_j < t_{j+1} < T \text{ for all } j \text{ and } (t_j, x_j) \rightarrow (T, \bar{x})$$

for some \bar{x} . Denote $\mathcal{J} = \{(t_j, x_j): j = 1, 2, \dots\}$ the set of all interaction points.

The situation described in items (a)–(c) above holds except in a finite number of interactions; let $\tau < T$ be the maximum time of these “exceptional” interactions. It is not restrictive to assume that $\tau < t_j < T$ for all $j = 1, 2, \dots$.

Starting from a point of \mathcal{J} , we “trace back” all the segments up to $t = 0$; we repeat the procedure for all the points of \mathcal{J} and call \mathcal{F} the set of the “traced segments” obtained in this way. In other words, a segment belongs to the set \mathcal{F} iff it can be joined forward in time to some point of \mathcal{J} by a continuous path along the wave fronts. The set \mathcal{F} is not empty: for instance, two segments interacting at the point (t_j, x_j) belong to \mathcal{F} for any $j = 1, 2, \dots$. Observe the following dichotomy property of \mathcal{F} which is used just below:

two interacting waves either both belong to \mathcal{F} or none of them does;
 moreover, if at least one of the outgoing waves belong to \mathcal{F} , then
 both the incoming waves must belong to \mathcal{F} .

We now partition all the interaction points of the algorithm that occur for times $t > \tau$ into the following sets:

- \mathcal{I}_0 : the interaction points where no ingoing wave belongs to \mathcal{F} ;
- \mathcal{I}_1 : the interaction points where both incoming waves belong to \mathcal{F} and at most one outgoing segment belongs to \mathcal{F} ;
- \mathcal{I}_2 : the interaction points where exactly two outgoing segments belong to \mathcal{F} ;
- \mathcal{I}_3 : the interaction points where three outgoing waves all belong to \mathcal{F} .

Because of the dichotomy property quoted above no outgoing wave in case \mathcal{I}_0 can belong to \mathcal{F} . On the contrary, both incoming waves in $\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_3$ must belong to \mathcal{F} .

Recall that we are considering times $t > \tau$. Therefore the maximum number of emitted waves in an interaction is three, and this happens only in the situation considered in \mathcal{I}_3 , that is, for interactions as in (a) above. The case of more than one emitted wave per family cannot occur, and so the outgoing waves in \mathcal{I}_2 belong to different families. By definition we have $\mathcal{J} \cap \mathcal{I}_0 = \emptyset$, and so $\mathcal{J} \subset \mathcal{I}_1 \cup \mathcal{I}_2 \cup \mathcal{I}_3$.

Let $\mathcal{V}(t)$ be the total number of wave fronts of the families 1, 2, and 3 that belong to \mathcal{F} at time t . The functional $\mathcal{V}(t)$ is nonincreasing, and it decreases at least by 1 across \mathcal{I}_1 . Then \mathcal{I}_1 is finite. As a consequence, all the interaction points of \mathcal{J} belong to $\mathcal{I}_2 \cup \mathcal{I}_3$, except at most a finite number. Let $\tau_1 \in [\tau, T)$ be a time such that all points in \mathcal{I}_1 lie in $t < \tau_1$.

Let $P = \{x_1, \dots, x_{N_1}\}$ be the set of points of the x -axis where a 2-wave is located. We consider two cases and make use of [3, Lemma 2.5].

$\boxed{\bar{x} \notin P}$. In this case we can choose a time $\tau_1 < T$ such that after that time no segment belonging to \mathcal{F} crosses a 2-wave. Then all the points in \mathcal{J} with $t_j > \tau_1$ belong to \mathcal{I}_2 .

Take a point $(t^*, x^*) \in \mathcal{J}$ with $t^* > \tau_1$, so that $(t^*, x^*) \in \mathcal{I}_2$. At (t^*, x^*) there are two outgoing segments, both of them belonging to \mathcal{F} ; for $t > t^*$ and t close to t^* we define $\gamma_\ell(t)$ to be the segment on the left and $\gamma_r(t)$ that on the right. When $\gamma_\ell(t)$ (resp., $\gamma_r(t)$) reaches an interaction point, we prolong it by the segment of \mathcal{F} outgoing from that point and located on the left (resp., on the right).

In this way we define recursively for any $t < T$ two paths: $\gamma_\ell(t)$ is made by segments of the families 1 or 3, while $\gamma_r(t)$ is made by segments of the families 3 or 4.

We claim that the speeds of the two paths are strictly separated. Indeed, if γ_r starts following a 3-segment, then γ_ℓ starts with a 1-segment and so will always follow 1-segments; therefore $\dot{\gamma}_\ell(t) \leq a_2 < 0 < b_1 \leq \dot{\gamma}_r(t)$. If γ_r starts following a 4-segment, then it will always follow 4-segments; in this case $\dot{\gamma}_\ell(t) \leq b_2 < \hat{s} = \dot{\gamma}_r(t)$.

Thus for $c = \min\{\hat{s} - b_2, b_1 - a_2\} > 0$ we have $\gamma_r(t) - \gamma_\ell(t) \geq c(t - t^*)$. Now set $\rho = c(T - t^*)$ and choose $t_n \in (t^*, T)$, with $(t_n, x_n) \in \mathcal{J}$, such that

$$(6.2) \quad c(t_n - t^*) > \frac{\rho}{2}, \quad (T - t_n)(|a_1| + \hat{s}) < \frac{\rho}{4}.$$

By definition of \mathcal{F} , the points $(t_n, \gamma_\ell(t_n))$ and $(t_n, \gamma_r(t_n))$ can be joined forward in time to points (t_h, x_h) , resp., (t_k, x_k) , of \mathcal{J} , with $h, k > n$. Then

$$x_h \leq \gamma_\ell(t_n) + \hat{s}(t_h - t_n), \quad x_k \geq \gamma_r(t_n) - |a_1|(t_k - t_n).$$

Thanks to (6.2), we get

$$x_k - x_h \geq c(t_n - t^*) - (|a_1| + \hat{s})(T - t_n) > \frac{\rho}{4}.$$

Since n can be taken arbitrarily large, the last inequality contradicts the convergence of the sequence x_j .

$\boxed{\bar{x} \in P}$. Consider case (a) above. The possibility that the outgoing waves belonging to \mathcal{F} are precisely one physical and one nonphysical wave may happen only a finite number of times, since the functional $\mathcal{V}(t)$ is nonincreasing. Therefore we can assume that, for $t > \tau_1$, in case (a) the two outgoing physical waves belong to \mathcal{F} .

Let $(t^*, x^*) \in \mathcal{J}$ with $t^* > \tau_1$. As before we define for $t \in [t^*, T)$ two continuous paths $\gamma_\ell(t)$, $\gamma_r(t)$ starting at (t^*, x^*) in the following way.

At (t^*, x^*) there are either two or three outgoing segments belonging to \mathcal{F} ; for times $t > t^*$ and sufficiently close to t^* we define $\gamma_\ell(t)$ to be the segment on the left and $\gamma_r(t)$ the one on the right. When $\gamma_\ell(t)$ ($\gamma_r(t)$) reaches an interaction point, it is prolonged by the segment of \mathcal{F} on the left (resp., on the right). Then the path $\gamma_\ell(t)$ is made by segment of the families 1, 2, or 3, while $\gamma_r(t)$ is made by segments of families 3 or 4; in fact the interaction of a 1-wave with a 2-wave always produces a 4-wave (except in a finite number of cases).

Now we prove that the speeds of the paths $\gamma_\ell(t)$, $\gamma_r(t)$ are strictly separated. Indeed, if $\gamma_r(t)$ starts following a 3-segment, then γ_ℓ starts with either a 2- or a 1-segment and, by the remark at the beginning of this subcase, $\dot{\gamma}_\ell(t) \leq 0 < b_1 \leq \dot{\gamma}_r(t)$. If $\gamma_r(t)$ starts following a 4-segment, then $\dot{\gamma}_\ell(t) \leq b_2 < \hat{s} = \dot{\gamma}_r(t)$.

Finally, for $c = \min\{\hat{s} - b_2, b_1\}$ the same argument exploited in the other case leads to a contradiction. \square

6.2. Control of the total size of nonphysical fronts. Assume as above that the assumptions of Proposition 5.14 hold. We assign inductively to each wave α a *generation order* k_α as in [9, page 140]. This is done according to the following procedure. First, at time $t = 0$ each wave has order 1. Second, assume that two waves α and β interact at time t ; if α and β belong to different families, the outgoing waves of those families keep the order of the incoming waves, and the other waves assume order $\max\{k_\alpha, k_\beta\} + 1$; if α and β belong to the same family, and the outgoing wave of that family takes the order $\min\{k_\alpha, k_\beta\}$, and the other waves are assigned order $\max\{k_\alpha, k_\beta\} + 1$.

When specialized to the current setting this has the following consequences:

- every 2-wave has order 1; when an i -wave, $i = 1, 3$, of order k interacts with a 2-wave, the outgoing i -wave has order k , and the other outgoing wave (of the family j , $j = 1, 3$, $j \neq i$, or a nonphysical wave) has order $k + 1$;
- in the interaction of a 1- with a 3-wave the waves cross without changing order; in the interaction of two waves α , β of the same family $i = 1, 3$, the outgoing wave of the family i takes order $\min\{k_\alpha, k_\beta\}$, and the wave of the family $j = 1, 3$, $j \neq i$, has order $\max\{k_\alpha, k_\beta\} + 1$;
- when a nonphysical wave interacts with any other wave, both waves cross without changing order; in particular a nonphysical wave keeps the order it has been assigned when generated.

For $t \geq 0$ not an interaction time and any $k = 1, 2, \dots$ define (see (5.2), (5.1))

$$V_k(t) = \sum_{\substack{\gamma > 0 \\ k_\gamma = k}} |\gamma| + \xi \sum_{\substack{\gamma < 0 \\ k_\gamma = k}} |\gamma| + K_{np} \sum_{\substack{\gamma \in \mathcal{NP} \\ k_\gamma = k}} |\gamma|,$$

$$Q_k(t) = \sum_{\substack{\gamma_3 \text{ at the left of } \delta_2 \\ k_{\gamma_3} = k}} |\gamma_3| |\delta_2| + \sum_{\substack{\gamma_1 \text{ at the right of } \delta_2 \\ k_{\gamma_1} = k}} |\delta_2| |\gamma_1|,$$

$$F_k(t) = V_k(t) + KQ_k(t),$$

and

$$\tilde{V}_k(t) = \sum_{\ell \geq k} V_\ell(t), \quad \tilde{Q}_k(t) = \sum_{\ell \geq k} Q_\ell(t), \quad \tilde{F}_k(t) = \tilde{V}_k(t) + K\tilde{Q}_k(t).$$

We remark that

$$\tilde{F}_1(0+) = L_\xi(0+) + KQ(0+), \quad \tilde{F}_k(0+) = 0 \quad \text{for } k \geq 2.$$

Observe that if a nonphysical front interacts with another wave, the functionals above do not change; the same holds for interactions between 3- and 1-waves. Then we focus on interactions of waves of the same i family, $i = 1, 3$ (as usual denote $j = 1, 3$, $j \neq i$), and on interactions between 1- or 3-waves with a 2-wave.

For $h \in \mathbb{N}$, denote by I_h the set of times t when an interaction occurs between two waves α and β of families 1 or 3 with $\max\{k_\alpha, k_\beta\} = h$; denote by J_h the set of interaction times t of a 1- or 3-wave of order h with a 2-wave. Finally, denote $\mathcal{T}_h = I_h \cup J_h$ and $I = \bigcup_{h \geq 1} I_h$, $J = \bigcup_{h \geq 1} J_h$, $\mathcal{T} = I \cup J$.

In order to control the total size of nonphysical fronts we must strengthen the assumptions (5.35), (5.36) required in Proposition 5.14. First, for any fixed $m > 0$, instead of (5.35) we require the stronger condition

$$(6.3) \quad 0 < A_o < \frac{1 - \sqrt{d}}{2 - \sqrt{d}}.$$

Then denote

$$(6.4) \quad \lambda \doteq \frac{1 + KA_o}{\xi}, \quad \lambda_2 \doteq \frac{\xi + KA_o}{K(2 - A_o) - \xi}, \quad \mu \doteq \max \left\{ \lambda, \lambda_2, \frac{K_{np}C_o}{K} \right\}.$$

We need $0 < \mu < 1$. From (5.36)₃ we have $K_{np}C_o < K$; moreover, we have $0 < \lambda < 1$ and $\lambda_2 > 0$ because of (5.36)₂. At last $\lambda_2 < 1$ holds iff $\frac{\xi}{1 - A_o} < K$; this condition is stronger than the left inequality in (5.36)₂. Hence, instead of (5.36) we assume

$$(6.5) \quad \frac{1 - A_o}{1 - 2A_o} < \xi < \frac{1}{\sqrt{d}}, \quad \frac{\xi}{1 - A_o} < K < \frac{\xi - 1}{A_o}, \quad K_{np} < \frac{K}{C_o}.$$

Remark the new upper bound required on ξ . As in the proof of Proposition 5.13, the interval where K varies is not empty if $A_o < \frac{1}{2}$ and $\xi > \frac{1 - A_o}{1 - 2A_o}$. In turn, we can find ξ satisfying (6.5)₁ if (6.3) holds; remark that $\frac{1 - \sqrt{d}}{2 - \sqrt{d}} \leq \frac{1}{2}$. The last condition in (6.5) coincides with that in (5.36).

Remark 6.4. Proposition 5.14 still holds under the stronger assumptions (6.3), (6.5) under the same condition (5.38), because the inequality on the right-hand side in (6.5)₂ has not changed; see (5.39).

PROPOSITION 6.5. *Fix $m > 0$ and assume (6.3), (6.5). We have the following:*

1. $\tau \in \mathcal{T}_h$, $h \leq k - 2$: then $\Delta \tilde{F}_k = \Delta F_k = 0$.
2. $\tau \in \mathcal{T}_{k-1}$: then $\Delta F_{k-1} < 0$, $\Delta \tilde{F}_k = \Delta F_k > 0$, and

$$(6.6) \quad [\Delta \tilde{F}_k]_+ \leq \mu \left([\Delta F_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta F_\ell]_+ \right).$$

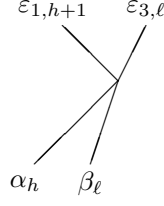
3. $\tau \in \mathcal{T}_h$, $h \geq k$: if $h = k$, then $\Delta F_k < 0$; in any case $\Delta \tilde{F}_k < 0$ and

$$(6.7) \quad \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+ < [\Delta \tilde{F}_k]_-.$$

Proof. As we pointed out above, for $\tau \in I$ only interactions of waves of the same family are taken into account; see Figure 6.1. Remark that by Proposition 5.14 we have

$$(6.8) \quad \sum_{\ell=1}^{k-1} \Delta F_\ell + \Delta \tilde{F}_k < 0.$$

1. If $h \leq k - 2$, no waves with order $\geq k$ are involved, and then $\Delta \tilde{F}_k = \Delta F_k = 0$.
2. Let $h = k - 1$. First, consider $\tau \in I_{k-1}$; then $\Delta \tilde{V}_k = \Delta V_k > 0$. We prove that

FIG. 6.1. *Interactions of 3-waves; $h \geq \ell$ denote generation orders.*

$$(6.9) \quad [\Delta \tilde{V}_k]_+ \leq \frac{1}{\xi} \left([\Delta V_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta V_\ell]_+ \right).$$

Indeed, from (5.30)–(5.32), we deduce that

$$(6.10) \quad \xi[\Delta \tilde{V}_k]_+ + \Delta V_{k-1} + \sum_{\ell=1}^{k-2} \Delta V_\ell \leq 0.$$

If $\min\{k_\alpha, k_\beta\} = k - 1$, then $\Delta V_\ell = 0$ for $\ell = 1, \dots, k - 2$ and (6.10) becomes $\xi[\Delta \tilde{V}_k]_+ + \Delta V_{k-1} < 0$; this implies $\Delta V_{k-1} < 0$ and hence $[\Delta \tilde{V}_k]_+ < (1/\xi)[\Delta V_{k-1}]_-$, that is, (6.9). Moreover, in this case, one easily finds that $\Delta Q_{k-1} \leq 0$, because of (5.21).

If $\min\{k_\alpha, k_\beta\} = \ell \leq k - 2$, then $\Delta V_{k-1} < 0$ and $\Delta Q_{k-1} < 0$, since no waves of order $k - 1$ are present after the interaction. Therefore the estimate (6.10) becomes

$$\xi[\Delta \tilde{V}_k]_+ - [\Delta V_{k-1}]_- + \Delta V_\ell \leq 0.$$

If $\Delta V_\ell \geq 0$, we get (6.9). If $\Delta V_\ell < 0$, we check directly that $\xi[\Delta \tilde{V}_k]_+ - [\Delta V_{k-1}]_- \leq \xi^2|\varepsilon_1| - |\alpha_3| < 0$ because of (5.20) and (6.5)₁. This completes the proof of (6.9).

From this proof we see also that $\Delta F_{k-1} = \Delta V_{k-1} + K\Delta Q_{k-1} < 0$, since both terms in the sum are negative. Then we use (6.9) and $0 \leq \Delta \tilde{Q}_k \leq A_o \Delta \tilde{V}_k$ to get

$$(6.11) \quad 0 < \Delta \tilde{F}_k \leq (1 + KA_o)[\Delta \tilde{V}_k]_+ \leq \lambda \left([\Delta V_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta V_\ell]_+ \right).$$

We now prove that

$$(6.12) \quad [\Delta Q_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta Q_\ell]_+ \geq 0.$$

We have only to consider the case in which $[\Delta Q_\ell]_+ > 0$ for some $\ell \leq k - 2$; but in this case $[\Delta Q_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta Q_\ell]_+ = L_{cd}^+(|\alpha_3| + |\beta_3| - |\varepsilon_3|) \geq 0$ because of (5.21).

Therefore (6.6) for $\tau \in I_{k-1}$ follows from (6.11) and (6.12).

Second, assume $\tau \in J_{k-1}$; we prove that

$$(6.13) \quad [\Delta \tilde{F}_k]_+ \leq \mu[\Delta F_{k-1}]_-.$$

Indeed, if the reflected wave is a physical wave, then, under the notation of Proposition 5.3,

$$\Delta V_{k-1} \leq \xi \frac{|\delta_1 \delta_2|}{2}, \quad \Delta Q_{k-1} \leq -|\delta_1 \delta_2| + \frac{|\delta_1 \delta_2|}{2} A_o = -\frac{|\delta_1 \delta_2|}{2} (2 - A_o)$$

so that $\Delta F_{k-1} \leq [\xi - K(2 - A_o)] |\delta_1 \delta_2| / 2 < 0$ because of (6.5). Then (6.13) follows since $\Delta \tilde{F}_k = \Delta F_k > 0$ and

$$[\Delta \tilde{F}_k]_+ \leq \frac{|\delta_1 \delta_2|}{2} (\xi + K A_o) = \frac{|\delta_1 \delta_2|}{2} [K(2 - A_o) - \xi] \cdot \lambda_2 \leq \lambda_2 [\Delta F_{k-1}]_- \leq \mu [\Delta F_{k-1}]_-.$$

If the reflected wave is a nonphysical wave, we have, under the notation of Proposition 5.12,

$$0 < \Delta F_k = \Delta V_k \leq K_{np} C_o |\delta \delta_2|, \quad \Delta V_{k-1} = 0, \quad \Delta Q_{k-1} = -|\delta \delta_2|,$$

and then

$$[\Delta F_{k-1}]_- = K |\delta \delta_2|, \quad [\Delta F_k]_+ \leq \frac{K_{np} C_o}{K} [\Delta F_{k-1}]_- \leq \mu [\Delta F_{k-1}]_-.$$

The estimate (6.13) is then completely proved. From (6.13) we get (6.6) since no 1-, 3-, or nonphysical waves of order $\leq k - 2$ are involved in the interaction.

3. Finally, let $h \geq k$. We first consider $\tau \in I_h$. If $\min\{k_\alpha, k_\beta\} \geq k$, then $\Delta \tilde{V}_k = \Delta L_\xi < 0$ and $\Delta \tilde{F}_k = \Delta F < 0$. If $\min\{k_\alpha, k_\beta\} \leq k - 1$, assume $k_\alpha \geq k$ and $k_\beta \leq k - 1$; then $\Delta \tilde{V}_k \leq \xi |\varepsilon_1| - |\alpha_3| \leq (\xi d - 1) |\alpha_3| < 0$ by (5.20) and (5.36)₁. Moreover, $\Delta \tilde{Q}_k = |\varepsilon_1| L_{cd}^- - |\alpha_3| L_{cd}^+ \leq |\varepsilon_1| A_o \leq d A_o |\alpha_3|$ again by (5.20). Then

$$\Delta \tilde{F}_k \leq [(\xi + K A_o) d - 1] \cdot |\alpha_3|,$$

and since, because of (6.5),

$$(\xi + K A_o) d - 1 < (2\xi - 1) d - 1 < \left(\frac{2}{\sqrt{d}} - 1 \right) d - 1 = -(1 - \sqrt{d})^2 < 0,$$

we proved that $\Delta \tilde{F}_k < 0$. Now from (6.8) we deduce that $\sum_{\ell=1}^{k-1} \Delta F_\ell < [\Delta \tilde{F}_k]_-$. Since at most one nonzero term is present in the first sum, (6.7) follows. If $h = k$, we are in the same situation as in case 2, so we have $\Delta F_k < 0$.

Now assume $\tau \in J_h$. Then $\Delta \tilde{F}_k = \Delta F < 0$ and no 1-, 3-, or nonphysical waves of order $< k$ are present, so (6.7) holds. If $h = k$ and the reflected wave is physical, from the proof of Proposition 5.3 and (6.5) we find that

$$\Delta F_k = \Delta V_k + K \Delta Q_k \leq \frac{|\delta_1 \delta_2|}{2} (\xi - 2K + K A_o) < 0.$$

If $h = k$ and the reflected wave is nonphysical, then $\Delta V_k = 0$, $\Delta Q_k < 0$, and $\Delta F_k < 0$. \square

Summarizing, for $\tau \in \mathcal{T}$ we have the following table:

$\mathcal{T}_h; h$	$\leq k - 2$	$k - 1$	k	$\geq k + 1$
ΔF_k	0	+	-	\pm
$\Delta \tilde{F}_k$	0	+	-	-

We write $\tilde{F}_k^\pm(t) = \sum_{\tau \leq t} [\Delta \tilde{F}_k(\tau)]_\pm$ for $k \geq 2$. For simplicity the time τ in such sums is omitted.

LEMMA 6.6. *Under the assumptions of Proposition 6.5 we have*

$$(6.14) \quad \tilde{F}_2^+(t) \leq \mu (L_\xi(0) + KQ(0)) + \sum_{\mathcal{T}_h, h \geq 2} [\Delta F_1]_+,$$

$$(6.15) \quad \tilde{F}_k^+(t) \leq \mu \left(\tilde{F}_{k-1}^+(t) + \sum_{\mathcal{T}_h, h \geq k} [\Delta F_{k-1}]_+ - \sum_{\mathcal{T}_{k-1}} \sum_{\ell=1}^{k-2} [\Delta F_\ell]_+ \right), \quad k \geq 3.$$

Proof. By Proposition 6.5 (see also the table above), the functional F_k increases at times $\tau \in \mathcal{T}_{k-1}$ and decreases at times $\tau \in \mathcal{T}_k$, while it does not have a given sign at times $\tau \in \mathcal{T}_h$, with $h \geq k+1$.

First, by summing up (6.6) we obtain

$$(6.16) \quad \tilde{F}_k^+(t) \leq \mu \sum_{\mathcal{T}_{k-1}} \left([\Delta F_{k-1}]_- - \sum_{\ell=1}^{k-2} [\Delta F_\ell]_+ \right)$$

for $k \geq 2$, where the last term in (6.16) is missing if $k = 2$.

Now recall that $F_1(0) = L_\xi(0) + KQ(0)$; therefore

$$F_1(t) \leq L_\xi(0) + KQ(0) - \sum_{\mathcal{T}_1} [\Delta F_1]_- + \sum_{\mathcal{T}_h, h \geq 2} [\Delta F_1]_+$$

and then

$$(6.17) \quad \sum_{\mathcal{T}_1} [\Delta F_1]_- \leq L_\xi(0) + KQ(0) + \sum_{\mathcal{T}_h, h \geq 2} [\Delta F_1]_+.$$

On the other hand, $F_k(0) = 0$ for $k \geq 2$; from Proposition 6.5 we have

$$F_k(t) \leq \sum_{\mathcal{T}_{k-1}} [\Delta F_k]_+ - \sum_{\mathcal{T}_k} [\Delta F_k]_- + \sum_{\mathcal{T}_h, h \geq k+1} [\Delta F_k]_+.$$

Moreover,

$$\sum_{\mathcal{T}_{k-1}} [\Delta F_k]_+ = \sum_{\mathcal{T}_{k-1}} [\Delta \tilde{F}_k]_+ = \tilde{F}_k^+(t)$$

and then

$$(6.18) \quad \sum_{\mathcal{T}_k} [\Delta F_k]_- \leq \tilde{F}_k^+(t) + \sum_{\mathcal{T}_h, h \geq k+1} [\Delta F_k]_+.$$

From (6.16), (6.17), (6.18) we get (6.14), (6.15). \square

PROPOSITION 6.7 (a contraction property). *Under the assumptions of Proposition 6.5, for any $t \geq 0$ and $k \geq 1$ we have*

$$(6.19) \quad \tilde{V}_k(t) \leq \tilde{F}_k(t) \leq \mu^{k-1} \cdot (L_\xi(0) + KQ(0)).$$

Proof. The estimate (6.19) holds for $k = 1$ because $\tilde{F}_1(t) = F(t) \leq L_\xi(0) + KQ(0)$. Next, we prove by induction on $k \geq 2$ that for any t

$$(6.20) \quad \tilde{F}_k^+(t) \leq \mu^{k-1} (L_\xi(0) + KQ(0)) + \sum_{\mathcal{T}_h, h \geq k} \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+.$$

Since by summing up (6.7) we obtain

$$(6.21) \quad \tilde{F}_k^-(t) \geq \sum_{\mathcal{T}_h, h \geq k} \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+,$$

then (6.19) will follow from (6.20) for any $k \geq 2$ because of (6.21).

Formula (6.20) for $k = 2$ reduces to (6.14). Next, assume that (6.20) holds for some $k \geq 2$. By (6.15) and the induction assumption

$$\begin{aligned} \tilde{F}_{k+1}^+(t) &\leq \mu \left(\tilde{F}_k^+(t) + \sum_{\mathcal{T}_h, h \geq k+1} [\Delta F_k]_+ - \sum_{\mathcal{T}_k} \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+ \right) \\ &\leq \mu^k (L_\xi(0) + KQ(0)) + \mu \left(\sum_{\mathcal{T}_h, h \geq k} \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+ + \sum_{\mathcal{T}_h, h \geq k+1} [\Delta F_k]_+ - \sum_{\mathcal{T}_k} \sum_{\ell=1}^{k-1} [\Delta F_\ell]_+ \right) \\ &\leq \mu^k (L_\xi(0) + KQ(0)) + \mu \sum_{\mathcal{T}_h, h \geq k+1} \sum_{\ell=1}^k [\Delta F_\ell]_+. \end{aligned}$$

Since $\mu < 1$, we get (6.20) for $k + 1$. \square

Remark 6.8. We now comment on the case $A_o = 0$. In this case system (1.1) reduces to the p -system with pressure law given by (2.1) for fixed λ . According to our front-tracking algorithm, stationary and nonphysical waves do not appear, and so $L_{cd} = Q = 0$; the algorithm reduces to the one introduced in [3]. Then Proposition 6.5 holds with \tilde{V}_k and $1/\xi$ replacing \tilde{F}_k and μ , respectively, and at last (6.19) reads

$$\tilde{V}_k(t) \leq \frac{1}{\xi^{k-1}} \cdot L_\xi(0).$$

Next, we conclude that the total strength of all nonphysical waves is small by proceeding as in [9, page 142]. Recall the notation in section 4. First, by Remark 5.16, the sequence $\{v^\nu\}$ is uniformly bounded from above and away from 0; then the eigenvalues e_1 and e_3 are bounded, and this makes possible the choice of a suitable \hat{s} . We have two more parameters η, ρ to be chosen. Fix $\eta > 0$ with the condition $\eta = \eta_\nu \rightarrow 0$ as $\nu \rightarrow \infty$ and estimate the total number of waves of order $< k$. We have

$$\begin{aligned} \sum_{\gamma \in \mathcal{NP}} |\gamma|(t) &\leq \tilde{V}_k(t) + \sum_{\gamma \in \mathcal{NP}, k_\gamma < k} |\gamma|(t) \\ &\leq \mu^{k-1} \cdot (L_\xi(0) + KQ(0)) + C_o \rho \cdot [\text{number of fronts of order } < k] \leq \frac{1}{\nu} \end{aligned}$$

by choosing k sufficiently large to have the first term $\leq 1/(2\nu)$ and then choose ρ small enough to have the second term $\leq 1/(2\nu)$.

We now accomplish the proof of Theorem 2.2. First define

$$(6.22) \quad k(m) = \frac{1 - \sqrt{d(m)}}{2 - \sqrt{d(m)}}.$$

From the properties of the function $d(m)$ stated in Remark 5.7 we see that $k(0) = 1/2$ and that $k(m)$ is decreasing, tending to 0 for $m \rightarrow +\infty$. The assumption (2.8) implies that (6.3) holds.

Now, by hypotheses (2.7) it follows that we can choose ξ such that

$$\frac{1}{2}\text{TV}\log(p_o) + \frac{1}{2\inf a_o}\text{TV}(u_o) < \frac{m}{2\xi - 1} < (1 - 2A_o)m$$

and that (6.5)₁ holds. Hence, using (3.12) and (i) in section 4, we have

$$L(0+) \leq \frac{1}{2}\text{TV}\log(p_o) + \frac{1}{2\inf a_o}\text{TV}(u_o) < \frac{m}{2\xi - 1}$$

so that the hypotheses (5.38) of Proposition 5.14 hold. Theorem 2.2 now follows along the lines of [9, section 7.4].

Appendix A. The weighted total variation. In this appendix we prove Proposition 2.1. Remark that the map $d(a, b) \doteq \frac{|a-b|}{a+b}$ is a distance on \mathbb{R}_+ , as one can easily prove.

We start with the proof of the inequality on the right in (2.5). It is enough to prove that

$$(A.1) \quad 2 \sum_{j=1}^n \frac{|f(x_j) - f(x_{j-1})|}{f(x_j) + f(x_{j-1})} \leq \sum_{j=1}^n |\log f(x_j) - \log f(x_{j-1})|.$$

We claim that

$$(A.2) \quad \log t \geq \frac{2(t-1)}{t+1} \quad \text{for } t \geq 1,$$

where the inequality is strict if $t > 1$. To prove the claim it is sufficient to notice that the function $\phi(t) = \log t - \frac{2t-1}{t+1}$ vanishes in 1 and $\phi'(t) = \frac{(t-1)^2}{t(t+1)^2} > 0$ if $t > 1$.

We apply (A.2) to $t = x/y$ for $0 < y \leq x$ and arguing by symmetry deduce that

$$|\log x - \log y| \geq \frac{2|x-y|}{x+y} \quad \text{for every } x, y > 0.$$

Then (A.1) follows. The proof of the inequality on the left in (2.5) is analogous, starting from the inequality

$$(A.3) \quad \frac{1}{t} \log t \leq \frac{2(t-1)}{t+1} \quad \text{for } t \geq 1$$

with strict inequality if $t > 1$.

Now assume that $f \in C(\mathbb{R})$. To show that $\text{WTV}(f) = \text{TV}(\log(f))$, we have to prove that the inequality

$$(A.4) \quad \text{TV}(\log(f)) \leq 2 \sup \sum_{j=1}^n \frac{|f(x_j) - f(x_{j-1})|}{f(x_j) + f(x_{j-1})}$$

holds for any $[a, b] \subset \mathbb{R}$, the supremum being taken on the set of all partitions $a = x_o < x_1 < \dots < x_n = b$, $n \in \mathbb{N}$. Consider any such partition; by the mean value theorem and by the intermediate value theorem applied to f , we get

$$\begin{aligned} |\log(f(x_j)) - \log(f(x_{j-1}))| &= \frac{|f(x_j) - f(x_{j-1})|}{\zeta_j} \\ &= \frac{f(x_j) + f(x_{j-1})}{2f(\eta_j)} \cdot 2 \frac{|f(x_j) - f(x_{j-1})|}{f(x_j) + f(x_{j-1})} \end{aligned}$$

for some ζ_j between $f(x_j)$ and $f(x_{j-1})$ and $\eta_j \in [x_{j-1}, x_j]$. We exploit again the continuity of f in $[a, b]$. On one hand, its image is compact; then $\min_{[a,b]} f = m > 0$. On the other hand, f is uniformly continuous in $[a, b]$, so that for any $\varepsilon > 0$ there exists $\delta_\varepsilon > 0$ such that $|f(x) - f(y)| < \varepsilon$ if $|x - y| < \delta_\varepsilon$ for $x, y \in [a, b]$.

Now fix any $\varepsilon > 0$; without loss of generality we can consider partitions of the interval $[a, b]$ of mesh less than δ_ε . Assume for instance $f(x_{j-1}) \leq f(x_j)$; then from the inequalities

$$f(x_j) - f(x_{j-1}) < \varepsilon, \quad f(x_{j-1}) \leq f(\eta_j) \leq f(x_j)$$

it follows that

$$\frac{f(x_j) + f(x_{j-1})}{2f(\eta_j)} \leq \frac{2f(x_{j-1}) + \varepsilon}{2f(x_{j-1})} \leq 1 + \frac{\varepsilon}{m}.$$

The inequality (A.4) follows by remarking that then for any partition of mesh less than δ_ε

$$\sum_{j=1}^n |\log(f(x_j)) - \log(f(x_{j-1}))| \leq \left(1 + \frac{\varepsilon}{m}\right) 2 \sum_{j=1}^n \frac{|f(x_j) - f(x_{j-1})|}{f(x_j) + f(x_{j-1})}.$$

The proof of Proposition 2.1 is complete.

Remark A.1. Observe that if $\text{TV}(\log(f)) < \infty$ and f is discontinuous, then the inequality on the right in (2.5) is strict, because of the strict inequality in (A.2). For example, if f has a single jump and assumes the values $c > 0$ and $d > 0$, then $\text{WTV}(f) = 2 \frac{|c-d|}{c+d} < |\log c - \log d| = \text{TV}(\log(f))$.

Remark, moreover, that WTV and TV are not equivalent, in the sense that there does not exist a positive constant C such that $C \cdot \text{TV}(\log(f)) \leq \text{WTV}(f)$. This follows from the fact that clearly the inequality $C \log t \leq \frac{2(t-1)}{t+1}$ does not hold for every $t \geq 1$.

Appendix B. Shock-rarefaction interactions. In this appendix we consider a particular case of Lemma 5.6. Actually, this is the only case needed in order to define a decreasing functional (see section 5); however, we needed further analysis for the control and treatment of the nonphysical waves.

LEMMA B.1 (the case $SR, RS \rightarrow SS$). *Consider the interaction of a shock α_i and a rarefaction β_i of the same family, $i = 1, 3$, producing two outgoing shocks $\varepsilon_1, \varepsilon_3$. Then there exists a smooth function B satisfying $|\alpha_i| \leq B(\alpha_i) \leq \min\{\sinh(|\alpha_i|), 2|\alpha_i|\}$ such that*

$$(B.1) \quad 0 < \beta_i \leq B(\alpha_i).$$

Moreover, assume

$$(B.2) \quad |\alpha_i| \leq m$$

for some $m > 0$ and denote $c = c(m) = \frac{\cosh(m)-1}{\cosh(m)+1}$. Then both the variation of shock waves and the reflected wave ε_j , $j \neq i$, $j = 1, 3$, are estimated by the interacting rarefaction as

$$(B.3) \quad |\varepsilon_1| + |\varepsilon_3| - |\alpha_i| \leq (2c - 1) \cdot |\beta_i|,$$

$$(B.4) \quad |\varepsilon_j| \leq c \cdot |\beta_i|.$$

Proof. We focus on the case $i = 3, j = 1$; see Figure 5.3(b). Therefore we consider $\alpha_3 < 0, \beta_3 > 0$ and $\varepsilon_1 < 0, \varepsilon_3 < 0$. Then (5.5), (5.6) become

$$(B.5) \quad |\varepsilon_1| - |\varepsilon_3| = |\beta_3| - |\alpha_3|,$$

$$(B.6) \quad \sinh(|\varepsilon_1|) + \sinh(|\varepsilon_3|) = \sinh(|\alpha_3|) - |\beta_3|.$$

From the second equation $|\varepsilon_1| < |\alpha_3|, |\varepsilon_3| < |\alpha_3|$, and $|\beta_3| < \sinh(|\alpha_3|)$; using the first equation and $|\varepsilon_3| < |\alpha_3|$, we get $|\varepsilon_1| < |\beta_3|$. Therefore in conclusion

$$(B.7) \quad |\varepsilon_1| < \min\{|\alpha_3|, |\beta_3|\}, \quad |\varepsilon_3| < |\alpha_3|, \quad |\beta_3| < \sinh(|\alpha_3|).$$

Step 1: notation. We set $x = |\beta_3|, y = |\varepsilon_1|, z = |\alpha_3|$, so that

$$(B.8) \quad |\varepsilon_3| = y - x + z;$$

see Figure B.1. Under this notation, (B.6) writes as

$$(B.9) \quad F(x, y; z) = \sinh y + \sinh(y - x + z) - \sinh z + x = 0$$

for $x \geq 0, y \geq 0, z \geq 0, y - x + z \geq 0$. By (B.7), any solution of (B.9) satisfies

$$(B.10) \quad y < z, \quad y < x, \quad x < \sinh(z).$$

Step 2: the threshold. Observe that, despite the last inequality in (B.7), we may well have $|\beta_3| > |\alpha_3|$, that is, $x > z$. Consider in fact the limit case of $\varepsilon_3 = 0$: we have $y = x - z > 0$ and $\sinh(y) = \sinh(z) - x$ that give

$$\sinh(x - z) = \sinh(z) - x, \quad x > z.$$

The last equality is the relation needed for β_3, α_3 in order to have that the shock and rarefaction cancel out exactly, giving rise only to a wave of the opposite family. Observe that the size of the rarefaction must be larger than the one of the shock. Under the notation above, the threshold curve separating the case of the outgoing waves S_1S_3 from the case S_1R_3 is given by

$$(B.11) \quad f(x, z) = \sinh(x - z) - \sinh(z) + x = 0.$$

Since $f_x = \cosh(x - z) + 1 > 0$ and $f(z, z) < 0$, the implicit equation $f(x, z) = 0$ is solved by $x = x_o(z) \geq z$ with $x'_o(z) = \frac{\cosh(x-z) + \cosh(z)}{\cosh(x-z) + 1} > 0$ for every $z \geq 0$; the curve has for tangent at $(0, 0)$ the line $z = x$. Observe that $x_o(z) \leq 2z$ because $f(2z, z) = z \geq 0$. In conclusion

$$(B.12) \quad z \leq x_o(z) \leq 2z;$$

see Figure B.2(a). This estimate and the third inequality in (B.7) prove (B.1) for $B(\alpha_i) = x_o(|\alpha_i|)$. We can prove more than (B.12), that is,

$$(B.13) \quad \lim_{z \rightarrow +\infty} (x_o(z) - 2z) = 0.$$

Indeed, we show that the inequality $x_o(z) > 2z - q$ holds for large z and $q > 0$. This follows from $f(2z - q, z) = \sinh(z - q) - \sinh z + 2z - q \sim e^z (e^{-q} - 1) / 2 \rightarrow -\infty$ for $z \rightarrow \infty$.

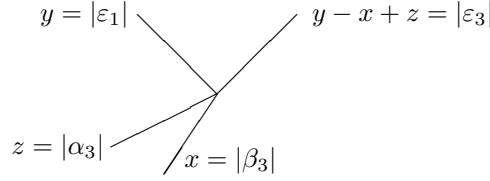


FIG. B.1. Interactions.

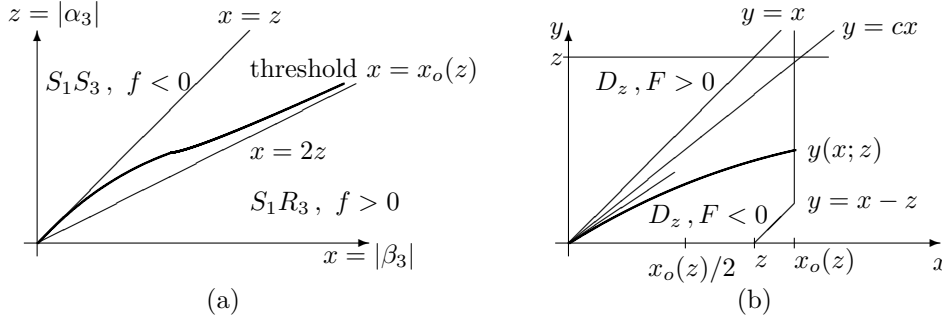


FIG. B.2. (a) the threshold curve $f(x, z) = \sinh(x - z) - \sinh(z) + x = 0$; (b) the domain D_z and the function $y = y(x; z)$.

Remark that the fact that ε_3 is a shock implies that

$$(B.14) \quad f(x, z) = \sinh(x - z) - \sinh(z) + x < 0.$$

Step 3: the amount of shocks can increase. From (B.5) we have

$$(B.15) \quad |\varepsilon_1| + |\varepsilon_3| - |\alpha_3| = 2|\varepsilon_1| - |\beta_3|.$$

We now prove that the inequality $|\varepsilon_1| + |\varepsilon_3| - |\alpha_3| < 0$, or equivalently $|\varepsilon_1| < \frac{1}{2}|\beta_3|$, does not hold if m is large.

The equation giving $|\varepsilon_1| = y$ in terms of $|\beta_3| = x$, for a given parameter $|\alpha_3| = z$, is (B.9), to be considered in the domain

$$D_z = \{(x, y): 0 \leq x \leq x_o(z), y \geq \max\{0, x - z\}\} \quad \text{for } z \geq 0,$$

where $x_o(z)$ satisfies (B.11); see Figure B.2(b). Since $F_y = \cosh y + \cosh(y - x + z) > 0$, the implicit equation (B.9) defines a function $y = y(x; z)$ with $0 \leq y(x; z) \leq x$ and $y(x; z) \leq z$; see (B.10). Remark that $F(x, x; z) = \sinh x + x > 0$. Moreover, $F(x, 0; z) = \sinh(z - x) - \sinh z + x$ so that $F(0, 0; z) = 0$; by $F_x = 1 - \cosh(z - x) < 0$, we deduce that $F(x, 0; z) < 0$ if $x > 0$. By the implicit function theorem we have $y'(x; z) \geq 0$,

$$(B.16) \quad y'(0; z) = \frac{\cosh(z) - 1}{\cosh(z) + 1} \in [0, 1),$$

and $y''(0; z) = -\frac{4 \sinh(z)}{(1 + \cosh(z))^2} \leq 0$. The function $\frac{\cosh(z) - 1}{\cosh(z) + 1}$ is increasing, and then for $z \in [0, m]$ its maximum is $\frac{\cosh(m) - 1}{\cosh(m) + 1}$; this quantity is strictly larger than $1/2$ if $m > \log(3 + 2\sqrt{2})$. Thus in general the estimate $y(x; z) < x/2$ cannot hold.

Step 4: proof of the estimate. From (B.15) we see that $|\varepsilon_1| + |\varepsilon_3| - |\alpha_3| \leq (2c - 1) \cdot |\beta_3| \iff |\varepsilon_1| \leq c \cdot |\beta_3|$, that is, that (B.3) and (B.4) are equivalent; we shall prove (B.4). To bypass the study of the function $y(x; z)$ we define

$$\Phi(x; z, c) = F(x, cx; z) = \sinh(cx) + \sinh(z - (1 - c)x) - \sinh z + x.$$

If $1/2 \leq c < 1$, then $z > (1 - c)x$ and $\Phi(0; z, c) = 0$,

$$\Phi_x(x; z, c) = 1 + c \cosh(cx) - (1 - c) \cosh(z - (1 - c)x),$$

$$\Phi_{xx}(x; z, c) = c^2 \sinh(cx) + (1 - c)^2 \sinh(z - (1 - c)x) > 0.$$

Therefore the function $x \rightarrow \Phi_x(x; z, c)$ is increasing, and then $\Phi(x; z, c) \geq 0$ if

$$\Phi_x(0; z, c) = (1 + c) - (1 - c) \cosh(z) > 0,$$

that is, if $\cosh(z) \leq \frac{1+c}{1-c}$; this is just (B.2). Then $y(x; z) \leq cx$ for all $x \in (0, x_o(z))$, and so (B.4) is proved. \square

Remark B.2. From (B.16) we see that condition (B.2) is equivalent to the geometric condition $y'(0; z) = (\cosh z - 1)/(\cosh z + 1) < c$. Moreover, as we noticed in the above proof, condition (B.2) is equivalent to

$$(B.17) \quad \cosh(|\alpha_i|) \leq \frac{1 + c}{1 - c},$$

which, in turn, is equivalent to

$$(B.18) \quad |\alpha_i| \leq \log \frac{(1 + \sqrt{c})^2}{1 - c}.$$

From the definition of the strength, one has that $|\alpha_i| = (1/2) \log(v_{max}/v_{min})$, where $v_{max} = \max\{v_\ell, v_r\}$, $v_{min} = \min\{v_\ell, v_r\}$, v_ℓ, v_r being, respectively, the left and right values of v for the wave of size α_i . Hence (B.18) is equivalent to

$$\sqrt{\frac{v_{max}}{v_{min}}} \leq \frac{(1 + \sqrt{c})^2}{1 - c}.$$

Remark B.3. In the proof above we showed that $\Delta L_{\text{shocks}} = |\varepsilon_1| + |\varepsilon_3| - |\alpha_3| = 2|\varepsilon_1| - |\beta_3|$ may be positive, differently from Nishida's paper, where it is always decreasing. This depends on the definition of the wave strengths, which was imposed to us in order to have good estimates when dealing with interactions with the 2-waves. In any case $\Delta L = \Delta L_{\text{shocks}} + \Delta L_{\text{rarefactions}} \leq 0$.

Remark B.4. Under the notation and assumptions of Lemma B.1 we verify that

$$(B.19) \quad |\varepsilon_j| \leq c \cdot |\alpha_i|.$$

This estimate, together with (B.4), allows us to obtain (in a special case) the analogue of (5.20) with c in place of d .

The proof makes use of a numerical computation. As in that lemma we consider the case $j = 1, i = 3$. Let $z > 0$ be fixed and $0 \leq x \leq x_o(z)$. From the proof of Lemma B.1 we deduce that $y = y(x, z)$; for z fixed the function $x \rightarrow y(x, z)$ is increasing. Then define

$$Y(z) = y(x_o(z), z) = x_o(z) - z.$$

In order to prove (B.19) it is sufficient to prove that

$$Y(z) \leq cz \quad \text{if} \quad \cosh(z) \leq \frac{1 + c}{1 - c}.$$

Remark that from (B.13) we know that $Y(z)/z \rightarrow 1$ for $z \rightarrow +\infty$; however, the constraint (B.17) implies that z is bounded. The inequality $Y(z) \leq cz$ is equivalent to $x_o(z) - z \leq cz$, i.e., $x_o(z) \leq (1+c)z$. Therefore we need to prove that

$$(B.20) \quad \phi(z, c) \doteq \sinh(cz) - \sinh(z) + (1+c)z \geq 0 \quad \text{if} \quad \cosh(z) \leq \frac{1+c}{1-c}.$$

Notice that $\cosh(z) \leq \frac{1+c}{1-c}$ means $0 \leq z \leq z_c$ for $z_c = \log\left(\frac{1+c}{1-c} + \sqrt{\left(\frac{1+c}{1-c}\right)^2 - 1}\right)$. Formula (B.20) is shown to hold true by numerical computations. Remark, however, that

$$\phi'(z, c) = c \cosh(cz) - \cosh(z) + 1 + c,$$

$$\phi''(z, c) = c^2 \sinh(cz) - \sinh(z) \leq 0,$$

so $\phi'(0, c) = 2c$, $\phi(\cdot, c)$ is concave, $\lim_{z \rightarrow +\infty} \phi(z, c) = -\infty$, and $\phi(\cdot, c)$ has a single point of maximum; at last $\phi'(z_c) = c(\cosh(cz_c) - \cosh(z_c)) < 0$. This concludes the proof of (B.19).

Remark that if $c = 1/2$, then $\phi(z, c) = \sinh(z/2) - \sinh z + \frac{3}{2}z$. If z is such that $\cosh z = 3$, then $\sinh z = 2\sqrt{2}$, $\sinh(z/2) = 1$, $z = \log(3 + \sqrt{8})$, and (B.20) holds with strict inequality.

Acknowledgment. The authors thank Graziano Guerra for stimulating remarks and Umberto Massari for hints on BV functions.

REFERENCES

- [1] R. ABEYARATNE AND J. K. KNOWLES, *Kinetic relations and the propagation of phase boundaries in solids*, Arch. Rational Mech. Anal., 114 (1991), pp. 119–154.
- [2] D. AMADORI AND A. CORLI, *A hyperbolic model of multiphase flow*, in Hyperbolic Problems: Theory, Numerics, Applications, S. Benzoni-Gavage and D. Serre, eds., Springer, Berlin, Heidelberg, 2008, pp. 407–414.
- [3] D. AMADORI AND G. GUERRA, *Global BV solutions and relaxation limit for a system of conservation laws*, Proc. Roy. Soc. Edinburgh Sect. A, 131 (2001), pp. 1–26.
- [4] F. ASAKURA, *Wave-front tracking for the equations of isentropic gas dynamics*, Quart. Appl. Math., 63 (2005), pp. 20–33.
- [5] F. ASAKURA, *Wave-front tracking for the equations of non-isentropic gas dynamics*, to appear.
- [6] S. BENZONI-GAVAGE, *Analyse numérique des modèles hydrodynamiques d'écoulements diphasiques instationnaires dans les réseaux de production pétrolière*, Ph.D. thesis, Ecole Normale Supérieure de Lyon, Lyon, France, 1991.
- [7] F. BEREUX, E. BONNETIER, AND P. G. LEFLOCH, *Gas dynamics system: Two special cases*, SIAM J. Math. Anal., 28 (1997), pp. 499–515.
- [8] S. BIANCHINI, *The semigroup generated by a Temple class system with non-convex flux function*, Differential Integral Equations, 13 (2000), pp. 1529–1550.
- [9] A. BRESSAN, *Hyperbolic Systems of Conservation Laws. The One-Dimensional Cauchy Problem*, Oxford University Press, Oxford, UK, 2000.
- [10] A. CORLI AND H. FAN, *The Riemann problem for reversible reactive flows with metastability*, SIAM J. Appl. Math., 65 (2004), pp. 426–457.
- [11] C. M. DAFERMOS, *Hyperbolic Conservation Laws in Continuum Physics*, 2nd ed., Grundlehren Math. Wiss., 325, Springer-Verlag, Berlin, 2005.
- [12] R. J. DIPIERNA, *Existence in the large for quasilinear hyperbolic conservation laws*, Arch. Rational Mech. Anal., 52 (1973), pp. 244–257.
- [13] H. FAN, *On a model of the dynamics of liquid/vapor phase transitions*, SIAM J. Appl. Math., 60 (2000), pp. 1270–1301.
- [14] J. GLIMM, *Solutions in the large for nonlinear hyperbolic systems of equations*, Comm. Pure Appl. Math., 18 (1965), pp. 697–715.

- [15] L. GOSSE, *Existence of L^∞ entropy solutions for a reacting Euler system*, Port. Math. (N.S.), 58 (2001), pp. 473–484.
- [16] T.-P. LIU, *Initial-boundary value problems for gas dynamics*, Arch. Rational Mech. Anal., 64 (1977), pp. 137–168.
- [17] T.-P. LIU, *Solutions in the large for the equations of nonisentropic gas dynamics*, Indiana Univ. Math. J., 26 (1977), pp. 147–177.
- [18] Y. LU, *Hyperbolic Conservation Laws and the Compensated Compactness Method*, Chapman Hall/CRC Monogr. Surv. Pure Appl. Math., 128, Chapman and Hall/CRC, Boca Raton, FL, 2003.
- [19] T. NISHIDA, *Global solution for an initial boundary value problem of a quasilinear hyperbolic system*, Proc. Japan Acad., 44 (1968), pp. 642–646.
- [20] T. NISHIDA AND J. A. SMOLLER, *Solutions in the large for some nonlinear hyperbolic conservation laws*, Comm. Pure Appl. Math., 26 (1973), pp. 183–200.
- [21] Y.-J. PENG, *Solutions faibles globales pour l'équation d'Euler d'un fluide compressible avec de grandes données initiales*, Comm. Partial Differential Equations, 17 (1992), pp. 161–187.
- [22] Y.-J. PENG, *Solutions faibles globales pour un modèle d'écoulements diphasiques*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 21 (1994), pp. 523–540.
- [23] F. POUPAUD, M. RASCLE, AND J.-P. VILA, *Global solutions to the isothermal Euler-Poisson system with arbitrarily large data*, J. Differential Equations, 123 (1995), pp. 93–121.
- [24] D. SERRE, *Systems of Conservation Laws. Vol. 1*, Cambridge University Press, Cambridge, UK, 1999.
- [25] D. SERRE, *Systems of Conservation Laws. Vol. 2*, Cambridge University Press, Cambridge, UK, 2000.
- [26] J. SMOLLER, *Shock Waves and Reaction-Diffusion Equations*, 2nd ed., Springer-Verlag, New York, 1994.
- [27] J. B. TEMPLE, *Solutions in the large for the nonlinear hyperbolic conservation laws of gas dynamics*, J. Differential Equations, 41 (1981), pp. 96–161.
- [28] D. H. WAGNER, *Equivalence of the Euler and Lagrangian equations of gas dynamics for weak solutions*, J. Differential Equations, 68 (1987), pp. 118–136.

ON THE GLOBAL WELL-POSEDNESS OF THE CRITICAL QUASI-GEOSTROPHIC EQUATION*

HAMMADI ABIDI[†] AND TAOUFIK HMIDI[†]

Abstract. We prove the global well-posedness of the critical dissipative quasi-geostrophic equation for large initial data belonging to the critical Besov space $\dot{B}_{\infty,1}^0(\mathbb{R}^2)$.

Key words. quasi-geostrophic equation, global existence, Besov spaces

AMS subject classifications. 76D03, 35B33, 35Q35, 76D05

DOI. 10.1137/070682319

1. Introduction. In this paper we are concerned with the initial value problem of the 2D (two-dimensional) dissipative quasi-geostrophic equation

$$(\text{QG})_{\alpha} \begin{cases} \partial_t \theta + v \cdot \nabla \theta + |\text{D}|^{2\alpha} \theta = 0, \\ \theta|_{t=0} = \theta^0, \end{cases}$$

where the scalar function θ represents the potential temperature and the parameter $\alpha \in]0, 1]$. The velocity $v = (v^1, v^2)$ is determined by Riesz transforms of θ ,

$$v = (-\partial_2 |\text{D}|^{-1} \theta, \partial_1 |\text{D}|^{-1} \theta) := (-R_2 \theta, R_1 \theta), \quad |\text{D}| = \sqrt{-\Delta}.$$

In addition to its intrinsic mathematical importance, this equation serves as a 2D model in geophysical fluid dynamics; for more details about the subject, see [6, 19].

This equation has been intensively investigated, and much attention is devoted to the problem of global well-posedness. For the subcritical case ($\alpha > \frac{1}{2}$) the theory seems to be in a satisfactory state. Indeed, global existence and uniqueness for arbitrary initial data are established in various function spaces (see, for example, [8, 20]). However, the critical and supercritical cases, corresponding respectively to $\alpha = \frac{1}{2}$ and $\alpha < \frac{1}{2}$, are harder to deal with. In the supercritical case, we have until now only global results for small initial data; see, for instance, [3, 5, 12, 13, 23, 24]. For the critical case, Constantin, Córdoba, and Wu showed in [7] the global existence in Sobolev space H^1 under a smallness assumption of the L^∞ norm of θ^0 . Many other relevant results can be found in [9, 14, 15, 17]. Very recently, Kiselev, Nazarov, and Volberg proved in [16] global well-posedness for arbitrary periodic smooth initial data by using an elegant argument of the modulus of continuity. In [2], Caffarelli and Vasseur established the global regularity of weak solutions associated with L^2 initial data.

The main goal of this work is to establish global well-posedness in the critical case when initial data belong to the homogeneous critical Besov space $\dot{B}_{\infty,1}^0(\mathbb{R}^2)$: we remove the periodic condition, and we weaken the initial regularity. Before giving our main result let us first specify our notion of critical spaces. Let θ be a solution of

*Received by the editors February 9, 2007; accepted for publication (in revised form) December 3, 2007; published electronically April 4, 2008.

<http://www.siam.org/journals/sima/40-1/68231.html>

[†]IRMAR, Université de Rennes 1, Campus de Beaulieu, 35042 Rennes cedex, France (hamadi.abidi@univ-rennes1.fr, thmidi@univ-rennes1.fr).

$(\text{QG})_\alpha$ and $\lambda > 0$; then $\theta_\lambda(t, x) = \theta(\lambda t, \lambda x)$ is also a solution. One class of scaling invariant spaces is the homogeneous Besov spaces $(\dot{B}_{p,r}^{2/p})$, with $p, r \in [1, \infty]$.

Our first main result reads as follows.

THEOREM 1.1. *Let $\theta^0 \in \dot{B}_{\infty,1}^0$; then there exists a unique global solution θ to $(\text{QG})_\alpha$ such that*

$$\theta \in \mathcal{C}(\mathbb{R}_+; \dot{B}_{\infty,1}^0) \cap L^1_{\text{loc}}(\mathbb{R}_+; \dot{B}_{\infty,1}^1).$$

The proof relies essentially on two facts: the first is the establishment of local existence, which is the major part of this paper, and the derivation of some smoothing effects of the solution, described in the next theorem. The second is the use of the modulus of continuity, as used in [16]. We mention that the property allowing us to remove the periodicity is the spatial decay of the solution. The key to our local existence result is some new estimates for the following transport-diffusion equation:

$$(\text{TD}) \begin{cases} \partial_t \theta + v \cdot \nabla \theta + |\text{D}| \theta = f, \\ \theta|_{t=0} = \theta^0, \end{cases}$$

where v and f are given and θ is the unknown scalar function. We now state our second main result.

THEOREM 1.2. *Let $s \in]-1, 1[$, $r, \bar{r} \in [1, +\infty]$ with $r \geq \bar{r}$, $f \in \tilde{L}^{\bar{r}}_{\text{loc}}(\mathbb{R}_+; \dot{B}_{\infty,1}^{s+\frac{1}{\bar{r}}-1})$, and v be a divergence-free vector field belonging to $L^1_{\text{loc}}(\mathbb{R}_+; \text{Lip}(\mathbb{R}^2))$. We consider a smooth solution θ of the transport-diffusion equation (TD); then there exists a constant C depending only on s such that for every $t \in \mathbb{R}_+$*

$$\|\theta\|_{\tilde{L}^r_t \dot{B}_{\infty,1}^{s+\frac{1}{r}}} \leq C e^{C \int_0^t \|\nabla v(\tau)\|_{L^\infty} d\tau} (\|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}^{\bar{r}}_t \dot{B}_{\infty,1}^{s+\frac{1}{\bar{r}}-1}}).$$

Additionally, if $v = \nabla^\perp |\text{D}|^{-1} \theta$, then we have for all $s \geq 1$

$$\|\theta\|_{\tilde{L}^r_t \dot{B}_{\infty,1}^{s+\frac{1}{r}}} \leq C e^{C \int_0^t (\|\nabla \theta(\tau)\|_{L^\infty} + \|\nabla v(\tau)\|_{L^\infty}) d\tau} (\|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}^{\bar{r}}_t \dot{B}_{\infty,1}^{s+\frac{1}{\bar{r}}-1}}).$$

We use for the proof a new approach based on Lagrangian coordinates combined with paradifferential calculus and a new commutator estimate. This idea has been recently used by the second author in [11, 12].

The rest of the paper is structured as follows. In section 2 we review some basic results of Littlewood–Paley theory, and we give some useful lemmas. Section 3 deals with a new commutator estimate which is needed for the proof of Theorem 1.2, done in section 4. Theorem 1.1 is proved in section 5.

2. Preliminaries. Throughout the paper, C stands for a constant which may be different in each occurrence. We shall sometimes use the notation $A \lesssim B$ instead of $A \leq CB$, and $A \approx B$ means that $A \lesssim B$ and $B \lesssim A$. We denote by $\mathcal{F}f$ the Fourier transform of f and by $[\eta]$ the whole part of η .

One starts with recalling a traditional result that will be frequently used.

LEMMA 2.1 (Bernstein). *Let $f \in L^p(\mathbb{R}^2)$, with $1 \leq p \leq \infty$ and $0 < r < R$. Then there exists a constant $C > 0$ such that $\forall k \in \mathbb{N}$ and $\lambda > 0$ we have*

$$\text{supp } \mathcal{F}f \subset B(0, \lambda r) \implies \sup_{|\beta|=k} \|\partial^\beta f\|_{L^q} \leq C^k \lambda^{k+2(\frac{1}{p}-\frac{1}{q})} \|f\|_{L^p},$$

$$\text{supp } \mathcal{F}f \subset \mathcal{C}(0, \lambda r, \lambda R) \implies C^{-k} \lambda^k \|f\|_{L^p} \leq \sup_{|\beta|=k} \|\partial^\beta f\|_{L^p} \leq C^k \lambda^k \|f\|_{L^p}.$$

These estimates hold true if we replace the derivation ∂^β by $|\mathbf{D}|^{|\beta|}$.

To define Besov spaces we need to recall the homogeneous Littlewood–Paley decomposition based on a dyadic unity partition. Let φ be a smooth function supported in the ring $\mathcal{C} := \{\xi \in \mathbb{R}^2, \frac{3}{4} \leq |\xi| \leq \frac{8}{3}\}$ and such that

$$\sum_{q \in \mathbb{Z}} \varphi(2^{-q}\xi) = 1 \quad \text{for } \xi \neq 0.$$

Now, for $u \in \mathcal{S}'$ we set

$$\forall q \in \mathbb{Z}, \quad \Delta_q u = \varphi(2^{-q}\mathbf{D})u \quad \text{and} \quad S_q u = \sum_{j \leq q-1} \Delta_j u.$$

We have the formal decomposition

$$u = \sum_{q \in \mathbb{Z}} \Delta_q u \quad \forall u \in \mathcal{S}'(\mathbb{R}^2)/\mathcal{P}[\mathbb{R}^2],$$

where $\mathcal{P}[\mathbb{R}^2]$ is the set of polynomials (see [18]). Moreover, the Littlewood–Paley decomposition satisfies the property of almost orthogonality:

$$(1) \quad \Delta_k \Delta_q u \equiv 0 \quad \text{if } |k - q| \geq 2, \quad \text{and} \quad \Delta_k(S_{q-1}u \Delta_q u) \equiv 0 \quad \text{if } |k - q| \geq 5.$$

We recall now the definition of Besov spaces. Letting $(p, m) \in [1, +\infty]^2$, $s \in \mathbb{R}$, and $u \in \mathcal{S}'$, we set

$$\|u\|_{\dot{B}_{p,m}^s} := \left(2^{qs} \|\Delta_q u\|_{L^p}\right)_{\ell^m}, \quad \dot{B}_{p,m}^s := \left\{u \in \mathcal{S} \mid \|u\|_{\dot{B}_{p,m}^s} < \infty\right\}.$$

- For $s < \frac{2}{p}$ (or $s \leq \frac{2}{p}$ if $m = 1$), we then define $\dot{B}_{p,m}^s$ as the completion of $\dot{B}_{p,m}^s$ for $\|\cdot\|_{\dot{B}_{p,m}^s}$.
- If $k \in \mathbb{N}$ and $\frac{2}{p} + k - 1 \leq s < \frac{2}{p} + k$ (or $s = \frac{2}{p} + k$ if $m = 1$), then $\dot{B}_{p,m}^s$ is defined as the subset of distributions $u \in \mathcal{S}'$ such that $\partial^\beta u \in \dot{B}_{p,m}^{s-k}$ whenever $|\beta| = k$.

Another characterization of the homogeneous Besov spaces that will be needed later is the following; see, for instance, [21]. For $s \in]0, 1[$, $p, m \in [1, \infty]$

$$(2) \quad \left(\int_{\mathbb{R}^2} \frac{\|u(\cdot - x) - u(\cdot)\|_{L^p}^m dx}{|x|^{sm}} \frac{dx}{|x|^2}\right)^{\frac{1}{m}} \approx \|u\|_{\dot{B}_{p,m}^s},$$

with the usual modification if $m = \infty$.

In our next study we require two kinds of coupled space-time Besov spaces. The first one is defined in the following manner: for $T > 0$ and $m \geq 1$, we denote by $L_T^r \dot{B}_{p,m}^s$ the set of all tempered distributions u satisfying

$$\|u\|_{L_T^r \dot{B}_{p,m}^s} := \left\| \left(2^{qs} \|\Delta_q u\|_{L^p}\right)_{\ell^m} \right\|_{L_T^r} < \infty.$$

The second mixed space is $\tilde{L}_T^r \dot{B}_{p,m}^s$, which is the set of tempered distributions u satisfying

$$\|u\|_{\tilde{L}_T^r \dot{B}_{p,m}^s} := \left(2^{qs} \|\Delta_q u\|_{L_T^r L^p}\right)_{\ell^m} < \infty.$$

We can define in the same way the spaces $L_T^r B_{p,m}^s$ and $\tilde{L}_T^r B_{p,m}^s$. The following embeddings are a direct consequence of Minkowski's inequality.

Let $s \in \mathbb{R}$, $r \geq 1$, and $(p, m) \in [1, \infty]^2$; then we have

$$(3) \quad \begin{aligned} L_T^r \dot{B}_{p,m}^s &\hookrightarrow \tilde{L}_T^r \dot{B}_{p,m}^s && \text{if } m \geq r \quad \text{and} \\ \tilde{L}_T^r \dot{B}_{p,m}^s &\hookrightarrow L_T^r \dot{B}_{p,m}^s && \text{if } r \geq m. \end{aligned}$$

The next lemma will be useful.

PROPOSITION 2.2. *The following results hold true:*

•

$$\dot{B}_{p,m}^s \hookrightarrow \dot{B}_{p_1,m_1}^{s-2(\frac{1}{p}-\frac{1}{p_1})} \quad \text{for } p \leq p_1 \quad \text{and } m \leq m_1.$$

- Let $|\mathbf{D}| := \sqrt{-\Delta}$ and $\sigma \in \mathbb{R}$; then the operator $|\mathbf{D}|^\sigma$ is an isomorphism from $\dot{B}_{p,m}^s$ to $\dot{B}_{p,m}^{s-\sigma}$.
- Let $\gamma \in]0, 1[$, $s_1, s_2 \in \mathbb{R}$ such that $s_1 < s_2$ and $u \in \dot{B}_{p,\infty}^{s_1} \cap \dot{B}_{p,\infty}^{s_2}$; then

$$\|u\|_{\dot{B}_{p,1}^{\gamma s_1 + (1-\gamma)s_2}} \lesssim \|u\|_{\dot{B}_{p,\infty}^{\gamma s_1}} \|u\|_{\dot{B}_{p,\infty}^{1-\gamma}}.$$

- For $s > 0$, $\dot{B}_{p,m}^s \cap L^\infty$ is an algebra.

We now recall some commutator estimates (see [4, 10] and the references therein).

LEMMA 2.3. *Let $p, r \in [1, \infty]$, $1 = \frac{1}{r} + \frac{1}{r'}$, $\rho_1 < 1$, $\rho_2 < 1$, and v be a divergence-free vector field of \mathbb{R}^2 . Assume in addition that*

$$\rho_1 + \rho_2 + 2 \min\{1, 2/p\} > 0 \quad \text{and} \quad \rho_1 + 2/p > 0.$$

Then we have

$$\sum_{q \in \mathbb{Z}} 2^{q(\frac{2}{p} + \rho_1 + \rho_2 - 1)} \|[\Delta_q, v \cdot \nabla]u\|_{L_t^1 L^p} \lesssim \|v\|_{\tilde{L}_t^r \dot{B}_{p,1}^{\frac{2}{p} + \rho_1}} \|u\|_{\tilde{L}_t^{r'} \dot{B}_{p,1}^{\frac{2}{p} + \rho_2}}.$$

Moreover, we have for $s \in]-1, 1[$

$$\sum_{q \in \mathbb{Z}} 2^{qs} \|[\Delta_q, v \cdot \nabla]u\|_{L^p} \lesssim \|\nabla v\|_{L^\infty} \|u\|_{\dot{B}_{p,1}^s}.$$

If $v = \nabla^\perp |\mathbf{D}|^{-1} \theta$, then the above estimate holds true for $s \geq 1$ if we replace $\|\nabla v\|_{L^\infty}$ by $\|\nabla v\|_{L^\infty} + \|\nabla u\|_{L^\infty}$.

The following result is due to Vishik [22].

LEMMA 2.4. *Let f be a function in the Schwartz class and ψ a diffeomorphism preserving Lebesgue measure; then $\forall p \in [1, +\infty]$ and $\forall j, q \in \mathbb{Z}$,*

$$\|\Delta_j (\Delta_q f \circ \psi)\|_{L^p} \leq C 2^{-|j-q|} \|\nabla \psi^{\epsilon(j,q)}\|_{L^\infty} \|\Delta_q f\|_{L^p},$$

with

$$\epsilon(j, q) = \text{sign}(j - q).$$

The following result is proved in [9].

PROPOSITION 2.5. *Let v be a smooth divergence-free vector field and f be a smooth function. We assume that θ is a smooth solution of the equation*

$$\partial_t \theta + v \cdot \nabla \theta + \kappa |\mathbf{D}|^{2\alpha} \theta = f, \quad \text{with } \kappa \geq 0 \quad \text{and } \alpha \in [0, 1].$$

Then for $p \in [1, +\infty]$ we have

$$\|\theta(t)\|_{L^p} \leq \|\theta(0)\|_{L^p} + \int_0^t \|f(\tau)\|_{L^p} d\tau.$$

We can find a proof of the next proposition in [12].

PROPOSITION 2.6. *Let \mathcal{C} be a ring and $\alpha \in \mathbb{R}_+$. There exists a positive constant C such that, for any $p \in [1, +\infty]$, for any pair (t, λ) of positive real numbers, we have*

$$\text{supp } \mathcal{F}u \subset \lambda \mathcal{C} \Rightarrow \|e^{-t|\mathbf{D}|^\alpha} u\|_{L^p} \leq C e^{-C^{-1}t\lambda^\alpha} \|u\|_{L^p}.$$

3. Commutator estimate. The main result of this section is the following estimate that will play a crucial role for the proof of Theorem 1.2.

PROPOSITION 3.1. *Let v be a divergence-free vector field belonging to $L^1_{\text{loc}}(\mathbb{R}_+; \text{Lip}(\mathbb{R}^2))$. For $q \in \mathbb{Z}$ we denote by ψ_q the flow of the regularized vector field $S_{q-1}v$. Then for $f \in \dot{B}^1_{\infty, \infty}$ and for $q \in \mathbb{Z}$ we have*

$$\| |\mathbf{D}|(\Delta_q f \circ \psi_q) - (|\mathbf{D}|\Delta_q f) \circ \psi_q \|_{L^\infty} \leq C e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \|\Delta_q f\|_{L^\infty},$$

where $V(t) = \|\nabla v\|_{L^1_t L^\infty(\mathbb{R}^2)}$ and C is an absolute constant.

Proof. We set $f_q := \Delta_q f$, and then it is obvious that

$$\begin{aligned} |\mathbf{D}|(f_q \circ \psi_q) - (|\mathbf{D}|f_q) \circ \psi_q &= |\mathbf{D}|^{\frac{1}{2}} \{ (|\mathbf{D}|^{\frac{1}{2}} f_q) \circ \psi_q \} - \{ |\mathbf{D}|^{\frac{1}{2}} (|\mathbf{D}|^{\frac{1}{2}} f_q) \} \circ \psi_q \\ &\quad + |\mathbf{D}|^{\frac{1}{2}} \{ |\mathbf{D}|^{\frac{1}{2}} (f_q \circ \psi_q) - (|\mathbf{D}|^{\frac{1}{2}} f_q) \circ \psi_q \} \\ &:= \text{I} + \text{II}. \end{aligned}$$

For the first term we apply Proposition 3.1 from [12], with $\alpha = \frac{1}{2}$ and $F_q = |\mathbf{D}|^{\frac{1}{2}} f_q$, yielding

$$\begin{aligned} \|\text{I}\|_{L^\infty} &\leq C e^{CV(t)} (e^{CV(t)} - 1) \|F_q\|_{\dot{B}^{\frac{1}{2}}_{\infty, 1}} \\ &\leq C e^{CV(t)} (e^{CV(t)} - 1) 2^q \|f_q\|_{L^\infty} \\ &\lesssim e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \|f_q\|_{L^\infty}. \end{aligned}$$

For the second term we use the following formula for the fractional Laplacian:

$$|\mathbf{D}|^{\frac{1}{2}} f(x) = C \int_{\mathbb{R}^2} \frac{f(x) - f(y)}{|x - y|^{\frac{5}{2}}} dy.$$

Since the flow ψ_q preserves Lebesgue measure, we easily get

$$\begin{aligned} |\mathbf{D}|^{\frac{1}{2}} (f_q \circ \psi_q)(x) - (|\mathbf{D}|^{\frac{1}{2}} f_q) \circ \psi_q(x) &= C \int_{\mathbb{R}^2} \frac{f_q(\psi_q(x)) - f_q(\psi_q(y))}{|x - y|^{\frac{5}{2}}} \\ &\quad \times \left(1 - \frac{|x - y|^{\frac{5}{2}}}{|\psi_q(x) - \psi_q(y)|^{\frac{5}{2}}} \right) dy. \end{aligned}$$

We denote $g_q(x) = f_q(\psi_q(x))$ and we set $h = x - y$:

$$|\mathbf{D}|^{\frac{1}{2}}(f_q \circ \psi_q)(x) - (|\mathbf{D}|^{\frac{1}{2}}f_q) \circ \psi_q(x) = C \int_{\mathbb{R}^2} \frac{g_q(x) - g_q(x-h)}{|h|^{\frac{5}{2}}} \bar{\psi}_q(x, h) dh,$$

with

$$\bar{\psi}_q(x, h) = 1 - \frac{|h|^{\frac{5}{2}}}{|\psi_q(x) - \psi_q(x-h)|^{\frac{5}{2}}}.$$

It follows from law products and the embedding $\dot{B}_{\infty,1}^0 \hookrightarrow L^\infty$ that

$$\begin{aligned} \|\mathbf{II}\|_{L^\infty} &\leq \| |\mathbf{D}|^{\frac{1}{2}}(f_q \circ \psi_q) - (|\mathbf{D}|^{\frac{1}{2}}f_q) \circ \psi_q \|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \\ &\leq C \|\bar{\psi}_q\|_{L^\infty(\mathbb{R}^4)} \int_{\mathbb{R}^2} |h|^{-\frac{5}{2}} \|g_q(\cdot) - g_q(\cdot-h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} dh \\ &\quad + C \sup_{h \in \mathbb{R}^2} \|\bar{\psi}_q(\cdot, h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \int_{\mathbb{R}^2} |h|^{-\frac{5}{2}} \|g_q(\cdot) - g_q(\cdot-h)\|_{L^\infty} dh \\ &= J_q^1 + J_q^2. \end{aligned}$$

We intend to estimate J_q^1 . It is plain from the mean value theorem that

$$\frac{1}{\|\nabla \psi\|_{L^\infty}^{\frac{5}{2}}} \leq \frac{|h|^{\frac{5}{2}}}{|\psi(x) - \psi(x-h)|^{\frac{5}{2}}} \leq \|\nabla \psi^{-1}\|_{L^\infty}^{\frac{5}{2}},$$

which easily gives the inequality

$$\|\bar{\psi}_q\|_{L^\infty(\mathbb{R}^4)} \leq \max\left(|1 - \|\nabla \psi_q^{-1}\|_{L^\infty}^{\frac{5}{2}}|; |1 - \|\nabla \psi_q\|_{L^\infty}^{\frac{5}{2}}|\right).$$

On the other hand, we have the classical estimates

$$(4) \quad \begin{aligned} e^{-C\|S_{q-1}\nabla v\|_{L_t^1 L^\infty}} &\leq \|\nabla \psi_q^{\mp 1}\|_{L^\infty} \leq e^{C\|S_{q-1}\nabla v\|_{L_t^1 L^\infty}} \\ &\text{and } \|S_{q-1}\nabla v\|_{L_t^1 L^\infty} \leq CV(t). \end{aligned}$$

We thus get

$$(5) \quad \|\bar{\psi}_q\|_{L^\infty(\mathbb{R}^4)} \leq C e^{CV(t)} (e^{CV(t)} - 1).$$

Using the definition of Besov spaces and the commutation of Δ_j with translation operators, one finds

$$\begin{aligned} &\int_{\mathbb{R}^2} |h|^{-\frac{5}{2}} \|g_q(\cdot) - g_q(\cdot-h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} dh \\ &\leq \sum_j 2^{\frac{1}{2}j} \int_{\mathbb{R}^2} |h|^{-\frac{1}{2}} \|\Delta_j g_q(\cdot) - (\Delta_j g_q)(\cdot-h)\|_{L^\infty} \frac{dh}{|h|^2}. \end{aligned}$$

Applying the characterization of Besov spaces (2) yields

$$\begin{aligned} \int_{\mathbb{R}^2} |h|^{-\frac{5}{2}} \|g_q(\cdot) - g_q(\cdot - h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} dh &\leq C \sum_j 2^{\frac{1}{2}j} \|\Delta_j g_q\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \\ &\leq C \sum_{|j-k|\leq 1} 2^{j\frac{1}{2}} 2^{\frac{1}{2}k} \|\Delta_j \Delta_k g_q\|_{L^\infty} \\ &\leq C \|g_q\|_{\dot{B}_{\infty,1}^1}. \end{aligned}$$

Now we use the following interpolation estimate:

$$\begin{aligned} \|g_q\|_{\dot{B}_{\infty,1}^1} &\lesssim \|g_q\|_{L^\infty}^{\frac{1}{2}} \|\Delta g_q\|_{L^\infty}^{\frac{1}{2}} \\ &\lesssim \|f_q\|_{L^\infty}^{\frac{1}{2}} \|\Delta g_q\|_{L^\infty}^{\frac{1}{2}}. \end{aligned}$$

It is easy to check from Leibniz rule that

$$\Delta g_q = \Delta(f_q \circ \psi_q) = \sum_{i=1}^d \langle (\nabla^2 f_q) \circ \psi_q \cdot \partial_i \psi_q, \partial_i \psi_q \rangle + (\nabla f_q) \circ \psi_q \cdot \Delta \psi_q.$$

Applying the Bernstein inequality, we get

$$\|\Delta g_q\|_{L^\infty} \lesssim e^{CV(t)} 2^{2q} \|f_q\|_{L^\infty} + 2^q \|f_q\|_{L^\infty} \|\Delta \psi_q\|_{L^\infty}.$$

The derivative of the flow equation with respect to x and the use of the Gronwall and Bernstein inequalities gives

$$\begin{aligned} \|\nabla^2 \psi_q(t)\|_{L^\infty} &\lesssim e^{CV(t)} \int_0^t \|\nabla^2 S_{q-1} v(\tau)\|_{L^\infty} d\tau \\ (6) \quad &\lesssim e^{CV(t)} 2^q. \end{aligned}$$

Combining both last estimates, we obtain

$$(7) \quad \|\Delta g_q\|_{L^\infty} \lesssim e^{CV(t)} 2^{2q} \|f_q\|_{L^\infty}.$$

Putting together (5) and (7), we conclude that

$$\|J_q^1(t)\|_{L^\infty} \lesssim e^{CV(t)} (e^{CV(t)} - 1) 2^q \|f_q\|_{L^\infty}.$$

Let us now turn to the second term J_q^2 . The integral term can be estimated from (2) as follows:

$$\int_{\mathbb{R}^2} |h|^{-\frac{5}{2}} \|g_q(\cdot) - g_q(\cdot - h)\|_{L^\infty} dh \lesssim \|g_q\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}}.$$

According to the classical composition result we write

$$\begin{aligned} \|g_q(t)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} &\lesssim \|\nabla \psi_q\|_{L^\infty}^{\frac{1}{2}} \|f_q\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \\ (8) \quad &\lesssim e^{CV(t)} 2^{q\frac{1}{2}} \|f_q\|_{L^\infty}. \end{aligned}$$

In order to estimate $\bar{\psi}_q$ we use the interpolation inequality

$$\|\bar{\psi}_q(\cdot, h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \lesssim \|\bar{\psi}_q(\cdot, h)\|_{L^\infty}^{\frac{1}{2}} \|\nabla_x \bar{\psi}_q(\cdot, h)\|_{L^\infty}^{\frac{1}{2}}.$$

This leads in view of (5) to

$$(9) \quad \|\bar{\psi}_q(\cdot, h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \leq C e^{CV(t)} (e^{CV(t)} - 1)^{\frac{1}{2}} \|\nabla_x \bar{\psi}_q(\cdot, h)\|_{L^\infty}^{\frac{1}{2}}.$$

The derivative of $\bar{\psi}_q$ with respect to x yields

$$\begin{aligned} |\nabla_x \bar{\psi}_q(x, h)| &\lesssim \frac{|h|^{\frac{7}{2}}}{|\psi_q(x) - \psi_q(x-h)|^{\frac{7}{2}}} \frac{|\nabla_x \psi_q(x) - \nabla_x \psi_q(x-h)|}{|h|} \\ &\lesssim \|\nabla \psi_q^{-1}\|_{L^\infty}^{\frac{7}{2}} \|\nabla^2 \psi_q\|_{L^\infty}. \end{aligned}$$

Combining (4) and (6), we obtain

$$(10) \quad \|\nabla_x \bar{\psi}_q(t)\|_{L^\infty(\mathbb{R}^4)} \lesssim e^{CV(t)} 2^q.$$

Plugging (10) into (9), we find

$$(11) \quad \|\bar{\psi}_q(\cdot, h)\|_{\dot{B}_{\infty,1}^{\frac{1}{2}}} \lesssim e^{CV(t)} V^{\frac{1}{2}}(t) 2^{\frac{q}{2}}.$$

Thus we deduce from (11) and (8) that

$$\|J_q^2(t)\|_{L^\infty} \leq C e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \|f_q(t)\|_{L^\infty}.$$

This achieves the proof. \square

4. Proof of Theorem 1.2. The Fourier localized function $\theta_q := \Delta_q \theta$ satisfies

$$(12) \quad \partial_t \theta_q + S_{q-1} v \cdot \nabla \theta_q + |\mathrm{D}| \theta_q = -[\Delta_q, v \cdot \nabla] \theta + (S_{q-1} v - v) \cdot \nabla \theta_q + f_q := \mathcal{R}_q.$$

Let ψ_q denote the flow of the velocity $S_{q-1} v$, and set

$$\bar{\theta}_q(t, x) = \theta_q(t, \psi_q(t, x)) \quad \text{and} \quad \bar{\mathcal{R}}_q(t, x) = \mathcal{R}_q(t, \psi_q(t, x)).$$

Since ψ_q is an homeomorphism, then

$$(13) \quad \|\bar{\mathcal{R}}_q\|_{L^\infty} \leq \|[\Delta_q, v \cdot \nabla] \theta\|_{L^\infty} + \|(S_{q-1} v - v) \cdot \nabla \theta_q\|_{L^\infty} + \|f_q\|_{L^\infty}.$$

It is not hard to check that the function $\bar{\theta}_q$ satisfies

$$(14) \quad \partial_t \bar{\theta}_q + |\mathrm{D}| \bar{\theta}_q = |\mathrm{D}|(\theta_q \circ \psi_q) - (|\mathrm{D}| \theta_q) \circ \psi_q + \bar{\mathcal{R}}_q := \bar{\mathcal{R}}_q^1.$$

From Proposition 3.1 we find that for $q \in \mathbb{Z}$

$$(15) \quad \| |\mathrm{D}|(\theta_q \circ \psi_q) - (|\mathrm{D}| \theta_q) \circ \psi_q \|_{L^\infty} \lesssim e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \|\theta_q(t)\|_{L^\infty},$$

where $V(t) := \|\nabla v\|_{L_t^1 L^\infty}$. Putting together (13) and (15) yields

$$\begin{aligned} \|\bar{\mathcal{R}}_q^1(t)\|_{L^\infty} &\lesssim \|f_q(t)\|_{L^\infty} + \|(S_{q-1} v - v) \cdot \nabla \theta_q\|_{L^\infty} + \|[\Delta_q, v \cdot \nabla] \theta\|_{L^\infty} \\ &\quad + e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \|\theta_q(t)\|_{L^\infty}. \end{aligned}$$

Applying the operator Δ_j to (14) and using Proposition 2.6, we obtain

$$\begin{aligned}
 (16) \quad \|\Delta_j \bar{\theta}_q(t)\|_{L^\infty} &\lesssim e^{-ct2^j} \|\Delta_j \theta_q^0\|_{L^\infty} + \int_0^t e^{-c(t-\tau)2^j} \|f_q(\tau)\|_{L^\infty} d\tau \\
 &\quad + e^{CV(t)} V^{\frac{1}{2}}(t) 2^q \int_0^t e^{-c(t-\tau)2^j} \|\theta_q(\tau)\|_{L^\infty} d\tau \\
 &\quad + \int_0^t e^{-c(t-\tau)2^j} \|[\Delta_q, v \cdot \nabla] \theta(\tau)\|_{L^\infty} d\tau \\
 (17) \quad &\quad + \int_0^t e^{-c(t-\tau)2^j} \|(S_{q-1}v - v) \cdot \nabla \theta_q(\tau)\|_{L^\infty} d\tau.
 \end{aligned}$$

Integrating this estimate with respect to the time and using the Young inequality, we get

$$\begin{aligned}
 \|\Delta_j \bar{\theta}_q\|_{L_t^r L^\infty} &\lesssim 2^{-\frac{j}{r}} (1 - e^{-crt2^j})^{\frac{1}{r}} \|\Delta_j \theta_q^0\|_{L^\infty} + 2^{-j(1+\frac{1}{r}-\frac{1}{p})} \|f_q\|_{L_t^r L^\infty} \\
 &\quad + e^{CV(t)} V^{\frac{1}{2}}(t) 2^{(q-j)} \|\theta_q\|_{L_t^r L^\infty} \\
 &\quad + 2^{-\frac{j}{r}} \int_0^t \|[\Delta_q, v \cdot \nabla] \theta(\tau)\|_{L^\infty} d\tau \\
 (18) \quad &\quad + 2^{-\frac{j}{r}} \int_0^t \|(S_{q-1}v - v) \cdot \nabla \theta_q(\tau)\|_{L^\infty} d\tau.
 \end{aligned}$$

Since the flow ψ is a homeomorphism, one writes

$$\begin{aligned}
 2^{q(s+\frac{1}{r})} \|\theta_q\|_{L_t^r L^\infty} &= 2^{q(s+\frac{1}{r})} \|\bar{\theta}_q\|_{L_t^r L^\infty} \\
 &\leq 2^{q(s+\frac{1}{r})} \left(\sum_{|j-q|>N} \|\Delta_j \bar{\theta}_q\|_{L_t^r L^\infty} + \sum_{|j-q|\leq N} \|\Delta_j \bar{\theta}_q\|_{L_t^r L^\infty} \right) \\
 &:= \text{I} + \text{II}.
 \end{aligned}$$

To estimate the term I we appeal to Lemma 2.4,

$$\begin{aligned}
 \|\Delta_j \bar{\theta}_q\|_{L_t^r L^\infty} &\lesssim 2^{-|q-j|} e^C \int_0^t \|\nabla v(\tau)\|_{L^\infty} d\tau \|\theta_q\|_{L_t^r L^\infty} \\
 &\leq C 2^{-|q-j|} e^{CV(t)} \|\theta_q\|_{L_t^r L^\infty}.
 \end{aligned}$$

Therefore we get

$$(19) \quad \text{I} \leq C 2^{-N} e^{CV(t)} 2^{q(s+\frac{1}{r})} \|\theta_q\|_{L_t^r L^\infty}.$$

In order to bound the second term II we use (18):

$$\begin{aligned}
 \text{II} &\lesssim (1 - e^{-crt2^q})^{\frac{1}{r}} 2^{qs} \|\theta_q^0\|_{L^\infty} + 2^{N(\frac{1}{r}+1-\frac{1}{p})} 2^{q(s+\frac{1}{p}-1)} \|f_q\|_{L_t^r L^\infty} \\
 &\quad + 2^N e^{CV(t)} V^{\frac{1}{2}}(t) 2^{q(s+\frac{1}{r})} \|\theta_q\|_{L_t^r L^\infty} \\
 &\quad + 2^{\frac{N}{r}} 2^{qs} \int_0^t \|[\Delta_q, v \cdot \nabla] \theta(\tau)\|_{L^\infty} d\tau \\
 (20) \quad &\quad + 2^{\frac{N}{r}} 2^{qs} \int_0^t \|(S_{q-1}v - v) \cdot \nabla \theta_q(\tau)\|_{L^\infty} d\tau.
 \end{aligned}$$

Denoting $Z_q^r(t) := 2^{q(s+\frac{1}{r})}\|\theta_q\|_{L_t^r L^\infty}$, we then obtain, in view of (19) and (20),

$$\begin{aligned} Z_q^r(t) &\leq C(1 - e^{-crt2^q})^{\frac{1}{r}} 2^{qs} \|\theta_q^0\|_{L^\infty} + C2^{N(\frac{1}{r}+1-\frac{1}{r})} 2^{q(s+\frac{1}{r}-1)} \|f_q\|_{L_t^r L^\infty} \\ &\quad + C\{2^N e^{CV(t)} V^{\frac{1}{2}}(t) + 2^{-N} e^{CV(t)}\} Z_q^r(t) \\ &\quad + C2^{\frac{N}{r}} 2^{qs} \int_0^t \|[\Delta_q, v \cdot \nabla]\theta(\tau)\|_{L^\infty} d\tau \\ &\quad + C2^{\frac{N}{r}} 2^{qs} \int_0^t \|(S_{q-1}v - v) \cdot \nabla\theta_q(\tau)\|_{L^\infty} d\tau. \end{aligned}$$

It is easy to check the existence of two absolute constants N and C_0 such that

$$V(t) \leq C_0 \Rightarrow C2^{-N} e^{CV(t)} + C2^N e^{CV(t)} V^{\frac{1}{2}}(t) \leq \frac{1}{2}.$$

Thus we obtain under this condition

$$\begin{aligned} Z_q^r(t) &\lesssim (1 - e^{-crt2^q})^{\frac{1}{r}} 2^{qs} \|\theta_q^0\|_{L^\infty} + 2^{q(s+\frac{1}{r}-1)} \|f_q\|_{L_t^r L^\infty} \\ (21) \quad &\quad + 2^{qs} \int_0^t \left(\|[\Delta_q, v \cdot \nabla]\theta(\tau)\|_{L^\infty} + \|(S_{q-1}v - v) \cdot \nabla\theta_q\|_{L^\infty} \right) d\tau. \end{aligned}$$

Summing over q and using Lemma 2.3 leads, for $V(t) \leq C_0$, to

$$\begin{aligned} \|\theta\|_{\tilde{L}_t^r \dot{B}_{\infty,1}^{s+\frac{1}{r}}} &\lesssim \|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_t^r \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}} + \int_0^t \|\nabla v(\tau)\|_{L^\infty} \|\theta(\tau)\|_{\dot{B}_{\infty,1}^s} d\tau \\ (22) \quad &\lesssim \|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_t^r \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}} + C_0 \|\theta\|_{L_t^\infty \dot{B}_{\infty,1}^s}. \end{aligned}$$

Let us show how to conclude the proof in the case of $r = \infty$. If C_0 is sufficiently small, then we obtain from (22) the desired estimate:

$$(23) \quad \|\theta\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^s} \lesssim \|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}}.$$

Now for an arbitrary positive time T we take a partition $(T_i)_{i=0}^M$ of $[0, T]$ such that $\int_{T_i}^{T_{i+1}} \|\nabla v(\tau)\|_{L^\infty} d\tau \approx C_0$. We can proceed analogously to the above calculus and obtain

$$\|\theta\|_{\tilde{L}_{[T_i, T_{i+1}]}^\infty \dot{B}_{\infty,1}^s} \lesssim \|\theta(T_i)\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_{[T_i, T_{i+1}]}^\infty \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}}.$$

An iteration argument leads to

$$\|\theta\|_{\tilde{L}_{[T_i, T_{i+1}]}^\infty \dot{B}_{\infty,1}^s} \leq C^{i+1} \left(\|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_{[0, T_{i+1}]}^\infty \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}} \right).$$

The triangle inequality and the fact that $C_0 M \simeq 1 + V(t)$ give

$$(24) \quad \|\theta\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^s} \leq C e^C \int_0^T \|\nabla v(\tau)\|_{L^\infty} \left(\|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}} \right).$$

Let us now turn to the case of finite r . Combining (22) and (23), we obtain under the assumption $V(t) \leq C_0$

$$(25) \quad \|\theta\|_{\tilde{L}_T^r \dot{B}_{\infty,1}^{s+\frac{1}{r}}} \lesssim \|\theta^0\|_{\dot{B}_{\infty,1}^s} + \|f\|_{\tilde{L}_T^r \dot{B}_{\infty,1}^{s+\frac{1}{r}-1}}.$$

This gives the result for a short time, and as for the case $r = \infty$, we obtain the required global estimate.

Concerning the last estimate of Theorem 1.2, we use in the commutator term of (21) the last part of Lemma 2.3. \square

5. Proof of Theorem 1.1. The proof is divided into two parts: in the first one we construct a local unique solution and give a criterion of global existence. However, in the second part we discuss the global existence by reproducing an idea of [16].

5.1. Local existence. We aim to prove the following result.

PROPOSITION 5.1. *Given any $\theta^0 \in \dot{B}_{\infty,1}^0$, there is $T > 0$ such that the (QG) $_{\alpha}$ equation has a unique solution θ with*

$$\theta \in \tilde{L}_T^{\infty} \dot{B}_{\infty,1}^0 \cap L_T^1 \dot{B}_{\infty,1}^1.$$

Moreover, for all $\beta \in \mathbb{R}_+$ we have $t^{\beta}\theta \in \tilde{L}_T^{\infty} \dot{B}_{\infty,1}^{\beta}$.

Proof. The existence is based on Theorem 1.2 and an iterative method. We denote $\theta_0(t, x) := e^{-t|\mathbb{D}|}\theta^0(x)$, $v_0 := (-R_2\theta_0, R_1\theta_0)$, and θ_{n+1} the solution of the linear system

$$\begin{cases} \partial_t \theta_{n+1} + v_n \cdot \nabla \theta_{n+1} + |\mathbb{D}|\theta_{n+1} = 0, \\ v_n = (-R_2\theta_n, R_1\theta_n), \\ \theta_{n+1}|_{t=0} = \theta^0. \end{cases}$$

Since $\theta_0 \in L^1(\mathbb{R}_+; \dot{B}_{\infty,1}^1)$ and from the continuity of Riesz transforms in the homogeneous Besov spaces we find $v_0 \in L^1(\mathbb{R}_+; \dot{B}_{\infty,1}^1)$. Thus by iteration and thanks to Theorem 1.2, one deduces that $\forall n \in \mathbb{N}$,

$$\theta_n \in \tilde{L}^{\infty}(\mathbb{R}_+; \dot{B}_{\infty,1}^0) \cap L^1(\mathbb{R}_+; \dot{B}_{\infty,1}^1).$$

Step 1: Uniform bounds. Now we intend to obtain uniform bounds, with respect to the parameter n , for some $T > 0$ independent of n .

By (21), we have for all $T \geq 0$ such that

$$(26) \quad \int_0^T \|\theta_n(\tau)\|_{\dot{B}_{\infty,1}^1} d\tau \leq C_1 (:= CC_0)$$

the following estimate:

$$\begin{aligned} \|\theta_{n+1}\|_{\tilde{L}_T^2 \dot{B}_{\infty,1}^{\frac{1}{2}}} + \|\theta_{n+1}\|_{L_T^1 \dot{B}_{\infty,1}^1} &\lesssim \sum_{q \in \mathbb{Z}} (1 - e^{-cT2^q})^{\frac{1}{2}} \|\Delta_q \theta^0\|_{L^{\infty}} \\ &+ \sum_{q \in \mathbb{Z}} \int_0^T \|[\Delta_q, v_n \cdot \nabla] \theta_{n+1}(\tau)\|_{L^{\infty}} d\tau \\ &+ \sum_{q \in \mathbb{Z}} \int_0^T \|(S_{q-1} v_n - v_n) \cdot \nabla \Delta_q \theta_{n+1}(\tau)\|_{L^{\infty}} d\tau. \end{aligned}$$

Since $\operatorname{div} v_n = 0$, then Lemma 2.3 combined with the continuity of Riesz transforms gives

$$\begin{aligned} \sum_{q \in \mathbb{Z}} \int_0^T \|\Delta_q v_n \cdot \nabla \theta_{n+1}(\tau)\|_{L^\infty} d\tau &\lesssim \|v_n\|_{\tilde{L}_T^2 \dot{B}_{\infty, \infty}^{\frac{1}{2}}} \|\theta_{n+1}\|_{\tilde{L}_T^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}} \\ &\lesssim \|\theta_n\|_{\tilde{L}_T^2 \dot{B}_{\infty, \infty}^{\frac{1}{2}}} \|\theta_{n+1}\|_{\tilde{L}_T^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}}. \end{aligned}$$

We deduce from the Hölder and Young inequalities that

$$\begin{aligned} \sum_{q \in \mathbb{Z}} \int_0^t \|(S_q v_n - v_n) \cdot \nabla \Delta_q \theta_{n+1}\|_{L^\infty} d\tau &\lesssim \sum_{q \in \mathbb{Z}} 2^q \|\Delta_q \theta_{n+1}\|_{L_t^2 L^\infty} \|S_q v_n - v_n\|_{L_t^2 L^\infty} \\ &\lesssim \sum_{q \in \mathbb{Z}} 2^{\frac{1}{2}q} \|\Delta_q \theta_{n+1}\|_{L_t^2 L^\infty} \sum_{k \geq q} 2^{\frac{1}{2}(q-k)} 2^{\frac{1}{2}k} \|\Delta_k v_n\|_{L_t^2 L^\infty} \\ &\lesssim \|\theta_{n+1}\|_{\tilde{L}_t^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}} \|\theta_n\|_{\tilde{L}_t^2 \dot{B}_{\infty, \infty}^{\frac{1}{2}}}. \end{aligned}$$

Therefore we obtain from the above inequalities

$$\begin{aligned} \|\theta_{n+1}\|_{\tilde{L}_t^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}} + \|\theta_{n+1}\|_{L_t^1 \dot{B}_{\infty, 1}^1} &\lesssim \sum_{q \in \mathbb{Z}} (1 - e^{-ct2^q})^{\frac{1}{2}} \|\Delta_q \theta^0\|_{L^\infty} \\ &\quad + \|\theta_{n+1}\|_{\tilde{L}_t^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}} \|\theta_n\|_{\tilde{L}_t^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}}. \end{aligned}$$

Thus there exists an absolute constant $\varepsilon_0 > 0$ such that, if

$$(27) \quad \sum_{q \in \mathbb{Z}} (1 - e^{-cT2^q})^{\frac{1}{2}} \|\Delta_q \theta^0\|_{L^\infty} \leq \varepsilon_0,$$

then

$$(28) \quad \|\theta_{n+1}\|_{\tilde{L}_T^2 \dot{B}_{\infty, 1}^{\frac{1}{2}}} + \|\theta_{n+1}\|_{L_T^1 \dot{B}_{\infty, 1}^1} \leq 2\varepsilon_0.$$

The existence of $T > 0$ is due to the Lebesgue theorem.

Hence, by using the estimate

$$\int_0^T \|\nabla v(\tau)\|_{L^\infty} d\tau \lesssim \int_0^T \|\theta(\tau)\|_{\dot{B}_{\infty, 1}^1} d\tau$$

and Theorem 1.2, we obtain

$$\|\theta_{n+1}\|_{\tilde{L}_T^\infty \dot{B}_{\infty, 1}^0} \lesssim \|\theta^0\|_{\dot{B}_{\infty, 1}^0}.$$

Thus we prove that the sequence $(v_n, \theta_n)_{n \in \mathbb{N}}$ is uniformly bounded in the space $\tilde{L}_T^\infty \dot{B}_{\infty, 1}^0 \cap L_T^1 \dot{B}_{\infty, 1}^1$.

Step 2: Strong convergence. We will prove that (v_n, θ_n) is a Cauchy sequence in $\tilde{L}_T^\infty \dot{B}_{\infty, 1}^0$. Let $(n, m) \in \mathbb{N}^2$, $\theta_{n,m} := \theta_{n+1} - \theta_{m+1}$, and $v_{n,m} := v_n - v_m$; then

$$\begin{cases} \partial_t \theta_{n,m} + v_n \cdot \nabla \theta_{n,m} + |\mathbf{D}| \theta_{n,m} = -v_{n,m} \cdot \nabla \theta_{m+1}, \\ \theta_{n,m}|_{t=0} = 0. \end{cases}$$

Applying Theorem 1.2 to this equation gives

$$(29) \quad \|\theta_{n,m}\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^0} \leq C e^{C\|\theta_n\|_{L_t^1 \dot{B}_{\infty,1}^1}} \int_0^t \|v_{n,m} \cdot \nabla \theta_{m+1}(\tau)\|_{\dot{B}_{\infty,1}^0} d\tau.$$

Thanks to Bony's decomposition [1], the embedding $\dot{B}_{\infty,1}^0 \hookrightarrow L^\infty$, and the fact that $\operatorname{div} v_{n,m} = 0$,

$$(30) \quad \|v_{n,m} \cdot \nabla \theta_{m+1}\|_{\dot{B}_{\infty,1}^0} \lesssim \|v_{n,m}\|_{\dot{B}_{\infty,1}^0} \|\theta_{m+1}\|_{\dot{B}_{\infty,1}^1}.$$

Since Riesz transforms continuously map $\dot{B}_{\infty,1}^0$ into itself, we get

$$(31) \quad \|v_{n,m} \cdot \nabla \theta_{m+1}\|_{\dot{B}_{\infty,1}^0} \lesssim \|\theta_{n-1,m-1}\|_{\dot{B}_{\infty,1}^0} \|\theta_{m+1}\|_{\dot{B}_{\infty,1}^1}.$$

Thus we infer

$$\|\theta_{n,m}\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^0} \leq C \|\theta_{n-1,m-1}\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^0} e^{C\|\theta_n\|_{L_t^1 \dot{B}_{\infty,1}^1}} \int_0^t \|\theta_{m+1}(\tau)\|_{\dot{B}_{\infty,1}^1} d\tau.$$

According to the inequality (28) one can choose ε_0 small such that

$$\|\theta_{n,m}\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0} \leq \eta \|\theta_{n-1,m-1}\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0}$$

with $\eta < 1$. Let us suppose that $n \geq m$; then by induction one finds

$$\|\theta_{n,m}\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0} \lesssim \eta^m \|\theta^0\|_{\dot{B}_{\infty,1}^0}.$$

Thus $(\theta_n)_{n \in \mathbb{N}}$ is a Cauchy sequence in $\tilde{L}_T^\infty \dot{B}_{\infty,1}^0$. Then there exists $\theta \in \tilde{L}_t^\infty \dot{B}_{\infty,1}^0$ such that θ_n converges strongly to θ in $\tilde{L}_t^\infty \dot{B}_{\infty,1}^0$. Moreover, the Fatou lemma and inequality (28) imply that $\theta \in L_t^1 \dot{B}_{\infty,1}^1$. These pieces of information allow us to pass to the limit into the equation.

Step 3: Uniqueness. Let $X_T := L_T^\infty \dot{B}_{\infty,1}^0 \cap L_T^1 \dot{B}_{\infty,1}^1$, and θ_i , $i = 1, 2$ (v_i the corresponding velocity), be two solutions of the $(\text{QG})_\alpha$ equation with the same initial data and belonging to the space X_T . We set $\theta_{1,2} = \theta_1 - \theta_2$ and $v_{1,2} = v_1 - v_2$; then it is plain that

$$\partial_t \theta_{1,2} + v^1 \cdot \nabla \theta_{1,2} + |\text{D}| \theta_{1,2} = -v_{1,2} \cdot \nabla \theta^2, \quad \theta_{1,2}|_{t=0} = 0.$$

Thanks to the inequalities (29) and (31), we have

$$\|\theta_{1,2}\|_{\tilde{L}_t^\infty \dot{B}_{\infty,1}^0} \leq C e^{C\|\theta_1\|_{L_t^1 \dot{B}_{\infty,1}^1}} \int_0^t \|\theta_{1,2}\|_{\tilde{L}_\tau^\infty \dot{B}_{\infty,1}^0} \|\theta_2(\tau)\|_{\dot{B}_{\infty,1}^1} d\tau.$$

Thus Gronwall's inequality gives the desired result.

Step 4: Smoothing effect. We will show the precise estimate: $\forall \beta \in \mathbb{R}_+$ we have

$$(32) \quad \|t^\beta \theta(t)\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^\beta} \leq C_\beta e^{C(\beta+1)\|\theta\|_{L_T^1 \dot{B}_{\infty,1}^1}} \|\theta\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0}.$$

It is clear that

$$\begin{cases} \partial_t(t^\beta \theta) + v \cdot \nabla(t^\beta \theta) + |\text{D}|(t^\beta \theta) = \beta t^{\beta-1} \theta, \\ (t^\beta \theta)|_{t=0} = 0. \end{cases}$$

We will proceed by induction and start the proof with the case $\beta \in \mathbb{N}$.

For $\beta = 1$, we apply Theorem 1.2 with $\bar{r} = +\infty$,

$$\|t\theta(t)\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^1} \lesssim e^{C\|\theta\|_{L_T^1 \dot{B}_{\infty,1}^1}} \|\theta\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0}.$$

Assume that (32) holds for degree n ; we will prove it for $n+1$.

Applying Theorem 1.2 to the equation of $t^{n+1}\theta$, we get

$$\begin{aligned} \|t^{n+1}\theta(t)\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^{n+1}} &\leq C(n+1)e^{C\|\theta\|_{L_T^1 \dot{B}_{\infty,1}^1}} \|t^n\theta\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^n} \\ &\leq C_n e^{C(n+2)\|\theta\|_{L_T^1 \dot{B}_{\infty,1}^1}} \|\theta\|_{\tilde{L}_T^\infty \dot{B}_{\infty,1}^0}. \end{aligned}$$

For $\beta \in \mathbb{R}_+$, we have $[\beta] \leq \beta < [\beta] + 1$, and by interpolation, one has

$$\|t^\beta\theta\|_{\tilde{L}_T^\infty (\dot{B}_{\infty,1}^\beta)} \lesssim \left\| t^{[\beta]}\theta \right\|_{\tilde{L}_T^\infty (\dot{B}_{\infty,1}^{[\beta]})}^{1+[\beta]-\beta} \left\| t^{[\beta]+1}\theta \right\|_{\tilde{L}_T^\infty (\dot{B}_{\infty,1}^{[\beta]+1})}^{\beta-[\beta]}.$$

This completes the proof. \square

5.2. Blowup criteria. The main result of this section is the following.

PROPOSITION 5.2. *Let T^* be the maximum local existence time of θ in $L_T^\infty \dot{B}_{\infty,1}^0 \cap L_T^1 \dot{B}_{\infty,1}^1$. There exists an absolute constant $\varepsilon_0 > 0$ such that if $T^* < \infty$, then*

$$\liminf_{t \rightarrow T^*} (T^* - t) \|\nabla\theta(t)\|_{L^\infty} \geq \varepsilon_0.$$

Proof. From local existence theory and especially (27) we see that if $T^* < \infty$, then necessarily

$$\liminf_{t \rightarrow T^*} \sum_{q \in \mathbb{Z}} (1 - e^{-c(T^*-t)2^q})^{\frac{1}{2}} \|\theta_q(t)\|_{L^\infty} \geq \varepsilon_0;$$

otherwise we can continue the solution over T^* . It follows that

$$\liminf_{t \rightarrow T^*} \sum_{q \in \mathbb{Z}} (1 - e^{-c(T^*-t)2^q})^{\frac{1}{2}} \sup_{t \leq T^*} \|\theta_q(t)\|_{L^\infty} \geq \varepsilon_0.$$

Consequently from the Lebesgue theorem we obtain

$$\|\theta\|_{\tilde{L}_{T^*}^\infty (\dot{B}_{\infty,1}^0)} = \infty.$$

Using the Bernstein inequality and the fact that $\|\theta_q\|_{L^\infty} \lesssim \|\theta^0\|_{L^\infty}$, we have

$$\begin{aligned} \varepsilon_0 &\leq \liminf_{t \rightarrow T^*} \left\{ \sum_{q \leq N} (1 - e^{-c(T^*-t)2^q})^{\frac{1}{2}} \|\theta_q(t)\|_{L^\infty} \right. \\ &\quad \left. + \sum_{q \geq N} (1 - e^{-c(T^*-t)2^q})^{\frac{1}{2}} \|\theta_q(t)\|_{L^\infty} \right\} \\ &\lesssim \liminf_{t \rightarrow T^*} \left\{ (T^* - t)^{\frac{1}{2}} \|\theta^0\|_{L^\infty} \sum_{q \leq N} 2^{q/2} + \|\nabla\theta(t)\|_{L^\infty} \sum_{q \geq N} 2^{-q} \right\} \\ &\lesssim \liminf_{t \rightarrow T^*} \left\{ (T^* - t) \|\theta^0\|_{L^\infty} 2^N + \|\nabla\theta(t)\|_{L^\infty} 2^{-N} \right\}. \end{aligned}$$

Choosing judiciously N , we obtain the desired result. \square

5.3. Global existence. We will use the idea of [16]. Let T^* be the maximal time existence of the solution in the space $\tilde{L}_{loc}^\infty([0, T^*[, \dot{B}_{\infty,1}^0) \cap L_{loc}^1([0, T^*[, \dot{B}_{\infty,1}^1)$. From the local existence, there exists $T_0 > 0$ such that

$$\forall t \in [0, T_0], \quad t \|\nabla \theta(t)\|_{L^\infty} \leq C \|\theta^0\|_{\dot{B}_{\infty,1}^0}.$$

Let λ be a real positive number that will be fixed later and $T_1 \in]0, T_0[$. We define the set

$$I := \left\{ T \in [T_1, T^*[, \forall t \in [T_1, T], \forall x \neq y \in \mathbb{R}^2, |\theta(t, x) - \theta(t, y)| < \omega_\lambda(|x - y|) \right\},$$

where

$$\omega : \mathbb{R}_+ \longrightarrow \mathbb{R}_+$$

is strictly nondecreasing, concave, $\omega(0) = 0$, $\omega'(0) < +\infty$, $\lim_{\xi \rightarrow 0^+} \omega''(\xi) = -\infty$, and

$$\omega_\lambda(|x - y|) = \omega(\lambda|x - y|).$$

The function ω is a modulus of continuity chosen as in [16]. We shall first check that I is nonempty. It suffices for this purpose to prove that T_1 belongs to I under suitable conditions over λ . Let C_0 be a large positive number such that

$$(33) \quad \omega(C_0) > 2\|\theta^0\|_{L^\infty}.$$

Since ω is a strictly nondecreasing function, then we get from the maximum principle that

$$\forall x, y, \quad \lambda|x - y| \geq C_0 \Rightarrow |\theta(T_1, x) - \theta(T_1, y)| \leq 2\|\theta^0\|_{L^\infty} < \omega_\lambda(|x - y|).$$

On the other hand, we have from the mean value theorem

$$|\theta(T_1, x) - \theta(T_1, y)| \leq |x - y| \|\nabla \theta(T_1)\|_{L^\infty}.$$

Let $0 < \delta_0 < C_0$. Then using the concavity of ω , one obtains

$$\lambda|x - y| \leq \delta_0 \Rightarrow \omega_\lambda(|x - y|) \geq \frac{\omega(\delta_0)}{\delta_0} \lambda|x - y|.$$

If we choose λ so that

$$\lambda > \frac{\delta_0}{\omega(\delta_0)} \|\nabla \theta(T_1)\|_{L^\infty},$$

then we get

$$0 < \lambda|x - y| \leq \delta_0 \Rightarrow |\theta(T_1, x) - \theta(T_1, y)| < \omega_\lambda(|x - y|).$$

Let us now move to the case $\delta_0 \leq \lambda|x - y| \leq C_0$. By an obvious computation we find

$$|\theta(T_1, x) - \theta(T_1, y)| \leq \frac{C_0}{\lambda} \|\nabla \theta(T_1)\|_{L^\infty} \quad \text{and}$$

$$\omega(\delta_0) \leq \omega(\lambda|x - y|).$$

Choosing λ such that

$$\lambda > \frac{C_0}{\omega(\delta_0)} \|\nabla\theta(T_1)\|_{L^\infty},$$

then we obtain

$$\delta_0 \leq \lambda|x - y| \leq C_0 \Rightarrow |\theta(T_1, x) - \theta(T_1, y)| < \omega_\lambda(|x - y|).$$

All the preceding conditions over λ can be obtained if we take

$$(34) \quad \lambda = \frac{\omega^{-1}(3\|\theta^0\|_{L^\infty})}{2\|\theta^0\|_{L^\infty}} \|\nabla\theta(T_1)\|_{L^\infty}.$$

From the construction, the set I is an interval of the form $[T_1, T_*)$. We have three possibilities. The first one is $T_* = T^*$, and in this case we must have $T^* = +\infty$ because the Lipschitz norm of θ does not blow up. The second one is $T_* \in I$, and we will show that is not possible. Indeed, let C_0 be as in (33); then $\forall t \in [T_1, T^*)$

$$\lambda|x - y| \geq C_0 \Rightarrow |\theta(t, x) - \theta(t, y)| < \omega_\lambda(|x - y|).$$

Since $\nabla\theta(t)$ belongs to $C([0, T^*]; \dot{B}_{\infty,1}^0)$, then for $\epsilon > 0$ there exist $\eta_0, R > 0$ such that

$$\forall t \in [T_*, T_* + \eta_0] \Rightarrow \|\nabla\theta(t)\|_{L^\infty} \leq \|\nabla\theta(T_*)\|_{L^\infty} + \epsilon/2 \quad \text{and}$$

$$\|\nabla\theta(T_*)\|_{L^\infty(B_{(0,R)})} \leq \epsilon/2,$$

where $B_{(0,R)}$ is the ball of radius R and with center the origin.

Hence for $\lambda|x - y| \leq C_0$ and x or $y \in B_{(0,R+\frac{C_0}{\lambda})}^c$ we have for $t \in [T_*, T_* + \eta_0]$

$$\begin{aligned} |\theta(t, x) - \theta(t, y)| &\leq |x - y| \|\nabla\theta(t)\|_{L^\infty(B_{(0,R)})} \\ &\leq \epsilon|x - y|. \end{aligned}$$

On the other hand, from the concavity of ω we have

$$\lambda|x - y| \leq C_0 \Rightarrow \frac{\omega(C_0)}{C_0} \lambda|x - y| \leq \omega_\lambda(|x - y|).$$

Thus if we take ϵ sufficiently small such that

$$\epsilon < \frac{\omega(C_0)}{C_0} \lambda,$$

then we find that

$$\lambda|x - y| \leq C_0, \quad x \text{ or } y \in B_{(0,R+\frac{C_0}{\lambda})}^c \Rightarrow |\theta(t, x) - \theta(t, y)| < \omega_\lambda(|x - y|).$$

It remains to study the case where $x, y \in B_{(0,R+\frac{C_0}{\lambda})}$. Since $\|\nabla^2\theta(T_*)\|_{L^\infty}$ is finite (see Proposition 5.1), we get for each $x \in \mathbb{R}^2$

$$|\nabla\theta(T_*, x)| < \lambda\omega'(0).$$

For the proof, see [16, p. 3]. From the continuity of $x \mapsto |\nabla\theta(T_*, x)|$ we obtain

$$\|\nabla\theta(T_*)\|_{L^\infty(B_{(0,R+\frac{c_0}{\lambda})})} < \lambda\omega'(0).$$

Let $\delta_0 \ll 1$; then using the continuity in time of the quantity $\|\nabla\theta(t)\|_{L^\infty}$, one can find $\eta_1 > 0$ such that $\forall t \in [T_*, T_* + \eta_1]$

$$\|\nabla\theta(t)\|_{L^\infty(B_{(0,R+\frac{c_0}{\lambda})})} < \lambda\frac{\omega(\delta_0)}{\delta_0}.$$

For $\lambda|x - y| \leq \delta_0$ and $x \neq y$ belonging together to $B_{(0,R+\frac{c_0}{\lambda})}$ we have

$$\begin{aligned} |\theta(t, x) - \theta(t, y)| &\leq |x - y| \|\nabla\theta(t)\|_{L^\infty(B_{(0,R+\frac{c_0}{\lambda})})} \\ &< \lambda|x - y| \frac{\omega(\delta_0)}{\delta_0} \leq \omega_\lambda(|x - y|). \end{aligned}$$

Now for the other case we have

$$\forall x, y \in B_{(0,R+\frac{c_0}{\lambda})}, \quad \delta_0 \leq \lambda|x - y|; \quad |\theta(T_*, x) - \theta(T_*, y)| < \omega_\lambda(|x - y|);$$

then we get from a standard compact argument the existence of $\eta_2 > 0$ such that $\forall t \in [T_*, T_* + \eta_2]$

$$\forall x, y \in B_{(0,R+\frac{c_0}{\lambda})}, \quad \delta_0 \leq \lambda|x - y|; \quad |\theta(t, x) - \theta(t, y)| < \omega_\lambda(|x - y|).$$

Taking $\eta = \min\{\eta_0, \eta_1, \eta_2\}$, we obtain that $T_* + \eta \in I$, which contradicts the fact that T_* is maximal.

The last case that we have to treat is that T_* does not belong to I . Thus we have by the time continuity of θ the existence of $x \neq y$ such that

$$\theta(T_*, x) - \theta(T_*, y) = \omega_\lambda(\xi), \quad \text{with } \xi = |x - y|.$$

We will show that this scenario cannot occur, and more precisely

$$f'(T_*) < 0, \quad \text{where } f(t) = \theta(t, x) - \theta(t, y).$$

This is impossible since $f(t) \leq f(T_*) \forall t \in [0, T_*]$. The proof is the same as in [16], but for the convenience of the reader we will outline the proof. From the regularity of the solution we see that $(\text{QG})_\alpha$ can be defined in the classical manner and

$$f'(T_*) = (u \cdot \nabla\theta)(T_*, x) - (u \cdot \nabla\theta)(T_*, y) + |\text{D}|\theta(T_*, x) - |\text{D}|\theta(T_*, y).$$

From [16] we have

$$(u \cdot \nabla\theta)(T_*, x) - (u \cdot \nabla\theta)(T_*, y) \leq \Omega_\lambda(\xi)\omega'_\lambda(\xi),$$

where

$$\Omega_\lambda(\xi) = C \left(\int_0^\xi \frac{\omega_\lambda(\eta)}{\eta} d\eta + \xi \int_\xi^\infty \frac{\omega_\lambda(\eta)}{\eta^2} d\eta \right) = \Omega(\lambda\xi).$$

Again from [16],

$$\begin{aligned} |\mathbf{D}|\theta(T_*, x) - \mathbf{D}|\theta(T_*, y)| &\leq \frac{1}{\pi} \int_0^{\frac{\xi}{2}} \frac{\omega_\lambda(\xi + 2\eta) + \omega_\lambda(\xi - 2\eta) - 2\omega_\lambda(\xi)}{\eta^2} d\eta \\ &\quad + \frac{1}{\pi} \int_{\frac{\xi}{2}}^\infty \frac{\omega_\lambda(2\eta + \xi) - \omega_\lambda(2\eta - \xi) - 2\omega_\lambda(\xi)}{\eta^2} d\eta \\ &\leq \lambda \mathcal{I}(\lambda\xi), \end{aligned}$$

where

$$\begin{aligned} \mathcal{I}(\xi) &= \frac{1}{\pi} \int_0^{\frac{\xi}{2}} \frac{\omega(\xi + 2\eta) + \omega(\xi - 2\eta) - 2\omega(\xi)}{\eta^2} d\eta \\ &\quad + \frac{1}{\pi} \int_{\frac{\xi}{2}}^\infty \frac{\omega(2\eta + \xi) - \omega(2\eta - \xi) - 2\omega(\xi)}{\eta^2} d\eta. \end{aligned}$$

Thus we get

$$f'(T_*) = \lambda(\Omega\omega' + \mathcal{I})(\lambda\xi).$$

Now, we choose the same function as [16] (see p. 5):

$$\omega(\xi) = \xi - \xi^{\frac{3}{2}} \quad \text{if } \xi \in [0, \delta],$$

and

$$\omega'(\xi) = \frac{\gamma}{\xi(4 + \log(\xi/\delta))} \quad \text{if } \xi > \delta,$$

where δ and γ are small numbers and satisfy $0 < \gamma < \delta$. It is shown in [16] that

$$\Omega(\xi)\omega'(\xi) + \mathcal{I}(\xi) < 0 \quad \forall \xi \neq 0.$$

This yields $f'(T_*) < 0$.

Finally we have $T^* = +\infty$ and

$$\forall t \in [T_1, +\infty), \quad \|\nabla\theta(t)\|_{L^\infty} \leq \lambda.$$

The value of λ is given by (34).

REFERENCES

- [1] J.-M. BONY, *Calcul symbolique et propagation des singularités pour les équations aux dérivées partielles non linéaires*, Ann. Ecole Norm. Sup., 14 (1981), pp. 209–246.
- [2] L. CAFFARELLI AND V. VASSEUR, *Drift Diffusion Equations with Fractional Diffusion and the Quasi-Geostrophic Equations*, preprint, <http://arXiv.org/abs/math/0608447>.
- [3] D. CHAE AND J. LEE, *Global well-posedness in the supercritical dissipative quasi-geostrophic equations*, Asymptot. Anal., 38 (2004), pp. 339–358.
- [4] J.-Y. CHEMIN, *Perfect Incompressible Fluids*, Clarendon Press, Oxford University Press, New York, 1998.
- [5] Q. CHEN, C. MIAO, AND Z. ZHANG, *A new Bernstein's inequality and the 2D dissipative quasi-geostrophic equation*, Comm. Math. Phys., 271 (2007), pp. 821–838.

- [6] P. CONSTANTIN, A. MAJDA, AND E. TABAK, *Formation of strong fronts in the 2D quasi-geostrophic thermal active scalar*, Nonlinearity, 7 (1994), pp. 1495–1533.
- [7] P. CONSTANTIN, D. CÓRDOBA, AND J. WU, *On the critical dissipative quasi-geostrophic equation*, Indiana Univ. Math. J., 50 (2001), pp. 97–107.
- [8] P. CONSTANTIN AND J. WU, *Behavior of solutions of 2D quasi-geostrophic equations*, SIAM J. Math. Anal., 30 (1999), pp. 937–948.
- [9] A. CÓRDOBA AND D. CÓRDOBA, *A maximum principle applied to quasi-geostrophic equations*, Comm. Math. Phys., 249 (2004), pp. 511–528.
- [10] R. DANCHIN, *Density-dependent incompressible viscous fluids in critical spaces*, Proc. Roy. Soc. Ed., 133 (2003), pp. 1311–1334.
- [11] T. HMIDI, *Régularité höldérienne des poches de tourbillon visqueuses*, J. Math. Pures Appl. (9), 84 (2005), pp. 1455–1495.
- [12] T. HMIDI AND S. KERAANI, *Global solutions of the super-critical 2D quasi-geostrophic equation in Besov spaces*, Adv. Math., 214 (2007), pp. 618–638.
- [13] N. JU, *Existence and uniqueness of the solution to the dissipative 2D quasi-geostrophic equations in the Sobolev space*, Comm. Math. Phys., 251 (2004), pp. 365–376.
- [14] N. JU, *On the two dimensional quasi-geostrophic equations*, Indiana Univ. Math. J., 54 (2005), pp. 897–926.
- [15] N. JU, *Global solutions to the two dimensional quasi-geostrophic equation with critical or super-critical dissipation*, Math. Ann., 334 (2006), pp. 627–642.
- [16] A. KISELEV, F. NAZAROV, AND A. VOLBERG, *Global well-posedness for the critical 2D dissipative quasi-geostrophic equation*, Invent. Math., 167 (2007), pp. 445–453.
- [17] F. MARCHAND AND P.-G. LEMARIÉ-RIEUSSET, *Solutions auto-similaires non radiales pour l'équation quasi-geostrophique dissipative critique*, C. R. Math. Acad. Sci. Paris, 341 (2005), pp. 535–538.
- [18] J. PEETRE, *New Thoughts on Besov Spaces*, Duke University Mathematical Series 1, Duke University Press, Durham, NC, 1976.
- [19] J. PEDLOSKY, *Geophysical Fluid Dynamics*, Springer-Verlag, New York, 1987.
- [20] S. RESNICK, *Dynamical Problem in Nonlinear Advective Partial Differential Equations*, Ph.D. thesis, University of Chicago, 1995.
- [21] H. TRIEBEL, *Theory of Function Spaces*, Leipzig, Germany, 1983.
- [22] M. VISHIK, *Hydrodynamics in Besov spaces*, Arch. Ration. Mech. Anal., 145 (1998), pp. 197–214.
- [23] J. WU, *Solutions to the 2D quasi-geostrophic equations in Hölder spaces*, Nonlinear Anal., 62 (2005), pp. 579–594.
- [24] J. WU, *Global solutions of the 2D dissipative quasi-geostrophic equation in Besov spaces*, SIAM J. Math. Anal., 36 (2005), pp. 1014–1030.

A Γ -CONVERGENCE RESULT FOR THIN MARTENSITIC FILMS IN LINEARIZED ELASTICITY*

PETER HORNING[†]

Abstract. The elastic energy of a thin film Ω_h of thickness h with displacement $u : \Omega_h \rightarrow \mathbb{R}^3$ is given by the functional $E^h(u) = \int_{\Omega_h} W(\nabla u)$. We consider materials whose energy density W is linearly frame indifferent and vanishes on two linearized wells which are compatible in the plane but incompatible in the thickness direction. We prove compactness of displacement sequences $u^{(h)} : \Omega_h \rightarrow \mathbb{R}^3$ satisfying $E^h(u^{(h)}) \leq Ch^2$, and we derive the Γ -limit of the functionals $\frac{1}{h^2} E^h$ as $h \rightarrow 0$.

Key words. singular variational problems, thin films, martensitic phase transitions, Γ -convergence

AMS subject classifications. 74G65, 49J45, 74N99

DOI. 10.1137/070683167

1. Introduction. The study of solid-solid phase transitions in thin elastic films leads to functionals of the form

$$(1) \quad E^h(v) = \int_{\Omega_h} W(\nabla v(x)) \, dx,$$

where $\Omega_h = S \times (-\frac{h}{2}, \frac{h}{2})$ is a cylindrical domain of thickness h , $S \subset \mathbb{R}^2$ is a bounded Lipschitz domain, $v : \Omega_h \rightarrow \mathbb{R}^3$ is the elastic deformation (in the nonlinear setting) or the displacement (in the linearized setting), and W is a free energy density with n energy minima F_i , i.e., $W(F_i) = 0$ for $i = 1, \dots, n$. In the context of nonlinear elasticity W is invariant under proper rotations, and in the context of linearized elasticity it is invariant under addition of skew symmetric matrices. In [10] it was observed (in the context of nonlinear elasticity) that for many materials which undergo austenite-martensite phase transitions, the low-energy states of very thin samples of material display a much richer variety of structures than bulk samples made of the same material. The reason is that three dimensional compatibility requires a plane on which two juxtaposed affine deformations coincide, i.e., that their gradients be rank-one connected. In contrast, two dimensional compatibility is already satisfied if there exists one in-plane vector on which the two deformations agree, so a rank-two connection between the gradients suffices. This fact leads to the existence of many nontrivial low-energy states, including laminates, tunnels, and tents; see, e.g., [11] for experimental results. This rich structure makes thin martensitic films particularly interesting for applications.

In this article we study the asymptotic behavior of thin martensitic films in the context of linearized elasticity. In our model the zero set of the energy density W consists of two linearized wells which are incompatible in bulk but compatible in the plane (see section 2 for details). We study the asymptotic behavior of the functionals (1) in the thin film-limit $h \rightarrow 0$. We prove compactness of displacement sequences

*Received by the editors February 20, 2007; accepted for publication (in revised form) December 3, 2007; published electronically April 4, 2008. This work was supported by the EU research and training network MULTIMAT MRTN-CT-2004-505226.

<http://www.siam.org/journals/sima/40-1/68316.html>

[†]Fachbereich Mathematik, Universität Duisburg-Essen, Lotharstrasse 65, 47057 Duisburg, Germany (peter.hornung@uni-due.de).

whose energy scales like h^2 , and we derive the Γ -limit of the functionals $\frac{1}{h^2} E^h$ as the film thickness h converges to zero. To our knowledge this is the first Γ -convergence result for thin martensitic films in which both the domain and the image space are three dimensional and no interfacial energy term is added to the elastic energy. In our model, the formation of interfaces is penalized in a natural way by the interplay of nonzero film thickness with incompatibility of the energy wells in the thickness direction.

Thin films of single-well materials have been studied, e.g., in [1, 3, 12, 17] (in a linearized setting) and in [34, 13, 28, 29] (in a nonlinear setting). Thin films of multiwell materials have been studied, e.g., in [19] (in a linearized setting) and in [37, 10, 15, 6] (in a nonlinear setting).

To state our main result let us introduce the functionals

$$I^h(u; S) = \begin{cases} \frac{1}{h^2} \int_{S \times (-\frac{h}{2}, \frac{h}{2})} W(\nabla u(x)) \, dx & \text{if } u \in W^{1,2}(S \times (-\frac{h}{2}, \frac{h}{2}); \mathbb{R}^3), \\ +\infty & \text{otherwise} \end{cases}$$

and

$$I^0(w; S) = \begin{cases} \int_J k(\nu(x)) d\mathcal{H}^1(x) & \text{if } w \in \mathcal{A}(S), \\ +\infty & \text{otherwise,} \end{cases}$$

where the class $\mathcal{A}(S)$ of admissible limiting displacements is given in (30) below, J denotes the jump set of $\text{sym } \nabla' w \in BV$, and ν denotes the normal to it, which can assume only two values (up to a sign). The function k is a ‘‘surface tension’’ which depends on the normal and which we define in (29) below. Let us write v' to denote the first two entries of $v \in \mathbb{R}^3$, and let us call a domain $S \subset \mathbb{R}^2$ strictly star-shaped if there is $z \in S$ such that for all $z' \in \bar{S}$ the open segment (z, z') is contained in S . Our main result is the following theorem.

THEOREM 1. *Let $A, B \in \mathbb{R}^{3 \times 3}$ satisfy (i)–(iv) from section 2, let W satisfy the conditions (4)–(6) below, and let $S \subset \mathbb{R}^2$ be a bounded strictly star-shaped Lipschitz domain. Then a Γ -type convergence $I^h(\cdot; S) \xrightarrow{\Gamma} I^0(\cdot; S)$ holds in the following sense:*

(i) *Ansatz-free lower bound. Let $w \in L^2(S; \mathbb{R}^2)$, $h_n \rightarrow 0$, let $v_n \in W^{1,2}(S \times (-\frac{h_n}{2}, \frac{h_n}{2}); \mathbb{R}^3)$, and set $w_n(x') = \frac{1}{h_n} \int_{-\frac{h_n}{2}}^{\frac{h_n}{2}} v'_n(x', x_3) \, dx_3$. If $w_n \rightarrow w$ in $L^2(S; \mathbb{R}^2)$, then $\liminf_{n \rightarrow \infty} I^{h_n}(v_n; S) \geq I^0(w; S)$.*

(ii) *Existence of recovery sequences. Let $w \in L^2(S; \mathbb{R}^2)$ and $h_n \rightarrow 0$. Then there is a sequence $v_n \in W^{1,2}(S \times (-\frac{h_n}{2}, \frac{h_n}{2}); \mathbb{R}^3)$ such that, setting $w_n(x') = \frac{1}{h_n} \int_{-\frac{h_n}{2}}^{\frac{h_n}{2}} v'_n(x', x_3) \, dx_3$, we have $w_n \rightarrow w$ strongly in $W^{1,2}(S; \mathbb{R}^2)$ and $\lim_{n \rightarrow \infty} I^{h_n}(v_n; S) = I^0(w; S)$.*

Remarks. (i) Theorem 1 is complemented by a compactness result for sequences v_n whose energy $E^{h_n}(v_n)$ scales like h_n^2 (Theorem 6 below).

(ii) Notice that in Theorem 1(ii) we state the existence of recovery sequences for any given sequence $h_n \rightarrow 0$.

(iii) The vertical average w_n can be interpreted as the in-plane displacement of the midplane S of the thin film; compare, e.g., [3].

(iv) The lower bound (i) is true for general (also non–star-shaped) Lipschitz domains S ; see Theorem 12. Star-shapedness is used in the proof of the upper bound (ii) to show that limiting displacements with finitely many well-separated interfaces are energy dense and to avoid the necessity of a lateral matching of two local recovery sequences. The same technical difficulty concerning non–star-shaped domains occurs

in [19, 20] and in [18]. A figure depicting the problematic situation can be found in [18, Figure 4].

(v) The Γ -limit obtained in Theorem 1 has the same structure as that derived in [19]. The reason is that the functionals I^h are related to singularly perturbed functionals of the form

$$(2) \quad J^{(\varepsilon)}(u; S) = \int_S \frac{1}{\varepsilon} W_{2D}(\nabla u(x')) + \varepsilon |\nabla^2 u(x')|^2 \, dx'.$$

The asymptotic behavior of singularly perturbed functionals has been extensively studied in the literature (see [35, 27, 26, 5, 38] and, for gradient phase transitions, see [25, 18, 19, 20]). Recently, Conti and Schweizer derived the Γ -limit of the functionals (2) both under the assumption of linearized frame indifference [19] and under nonlinear frame indifference [20] of W_{2D} . Their results are restricted to two dimensions.

There are two crucial differences between functionals of the form (2) and the model studied in this article: In the former the domain and the image space are two dimensional and an extra term $\int |\nabla^2 u|^2$, weighted with some small parameter ε^2 , is added to the elastic energy. The role of this term is to penalize the formation of phase interfaces. In our model (1), the domain and the image space are genuinely three dimensional and the energy functional does not involve higher derivatives: It is a key property of our model that no interfacial energy contribution is added to the elastic energy (this also contrasts with other thin film models [10, 37, 6]). The formation of phase interfaces is naturally penalized by the interplay of nonzero film thickness with the fact that the zero energy displacements $A + \text{Skew}$ and $B + \text{Skew}$ are incompatible in the thickness direction.

Lemma 14 makes the relation between (1) and (2) more precise. This requires a subtle mollification argument since (2) requires control on the second derivatives. Lemma 14 suggests that the small parameter ε should be interpreted as the film thickness h .

(vi) No nonlinear version of Theorem 1 has yet been proven. The only result in this direction is [15], where it is shown that the energy of thin-film deformations consisting of two phases scales like h^2 . Notice that, in contrast to the linearized setting, the model considered in [20] would not be appropriate to describe thin martensitic films since it is too rigid. Nonlinearly elastic rods of multiphase materials were studied in [36].

This article is organized as follows. In section 2 we introduce some definitions and reduce the problem to a canonical form. Then we prove a two-well analogue of Korn's inequality, Theorem 3, which applies to incompatible linear wells. Then we apply this result to deduce the compactness result Theorem 6. In section 3 we obtain the lower bound, Theorem 12, by an abstract scaling argument. Finally, in section 4 we derive the upper bound by constructing three dimensional recovery sequences. The proof of Theorem 1 closes section 4.

Notation. We use the letter C to denote constants depending only on the domain and on W . Within an expression the explicit value of C may change from line to line. A bar above a given 3×3 matrix denotes its upper left 2×2 submatrix, and in general we use barred letters to denote 2×2 matrices. Primes on 3-vectors will denote the 2-vector consisting of the first two entries, so in particular $x = (x', x_3)$. For a matrix A we write $\text{sym } A = \frac{1}{2}(A + A^T)$, $\text{skew } A = \frac{1}{2}(A - A^T)$, and $|A|^2 = \text{Tr}(A^T A)$, where Tr denotes the trace. By a subscript i we will denote the partial derivative with respect to the x_i -variable. By ∇' we denote the in-plane gradient, that is, $\nabla' w = (w_{,1}|w_{,2})$. For $h > 0$ we set $I_h = (-\frac{h}{2}, \frac{h}{2})$. All intervals in this article are implicitly assumed to

be nonempty and bounded. We use a dashed integral sign $\overline{\int}$ to denote the average. Often we will simply write $\{f = a\}$ instead of $\{x \in S : f(x) = a\}$. For $\rho > 0$ we set $[\rho] = \max\{n \in \mathbb{N} : n \leq \rho\}$. If $E \subset \mathbb{R}^n$, then $|E|$ denotes its n dimensional Lebesgue measure and $\mathcal{H}^k(E)$ its k dimensional Hausdorff measure [24]. For $j \in \{1, 2\}$ we denote by e_j the j th unit vector, by $\pi_j : \mathbb{R}^2 \rightarrow \text{span}\{e_j\}$ the orthogonal projection onto $\text{span}\{e_j\}$, and by $\pi_j^\perp : \mathbb{R}^2 \rightarrow \{e_j\}^\perp$ the orthogonal projection onto $\{e_j\}^\perp$. If $E \subset \mathbb{R}^n$, then $B_\varepsilon(E) = \{x \in \mathbb{R}^n : \text{dist}_E(x) < \varepsilon\}$.

2. Preliminaries and compactness. We consider the functional, defined for any Lipschitz domain $U \subset \mathbb{R}^2$,

$$(3) \quad I^h(u; U) = \begin{cases} \frac{1}{h^2} \int_{U \times I_h} W(\nabla u(x)) \, dx & \text{if } u \in W^{1,2}(U \times I_h; \mathbb{R}^3), \\ +\infty & \text{otherwise.} \end{cases}$$

Throughout this article, W is assumed to satisfy the following conditions:

- (4) $W : \mathbb{R}^{3 \times 3} \rightarrow \mathbb{R}$ is continuous,
- (5) linearized frame indifference: $W(F) = W(\text{sym } F)$ for all $F \in \mathbb{R}^{3 \times 3}$,
- (6) quadratic growth and coercivity: $c_0 W_0(F) \leq W(F) \leq C_0 W_0(F)$,

where c_0, C_0 are positive constants. Here we have introduced the standard energy density $W_0(F) = \text{dist}^2(\text{sym } F, \{A, B\})$, where A and B are symmetric 3×3 matrices to be specified below. We define the reduced functional

$$I_{2D}^h(w; U) = \begin{cases} \int_U \frac{1}{h} W_{2D}(\nabla' w) + h |\nabla'^2 w|^2 & \text{if } w \in W^{2,2}(U; \mathbb{R}^2), \\ +\infty & \text{otherwise,} \end{cases}$$

where $W_{2D}(\bar{F}) = \text{dist}^2(\text{sym } \bar{F}, \{\bar{A}, \bar{B}\})$. We make the following assumptions on the wells A and B :

- (i) $A \in \mathbb{R}^{3 \times 3}$ and $B \in \mathbb{R}^{3 \times 3}$ are symmetric; i.e., $A = \text{sym } A$ and $B = \text{sym } B$.
- (ii) *Incompatibility in bulk.* $\text{rank}(A - B + T) \geq 2$ for all $T \in \mathbb{R}^{3 \times 3}$ with $\text{sym } T = 0$.
- (iii) *Compatibility in the plane.* There exists $\bar{T} \in \mathbb{R}^{2 \times 2}$ with $\text{sym } \bar{T} = 0$ such that $\text{rank}(\bar{A} - \bar{B} + \bar{T}) \leq 1$.
- (iv) *Nondegeneracy.* $\det(\bar{A} - \bar{B}) \neq 0$.

Item (iii) is satisfied if and only if there exists $t \in \mathbb{R}$ such that

$$0 = \det \left(\bar{A} - \bar{B} + \begin{pmatrix} 0 & t \\ -t & 0 \end{pmatrix} \right) = \det(\bar{A} - \bar{B}) + t^2,$$

whence (iii) is equivalent to $\det(\bar{A} - \bar{B}) \leq 0$ with equality if and only if \bar{A} and \bar{B} are rank-one connected. Thus (iii) and (iv) together are equivalent to $\det(\bar{A} - \bar{B}) < 0$. Table 11.1 in [7] shows that conditions (i)–(iv) are generically satisfied by real materials (in a linearized framework).

Let us now reduce the set of all matrices satisfying (i)–(iv) to a canonical form. Let \tilde{A}, \tilde{B} satisfy conditions (i)–(iv) but be arbitrary otherwise. Then there is an orthogonal matrix $R \in O(3)$ with $Re_3 = e_3$ such that

$$R^T(\tilde{B} - \tilde{A})R = \lambda_1 e_1 \otimes e_1 + \lambda_2 e_2 \otimes e_2 + \sum_{i=1}^3 \tilde{\mu}_i \frac{e_i \otimes e_3 + e_3 \otimes e_i}{2},$$

where λ_i are the eigenvalues of the matrix $\bar{B} - \bar{A}$ and $\tilde{\mu}_i$ are some real numbers. By possibly choosing R differently (by interchanging the first two columns), we may assume that $\lambda_1 \geq \lambda_2$, so since $\det(\bar{A} - \bar{B}) < 0$, we must in fact have $\lambda_1 > 0 > \lambda_2$. Let $Q = \text{diag}(|\lambda_1|^{-\frac{1}{2}}, |\lambda_2|^{-\frac{1}{2}}, 1)$ and set $\hat{B} = QR^T(\bar{B} - \bar{A})RQ$. This gives $\hat{B} = e_1 \otimes e_1 - e_2 \otimes e_2 + \sum_{i=1}^3 \hat{\mu}_i \frac{e_i \otimes e_3 + e_3 \otimes e_i}{2}$, where $\hat{\mu}_i$ are related to $\tilde{\mu}_i$ and λ_i . Now we can find a proper rotation $\hat{Q} \in SO(3)$ with eigenvector e_3 such that

$$(7) \quad B = \hat{Q}^T \hat{B} \hat{Q} = e_1 \otimes e_2 + e_2 \otimes e_1 + \sum_{i=1}^3 \mu_i \frac{e_i \otimes e_3 + e_3 \otimes e_i}{2}$$

for some $\mu_i \in \mathbb{R}$. Since the structural assumptions on the energy density W and on the shape of the domain (i.e., strict star-shapedness with respect to the origin and a cylindrical form $S \times I_h$) are invariant under the transformations introduced above, we obtain the following lemma.

LEMMA 2. *If Theorem 1 is shown for the special pairs A, B given by $A = 0$ and B as in (7), then it holds for all possible choices of A and B which satisfy conditions (i)–(iv).*

2.1. Korn’s inequality for two incompatible strains. The following theorem provides a generalization of Korn’s inequality to the case of two incompatible linearized wells. A nonquantitative version of this result can be found in [22]; compare also [39, 9, 8, 21]. In [19] an example is provided which shows that no Korn-type rigidity like the one derived here can be expected in the case of two compatible wells.

THEOREM 3. *Let $\Omega \subset \mathbb{R}^n$ be a bounded connected Lipschitz domain, $n \geq 2$, and $K = (A + \text{Skew}) \cup (B + \text{Skew})$, where A and B are incompatible strains, i.e., $(B - A) + \text{Skew}$ does not contain rank-one matrices. Then there exists a positive constant $C(\Omega, A, B)$ with the following property: For every $u \in W^{1,2}(\Omega; \mathbb{R}^n)$, there exists an associated $R \in K$ such that $\|\nabla u - R\|_{L^2(\Omega; \mathbb{R}^{n \times n})} \leq C(\Omega, A, B) \|\text{dist}(\nabla u, K)\|_{L^2(\Omega; \mathbb{R}^{n \times n})}$.*

This theorem will follow from the interior estimate provided by the following lemma.

LEMMA 4. *With assumptions as in Theorem 3 and $\bar{U} \subset \Omega$, where U is Lipschitz and connected, there is a constant $C(U, \Omega, A, B)$ such that the following holds: For every $u \in W^{1,2}(\Omega; \mathbb{R}^n)$, there exists an associated $R \in K$ such that*

$$(8) \quad \|\nabla u - R\|_{L^2(U; \mathbb{R}^{n \times n})} \leq C(U, \Omega, A, B) \|\text{dist}(\nabla u, K)\|_{L^2(\Omega; \mathbb{R}^{n \times n})}.$$

Proof. From (15) on this proof follows [14] rather closely with some minor changes. Define $d(F) = \text{dist}(F, \{A, B\})$ and set $\varepsilon^2 = \int_{\Omega} \text{dist}^2(\nabla u, K)$. Notice that $\text{dist}^2(F, K) = d^2(\text{sym } F)$ for all $F \in \mathbb{R}^{n \times n}$. Since $|F - A| \leq d(F) + |A - B|$, by Korn’s inequality there exists $C > 0$ with the property that (8) is satisfied whenever u is such that $\varepsilon \geq 1$. Hence we may assume without loss of generality that $\varepsilon < 1$. By setting $\tilde{B} = B - A$ and applying the lemma to $\tilde{u}(x) = u(x) - Ax$ we may also assume without loss of generality that $A = 0$.

Denote by $P : \mathbb{R}^{n \times n} \rightarrow ((\text{span}\{B\}) \oplus \text{Skew})^\perp$ the orthogonal projection onto the orthogonal complement of the subspace $(\text{span}\{B\}) \oplus \text{Skew}$. Since $B + \text{Skew}$ does not contain rank-one matrices, the only rank-one matrix contained in $(\text{span}\{B\}) \oplus \text{Skew}$ is the zero matrix. Thus $|P(a \otimes b)|^2 > 0$ for all $a, b \neq 0$. Hence, by continuity and by compactness of the sphere, P satisfies the Legendre–Hadamard ellipticity condition $\Lambda |a|^2 |b|^2 \geq |P(a \otimes b)|^2 \geq \lambda |a|^2 |b|^2$ for some $\Lambda > \lambda > 0$. Now let $w \in W^{1,2}(\Omega; \mathbb{R}^n)$ be

a weak solution of the linear elliptic system with constant coefficients

$$(9) \quad \begin{aligned} \operatorname{div} P(\nabla w) &= 0 \text{ in } \Omega, \\ w &= u \text{ on } \partial\Omega. \end{aligned}$$

Set $z = u - w$. Then $z \in W_0^{1,2}(\Omega; \mathbb{R}^n)$ is a weak solution of $\operatorname{div} P(\nabla z) = \operatorname{div} P(\nabla u)$. Testing this with z gives

$$(10) \quad \int_{\Omega} P(\nabla z) : \nabla z = \int_{\Omega} P(\nabla u) : \nabla z \leq \left(\int_{\Omega} |P(\nabla u)|^2 \right)^{\frac{1}{2}} \left(\int_{\Omega} |\nabla z|^2 \right)^{\frac{1}{2}}.$$

Since by ellipticity the left-hand side of (10) is greater than $\int_{\Omega} \lambda |\nabla z|^2$, we conclude

$$(11) \quad \int_{\Omega} |\nabla z|^2 \leq C \int_{\Omega} |P(\nabla u)|^2 = C \int_{\Omega} \operatorname{dist}^2(\operatorname{sym} \nabla u, \operatorname{span}\{B\}) \leq C\varepsilon^2.$$

Thus it remains to prove that there exists $R \in K$ such that $\int_{\Omega} |\nabla w(x) - R|^2 dx \leq C\varepsilon^2$, where C is independent of w .

Set $e_w = \operatorname{sym} \nabla w$ and let $y \in \Omega$ be such that $B(y, 2r) \subset \Omega$. By Korn's inequality there is a $C = C(n)$ (which by scaling invariance is independent of r), and $T \in \mathbb{R}^{n \times n}$ with $\operatorname{sym} T = 0$ such that

$$(12) \quad \int_{B(y, 2r)} |\nabla w - T|^2 \leq C \int_{B(y, 2r)} |e_w|^2.$$

Since by $P(M) = P(\operatorname{sym} M)$ we have $P(T) = 0$, the mapping $v(x) = w(x) - Tx$ is a weak solution of

$$(13) \quad \begin{aligned} \operatorname{div} P(\nabla v) &= 0 \text{ in } \Omega, \\ v &= u - Tx \text{ on } \partial\Omega. \end{aligned}$$

By standard elliptic regularity for linear systems with constant coefficients (see, e.g., [30]), we obtain the inequality

$$(14) \quad \int_{B(y, r)} |\nabla^2 v|^2 \leq \frac{C}{r^2} \int_{B(y, 2r)} |\nabla v|^2 = \frac{C}{r^2} \int_{B(y, 2r)} |\nabla w - T|^2.$$

We have $|\nabla e_w|^2 = \frac{1}{4} \sum_{i,j,k} (w_{i,jk} + w_{j,ik})^2 \leq |\nabla^2 w|^2$. Hence by the choice of T and since $|\nabla^2 w|^2 = |\nabla^2 v|^2$ on $B(y, 2r)$, we conclude from (12) and (14) that

$$(15) \quad \int_{B(y, r)} |\nabla e_w|^2 \leq \frac{C}{r^2} \int_{B(y, 2r)} |e_w|^2.$$

This inequality holds for all $y \in \Omega$ with $B(y, 2r) \subset \Omega$.

Fix $r_0 \in (0, \frac{\operatorname{dist}(U, \partial\Omega)}{4})$ such that there exists $c_0 > 0$ with the property that $|B_r(x) \cap U| \geq c_0 |B_r|$ for all $x \in U$ and for all $r \leq r_0$. (Here and in what follows we will sometimes omit the center of the ball when denoting its volume.) Existence of such an r_0 follows from the Lipschitz property of U , and c_0 will depend on U . Covering \bar{U} with finitely many balls of radius $\frac{1}{3} \operatorname{dist}(U, \partial\Omega)$ and applying (15) shows that $\int_U |\nabla e_w|^2$ is bounded by a constant independent of u (since $|e_w| \leq d(e_w) + C$

and $\int_{\Omega} d^2(e_w) \leq C\varepsilon^2$ by (11)). Hence, by applying Lemma 5 below with $K_1 = \{0\}$, $K_2 = \{B\}$, and $F = e_w$, we obtain

$$(16) \quad \min \left\{ \int_U |e_w - B|^2, \int_U |e_w|^2 \right\} \leq C \left(\int_U d^2(e_w) \int_U |\nabla e_w|^2 \right)^{\frac{n}{2(n-1)}} + \int_U d^2(e_w) \leq C(\varepsilon^{\frac{n}{n-1}} + \varepsilon^2).$$

Set $\rho = \frac{|B|}{2}$, and let us assume that B is the minority phase in U ; i.e., the set $E = \{x \in U : |e_w(x) - B|^2 \leq \rho^2\}$ satisfies $|E| \leq |\{x \in U : |e_w(x)|^2 \leq \rho^2\}|$. (The case when A is the minority phase is treated similarly.) In particular, this implies $|E| \leq |U \setminus E|$ by the choice of ρ . We have $\rho^2|E| \leq \int_U |e_w|^2$ because $|e_w| \geq \rho$ on E , and similarly $\rho^2|E| \leq \rho^2|U \setminus E| \leq \int_U |e_w - B|^2$. Thus by (16), whenever $\varepsilon < 1$,

$$(17) \quad |E| \leq C_1 \left(\varepsilon^{\frac{n}{n-1}} + \varepsilon^2 \right)$$

for some constant C_1 independent of u . Now we fix $\varepsilon_0 \in (0, 1)$ such that $C_1(\varepsilon_0^{\frac{n}{n-1}} + \varepsilon_0^2) < \frac{c_0}{2}|B_{r_0}|$. From now on we assume that $\varepsilon \leq \varepsilon_0$; the other case is treated at the end of this proof. From (17) we deduce that $|E| < \frac{c_0}{2}|B_{r_0}|$. Our aim is to show that

$$(18) \quad |E| \leq C\varepsilon^2$$

for a constant C independent of u . Using Lemma 5 (notice that the constant in its conclusion is invariant under a rescaling of the domain) as in (16) with B_r replacing U , and then dividing through $|B_r|$ and applying (15), one obtains

$$(19) \quad \min \left\{ \int_{B_r(x)} |e_w|^2, \int_{B_r(x)} |e_w - B|^2 \right\} \leq C \left[\left(\mathcal{M}(|e_w|^2)(x) \int_{B_r(x)} d^2(e_w) \right)^{\frac{n}{2(n-1)}} + \int_{B_r(x)} d^2(e_w) \right]$$

for all $x \in \Omega$ and for all $r > 0$ such that $B_{2r}(x) \subset \Omega$. Here \mathcal{M} denotes the Hardy–Littlewood maximal function, $\mathcal{M}(f)(x) = \sup_{r>0} \int_{B_r(x)} |f|$. Above and in what follows we extend e_w by zero outside Ω .

Claim 1. The set $A_{\infty} = \{x \in \Omega : \mathcal{M}(|e_w|^2)(x) \geq 10|B|^2\}$ satisfies $|A_{\infty}| \leq C\varepsilon^2$.

In fact, $\mathcal{M}(|e_w|^2) \leq \mathcal{M}(|e_w|^2 - 5|B|^2)_+ + 5|B|^2$, whence $x \in A_{\infty}$ implies $\mathcal{M}(|e_w|^2 - 5|B|^2)_+(x) \geq 5|B|^2$. Since $|e_w|^2 \leq 2(d^2(e_w) + |B|^2)$, we have $d^2(e_w) \geq \frac{1}{2}(|e_w|^2 - 5|B|^2)_+$. We conclude that $A_{\infty} \subset \{\mathcal{M}(d^2(e_w)) \geq \frac{5}{2}|B|^2\}$. Thus, by the Hardy–Littlewood maximal theorem [31, Chapter 4], $|A_{\infty}| \leq |\{\mathcal{M}(d^2(e_w)) \geq \frac{5}{2}|B|^2\}| \leq C \int_{\Omega} d^2(e_w)$, which proves Claim 1.

For almost every $x \in E \setminus A_{\infty}$ there is an $r_x \leq r_0$ such that

$$(20) \quad \frac{|E \cap B_{r_x}(x)|}{|B_{r_x}(x)|} = \frac{c_0}{2}.$$

In fact, by the Lebesgue point theorem, for almost all $x \in E \setminus A_{\infty}$, we have $\frac{|E \cap B_r(x)|}{|B_r(x)|} \rightarrow 1$ as $r \rightarrow 0$. On the other hand, $\frac{|E \cap B_r(x)|}{|B_r(x)|} \leq \frac{|E|}{|B_r|}$, which is strictly less than $c_0/2$ for $r > r_0$ by the choice of ε_0 . In particular, $B_{2r_x}(x) \subset \Omega$ for every x as above, by the choice of r_0 .

By Vitali's covering theorem [24, Theorem 1, section 1.5] we can choose countably many $x_i \in E \setminus A_\infty$ satisfying (20) and such that

$$(21) \quad |E \setminus A_\infty| \leq C \sum |B_{r_{x_i}}(x_i)|$$

with pairwise disjoint balls on the right-hand side. By (20) and since $|e_w|^2 \geq \rho^2$ on E and $|e_w - B|^2 \geq \rho^2$ on $U \setminus E$, for every i we have

$$(22) \quad \begin{aligned} \frac{c_0 \rho^2}{2} &\leq \frac{1}{|B_{r_{x_i}}|} \min \left\{ \int_{B_{r_{x_i}}(x_i) \cap E} |e_w|^2, \int_{B_{r_{x_i}}(x_i) \cap U \setminus E} |e_w - B|^2 \right\} \\ &\leq C \left[\left(\int_{B_{r_{x_i}}(x_i)} d^2(e_w) \right)^{\frac{n}{2n-2}} + \int_{B_{r_{x_i}}(x_i)} d^2(e_w) \right]. \end{aligned}$$

For the first inequality we have used that $x_i \in U$, whence $|B_{r_{x_i}}(x_i) \cap U| \geq c_0 |B_{r_{x_i}}|$ by definition of c_0 , and so $|B_{r_{x_i}}(x_i) \cap U \setminus E| \geq \frac{c_0}{2} |B_{r_{x_i}}|$ by (20). For the second inequality we have used (19) and the definition of A_∞ . From (22) we have $\frac{c_0 \rho^2}{2} \leq 2C \max \left\{ \left(\int_{B_{r_{x_i}}(x_i)} d^2(e_w) \right)^{\frac{n}{2n-2}}, \int_{B_{r_{x_i}}(x_i)} d^2(e_w) \right\}$. We conclude that $|B_{r_{x_i}}(x_i)| \leq \max \left\{ \frac{4C}{c_0 \rho^2}, \left(\frac{4C}{c_0 \rho^2} \right)^{\frac{2n-2}{n}} \right\} \int_{B_{r_{x_i}}(x_i)} d^2(e_w)$. Summing over i , from the disjointedness of the $B_{r_{x_i}}(x_i)$ and from (21) we conclude that $|E \setminus A_\infty| \leq C \varepsilon^2$, since $\int_\Omega d^2(e_w) \leq C \varepsilon^2$ by (11). By Claim 1 this implies (18) in the case $\varepsilon \leq \varepsilon_0$. But (18) also holds when $\varepsilon > \varepsilon_0$ by (17) (e.g., by choosing $C = \frac{2C_1}{\varepsilon_0^2}$ in (18); recall that $\varepsilon < 1$). Using (18) we can finally estimate

$$\int_U |e_w|^2 = \int_{U \setminus E} |e_w|^2 + \int_E |e_w|^2 \leq C \left[\int_{U \setminus E} d^2(e_w) + |E| + \int_E d^2(e_w) \right] \leq C \varepsilon^2.$$

The desired estimate now follows from Korn's inequality. \square

The proof of Theorem 3 is completed using a cube decomposition of Ω and applying a weighted Poincaré inequality exactly as in the proof of Theorem 2 in [14]. We have used the following lemma, the proof of which is the same as that of Lemma 2.4 in [14], where one can replace ∇w throughout by an arbitrary matrix-valued $W^{1,2}$ -function F .

LEMMA 5. *Let $n \geq 2$, let $\Omega \subset \mathbb{R}^n$ be a bounded and connected Lipschitz domain, and let K_1, K_2 be compact disjoint subsets of $\mathbb{R}^{n \times n}$, $K = K_1 \cup K_2$. Then there is a constant $C = C(K, \Omega)$, such that for any $F \in W^{1,2}(\Omega; \mathbb{R}^{n \times n})$*

$$\begin{aligned} \min_{i=1,2} \int_\Omega \text{dist}^2(F, K_i) &\leq C(K, \Omega) \left(\int_\Omega \text{dist}^2(F, K) \int_\Omega |\nabla F|^2 \right)^{\frac{n}{2(n-1)}} \\ &\quad + C(K, \Omega) \int_\Omega \text{dist}^2(F, K). \end{aligned}$$

2.2. Compactness. The following theorem provides the compactness result which complements the Γ -convergence result of Theorem 1. Its proof uses some facts which were derived in [15] (in order to prove a lower scaling bound in a nonlinearly elastic setting) and which in spirit are close to [28]. It is different from the Young

measure arguments used in the literature of singularly perturbed functionals (e.g., [27, 18, 19, 20]).

THEOREM 6. *Let $S \subset \mathbb{R}^2$ be a bounded Lipschitz domain, let $A, B \in \mathbb{R}^{2 \times 2}$ be symmetric and such that $(B - A) + \text{Skew}$ does not contain rank-one matrices, and set $K = (A + \text{Skew}) \cup (B + \text{Skew})$. Let $h_n \rightarrow 0$, set $\Omega_{h_n} = S \times I_{h_n}$, and suppose that a sequence $u_n \in W^{1,2}(\Omega_{h_n}; \mathbb{R}^3)$ satisfies*

$$(23) \quad \limsup_{n \rightarrow \infty} \frac{1}{h_n^2} \int_{\Omega_{h_n}} \text{dist}^2(\nabla u_n, K) < \infty.$$

Set $w_n(x') = \int_{I_{h_n}} (u_n(x', x_3))' dx_3$. Then there exist a subsequence (not relabeled) and affine mappings $f_n : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla' f_n = 0$, and there is $w_0 \in W^{1,2}(S; \mathbb{R}^2)$ satisfying $\text{sym } \nabla' w_0 \in BV(S; \{\bar{A}, \bar{B}\})$ such that $w_n + f_n \rightarrow w_0$ strongly in $W^{1,2}(S; \mathbb{R}^2)$.

Proof. For $h > 0$ we consider a lattice of squares $S_{a,h} = a + (-\frac{h}{2}, \frac{h}{2})^2$, $a \in h\mathbb{Z}^2$, and we let $S'_h = \bigcup_{\{a \in h\mathbb{Z}^2 : S_{a,h} \subset S\}} S_{a,h}$. Now apply Theorem 3 to $u^{(h)}$ (here and in what follows we write $u^{(h)}$ instead of u_n and h instead of h_n to avoid cumbersome notation) restricted to each cube $a + (-\frac{h}{2}, \frac{h}{2})^3$ with $a \in h\mathbb{Z}^2$. This yields a piecewise constant map $R^{(h)} : S'_h \rightarrow K$ such that

$$(24) \quad \int_{S_{a,h} \times I_h} |\nabla u^{(h)}(x) - R^{(h)}(x')|^2 dx \leq C \int_{S_{a,h} \times I_h} \text{dist}^2(\nabla u^{(h)}(x), K) dx.$$

Define the piecewise constant map $L^{(h)} : S'_h \rightarrow \{A, B\}$ by setting $L^{(h)}(x) = \text{sym } R^{(h)}(x)$. Let $\varepsilon > 0$ be sufficiently small (to be fixed below). We divide the family of squares $\{S_{a,h} : a \in h\mathbb{Z}^2 \text{ and } S_{a,h} \subset S\}$ into three different groups:

$$(25) \quad a \in \mathcal{A}_0^h \text{ if and only if } \int_{S_{a,h} \times I_h} \text{dist}^2(\nabla u^{(h)}(x), K) dx \geq \varepsilon h^3.$$

If $a \notin \mathcal{A}_0^h$, then the matrix $L^{(h)}(a) \in \{A, B\}$ is such that $\frac{1}{h^3} \int_{S_{a,h} \times I_h} |\text{sym } \nabla u^{(h)} - L^{(h)}(a)|^2 \leq C\varepsilon$. This follows from (24) and (25) by the definition of $L^{(h)}(a)$. Now define

$$a \in \mathcal{A}_1^h \text{ if and only if } a \notin \mathcal{A}_0^h \text{ and } L^{(h)}(a) = A,$$

$$(26) \quad a \in \mathcal{A}_2^h \text{ if and only if } a \notin \mathcal{A}_0^h \text{ and } L^{(h)}(a) = B.$$

For ε small enough, each square $S_{a,h}$ belongs to exactly one of these three groups. Thus the sets

$$(27) \quad \Omega_i^h = \text{int} \left(\bigcup_{a \in \mathcal{A}_i^h} \bar{S}_{a,h} \right),$$

$i = 0, 1, 2$, are disjoint and cover S'_h up to an \mathcal{H}^2 null set. As in [15] one can prove that, for ε small enough, the following implication holds:

$$(28) \quad a \in \mathcal{A}_i^h, a' \in \mathcal{A}_j^h, \text{ and } \mathcal{H}^1(\bar{S}_{a',h} \cap \bar{S}_{a,h}) > 0 \implies j \in \{0, i\},$$

that is, a square of type \mathcal{A}_1^h can only have neighboring squares of type \mathcal{A}_1^h or \mathcal{A}_0^h (never of type \mathcal{A}_2^h), and the analogous statement holds with \mathcal{A}_1^h and \mathcal{A}_2^h swapped. But from (23) and from (25) we deduce $\#\mathcal{A}_0^h \leq \frac{C}{h}$. Since the side-length of each square $S_{a,h}$ is h , this leads to the estimate $\mathcal{H}^1(\partial\Omega_1^h \setminus \partial S) \leq C$; compare [15]. This implies that the characteristic functions $\chi_{\Omega_1^h}$ are bounded in $BV(S)$, whence they have a subsequence converging strongly in $L^1(S)$ and hence (by interpolation) in all $L^p(S)$ with $p < \infty$. Since $\#\mathcal{A}_0^h \leq \frac{C}{h}$, we have $|\Omega_0^h| \leq Ch$, whence $\chi_{\Omega_0^h} \rightarrow 0$ in $L^1(\mathbb{R}^2)$. Hence we also have strong $L^1(S)$ -convergence of $\chi_{\Omega_2^h}$. Note that the respective limit functions, which we denote χ_{Ω_1} and χ_{Ω_2} , both belong to $BV(S)$.

On the other hand, $L^{(h)} = A\chi_{\Omega_1^h} + B\chi_{\Omega_2^h} + L^{(h)}\chi_{\Omega_0^h}$. Let us extend $L^{(h)}$ by zero to all of S . By the convergence $\chi_{S'_h} \rightarrow 1$ in $L^1(S)$ we obtain that $L^{(h)} \rightarrow L$ strongly in $L^2(S; \mathbb{R}^{2 \times 2})$, where $L = \chi_{\Omega_1}A + \chi_{\Omega_2}B \in BV(S; \{A, B\})$. By (24), (23), and Jensen's inequality we have $\int_{S'_h} |\text{sym } \nabla' w^{(h)} - \bar{L}^{(h)}|^2 \leq Ch$. Using $L^{(h)} = 0$ on $S \setminus S'_h$ and applying Jensen's inequality, we find $\int_{S \setminus S'_h} |\text{sym } \nabla' w^{(h)} - \bar{L}^{(h)}|^2 \leq C|S \setminus S'_h| + Ch$. We conclude that $\text{sym } \nabla' w^{(h)} \rightarrow \bar{L}$ strongly in $L^2(S; \mathbb{R}^{2 \times 2})$. Since the subspace of symmetrized gradients is strongly closed in $L^2(S; \mathbb{R}^{2 \times 2})$, there is a $w_0 \in W^{1,2}(S; \mathbb{R}^2)$ such that $\text{sym } \nabla' w_0 = \bar{L} \in BV(S; \{A, B\})$. An application of Korn's and of Poincaré's inequalities on S for each h yields affine mappings $f^{(h)} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla' f^{(h)} = 0$ (explicitly, e.g., $\nabla' f^{(h)} = \text{skew } \int_S \nabla'(w^{(h)} - w_0)$) such that $\|w^{(h)} + f^{(h)} - w_0\|_{W^{1,2}(S; \mathbb{R}^2)}^2 \leq C \int_S |\text{sym } \nabla' w^{(h)} - \text{sym } \nabla' w_0|^2$. This converges to zero because $\text{sym } \nabla' w^{(h)} \rightarrow \bar{L}$ in $L^2(S; \mathbb{R}^{2 \times 2})$. \square

Remark. Let $y^{(h)}(x', x_3) = u^{(h)}(x', hx_3)$ be the rescaled displacements defined on $\Omega = S \times (-\frac{1}{2}, \frac{1}{2})$ and set $\nabla_h y = (\nabla' y | \frac{1}{h} y_{,3})$. Then the proof of Theorem 6 shows that $\text{sym } \nabla_h y^{(h)} \rightarrow L$ strongly in $L^2(\Omega; \mathbb{R}^{3 \times 3})$. Thus $\nabla'(y^{(h)})' \rightarrow \bar{L}$ and $y_{,3} \rightarrow 0$ strongly in L^2 . If, using a scaling analogous to that in [17, section 1.3], we set $(\bar{y}^{(h)})' = (y^{(h)})'$ and $(\bar{y}^{(h)})_3 = h(y^{(h)})_3$, then Korn's inequality on Ω implies that there are affine mappings $F^{(h)} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ with $\text{sym } \nabla F^{(h)} = 0$ and such that $\bar{y}^{(h)} + F^{(h)} \rightarrow F_L$ strongly in $W^{1,2}(\Omega; \mathbb{R}^3)$, where $F_L(x) = (\bar{L}x'_0)$.

In [19, Proposition 2.2], the following characterization is provided for functions whose symmetrized gradient has bounded variation and is supported on two incompatible matrices \bar{A}, \bar{B} . (For earlier results in this direction, compare [23].)

PROPOSITION 7. *Let $S \subset \mathbb{R}^2$ be a bounded Lipschitz domain. Let \bar{A}, \bar{B} satisfy (iii)–(iv) from the beginning of section 2, let ν_1, ν_2 be linearly independent solutions to $\bar{A} - \bar{B} + t(e_1 \otimes e_2 - e_2 \otimes e_1) = a \otimes \nu_i$, where $a \in \mathbb{R}^3$ and $t \in \mathbb{R}$, and let $w \in W^{1,2}(S; \mathbb{R}^2)$ satisfy $\text{sym } \nabla' w \in BV(S; \{\bar{A}, \bar{B}\})$. Then the jump set J of $\text{sym } \nabla' w$ consists of countably many disjoint segments whose endpoints belong to ∂S and which have normal directions ν_1 or ν_2 . In addition, $\nabla' w$ is constant on each connected component of $S \setminus J$.*

3. Lower bound. In this section we prove part (i) of Theorem 1. From now on we assume that $A = 0$ and B is as in (7). This choice allows exactly two different directions for the interface normal; the directions are orthogonal to each other: Setting $T_1 = e_2 \otimes e_1 - e_1 \otimes e_2$ and $T_2 = e_1 \otimes e_2 - e_2 \otimes e_1$, we have $\bar{B} + T_1 = 2e_2 \otimes e_1$, giving the normal $\nu_1 = e_1$, and $\bar{B} + T_2 = 2e_1 \otimes e_2$, giving the normal $\nu_2 = e_2$. Notice that,

if $\bar{T} \in \mathbb{R}^{2 \times 2} \setminus \{T_1, T_2\}$ with $\text{sym } \bar{T} = 0$, then $\det(\bar{B} + \bar{T}) \neq 0$. Define the piecewise constant mappings $F_i^\pm : \mathbb{R}^2 \rightarrow \mathbb{R}^{2 \times 2}$ by

$$F_i^\pm(x') = \begin{cases} 0 & \text{for } \pm x' \cdot \nu_i < 0, \\ \bar{B} + T_i & \text{otherwise,} \end{cases}$$

and set $w_i^\pm(x') = F_i^\pm(x')x'$; note that the jump set of $\nabla' w_i^\pm$ agrees with $\{\nu_i\}^\perp$.

DEFINITION 8. Let $\sigma \in \{-, +\}$, $i \in \{1, 2\}$, let J_1, J_2 be open intervals, and set $S = J_1 \times J_2$.

(i) We set $w_{i,S}^\sigma(x') = w_i^\sigma(x' - \xi e_i)$, where $\xi = \frac{1}{2}(\sup J_i - \inf J_i)$. The set $\mathcal{J}_{i,S} = S \cap (\xi \nu_i + \{\nu_i\}^\perp)$ is called the interface of $w_{i,S}^\sigma$.

(ii) A pair of sequences (u_n, h_n) is called (σ, i) -admissible on S , provided that $h_n \in \mathbb{R}_+$, $h_n \rightarrow 0$, $u_n \in W^{1,2}(S \times I_{h_n}; \mathbb{R}^3)$, and $\int_{I_{h_n}} u_n' dx_3 \rightarrow w_{i,S}^\sigma$ in $L^2(S; \mathbb{R}^2)$.

(iii) If $J_1 = \emptyset$ or $J_2 = \emptyset$, then we set $\mathcal{F}_i^\sigma(S) = 0$. Otherwise, we define

$$\mathcal{F}_i^\sigma(S) = \inf \left\{ \liminf_{n \rightarrow \infty} I^{h_n}(u_n; S) : \text{There exist } (h_n) \text{ and } (u_n) \text{ such that } (u_n, h_n) \text{ is } (\sigma, i)\text{-admissible on } S \right\}.$$

We define

$$(29) \quad k(\nu_i) = \mathcal{F}_i^+ \left(\left(-\frac{1}{2}, \frac{1}{2} \right) \times \left(-\frac{1}{2}, \frac{1}{2} \right) \right).$$

(iii) A pair of sequences (u_n, h_n) is called a (σ, i) -recovery sequence on S , provided that (u_n, h_n) is (σ, i) -admissible on S and $\lim_{n \rightarrow \infty} I^{h_n}(u_n; S) = \mathcal{F}_i^\sigma(S)$.

Remarks. (i) Our definition of \mathcal{F}_i^σ differs slightly from the usual one. If one sets $\tilde{\mathcal{F}}_1^\sigma(J; \varepsilon) = \mathcal{F}_1^\sigma((-\varepsilon, \varepsilon) \times J)$ and $\tilde{\mathcal{F}}_2^\sigma(J; \varepsilon) = \mathcal{F}_2^\sigma(J \times (-\varepsilon, \varepsilon))$, then the $\tilde{\mathcal{F}}_i^\sigma$ correspond to the \mathcal{F}_i^σ as defined, e.g., in [18, formula (4.2)].

(ii) In Lemma 13 we will show that for any fixed sequence $h_n \rightarrow 0$ there exist u_n such that (u_n, h_n) is a (σ, i) -recovery sequence.

(iii) When it is clear from the context which sequence $h_n \rightarrow 0$ is meant, then we will often just say that u_n is a (σ, i) -recovery sequence. Also, we will drop the prefix (σ, i) when it is clear from the context.

LEMMA 9. Let J_1, J_2 be open intervals, set $S = J_1 \times J_2$, and let $i \in \{1, 2\}$. Then $\mathcal{F}_i^+(S) = \mathcal{F}_i^-(S) = k(\nu_i) |\pi_i^\perp(S)|$. (Recall that π_i^\perp denotes the orthogonal projection onto $\{e_i\}^\perp = \{\nu_i\}^\perp$.)

Proof. Similar to [18, Lemma 4.3] or [19, Lemma 3.2], one can prove the lemma by showing the following facts. Let $J'_1 \subset J_1$, $J'_2 \subset J_2$ be open intervals and set $S' = J'_1 \times J'_2$. Then

(i) $\mathcal{F}_i^+(S') = \mathcal{F}_i^-(S') =: \mathcal{F}_i(S')$.

(ii) Behavior under homotheties. $\mathcal{F}_i(x' + \lambda S) = \lambda \mathcal{F}_i(S)$ for all $x' \in \mathbb{R}^2$ and all $\lambda > 0$.

(iii) Monotonicity. $\mathcal{F}_i(S') \leq \mathcal{F}_i(S)$.

(iv) Concentration. $\mathcal{F}_i(S') = \mathcal{F}_i(S)$ if $|\pi_i^\perp(S')| = |\pi_i^\perp(S)|$. \square

LEMMA 10. Let $i \in \{1, 2\}$, $\sigma \in \{+, -\}$, let J_1, J_2 be open intervals, let $\lambda > 0$, $y \in \mathbb{R}^2 \times \{0\}$, and let (v_n, h_n) be a (σ, i) -recovery sequence on $S = J_1 \times J_2$. Set $\hat{v}_n(x) = \lambda v_n(\frac{x-y}{\lambda})$ for all $n \in \mathbb{N}$. Then $(\hat{v}_n, \lambda h_n)$ is a (σ, i) -recovery sequence on $y' + \lambda S$.

Proof. Since from Lemma 9 we have $\mathcal{F}_i^\sigma(y' + S) = \mathcal{F}_i^\sigma(S)$, we may assume without loss of generality that $y = 0$ and that S is centered about the origin. Since (v_n, h_n)

is (σ, i) -admissible on S , $(\hat{v}_n, \lambda h_n)$ is also (σ, i) -admissible on λS . By Lemma 9 and a change of variables we have $\limsup_{n \rightarrow \infty} I^{\lambda h_n}(\hat{v}_n; \lambda S) = \limsup_{n \rightarrow \infty} \lambda I^{h_n}(v_n; S) = \lambda \mathcal{F}_i^\sigma(S) = \mathcal{F}_i^\sigma(\lambda S)$. \square

LEMMA 11. Let $i \in \{1, 2\}$, $\sigma \in \{+, -\}$, let J_j, J'_j be open intervals with $J'_j \subset J_j$, $j = 1, 2$, and set $S = J_1 \times J_2$, $S' = J'_1 \times J'_2$. Let (v_n, h_n) be a (σ, i) -recovery sequence on S . Then the following hold:

- (i) If $\overline{S'} \cap \mathcal{J}_{i,S} = \emptyset$, then $\lim_{n \rightarrow \infty} I^{h_n}(v_n; S') = 0$.
- (ii) If $S' \cap \mathcal{J}_{i,S} \neq \emptyset$, then $\lim_{n \rightarrow \infty} I^{h_n}(v_n; S') = \mathcal{F}_i^\sigma(S')$.

Proof. By translation invariance we may assume without loss of generality that S is centered about the origin, so $\mathcal{J}_{i,S} = S \cap \{\nu_i\}^\perp$. Statement (i) is an immediate consequence of Lemma 9. In fact, let $\varepsilon > 0$ be so small that $S_\varepsilon = \{x' \in S : |x' \cdot \nu_i| < \varepsilon\}$ satisfies $S_\varepsilon \cap S' = \emptyset$. Then $\mathcal{F}_i^\sigma(S_\varepsilon) = \mathcal{F}_i^\sigma(S)$ by Lemma 9 because $|\pi_i^\perp(S_\varepsilon)| = |\pi_i^\perp(S)|$. Hence $\lim_{n \rightarrow \infty} I^{h_n}(v_n; S \setminus S_\varepsilon) = 0$.

To prove statement (ii) notice that by Lemma 11(i) we may assume without loss of generality that $\pi_i(S') = \pi_i(S)$ (recall that π_i denotes the orthogonal projection onto the subspace spanned by $e_i = \nu_i$). In other words, $J'_i = J_i$, so S' is a stripe of width $|\pi_i^\perp(S')|$ perpendicular to the interface $\mathcal{J}_{i,S}$. Assume that (ii) is false, and so, since v_n is admissible on S' , we have $\limsup_{n \rightarrow \infty} I^{h_n}(v_n; S') > \mathcal{F}_i^\sigma(S')$. If $S' = S$, then this would contradict the fact that v_n is a recovery sequence on S . Otherwise, denote by S_1, S_2 the two connected components of $S \setminus S'$. (If $S \setminus S'$ consists of only one connected component, then we call it S_1 and set $S_2 = \emptyset$.) After passing to subsequences (not relabeled) we may assume that $I^{h_n}(v_n; S')$ and $I^{h_n}(v_n; S_j)$, $j = 1, 2$, converge. Hence we obtain the contradiction

$$\begin{aligned} \mathcal{F}_i^\sigma(S) &= \mathcal{F}_i^\sigma(S') + \mathcal{F}_i^\sigma(S_1) + \mathcal{F}_i^\sigma(S_2) \\ &< \lim_{n \rightarrow \infty} I^{h_n}(v_n; S') + \lim_{n \rightarrow \infty} I^{h_n}(v_n; S_1) + \lim_{n \rightarrow \infty} I^{h_n}(v_n; S_2) \\ &= \lim_{n \rightarrow \infty} I^{h_n}(v_n; S) = \mathcal{F}_i^\sigma(S). \end{aligned}$$

The first equality follows from Lemma 9; the strict inequality holds because $(v_n|_{S_j}, h_n)$ are admissible on S_j and because by assumption we have $\lim I^{h_n}(v_n; S') > \mathcal{F}_i^\sigma(S')$. The last equality holds because (v_n, h_n) is a recovery sequence on S . \square

Now we define the set of admissible limiting functions as

$$(30) \quad \mathcal{A}(S) = \left\{ w \in W^{1,2}(S; \mathbb{R}^2) : \text{sym } \nabla' w \in BV(S; \{0, \bar{B}\}) \right\}$$

and the limiting functional

$$(31) \quad I^0(w; S) = \begin{cases} \int_J k(\nu(x)) d\mathcal{H}^1(x) & \text{if } w \in \mathcal{A}(S), \\ +\infty & \text{otherwise.} \end{cases}$$

Here J denotes the jump set of $\text{sym } \nabla' w$, also called the phase interface, and ν denotes the normal (the sign does not matter), which up to a sign can assume only the values $\nu_1 = e_1$ and $\nu_2 = e_2$.

THEOREM 12. Let $S \subset \mathbb{R}^2$ be a bounded Lipschitz domain and $w \in L^2(S; \mathbb{R}^2)$. Then, for all $h_n \rightarrow 0$ and all $u_n \in L^2(S \times I_{h_n}; \mathbb{R}^3)$ satisfying $\int_{I_{h_n}} u'_n dx_3 \rightarrow w$ in $L^2(S; \mathbb{R}^2)$, one has $\liminf_{n \rightarrow \infty} I^{h_n}(u_n; S) \geq I^0(w; S)$.

Proof. If $\liminf_{n \rightarrow \infty} I^{h_n}(u_n; S) = \infty$, then there is nothing to prove. Otherwise, by passing to a subsequence (not relabeled) we may assume that the sequence

$I^{h_n}(u_n; S)$ converges, so, in particular, $\limsup_{n \rightarrow \infty} I^{h_n}(u_n; S) < \infty$. After passing to a further subsequence, Theorem 6 implies that there is a sequence of affine mappings $f_n : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla f_n = 0$ such that $w_n + f_n \rightarrow w_0$ in $W^{1,2}(S; \mathbb{R}^2)$ for some $w_0 \in \mathcal{A}(S)$, where we have set $w_n = \int_{I_{h_n}} u'_n dx_3$. Since $w_n \rightarrow w$ in $L^2(S; \mathbb{R}^2)$, we deduce that f_n converges in $L^2(S; \mathbb{R}^2)$, whence there is $\bar{T} \in \mathbb{R}^{2 \times 2}$ with $\text{sym } \bar{T} = 0$ and a vector $c \in \mathbb{R}^2$ such that $f_n(x') \rightarrow c + \bar{T}x'$ for all $x' \in S$. Hence $w = w_0 - \bar{T}x' - c$, and, in particular, we have $w \in \mathcal{A}(S)$. By the strong $W^{1,2}$ -convergence of both $w_n + f_n$ and f_n we have $w_n \rightarrow w$ in $W^{1,2}(S; \mathbb{R}^2)$. By Proposition 7 the jump set of $\text{sym } \nabla' w$ consists of a countable union of disjoint segments \mathcal{J}_k with normal ν_1 or ν_2 . The rest of the proof is standard: One covers each \mathcal{J}_k with a box, applies Lemma 9 to each box separately, and uses the minimality of \mathcal{F}_i^\pm (see, e.g., the proof of Proposition 3.1 in [19] for the details). \square

4. Upper bound. In this section we will show that for any admissible limit function w and for any given sequence $h_n \rightarrow 0$ one can find a sequence (v_n, h_n) such that $\int_{I_{h_n}} v'_n dx_3 \rightarrow w$ strongly in $W^{1,2}(S; \mathbb{R}^2)$ and $I^{h_n}(v_n; S) \rightarrow I^0(w; S)$. We will first show that given (σ, i) and any sequence $h_n \rightarrow 0$ one can find a sequence (v_n) such that (v_n, h_n) is a (σ, i) -recovery sequence on S . A key difference from the proof of the analogous Proposition 5.5 in [19] is that we do not rely on the existence of special recovery sequences which are affine away from the interface but work directly with an arbitrary recovery sequence.

LEMMA 13. *Let J_1, J_2 be open intervals and set $S = J_1 \times J_2$. Let $\sigma \in \{+, -\}$, $i \in \{1, 2\}$, and $H_n \rightarrow 0$ be given. Then we have*

$$\mathcal{F}_i^\sigma(S) = \inf \left\{ \liminf_{n \rightarrow \infty} I^{H_n}(u_n; S) : \text{there is } (u_n) \text{ such that } (u_n, H_n) \right.$$

$\left. \text{is } (\sigma, i)\text{-admissible on } S \right\}$.

Proof. Clearly we must prove only the “ \geq ”-inequality, and by Lemmas 10 and 11 we may assume without loss of generality that $S = (-\frac{1}{2}, \frac{1}{2})^2$. In fact, suppose Lemma 13 is shown for this particular case. Now let $H_n \rightarrow 0$ and an arbitrary S be given and assume without loss of generality (by translation invariance) that S is centered about the origin. Then there is $\lambda > 0$ such that $\lambda S \subset (-\frac{1}{2}, \frac{1}{2})^2$. By the special case of Lemma 13 there is a recovery sequence $(v_n, \lambda H_n)$ on $(-\frac{1}{2}, \frac{1}{2})^2$. By Lemma 11(ii) we have that $(v_n|_{\lambda S}, \lambda H_n)$ is a recovery sequence on λS . By Lemma 10, setting $\hat{v}_n(x) = \frac{1}{\lambda} v_n(\lambda x)$, we conclude that (\hat{v}_n, H_n) is a recovery sequence on S .

So suppose that $S = (-\frac{1}{2}, \frac{1}{2})^2$ and let $H_n \rightarrow 0$ be given. Let us restrict our attention to the case $\sigma = +$ and $i = 1$, so the interface normal is $\nu_1 = e_1$, the phase “0” is used on the left, $\{\nabla' w_1^+ = 0\} = S \cap \{x_1 < 0\}$, and $\{\nabla' w_1^+ = \bar{B} + T_1\} = S \cap \{x_1 > 0\}$ up to a null set; the other cases are similar. Note that the infimum in the definition of $\mathcal{F}_i^\sigma(S) = k(\nu_i)$ is attained; i.e., there is a sequence $h_n \rightarrow 0$ and $v_n \in W^{1,2}(S \times I_{h_n}; \mathbb{R}^3)$ such that $\int_{I_{h_n}} v'_n dx_3 \rightarrow w_1^+$ in $L^2(S; \mathbb{R}^2)$ and $\lim_{n \rightarrow \infty} I^{h_n}(v_n; S) = k(\nu_1)$. Since after passing to subsequences this equality remains valid, we may assume without loss of generality that $h_n \ll H_n$, i.e., $\alpha_n = \frac{H_n}{h_n} \rightarrow \infty$.

Set $y_1^{(n)} = \frac{1}{2\alpha_n} - \frac{1}{2}$ and $y_{m+1}^{(n)} = y_m^{(n)} + \frac{1}{\alpha_n}$, $m = 1, \dots, [\alpha_n] - 1$, and let $S_m^{(n)} = (-\frac{1}{2}, \frac{1}{2}) \times (y_m^{(n)} - \frac{1}{2\alpha_n}, y_m^{(n)} + \frac{1}{2\alpha_n})$. We recall (25)–(28) from the proof of Theorem 6 and apply them to v_n instead of $u^{(h)}$. Define $S_{m,1}^{(n)} = S_m^{(n)} \cap \Omega_1^{h_n} \cap \{\nabla' w_1^+ = 0\}$ and $S_{m,2}^{(n)} = S_m^{(n)} \cap \Omega_2^{h_n} \cap \{\nabla' w_1^+ = \bar{B} + T_1\}$. It follows from the proof of Theorem 6 applied to v_n that $\chi_{\Omega_1^{h_n}} \rightarrow \chi_{\{\nabla' w_1^+ = 0\} \cap S}$ and $\chi_{\Omega_2^{h_n}} \rightarrow \chi_{\{\nabla' w_1^+ = \bar{B} + T_1\} \cap S}$ strongly in $L^1(\mathbb{R}^2)$

(notice that, by uniqueness of the limits, the full sequence indeed converges). Now denote by $G_n = \{m = 1, \dots, [\alpha_n] : |S_{m,1}^{(n)}| > \frac{1}{3\alpha_n} \text{ and } |S_{m,2}^{(n)}| > \frac{1}{3\alpha_n}\}$ the index set of “good” stripes, and for $i = 1, 2$ set $M_{n,i} = \{m = 1, \dots, [\alpha_n] : |S_{m,i}^{(n)}| \leq \frac{1}{3\alpha_n}\}$. We claim that

$$(32) \quad \frac{\#G_n}{\alpha_n} \rightarrow 1 \text{ as } n \rightarrow \infty.$$

Indeed, if (32) were false, then by definition of G_n there would exist $i \in \{1, 2\}$, $\gamma > 0$, and a subsequence (not relabeled) such that $\frac{\#M_{n,i}}{\alpha_n} \in (\frac{5}{6}\gamma, \gamma)$ for all n . For definiteness suppose that $i = 1$; the case $i = 2$ is similar. But for all n we have $\Omega_1^{h_n} \cap \{\nabla' w_1^+ = 0\} \subset (S \setminus \bigcup_{m=1}^{[\alpha_n]} S_m^{(n)}) \cup \bigcup_{m \notin M_{n,1}} S_{m,1}^{(n)} \cup \bigcup_{m \in M_{n,1}} S_{m,1}^{(n)}$. Taking measures on both sides, we find (notice that by definition $|S_{m,i}^{(n)}| \leq \frac{1}{2\alpha_n}$)

$$\begin{aligned} |\Omega_1^{h_n} \cap \{\nabla' w_1^+ = 0\}| &\leq 1 - \frac{[\alpha_n]}{\alpha_n} + \frac{1}{2\alpha_n}([\alpha_n] - \#M_{n,1}) + \frac{1}{3\alpha_n} \#M_{n,1} \\ &\leq 1 - \frac{[\alpha_n]}{\alpha_n} + \frac{1}{2} \left(1 - \frac{5}{6}\gamma\right) + \frac{\gamma}{3}. \end{aligned}$$

As $n \rightarrow \infty$, the left-hand side converges to $\frac{1}{2}$ and the right-hand side converges to $\frac{1}{2} - \frac{1}{12}\gamma$. This contradiction proves (32).

From (32) one deduces that

$$(33) \quad \limsup_{n \rightarrow \infty} (\alpha_n \min_{m \in G_n} I^{h_n}(v_n; S_m^{(n)})) \leq \limsup_{n \rightarrow \infty} I^{h_n}(v_n; S).$$

In fact, if (33) were false, then (after passing to an unlabeled subsequence) there would exist $\gamma > 0$ such that $I^{h_n}(v_n; S_m^{(n)}) - \frac{1}{\alpha_n} I^{h_n}(v_n; S) \geq \frac{\gamma}{\alpha_n}$ for all n and for all $m \in G_n$. Summing over $m \in G_n$ we would find $(1 - \frac{\#G_n}{\alpha_n}) I^{h_n}(v_n; S) \geq I^{h_n}(v_n; \bigcup_{m \in G_n} S_m^{(n)}) - \frac{\#G_n}{\alpha_n} I^{h_n}(v_n; S) \geq \frac{\#G_n}{\alpha_n} \gamma$. By (32) the left-hand side converges to zero as $n \rightarrow \infty$, while the right-hand side converges to γ , which is a contradiction proving (33).

Now choose $m_n \in G_n$ such that $I^{h_n}(v_n; S_{m_n}^{(n)}) = \min_{m \in G_n} I^{h_n}(v_n; S_m^{(n)})$, and set $\hat{y}_n = y_{m_n}^{(n)}$. Let $r^{(n)} = (-\frac{l}{2\alpha_n}, \frac{l}{2\alpha_n}) \times (\hat{y}_n - \frac{1}{2\alpha_n}, \hat{y}_n + \frac{1}{2\alpha_n})$, where $l = 1 + \frac{k(\nu_1)}{k(\nu_2)}$. Consider the functions $g_n(x_1) = |\Omega_1^{h_n} \cap ((x_1, 0) + r^{(n)})|$. Since $m_n \in G_n$, there is $x_1 \leq 0$ such that $g_n(x_1) > \frac{|r^{(n)}|}{2}$. The existence of such x_1 can be seen, e.g., by the following argument (another argument uses Fubini’s theorem): Set $x_1^k = \frac{l}{2\alpha_n} - \frac{1}{2} + \frac{kl}{\alpha_n}$, $k = 0, \dots, N_n$, where $N_n = \lfloor \frac{\alpha_n}{2l} - \frac{1}{2} \rfloor$. Then $x_1^k \leq 0$ for all $k = 0, \dots, N_n$ and

$$(34) \quad S_{m_n}^{(n)} \cap \{\nabla' w_1^+ = 0\} \subset \hat{r}^{(n)} \cup \bigcup_{k=0}^{N_n} ((x_1^k, 0) + r^{(n)}),$$

where $\hat{r}^{(n)}$ is a rectangle with $|\hat{r}^{(n)}| \leq |r^{(n)}|$. If $g_n(x_1^k) \leq \frac{|r^{(n)}|}{2}$ for all $k = 0, \dots, N_n$, then $|\Omega_1^{h_n} \cap ((x_1^k, 0) + r^{(n)})| \leq \frac{|r^{(n)}|}{2}$ for all k . Intersecting (34) with $\Omega_1^{h_n}$, taking measures, and multiplying by α_n , we find $\alpha_n |S_{m_n,1}^{(n)}| \leq \alpha_n |\hat{r}^{(n)}| + \alpha_n (N_n + 1) \frac{|r^{(n)}|}{2} \leq \frac{C}{\alpha_n} + \frac{1}{4}$ since $|r^{(n)}| = \frac{l}{\alpha_n}$. As $n \rightarrow \infty$, the right-hand side converges to $\frac{1}{4}$, while the left-hand side is greater than $\frac{1}{3}$ because $m_n \in G_n$, which is a contradiction. Similarly, one proves existence of $x_1 \geq 0$ such that $g_n(x_1) < \frac{|r^{(n)}|}{2}$.

Since g_n is continuous, we conclude that there is some \hat{x}_n with $g_n(\hat{x}_n) = \frac{|r^{(n)}|}{2}$. By (33) and the choice of m_n the rectangle $\hat{S}_n = (\hat{x}_n - \frac{l}{2\alpha_n}, \hat{x}_n + \frac{l}{2\alpha_n}) \times (\hat{y}_n - \frac{1}{2\alpha_n}, \hat{y}_n + \frac{1}{2\alpha_n})$ satisfies $\limsup \alpha_n I^{h_n}(v_n; \hat{S}_n) \leq k(\nu_1)$. Recalling (25) we set $\hat{\mathcal{A}}_0^{(n)} = \{a \in \mathcal{A}_0^{h_n} : a + (-\frac{h_n}{2}, \frac{h_n}{2})^2 \subset \hat{S}_n\}$. Then by (25) we have $\#\hat{\mathcal{A}}_0^{(n)} \leq \frac{C}{h_n^3} \int_{\hat{S}_n} W(\nabla v_n) \leq \frac{C}{h_n \alpha_n}$. Hence $|\Omega_0^{h_n} \cap \hat{S}_n| \leq |\bigcup_{a \in \hat{\mathcal{A}}_0^{(n)}} (a + (-\frac{h_n}{2}, \frac{h_n}{2})^2)| + h_n \cdot \mathcal{H}^1(\partial \hat{S}_n) \leq \frac{Ch_n}{\alpha_n}$. But $\frac{Ch_n}{\alpha_n} \ll \frac{l}{\alpha_n^2} = |\hat{S}_n|$ because $H_n \rightarrow 0$. Hence $\frac{|\Omega_0^{h_n} \cap \hat{S}_n|}{|\hat{S}_n|} \rightarrow 0$, so we conclude that

$$(35) \quad \frac{|\Omega_i^{h_n} \cap \hat{S}_n|}{|\hat{S}_n|} \rightarrow \frac{1}{2} \text{ for } i = 1, 2 \text{ as } n \rightarrow \infty.$$

In fact, we have $|\hat{S}_n \cap \Omega_1^{h_n}| = \frac{|\hat{S}_n|}{2}$ by the choice of \hat{x}_n . Now set $V_n(x) = \alpha_n v_n(\frac{x}{\alpha_n} + (\hat{x}_n, \hat{y}_n))$ and $S' = (-l/2, l/2) \times (-\frac{1}{2}, \frac{1}{2})$. Then $V_n \in W^{1,2}(S' \times (-\frac{H_n}{2}, \frac{H_n}{2}); \mathbb{R}^3)$ and $\limsup I^{H_n}(V_n; S') \leq k(\nu_1)$. Set $W_n = \int_{I_{H_n}} V_n dx_3$. After possibly reflecting V_n about $\{\nu_1\}^\perp$ we may assume that

$$(36) \quad \left| \left\{ x \in S' : x_1 < 0, \text{sym } \nabla' W_n(x) \leq \frac{|B|}{4} \right\} \right| \\ \leq \left| \left\{ x \in S' : x_1 > 0, \text{sym } \nabla' W_n(x) \leq \frac{|B|}{4} \right\} \right|$$

for all n . By Theorem 6, for every subsequence there is a further subsequence, labeled with an index m , and there are affine mappings $F_m : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla F_m = 0$ such that $W_m + F_m \rightarrow W_0$ strongly in $W^{1,2}(S'; \mathbb{R}^2)$, where $W_0 \in \mathcal{A}(S')$. Moreover, we may assume that $\lim_{m \rightarrow \infty} I^{H_m}(V_m; S') = k(\nu_1)$. Theorem 12 implies that $I^0(W_0; S') \leq k(\nu_1)$. Since S' is a rectangle with sides parallel to e_1 and e_2 , Proposition 7 implies that W_0 has either only interfaces with normal $\nu_1 = e_1$ or only interfaces with normal $\nu_2 = e_2$. If it had an interface of the latter type, we would obtain the contradiction $I^0(W_0; S') \geq lk(\nu_2) > k(\nu_1)$, by the choice of l . Thus W_0 has only interfaces with normal ν_1 , and since $\limsup I^{H_n}(V_n; S') \leq k(\nu_1)$, there can be at most one such interface. On the other hand, by (35) we have $|\{\text{sym } \nabla' W_0 = 0\}| = |\{\text{sym } \nabla' W_0 = \bar{B}\}| = \frac{|S'|}{2}$, so there must be at least one interface. Hence W_0 has exactly one interface, it has normal ν_1 , and by (35) it lies in $\{\nu_1\}^\perp$. By (36) we have $(-\frac{l}{2}, 0) \times (-\frac{1}{2}, \frac{1}{2}) = \{\text{sym } \nabla' W_0 = 0\}$ up to a null set, so the mapping $w_1^+ - W_0$ is affine with $\text{sym } \nabla'(w_1^+ - W_0) = 0$. We conclude that $\text{sym } \nabla' W_m \rightarrow \text{sym } \nabla' w_1^+$ strongly in $L^2(S'; \mathbb{R}^{2 \times 2})$. Since the same limit is obtained for all subsequences, we conclude that the full sequence satisfies $\text{sym } \nabla' W_n \rightarrow \text{sym } \nabla' w_1^+$. By Korn's inequality on S' , there exist affine mappings $\tilde{F}_n : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla \tilde{F}_n = 0$ such that $W_n + \tilde{F}_n \rightarrow w_1^+$ strongly in $W^{1,2}(S'; \mathbb{R}^2)$. Denote by u_n the restriction of $x \mapsto V_n(x) + (\frac{\tilde{F}_n(x')}{0})$ to S . This sequence satisfies $\int_{I_{H_n}} u_n' dx_3 \rightarrow w_1^+$ in $W^{1,2}(S; \mathbb{R}^2)$ and $\limsup I^{H_n}(u_n; S) = k(\nu_1)$. \square

Notation. From now on we will drop the index n when dealing with sequences $h_n \rightarrow 0$ because, in view of Lemma 13, there exists a recovery sequence (v_n, h_n) for a particular sequence $h_n \rightarrow 0$ if and only if there exists one for *every* sequence $h_n \rightarrow 0$. We say that there exists a (σ, i) -recovery sequence $u^{(h)}$ on S if for all (h_n) there exist (u_n) such that (u_n, h_n) is a (σ, i) -recovery sequence on S .

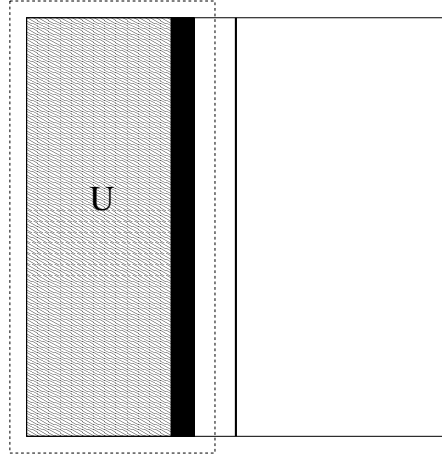


FIG. 1. Proof of Lemma 14: The black region represents the interpolation layer. The dashed rectangle denotes the boundary of V .

In a first modification step we will now change the recovery sequence furnished by Definition 8 and Lemma 13 in such a way that its vertical averages become smooth away from the interface. Lemma 14 provides a link between (1) and (2).

LEMMA 14. Let J_1, J_2 be open intervals, set $S = J_1 \times J_2$, and let $\sigma \in \{+, -\}$, $i \in \{1, 2\}$. Let $\varepsilon > 0$ and set $Q = S \cap B_\varepsilon(\mathcal{J}_{i,S})$. Then there exists a (σ, i) -recovery sequence $u^{(h)} \in W^{1,2}(S \times I_h; \mathbb{R}^3)$ on S such that $w^{(h)}(x') = \int_{I_h} (u^{(h)})'(x) dx_3$ and $\tau^{(h)}(x') = \int_{I_h} u_3^{(h)}(x) dx_3$ satisfy $w^{(h)} \in C^\infty(S \setminus \bar{Q}; \mathbb{R}^2)$ and $\tau^{(h)} \in C^\infty(S \setminus \bar{Q})$. Moreover,

$$(37) \quad \lim_{h \rightarrow 0} \left(I_{2D}^h(w^{(h)}; S \setminus Q) + h \int_{S \setminus Q} |\nabla'^2 \tau^{(h)}(x')|^2 dx' \right) = 0.$$

Proof. Arguing as at the beginning of the proof of Lemma 13, we may assume without loss of generality that $S = (-\frac{1}{2}, \frac{1}{2})^2$. Moreover, we will prove the statement for $i = 1, \sigma = +$ only, the other cases being analogous. Recall that $\{\nabla' w_1^+ = 0\} = \{x \in \mathbb{R}^2 : x_1 < 0\}$ up to a null set. Let $v^{(h)} \in W^{1,2}((2S) \times I_h; \mathbb{R}^3)$ be a $(+, 1)$ -recovery sequence on $2S$, so by Lemma 11(ii) we have $\lim_{h \rightarrow 0} I^h(v^{(h)}; S) = k(\nu_1)$. Fix $a > 0$ satisfying $(-9a, 9a) \times (-\frac{1}{2}, \frac{1}{2}) \subset Q$, set $V = (-\frac{1+a}{2}, -\frac{a}{2}) \times (-\frac{1+a}{2}, \frac{1+a}{2})$, and let $U = (-\frac{1}{2}, -a) \times (-\frac{1}{2}, \frac{1}{2})$, so $\bar{U} \subset V \subset \{\nabla' w_1^+ = 0\} \cap (2S)$. The situation is depicted in Figure 1. To obtain mappings defined on the full plate thickness, we mollify slicewise in horizontal planes: Let ψ be a standard mollifier supported on $(-\frac{1}{2}, \frac{1}{2})^2$, and set $\psi_h(x') = \frac{1}{h^2} \psi(\frac{x'}{h})$. Set

$$\tilde{v}^{(h)}(x) = (\psi_h * v^{(h)}(\cdot, x_3))(x') = \int_{I_h^2} \psi\left(\frac{y'}{h}\right) v^{(h)}(x' - y', x_3) dy',$$

which for h small enough is well defined on $S \times I_h$ (recall that $v^{(h)}$ is defined on $(2S) \times I_h$). Recall the definition of S'_h and of $R^{(h)} : S'_h \rightarrow \text{Skew} \cup (B + \text{Skew})$ introduced before (24) in the proof of Theorem 6. Adopting the notation introduced there, we have $S \subset (2S)'_h$ for h small enough, so $R^{(h)}$ is defined everywhere on S .

Since $\nabla \tilde{v}^{(h)}(x) = (\psi_h * \nabla v^{(h)}(\cdot, x_3))(x')$, we can estimate

$$\begin{aligned}
& \int_{U \times I_h} |\nabla \tilde{v}^{(h)}(x) - R^{(h)}(x')|^2 dx \\
& \leq \frac{C}{h^2} \int_{\text{spt } \psi_h} dy' \int_{U \times I_h} \left(|\nabla v^{(h)}(x' - y', x_3) - R^{(h)}(x' - y')|^2 \right. \\
& \quad \left. + |R^{(h)}(x' - y') - R^{(h)}(x')|^2 \right) dx \\
(38) \quad & \leq C \int_{V \times I_h} W(\nabla v^{(h)}(x)) dx.
\end{aligned}$$

In the first step we have applied Jensen's inequality and have added and subtracted $R^{(h)}(x' - y')$. In the last step we used that $\text{spt } \psi_h \subset (-\frac{h}{2}, \frac{h}{2})^2$ and applied the estimate

$$(39) \quad \int_U |R^{(h)}(x' + \zeta) - R^{(h)}(x')|^2 dx' \leq C \int_{V \times I_h} W(\nabla v^{(h)}(x)) dx,$$

which holds for all $\zeta \in \mathbb{R}^2$ with $|\zeta_1|, |\zeta_2| \leq h$. (The estimate (39) can be derived by arguments similar to the first part of the proof of Theorem 4.1 in [28], with our Theorem 3 replacing their Theorem 3.1.) From (38) we deduce $I^h(\tilde{v}^{(h)}; U) \leq CI^h(v^{(h)}; V)$ by (6) since $\text{sym } R^{(h)} \in \{A, B\}$. Let $w^{(h)} = \int_{I_h} (v^{(h)})' dx_3$, $\tau^{(h)} = \int_{I_h} (v^{(h)})_3 dx_3$ and $\tilde{w}^{(h)} = \int_{I_h} (\tilde{v}^{(h)})' dx_3$, $\tilde{\tau}^{(h)} = \int_{I_h} (\tilde{v}^{(h)})_3 dx_3$. Since $\text{sym } \bar{R}^{(h)} \in \{\bar{A}, \bar{B}\}$, by Jensen's inequality and by (38) we can estimate

$$\begin{aligned}
& \int_{U \times I_h} W_{2D}(\nabla' \tilde{w}^{(h)}(x')) dx \leq \int_{U \times I_h} |\text{sym } \nabla' \tilde{w}^{(h)}(x') - \text{sym } \bar{R}^{(h)}(x')|^2 dx \\
& \leq \int_{U \times I_h} |\nabla \tilde{v}^{(h)}(x) - R^{(h)}(x')|^2 dx \leq C \int_{V \times I_h} W(\nabla v^{(h)}(x)) dx.
\end{aligned}$$

Since $\nabla \tilde{v}^{(h)}(x) = (\psi_h * \nabla v^{(h)}(\cdot, x_3))(x')$, for $\alpha \in \{1, 2\}$ we have $(\nabla \tilde{v}^{(h)})_{,\alpha}(x) = \int \psi_{,\alpha}(y') (\nabla v(x' - y', x_3) - \nabla v(x', x_3)) dy'$ (where we have added a term which is zero by $\int \nabla' \psi_h = 0$). Using this together with Jensen's inequality and the fact that $|\nabla' \psi_h|^2 \leq \frac{C}{h^2}$ while $|\text{spt } \psi_h| \leq h^2$, an argument similar to (38) leads to $h \int_U |\nabla'^2 \tilde{w}^{(h)}(x')|^2 dx' + h \int_U |\nabla'^2 \tilde{\tau}^{(h)}(x')|^2 dx' \leq CI^h(v^{(h)}; V)$. Summarizing, we have shown that

$$(40) \quad I_{2D}^h(\tilde{w}^{(h)}; U) + h \int_U |\nabla'^2 \tilde{\tau}^{(h)}(x')|^2 dx' \leq CI^h(v^{(h)}; V).$$

For $\kappa \in (0, a)$ let ϕ be a smooth cutoff function that decreases from one to zero within the transition layer $(-a - \kappa, -a)$ with $\|\phi'\|_\infty \leq \frac{2}{\kappa}$. Consider the linear interpolation $u_\kappa^{(h)}(x) = v^{(h)}(x) + \phi(x_1)(\tilde{v}^{(h)}(x) - v^{(h)}(x))$. Since $\int_{I_h} (\tilde{v}^{(h)})' dx_3 \rightarrow w_1^+$ and $\int_{I_h} (v^{(h)})' dx_3 \rightarrow w_1^+$ in $L^2(U; \mathbb{R}^2)$, we also have $\int_{I_h} (u_\kappa^{(h)})' dx_3 \rightarrow w_1^+$ in $L^2(U; \mathbb{R}^2)$. Moreover, by (6) the energy $\int_{T_h} W(\nabla u_\kappa^{(h)})$ on the transition layer

$T_h = (-a - \kappa, -a) \times (-\frac{1}{2}, \frac{1}{2}) \times I_h$ is bounded by

$$\begin{aligned}
 C \int_{T_h} W_0(\nabla u_\kappa^{(h)}(x)) \, dx &\leq C \int_{T_h} W_0(\nabla v^{(h)}(x)) + \frac{1}{\kappa^2} |\tilde{v}^{(h)}(x) - v^{(h)}(x)|^2 \\
 &\quad + |\nabla \tilde{v}^{(h)}(x) - \nabla v^{(h)}(x)|^2 \, dx \\
 &\leq \frac{C}{\kappa^2} \int_{T_h} W_0(\nabla v^{(h)}(x)) + |\nabla \tilde{v}^{(h)}(x) - \nabla v^{(h)}(x)|^2 \, dx \\
 (41) \qquad \qquad \qquad &\leq \frac{C}{\kappa^2} \int_{V \times I_h} W(\nabla v^{(h)}(x)) \, dx.
 \end{aligned}$$

In passing to the second line we have assumed, by possibly adding a constant $c^{(h)}$ to $\tilde{v}^{(h)}$, that $\int_{T_h} (\tilde{v}^{(h)}(x) - v^{(h)}(x)) \, dx = 0$, so we could apply Poincaré’s inequality to estimate the term involving $|\tilde{v}^{(h)} - v^{(h)}|^2$ (the varying domain causes no problem in the application of Poincaré’s inequality; one could, e.g., apply it separately in the in-plane and in the x_3 -directions). Note that $c^{(h)} \rightarrow 0$, since $w^{(h)}$ and $\tilde{w}^{(h)}$ converge to the same limit w_1^+ in $W^{1,2}(S; \mathbb{R}^2)$. In the last step in (41) we have used the fact that by (38) and since $\int_{U \times I_h} |\nabla v^{(h)} - R^{(h)}|^2 \leq C \int_{V \times I_h} W(\nabla v^{(h)})$ (by the definition of $R^{(h)}$), we have $\int_{V \times I_h} |\nabla \tilde{v}^{(h)} - \nabla v^{(h)}|^2 \leq C \int_{V \times I_h} W(\nabla v^{(h)})$. After applying the analogous construction to the right of the interface (adding a different constant $c_2^{(h)}$ to the corresponding $\tilde{v}^{(h)}$), the lemma follows because by Lemma 11(i) we have $\frac{1}{h^2} \int_{V \times I_h} W(\nabla v^{(h)}) \rightarrow 0$ as $h \rightarrow 0$. \square

In the next lemma we further modify the recovery sequence such that the resulting functions are affine away from the interface. This is achieved via a two-step interpolation depicted in Figure 2. In the first step the recovery sequence is modified in such a way that it uses only one well away from the interface—namely, the one

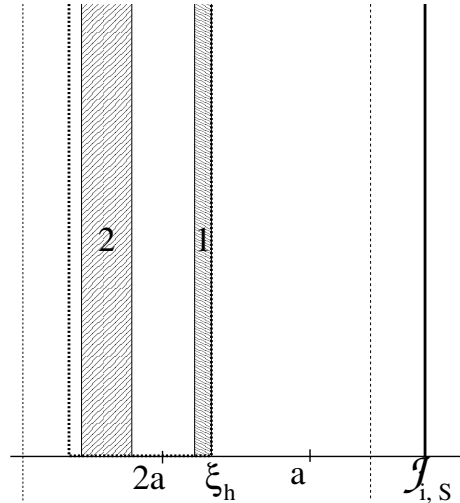


FIG. 2. The shaded regions represent the interpolation layers whose numbers correspond to the steps in Lemma 15. The bold dashed lines belong to ∂U_h , the thin dashed lines to ∂U , and the solid horizontal line to ∂S .

which is being used by the limiting mapping on that region. In a second step, it is further modified to become affine with gradient in the corresponding well.

LEMMA 15. *Let J_1, J_2 be open intervals, set $S = J_1 \times J_2$, and let $\sigma \in \{+, -\}$, $i \in \{1, 2\}$. Then there is a (σ, i) -recovery sequence $v^{(h)}$ on S with the following property: For any $\varepsilon > 0$ there is $h_0 > 0$ such that, writing $Q = S \cap B_\varepsilon(\mathcal{J}_{i,S})$, the following holds: For $h \in (0, h_0)$ the mapping $v^{(h)}$ is affine on each connected component of $(S \setminus Q) \times I_h$, and $\nabla v^{(h)} \in \text{Skew} \cup (B + \text{Skew})$ on $(S \setminus Q) \times I_h$.*

Proof. Arguing as at the beginning of the proof of Lemma 13, we may assume without loss of generality that $S = (-\frac{1}{2}, \frac{1}{2})^2$. We perform the construction only for $\sigma = +$ and $i = 1$ and only on the left side of the interface; the other cases are similar.

Set $\hat{S} = (-1, 1)^2$ and fix $a \in (0, \frac{1}{100})$. Let $v^{(h)}$ be a $(+, 1)$ -recovery sequence on \hat{S} as provided by Lemma 14, whose vertical averages $w^{(h)}(x') = \int_{I_h} (v^{(h)}(x))' dx_3$ and $\tau^{(h)}(x') = \int_{I_h} v_3^{(h)}(x) dx_3$ satisfy $w^{(h)} \in C^\infty(\hat{S} \setminus \{|x_1| < \frac{a}{10}\}; \mathbb{R}^2)$ and $\tau^{(h)} \in C^\infty(\hat{S} \setminus \{|x_1| < \frac{a}{10}\})$. Set $U = (-5a, -\frac{a}{2}) \times (-1, 1)$. We may assume without loss of generality that $h < \frac{a}{100}$, and by Lemma 11(i) we have $I^h(v^{(h)}; U) \rightarrow 0$ as $h \rightarrow 0$. By the strong convergence $w^{(h)} \rightarrow w_1^+$ in $W^{1,2}(\hat{S}; \mathbb{R}^2)$ and since $U \subset \{\nabla' w_1^+ = 0\}$ we have $\int_U |\nabla' w^{(h)}|^2 \rightarrow 0$. Hence, using (37) in the conclusion of Lemma 14, we can write

$$(42) \quad I^h(v^{(h)}; U) + I_{2D}^h(w^{(h)}; U) + h \int_U |\nabla'^2 \tau^{(h)}(x')|^2 dx' + \int_U |\nabla' w^{(h)}(x')|^2 dx' = \eta_h,$$

where $\eta_h \rightarrow 0$ as $h \rightarrow 0$.

Step 1. Interpolation to a displacement with low one-well energy. As in the proof of Theorem 6 set $S_{z,h} = z + (-\frac{h}{2}, \frac{h}{2})^2$ and $\hat{S}'_h = \bigcup_{\{z \in h\mathbb{Z}^2: S_{z,h} \subset \hat{S}\}} S_{z,h}$, and define the mapping $R^{(h)} : S'_h \rightarrow K$ to be constant on each $S_{z,h}$ with $z \in h\mathbb{Z}^2$ and such that $\int_{S_{z,h} \times I_h} |\nabla v^{(h)} - R^{(h)}|^2 \leq C \int_{S_{z,h} \times I_h} \text{dist}^2(\nabla v^{(h)}, K)$. (Here C is a universal constant given by applying Theorem 3 to a cube.) Let $G^{(h)} = h\mathbb{Z} \cap (-2a, -a)$ and set $N_h = \#G^{(h)}$. To every $\xi \in G^{(h)}$ define the column $Z_h(\xi) = ((\xi - \frac{h}{2}, \xi + \frac{h}{2}) \times (-1, 1)) \cap \hat{S}'_h$. By definition of U and since $h < \frac{a}{2}$, we have $Z_h(\xi) \subset U$ for all $\xi \in G^{(h)}$. Hence

$$(43) \quad \sum_{\xi \in G^{(h)}} \int_{Z_h(\xi)} |\nabla' w^{(h)}(x')|^2 dx' \leq \int_U |\nabla' w^{(h)}(x')|^2 dx' \leq \eta_h.$$

Now fix $\rho \in (0, 1)$ such that $(1 - 4\rho)^2 > \frac{2}{3}$ and denote by $G_1^{(h)}$ the set of all $\xi \in G^{(h)}$ with the property that

$$(44) \quad \int_{Z_h(\xi)} |\nabla' w^{(h)}(x')|^2 dx' \leq \frac{\eta_h}{[\rho N_h]}.$$

The estimate (43) implies that $\#G_1^{(h)} \geq (1 - \rho)N_h$. Notice that by (44) and since $N_h \geq \frac{a}{2h}$, for $\xi \in G_1^{(h)}$ we have

$$(45) \quad \int_{Z_h(\xi)} |\nabla' w^{(h)}(x')|^2 dx' \leq C\eta_h h.$$

Now set $L^{(h)} = \text{sym } R^{(h)}$ and define \mathcal{A}_i^h as in (25)–(26) and Ω_i^h as in (27), $i = 0, 1, 2$. By (42) and (25) we have $h \cdot \#(\mathcal{A}_0^h \cap U) \rightarrow 0$ as $h \rightarrow 0$. Hence the set $G_2^{(h)}$ of all $\xi \in G^{(h)}$ with the property that $Z_h(\xi) \cap \mathcal{A}_0^h = \emptyset$ satisfies $\frac{\#G_2^{(h)}}{N_h} \geq \frac{N_h - \#(\mathcal{A}_0^h \cap U)}{N_h} \rightarrow 1$

as $h \rightarrow 0$ because $hN_h \geq \frac{a}{2} > 0$. On the other hand, by the L^1 -convergence $\chi_{\Omega_1^h} \rightarrow \chi_{\hat{S} \cap \{\nabla' w_1^+ = 0\}}$, the set $G_3^{(h)}$ of all $\xi \in G^{(h)}$ with the property that $Z_h(\xi) \cap \mathcal{A}_1^h \neq \emptyset$ satisfies $\frac{\#G_3^{(h)}}{N_h} \rightarrow 1$ as well. By these two cardinality estimates, for h small enough we have $\#(G_2^{(h)} \cap G_3^{(h)}) \geq (1 - \rho)N_h$, and since also $\#G_1^{(h)} \geq (1 - \rho)N_h$, we conclude that $\#(G_1^{(h)} \cap G_2^{(h)} \cap G_3^{(h)}) \geq (1 - 4\rho)N_h$. On the other hand, $N_h \geq \frac{a}{h} - 2 \geq (1 - 4\rho)\frac{a}{h}$ for small h . Hence by the choice of ρ we conclude

$$(46) \quad \#(G_1^{(h)} \cap G_2^{(h)} \cap G_3^{(h)}) > \frac{2a}{3h}$$

for small h . From (28) we deduce that if $\xi \in G_2^{(h)} \cap G_3^{(h)}$, then $Z_h(\xi) \subset \Omega_1^h$. Hence using that by definition $\text{sym } R^{(h)} = 0$ on Ω_1^h , from (24), (42), and the definition of $R^{(h)}$ we conclude that $\xi \in G_2^{(h)} \cap G_3^{(h)}$ implies

$$(47) \quad \begin{aligned} \int_{Z_h(\xi) \times I_h} |\text{sym } \nabla v^{(h)}(x)|^2 dx &= \int_{Z_h(\xi) \times I_h} |\text{sym } \nabla v^{(h)}(x) - \text{sym } R^{(h)}(\xi, x_2)|^2 dx \\ &\leq Ch^2 I^h(v^{(h)}; U) \leq C\eta_h h^2. \end{aligned}$$

Let $\tilde{J}_2^{(h)}$ be the set of all $\xi \in (-2a, -a)$ satisfying the property (P_h) (defined in the statement of Lemma 18 in the appendix) for $w^{(h)}$. Applying Lemma 18 on the domain $(-2a, -a) \times (-1, 1)$, for h small enough, we have $\mathcal{H}^1(\tilde{J}_2^{(h)}) \geq \frac{a}{2}$. For every $\xi \in \tilde{J}_2^{(h)}$ there is $\tilde{w}^{(h)} \in W^{1,2}((\xi - a, \xi) \times (-\frac{1}{2}, \frac{1}{2}); \mathbb{R}^2)$ satisfying $w^{(h)}(\xi, \cdot) = \tilde{w}^{(h)}(\xi, \cdot)$ in the trace sense and

$$(48) \quad \begin{aligned} \frac{1}{h} \int_{(\xi-a, \xi) \times (-\frac{1}{2}, \frac{1}{2})} |\text{sym } \nabla' \tilde{w}^{(h)}(x')|^2 dx' &\leq C \left(I_{2D}^h(w^{(h)}; U) + \int_U |\nabla' w^{(h)}(x')|^2 dx' \right) \\ &\leq C\eta_h \end{aligned}$$

with C independent of ξ and h . We restrict the further construction to the domain of interest $S = (-\frac{1}{2}, \frac{1}{2})^2$. For $\xi \in (-2a, -a)$ define the columns

$$(49) \quad Z'_h(\xi) = \left(\xi - \frac{h}{8}, \xi \right) \times \left(-\frac{1}{2}, \frac{1}{2} \right)$$

and consider $\tilde{J}_1^{(h)} = \bigcup_{\xi \in \bigcap_{i=1}^3 G_i^{(h)}} (\xi - \frac{3h}{8}, \xi + \frac{3h}{8})$. By (46) we have $\mathcal{H}^1(\tilde{J}_1^{(h)} \cap (-2a, -a)) = \frac{3}{4}h \cdot (\# \bigcap_{i=1}^3 G_i^{(h)}) > \frac{a}{2}$. Since $\tilde{J}_2^{(h)} \subset (-2a, -a)$ and $\mathcal{H}^1(\tilde{J}_2^{(h)}) \geq \frac{a}{2}$ we conclude that there is $\xi_h \in \tilde{J}_1^{(h)} \cap \tilde{J}_2^{(h)}$ with $\xi_h \in (-2a, -a)$. Note that, in general, $\xi_h \notin G^{(h)}$ but that $Z'_h(\xi_h) \subset Z_h(\xi)$ for some $\xi \in \bigcap_{i=1}^3 G_i^{(h)}$ by the definitions of $\tilde{J}_1^{(h)}$ and $Z'_h(\xi)$. Let us introduce the set $U_h = (\xi_h - a, \xi_h) \times (-\frac{1}{2}, \frac{1}{2})$, which satisfies $U_h \subset (-3a, -a) \times (-\frac{1}{2}, \frac{1}{2}) \subset U$ for all h , since $\xi_h \in (-2a, -a)$. Setting

$$(50) \quad W_h = \text{skew} \int_{U_h} \nabla' \tilde{w}^{(h)}(x') dx',$$

we apply Korn's inequality in the plane to deduce that there is an affine mapping $f^{(h)}$ with $\nabla f^{(h)} = W_h$ such that

$$(51) \quad \int_{U_h} |\tilde{w}^{(h)} - f^{(h)}|^2 + |\nabla' \tilde{w}^{(h)} - W_h|^2 dx' \leq C \int_{U_h} |\text{sym } \nabla' \tilde{w}^{(h)}|^2 dx' \leq Ch\eta_h,$$

where the last estimate holds by (48). Notice that C in (51) is independent of h because the constant appearing in Korn's inequality is invariant under translation of the domain. We claim that

$$(52) \quad W_h \rightarrow 0 \text{ in } \mathbb{R}^{2 \times 2}.$$

Indeed, consider any subsequence. By the trace inequality and the fact that \tilde{w} and w agree on $\{x \in S : x_1 = \xi_h\}$ we have

$$\begin{aligned} \int_{-\frac{1}{2}}^{\frac{1}{2}} |f^{(h)}(\xi_h, x_2)|^2 dx_2 &\leq C \int_{-\frac{1}{2}}^{\frac{1}{2}} |f^{(h)}(\xi_h, x_2) - \tilde{w}^{(h)}(\xi_h, x_2)|^2 dx_2 \\ &\quad + C \int_{-\frac{1}{2}}^{\frac{1}{2}} |w^{(h)}(\xi_h, x_2)|^2 dx_2 \\ &\leq C \int_{U_h} |f^{(h)}(x') - \tilde{w}^{(h)}(x')|^2 + |W_h - \nabla' \tilde{w}^{(h)}(x')|^2 \\ &\quad + |w^{(h)}(x')|^2 + |\nabla' w^{(h)}(x')|^2 dx', \end{aligned}$$

which tends to zero by (51) and since $w^{(h)} \rightarrow w_1^+$ in $W^{1,2}(S; \mathbb{R}^2)$. Hence (after passing to a subsequence) $f^{(h)}(\xi_h, x_2) \rightarrow 0$ for all $x_2 \in (-\frac{1}{2}, \frac{1}{2})$. From this and using that $\text{sym } W_h = 0$ we deduce (52).

Now we extend $\tilde{w}^{(h)}$ to a three dimensional displacement $\tilde{v}^{(h)}$ by defining

$$(53) \quad \tilde{v}^{(h)}(x) = \begin{pmatrix} \tilde{w}^{(h)}(x') \\ \tau^{(h)}(x') \end{pmatrix} + x_3 \begin{pmatrix} -\tau_{,1}^{(h)}(x') \\ -\tau_{,2}^{(h)}(x') \\ 0 \end{pmatrix}.$$

By (42) and (48) we have

$$(54) \quad \begin{aligned} \int_{U_h \times I_h} |\text{sym } \nabla \tilde{v}^{(h)}(x)|^2 dx &\leq h^3 \int_{U_h} |\nabla'^2 \tau^{(h)}(x')|^2 dx' + h \int_{U_h} |\text{sym } \nabla' \tilde{w}^{(h)}(x')|^2 dx' \\ &\leq C \eta_h h^2. \end{aligned}$$

(Later we will repeat this construction on the other side of the interface. Then one must replace the second summand in (53) by $x_3(\mu_2 - \tau_{,1}^{(h)}(x'), \mu_1 - \tau_{,2}^{(h)}(x'), \mu_3)^T$.) Now consider the interpolation

$$(55) \quad u^{(h)}(x) = v^{(h)}(x) + \phi^{(h)}(x_1)(\tilde{v}^{(h)}(x) - v^{(h)}(x)),$$

where $\phi^{(h)} : \mathbb{R} \rightarrow [0, 1]$ denotes a smooth cutoff function that decreases from one to zero within the interval $(\xi_h - \frac{h}{8}, \xi_h)$, so $\text{spt } (\phi^{(h)})' \subset (\xi_h - \frac{h}{8}, \xi_h)$. We claim that

$$(56) \quad \int_{U_h \times I_h} |\text{sym } \nabla u^{(h)}(x)|^2 dx \leq C \tilde{\eta}_h h^2,$$

where $\tilde{\eta}_h = \eta_h + |W_h|^2$ converges to zero as $h \rightarrow 0$.

To prove (56), recall that by the definition of U_h and (49) we have $Z'_h(\xi_h) \subset U_h$. Now notice that $u^{(h)} = \tilde{v}^{(h)}$ on $(U_h \times I_h) \setminus (Z'_h(\xi_h) \times I_h)$, whence by (54) we have

$\int_{(U_h \setminus Z'_h(\xi_h)) \times I_h} |\text{sym } \nabla u^{(h)}|^2 \leq C\eta_h h^2$. It remains to prove (56) with the integration domain $Z'_h(\xi_h) \times I_h$ replacing $U_h \times I_h$. We make a standard calculation to obtain

$$(57) \quad \int_{Z'_h(\xi_h) \times I_h} |\text{sym } \nabla u^{(h)}(x)|^2 dx \leq C \int_{Z'_h(\xi_h) \times I_h} |\text{sym } \nabla v^{(h)}(x)|^2 + |\text{sym } \nabla \tilde{v}^{(h)}(x)|^2 + \frac{1}{h^2} |\tilde{v}^{(h)}(x) - v^{(h)}(x)|^2 dx.$$

Since $\xi_h \in G_2^{(h)} \cap G_3^{(h)}$, the first term on the right-hand side is estimated by (47). By (54) the second term in (57) satisfies $\int_{Z'_h(\xi_h) \times I_h} |\text{sym } \nabla \tilde{v}(x)|^2 dx \leq C\eta_h h^2$. Let us estimate the third term in (57). Since

$$\begin{pmatrix} w^{(h)}(x') \\ \tau^{(h)}(x') \end{pmatrix} = \int_{I_h} v^{(h)}(x) dx_3 = \int_{I_h} v^{(h)}(x) + x_3 (\nabla' \tau^{(h)}(x'))^T dx_3$$

and since $w^{(h)} = \tilde{w}^{(h)}$ on the line $\{x \in S : x_1 = \xi_h\}$, we can apply a Poincaré inequality (see, e.g., [16, Theorem 6.1-8]) to estimate the second term in the last step in (58) below. The first term in that step is estimated by the usual Poincaré inequality in the x_3 -direction:

$$(58) \quad \begin{aligned} & \int_{Z'_h(\xi_h) \times I_h} \frac{1}{h^2} |v^{(h)}(x) - \tilde{v}^{(h)}(x)|^2 dx \\ & \leq \frac{C}{h^2} \int_{Z'_h(\xi_h) \times I_h} \left| v^{(h)}(x) + x_3 \begin{pmatrix} \tau_{,1}^{(h)}(x') \\ \tau_{,2}^{(h)}(x') \\ 0 \end{pmatrix} - \begin{pmatrix} w^{(h)}(x') \\ \tau^{(h)}(x') \end{pmatrix} \right|^2 + |w^{(h)}(x') - \tilde{w}^{(h)}(x')|^2 dx \\ & \leq C \int_{Z'_h(\xi_h) \times I_h} \left| v_{,3}^{(h)}(x) + \begin{pmatrix} \tau_{,1}^{(h)}(x') \\ \tau_{,2}^{(h)}(x') \\ 0 \end{pmatrix} \right|^2 + |\nabla' w^{(h)}(x') - \nabla' \tilde{w}^{(h)}(x')|^2 dx. \end{aligned}$$

To estimate the first term in (58), we observe that $(\tau_{,1}^{(h)}(x'), \tau_{,2}^{(h)}(x'), 0) = \nabla \int_{I_h} v_3^{(h)}(x) dx_3$, so we have

$$\begin{aligned} \int_{Z'_h(\xi_h) \times I_h} \left| v_{,3}^{(h)}(x) + \begin{pmatrix} \tau_{,1}^{(h)}(x') \\ \tau_{,2}^{(h)}(x') \\ 0 \end{pmatrix} \right|^2 dx & \leq C \int_{Z'_h(\xi_h) \times I_h} \left| v_{,3}^{(h)}(x) - \int_{I_h} v_3^{(h)}(x', z) dz \right|^2 \\ & \quad + \left| \int_{I_h} v_3^{(h)}(x', z) dz + \left(\nabla \int_{I_h} v_3^{(h)}(x', z) dz \right)^T \right|^2 dx. \end{aligned}$$

The second term is bounded by $\int_{Z'_h(\xi_h) \times I_h} |\text{sym } \nabla v^{(h)}|^2$ and can therefore be estimated by (47). By the x_3 -independence of $R^{(h)}$, by Jensen's inequality, by definition of $R^{(h)}$,

and by (42), the first term can be estimated as follows:

$$\begin{aligned} & \int_{Z'_h(\xi_h) \times I_h} \left| v_{,3}^{(h)}(x) - \int_{I_h} v_{,3}^{(h)}(x', z) dz \right|^2 dx \\ & \leq C \int_{Z'_h(\xi_h) \times I_h} |v_{,3}^{(h)}(x) - R_3^{(h)}(x')|^2 + \left| \int_{I_h} R_3^{(h)}(x') - v_{,3}^{(h)}(x', z) dz \right|^2 dx \\ & \leq C \eta_h h^2. \end{aligned}$$

Finally, the second term in (58) is bounded by

$$\begin{aligned} & C \int_{Z'_h(\xi_h) \times I_h} |\nabla' w^{(h)}(x')|^2 + |\nabla' \tilde{w}^{(h)}(x') - W_h|^2 + |W_h|^2 dx \\ & \leq C(\eta_h h^2 + |Z'_h(\xi_h) \times I_h| |W_h|^2). \end{aligned}$$

We have applied (52) and (45) multiplied by h (recall that $\xi_h \in G_1^{(h)}$), since here we are integrating over the thickness on the left-hand side. This proves (56) and finishes the first interpolation step.

Step 2. Interpolation to an affine displacement. We apply Lemma 16 to the mapping $u^{(h)}$ defined in (55) with $J_h = (\xi_h - a, \xi_h)$ instead of J_1 , (so $J_h \times (-\frac{1}{2}, \frac{1}{2}) = U_h$) and with $t = \xi_h - \frac{a}{2}$ and $b = \frac{a}{4}$. Lemma 16 furnishes a mapping $\tilde{u}^{(h)}$ which agrees with $u^{(h)}$ on $\{x \in S \times I_h : x_1 > \xi_h - \frac{a}{2}\}$ and on $\{x \in S \times I_h : x_1 < \xi_h - \frac{3a}{4}\}$ agrees with an affine function $f^{(h)}$ with $\text{sym } \nabla f^{(h)} = 0$ (the mapping $\tilde{u}^{(h)}$ is at first not defined on $\{x_1 < \xi_h - a\}$, but since it is affine on $\{x \in S \times I_h : x_1 \in (\xi_h - a, \xi_h - \frac{3a}{4})\}$ we can extend it affinely). Moreover, $\tilde{u}^{(h)}$ satisfies $\int_{U_h \times I_h} |\text{sym } \nabla \tilde{u}^{(h)}|^2 \leq \frac{C}{a^4} \int_{U_h \times I_h} |\text{sym } \nabla u^{(h)}|^2$. Combining this with (56) and with the fact that $\tilde{u}^{(h)} = v^{(h)}$ on $\{x \in S \times I_h : x_1 > \xi_h\}$ and $\text{sym } \nabla \tilde{u}^{(h)} = 0$ on $\{x \in S \times I_h : x_1 < \xi_h - \frac{3a}{4}\}$, we conclude that

$$(59) \quad I^h(\tilde{u}^{(h)}; S) \rightarrow k(\nu_1).$$

Step 3. Convergence. Now we apply Steps 1 and 2 with obvious modifications also on the other side of the interface. Let us denote the resulting mappings by $\tilde{u}_a^{(h)}$. Hence $\tilde{u}_a^{(h)}$ is a $(+, 1)$ -recovery sequence on S which is affine on $\{x \in S \times I_h : |x_1| > 5a\}$ with $\text{sym } \nabla \tilde{u}_a^{(h)} = 0$ on $\{x \in S \times I_h : x_1 < -5a\}$ and $\text{sym } \nabla \tilde{u}_a^{(h)} = B$ on $\{x \in S \times I_h : x_1 > 5a\}$. By Proposition 19 there exists a sequence $a^{(h)} \rightarrow 0$ such that

$$(60) \quad \limsup_{h \rightarrow 0} I^h(\tilde{u}_{a^{(h)}}^{(h)}; S) = \limsup_{a \rightarrow 0} \limsup_{h \rightarrow 0} I^h(\tilde{u}_a^{(h)}; S) = k(\nu_1).$$

Theorem 6 implies that there exist affine mappings $f^{(h)}$ with $\text{sym } \nabla f^{(h)} = 0$ and $w \in \mathcal{A}(S)$ such that, after passing to an unlabeled subsequence,

$$(61) \quad \bar{w}^{(h)} + f^{(h)} \rightarrow w \text{ strongly in } W^{1,2}(S; \mathbb{R}^2),$$

where we have set $\bar{w}^{(h)}(x') = \int_{I_h} (\tilde{u}_{a^{(h)}}^{(h)}(x', x_3))' dx_3$. By Theorem 12 and (60) the limiting function w satisfies $I^0(w; S) \leq k(\nu_1)$. But by (61) and the properties of $\tilde{u}_{a^{(h)}}^{(h)}$ necessarily $\text{sym } \nabla' w = 0$ on $\{x \in S : x_1 < 0\}$ and $\text{sym } \nabla' w = \bar{B}$ on $\{x \in S : x_1 > 0\}$. Hence after possibly adding $w_1^+ - w$ to all $f^{(h)}$ (notice that $w_1^+ - w$ is affine with $\text{sym } \nabla'(w_1^+ - w) = 0$) we may assume that (61) holds with $w = w_1^+$ on the right-hand side. Since the same limit is obtained for every subsequence, (61) holds for the full

sequence with $w = w_1^+$, whence $x \mapsto \tilde{u}_{a^{(h)}}^{(h)}(x) + (f^{(h)})'_0(x')$ is the sought recovery sequence. \square

LEMMA 16. *Let J_1, J_2 be open intervals, set $U = J_1 \times J_2$, let $i \in \{1, 2\}$, and let $D \in \mathbb{R}^{3 \times 3}$ be a symmetric matrix. Then there is a constant $C > 0$ such that the following holds: For every $h \in (0, 1)$, for every $u^{(h)} \in W^{1,2}(U \times I_h; \mathbb{R}^3)$, for every $t \in J_i$, and for every $b \in (0, 1)$ satisfying $t + b \in J_i$ (resp., $t - b \in J_i$), there exist $c^{(h)} \in \mathbb{R}^3$, $T^{(h)} \in \mathbb{R}^{3 \times 3}$ with $\text{sym } T^{(h)} = 0$, and $\tilde{u}^{(h)} \in W^{1,2}(U \times I_h; \mathbb{R}^3)$ such that $\tilde{u}^{(h)} = u^{(h)}$ on $\{x \in U \times I_h : x_i < t\}$ (resp., $\{x \in U \times I_h : x_i > t\}$) and $\tilde{u}^{(h)} = Dx + T^{(h)}x + c^{(h)}$ on $\{x \in U \times I_h : x_i > t + b\}$ (resp., $\{x \in U \times I_h : x_i < t - b\}$) and such that $\int_{U \times I_h} |\text{sym } \nabla \tilde{u}^{(h)} - D|^2 \leq \frac{C}{b^4} \int_{U \times I_h} |\text{sym } \nabla u^{(h)} - D|^2$. One can take*

$$(62) \quad T^{(h)} = \text{skew} \int_{U \times I_h} \nabla u^{(h)}(x) \, dx \text{ and } c^{(h)} = \int_{U \times I_h} u^{(h)}(x) \, dx.$$

Proof. We may assume without loss of generality that $D = 0$. (In fact, if $D \neq 0$, then apply the lemma with $D = 0$ to $\hat{u}^{(h)}(x) = u^{(h)}(x) - Dx$ instead of $u^{(h)}$ to obtain $\hat{u}^{(h)}$. Then define $\tilde{u}^{(h)}(x) = \hat{u}^{(h)}(x) + Dx$.) Moreover, we prove only the case when $i = 1$ and $t + b \in J_1$. By Proposition 17(ii) and Poincaré’s inequality the mappings $f^{(h)}(x) = T^{(h)}x + c^{(h)}$ with $T^{(h)}$ and $c^{(h)}$ as in (62) satisfy

$$(63) \quad \int_{U \times I_h} |u^{(h)}(x) - f^{(h)}|^2 + |\nabla u^{(h)}(x) - T^{(h)}|^2 \, dx \leq \frac{C}{h^2} \int_{U \times I_h} |\text{sym } \nabla u^{(h)}(x)|^2 \, dx.$$

Also, by Proposition 17(iii)

$$(64) \quad \int_{U \times I_h} |(u^{(h)})'(x) - (f^{(h)})'(x)|^2 \, dx \leq C \int_{U \times I_h} |\text{sym } \nabla u^{(h)}(x)|^2 \, dx.$$

Fix a smooth cutoff function $\phi(x_1)$ which decreases from one to zero within the interval $(t + \frac{b}{4}, t + \frac{3b}{4})$. Set

$$\tilde{u}^{(h)}(x) = f^{(h)}(x) + \phi(x_1)(u^{(h)}(x) - f^{(h)}(x)) - x_3 \phi'(x_1)(u_3^{(h)}(x) - f_3^{(h)}(x))e_1.$$

Then

$$\begin{aligned} \nabla \tilde{u}^{(h)}(x) &= T^{(h)} + \phi(x_1)(\nabla u^{(h)}(x) - T^{(h)}) + \begin{pmatrix} (u^{(h)} - f^{(h)})'(x) \\ 0 \end{pmatrix} \otimes e_1 \phi'(x_1) \\ &+ (u^{(h)} - f^{(h)})_3(x) e_3 \otimes e_1 \phi'(x_1) - (u^{(h)} - f^{(h)})_3(x) e_1 \otimes e_3 \phi'(x_1) \\ &- x_3 \left(\phi''(x_1)(u^{(h)} - f^{(h)})_3(x) e_1 \otimes e_1 + \phi'(x_1) e_1 \otimes (\nabla u_3^{(h)}(x) - (T^{(h)})^T e_3) \right). \end{aligned}$$

Upon taking the symmetric part of the above expression, the second line cancels, so we obtain

$$\begin{aligned} \int_{U \times I_h} |\text{sym } \nabla \tilde{u}^{(h)}|^2 \, dx &\leq C \int_{U \times I_h} |\text{sym } \nabla u^{(h)}|^2 + \frac{1}{b^2} |(u^{(h)} - f^{(h)})'|^2 \, dx \\ &+ h^2 \int_{U \times I_h} \left(\frac{1}{b^4} |u^{(h)} - f^{(h)}|^2 + \frac{1}{b^2} |\nabla u^{(h)} - T^{(h)}|^2 \right) \, dx, \end{aligned}$$

since $|\phi'| \leq \frac{C}{b}$ and $|\phi''| \leq \frac{C}{b^2}$. The last term is controlled by (63) and the $(u^{(h)} - f^{(h)})'$ -term is controlled by (64). \square

Proof of Theorem 1. By Lemma 2 we must prove the theorem only for the special case $A = 0$ and B as in (7). Statement (i) just rephrases the content of Theorem 12. The proof of (ii) is similar to that of Proposition 5.1 in [19]; compare also Step 2 in the proof of Theorem 5.6 in [18]. If $I^0(w) = \infty$, then the proof is trivial. Otherwise, $w \in \mathcal{A}(S)$, so by Proposition 7, w is affine on each connected component of $S \setminus \bigcup_{i=1}^M \mathcal{J}_i$, where the interfaces $\mathcal{J}_i \subset S$, $i = 1, \dots, M$ ($M \in \mathbb{N} \cup \{\infty\}$), are straight line segments parallel to e_1 or e_2 which intersect ∂S at their ends. By bounded variation we have $\sum_{i=1}^M \mathcal{H}^1(\mathcal{J}_i) < \infty$, so if $M = \infty$, then $\mathcal{H}^1(\mathcal{J}_i) \rightarrow 0$ as $i \rightarrow \infty$, so the \mathcal{J}_i can accumulate only at ∂S . By translation invariance we may assume without loss of generality that S is strictly star-shaped with respect to the origin, so $\bar{S} \subset \eta S$ for any $\eta > 1$. Hence the restriction of $w_\eta(x) = \eta w(\frac{x}{\eta})$ to S is well defined and satisfies $w_\eta|_S \in \mathcal{A}(S)$. The mapping $w_\eta|_S$ has only finitely many interfaces \mathcal{J}_i^η , $i = 1, \dots, M_\eta$ (recall that the \mathcal{J}_i accumulate only near ∂S), and they satisfy $\text{dist}(\mathcal{J}_i^\eta, \mathcal{J}_k^\eta) > 0$ whenever $k \neq i$. By construction we have $\mathcal{J}_i^\eta = \tilde{\mathcal{J}}_i^\eta \cap S$, where $\tilde{\mathcal{J}}_i^\eta \subset \eta S$ are the interfaces of $w_\eta \in \mathcal{A}(\eta S)$. For $i = 1, \dots, M_\eta$, let $k_i \in \{1, 2\}$ be such that ν_{k_i} is normal to \mathcal{J}_i^η . Let $\varepsilon > 0$ be small (fixed below) and define the rectangles $S_i = \{x \in \mathbb{R}^2 : \pi_{k_i}(x) \in \pi_{k_i}(\tilde{\mathcal{J}}_i^\eta) + (-\varepsilon, \varepsilon) \text{ and } \pi_{k_i}^\perp(x) \in \pi_{k_i}^\perp(\tilde{\mathcal{J}}_i^\eta)\}$. Since the endpoints of $\tilde{\mathcal{J}}_i^\eta$ lie in $\partial(\eta S)$ and since the $\tilde{\mathcal{J}}_i^\eta$ do not intersect in ηS , for ε small enough the sides of ∂S_i which are parallel to ν_{k_i} do not intersect \bar{S} and the S_i are pairwise disjoint. Define $\sigma_i \in \{-, +\}$ by the requirement that $w_\eta = w_{k_i, S_i}^{\sigma_i}$ on $S_i \cap \eta S$.

Now let $h_n \rightarrow 0$ be given and for all $i = 1, \dots, M_\eta$ let $(u_n^{(i)}, h_n)$ be a (σ_i, k_i) -recovery sequence on S_i . By Lemma 15 we may assume that $u_n^{(i)}$ is affine with $\text{sym } \nabla u_n^{(i)} \in \{A, B\}$ in a neighborhood (relatively open in $S \times I_{h_n}$) of $(S \cap \partial S_i) \times I_{h_n}$. Hence we can extend each $u_n^{(i)}$ to all of $S \times I_{h_n}$ in such a way that $u_n^{(i)}$ is affine on each connected component of $(S \setminus S_i) \times I_{h_n}$. For $i, j = 1, \dots, M_\eta$ let us write $i \sim j$ whenever there exists a connected component S'_{ij} of $S \setminus \bigcup_{i=1}^{M_\eta} \mathcal{J}_i^\eta$ satisfying $S \cap \partial S'_{ij} = \mathcal{J}_i^\eta \cup \mathcal{J}_j^\eta$ (i.e., \mathcal{J}_i^η and \mathcal{J}_j^η are neighbors). By star-shapedness and the fact that $M_\eta < \infty$, and recalling that $u_n^{(i)}$ is affine on $(S \setminus S_i) \times I_{h_n}$, one inductively finds affine mappings $f_n^{(i)} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ with $\text{sym } \nabla f_n^{(i)} = 0$, $i = 1, \dots, M_\eta$, such that $f_n^{(i)} + u_n^{(i)} = f_n^{(j)} + u_n^{(j)}$ on $S'_{ij} \setminus (S_i \cup S_j)$ whenever $i \sim j$. Now, for every $i = 1, \dots, M_\eta$, we set $v_n^\eta = u_n^{(i)} + f_n^{(i)}$ on S_i and extend it affinely to $S \times I_{h_n}$. Then $v_n^\eta \in W^{1,2}(S \times I_{h_n}; \mathbb{R}^3)$ is well defined by the choice of the $f_n^{(i)}$. Moreover, $\limsup_{n \rightarrow \infty} I^{h_n}(v_n^\eta; S) \leq I^0(w_\eta; S) + \rho_\eta$, where $\rho_\eta \downarrow 0$ as $\eta \downarrow 1$. By a diagonal sequence argument, this implies part (ii) of Theorem 1, since $w_\eta|_S \rightarrow w$ in $W^{1,2}(S; \mathbb{R}^2)$ and $I^0(w_\eta; S) \rightarrow I^0(w; S)$ as $\eta \downarrow 1$ (i.e., limiting displacements which arise as rescalings w_η of some $w \in \mathcal{A}(S)$ are energy dense; compare [2, p. 3]). \square

Appendix. The following proposition was used in the proof of Lemma 16. Statement (ii) is Korn's inequality for thin films as presented in [3]; see also Problem 1.12 in [17].

PROPOSITION 17. *Let $S \subset \mathbb{R}^2$ be a bounded Lipschitz domain and let $A, B \in \mathbb{R}^{3 \times 3}$ be such that $\text{rank}(A - B + F) \geq 2$ for all $F \in \mathbb{R}^{3 \times 3}$ with $\text{sym } F = 0$. Then there is a constant $C(S)$ such that for all $h \in (0, 1)$ and for all $v^{(h)} \in W^{1,2}(S \times I_h; \mathbb{R}^3)$ the following hold:*

(i) *There exists a matrix $\text{sym } T^{(h)} \in \{A, B\}$ such that*

$$\int_{S \times I_h} |\nabla v^{(h)}(x) - T^{(h)}|^2 dx \leq \frac{C(S)}{h^2} \int_{S \times I_h} \text{dist}^2(\text{sym } \nabla v^{(h)}(x), \{A, B\}) dx.$$

(ii) *The estimate*

$$\int_{S \times I_h} |\nabla v^{(h)}(x) - T^{(h)}|^2 dx \leq \frac{C(S)}{h^2} \int_{S \times I_h} |\text{sym } \nabla v^{(h)}(x)|^2 dx$$

holds for $T^{(h)} = \text{skew } \int_{S \times I_h} \nabla v^{(h)}$.

(iii) *There exists $c^{(h)} \in \mathbb{R}^2$ such that*

$$\int_{S \times I_h} |(v^{(h)})'(x) - (T^{(h)}x)' - c^{(h)}|^2 dx \leq C(S) \int_{S \times I_h} |\text{sym } \nabla v^{(h)}(x)|^2 dx$$

for the same $T^{(h)}$ as in (ii).

Proof. The proof of (i) is analogous to that of Theorem 10 in [29], with Theorem 3 replacing their geometric rigidity theorem. Statement (ii) can be proven in the same way, with Korn’s inequality for one well replacing their geometric rigidity theorem. Another proof is given in [3] and [33]. Notice that if (ii) holds for some skew matrix, then it will also hold for the special choice $T^{(h)} = \text{skew } \int_{S \times I_h} \nabla v^{(h)} dx$.

To prove statement (iii), set $w^{(h)}(x') = \int_{I_h} (v^{(h)})'(x', x_3) dx_3$. From Korn’s inequality in the plane and from Jensen’s inequality we obtain

$$\begin{aligned} \int_S \left| \nabla' w^{(h)}(x') - \bar{T}^{(h)} \right|^2 dx' &\leq C \int_S \left| \text{sym } \nabla' w^{(h)}(x') \right|^2 dx' \\ (65) \qquad \qquad \qquad &\leq \frac{C}{h} \int_{S \times I_h} |\text{sym } \nabla v^{(h)}(x)|^2 dx. \end{aligned}$$

With $c^{(h)} = \int_S (w^{(h)}(x') - \bar{T}^{(h)}x') dx'$ we obtain

$$\begin{aligned} &\int_{S \times I_h} \left| (v^{(h)})'(x) - (T^{(h)}x)' - c^{(h)} \right|^2 dx \\ &\leq C \int_{S \times I_h} \left| (v^{(h)})'(x) - (T_3^{(h)})'x_3 - w^{(h)}(x') \right|^2 + \left| w^{(h)}(x') - \bar{T}^{(h)}x' - c^{(h)} \right|^2 dx, \end{aligned}$$

where $T_3^{(h)}$ denotes the third column of $T^{(h)}$. The second term is estimated by applying Poincaré’s inequality on S and then (65). To estimate the first term, notice that since the integration domain is symmetric, we have $w^{(h)}(x') = \int_{I_h} (v'(x) - (T_3^{(h)})'x_3) dx_3$. Applying Poincaré’s inequality in the x_3 -direction for almost every x' and subsequently using (ii) shows that the first term is controlled by $\int_{S \times I_h} |\text{sym } \nabla v^{(h)}|^2$. \square

The following lemma is a corollary of Proposition 4.1 in [19]. Notice that their ε corresponds to our h .

LEMMA 18. *Let $a, l, d > 0$, let $U = (-l, l) \times (-d, d)$, let $\bar{A} = 0$ and $\bar{B} = e_1 \otimes e_2 + e_2 \otimes e_1$, and let $\bar{F} \in \{\bar{A}, \bar{B}\}$. Then there are constants $\eta_0, C_0 > 0$ such that for every $h \in (0, 1)$ and $w \in W^{2,2}(U; \mathbb{R}^2)$ with*

$$I_{2D}^h(w; U) \leq \eta_0 \text{ and } \int_U |\text{sym } \nabla' w(x') - \bar{F}|^2 dx' \leq \eta_0$$

the set of $\xi \in (-l, l)$ satisfying property (P_h) for w has measure not smaller than l .

We say $\xi \in (-l, l)$ satisfies property (P_h) for $w \in W^{2,2}(U; \mathbb{R}^2)$ if, setting $U_1 = (\xi - a, \xi) \times (-\frac{d}{2}, \frac{d}{2})$ and $U_2 = (\xi, \xi + a) \times (-\frac{d}{2}, \frac{d}{2})$, for each $i = 1, 2$ there exist $\tilde{w}_i \in W^{1,2}(U_i; \mathbb{R}^2)$ with $\tilde{w}(\xi, \cdot) = w(\xi, \cdot)$ on $(-d/2, d/2)$ and

$$\frac{1}{h} \int_{U_i} |\text{sym } \nabla' \tilde{w}_i(x') - \bar{F}|^2 dx' \leq C_0 \left(I_{2D}^h(w; U) + \int_U |\text{sym } \nabla' w(x') - \bar{F}|^2 dx' \right).$$

An analogous result holds for lines of the form $\{x_2 = \xi\}$.

Proof. Set $\eta = I_{2D}^h(w; U) + \int_U |\text{sym } \nabla' w - \bar{F}|^2$. By Proposition 4.1 in [19] there exists a Borel set $\Sigma \subset (-l, l)$ with $|\Sigma| \geq l$ and such that for all $\xi \in \Sigma$ there exists an affine mapping $w_\xi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $\text{sym } \nabla' w_\xi = \bar{F}$ and

$$(66) \quad \|w(\xi, \cdot) - w_\xi(\xi, \cdot)\|_{H^{1/2}((-\frac{d}{2}, \frac{d}{2}); \mathbb{R}^2)}^2 \leq Ch\eta.$$

By the properties of the $H^{1/2}$ -norm (see, e.g., the appendix of [32] for a review), for $i = 1, 2$ there exist $v_i \in W^{1,2}(U_i; \mathbb{R}^2)$ such that

$$(67) \quad \int_{U_i} |\nabla' v_i(x')|^2 dx' \leq \|w(\xi, \cdot) - w_\xi(\xi, \cdot)\|_{H^{1/2}((-\frac{d}{2}, \frac{d}{2}); \mathbb{R}^2)}^2$$

and $v_i(\xi, \cdot) = w(\xi, \cdot) - w_\xi(\xi, \cdot)$ on $(-\frac{d}{2}, \frac{d}{2})$ in the trace sense. Setting $\tilde{w}_i = v_i + w_\xi$, we find

$$\frac{1}{h} \int_{U_i} |\text{sym } \nabla' \tilde{w}_i(x') - \bar{F}|^2 dx' \leq \frac{C}{h} \int_{U_i} |\nabla' v_i(x')|^2 + |\text{sym } \nabla' w_\xi(x') - \bar{F}|^2 dx' \leq C\eta$$

and $\tilde{w}_i(\xi, \cdot) = w(\xi, \cdot)$ on $(-\frac{d}{2}, \frac{d}{2})$ in the trace sense. We have used (66)–(67) and the fact that $\text{sym } \nabla' w_\xi = \bar{F}$. \square

The following proposition is a standard diagonalization lemma (compare [4, Corollary 1.16] or [13, Lemma 7.2]).

PROPOSITION 19. *Let $a_{k,j}$ be a doubly indexed sequence of real numbers, $k, j \rightarrow \infty$. Then there exists a subsequence $k_j \rightarrow \infty$ such that*

$$\limsup_{j \rightarrow \infty} a_{k_j, j} = \limsup_{k \rightarrow \infty} \limsup_{j \rightarrow \infty} a_{k, j}.$$

Acknowledgments. This work is part of my Ph.D. thesis supervised by Prof. S. Müller (Max-Planck-Institute for Mathematics in the Sciences, Leipzig), whom I thank for his steady advice. I also wish to thank Prof. S. Conti for interesting discussions on the topic and for a lot of helpful comments on a preliminary version of this work.

REFERENCES

- [1] E. ACERBI, G. BUTTAZZO, AND D. PERCIVALE, *A variational definition of an elastic string*, J. Elasticity, 25 (1991), pp. 137–148.
- [2] G. ALBERTI, *Variational models for phase transitions: An approach via Γ -convergence*, in Calculus of Variations and Partial Differential Equations, Topics on Geometrical Evolution Problems and Degree Theory, G. Buttazzo et al., eds., Springer-Verlag, Berlin, 2000, pp. 95–114.
- [3] G. ANZELLOTTI, S. BALDO, AND D. PERCIVALE, *Dimension reduction in variational problems, asymptotic development in Γ -convergence and thin structures in elasticity*, Asymptot. Anal., 9 (1994), pp. 61–100.
- [4] H. ATTOUCH, *Variational Convergence for Functions and Operators*, Pitman, Boston, 1984.

- [5] A. C. BARROSO AND I. FONSECA, *Anisotropic singular perturbations—the vectorial case*, Proc. Roy. Soc. Edinburgh Sect. A, 124 (1994), pp. 527–571.
- [6] P. BÉLIK AND M. LUSKIN, *The Γ -convergence of a sharp interface thin film model with non-convex elastic energy*, SIAM J. Math. Anal., 38 (2006), pp. 414–433.
- [7] K. BHATTACHARYA, *Microstructure of Martensite*, Oxford Series on Materials Modelling, Oxford University Press, Oxford, UK, 2003.
- [8] K. BHATTACHARYA AND G. DOLZMANN, *Relaxation of some multi-well problems*, Proc. Roy. Soc. Edinburgh Sect. A, 131 (2001), pp. 279–320.
- [9] K. BHATTACHARYA, N. B. FIROOZY, R. D. JAMES, AND R. V. KOHN, *Restrictions on microstructure*, Proc. Roy. Soc. Edinburgh Sect. A, 124 (1994), pp. 843–878.
- [10] K. BHATTACHARYA AND R. D. JAMES, *A theory of thin films of martensitic materials with applications to microactuators*, J. Mech. Phys. Solids, 47 (1999), pp. 531–576.
- [11] K. BHATTACHARYA AND R. D. JAMES, *The material is the machine*, Science, 307 (2005), pp. 53–54.
- [12] F. BOURQUIN, P. G. CIARLET, G. GEYMONAT, AND A. RAOULT, *Γ -convergence et analyse asymptotique des plaques minces*, C. R. Acad. Sci. Paris Sér. I Math., 315 (1992), pp. 1017–1024.
- [13] A. BRAIDES, I. FONSECA, AND G. FRANCFORT, *3D-2D asymptotic analysis for inhomogeneous thin films*, Indiana Univ. Math. J., 49 (2000), pp. 1367–1403.
- [14] N. CHAUDHURI AND S. MÜLLER, *Rigidity estimate for two incompatible wells*, Calc. Var. Partial Differential Equations, 19 (2004), pp. 379–390.
- [15] N. CHAUDHURI AND S. MÜLLER, *Scaling of the energy for thin martensitic films*, SIAM J. Math. Anal., 38 (2006), pp. 468–477.
- [16] P. G. CIARLET, *Mathematical Elasticity Vol. I*, Studies in Mathematics and Its Applications 20, North-Holland, Amsterdam, 1988.
- [17] P. G. CIARLET, *Mathematical Elasticity Vol. II*, Studies in Mathematics and Its Applications 27, North-Holland, Amsterdam, 1997.
- [18] S. CONTI, I. FONSECA, AND G. LEONI, *A Γ -convergence result for the two-gradient theory of phase transitions*, Comm. Pure Appl. Math., 55 (2002), pp. 857–936.
- [19] S. CONTI AND B. SCHWEIZER, *A sharp-interface limit for a two-well problem*, Arch. Ration. Mech. Anal., 179 (2006), pp. 413–452.
- [20] S. CONTI AND B. SCHWEIZER, *Rigidity and Gamma convergence for solid-solid phase transitions with $SO(2)$ -invariance*, Comm. Pure Appl. Math., 59 (2006), pp. 830–868.
- [21] C. DELELLIS AND L. SZÉKELYHIDI, JR., *Simple proof of two-well rigidity*, C. R. Math. Acad. Sci. Paris, 343 (2006), pp. 367–370.
- [22] A. DE SIMONE AND G. FRIESECKE, *On the problem of two linearized wells*, Calc. Var. Partial Differential Equations, 4 (1996), pp. 293–304.
- [23] G. DOLZMANN AND S. MÜLLER, *Microstructures with finite surface energy: The two-well problem*, Arch. Ration. Mech. Anal., 132 (1995), pp. 101–141.
- [24] L. C. EVANS AND R. F. GARIEPY, *Measure Theory and Fine Properties of Functions*, CRC Press, Boca Raton, FL, 1992.
- [25] I. FONSECA, *Phase transitions of elastic solid materials*, Arch. Ration. Mech. Anal., 107 (1989), pp. 195–223.
- [26] I. FONSECA AND C. MANTEGAZZA, *Second order singular perturbation models for phase transitions*, SIAM J. Math. Anal., 31 (2000), pp. 1121–1143.
- [27] I. FONSECA AND L. TARTAR, *The gradient theory of phase transitions for systems with two potential wells*, Proc. Roy. Soc. Edinburgh Sect. A, 111 (1989), pp. 89–102.
- [28] G. FRIESECKE, R. D. JAMES, AND S. MÜLLER, *A theorem on geometric rigidity and the derivation of nonlinear plate theory from three dimensional elasticity*, Comm. Pure Appl. Math., 55 (2002), pp. 1461–1506.
- [29] G. FRIESECKE, R. D. JAMES, AND S. MÜLLER, *A hierarchy of plate models derived from nonlinear elasticity by Gamma-convergence*, Arch. Ration. Mech. Anal., 180 (2006), pp. 183–236.
- [30] M. GIAQUINTA, *Multiple Integrals in the Calculus of Variations and Nonlinear Elliptic Systems*, Princeton University Press, Princeton, NJ, 1983.
- [31] M. GIAQUINTA, *Introduction to Regularity Theory for Nonlinear Elliptic Systems*, Birkhäuser, Basel, 1993.
- [32] R. V. KOHN AND S. MÜLLER, *Surface energy and microstructure in coherent phase transitions*, Comm. Pure Appl. Math., 47 (1994), pp. 405–435.
- [33] R. V. KOHN AND M. VOGELIUS, *A new model for thin plates with rapidly varying thickness II: A convergence proof*, Quart. Appl. Math., 43 (1985), pp. 1–22.
- [34] H. LE DRET AND A. RAOULT, *The nonlinear membrane model as variational limit of nonlinear three-dimensional elasticity*, J. Math. Pures Appl., 74 (1995), pp. 549–578.

- [35] L. MODICA AND S. MORTOLA, *Un esempio di Γ -convergenza*, Boll. Un. Mat. Ital. B (5), 14 (1977), pp. 285–299.
- [36] M. G. MORA AND S. MÜLLER, *Derivation of a rod theory for multiphase materials*, Calc. Var. Partial Differential Equations, 28 (2007), pp. 161–178.
- [37] Y. C. SHU, *Heterogeneous thin films of martensitic materials*, Arch. Ration. Mech. Anal., 153 (2000), pp. 39–90.
- [38] P. STERNBERG, *The effect of a singular perturbation on nonconvex variational problems*, Arch. Ration. Mech. Anal., 101 (1988), pp. 209–260.
- [39] K. ZHANG, *Isolated microstructures on linear elastic strains*, Proc. R. Soc. Lond. A Ser. A math. Phys. Eng. Sci., 460 (2004), pp. 2993–3011.

DERIVATION OF A MACROSCOPIC RECEPTOR-BASED MODEL USING HOMOGENIZATION TECHNIQUES*

ANNA MARCINIAK-CZOCHRA[†] AND MARIYA PTASHNYK[‡]

Abstract. We study the problem of diffusive transport of biomolecules in the intercellular space, modeled as porous medium, and of their binding to the receptors located on the surface membranes of the cells. Cells are distributed periodically in a bounded domain. To describe this process we introduce a reaction-diffusion equation coupled with nonlinear ordinary differential equations on the boundary. We prove existence and uniqueness of the solution of this problem. We consider the limit, when the number of cells tends to infinity and at the same time their size tends to zero, while the volume fraction of the cells remains fixed. Using the homogenization technique of two-scale convergence, we show that the sequence of solutions of the original problem converges to the solution of the so-called macroscopic problem. To show the convergence of the nonlinear terms on the surfaces we use the unfolding method (periodic modulation). We discuss applicability of the result to mathematical description of membrane receptors of biological cells and compare the derived model with those previously considered.

Key words. homogenization, two-scale convergence, intercellular communication, receptor-ligand binding, reaction-diffusion equations, unfolding method (periodic modulation)

AMS subject classifications. 35B27, 74Q10, 74Q15, 35K57, 35K60

DOI. 10.1137/050645269

1. Introduction. Regulatory and signaling molecules (ligands) act by binding and activating receptor molecules. Receptors are usually located in the cell membrane, with some exceptions such as lipophilic ligands, which are located in the cytoplasm [24, 33, 34]. Some receptors interact with surface-bound ligands, such as adhesion proteins and extracellular matrix components. Other receptors bind soluble ligands, such as growth factors and cytokines. There are also many ligands which are present in both forms. As an example, antibodies, which are secreted by B cells as soluble molecules, become surface-bound ligands for the Fc receptors upon binding to antigens deposited on the surface [25].

Soluble molecules which are secreted to the intercellular space and transported via diffusion provide cell-to-cell communication, which results in the activation of processes in cells at a distance from the original signal. This happens, for example, in the case of the bystander effect. There is strong evidence that unirradiated bystander cells respond to signals emitted by irradiated cells [28]. In another case, the interplay between the spatial transport of virions and interferons results in the formation of patterns of infected and resistant cells [10]. Intercellular signaling can also lead to the formation of spatially nonhomogeneous structures, which is especially evident in developmental processes [33, 34]. The effects of the spatial transport of the soluble molecules are even visible in experiments in which only spatial averages are measured in order to understand the time dynamics of a signaling pathway. There is evidence

*Received by the editors November 15, 2005; accepted for publication (in revised form) August 31, 2007; published electronically April 16, 2008.

<http://www.siam.org/journals/sima/40-1/64526.html>

[†]Center for Modeling and Simulation in the Biosciences (BIOMS), Institute of Applied Mathematics, University of Heidelberg, Im Neuenheimer Feld 294, 69120 Heidelberg, Germany (anna.marciniak@iwr.uni-heidelberg.de).

[‡]Institute of Applied Mathematics, University of Heidelberg, Im Neuenheimer Feld 294, 69120 Heidelberg, Germany (mariya.ptashnyk@iwr.uni-heidelberg.de).

that different mixing conditions strongly influence quantitative and also qualitative results of such experiments [17]. Therefore, there arises a need to explain how the intercellular transport of the molecules should be described on the macroscopic level.

Models proposed so far are mainly phenomenological and describe all the processes on the macroscale level represented by a two- or three-dimensional sheet of cells [29, 30, 31, 32, 42, 46]. However, the real geometry is much more complicated, and binding a soluble ligand to a cell surface receptor requires interaction of molecules diffusing in three-dimensional space with some molecules attached to a two-dimensional surface. Since the size of a cell is very small compared to the dimension of the whole tissue, systems that include cells have to be treated as multiscale systems.

The aim of the present work is to derive a macroscopic model of receptor-ligand binding, based on a microscopic description, using methods of asymptotic analysis. Such an approach is called homogenization and it hinges on demonstrating the convergence of solutions of a sequence of microscopic problems to the solution of the macroscopic problem in properly chosen function spaces. We use here the two-scale convergence, which was introduced in [2] and [36] for sequences of functions $\{u^\varepsilon\}$ bounded in L^2 or in H^1 on an ε -periodic domain. Then, in [37] and [3], the definition of two-scale convergence was extended to sequences of functions defined on ε -periodic hypersurfaces, with dependence on parameters. This extension was used to homogenize a diffusion-reaction process in a catalyst consisting of distributed bars [37]. A similar problem with convection was studied in [19] using the standard homogenization technique, the energy method. A model describing processes of diffusion, convection, and nonlinear reactions in a periodic array of cells was studied in [20]. In that paper, the convergence of the nonlinear terms was shown using their monotonicity. Homogenization of models of chemical reactive flows in domains with periodically distributed reactive solid grains was also recently studied by Conca et al. [9]. They considered a stationary reaction-diffusion model with nonlinear, fast growing but monotone kinetics on the the surface of reactive solid grains and a model of reaction-diffusion processes both inside and outside of grains. Homogenization of the reaction-diffusion-convection processes with linear reactions on the surface of microstructures was also considered by Hornung in [18].

The model presented in this paper includes the dynamics of molecule concentrations on the surface of microstructures described by nonlinear ordinary differential equations. Therefore, we apply the concept of two-scale convergence of functions from L^∞ on ε -periodic hypersurfaces. To show convergence of the nonlinear terms on the surface of microstructures we use the unfolding method (periodic modulation); see [5, 6, 7].

Our paper is organized as follows. First, we present a precise description of the considered ε -periodic geometry (section 2) and of the equations describing the microscopic nature of the receptor-ligand binding process (section 3). These equations are spatially scaled by ε . Then we show existence and uniqueness of solutions of the microscopic model (section 3.2) and a priori estimates (section 3.3). In section 4, after extension of the solutions from the porous domain to the whole domain, using a priori estimates, we show the convergence of solutions of the microscopic problem to the solutions of a macroscopic homogenized model. Effective macroscopic equations are derived in section 4.2 and formulated in Theorem 4.4. In section 5 we compare a derived macroscopic model of the receptor-ligand binding on cells surfaces with the phenomenological models previously discussed in the literature.

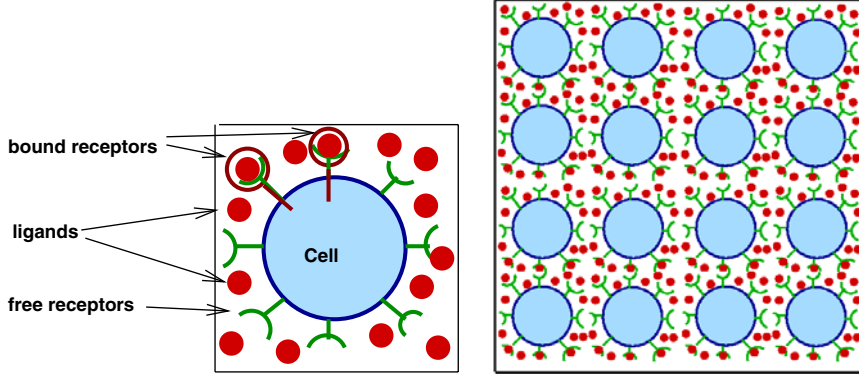


FIG. 1. *Geometry of the model. The array of the cells (on the right-hand side) consists of periodic repetition of the so-called standard cell, $Z = [0, 1]^3$ (on the left-hand side), which corresponds to a single biological cell with the surrounding intercellular space.*

2. Problem formulation. We consider a model involving a system of cells, periodically distributed in a three-dimensional cube $\Omega = [a, b]^3$, $a, b \in \mathbb{R}$, $a < b$, with boundary Γ^N . For the mathematical formulation of the problem we consider the so-called standard cell, $Z = [0, 1]^3$, periodically repeated over \mathbb{R}^3 with $Y_0 \subset Z$, an open subset with a smooth boundary Γ ; $Y = Z \setminus \bar{Y}_0$; and ν , the outer normal of Y (see Figure 1).

Let $\varepsilon > 0$ be a given scale factor such that $\varepsilon = \frac{b-a}{n}$, $n \in \mathbb{N}$, denoting the ratio between the size of the cells and the size of the whole domain Ω . Then the geometric structure within the fixed domain Ω is obtained by intersecting the ε -multiple εZ with Ω . We define, for $k \in \mathbb{Z}^3$, a triple of integers; and e_i , unit vectors, $\Gamma^k = \Gamma + \sum_{i=1}^3 k_i e_i$, $Y_0^k = Y_0 + \sum_{i=1}^3 k_i e_i$, $Z^k = Z + \sum_{i=1}^3 k_i e_i$, $\Gamma^* = \cup\{\Gamma^k, k \in \mathbb{Z}^3\}$, $Z^* = \cup\{Z^k, k \in \mathbb{Z}^3\}$. We further define $\Omega_0^\varepsilon = \cup\{\varepsilon Y_0^k | \varepsilon Z^k \subset \Omega, k \in \mathbb{Z}^3\}$, $\Omega^\varepsilon = \Omega \setminus \Omega_0^\varepsilon$, $\Gamma^\varepsilon = \cup\{\varepsilon \Gamma^k | \varepsilon Z^k \subset \Omega, k \in \mathbb{Z}^3\}$.

Remark 2.1. The geometry defined above fulfills the assumptions that

1. cells (holes in the domain) do not touch the boundary $\partial\Omega$;
2. cells do not touch each other;
3. cells have smooth boundary.

These assumptions allow for the definition of the functions on the cell boundaries using periodic repetition, and the definition of extension as proposed in [8]. Therefore, these assumptions are important for the methods applied in this paper. Homogenization of the Neumann problem in domains with more complicated geometry was considered in [1] and [4].

We assume that new ligands and new free receptors are produced on the cell surface through a combination of recycling (dissociation of bound receptors) and *de novo* production within the cell. Free receptors exist only on the surfaces, while ligands are transported by diffusion within the intercellular space, which is a porous medium. A ligand reversibly binds to a free receptor, which results in a bound receptor that can be internalized into the cell. Bound receptors also dissociate. Both ligands and free receptors undergo natural decay. We denote the concentration of ligands by $l^\varepsilon : (0, T) \times \Omega^\varepsilon \rightarrow \mathbb{R}$. Bound and free receptor densities are denoted by $r_b^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$ and $r_f^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$, respectively. For simplicity we assume that all binding processes are governed by the law of mass action without saturation effects.

3. Microscopic model.

3.1. Model assumptions. The microscopic model consists of the following equations:

Diffusion equation for ligands in the intercellular space,

$$\begin{aligned} \frac{\partial}{\partial t} l^\varepsilon(t, x) &= \nabla \cdot (D^\varepsilon(t, x) \nabla l^\varepsilon(t, x)) - \mu_l^\varepsilon(t, x) l^\varepsilon(t, x) + p_l^\varepsilon(t, x, l^\varepsilon(x, t)) \quad \text{in } (0, T) \times \Omega^\varepsilon, \\ \nu^\varepsilon \cdot \nabla_x l^\varepsilon(t, x) &= 0 \quad \text{on } (0, T) \times \Gamma^N, \\ (1) \quad l^\varepsilon(x, t) &= l_0(x), \quad t = 0, \quad x \in \Omega^\varepsilon. \end{aligned}$$

Binding equation on the surfaces,

$$(2) \quad -D^\varepsilon(t, x) \nabla l^\varepsilon(t, x) \cdot \nu^\varepsilon = \varepsilon(b^\varepsilon(t, x) l^\varepsilon(t, x) r_f^\varepsilon(t, x) - d^\varepsilon(t, x) r_b^\varepsilon(t, x)) \quad \text{on } (0, T) \times \Gamma^\varepsilon.$$

Reaction equations for receptors on the surfaces,

$$\begin{aligned} \frac{\partial}{\partial t} r_f^\varepsilon(x, t) &= -\mu_f^\varepsilon(t, x) r_f^\varepsilon(x, t) + p_r^\varepsilon(t, x, r_b^\varepsilon(x, t)) - b^\varepsilon(t, x) r_f^\varepsilon(x, t) l^\varepsilon(x, t) \\ (3) \quad &+ d^\varepsilon(t, x) r_b^\varepsilon(x, t), \\ (4) \quad \frac{\partial}{\partial t} r_b^\varepsilon(x, t) &= -\mu_b^\varepsilon(t, x) r_b^\varepsilon(x, t) + b^\varepsilon(t, x) r_f^\varepsilon(x, t) l^\varepsilon(x, t) - d^\varepsilon(t, x) r_b^\varepsilon(x, t), \end{aligned}$$

with initial conditions

$$(5) \quad r_f^\varepsilon(x, t) = r_{f0}(x), \quad t = 0, \quad x \in \Gamma^\varepsilon,$$

$$(6) \quad r_b^\varepsilon(x, t) = r_{b0}(x), \quad t = 0, \quad x \in \Gamma^\varepsilon.$$

The following is a list of functional coefficients in these equations:

$\mu_l^\varepsilon : (0, T) \times \Omega \rightarrow \mathbb{R}$	rate of decay of ligands,
$p_l^\varepsilon : (0, T) \times \Omega \times \mathbb{R} \rightarrow \mathbb{R}$	production of ligands,
$D^\varepsilon : (0, T) \times \Omega \rightarrow \mathbb{R}^{3 \times 3}$	diffusion coefficient for ligands,
$p_r^\varepsilon : (0, T) \times \Gamma^\varepsilon \times \mathbb{R} \rightarrow \mathbb{R}$	production of new free receptors,
$\mu_f^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$	rate of decay of free receptors,
$\mu_b^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$	rate of decay of bound receptors,
$d^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$	rate of dissociation of bound receptors,
$b^\varepsilon : (0, T) \times \Gamma^\varepsilon \rightarrow \mathbb{R}$	rate of binding of ligands and free receptors,

where functions on Ω or Γ^ε are defined by Z -periodic function: $D_{i,j}^\varepsilon(t, x) = D_{i,j}(t, \frac{x}{\varepsilon})$, $p_l^\varepsilon(t, x, \xi) = p_l(t, \frac{x}{\varepsilon}, \xi)$, $\mu_l^\varepsilon(t, x) = \mu_l(t, \frac{x}{\varepsilon})$, $\mu_f^\varepsilon(t, x) = \mu_f(t, \frac{x}{\varepsilon})$, $\mu_b^\varepsilon(t, x) = \mu_b(t, \frac{x}{\varepsilon})$, $b^\varepsilon(t, x) = b(t, \frac{x}{\varepsilon})$, $d^\varepsilon(t, x) = d(t, \frac{x}{\varepsilon})$, $p_r^\varepsilon(t, x, \xi) = p_r(t, \frac{x}{\varepsilon}, \xi)$, defined on Z^* and Γ^* , respectively.

We assume that decay processes are linear and that binding is a product of the density of ligands and free receptors. The proposed functions are the simplest functions usually used to describe decay or binding processes (see the models described in [35]), modeled by the law of mass action. We assume that *de novo* production of free receptors, denoted by p_r , is regulated by bound receptors. We assume that p_r

is a bounded Lipschitz continuous function in r_b and is nonnegative for nonnegative values of r_b , for example a Michaelis–Menten function $p_r = \frac{m_1 r_b}{1+r_b}$. In addition, we assume that the production of ligands depends on their density. It could be regulated via some other receptors not considered in our model. Thus, we assume that p_l is a Lipschitz continuous function in l , nonnegative for nonnegative values of l .

Assumption 3.1.

1. $D \in L^\infty((0, T) \times Z)^{3 \times 3}$, $\partial_t D \in L^\infty((0, T) \times Z)^{3 \times 3}$, $(D(t, x)\xi, \xi) \geq d_0 |\xi|^2$ for some $d_0 > 0$, for every $\xi \in \mathbb{R}^3$, a.a. $(t, x) \in (0, T) \times Z$.
2. $\mu_l \in L^\infty((0, T) \times Z)$ and $\mu_l \geq 0$ a.e. in $(0, T) \times Z$.
3. p_l is measurable in t and x , sublinear, i.e., $|p_l(t, x, \xi)| \leq c_1 + c_2 |\xi|$ for a.a. $(t, x) \in (0, T) \times Z$, Lipschitz continuous in ξ , and $p_l(t, x, \xi) \geq 0$ for $\xi \geq 0$.
4. $b \in C([0, T]; C^{0, \alpha}(\Gamma))$, $b \geq 0$, in $[0, T] \times \Gamma$, $\partial_t b \in L^\infty((0, T) \times \Gamma)$.
5. $d \in C([0, T]; C^{0, \alpha}(\Gamma))$, $d \geq 0$, in $[0, T] \times \Gamma$, $\partial_t d \in L^\infty((0, T) \times \Gamma)$.
6. $\mu_f, \mu_b \in C([0, T]; C^{0, \alpha}(\Gamma))$, $\mu_f \geq 0$, $\mu_b \geq 0$, in $[0, T] \times \Gamma$.
7. $p_r(\xi) \in C([0, T]; C^{0, \alpha}(\Gamma))$ for all $\xi \in \mathbb{R}$, $p_r(t, x, \xi) \geq 0$ for $\xi \geq 0$, p_r is bounded, i.e., $|p_r(t, x, \xi)| \leq m_1$ for all $(t, x, \xi) \in (0, T) \times \Gamma \times \mathbb{R}$ and is Lipschitz continuous in ξ .

3.2. Existence of the solutions of the microscopic model. We start with a weak formulation of the microscopic model.

DEFINITION 3.2. *The triple $(l^\varepsilon, r_f^\varepsilon, r_b^\varepsilon)$ is a solution of problem (1)–(6) if $l^\varepsilon \in L^2((0, T); H^1(\Omega^\varepsilon))$, $\partial_t l^\varepsilon \in L^2((0, T) \times \Omega^\varepsilon)$, $l^\varepsilon \in L^\infty((0, T) \times \Omega^\varepsilon)$, $r_f^\varepsilon, r_b^\varepsilon \in L^\infty((0, T) \times \Gamma^\varepsilon)$, $\partial_t r_f^\varepsilon, \partial_t r_b^\varepsilon \in L^\infty((0, T) \times \Gamma^\varepsilon)$ such that*

1.

$$(7) \quad \begin{aligned} (\partial_t l^\varepsilon, \phi)_{(0, T) \times \Omega^\varepsilon} &= -(D^\varepsilon \nabla l^\varepsilon, \nabla \phi)_{(0, T) \times \Omega^\varepsilon} - (\mu_l^\varepsilon l^\varepsilon, \phi)_{(0, T) \times \Omega^\varepsilon} \\ &+ (d^\varepsilon r_b^\varepsilon - b^\varepsilon r_f^\varepsilon l^\varepsilon, \phi)_{(0, T) \times \Gamma^\varepsilon} + (p_l^\varepsilon(l^\varepsilon), \phi)_{(0, T) \times \Omega^\varepsilon} \end{aligned}$$

for all $\phi \in L^2((0, T); H^1(\Omega^\varepsilon))$;

2. l^ε satisfies the initial condition, i.e., $l^\varepsilon \rightarrow l_0$ in $L^2(\Omega^\varepsilon)$ as $t \rightarrow 0$;

3.

$$(8) \quad \begin{cases} \frac{\partial}{\partial t} r_f^\varepsilon(x, t) = -\mu_f^\varepsilon r_f^\varepsilon(x, t) + p_r^\varepsilon(t, x, r_b^\varepsilon(x, t)) \\ \quad - b^\varepsilon r_f^\varepsilon(x, t) l^\varepsilon(x, t) + d^\varepsilon r_b^\varepsilon(x, t), \\ \frac{\partial}{\partial t} r_b^\varepsilon(x, t) = -\mu_b^\varepsilon r_b^\varepsilon(x, t) + b^\varepsilon r_f^\varepsilon(x, t) l^\varepsilon(x, t) \\ \quad - d^\varepsilon r_b^\varepsilon(x, t) \end{cases}$$

a.e. $(0, T) \times \Gamma^\varepsilon$;

4. $r_f^\varepsilon, r_b^\varepsilon$ satisfy the initial conditions (5)–(6).

Here $(u, v)_{(0, T) \times \Omega^\varepsilon} = \int_0^T \int_{\Omega^\varepsilon} u v \, dx \, dt$ and $(u, v)_{(0, T) \times \Gamma^\varepsilon} = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} u v \, d\gamma_x \, dt$.

THEOREM 3.3. *Let Assumption 3.1 be satisfied and*

$$l_0 \in C^{0, \alpha}(\overline{\Omega}), \quad l_0 \in H^1(\Omega), \quad l_0 \geq 0,$$

$$r_{f0}, r_{b0} \in C^{0, \alpha}(\overline{\Omega}), \quad r_{f0} \geq 0, \quad r_{b0} \geq 0.$$

Then there exists a unique solution $(l^\varepsilon, r_f^\varepsilon, r_b^\varepsilon)$ of problem (1)–(6), such that

$$l^\varepsilon \in H^1(0, T; L^2(\Omega^\varepsilon)), \quad l^\varepsilon \in L^2(0, T; H^1(\Omega^\varepsilon)),$$

$$l^\varepsilon \in C^{0, \beta/2}([0, T]; C^{0, \beta}(\overline{\Omega^\varepsilon})),$$

$$r_f^\varepsilon, r_b^\varepsilon \in C^1([0, T]; C^{0, \beta}(\Gamma^\varepsilon)), \quad \text{where } \beta \in (0, \alpha],$$

$$\text{and } l^\varepsilon \geq 0, r_f^\varepsilon \geq 0, r_b^\varepsilon \geq 0.$$

Proof. Existence. The existence of a solution of the system (1), (3), (4) will be proved by showing the existence of a fix point of the operator K defined on $C([0, T] \times \overline{\Omega}^\varepsilon)$ by $l^{n,\varepsilon} = K(l^{n-1,\varepsilon})$ with $l^{n,\varepsilon}$ given by

$$(9) \quad \begin{cases} \partial_t l^{n,\varepsilon} = \nabla \cdot (D^\varepsilon \nabla l^{n,\varepsilon}) - \mu_l^\varepsilon l^{n,\varepsilon} + p_l^\varepsilon(l^{n-1,\varepsilon}), & t > 0, x \in \Omega^\varepsilon, \\ \nabla l^{n,\varepsilon} \cdot \nu^\varepsilon = 0, & t > 0, x \in \Gamma^N, \\ l^{n,\varepsilon} = l_0, & t = 0, x \in \Omega^\varepsilon, \\ -D^\varepsilon \nabla l^{n,\varepsilon} \cdot \nu^\varepsilon = \varepsilon(b^\varepsilon l^{n,\varepsilon} r_f^{n,\varepsilon} - d^\varepsilon r_b^{n,\varepsilon}), & t > 0, x \in \Gamma^\varepsilon, \end{cases}$$

$$(10) \quad \begin{cases} \partial_t r_f^{n,\varepsilon} = -\mu_f^\varepsilon r_f^{n,\varepsilon} + p_r^\varepsilon(r_b^{n,\varepsilon}) - b^\varepsilon r_f^{n,\varepsilon} l^{n-1,\varepsilon} + d^\varepsilon r_b^{n,\varepsilon}, & t > 0, x \in \Gamma^\varepsilon, \\ \partial_t r_b^{n,\varepsilon} = -\mu_b^\varepsilon r_b^{n,\varepsilon} + b^\varepsilon r_f^{n,\varepsilon} l^{n-1,\varepsilon} - d^\varepsilon r_b^{n,\varepsilon}, & t > 0, x \in \Gamma^\varepsilon, \\ r_f^{n,\varepsilon} = r_{f0}, & t = 0, x \in \Gamma^\varepsilon, \\ r_b^{n,\varepsilon} = r_{b0}, & t = 0, x \in \Gamma^\varepsilon. \end{cases}$$

For a given $l^{n-1,\varepsilon} \in C([0, T] \times \overline{\Omega}^\varepsilon)$, $l^{n-1,\varepsilon} \geq 0$ on $[0, T] \times \overline{\Omega}^\varepsilon$, there exists a unique solution of system (10), $r_f^{n,\varepsilon}, r_b^{n,\varepsilon} \in C^1([0, T]; C(\Gamma^\varepsilon))$, because the right-hand side of the system of ordinary differential equations (10) is Lipschitz continuous [45]. Since p_r is a nonnegative function for nonnegative values of r_b and $l^{n-1,\varepsilon} \geq 0$ on $[0, T] \times \Gamma^\varepsilon$ and $r_{f0} \geq 0, r_{b0} \geq 0$, we deduce that $r_f^{n,\varepsilon} \geq 0, r_b^{n,\varepsilon} \geq 0$ on $[0, T] \times \Gamma^\varepsilon$.

Using the Galerkin method and a priori estimates similar to the estimates in Lemma 3.4, we obtain the existence of a weak solution of (9), $l^{n,\varepsilon} \in L^2(0, T; H^1(\Omega^\varepsilon))$, $\partial_t l^{n,\varepsilon} \in L^2(0, T; L^2(\Omega^\varepsilon))$; see [23]. Since $l_0 \in C^{0,\alpha}(\overline{\Omega})$, there exists $\max_{\overline{\Omega}^\varepsilon} |l_0| = M$. In addition, $r_f^{n,\varepsilon} \geq 0$ and $|r_b^{n,\varepsilon}| \leq C$. Thus, we may apply the result from [26] (Theorem 6.40) stating that for parabolic equations with uniformly elliptic operator, sublinear terms of lower order, bounded free terms, and bounded coefficients of Robin boundary conditions, the boundedness of the initial conditions implies the boundedness of the supremum of a solution. From this, we conclude that $\sup_{(0,T) \times \Omega^\varepsilon} |l^{n,\varepsilon}| \leq M_1$. Then, since $l_0 \in C^{0,\alpha}(\overline{\Omega})$, $r_f^{n,\varepsilon} \geq 0$, and $r_b^{n,\varepsilon} \in C^1([0, T]; C(\Gamma^\varepsilon))$, we conclude also that $l^{n,\varepsilon} \in C^{0,\beta/2}([0, T]; C^{0,\beta}(\overline{\Omega}^\varepsilon))$ (see Theorem III.10.1 in [23], generalized for Robin boundary conditions, or [11] and [26]). Using the maximum principle and the continuity of $l^{n,\varepsilon}$, we obtain that $l^{n,\varepsilon} \geq 0$ in $[0, T] \times \overline{\Omega}^\varepsilon$ [12].

The space $C^{0,\beta/2}([0, T]; C^{0,\beta}(\overline{\Omega}^\varepsilon))$ is compact embedded in $C([0, T] \times \overline{\Omega}^\varepsilon)$. Then, by virtue of the Schauder theorem, there exists a fixed point of K , a solution of the microscopic problem $l^\varepsilon, r_f^\varepsilon$, and r_b^ε . In addition, we obtain that $l^\varepsilon \geq 0, r_f^\varepsilon \geq 0$, and $r_b^\varepsilon \geq 0$. Since $r_{f0}, r_{b0} \in C^{0,\alpha}(\Omega)$ and $l^\varepsilon \in C^{0,\beta/2}([0, T]; C^{0,\beta}(\overline{\Omega}^\varepsilon))$, we conclude also that $r_f^\varepsilon, r_b^\varepsilon \in C^1([0, T]; C^{0,\beta}(\Gamma^\varepsilon))$.

Uniqueness. Suppose there are two solutions of the problem $(l^{1,\varepsilon}, r_f^{1,\varepsilon}, r_b^{1,\varepsilon})$ and $(l^{2,\varepsilon}, r_f^{2,\varepsilon}, r_b^{2,\varepsilon})$. We denote $l^\varepsilon = l^{1,\varepsilon} - l^{2,\varepsilon}$ and choose $\phi = l^\varepsilon$. We calculate

$$\begin{aligned} & \frac{1}{2} \int_0^\tau \int_{\Omega^\varepsilon} \left(\partial_t |l^\varepsilon|^2 + (D^\varepsilon \nabla l^\varepsilon, \nabla l^\varepsilon) + \mu_l^\varepsilon |l^\varepsilon|^2 \right) dx dt = \int_0^\tau \int_{\Omega^\varepsilon} (p_l^\varepsilon(l^{1,\varepsilon}) - p_l^\varepsilon(l^{2,\varepsilon})) l^\varepsilon dx dt \\ & + \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} ((d^\varepsilon r_b^{1,\varepsilon} - b^\varepsilon r_f^{1,\varepsilon} l^{1,\varepsilon}) - (d^\varepsilon r_b^{2,\varepsilon} - b^\varepsilon r_f^{2,\varepsilon} l^{2,\varepsilon})) l^\varepsilon d\gamma dt \end{aligned}$$

for any $\tau \in [0, T]$. For r_f^ε and r_b^ε we obtain

$$\begin{aligned} \frac{\partial}{\partial t}(r_f^{1,\varepsilon} - r_f^{2,\varepsilon}) &= -\mu_f^\varepsilon(r_f^{1,\varepsilon} - r_f^{2,\varepsilon}) + (p_r^\varepsilon(r_b^{1,\varepsilon}) - p_r^\varepsilon(r_b^{2,\varepsilon})) - b^\varepsilon(r_f^{1,\varepsilon}l^{1,\varepsilon} - r_f^{2,\varepsilon}l^{2,\varepsilon}) \\ &\quad + d^\varepsilon(r_b^{1,\varepsilon} - r_b^{2,\varepsilon}), \end{aligned}$$

$$\frac{\partial}{\partial t}(r_b^{1,\varepsilon} - r_b^{2,\varepsilon}) = -\mu_b^\varepsilon(r_b^{1,\varepsilon} - r_b^{2,\varepsilon}) + b^\varepsilon(r_f^{1,\varepsilon}l^{1,\varepsilon} - r_f^{2,\varepsilon}l^{2,\varepsilon}) - d^\varepsilon(r_b^{1,\varepsilon} - r_b^{2,\varepsilon}).$$

Integrating by parts with respect to time and summing up side by side the last two equations, we obtain

$$\begin{aligned} |r_f^{1,\varepsilon} - r_f^{2,\varepsilon}| + |r_b^{1,\varepsilon} - r_b^{2,\varepsilon}| &\leq \int_0^\tau \left(\mu_f^1 |r_f^{1,\varepsilon} - r_f^{2,\varepsilon}| + \mu_b^1 |r_b^{1,\varepsilon} - r_b^{2,\varepsilon}| + c_r |r_b^{1,\varepsilon} - r_b^{2,\varepsilon}| \right) dt \\ &+ \int_0^\tau \left(2b_1 \max_{[0,T] \times \Gamma^\varepsilon} |l^{1,\varepsilon}| |r_f^{1,\varepsilon} - r_f^{2,\varepsilon}| + 2b_1 \max_{[0,T] \times \Gamma^\varepsilon} |r_f^{2,\varepsilon}| |l^{1,\varepsilon} - l^{2,\varepsilon}| + 2d_1 |r_b^{1,\varepsilon} - r_b^{2,\varepsilon}| \right) dt, \end{aligned}$$

where c_r is the Lipschitz constant of p_r , $\mu_f^1 = \sup_{[0,T] \times \Gamma^\varepsilon} |\mu_f^\varepsilon|$, $\mu_b^1 = \sup_{[0,T] \times \Gamma^\varepsilon} |\mu_b^\varepsilon|$, $b_1 = \sup_{[0,T] \times \Gamma^\varepsilon} |b^\varepsilon|$, $d_1 = \sup_{[0,T] \times \Gamma^\varepsilon} |d^\varepsilon|$. The Gronwall lemma implies

$$(11) \quad |r_f^{1,\varepsilon} - r_f^{2,\varepsilon}| + |r_b^{1,\varepsilon} - r_b^{2,\varepsilon}| \leq C \int_0^\tau |l^{1,\varepsilon} - l^{2,\varepsilon}| dt.$$

Using the above estimate and nonnegativity of b^ε and $r_f^{2,\varepsilon}$, we obtain

$$\begin{aligned} &\frac{1}{2} \int_0^\tau \int_{\Omega^\varepsilon} \partial_t |l^\varepsilon|^2 dx dt + d_0 \int_0^\tau \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt + \int_0^\tau \int_{\Omega^\varepsilon} \mu_i^\varepsilon |l^\varepsilon|^2 dx dt \\ &\leq C d_1^2 \varepsilon \frac{1}{2\delta} \int_0^\tau \int_{\Gamma^\varepsilon} \int_0^t |l^\varepsilon|^2 ds d\gamma dt + \varepsilon \frac{\delta}{2} \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma dt + c_l \int_0^\tau \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx dt \\ &+ C b_1 \varepsilon \max_{[0,T] \times \Gamma^\varepsilon} |l^{1,\varepsilon}| \int_0^\tau \int_{\Gamma^\varepsilon} \int_0^t |l^\varepsilon|^2 ds d\gamma dt + C b_1 \varepsilon \max_{[0,T] \times \Gamma^\varepsilon} |l^{1,\varepsilon}| \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma dt, \end{aligned}$$

where c_l is the Lipschitz constant of p_l . Furthermore, using the estimate

$$(12) \quad \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma dt \leq c \int_0^\tau \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx dt + c\varepsilon^2 \int_0^\tau \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt,$$

we obtain

$$\begin{aligned} &\frac{1}{2} \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx + (d_0 - \varepsilon^2 \delta) \int_0^\tau \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt + \int_0^\tau \int_{\Omega^\varepsilon} \mu_i^\varepsilon |l^\varepsilon|^2 dx dt \\ &\leq C \frac{1}{\delta} \int_0^\tau \int_0^t \int_{\Omega^\varepsilon} (|l^\varepsilon|^2 + |\nabla l^\varepsilon|^2) dx ds dt + c_l \int_0^\tau \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx dt. \end{aligned}$$

From the Gronwall lemma and $\mu_i^\varepsilon \geq 0$, taking the supremum over $\tau \in [0, T]$, we conclude that

$$\int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx + C \int_0^T \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt \leq 0$$

and, therefore, $l^{1,\varepsilon} = l^{2,\varepsilon}$ in $(0, T) \times \Omega^\varepsilon$. Due to (11), also $r_f^{1,\varepsilon} = r_f^{2,\varepsilon}$ and $r_b^{1,\varepsilon} = r_b^{2,\varepsilon}$ on $[0, T] \times \Gamma^\varepsilon$. \square

3.3. A priori estimates for the microscopic solutions.

LEMMA 3.4. *For any solution of problem (1)–(6) from Theorem 3.3 the following estimates hold:*

$$\begin{aligned} \|l^\varepsilon\|_{L^2(0,T;H^1(\Omega^\varepsilon))} &\leq C, & \|\partial_t l^\varepsilon\|_{L^2(0,T;L^2(\Omega^\varepsilon))} &\leq C, \\ \|r_f^\varepsilon\|_{L^\infty((0,T)\times\Gamma^\varepsilon)} &\leq C, & \|r_b^\varepsilon\|_{L^\infty((0,T)\times\Gamma^\varepsilon)} &\leq C, \\ \|\partial_t r_f^\varepsilon\|_{L^2((0,T)\times\Gamma^\varepsilon)} &\leq C, & \|\partial_t r_b^\varepsilon\|_{L^2((0,T)\times\Gamma^\varepsilon)} &\leq C, \end{aligned}$$

where C is a constant independent on ε .

Proof. To show the estimates for r_f^ε and r_b^ε we add (3) and (4) side by side and obtain

$$\partial_t(r_f^\varepsilon + r_b^\varepsilon) \leq m_1.$$

Since r_f^ε and r_b^ε are nonnegative (see Theorem 3.3), we conclude that

$$\|r_f^\varepsilon\|_{L^\infty((0,T)\times\Gamma^\varepsilon)} \leq C \quad \text{and} \quad \|r_b^\varepsilon\|_{L^\infty((0,T)\times\Gamma^\varepsilon)} \leq C.$$

Now we show the estimates for l^ε . We choose $\phi = l^\varepsilon$ as a test function in (7) and calculate

$$\begin{aligned} &\frac{1}{2} \int_0^\tau \int_{\Omega^\varepsilon} \partial_t |l^\varepsilon|^2 dx dt + \int_0^\tau \int_{\Omega^\varepsilon} (D^\varepsilon \nabla l^\varepsilon, \nabla l^\varepsilon) dx dt + \int_0^\tau \int_{\Omega^\varepsilon} \mu_l^\varepsilon |l^\varepsilon|^2 dx dt \\ &= \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} (d^\varepsilon r_b^\varepsilon - b^\varepsilon r_f^\varepsilon l^\varepsilon) l^\varepsilon d\gamma dt + \int_0^\tau \int_{\Omega^\varepsilon} p_l^\varepsilon(l^\varepsilon) l^\varepsilon dx dt \end{aligned}$$

for any $\tau \in [0, T]$. Applying the Young inequality we obtain

$$\begin{aligned} &\frac{1}{2} \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx + \int_0^\tau \int_{\Omega^\varepsilon} d_0 |\nabla l^\varepsilon|^2 dx dt + \int_0^\tau \int_{\Omega^\varepsilon} \mu_l^\varepsilon |l^\varepsilon|^2 dx dt \\ &\leq \frac{\varepsilon d_1}{2\delta} \int_0^\tau \int_{\Gamma^\varepsilon} |r_b^\varepsilon|^2 d\gamma dt + \varepsilon \frac{\delta}{2} \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma dt \\ &\quad - \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon |l^\varepsilon|^2 d\gamma dt + c_1 \int_0^\tau \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx dt + \frac{1}{2} \int_{\Omega^\varepsilon} |l_0^\varepsilon|^2 dx. \end{aligned}$$

Now we use (12), $\mu_l^\varepsilon \geq 0$, $b^\varepsilon \geq 0$, and $r_f^\varepsilon \geq 0$ and obtain

$$\begin{aligned} &\frac{1}{2} \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx + \int_0^\tau \int_{\Omega^\varepsilon} \left(d_0 - \frac{\delta \varepsilon^2}{2} \right) |\nabla l^\varepsilon|^2 dx dt \\ &\leq \frac{\varepsilon}{2\delta} \int_0^\tau \int_{\Gamma^\varepsilon} |r_b^\varepsilon|^2 d\gamma dt + c_1 \int_0^\tau \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx dt + \frac{1}{2} \int_{\Omega^\varepsilon} |l_0^\varepsilon|^2 dx. \end{aligned}$$

Then, from the Gronwall lemma and the estimate for r_b^ε , it follows that

$$\int_{\Omega} |l^\varepsilon|^2 dx + \int_0^T \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt \leq C.$$

Using the estimates for l^ε , r_f^ε , and r_b^ε , we conclude from (8) that

$$\begin{aligned} \|\partial_t r_f^\varepsilon\|_{L^2((0,T)\times\Gamma^\varepsilon)} &\leq C, \\ \|\partial_t r_b^\varepsilon\|_{L^2((0,T)\times\Gamma^\varepsilon)} &\leq C. \end{aligned}$$

To obtain the estimates for $\partial_t l^\varepsilon$ we choose $\phi = \partial_t l^\varepsilon$ as a test function and calculate

$$\begin{aligned} &\int_0^\tau \int_{\Omega^\varepsilon} |\partial_t l^\varepsilon|^2 dx dt + \frac{1}{2} \int_0^\tau \int_{\Omega^\varepsilon} \left(\partial_t (D^\varepsilon \nabla l^\varepsilon, \nabla l^\varepsilon) - (\partial_t D^\varepsilon \nabla l^\varepsilon, \nabla l^\varepsilon) \right) dx dt \\ &= \int_0^\tau \int_{\Omega^\varepsilon} (p_l^\varepsilon(l^\varepsilon) - \mu_l^\varepsilon l^\varepsilon) \partial_t l^\varepsilon dx dt + \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} \left(\partial_t (d^\varepsilon r_b^\varepsilon l^\varepsilon) - d^\varepsilon \partial_t r_b^\varepsilon l^\varepsilon - \partial_t d^\varepsilon r_b^\varepsilon l^\varepsilon \right) d\gamma dt \\ &+ \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} \left(-\partial_t (b^\varepsilon r_f^\varepsilon |l^\varepsilon|^2) + b^\varepsilon \partial_t r_f^\varepsilon |l^\varepsilon|^2 + \partial_t b^\varepsilon r_f^\varepsilon |l^\varepsilon|^2 \right) d\gamma dt. \end{aligned}$$

Using the Young inequality we obtain

$$\begin{aligned} &(1 - \delta) \int_0^\tau \int_{\Omega^\varepsilon} |\partial_t l^\varepsilon|^2 dx dt + \frac{d_0}{2} \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx \\ &\leq \frac{\varepsilon}{2\delta} \int_{\Gamma^\varepsilon} d_1^2 |r_b^\varepsilon|^2 d\gamma + D_2 \int_0^\tau \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx dt \\ &+ \varepsilon \frac{\delta}{2} \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma + \frac{\varepsilon}{2} \int_0^\tau \int_{\Gamma^\varepsilon} (d_1^2 |\partial_t r_b^\varepsilon|^2 + |\partial_t d| |r_b^\varepsilon|^2) d\gamma dt \\ &+ \frac{\varepsilon}{2} \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^2 d\gamma dt - \varepsilon \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon |l^\varepsilon|^2 d\gamma \\ &+ \frac{1}{2\delta} \int_0^\tau \int_{\Omega^\varepsilon} (|p_l^\varepsilon(l^\varepsilon)|^2 + \mu_l^\varepsilon |l^\varepsilon|^2) dx dt + \varepsilon \int_{\Gamma^\varepsilon} (d_1 r_{b0} l_0 + b_1 r_{f0} |l_0|^2) d\gamma \\ &+ D_1 \int_{\Omega^\varepsilon} |\nabla l_0|^2 dx + \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} (\partial_t b^\varepsilon r_f^\varepsilon + b^\varepsilon |\partial_t r_f^\varepsilon|) |l^\varepsilon|^2 d\gamma dt, \end{aligned}$$

where $D_1 = \sup_{(0,T)\times\Omega} |D^\varepsilon|$, $D_2 = \sup_{(0,T)\times\Omega} |\partial_t D^\varepsilon|$. For the estimate of the last integral we use the embedding for a space of dimension $n = 3$, i.e., $L^\infty(0, T; H^1(\Omega^\varepsilon)) \subset L^4((0, T) \times \Gamma^\varepsilon)$,

$$\begin{aligned} \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} b^\varepsilon |\partial_t r_f^\varepsilon| |l^\varepsilon|^2 d\gamma dt &\leq \frac{b_1^2 \varepsilon}{2\delta} \int_0^\tau \int_{\Gamma^\varepsilon} |\partial_t r_f^\varepsilon|^2 d\gamma dt + \frac{\delta \varepsilon}{2} \int_0^\tau \int_{\Gamma^\varepsilon} |l^\varepsilon|^4 d\gamma dt \\ &\leq \frac{b_1^2 \varepsilon}{2\delta} \int_0^\tau \int_{\Gamma^\varepsilon} |\partial_t r_f^\varepsilon|^2 d\gamma dt + \frac{\delta}{2} \sup_{[0,T]} \int_{\Omega^\varepsilon} |l^\varepsilon|^2 dx + \frac{\delta \varepsilon^2}{2} \sup_{[0,T]} \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx. \end{aligned}$$

Using estimate (12) and the positivity of b^ε and r_f^ε we obtain

$$\int_0^\tau \int_{\Omega^\varepsilon} |\partial_t l^\varepsilon|^2 dx dt + \sup_{[0,T]} \int_{\Omega^\varepsilon} |\nabla l^\varepsilon|^2 dx \leq C. \quad \square$$

To obtain a priori estimates for functions defined in the domain independent of ε , we extend functions l^ε defined on Ω^ε to functions \bar{l}^ε defined on the whole Ω .

3.4. Extension of l^ε . Since l^ε is defined only on Ω^ε , we extend it onto Ω ; see [8] or [19] for the proof.

LEMMA 3.5. 1. For $l \in H^1(Y)$ there exists an extension \tilde{l} to Z , such that

$$\|\tilde{l}\|_{L^2(Z)} \leq c\|l\|_{L^2(Y)} \quad \text{and} \quad \|\nabla\tilde{l}\|_{L^2(Z)} \leq c\|\nabla l\|_{L^2(Y)}.$$

2. For $l^\varepsilon \in H^1(\Omega^\varepsilon)$ there exists an extension \tilde{l}^ε to Ω , such that

$$\|\tilde{l}^\varepsilon\|_{H^1(\Omega)} \leq c\|l^\varepsilon\|_{H^1(\Omega^\varepsilon)}.$$

Remark 3.1. For $l^\varepsilon \in L^2(0, T; H^1(\Omega^\varepsilon))$ we define $\bar{l}^\varepsilon(\cdot, t) := \tilde{l}^\varepsilon(\cdot, t)$ for a.a. t . Since the extension operator is linear, then $\bar{l}^\varepsilon \in L^2(0, T; H^1(\Omega))$.

We identify l^ε with the extension \bar{l}^ε . For the extended functions, we obtain a priori estimate of the supremum norm of l^ε .

LEMMA 3.6. For any solution of problem (1)–(6), the following estimate holds:

$$(13) \quad \|l^\varepsilon\|_{L^\infty((0,T)\times\Omega)} \leq C,$$

where C is a constant independent on ε .

Estimate (13) follows from the nonnegativity of l^ε , r_f^ε , r_b^ε , the boundedness of r_b^ε and l_0 , and the estimate in Lemma 3.5; see Theorem 6.40 in [26] (for the sketch of proof see Appendix 6.1).

4. Convergence of solutions of microscopic problem.

4.1. Convergence of l^ε , r_f^ε , and r_b^ε . To show the convergence results we apply the method of two-scale convergence, introduced in [2] and [36], and extended further in [3, 37]. The definition and theorems concerning the two-scale convergence, used in this section are outlined in Appendix 6.2.

To show the compactness of l^ε we use the following Hilbert space.

DEFINITION 4.1 (see [47]). Let $W^{\beta,2}(\Omega)$ with $\beta \in \mathbb{R}$, $\beta > 0$ be a Hilbert space defined as the completion of $C^\infty(\Omega)$ with respect to the norm

$$\|u\|_{W^{\beta,2}(\Omega)} = \|u\|_{W^{k,2}(\Omega)} + \left\{ \int_{\Omega} \int_{\Omega} \frac{|u(x) - u(y)|^2}{|x - y|^{n+2(\beta-k)}} dx dy \right\}^{\frac{1}{2}},$$

where $k = [\beta]$.

LEMMA 4.2. 1. For a function $v^\varepsilon \in H^1(\Omega^\varepsilon)$ the following estimate holds:

$$\varepsilon \int_{\Gamma^\varepsilon} |v^\varepsilon|^2 d\gamma_x \leq C \int_{\Omega^\varepsilon} |v^\varepsilon|^2 dx + C\varepsilon^2 \int_{\Omega^\varepsilon} |\nabla v^\varepsilon|^2 dx,$$

where C is a constant independent on ε .

2. For a function $v^\varepsilon \in W^{\beta,2}(\Omega^\varepsilon)$, where $\frac{1}{2} < \beta < 1$, the following estimate holds:

$$\varepsilon \int_{\Gamma^\varepsilon} |v^\varepsilon|^2 d\gamma_x \leq C \int_{\Omega^\varepsilon} |v^\varepsilon|^2 dx + C\varepsilon^{2\beta} \int_{\Omega^\varepsilon} \int_{\Omega^\varepsilon} \frac{|v^\varepsilon(x_1) - v^\varepsilon(x_2)|^2}{|x_1 - x_2|^{n+2\beta}} dx_1 dx_2,$$

where C is a constant independent on ε .

Proof. 1. For the proof see [19, Lemma 3].

2. For a function $v \in W^{\beta,2}(Y)$ the trace theorem implies

$$\int_{\Gamma} |v|^2 d\gamma_y \leq C \int_Y |v|^2 dy + C \int_Y \int_Y \frac{|v(y_1) - v(y_2)|^2}{|y_1 - y_2|^{n+2\beta}} dy_1 dy_2.$$

Changing variables, $y = x/\varepsilon$, we obtain

$$\int_{\varepsilon\Gamma_i} |v^\varepsilon|^2 \frac{d\gamma_x}{\varepsilon^{n-1}} \leq C \int_{\varepsilon Y_i} |v^\varepsilon|^2 \frac{dx}{\varepsilon^n} + C \int_{\varepsilon Y_i} \int_{\varepsilon Y_i} \frac{|v^\varepsilon(x_1) - v^\varepsilon(x_2)|^2}{|x_1 - x_2|^{n+2\beta}} \varepsilon^{n+2\beta} \frac{dx_1}{\varepsilon^n} \frac{dx_2}{\varepsilon^n}.$$

Multiplying the inequality side by side with ε^{-n} and summing up over i from 1 to N implies the estimate of the lemma. \square

Using a priori estimates derived in section 3.3 and the concept of the two-scale convergence, we obtain the following compactness result.

LEMMA 4.3. *There exist functions l , r_f , and r_b such that*

1. $l^\varepsilon \rightharpoonup l$ in $L^2(0, T; H^1(\Omega))$, $\partial_t l^\varepsilon \rightharpoonup \partial_t l$ in $L^2((0, T) \times \Omega)$, $l^\varepsilon \overset{*}{\rightharpoonup} l$ in $L^\infty((0, T) \times \Omega)$,
2. $l^\varepsilon \rightharpoonup l$ in $L^2(0, T; W^{\beta,2}(\Omega))$ for $\frac{1}{2} < \beta < 1$ and $\lim_{\varepsilon \rightarrow 0} \|l^\varepsilon - l\|_{L^2((0,T) \times \Gamma^\varepsilon)} = 0$,
3. $l^\varepsilon \rightharpoonup l$ two-scale, $\nabla l^\varepsilon \rightharpoonup \nabla_x l + \nabla_y l_1$ two-scale, $l_1 \in L^2((0, T) \times \Omega; H^1_{per}(Z)/\mathbb{R})$,
4. $r_f^\varepsilon \rightharpoonup r_f$, $r_b^\varepsilon \rightharpoonup r_b$ two-scale and $r_f, r_b \in L^\infty((0, T) \times \Omega \times \Gamma)$,
5. $\partial_t r_f^\varepsilon \rightharpoonup \partial_t r_f$, $\partial_t r_b^\varepsilon \rightharpoonup \partial_t r_b$ two-scale, and $\partial_t r_f, \partial_t r_b \in L^2((0, T) \times \Omega \times \Gamma)$.

Proof. From the a priori estimates in Lemma 3.4, we obtain weak convergence $l^\varepsilon \rightharpoonup l$ in $L^2(0, T; H^1(\Omega))$, $\partial_t l^\varepsilon \rightharpoonup \partial_t l$ in $L^2((0, T) \times \Omega)$, and $l^\varepsilon \overset{*}{\rightharpoonup} l$ in $L^\infty((0, T) \times \Omega)$.

To obtain strong convergence of l^ε in $L^2((0, T), W^{\beta,2}(\Omega))$, $\frac{1}{2} < \beta < 1$, we use the compact embedding of $W^{\beta,2}(\Omega)$ in $H^1(\Omega)$ and apply the Lions–Aubin lemma [27] with $B = W^{\beta,2}(\Omega)$. Applying Lemma 4.2 we obtain the inequality

$$\|l^\varepsilon\|_{\Gamma^\varepsilon}^2 \leq c \|l^\varepsilon\|_{W^{\beta,2}(\Omega^\varepsilon)}^2.$$

It follows that

$$\|l^\varepsilon - l\|_{L^2((0,T) \times \Gamma^\varepsilon)} \leq c \|l^\varepsilon - l\|_{L^2(0,T;W^{\beta,2}(\Omega^\varepsilon))} \leq c \|l^\varepsilon - l\|_{L^2(0,T;W^{\beta,2}(\Omega))} \rightarrow 0 \text{ for } \varepsilon \rightarrow 0.$$

Since l^ε is bounded in $L^2(0, T; H^1(\Omega))$, the compactness theorem (see Theorem 6.3 in Appendix 6.2) implies the two-scale convergence of l^ε to the same function l and the existence of a function $l_1 \in L^2((0, T) \times \Omega; H^1_{per}(Z)/\mathbb{R})$ such that, up to a subsequence, ∇l^ε two-scale converges to $\nabla_x l(x) + \nabla_y l_1(x, y)$.

Invoking Theorem 6.5 (see Appendix 6.2) we obtain the two-scale convergence of r_f^ε and r_b^ε to functions in $L^\infty((0, T) \times \Omega \times \Gamma)$. Due to $\|\partial_t r_f^\varepsilon\|_{L^2((0,T) \times \Gamma^\varepsilon)} \leq C$ and [37, Theorem 2.2], we conclude that $\partial_t r_f^\varepsilon \rightharpoonup v$ two-scale and $v \in L^2((0, T) \times \Omega \times \Gamma)$. Then

$$\begin{aligned} \int_0^T \int_{\Gamma \times \Omega} v \phi \, dx \, d\gamma_y \, dt &= \lim_{\varepsilon \rightarrow 0} \int_0^T \int_{\Gamma^\varepsilon} \partial_t r_f^\varepsilon \phi \, d\gamma_x \, dt \\ &= - \lim_{\varepsilon \rightarrow 0} \int_0^T \int_{\Gamma^\varepsilon} r_f^\varepsilon \partial_t \phi \, d\gamma_x \, dt = - \int_0^T \int_{\Gamma \times \Omega} r_f \partial_t \phi \, dx \, d\gamma_y \, dt. \end{aligned}$$

Consequently, we conclude that $v = \partial_t r_f$. Analogously we obtain the two-scale convergence of $\partial_t r_b^\varepsilon$ to $\partial_t r_b$. \square

4.2. Macroscopic equations.

THEOREM 4.4. *As $\varepsilon \rightarrow 0$, the sequence of solutions of the microscopic problem (1)–(6) converges to the weak solution (l, r_f, r_b) , $l \in H^1(0, T; L^2(\Omega))$, $l \in L^2(0, T; H^1(\Omega))$, $l \in L^\infty((0, T) \times \Omega)$, $r_f, r_b \in H^1(0, T; L^2(\Omega \times \Gamma))$, $r_f, r_b \in L^\infty((0, T) \times \Omega \times \Gamma)$,*

of the following macroscopic problem:

$$(14) \quad \begin{cases} \partial_t l(t, x) = -\frac{1}{|Y|} \int_{\Gamma} (b(t, y)r_f(t, x, y)l(t, x) - d(t, y)r_b(t, x, y))d\gamma_y \\ \quad + (\nabla(S(t)\nabla l(t, x))) + \tilde{p}_l(t, l(t, x)) - \tilde{\mu}_l(t)l(t, x), & t > 0, x \in \Omega, \\ \nabla l(t, x) \cdot \nu = 0, & t > 0, x \in \Gamma^N, \\ l(t, x) = l_0(x), & t = 0, x \in \Omega, \end{cases}$$

$$(15) \quad \begin{cases} \partial_t r_f(t, x, y) = p_r(t, y, r_b(t, x, y)) - b(t, y)r_f(t, x, y)l(t, x) \\ \quad + d(t, y)r_b(t, x, y) - \mu_f(t, y)r_f(t, x, y), & y \in \Gamma, x \in \Omega, \\ \partial_t r_b(t, x, y) = b(t, y)r_f(t, x, y)l(t, x) - d(t, y)r_b(t, x, y) \\ \quad - \mu_b(t, y)r_b(t, x, y), & y \in \Gamma, x \in \Omega, \\ r_f(t, x, y) = r_{f0}(x, y), & t = 0, y \in \Gamma, x \in \Omega, \\ r_b(t, x, y) = r_{b0}(x, y), & t = 0, y \in \Gamma, x \in \Omega, \end{cases}$$

where $\tilde{\mu}_l(t) = \frac{1}{|Y|} \int_Y \mu_l(t, y) dy$, $\tilde{p}(t, l) = \frac{1}{|Y|} \int_Y p(t, y, l) dy$, and the matrix S is defined as $s_{ij} = \frac{1}{|Y|} \sum_{k=1}^3 \int_Y (D_{ij}(t, y) + D_{ik}(t, y)\partial_{y_k} w_j) dy$ with w_i being the solutions of the cell problem

$$-\nabla_y(D(t, y)\nabla_y w_i) = \sum_{k=1}^3 \partial_{y_k} D_{ki}(t, y) \text{ in } Y, \quad -D(t, y)\frac{\partial w_i}{\partial \nu} = \sum_{k=1}^3 D_{ki}(t, y)\nu_k \text{ on } \Gamma.$$

Proof. To derive a limit equation for l^ε we apply a standard two-scale convergence method and strong convergence of l^ε . Using in (7) a test function of the form $\phi(t, x) = \psi_0(t, x) + \varepsilon\psi_1(t, x, \frac{x}{\varepsilon})$, $\psi_0 \in C^\infty((0, T) \times \Omega)$, $\psi_1 \in C^\infty((0, T) \times \Omega; C_{per}^\infty(Z))$ and passing to the two-scale limit applying Lemma 4.3 yields

$$\begin{aligned} & |Y| \int_0^T \int_{\Omega} \partial_t l \psi_0(t, x) dx dt + |Y| \int_0^T \int_{\Omega} \tilde{\mu}_l(t) l(t, x) \psi_0(t, x) dx dt \\ & + \int_0^T \int_{\Omega} \int_Y D(t, y)(\nabla_x l(t, x) + \nabla_y l_1(t, x, y))(\nabla_x \psi_0 + \nabla_y \psi_1) dy dx dt \\ & = - \int_0^T \int_{\Omega} \int_{\Gamma} [b(t, y)r_f(t, x, y)l(t, x) - d(t, y)r_b(t, x, y)]\psi_0(t, x) d\gamma_y dx dt \\ & + |Y| \int_0^T \int_{\Omega} \tilde{p}_l(t, l) \psi_0 dx dt. \end{aligned}$$

To show the convergence of the nonlinear term $b^\varepsilon r_f^\varepsilon l^\varepsilon$ of the boundary integral, we rewrite this integral as a sum of two integrals,

$$\begin{aligned} & \varepsilon \int_0^T \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon l^\varepsilon \left(\psi_0(t, x) + \varepsilon\psi_1 \left(t, x, \frac{x}{\varepsilon} \right) \right) d\gamma_x dt \\ & = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon l \left(\psi_0(t, x) + \varepsilon\psi_1 \left(t, x, \frac{x}{\varepsilon} \right) \right) d\gamma_x dt \\ & \quad + \varepsilon \int_0^T \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon (l^\varepsilon - l) \left(\psi_0(t, x) + \varepsilon\psi_1 \left(t, x, \frac{x}{\varepsilon} \right) \right) d\gamma_x dt. \end{aligned}$$

The first integral converges to $\int_0^T \int_\Omega \int_\Gamma b(t, y) r_f(t, x, y) l(t, x) \psi_0(t, x) d\gamma_y dx dt$ due to the two-scale convergence of r_f^ε . Since $\|l^\varepsilon - l\|_{L^2((0,T) \times \Gamma^\varepsilon)} \rightarrow 0$ as $\varepsilon \rightarrow 0$, we obtain for the second integral

$$\begin{aligned} & \varepsilon \int_0^T \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon (l^\varepsilon - l) \left(\psi_0(t, x) + \varepsilon \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) \right) d\gamma_x dt \\ & \leq \varepsilon \left(\int_0^T \int_{\Gamma^\varepsilon} |b^\varepsilon r_f^\varepsilon \psi_0|^2 d\gamma_x dt \right)^{1/2} \left(\int_0^T \int_{\Gamma^\varepsilon} |l^\varepsilon - l|^2 d\gamma_x dt \right)^{1/2} \\ & + \varepsilon^2 \left(\int_0^T \int_{\Gamma^\varepsilon} |b^\varepsilon r_f^\varepsilon \psi_1|^2 d\gamma_x dt \right)^{1/2} \left(\int_0^T \int_{\Gamma^\varepsilon} |l^\varepsilon - l|^2 d\gamma_x dt \right)^{1/2} \rightarrow 0 \text{ as } \varepsilon \rightarrow 0. \end{aligned}$$

To determinate the unknown function $l_1 \in L^2((0, T) \times \Omega; H^1_{\text{per}}(Y)/\mathbb{R})$, we set $\psi_0 = 0$ and obtain the equation

$$\int_0^T \int_{\Omega \times Y} D(t, y) (\nabla_x l(t, x) + \nabla_y l_1(t, x, y)) \nabla_y \psi_1(t, x, y) dt dx dy = 0$$

for all ψ_1 . From this it follows that l_1 depends linearly on $\nabla_x l$, and it can be written in the form

$$l_1 = \sum_{i=1}^n \frac{\partial l}{\partial x_i} \cdot w_i,$$

where the functions w_i are defined as solutions of the cell problem

$$-\nabla(D(t, y) \nabla w_i) = \sum_{k=1}^3 \partial_{y_k} D_{ki}(t, y) \text{ in } Y, \quad -D(t, y) \frac{\partial w_i}{\partial \nu} = \sum_{k=1}^3 D_{ki}(t, y) \nu_k \text{ on } \Gamma.$$

Next, setting $\psi_1 = 0$, we obtain

$$\begin{aligned} & \int_0^T \int_\Omega \int_Y \sum_{i,j=1}^n D_{ij}(t, y) (\partial_{x_i} l(t, x) + \sum_{k=1}^n \partial_{y_i} w_k \partial_{x_k} l(t, x)) \partial_{x_j} \psi_0(t, x) dy dx dt \\ & = |Y| \int_0^T \int_\Omega \sum_{i,j=1}^n s_{ij} \partial_{x_i} \psi_0(t, x) \partial_{x_j} l(t, x) dy dx dt \end{aligned}$$

with $s_{ij} = \frac{1}{|Y|} \sum_{k=1}^3 \int_Y (D_{ij}(t, y) + D_{ik}(t, y) \partial_{y_k} w_j) dy$.

The difficulty arises in passing to the limit in nonlinear terms in the ordinary differential equations on the surface of microstructures. We have to show that $p_r^\varepsilon(t, x, r_b^\varepsilon(t, x)) \rightarrow p_r(t, y, r_b(t, x, y))$ in the two-scale sense. To cope with this difficulty we apply the unfolding method (periodic modulation), developed in [7, 5, 6]. Following [5] and [6], we define a dilation operator.

DEFINITION 4.5. For a given $\varepsilon > 0$, we define a dilation operator D^ε mapping measurable functions on $(0, T) \times \Gamma^\varepsilon$ to measurable functions on $(0, T) \times \Omega \times \Gamma$ by

$$D^\varepsilon u(t, x, y) = u(t, c^\varepsilon(x) + \varepsilon y), \quad y \in \Gamma, \quad (t, x) \in (0, T) \times \Omega,$$

where $c^\varepsilon(x)$ denotes the lattice translation point of the ε -cell domain containing x , $c^\varepsilon(x) = \varepsilon[\frac{x}{\varepsilon}]$. We extend $D^\varepsilon u$ from Γ to $\bigcup_k(\Gamma + k)$ periodically.

Remark 4.1. The dilation operator D^ε is well defined for all $(t, x, y) \in (0, T) \times \Omega \times \Gamma$ under the assumption on the geometry of domain Ω^ε (cf. Remark 2.1).

To proceed, we have to establish the link between the two-scale convergence and the weak convergence of the dilated sequences. Following [6], we formulate the lemma on the convergence of $D^\varepsilon u^\varepsilon : (0, T) \times \Omega \times \Gamma \rightarrow \mathbb{R}$. We define $L^2_{\text{per}}(\Gamma)$ as the space of functions $f \in L^2(\Gamma)$ defined on Γ and periodically extended to $\Gamma^* = \bigcup_k(\Gamma + k)$.

LEMMA 4.6. *If $D^\varepsilon u^\varepsilon \rightharpoonup u^*$ weakly in $L^2((0, T) \times \Omega; L^2_{\text{per}}(\Gamma))$ and $u^\varepsilon \rightarrow u$ two-scale, then $u^* = u$ a.e. in $(0, T) \times \Omega \times \Gamma$.*

Proof. Let u^* be a weak limit of $D^\varepsilon u^\varepsilon$. Then, for a test function $\psi(t, x)h(y)$, where $\psi \in C^\infty_0((0, T) \times \Omega)$ and $h \in C^\infty_{\text{per}}(\Gamma)$, we obtain

$$\begin{aligned} & \int_0^T \int_{\Omega \times \Gamma} D^\varepsilon u^\varepsilon(t, x, y) \psi(t, x) h(y) d\gamma_y dx dt \\ & \rightarrow \int_0^T \int_{\Omega \times \Gamma} u^*(t, x, y) \psi(t, x) h(y) d\gamma_y dx dt \quad \text{as } \varepsilon \rightarrow 0. \end{aligned}$$

On the other hand, we have

$$\begin{aligned} & \int_0^T \int_{\Omega \times \Gamma} D^\varepsilon u^\varepsilon(t, x, y) \psi(t, x) h(y) d\gamma_y dx dt \\ & = \int_0^T \int_{\Omega \times \Gamma} u^\varepsilon(t, \varepsilon y + c^\varepsilon(x)) \psi(t, x) h(y) d\gamma_y dx dt \\ & = \sum_{k=1}^N \int_0^T \int_{\varepsilon(Z+k)} \int_{\Gamma} u^\varepsilon(t, \varepsilon y + c^\varepsilon(x)) \psi(t, x) h(y) d\gamma_y dx dt. \end{aligned}$$

Changing variables $z = \varepsilon(y + k)$, where $c^\varepsilon(x) = \varepsilon[\frac{x}{\varepsilon}] = \varepsilon k$, and using the periodicity of h , we obtain

$$\begin{aligned} & \int_0^T \sum_{k=1}^N \varepsilon^{-2} \int_{\varepsilon(\Gamma+k)} u^\varepsilon(t, z) h\left(\frac{z}{\varepsilon}\right) \int_{\varepsilon(Z+k)} \psi(t, x) dx d\gamma_z dt \\ & = \varepsilon \int_0^T \sum_{k=1}^N \int_{\varepsilon(\Gamma+k)} u^\varepsilon(t, z) h\left(\frac{z}{\varepsilon}\right) \psi(t, z) d\gamma_z dt + c\varepsilon^2 \\ & \rightarrow \int_0^T \int_{\Omega} \int_{\Gamma} u(t, x, y) h(y) \psi(t, x) d\gamma_y dx dt, \end{aligned}$$

since from the continuity of ψ we have the estimate

$$|\varepsilon^{-3} \int_{\varepsilon(Z+k)} (\psi(t, x) - \psi(t, z)) dx| \leq c\varepsilon \text{ for } z \in \varepsilon(\Gamma + k).$$

Therefore, we conclude that $u^* = u$ a.e. in $(0, T) \times \Omega \times \Gamma$. □

In analogy to the above lemma and Lemma 2 in [5], we can prove the following properties of the dilation operator for oscillating surfaces.

LEMMA 4.7. For $u \in L^2((0, T) \times \Gamma^\varepsilon)$

$$\|D^\varepsilon u\|_{L^2(\Omega \times \Gamma)} = \|u\|_{L^2(\Gamma^\varepsilon)}.$$

If $u \in L^2(\Omega \times \Gamma)$ is constant in y , then $D^\varepsilon u \rightarrow u$ as $\varepsilon \rightarrow 0$ strongly in $L^2(\Omega \times \Gamma)$.

Changing variables, $\Gamma^\varepsilon \ni x \rightarrow \varepsilon y + c^\varepsilon(x)$, $c^\varepsilon(x) = \varepsilon k$ for $x \in \Gamma^\varepsilon$, we obtain equations on the fixed domain $(0, T) \times \Omega \times \Gamma$,

$$\begin{aligned} \frac{\partial}{\partial t} D^\varepsilon r_f^\varepsilon(t, x, y) &= -\mu_f(t, y) D^\varepsilon r_f^\varepsilon(t, x, y) + p_r(t, y, D^\varepsilon r_b^\varepsilon(t, x, y)) \\ &\quad - b(t, y) D^\varepsilon r_f^\varepsilon(t, x, y) D^\varepsilon l^\varepsilon(t, x, y) + d(t, y) D^\varepsilon r_b^\varepsilon(t, x, y), \\ \frac{\partial}{\partial t} D^\varepsilon r_b^\varepsilon(t, x, y) &= -\mu_b(t, y) D^\varepsilon r_b^\varepsilon(t, x, y) \\ &\quad + b(t, y) D^\varepsilon r_f^\varepsilon(t, x, y) D^\varepsilon l^\varepsilon(t, x, y) - d(t, y) D^\varepsilon r_b^\varepsilon(t, x, y). \end{aligned}$$

Applying the estimates for r_f^ε and r_b^ε , we obtain the estimates for $D^\varepsilon r_f^\varepsilon$ and $D^\varepsilon r_b^\varepsilon$ and the weak convergence of $D^\varepsilon r_f^\varepsilon$ to r_f and $D^\varepsilon r_b^\varepsilon$ to r_b in $L^2((0, T) \times \Omega; L^2_{\text{per}}(\Gamma))$ (see Lemma 4.6). Since $\sup_{[0, T] \times \bar{\Omega}} |l^\varepsilon| \leq C$, we conclude that $\sup_{[0, T] \times \Omega \times \Gamma} |D^\varepsilon l^\varepsilon| \leq C$.

Now we prove the strong convergence of $D^\varepsilon r_f^\varepsilon$ and $D^\varepsilon r_b^\varepsilon$ in $L^2((0, T) \times \Omega; L^2_{\text{per}}(\Gamma))$. For this we show that $D^\varepsilon r_f^\varepsilon$ and $D^\varepsilon r_b^\varepsilon$ are Cauchy sequences. We consider the equations for $D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}$ and $D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}$, with $n > m$, multiply them side by side with $D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}$ and $D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}$, respectively, and integrate over $\Omega \times \Gamma$.

$$\begin{aligned} \frac{\partial}{\partial t} \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma &= - \int_{\Omega \times \Gamma} \mu_f(t, y) |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma \\ &\quad + \int_{\Omega \times \Gamma} (p_r(t, y, D^{\varepsilon_n} r_b^{\varepsilon_n}) - p_r(t, y, D^{\varepsilon_m} r_b^{\varepsilon_m})) (D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}) dx d\gamma \\ &\quad - \int_{\Omega \times \Gamma} b(t, y) (D^{\varepsilon_n} r_f^{\varepsilon_n} D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m} D^{\varepsilon_m} l^{\varepsilon_m}) (D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}) dx d\gamma \\ &\quad - \int_{\Omega \times \Gamma} d(t, y) (D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}) (D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}) dx d\gamma, \\ \frac{\partial}{\partial t} \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma &= - \int_{\Omega \times \Gamma} \mu_b(t, y) |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma \\ &\quad + \int_{\Omega \times \Gamma} b(t, y) (D^{\varepsilon_n} r_f^{\varepsilon_n} D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m} D^{\varepsilon_m} l^{\varepsilon_m}) (D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}) dx d\gamma \\ &\quad - \int_{\Omega \times \Gamma} d(t, y) |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma. \end{aligned}$$

Using the Young inequality, we obtain

$$\begin{aligned}
(16) \quad & \frac{\partial}{\partial t} \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma \\
& \leq C_1 \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma \\
& \quad + C_2 \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma \\
& \quad + b_1 \sup_{(0,T) \times \Omega \times \Gamma} |D^{\varepsilon_n} l^{\varepsilon_n}| \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma \\
& \quad + C_3 \int_{\Omega \times \Gamma} |D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_m} l^{\varepsilon_m}|^2 dx d\gamma,
\end{aligned}$$

$$\begin{aligned}
(17) \quad & \frac{\partial}{\partial t} \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma \\
& \leq C_4 \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}|^2 dx d\gamma \\
& \quad + b_1 \sup_{(0,T) \times \Omega \times \Gamma} |D^{\varepsilon_n} l^{\varepsilon_n}| \int_{\Omega \times \Gamma} |D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}|^2 dx d\gamma \\
& \quad + C_5 \int_{\Omega \times \Gamma} |D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_m} l^{\varepsilon_m}|^2 dx d\gamma.
\end{aligned}$$

Due to Lemma 4.7 and strong convergence of l^ε on Γ^ε , we obtain

$$\int_0^T \int_{\Omega \times \Gamma} |D^\varepsilon l^\varepsilon - D^\varepsilon l|^2 d\gamma dx dt = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} |l^\varepsilon - l|^2 d\gamma_x dt \leq C\varepsilon.$$

Therefore, since $D^{\varepsilon_n} l \rightarrow l$ strongly in $L^2((0, T) \times \Omega \times \Gamma)$ (see Lemma 4.7),

$$\begin{aligned}
& \int_0^T \int_{\Omega \times \Gamma} |D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_m} l^{\varepsilon_m}|^2 d\gamma dx dt \\
& \leq \int_0^T \int_{\Omega \times \Gamma} \left(|D^{\varepsilon_n} l^{\varepsilon_n} - D^{\varepsilon_n} l|^2 + |D^{\varepsilon_n} l - l|^2 \right) d\gamma dx dt \\
& \quad + \int_0^T \int_{\Omega \times \Gamma} \left(|D^{\varepsilon_m} l - l|^2 + |D^{\varepsilon_m} l^{\varepsilon_m} - D^{\varepsilon_m} l|^2 \right) d\gamma dx dt \\
& \leq \varepsilon_n \int_0^T \int_{\Gamma^{\varepsilon_n}} |l^{\varepsilon_n} - l|^2 d\gamma_x dt \\
& \quad + \varepsilon_m \int_0^T \int_{\Gamma^{\varepsilon_m}} |l^{\varepsilon_m} - l|^2 d\gamma_x dt + \int_0^T \int_{\Omega \times \Gamma} \left(|D^{\varepsilon_n} l - l|^2 + |D^{\varepsilon_m} l - l|^2 \right) d\gamma dx dt \\
& \leq C(\varepsilon_n + \varepsilon_m).
\end{aligned}$$

We add (16) and (17) side by side and integrate with respect to time. Using additionally the boundedness of $D^\varepsilon l^\varepsilon$ on $(0, T) \times \Omega \times \Gamma$, we obtain

$$\begin{aligned} & \|D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}\|^2 + \|D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}\|^2 \\ & \leq C_1 \int_0^T (\|D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}\|^2 + \|D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}\|^2) dt + C_2 \frac{1}{n}, \end{aligned}$$

where $C_1 = C_1(\sup_{(0,T) \times \Omega} |l^\varepsilon|, \sup_{(0,T) \times \Gamma} |\mu_f|, \sup_{(0,T) \times \Gamma} |\mu_b|, \sup_{(0,T) \times \Gamma} |b|, \sup_{(0,T) \times \Gamma} |d|, \sup_{(0,T) \times \Gamma^\varepsilon} |r_f^\varepsilon|)$. Then the Gronwall lemma yields

$$\begin{aligned} \|D^{\varepsilon_n} r_f^{\varepsilon_n} - D^{\varepsilon_m} r_f^{\varepsilon_m}\|_{L^2(\Omega \times \Gamma)} & \leq C \frac{1}{n}, \\ \|D^{\varepsilon_n} r_b^{\varepsilon_n} - D^{\varepsilon_m} r_b^{\varepsilon_m}\|_{L^2(\Omega \times \Gamma)} & \leq C \frac{1}{n}. \end{aligned}$$

Using strong convergence of $D^\varepsilon r_b^\varepsilon$, continuity of p_r , and weak convergence of $p_r(t, y, D^\varepsilon r_b^\varepsilon)$, which results from the boundedness of p_r , we obtain that $p_r(t, y, D^\varepsilon r_b^\varepsilon)$ weakly converges to $p_r(t, y, r_b(t, x, y))$ in $L^2((0, T) \times \Omega; L^2_{\text{per}}(\Gamma))$.

Now we can take the two-scale limit in the equations on the boundary,

$$\begin{aligned} & \varepsilon \int_0^T \int_{\Gamma^\varepsilon} \partial_t r_f^\varepsilon \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) d\gamma_x dt = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} p_r^\varepsilon(t, x, r_b^\varepsilon(t, x)) \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) d\gamma_x dt \\ & + \varepsilon \int_0^T \int_{\Gamma^\varepsilon} \left(-b^\varepsilon r_f^\varepsilon(t, x) l^\varepsilon(t, x) + d^\varepsilon(t, x) r_b^\varepsilon(t, x) - \mu_f^\varepsilon r_f^\varepsilon(t, x) \right) \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) d\gamma_x dt, \\ & \varepsilon \int_0^T \int_{\Gamma^\varepsilon} \partial_t r_b^\varepsilon(t, x) \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) d\gamma_x dt = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} b^\varepsilon r_f^\varepsilon(t, x) l^\varepsilon(t, x) \psi_1 d\gamma_x dt \\ & + \varepsilon \int_0^T \int_{\Gamma^\varepsilon} \left(-d^\varepsilon r_b^\varepsilon(t, x) - \mu_b^\varepsilon r_b^\varepsilon(t, x) \right) \psi_1 \left(t, x, \frac{x}{\varepsilon} \right) d\gamma_x dt. \end{aligned}$$

The linear terms converge two-scale to their limit functions. The proof of convergence for the nonlinear term $b^\varepsilon r_f^\varepsilon(t, x) l^\varepsilon(t, x)$ is the same as in the equation for l^ε . Due to boundedness of p_r^ε and Lemma 4.6, $p_r^\varepsilon(t, x, r_b^\varepsilon)$ converges two-scale to $p_r(t, y, r_b(t, x, y))$. Therefore, we obtain the macroscopic equations for r_f and r_b . \square

The uniqueness of the solution of the macroscopic problem can be proved in the same way as for the microscopic problem.

Remark 4.2. Properties of the macroscopic model: Using the framework of bounded invariant rectangles (see [44]) we can show that solutions of system (14)–(15) remain positive for positive initial conditions and that they are also uniformly bounded. This results from the assumption of the nonnegativity of the model parameters and their boundedness independent of time. Methods outlined in [44, Chapter 14] can be used without major modifications.

5. Discussion. In this work, using homogenization techniques, we studied the macroscopic limit of the microscopic model describing receptor-ligand dynamics on cell membranes and in the intercellular space. We tried to answer the question of how processes which take place in different “spaces,” such as cells membranes, intercellular space, and also intracellular space, can be described by macroscopic models operating

in homogenized space. On one hand, this work provides a justification of previously proposed models, and on the other hand it is a starting point for further models.

Comparison of the macroscopic model (14)–(15) to the previously considered receptor-based model of the form

$$(18) \quad \begin{aligned} \frac{\partial}{\partial t} r_f &= -\mu_f r_f + p_r(r_b) - br_f l + dr_b, \\ \frac{\partial}{\partial t} r_b &= -\mu_b r_b + br_f l - dr_b, \\ \frac{\partial}{\partial t} l &= \frac{1}{\gamma} \frac{\partial^2}{\partial x^2} l - \mu_l l - br_f l + p_l(l) + dr_b, \end{aligned}$$

defined on the macroscopic domain Ω , shows in which cases the “older” models can be derived from the microscopic description. Model (14)–(15) is equivalent to model (18) in the case when neither the model parameters nor the initial conditions for r_f and r_b depend on the surface variable y . It means that the processes described are homogeneous within each cell and that there is no heterogeneity in the dissociation or binding processes on the cell surfaces. For nonadherent cells one can consider receptor production, binding, dissociation, or decay to be uniformly distributed on the cell surface, which results in model coefficients being constant with respect to the surface variable y . Under such assumptions we obtain a macroscopic model, in which the integral in the equation for the ligands disappears and the only difference with respect to model (18) is that the kinetics are multiplied by a coefficient $\int_{\Gamma} d\gamma_y/|Y|$. However, there is now considerable evidence of the existence of lipid raft microdomains, called membrane rafts, which organize the membrane into specialized functional units [14, 15, 38, 39, 41]. Rafts were described mainly for T-cells and T-cell receptor [15, 39, 40], but now it is clear that they play an important role for many different receptor classes [13, 41]. There are observations that the structure of lipid rafts could control cellular processes such as signaling cascades [40, 43] and receptor synthesis and trafficking [21] as well as cell adhesion and migration [16]. Some membrane proteins are located preferentially on the raft domains, whereas others are excluded from them [14]. Such a situation corresponds to the nonhomogeneous initial distribution of receptors on the cell surface and also *de novo* production terms depending on the surface variable. Our studies show that in such a case the “older” type of receptor-based model is not relevant.

Another example of cells with nonhomogeneous membrane properties are adherent cells. In the case of adherent cells there are two types of polarity, top-bottom and front-back, and it is not easy for the ligand to get in contact with the bottom of the cell. One can imagine that receptors may be concentrated on the frontal end of the cell (this determines cell motility in the case of chemotaxis), and, therefore, all the receptor-ligand processes are nonhomogeneous within the membrane [22].

6. Appendix.

6.1. Supremum estimate for l^ε . We present here a sketch of the proof of Lemma 3.6 used in section 4.

LEMMA. *For any solution of problem (1)–(6), the following estimate holds:*

$$\|l^\varepsilon\|_{L^\infty((0,T)\times\Omega)} \leq C + 2k,$$

where C is a constant independent on ε and $k = \max\{1, \sup_{\Omega} |l_0|\}$.

Proof. To show the boundedness of l^ε we use the Moser iteration technique, described in the proof of [26, Theorem 6.15]. We choose as a test function $v = \psi(l^\varepsilon)(l^\varepsilon - k)_+$, where $\psi \geq 0$ is a bounded $C^1(\mathbb{R})$ function and which satisfies for $s > k$

$$0 \leq \frac{\psi'(s)(s - k)}{\psi(s)} \leq k_1.$$

Due to the fact that $l_0 \leq k$, we obtain

$$\begin{aligned} (19) \quad & \int_{\Omega} \int_0^{l^\varepsilon} \psi(s)(s - k)_+ ds \chi^\varepsilon dx + \int_0^\tau \int_{\Omega} (D^\varepsilon \nabla l^\varepsilon, \psi(l^\varepsilon) \nabla l^\varepsilon) \chi^\varepsilon dx dt \\ & + \int_0^\tau \int_{\Omega} (D^\varepsilon \nabla l^\varepsilon, \psi'(l^\varepsilon)(l^\varepsilon - k)_+ \nabla l^\varepsilon) \chi^\varepsilon dx dt \\ & + \int_0^\tau \int_{\Omega} \mu_i^\varepsilon l^\varepsilon \psi(l^\varepsilon)(l^\varepsilon - k)_+ \chi^\varepsilon dx dt \\ & = \int_0^\tau \int_{\Omega} p_l^\varepsilon(l^\varepsilon) \psi(l^\varepsilon)(l^\varepsilon - k)_+ \chi^\varepsilon dx dt \\ & + \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} (d^\varepsilon r_b^\varepsilon - b^\varepsilon r_f^\varepsilon l^\varepsilon) \psi(l^\varepsilon)(l^\varepsilon - k)_+ d\gamma_x dt, \end{aligned}$$

where χ is a characteristic function of Y periodically extended to Z^* , and $\chi^\varepsilon(x) = \chi(\frac{x}{\varepsilon})$. From the properties of ψ , we obtain

$$\int_k^s \psi(t)(t - k) dt \geq \frac{1}{2 + k_1} \psi(s)(s - k) \quad \text{for } s \geq k.$$

The third and fourth terms on the left-hand side of (19) are nonnegative; the third term on the right-hand side is nonpositive. Using $l^\varepsilon \geq k \geq 1$ and Lemma 4.2, we obtain the estimate

$$\begin{aligned} & \varepsilon \int_0^\tau \int_{\Gamma^\varepsilon} \psi(l^\varepsilon)(l^\varepsilon - k)_+ d\gamma_x dt \\ & \leq C \int_0^\tau \int_{\Omega^\varepsilon} \left(\psi(l^\varepsilon)(l^\varepsilon - k)_+ + \varepsilon^2 \nabla(\psi(l^\varepsilon)(l^\varepsilon - k)_+) \right) dx dt \\ & \leq C(1 + \frac{1}{\delta})(1 + k_1) \int_0^\tau \int_{\Omega} \psi(l^\varepsilon) |l^\varepsilon|^2 \chi^\varepsilon dx dt + C\varepsilon^2 \delta \int_0^\tau \int_{\Omega} \psi(l^\varepsilon) |\nabla l^\varepsilon|^2 \chi^\varepsilon dx dt. \end{aligned}$$

Then boundedness of coefficients and sublinearity of p_l yields

$$\begin{aligned} & \int_{\Omega} \psi(l^\varepsilon)(l^\varepsilon - k)^2 \chi^\varepsilon dx + d_0 \int_0^\tau \int_{\Omega} \psi(l^\varepsilon) |\nabla l^\varepsilon|^2 \chi^\varepsilon dx dt \\ & \leq C(1 + k_1)^2 \int_0^\tau \int_{\Omega} \psi(l^\varepsilon) |l^\varepsilon|^2 \chi^\varepsilon dx dt. \end{aligned}$$

Choosing $\psi(s) = (\min\{s, Z\})(1 - k/s)_+^q$, where q, Z are positive constants, applying the Gronwall and Young inequalities, and taking $Z \rightarrow \infty$ leads to

$$\int_0^T \int_{\Omega} |l^\varepsilon|^{q+2} \chi^\varepsilon dx dt \leq C(q) |\Omega| T k^{q+2}$$

for any positive q . Thus, for a fixed $q > 1$ we can choose $\psi(s) = s^{2q-2} \left(1 - \frac{k}{s}\right)_+^{(n+2)(q-1)}$ and conclude that

$$\begin{aligned} & \int_{\Omega} (l^\varepsilon)^{2q} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)q-n} \chi^\varepsilon dx + c(d_0) \int_0^\tau \int_{\Omega} (l^\varepsilon)^{2q-2} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)(q-1)} |\nabla l^\varepsilon|^2 \chi^\varepsilon dx dt \\ & \leq Cq^2 \int_0^\tau \int_{\Omega} (l^\varepsilon)^{2q} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)(q-1)} \chi^\varepsilon dx dt. \end{aligned}$$

Setting $h = (l^\varepsilon)^q \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)q-n} 1/2$ gives $|\nabla h|^2 \leq c(n)q^2 (l^\varepsilon)^{2q-2} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)(q-1)} |\nabla l^\varepsilon|^2$. Using the property of extended function of l^ε , i.e., $\|l^\varepsilon\|_{H^1(\Omega)} \leq C\|l^\varepsilon\|_{H^1(\Omega^\varepsilon)}$, with a constant C independent of ε , yields

$$\sup_{(0,T)} \int_{\Omega} h^2 dx + \int_0^T \int_{\Omega} |\nabla h|^2 dx dt \leq Cq^4 \int_0^T \int_{\Omega} (l^\varepsilon)^{2q} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)(q-1)} dx dt.$$

Invoking the Sobolev embedding theorem on $(0, T) \times \Omega$, we obtain

$$\left(\int_0^T \int_{\Omega} h^{2\kappa} dx dt \right)^{1/\kappa} \leq Cq^4 \int_0^T \int_{\Omega} (l^\varepsilon)^{2q} \left(1 - \frac{k}{l^\varepsilon}\right)_+^{(n+2)q-n-2} dx dt,$$

where $\kappa = (n + 2)/n$. Iterating the last inequality for $q = 1, \kappa, \kappa^2, \dots$, as in [26], implies that

$$\sup_{(0,T) \times \Omega} |l^\varepsilon|^2 \left(1 - \frac{k}{l^\varepsilon}\right)_+^{n+2} \leq C \int_0^T \int_{\Omega} |l^\varepsilon|^2 dx dt.$$

Considering separately the cases $\sup_{(0,T) \times \Omega} l^\varepsilon \leq 2k$ and $\sup_{(0,T) \times \Omega} l^\varepsilon \geq 2k$ results in the estimate of the lemma. \square

6.2. Two-scale convergence with parameters. We recall here the definition of two-scale convergence for functions dependent on parameters and several important results concerning this notion presented in [37]. The proofs are straightforward modifications of the proofs for the standard two-scale convergence method presented in [2].

DEFINITION 6.1. Let (u_ε) be a sequence in $L^2(\Lambda \times \Omega)$, where ε is a sequence of strictly positive numbers, which tends to zero. (u_ε) is said to two-scale converge to a (unique) limit $u_0 \in L^2(\Lambda \times \Omega \times Z)$ iff for any $\phi \in \mathcal{D}(\Lambda \times \Omega, C_{per}^\infty(Z))$ we have

$$\lim_{\varepsilon \rightarrow 0} \int_{\Lambda} \int_{\Omega} u_\varepsilon(\lambda, x) \phi\left(\lambda, x, \frac{x}{\varepsilon}\right) dx d\lambda = \int_{\Lambda} \int_{\Omega} \int_Z u_0(\lambda, x, y) \phi(\lambda, x, y) dx dy d\lambda.$$

THEOREM 6.2. From each bounded sequence (u_ε) in $L^2(\Lambda \times \Omega)$ we can extract a subsequence which two-scale converges to $u_0 \in L^2(\Lambda \times \Omega \times Z)$.

THEOREM 6.3. 1. Let (u_ε) be a bounded sequence in $L^2(\Lambda, H^1(\Omega))$, which converges weakly to a limit function $u \in L^2(\Lambda, H^1(\Omega))$. Then there exists $u_1 \in L^2(\Lambda \times \Omega, H_{per}^1(Z))$ such that, up to a subsequence, u_ε two-scale converges to u and ∇u_ε two-scale converges to $\nabla u(\lambda, x) + \nabla_y u_1(\lambda, x, y)$.

2. Let (u_ε) and $(\varepsilon \nabla u_\varepsilon)$ be bounded sequences in $L^2(\Lambda \times \Omega)$. Then there exists $u_0 \in L^2(\Lambda \times \Omega, H_{per}^1(Z))$ such that, up to a subsequence, u_ε and $\varepsilon \nabla u_\varepsilon$ two-scale converge to $u_0(\lambda, x, y)$ and $\nabla_y u_0(\lambda, x, y)$, respectively.

Now, we transfer the compactness results to the case of a sequence u_ε defined on an $(n - 1)$ -dimensional ε -periodic manifold $\Gamma^\varepsilon \in \Omega$. Let $\Gamma \in Z$ be a smooth $(n - 1)$ -dimensional manifold (in our application a sphere, $n = 3$). Then Γ^ε is the union of all $\varepsilon\Gamma$. For each Γ^ε we consider the space $L^2(\Gamma^\varepsilon)$ equipped with the scalar product $(u, v)_{\Gamma^\varepsilon} := \varepsilon \int_{\Gamma^\varepsilon} u(x)v(x)dx$.

DEFINITION 6.4 (see [37]). *A sequence of functions $(w_\varepsilon) \in L^2(\Lambda \times \Gamma^\varepsilon)$ is said to two-scale converge to a limit $w \in L^2(\Lambda \times \Omega \times \Gamma)$ iff for any $\psi \in \mathcal{D}(\Lambda \times \Omega, C_{per}^\infty(\Gamma))$ we have*

$$\lim_{\varepsilon \rightarrow 0} \varepsilon \int_{\Lambda} \int_{\Gamma^\varepsilon} w^\varepsilon(\lambda, x) \psi \left(\lambda, x, \frac{x}{\varepsilon} \right) d\gamma_x d\lambda = \int_{\Lambda} \int_{\Omega} \int_{\Gamma} w(\lambda, x, y) \psi(\lambda, x, y) dx d\gamma_y d\lambda.$$

THEOREM 6.5. 1. *From each bounded sequence (w_ε) in $L^2(\Lambda \times \Gamma^\varepsilon)$ we can extract a subsequence which two-scale converges to $w \in L^2(\Lambda \times \Omega \times \Gamma)$.*

2. *If the sequence (w_ε) is bounded in $L^\infty(\Lambda \times \Gamma^\varepsilon)$, then the limit w belongs to $L^\infty(\Lambda \times \Omega \times \Gamma)$.*

Proof. For the proof of 1, see [37].

2. We know that if w^ε is bounded in $L^2((0, T) \times \Gamma^\varepsilon)$, then there exists $w \in L^2((0, T) \times \Omega \times \Gamma)$ such that $w^\varepsilon \rightarrow w$ two-scale [37]. Now we use the proof of that theorem and show that if w^ε is bounded in $L^\infty((0, T) \times \Gamma^\varepsilon)$, then $w^\varepsilon \rightarrow w$ two-scale and $w \in L^\infty((0, T) \times \Omega \times \Gamma)$. We define $\mu_\varepsilon(\phi) = \varepsilon \int_0^T \int_{\Gamma^\varepsilon} w^\varepsilon(t, x) \phi(t, x, \frac{x}{\varepsilon}) d\gamma_x^\varepsilon dt$ and obtain

$$|\mu_\varepsilon(\phi)| \leq \|w^\varepsilon\|_{L^2((0, T) \times \Gamma^\varepsilon)} \left(\int_0^T \int_{\Gamma^\varepsilon} \varepsilon \left| \phi \left(x, \frac{x}{\varepsilon} \right) \right|^2 d\gamma_x^\varepsilon dt \right)^{\frac{1}{2}} \leq c \|\phi\|_{C^0([0, T] \times \bar{\Omega}; C_{per}^0(\Gamma))}.$$

Therefore, $\{\mu_\varepsilon\}$ is a bounded sequence of functionals on $C^0([0, T] \times \bar{\Omega}; C_{per}^0(\Gamma))$. Since this space is a separable Banach space, there exists a subsequence of μ_ε that converges weakly* to μ . Using the boundedness of w^ε and a variant of the oscillation lemma [2], we obtain

$$|\mu(\phi)| = \lim_{\varepsilon \rightarrow 0} |\mu_\varepsilon(\phi)| \leq C \lim_{\varepsilon \rightarrow 0} \left(\int_0^T \int_{\Gamma^\varepsilon} \varepsilon \left| \phi \left(x, \frac{x}{\varepsilon} \right) \right|^2 d\gamma_x^\varepsilon dt \right)^{\frac{1}{2}} = c \|\phi\|_{L^2((0, T) \times \Omega \times \Gamma)}.$$

Therefore, μ is a bounded functional on $L^2((0, T) \times \Omega \times \Gamma)$. The Riesz representation theorem implies the existence of a function $w \in L^2((0, T) \times \Omega \times \Gamma)$. Furthermore, $\|w^\varepsilon\|_{L^\infty((0, T) \times \Gamma^\varepsilon)} \leq C$ yields

$$|\mu(\phi)| = \lim_{\varepsilon \rightarrow 0} |\mu_\varepsilon(\phi)| \leq C \lim_{\varepsilon \rightarrow 0} \int_0^T \int_{\Gamma^\varepsilon} \varepsilon \left| \phi \left(x, \frac{x}{\varepsilon} \right) \right| d\gamma_x^\varepsilon dt = c \|\phi\|_{L^1((0, T) \times \Omega \times \Gamma)}.$$

Finally, we conclude

$$\begin{aligned} \|w\|_{L^\infty((0, T) \times \Omega \times \Gamma)} &= \frac{\langle w, \phi \rangle}{\|\phi\|_{L^1((0, T) \times \Omega \times \Gamma)}} \\ &= \frac{|\mu(\phi)|}{\|\phi\|_{L^1((0, T) \times \Omega \times \Gamma)}} \leq \frac{C \|\phi\|_{L^1((0, T) \times \Omega \times \Gamma)}}{\|\phi\|_{L^1((0, T) \times \Omega \times \Gamma)}} = C. \quad \square \end{aligned}$$

THEOREM 6.6 (see [37]). *Let (u_ε) and $(\varepsilon \nabla u_\varepsilon)$ be bounded sequences in $L^2(\Lambda \times \Gamma^\varepsilon)$. Then there exists $u_0 \in L^2(\Lambda \times \Omega, H_{per}^1(\Gamma))$ such that, up to a subsequence, u_ε and $\varepsilon \nabla u_\varepsilon$, two-scale converge to $u_0(\lambda, x, y)$ and $\nabla_y u_0(\lambda, x, y)$, respectively.*

REFERENCES

- [1] E. ACERBI, V. CHIADO PIAT, G. DAL MASO, AND D. PERCIVALE, *An extension theorem from connected sets, and homogenization in general periodic domains*, *Nonlinear Anal.*, 18 (1992), pp. 481–496.
- [2] G. ALLAIRE, *Homogenization and two-scale convergence*, *SIAM J. Math. Anal.*, 23 (1992), pp. 1482–1518.
- [3] G. ALLAIRE, A. DAMLAMIAN, AND U. HORNUNG, *Two-scale convergence on periodic surfaces and applications*, in *Proceedings of the International Conference on Mathematical Modelling of Flow through Porous Media*, A. Bourgeat et al., eds., World Scientific, Singapore, 1996, pp. 15–25.
- [4] G. ALLAIRE AND F. MURAT, *Homogenization of the Neumann problem with nonisolated holes*, *Asymptotic Anal.*, 7 (1993), pp. 81–95.
- [5] T. ARBOGAST, J. DOUGLAS, JR., AND U. HORNUNG, *Derivation of the double porosity model of single phase flow via homogenization theory*, *SIAM J. Math. Anal.*, 21 (1990), pp. 823–836.
- [6] A. BOURGEAT, S. LUCKHAUS, AND A. MIKELIĆ, *Convergence of the homogenization process for a double-porosity model of immiscible two-phase flow*, *SIAM J. Math. Anal.*, 27 (1996), pp. 1520–1543.
- [7] D. CIORANESCU, A. DAMLAMIAN, AND G. GRISO, *Periodic unfolding and homogenization*, *C. R. Acad. Sci. Paris Sér. I Math.*, 335 (2002), pp. 99–104.
- [8] D. CIORANESCU AND J. SAINT JEAN PAULIN, *Homogenization in open sets with holes*, *J. Math. Anal. Appl.*, 71 (1979), pp. 590–607.
- [9] C. CONCA, J. I. DIAZ, A. LINAN, AND C. TIMOFTE, *Homogenization in chemical reactive flows*, *Electron. J. Differential Equations*, 40 (2004), pp. 1–22.
- [10] K. A. DUCA, V. LAM, I. KEREN, E. E. ENDLER, G. J. LETCHWORTH, I. S. NOVELLA, AND J. YIN, *Quantifying viral propagation in vitro: Toward a method for characterization of complex phenotypes*, *Biotechnology Progress*, 17 (2001), pp. 1156–1165.
- [11] L. DUNG, *Remarks on Hölder continuity for parabolic equations and convergence to global attractors*, *Nonlinear Anal.*, 41 (2000), pp. 921–941.
- [12] L. DUNG AND H. SMITH, *Strong positivity of solutions to parabolic and elliptic equations on nonsmooth domains*, *J. Math. Anal. Appl.*, 275 (2002), pp. 208–221.
- [13] J. I. ELLIOTT, A. SURPRENANT, F. M. MARELLI-BERG, J. C. COOPER, R. CASSADY-CAIN, C. WOODING, K. LINTON, D. R. ALEXANDER, AND C. F. HIGGINS, *Membrane phosphatidylserine distribution as a non-apoptotic signalling mechanism in lymphocytes*, *Nature Cell Biol.*, 7 (2005), pp. 808–816.
- [14] K. GAUS, M. RODRIGUEZ, K. R. RUBERU, I. GELISSEN, T. M. SLOANE, L. KRITHARIDES, AND W. JESSUP, *Domain-specific lipid distribution in macrophage plasma membranes*, *J. Lipid Res.*, 46 (2005), pp. 1526–1538.
- [15] O. GLEBOV AND B. J. NICHOLS, *Lipid raft proteins have a random distribution during localized activation of the T-cell receptor*, *Nature Cell Biol.*, 6 (2004), pp. 238–243.
- [16] C. GOMEZ-MOUNTON, J. L. ABAD, E. MIRA, R. A. LACALLE, E. GALLARO, S. JIMENEZ-BARANDA, I. ILLA, A. BERNARD, S. MANES, AND A. C. MARTINEZ, *Segregation of leading-edge and uropod components into specific lipid rafts during T cell polarization*, *Proc. Natl. Acad. Sci. USA*, 98 (2001), pp. 9642–9647.
- [17] A. GRATCHEV, *Personal communication*, Medical Center in Mannheim, University of Heidelberg, 2005.
- [18] U. HORNUNG, *Homogenization and Porous Media*, Springer-Verlag, New York, 1997.
- [19] U. HORNUNG AND W. JÄGER, *Diffusion, convection, adsorption and reaction of chemicals in porous media*, *J. Differential Equations*, 92 (1991), pp. 199–225.
- [20] U. HORNUNG, W. JÄGER, AND A. MIKELIĆ, *Reactive transport through an array of cells with semi-permeable membranes*, *RAIRO Modél. Math. Anal. Numér.*, 28 (1994), pp. 59–94.
- [21] E. IKONEN, *Roles of lipid rafts in membrane transport*, *Curr. Opin. Cell Biol.*, 13 (2001), pp. 470–477.
- [22] B. JOHNSTON AND E. C. BUTCHER, *Chemokines in rapid leukocyte adhesion triggering and migration*, *Immunology*, 14 (2002), pp. 83–92.
- [23] O. A. LADYZENSKAJA, V. A. SOLONNIKOV, AND N. N. URALCEVA, *Linear and Quasi-Linear Equations of Parabolic Type*, AMS, Providence, RI, 1968.
- [24] D. A. LAUFFENBURGER AND J. J. LINDERMAN, *Receptors. Models for Binding, Trafficking, and Signaling*, Oxford University Press, New York, 1993.
- [25] P. LI, P. SELVARAJ, AND C. ZHU, *Analysis of Competition binding between soluble and membrane-bound ligands for cell surface receptors*, *Biophys. J.*, 77 (1999), pp. 3394–3408.

- [26] G. M. LIEBERMAN, *Second Order Parabolic Differential Equations*, World Scientific, Singapore, 1996.
- [27] J. L. LIONS, *Quelques méthodes de résolution des problèmes aux limites non linéaires*, Dunod, Paris, 1969.
- [28] J. LITTLE, *Genomic instability and bystander effects: A historical perspective*, *Oncogene*, 22 (2003), pp. 6978–6987.
- [29] A. MARCINIAK-CZUCHRA, *Receptor-based models with diffusion-driven instability for pattern formation in hydra*, *J. Biol. Sys.*, 11 (2003), pp. 293–324.
- [30] A. MARCINIAK-CZUCHRA, *Receptor-based models with hysteresis for pattern formation in hydra*, *Math. Biosci.*, 199 (2005), pp. 97–119.
- [31] J. L. MARTIEL AND A. GOLDBETER, *A model based on receptor desensitization for cyclic AMP signaling in Dictyostelium cells*, *Biophys. J.*, 57 (1987), pp. 807–828.
- [32] P. B. MONK AND H. G. OTHMER, *Cyclic AMP oscillations in suspensions of Dictyostelium Discoideum*, *Philos. Trans. Roy. Soc. London Ser. B*, 323 (1989), pp. 185–224.
- [33] W. A. MÜLLER, *Pattern control in hydra: Basic experiments and concepts*, in *Experimental and Theoretical Advances in Biological Pattern Formation*, H. G. Othmer, P. K. Maini, and J. D. Murray, eds., Plenum Press, New York, 1993.
- [34] W. A. MÜLLER, *Developmental Biology*, Springer-Verlag, New York, 1997.
- [35] J. MURRAY, *Mathematical Biology*, Springer-Verlag, Berlin, 2003.
- [36] G. NGUETSENG, *A general convergence result for a functional related to the theory of homogenization*, *SIAM J. Math. Anal.*, 20 (1989), pp. 608–623.
- [37] M. NEUSS-RADU, *Some extensions of two-scale convergence*, *C. R. Acad. Sci. Paris Sér. I Math.*, 332 (1996), pp. 899–904.
- [38] P. O'SHEA, *Physical landscapes in biological membranes: Physico-chemical terrains for spatio-temporal control of biomolecular interactions and behaviour*, *Philos. Trans. Roy. Soc. London Ser. A*, 363 (2005), pp. 575–588.
- [39] I. PECHT AND D. M. GAKAMSKY, *Spatial coordination of CD8 and TCR molecules controls antigen recognition by CD8⁺ T-cells*, *FEBS Lett.*, 579 (2005), pp. 3336–3341.
- [40] T. M. RAZZAQ, P. OZEGBE, E. C. JURY, P. SEMBI, N. M. BLACKWELL, AND P. S. KABOURIDIS, *Regulation of T-cell receptor signalling by membrane microdomains*, *Immunology*, 113 (2004), pp. 413–426.
- [41] P. B. SEHGAL, *Plasma membrane rafts and chaperones in cytokine STAT signaling*, *Acta Bioch. Pol.*, 50 (2003), pp. 583–594.
- [42] J. A. SHERRATT, P. K. MAINI, W. JÄGER, AND W. MÜLLER, *A receptor-based model for pattern formation in hydra*, *Forma*, 10 (1995), pp. 77–95.
- [43] K. SIMONS AND D. TOOMRE, *Lipid rafts and signal transduction*, *Nat. Rev. Mol. Cell Biol.*, 1 (2000), pp. 31–39.
- [44] J. SMOLLER, *Shock-Waves and Reaction-Diffusion Equations*, Springer-Verlag, New York, 1994.
- [45] W. WALTER, *Gewöhnliche Differentialgleichungen*, Springer-Verlag, Berlin, 1996.
- [46] H. WEARING AND J. A. SHERRATT, *Keratinocyte growth factor signalling: A mathematical model of dermal-epidermal interaction in epidermal wound healing*, *Math. Biosci.*, 165 (2000), pp. 41–62.
- [47] J. WLOKA, *Partielle Differentialgleichungen*, Teubner Verlag, Stuttgart, 1982.

INVERSE PROBLEM OF DETERMINING THE DENSITY AND TWO LAMÉ COEFFICIENTS BY BOUNDARY DATA*

M. BELLASSOUED[†], O. IMANUVILOV[‡], AND M. YAMAMOTO[§]

Abstract. In this paper we study an inverse problem of determining spatially varying density and two Lamé coefficients by a single measurement of solution in a subboundary over a time interval. By assuming that, in a neighborhood of the boundary of the spatial domain, the density and the Lamé coefficients are known, we prove a logarithmic stability estimate for the inverse problem with a single measurement of data on an arbitrarily given subboundary.

Key words. Carleman estimate, Lamé system, inverse problem, Fourier–Bros–Iagolnitzer transform

AMS subject classifications. 35B60, 35R25, 35R30, 74B05

DOI. 10.1137/070679971

1. Introduction. This paper is concerned with the global stability in determining density and two Lamé coefficients in the classical isotropic elasticity system from data of the solution on a subboundary over a time interval. We will formulate our problem as follows: In a bounded domain $\Omega \subset \mathbb{R}^3$ with sufficiently smooth boundary $\Gamma = \partial\Omega$, we consider the isotropic elasticity system

$$(1.1) \quad \rho(x)\partial_t^2 \mathbf{u}(x, t) - \mathcal{L}_{\lambda, \mu}(x, \partial_x) \mathbf{u}(x, t) = 0, \quad (x, t) \in Q = \Omega \times (-T, T),$$

where

$$\begin{aligned} \mathcal{L}_{\lambda, \mu}(x, \partial_x) \mathbf{v}(x) &\equiv \mu(x)\Delta \mathbf{v}(x) + (\mu(x) + \lambda(x)) (\nabla \operatorname{div} \mathbf{v}(x)) \\ &\quad + (\operatorname{div} \mathbf{v}(x)) \nabla \lambda(x) + (\nabla \mathbf{v} + (\nabla \mathbf{v})^T) \nabla \mu(x), \quad x \in \Omega. \end{aligned}$$

Throughout this paper, t and $x = (x_1, x_2, x_3)$ denote the time variable and the spatial variable, respectively, and $\mathbf{u} = (u_1, u_2, u_3)^T$ denotes the displacement at the location x and the time t , where \cdot^T denotes the transpose of matrices. We will assume that the density ρ and the Lamé parameters μ and λ satisfy

$$\rho, \lambda, \mu \in C^3(\bar{\Omega}), \quad \rho(x) > 0, \quad \mu(x) > 0, \quad \lambda(x) + \mu(x) > 0 \quad \text{for } x \in \bar{\Omega}.$$

To system (1.1), we attach initial and boundary conditions:

$$(1.2) \quad \mathbf{u} = \Phi, \quad \partial_t \mathbf{u} = \Psi \quad \text{on } \Omega \times \{0\}$$

*Received by the editors January 12, 2007; accepted for publication (in revised form) September 5, 2007; published electronically April 16, 2008.

<http://www.siam.org/journals/sima/40-1/67997.html>

[†]Department of Mathematics, Faculty of Sciences of Bizerte, 7021 Jarzouna Bizerte, Tunisia (mourad.bellassoued@fsb.rnu.tn).

[‡]Department of Mathematics, Colorado State University, 101 Weber Building, Fort Collins, CO 80523 (oleg@math.colostate.edu). This author was supported by NSF grant DMS 0205148.

[§]Department of Mathematical Sciences, The University of Tokyo, 3-8-1 Komaba, Meguro, Tokyo 153, Japan (myama@ms.u-tokyo.ac.jp). This author was supported partially by grant 15340027 from the Japan Society for the Promotion of Science and grant 17654019 from the Ministry of Education, Culture, Sports and Technology.

and

$$(1.3) \quad \mathbf{u} = \mathbf{g} \quad \text{on} \quad \Sigma \equiv \Gamma \times (-T, T).$$

There exists a unique weak solution $\mathbf{u} \in C([-T, T]; \mathbf{H}_0^1(\Omega))$ for suitable Φ, Ψ, \mathbf{g} under sufficient compatibility conditions, and by $\mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g})$ we denote the solution to (1.1)–(1.3).

The main subject of this paper is the inverse problem of determining $\rho = \rho(x)$, $\lambda = \lambda(x)$, and $\mu = \mu(x)$ from observed data of the solution on a part of the boundary. It is an important problem, for example, in the geophysics to determine ρ, λ , and μ inside an elastic body from measurements on a subboundary. Thus our inverse problem is physically motivated.

1.1. Inverse problem. Let $\Gamma_1 \subset \Gamma$ be given arbitrarily, and let $\Phi_j, \Psi_j, \mathbf{g}_j$, $1 \leq j \leq \mathcal{N}$, be appropriately given. Then we want to determine $\lambda(x), \mu(x), \rho(x)$, $x \in \Omega$, by measurements

$$\partial_\nu \mathbf{u}(\lambda, \mu, \rho, \Phi_j, \Psi_j, \mathbf{g}_j)(x, t), \quad (x, t) \in \Sigma_1 \equiv \Gamma_1 \times (-T, T), \quad 1 \leq j \leq \mathcal{N}.$$

Here $\nu = \nu(x)$ denotes the unit outward normal vector and $\partial_\nu = \nabla \cdot \nu$.

Our formulation of the inverse problem requires only a finite number of observations (i.e., $\mathcal{N} < \infty$). As for inverse problems for a nonstationary Lamé system by infinitely many boundary observations (i.e., Dirichlet-to-Neumann map), we refer to Rachele [42], for example. Moreover see a monograph by Yakhno [48] for inverse problems for the Lamé system.

For the formulation with a finite number of observations, Bukhgeim and Klibanov [10] proposed a remarkable method based on a Carleman estimate and established the uniqueness for similar inverse problems for scalar partial differential equations. See also Baudouin and Puel [2], Bellassoued [4], [5], Bellassoued and Yamamoto [6], [7], Bukhgeim [8], Bukhgeim, Cheng, Isakov, and Yamamoto [9], Imanuvilov and Yamamoto [20], [21], [22], [23], [24], Isakov [25], [26], Isakov and Yamamoto [27], Khaidarov [29], Klibanov [30], [31], Klibanov and Timonov [33], Klibanov and Yamamoto [34], Kubo [35], Li [39], Puel and Yamamoto [40], [41], and Yamamoto [49].

A Carleman estimate is an inequality for a solution to a partial differential equation with weighted L^2 -norm and is a strong tool also for proving the uniqueness in the Cauchy problem or the unique continuation for a partial differential equation with nonanalytic coefficients. Moreover Carleman estimates have been applied essentially for estimating the energy (e.g., Kazemi and Klibanov [28], Klibanov and Malinsky [32]), while we refer to [1] as another method for the energy estimate, which is, however, not applicable to our inverse problem.

As a pioneering work concerning a Carleman estimate, we refer to Carleman's paper [11], which proved what is now called a Carleman estimate and applied it for proving the uniqueness in the Cauchy problem for a two-dimensional elliptic equation. Since [11], the theory of Carleman estimates has been studied extensively. We refer to a general theory by Hörmander [14] in the case where the symbol of a partial differential equation is isotropic and functions under consideration have compact supports (that is, they and their derivatives of suitable orders vanish on the boundary of a domain). Later Carleman estimates for functions with compact supports have been obtained for partial differential operators with anisotropic symbols by Isakov [26]. For Carleman estimates for functions without compact supports, see Imanuvilov [17] and Tataru [46]. We further refer to Fursikov and Imanuvilov [13] and Imanuvilov [16]. As

for a direct derivation of pointwise Carleman estimates for hyperbolic equations which are applicable to functions without compact supports, see Klivanov and Timonov [33] and Lavrent'ev, Romanov, and Shishat'skiĭ [36].

The Carleman estimate for the nonstationary Lamé system was obtained for functions with compact supports, by Eller, Isakov, Nakamura, and Tataru [12], Ikehata, Nakamura, and Yamamoto [15], Imanuvilov, Isakov, and Yamamoto [18], and Isakov [25]. Lemma 2.3 is our Carleman estimate for the Lamé system whose solutions have not necessarily compact supports, and the Carleman estimate bounds also the second-order x -derivatives, which is possible because the right-hand side is estimated in the weighted H^1 -norm in x , and we can use the identity $\Delta v = -\operatorname{rot} \operatorname{rot} v + \nabla(\operatorname{div} v)$ and the a priori estimate for the Dirichlet problem for the Laplace equation. By the methodology by [10] or [22] with such Carleman estimates, several uniqueness and stability results are available for the inverse problem for the Lamé system (1.1). That is, in [25] Isakov established the uniqueness in determining a single coefficient $\rho(x)$, using four measurements (i.e., $\mathcal{N} = 4$). Later [15] reduced the number of measurements to three. Recently [18] proved conditional stability and the uniqueness in the determination of all three functions λ , μ , and ρ , with two measurements, and Imanuvilov and Yamamoto [23], [24] proved conditional stability results with a single measurement, provided that initial data satisfy some nondegeneracy condition.

In all of the works [15], [18], [23], [24], [25], the authors assume some geometric condition of the observation subboundary. For such a kind of inverse problems, the uniqueness as well as the stability with boundary measurement on an arbitrary part of Γ are open problems. As for the corresponding unique continuation, we can refer to Bellassoued [3], Robbiano [43], and Tataru [46]. However, their methods are not applicable to the inverse problem. In our paper, by assuming that coefficients under consideration are given in a neighborhood of Γ , we will prove a stability result in the inverse problem. The coincidence of coefficients near the boundary is technically restrictive but acceptable from practical viewpoints, because one can directly know physical properties near the boundary.

Our main achievement of this paper is that we can take an arbitrary observation subboundary Γ_1 for the stability estimate. The key idea is a combination of the method (e.g., [10], [23]) by the Carleman estimates and the Fourier–Bros–Iagolnitzer (FBI) transformation which was used for sharp unique continuation by Robbiano [44]. More precisely, we apply the FBI transformation to change the problem near the boundary into a problem to which elliptic estimates can be applied.

1.2. Notation and statement of main results. In order to formulate our results, we need to introduce some notation. Let $x_0 \in \mathbb{R}^3 \setminus \bar{\Omega}$, $M_0 \geq 0$, $0 < \theta_0 \leq 1$, and $\theta_1 > 0$ be arbitrarily fixed, and let us introduce the conditions on a scalar function p :

$$(1.4) \quad \begin{cases} p(x) \geq \theta_1 > 0, & x \in \bar{\Omega}, \\ \|p\|_{C^3(\bar{\Omega})} \leq M_0, & \frac{(\nabla p(x) \cdot (x - x_0))}{2p(x)} \leq 1 - \theta_0, \quad x \in \overline{\Omega \setminus \omega}. \end{cases}$$

Next we define an admissible set of unknown coefficients λ , μ , ρ . Let $\omega \subset \Omega$ be a given arbitrary neighborhood of the boundary Γ . For fixed functions ρ_0 , λ_0 , μ_0 on ω

and Φ, Ψ in Ω , and a given constant $M_1 > 0$, we set $\Lambda = \Lambda_{M_0, M_1, \theta_0, \theta_1}$

$$(1.5) \quad \Lambda = \left\{ (\lambda, \mu, \rho) \in (C^3(\bar{\Omega}))^3; (\rho, \lambda, \mu) = (\rho_0, \lambda_0, \mu_0) \text{ in } \omega, \right. \\ \left. \left(\frac{\lambda + 2\mu}{\rho} \right), \left(\frac{\mu}{\rho} \right) \text{ satisfy (1.4), } \|\mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g})\|_{W^{8, \infty}(Q)} \leq M_1 \right\}.$$

Throughout this paper, let \mathbf{I}_3 be the 3×3 identity matrix. We note that $\mathcal{L}_{\lambda, \mu}(x, \partial_x)\Phi(x)$ is the 3-column vector for a 3-column vector Φ . Moreover by $\{\mathbf{a}\}_j$ we denote the matrix (or vector) obtained from \mathbf{a} after deleting the j th row, $\text{det}_j A$ means $\text{det} \{A\}_j$ for a square matrix A , and $\langle A \rangle_j$ is the matrix which is obtained from A by deleting the j th column of A . Furthermore we assume that the observation data are measured by the norm:

$$(1.6) \quad \epsilon(\Sigma_1) = \sum_{|\alpha|=1}^2 \|\partial_x^\alpha \mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g}) - \partial_x^\alpha \mathbf{u}(\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}, \Phi, \Psi, \mathbf{g})\|_{L^2(\Sigma_1)}^2$$

for $(\lambda, \mu, \rho), (\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}) \in \Lambda$.

1.3. Hypotheses (\mathcal{H}_1) – (\mathcal{H}_2) . Let (λ, μ, ρ) be an arbitrarily fixed element of Λ . For $\Phi = (\phi_1, \phi_2, \phi_3)^T$ and $\Psi = (\psi_1, \psi_2, \psi_3)^T$, we assume that there exist $j_1, j_2, \in \{1, \dots, 6\}$ such that for all $x \in \bar{\Omega}$

$$(\mathcal{H}_1) \quad \text{det}_{j_1} \begin{pmatrix} \mathcal{L}_{\lambda, \mu}(x, \partial_x)\Phi(x) & (\text{div } \Phi(x))\mathbf{I}_3 & (\nabla\Phi(x) + (\nabla\Phi(x))^T)(x - x_0) \\ \mathcal{L}_{\lambda, \mu}(x, \partial_x)\Psi(x) & (\text{div } \Psi(x))\mathbf{I}_3 & (\nabla\Psi(x) + (\nabla\Psi(x))^T)(x - x_0) \end{pmatrix} \neq 0,$$

$$(\mathcal{H}_2) \quad \text{det}_{j_2} \begin{pmatrix} \mathcal{L}_{\lambda, \mu}(x, \partial_x)\Phi(x) & \nabla\Phi(x) + (\nabla\Phi(x))^T & (\text{div } \Phi)(x - x_0) \\ \mathcal{L}_{\lambda, \mu}(x, \partial_x)\Psi(x) & \nabla\Psi(x) + (\nabla\Psi(x))^T & (\text{div } \Psi)(x - x_0) \end{pmatrix} \neq 0.$$

Now we are ready to state the main result, which proves that only one observation with a suitable initial value yields the logarithmic conditional stability for our inverse problem.

THEOREM 1.0. *Let $T > 0$ be sufficiently large for Ω, ω , and let $\Lambda = \Lambda_{M_0, M_1, \theta_0, \theta_1}$ be defined by (1.5). Moreover let (Φ, Ψ) satisfy the conditions (\mathcal{H}_1) – (\mathcal{H}_2) . Then there exist constants $C > 0$ and $\kappa \in (0, 1)$ such that the following estimate holds:*

$$\|\tilde{\lambda} - \lambda\|_{H^2(\Omega)} + \|\tilde{\mu} - \mu\|_{H^2(\Omega)} + \|\tilde{\rho} - \rho\|_{H^1(\Omega)} \leq C \left[\log \left(2 + \frac{C}{\epsilon(\Sigma_1)} \right) \right]^{-\kappa}$$

for any $(\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}) \in \Lambda$.

Here we note that $\epsilon(\Sigma_1)$ is given by (1.6) and the constants C and $\kappa \in (0, 1)$ are dependent on, $\Omega, \omega, T, M_0, M_1$ and independent of $(\lambda, \mu, \rho) \in \Lambda$.

Our stability result requires only one measurement, $\mathcal{N} = 1$, and the stability result is of logarithmic rate and weaker than any Hölder stability. We notice that, with a suitable geometric condition of the observation subboundary, we can prove Hölder (or the Lipschitz) stability by means of the method in Imanuvilov and Yamamoto [23], [24]. By Theorem 1.0, we can readily derive the uniqueness in the inverse problem.

COROLLARY 1.1 (uniqueness). *Under the assumptions in Theorem 1, for all $(\lambda, \mu, \rho), (\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}) \in \Lambda$, we have the uniqueness*

$$\partial_\nu \mathbf{u}(x, t) = \partial_\nu \tilde{\mathbf{u}}(x, t), \quad (x, t) \in \Sigma_1 \quad \text{implies} \quad (\lambda, \mu, \rho) = (\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}) \quad \text{in } \Omega.$$

For the determination of the three coefficients by a single measurement, we have to choose initial data satisfying conditions (\mathcal{H}_1) – (\mathcal{H}_2) . Thus conditions (\mathcal{H}_1) – (\mathcal{H}_2) are not generic properties, and we should satisfy them artificially and a posteriori. Moreover, as the following example shows, we can take such Φ and Ψ .

Example of Ω, Φ, Ψ meeting (\mathcal{H}_1) – (\mathcal{H}_2) . For simplicity, we assume that $x_0 = (0, 0, 0)$, $\bar{\Omega}$ does not intersect any of $\{x_1 = 0\}, \{x_2 = 0\}, \{x_3 = 0\}$, and $\{x_1 + x_3 = 0\}$, and λ, μ are positive constants. For example, we take

$$\Phi(x) = \begin{pmatrix} 0 \\ x_1 x_2 \\ 0 \end{pmatrix}, \quad \Psi(x) = \begin{pmatrix} x_2^2 \\ 0 \\ x_2^2 \end{pmatrix}.$$

Then, by choosing $j_1 = j_2 = 6$, we can verify that (\mathcal{H}_1) – (\mathcal{H}_2) are satisfied.

Thanks to the extra information $(\lambda, \mu, \rho) = (\tilde{\lambda}, \tilde{\mu}, \tilde{\rho})$ in a neighborhood ω of $\partial\Omega$, the sharp unique continuation by Bellassoued [3], implies $\mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g})(x, t) = \mathbf{u}(\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}, \Phi, \Psi, \mathbf{g})(x, t)$ on $\partial(\Omega \setminus \bar{\omega}) \times (-T, T)$, provided that $T > 0$ is sufficiently large. Therefore the method in Imanuvilov, Isakov, and Yamamoto[18], and Imanuvilov and Yamamoto [24] directly yields the uniqueness in our inverse problem. However, our main result is concerned with the stability in the inverse problem, and the direct combination of the existing results in [18], [21], [22] does not work. For our purpose, we will use the FBI transformation according to Robbiano [44], [43].

The remainder of the paper is organized as follows. In section 2, we give key estimates. In section 3, we prove Theorem 1.0 on the basis of the weak observation estimate, that is, an estimate of $\mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g}) - \mathbf{u}(\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}, \Phi, \Psi, \mathbf{g})$ by data on any small part of the boundary. Section 4 is devoted to the proof of the weak observation estimate.

2. Preliminaries and Carleman estimates. We set

$$\omega(\epsilon) = \{x \in \Omega; \text{dist}(x, \partial\Omega) \leq \epsilon\}$$

and

$$\omega(\epsilon_1, \epsilon_2) = \{x \in \Omega; \epsilon_1 \leq \text{dist}(x, \partial\Omega) \leq \epsilon_2\},$$

with $0 < \epsilon_1 < \epsilon_2$ and $\epsilon > 0$.

In this section we first derive several estimates. We choose $\epsilon_0, \epsilon_1, \epsilon_2 > 0$ such that

$$(2.1) \quad \omega(3\epsilon_0) \subset \omega$$

and

$$(2.2) \quad \omega(\epsilon_1, \epsilon_2) \subset \omega, \quad \epsilon_1 < \epsilon_2 < 8\epsilon_0.$$

We set

$$\omega_T(\epsilon) = \omega(\epsilon) \times [-T, T], \quad \omega_T(\epsilon_1, \epsilon_2) = \omega(\epsilon_1, \epsilon_2) \times [-T, T].$$

For δ_0 such that $0 < \delta_0 < T$, we set

$$Q_{\delta_0} = \Omega \times [-T + \delta_0, T - \delta_0], \quad Q_{\delta_0}(\epsilon) = (\Omega \setminus \omega(\epsilon)) \times [-T + \delta_0, T - \delta_0].$$

For formulating our Carleman estimate, we need some notation. Set

$$(2.3) \quad d = \left(\sup_{x \in \Omega} |x - x_0|^2 - \inf_{x \in \Omega} |x - x_0|^2 \right)^{\frac{1}{2}},$$

where $x_0 \notin \overline{\Omega}$ is arbitrarily fixed. We choose $\theta > 0$ such that

$$(2.4) \quad \theta + \frac{M_0 d}{\sqrt{\theta_1}} \sqrt{\theta} < \theta_0 \theta_1, \quad \theta_1 \inf_{x \in \Omega} |x - x_0|^2 - \theta \sup_{x \in \Omega} |x - x_0|^2 > 0.$$

Here we note that, since $x_0 \notin \overline{\Omega}$, such $\theta > 0$ exists.

We introduce two functions $\psi, \varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ of class C^1 by setting

$$(2.5) \quad \begin{aligned} \psi(x, t) &= |x - x_0|^2 - \theta |t|^2 \quad \text{for all } x \in \Omega, \quad -T \leq t \leq T, \\ \varphi(x, t) &= e^{\beta \psi(x, t)}, \quad \beta > 0, \end{aligned}$$

where $T > \frac{d}{\sqrt{\theta}}$. Therefore, by (2.4) and (2.5), we have

$$(2.6) \quad \varphi(x, 0) \geq d_0, \quad \varphi(x, \pm T) < d_0,$$

with $d_0 = \exp(\beta \inf_{x \in \Omega} |x - x_0|)$. Thus, for a given sufficiently small $\eta > 0$, we can choose a sufficiently small $\delta_0 = \delta_0(\eta)$ such that

$$(2.7) \quad \varphi(x, t) \leq d_0 - \eta \quad \text{for all } (x, t) \in Q \setminus \overline{Q_{2\delta_0}}.$$

We set $\nabla_{x,t} v(t, x) = \left(\frac{\partial v}{\partial x_1}, \frac{\partial v}{\partial x_2}, \frac{\partial v}{\partial x_3}, \frac{\partial v}{\partial t} \right) = (\nabla v, \partial_t v)$, and we shall use the weighted norm:

$$\|u\|_{H^{1,\tau}(Q)}^2 = \tau^2 \|u\|_{L^2(Q)}^2 + \|\nabla_{x,t} u\|_{L^2(Q)}^2.$$

Moreover $H_x^{k,\tau}(Q)$ is the Sobolev space equipped with the norm

$$\|u\|_{H_x^{k,\tau}(Q)}^2 = \sum_{|\alpha| \leq k} \tau^{2(k-|\alpha|)} \|\partial_x^\alpha u\|_{L^2(Q)}^2.$$

In what follows, $C > 0$ denote generic constants depending on $\widehat{\beta}$, τ_0 , M_0 , M_1 , θ_0 , θ_1 , Ω , T , x_0 , ω , χ and Φ , Ψ , ϵ , δ but independent of $\tau > \widehat{\tau}$.

2.1. Preliminary estimates. The first lemma is a classical Carleman estimate for a scalar hyperbolic equation (e.g., [14], [17], [25], [26]). See also Triggiani and Yao [47].

LEMMA 2.1. *Let φ be defined by (2.5). If $\frac{a}{\rho}$ satisfies (1.4), then there exists $\widehat{\beta} > 0$ such that for any $\beta > \widehat{\beta}$ we can choose $\widehat{\tau}(\beta) > 0$ such that the following estimate holds true:*

$$\tau \|e^{\tau \varphi} u\|_{H^{1,\tau}(Q)}^2 \leq C \|e^{\tau \varphi} (\rho \partial_t^2 - a \Delta) u\|_{L^2(Q)}^2,$$

whenever a function $u \in H^2(Q)$ is supported in Q and $\tau > \widehat{\tau}$.

The next lemma is a weighted estimate for the elliptic equation which follows from the formula

$$\begin{aligned} |\Delta(\mathbf{v}e^{\tau\varphi})| &= O(\tau^2)e^{\tau\varphi}|\mathbf{v}| + O(\tau)e^{\tau\varphi}|\nabla\mathbf{v}| \\ &\quad + e^{\tau\varphi}|\operatorname{rot}(\operatorname{rot}\mathbf{v}) - \nabla(\operatorname{div}\mathbf{v})| \end{aligned}$$

in Q and the standard a priori estimate for the Dirichlet problem for the Poisson equation.

LEMMA 2.2. *There exists a constant $C > 0$ such that*

$$\frac{1}{\tau} \|e^{\tau\varphi}\mathbf{v}\|_{H_x^{2,\tau}(Q)}^2 \leq C(\tau \|e^{\tau\varphi}\mathbf{v}\|_{H_x^{1,\tau}(Q)}^2 + \|e^{\tau\varphi}(\nabla\operatorname{div}\mathbf{v})\|_{L^2(Q)}^2 + \|e^{\tau\varphi}(\nabla\operatorname{rot}\mathbf{v})\|_{L^2(Q)}^2),$$

whenever $\mathbf{v} \in H^2(Q)$ and $\mathbf{v}|_{\Sigma} = 0$.

2.2. Carleman estimate for the Lamé system. In order to prove a Carleman estimate, we have to assume a condition called the pseudoconvexity (e.g., [14]) where the coefficient of the principal term is involved. Since such a coefficient is unknown in our inverse problem, we need to establish a Carleman estimate with one possible explicit characterization (1.4) of the coefficients for the pseudoconvexity, and we will argue similarly to Bellassoued [4]. Moreover for our stability estimates, unlike [12], [26], we require a Carleman estimate for functions which do not have compact supports.

Now we will consider the three-dimensional isotropic nonstationary Lamé system

$$(2.8) \quad P(x, \partial_x, \partial_t)\mathbf{y}(x, t) = \rho(x)\partial_t^2\mathbf{y}(x, t) - \mathcal{L}_{\lambda,\mu}(x, \partial_x)\mathbf{y}(x, t) = \mathbf{f}(x, t), \quad (x, t) \in Q.$$

We have a Carleman estimate.

LEMMA 2.3. *There exists $\hat{\beta} > 0$ such that for any $\beta > \hat{\beta}$ we can choose $\tau_0 = \tau_0(\beta) > 0$ such that for any solution $\mathbf{y} \in H^2(Q)$ to problem (2.8) the following estimate holds true:*

$$\begin{aligned} \frac{1}{\tau} \|e^{\tau\varphi}\mathbf{y}\|_{H_x^{2,\tau}(Q_{\delta_0}(3\epsilon_0))}^2 &\leq C(\|e^{\tau\varphi}\mathbf{f}\|_{L^2(Q)}^2 + \|e^{\tau\varphi}\nabla\mathbf{f}\|_{L^2(Q)}^2 + e^{C\tau} \|\mathbf{y}\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2 \\ &\quad + e^{2\tau(d_0-\eta)} \|\mathbf{y}\|_{H^2(Q)}^2) \end{aligned}$$

for any $\tau \geq \tau_0$, where the constant $C = C(\beta) > 0$ is independent of τ .

Proof. We introduce a cutoff function χ satisfying $0 \leq \chi \leq 1$, $\chi \in C_0^\infty(\mathbb{R}^3 \times \mathbb{R})$, $\chi_1 \in C_0^\infty(\mathbb{R}^3)$, $\chi_2 \in C_0^\infty(\mathbb{R})$, and

$$\chi(x, t) = \chi_1(x)\chi_2(t), \quad \chi_1(x) = \begin{cases} 0, & x \in \omega(\epsilon_0), \\ 1, & x \in \Omega \setminus \omega(3\epsilon_0), \end{cases} \quad \chi_2(t) = \begin{cases} 0, & |t| > T - \delta_0, \\ 1, & |t| < T - 2\delta_0. \end{cases}$$

Set $\mathbf{v}(x, t) = \chi(x, t)\mathbf{y}(x, t)$. Then we have

$$(2.9) \quad \partial_t^2\mathbf{v} - \frac{1}{\rho}\mathcal{L}_{\lambda,\mu}(x, \partial_x)\mathbf{v} = \frac{1}{\rho}\tilde{\mathbf{f}} \quad \text{in } Q, \quad \mathbf{v} = 0 \quad \text{in } \omega(\epsilon_0),$$

where

$$(2.10) \quad \tilde{\mathbf{f}}(x, t) = \chi(x, t)\mathbf{f} + [P, \chi_1]\chi_2\mathbf{y} + \chi_1[P, \chi_2]\mathbf{y}.$$

Let $v = \operatorname{div} \mathbf{v}$ and $\mathbf{w} = \operatorname{rot} \mathbf{v}$. Therefore, by [12], for example, we apply rot and div to (2.9) and obtain

$$\begin{aligned} \rho(x)\partial_t^2 \mathbf{v} - \mu(x)\Delta \mathbf{v} + A_1(\mathbf{v}, v) &= \tilde{\mathbf{f}}, \\ \rho(x)\partial_t^2 v - (2\mu(x) + \lambda(x))\Delta v + A_2(\mathbf{v}, v, \mathbf{w}) &= R_1 \tilde{\mathbf{f}}, \\ \rho(x)\partial_t^2 \mathbf{w} - \mu(x)\Delta \mathbf{w} + A_3(\mathbf{v}, v, \mathbf{w}) &= R_2 \tilde{\mathbf{f}}, \end{aligned}$$

where A_j and R_j are linear differential operators of the first order with bounded coefficients in Ω .

Since $\frac{\mu}{\rho}$ and $\frac{2\mu+\lambda}{\rho}$ satisfy (1.4), we have by Lemma 2.1

$$\begin{aligned} &\tau(\|e^{\tau\varphi} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} v\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \mathbf{w}\|_{H^{1,\tau}(Q)}^2) \\ &\leq C(\|\tilde{\mathbf{f}}e^{\tau\varphi}\|_{L^2(Q)}^2 + \|(\nabla \tilde{\mathbf{f}})e^{\tau\varphi}\|_{L^2(Q)}^2 + \|\mathbf{v}e^{\tau\varphi}\|_{L^2(Q)}^2 + \|(\nabla \mathbf{v})e^{\tau\varphi}\|_{L^2(Q)}^2 \\ &\quad + \|ve^{\tau\varphi}\|_{L^2(Q)}^2 + \|(\nabla v)e^{\tau\varphi}\|_{L^2(Q)}^2 + \|\mathbf{w}e^{\tau\varphi}\|_{L^2(Q)}^2 + \|(\nabla \mathbf{w})e^{\tau\varphi}\|_{L^2(Q)}^2). \end{aligned}$$

By taking $\tau > 0$ sufficiently large, we have

$$\begin{aligned} &\tau(\|e^{\tau\varphi} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} v\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \mathbf{w}\|_{H^{1,\tau}(Q)}^2) \\ (2.11) \quad &\leq C(\|\tilde{\mathbf{f}}e^{\tau\varphi}\|_{L^2(Q)}^2 + \|(\nabla \tilde{\mathbf{f}})e^{\tau\varphi}\|_{L^2(Q)}^2). \end{aligned}$$

On the other hand, since $[P, \chi_1]$ is a first-order differential operator which is supported in $\omega_T(\epsilon_0, 3\epsilon_0)$ and $[P, \chi_2]$ is a first-order differential operator which is supported in $Q \setminus Q_{2\delta_0}$, we obtain by (2.7) and (2.10)

$$\begin{aligned} \|e^{\tau\varphi} \tilde{\mathbf{f}}\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \nabla \tilde{\mathbf{f}}\|_{L^2(Q)}^2 &\leq C(\|e^{\tau\varphi} \mathbf{f}\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \nabla \mathbf{f}\|_{L^2(Q)}^2 + e^{C\tau} \|\mathbf{y}\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2 \\ &\quad + e^{2\tau(d_0-\eta)} \|\mathbf{y}\|_{H^2(Q)}^2). \end{aligned}$$

Thus, in terms of (2.11), we have

$$\begin{aligned} &\tau(\|e^{\tau\varphi} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \operatorname{div} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \operatorname{rot} \mathbf{v}\|_{H^{1,\tau}(Q)}^2) \\ &\leq C(\|e^{\tau\varphi} \mathbf{f}\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \nabla \mathbf{f}\|_{L^2(Q)}^2 + e^{C\tau} \|\mathbf{y}\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2 + e^{2\tau(d_0-\eta)} \|\mathbf{y}\|_{H^2(Q)}^2). \end{aligned}$$

By applying Lemma 2.2, we obtain

$$\begin{aligned} &\frac{1}{\tau} \|e^{\tau\varphi} \mathbf{v}\|_{H_x^{2,\tau}(Q)}^2 \\ &\leq C(\tau \|e^{\tau\varphi} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \operatorname{div} \mathbf{v}\|_{H^{1,\tau}(Q)}^2 + \|e^{\tau\varphi} \operatorname{rot} \mathbf{v}\|_{H^{1,\tau}(Q)}^2) \\ &\leq C(\|e^{\tau\varphi} \mathbf{f}\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \nabla \mathbf{f}\|_{L^2(Q)}^2 + e^{C\tau} \|\mathbf{y}\|_{H^3(\omega_T(\epsilon_0, 3\epsilon_0))}^2 + e^{2\tau(d_0-\eta)} \|\mathbf{y}\|_{H^2(Q)}^2). \end{aligned}$$

Since $\mathbf{v} = \chi \mathbf{y}$ and $\chi = 1$ in $Q_{\delta_0}(3\epsilon_0)$, we can replace \mathbf{v} by \mathbf{y} on the left-hand side, so that the proof of Lemma 2.3 is complete. \square

2.3. Carleman estimate for a first-order partial differential operator.

We consider a first-order partial differential equation

$$(2.12) \quad Av = \sum_{j=1}^n a_j(x) \partial_j v + a_0(x)v \equiv f(x), \quad x \in \Omega,$$

where

$$(2.13) \quad a_0 \in C(\overline{\Omega}), \quad a = (a_1, \dots, a_n) \in [C^1(\overline{\Omega})]^n,$$

and

$$(2.14) \quad |a(x) \cdot (x - x_0)| \geq c_0 > 0 \quad \text{on } \overline{\Omega},$$

with a constant $c_0 > 0$. We set

$$\varphi_0(x) = \varphi(x, 0), \quad x \in \Omega.$$

We have the following.

LEMMA 2.4. *In addition to (2.14), we assume that $\|a_0\|_{C^2(\overline{\Omega})} \leq M$ and $\|a_i\|_{C^2(\overline{\Omega})} \leq M$, $1 \leq i \leq 3$. Then there exists a constant $\widehat{\beta} > 0$ such that for all $\beta > \widehat{\beta}$ there exist $\widehat{\tau} = \widehat{\tau}(\beta) > 0$ and $C = C(\widehat{\tau}, \widehat{\beta}, \Omega, \omega) > 0$ such that*

$$\tau^2 \sum_{|\alpha| \leq 2} \int_{\Omega} |\partial_x^\alpha v|^2 e^{2\tau\varphi_0(x)} dx \leq C \sum_{|\alpha| \leq 2} \int_{\Omega} |\partial_x^\alpha f(x)|^2 e^{2\tau\varphi_0(x)} dx$$

for all $\tau > \widehat{\tau}$ and $v \in C_0^2(\Omega)$.

Proof. We will repeat the argument of Lemma 3.2 in [23], and, for completeness, we give the proof. We multiply both sides of (2.12) by $v(x)e^{2\tau\varphi_0(x)}$, and, by using the divergence theorem and $v \in C_0^2(\Omega)$, we obtain

$$\begin{aligned} & \int_{\Omega} Av(x) \cdot v(x) e^{2\tau\varphi_0(x)} dx = \int_{\Omega} \nabla v(x) \cdot \left(e^{2\tau\varphi_0(x)} v(x) a(x) \right) dx + \int_{\Omega} a_0(x) |v(x)|^2 dx \\ & = - \int_{\Omega} v(x) \operatorname{div} \left(e^{2\tau\varphi_0(x)} v(x) a(x) \right) dx + \int_{\Omega} a_0(x) |v(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & = - \int_{\Omega} |v|^2 e^{2\tau\varphi_0(x)} \operatorname{div}(a(x)) dx - 2\tau \int_{\Omega} |v(x)|^2 \nabla \varphi_0 \cdot a(x) e^{2\tau\varphi_0(x)} dx \\ & \quad - \int_{\Omega} e^{2\tau\varphi_0(x)} v(x) \nabla v(x) \cdot a(x) dx + \int_{\Omega} a_0(x) |v(x)|^2 e^{2\tau\varphi_0(x)} dx. \end{aligned}$$

By (2.14), we obtain

$$|\nabla \varphi_0(x) \cdot a(x)| \geq 2\beta c_0, \quad \nabla v(x) \cdot a(x) = Av - a_0(x)v(x), \quad x \in \Omega,$$

so that the Cauchy–Schwarz inequality yields

$$\begin{aligned} \tau \int_{\Omega} |v(x)|^2 e^{2\tau\varphi_0(x)} dx & \leq C \int_{\Omega} |(Av(x) \cdot v(x))| e^{2\tau\varphi_0(x)} dx + C \int_{\Omega} |v(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & \leq \frac{C\varepsilon}{\tau} \int_{\Omega} |Av(x)|^2 e^{2\tau\varphi_0(x)} dx + \varepsilon \tau \int_{\Omega} |v(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & \quad + C \int_{\Omega} |v(x)|^2 e^{2\tau\varphi_0(x)} dx. \end{aligned}$$

By choosing τ large and ε small, we obtain

$$(2.15) \quad \tau^2 \int_{\Omega} |v(x)|^2 e^{2\tau\varphi_0(x)} dx \leq C \int_{\Omega} |Av(x)|^2 e^{2\tau\varphi_0(x)} dx.$$

Since $A(\partial_j v) = \partial_j f(x) - \sum_{k=1}^n (\partial_j a_k) \partial_k v - (\partial_j a_0) v$ and $\partial_j v|_{\partial\Omega} = 0$, we apply (2.15) to $\partial_j v$, so that

$$\begin{aligned} \tau^2 \int_{\Omega} |\partial_j v(x)|^2 e^{2\tau\varphi_0} dx &\leq C \int_{\Omega} (|v(x)|^2 + |\nabla v(x)|^2) e^{2\tau\varphi_0} dx + C \int_{\Omega} |\partial_j f(x)|^2 e^{2\tau\varphi_0} dx \\ &\leq C \int_{\Omega} (|f(x)|^2 + |\partial_j f(x)|^2) e^{2\tau\varphi_0} dx + C \int_{\Omega} |\nabla v(x)|^2 e^{2\tau\varphi_0} dx. \end{aligned}$$

Therefore

$$\tau^2 \int_{\Omega} |\nabla v(x)|^2 e^{2\tau\varphi_0} dx \leq C \int_{\Omega} (|f(x)|^2 + |\nabla f(x)|^2) e^{2\tau\varphi_0} dx + C \int_{\Omega} |\nabla v(x)|^2 e^{2\tau\varphi_0} dx.$$

By taking $\tau_0 > 0$ sufficiently large, we have

$$\tau^2 \int_{\Omega} |\nabla v(x)|^2 e^{2\tau\varphi_0} dx \leq C \int_{\Omega} (|f(x)|^2 + |\nabla f(x)|^2) e^{2\tau\varphi_0} dx.$$

Next we have

$$A(\partial_k \partial_\ell v) = \partial_k \partial_\ell f - \sum_{j=1}^3 (\partial_k a_j) (\partial_\ell \partial_j v) + (\partial_\ell a_j) (\partial_k \partial_j v) + K(v, \nabla v),$$

where K is a linear operator of v and ∇v with bounded coefficients in Ω . Noting that $\partial_k \partial_\ell v = 0$ on $\partial\Omega$, we apply (2.15) to $\partial_k \partial_\ell v$, and we can complete the proof of Lemma 2.4. \square

2.4. Weak observation estimate. Let \mathbf{v} satisfy

$$\rho(x) \partial_t^2 \mathbf{v} - \mathcal{L}_{\lambda, \mu}(x, \partial_x) \mathbf{v} = R(x, t) \quad \text{in } Q \equiv \Omega \times (-T, T)$$

and

$$\mathbf{v}(x, t) = 0 \quad \text{on } \Sigma \equiv \Gamma \times (-T, T),$$

where we assume that

$$R(x, t) = 0, \quad (x, t) \in \omega \times (-T, T).$$

The following proposition shows the stability in the unique continuation of solutions of the Lamé system from lateral boundary data on an arbitrarily small part Γ_1 of $\partial\Omega$.

LEMMA 2.5. *For sufficiently large $T > 0$, there exists a constant $C > 0$ such that*

$$\|\mathbf{v}\|_{H^2(\omega_T(\varepsilon_0, 3\varepsilon_0))}^2 \leq C \left[\log \left(2 + \frac{C}{\sum_{|\alpha|=1}^2 \|\partial_x^\alpha \mathbf{v}\|_{L^2(\Sigma_1)}^2} \right) \right]^{-1}.$$

In our case, the corresponding uniqueness is already proved in Bellassoued [3] and Eller, Isakov, Nakamura, and Tataru [12]. For similar stability results for a scalar hyperbolic equation, see Bellassoued and Yamamoto [6] and Robbiano [44]. The proof of the lemma is given in section 4.

3. Proof of the main result. This section is devoted to the proof of Theorem 1.0. The key is the combination of Lemma 2.5 and the existing method (e.g., [4]).

3.1. Notation and preliminary estimates. For simplicity, we set

$$\mathbf{u} = \mathbf{u}(\lambda, \mu, \rho, \Phi, \Psi, \mathbf{g}), \quad \tilde{\mathbf{u}} = \mathbf{u}(\tilde{\lambda}, \tilde{\mu}, \tilde{\rho}, \Phi, \Psi, \mathbf{g}),$$

and

$$\mathbf{u}_* = \mathbf{u} - \tilde{\mathbf{u}}, \quad \rho_* = \rho - \tilde{\rho}, \quad \lambda_* = \lambda - \tilde{\lambda}, \quad \mu_* = \mu - \tilde{\mu}.$$

Then we obtain

$$(3.1) \quad \tilde{\rho} \partial_t^2 \mathbf{u}_*(x, t) = \mathcal{L}_{\tilde{\lambda}, \tilde{\mu}}(x, \partial_x) \mathbf{u}_*(x, t) + \mathbf{f}(x, t) \quad \text{in } Q,$$

$$(3.2) \quad \mathbf{u}_*(x, 0) = \partial_t \mathbf{u}_*(x, 0) = 0, \quad x \in \Omega,$$

and

$$(3.3) \quad \mathbf{u}_* = 0 \quad \text{on } \Sigma = \Gamma \times (-T, T).$$

Here we set

$$\begin{aligned} \mathbf{f}(x, t) &= -\rho_*(x) \partial_t^2 \mathbf{u}(x, t) + (\lambda_*(x) + \mu_*(x)) \nabla(\operatorname{div} \mathbf{u})(x, t) + \mu_*(x) \Delta \mathbf{u}(x, t) \\ &\quad + (\operatorname{div} \mathbf{u})(x, t) \nabla \lambda_*(x) + (\nabla \mathbf{u}(x, t) + (\nabla \mathbf{u}(x, t))^T) \nabla \mu_*(x) \\ &= -(\rho_*(x) \partial_t^2 \mathbf{u} - \mathcal{L}_{\lambda_*, \mu_*}(x, \partial_x) \mathbf{u}). \end{aligned}$$

Moreover we set

$$\mathbf{y}_1(x, t) = \partial_t^2 \mathbf{u}_*(x, t), \quad \mathbf{y}_2(x, t) = \partial_t^3 \mathbf{u}_*(x, t), \quad \mathbf{y}_3(x, t) = \partial_t^4 \mathbf{u}_*(x, t).$$

Then we have

$$(3.4) \quad \tilde{\rho} \partial_t^2 \mathbf{y}_j - \mathcal{L}_{\tilde{\lambda}, \tilde{\mu}}(x, \partial_x) \mathbf{y}_j = \partial_t^{j+2} \mathbf{f}.$$

For simplicity, it is convenient to use the following notation:

$$\mathcal{D} = \sum_{j=2}^4 \|\partial_t^j \mathbf{u}_*\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2 = \sum_{j=1}^3 \|\mathbf{y}_j\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2$$

and

$$\mathcal{E} = \sum_{|\alpha| \leq 1} \|e^{\tau\varphi} \partial_x^\alpha \rho_*\|_{L^2(Q)}^2 + \sum_{|\alpha| \leq 2} (\|e^{\tau\varphi} \partial_x^\alpha \lambda_*\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \partial_x^\alpha \mu_*\|_{L^2(Q)}^2),$$

$$\mathbf{z}_1(x) = \mathbf{y}_1(x, 0), \quad \mathbf{z}_2(x) = \mathbf{y}_2(x, 0).$$

We have the following.

LEMMA 3.1. *There exists a constant $\hat{\beta} > 0$ such that for all $\beta \geq \hat{\beta}$ there exist $\hat{\tau}$ and $C > 0$ such that*

$$\sum_{|\alpha| \leq 2} \int_{\Omega \setminus \omega(3\epsilon_0)} \tau^{4-2|\alpha|} \left(|\partial_x^\alpha \mathbf{z}_1(x)|^2 + |\partial_x^\alpha \mathbf{z}_2(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \leq C(\tau^4 e^{2\tau(d_0 - \eta)} + \tau^2 \mathcal{E} + e^{C\tau} \mathcal{D})$$

for all large $\tau > \hat{\tau}$.

Here we recall that $\varphi_0(x) = \varphi(x, 0)$, $x \in \Omega$.

Proof. Noting that $\mathbf{u}_* \in W^{8,\infty}(Q)$, we can apply Lemma 2.3 to (3.4), so that we have

$$(3.5) \quad \frac{1}{\tau} \|e^{\tau\varphi} \mathbf{y}_j\|_{H_x^{2,\tau}(Q_{\delta_0}(3\epsilon_0))}^2 \leq C(\|e^{\tau\varphi} \partial_t^{j+2} \mathbf{f}\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \nabla(\partial_t^{j+2} \mathbf{f})\|_{L^2(Q)}^2) \\ + e^{C\tau} \|\mathbf{y}_j\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^2 + e^{2\tau(d_0-\eta)} \|\mathbf{y}_j\|_{H^2(Q)}^2, \quad j = 1, 2, 3.$$

By the definition of \mathbf{f} , we have

$$\left| \nabla(\partial_t^{j+2} \mathbf{f}(x, t)) \right|^2 \leq C \left(\sum_{|\alpha| \leq 1} |\partial_x^\alpha \rho_*(x)|^2 + \sum_{|\alpha| \leq 2} (|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2) \right)$$

and

$$\left| \partial_t^{j+2} \mathbf{f}(x, t) \right|^2 \leq C \left(|\rho_*(x)|^2 + \sum_{|\alpha| \leq 1} (|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2) \right) \quad \text{in } Q.$$

Therefore

$$\sum_{j=1}^5 \left(\|e^{\tau\varphi} \nabla(\partial_t^j \mathbf{f})\|_{L^2(Q)}^2 + \|e^{\tau\varphi} \partial_t^j \mathbf{f}\|_{L^2(Q)}^2 \right) \leq C\mathcal{E}.$$

Thus (3.5) implies that

$$(3.6) \quad \frac{1}{\tau} \sum_{j=1}^3 \|e^{\tau\varphi} \mathbf{y}_j\|_{H_x^{2,\tau}(Q_{\delta_0}(3\epsilon_0))}^2 \leq C(e^{2\tau(d_0-\eta)} + \mathcal{E} + e^{C\tau}\mathcal{D}).$$

On the other hand, we introduce a cutoff function $\chi_2 \in C_0^\infty(\mathbb{R})$ satisfying $0 \leq \chi_2 \leq 1$ and

$$\chi_2(t) = \begin{cases} 0, & |t| > T - \delta_0, \\ 1, & |t| < T - 2\delta_0. \end{cases}$$

Then we have

$$\int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha \mathbf{z}_j(x)|^2 e^{2\tau\varphi_0(x)} dx = \int_{-T}^0 \frac{\partial}{\partial t} \left(\int_{\Omega \setminus \omega(3\epsilon_0)} |(\partial_x^\alpha \mathbf{y}_j)(x, t)|^2 \chi_2(t)^2 e^{2\tau\varphi} dx \right) dt \\ = 2 \int_{-T}^0 \int_{\Omega \setminus \omega(3\epsilon_0)} (\partial_x^\alpha \mathbf{y}_{j+1} \cdot \partial_x^\alpha \mathbf{y}_j) \chi_2^2(t) e^{2\tau\varphi} dx dt \\ + 2\tau \int_{-T}^0 \int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha \mathbf{y}_j|^2 \chi_2^2(\partial_t \varphi) e^{2\tau\varphi} dx dt \\ + 2 \int_{-T}^0 \int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha \mathbf{y}_j|^2 \chi_2'(t) \chi_2(t) e^{2\tau\varphi} dx dt, \quad j = 1, 2.$$

Since $\chi_2'(t)$ is supported in $[-T, -T + 2\delta_0] \cup [T - 2\delta_0, T]$, by the Cauchy-Schwarz

inequality we obtain

$$(3.7) \quad \begin{aligned} & \sum_{j=1}^2 \sum_{|\alpha| \leq 2} \int_{\Omega \setminus \omega(3\epsilon_0)} \tau^{4-2|\alpha|} |\partial_x^\alpha \mathbf{z}_j(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & \leq C \left(\tau \sum_{j=1}^3 \|e^{\tau\varphi} \mathbf{y}_j\|_{H_x^{2,\tau}(Q_{\delta_0}(3\epsilon_0))}^2 + \tau^4 e^{2\tau(d_0-\eta)} \right). \end{aligned}$$

Therefore (3.6) completes the proof of the lemma. \square

3.2. Estimation for the two Lamé coefficients. We will consider a first-order partial differential equations in λ_* , μ_* , and ρ_* . That is, by (3.1), (3.2), and $\mathbf{u}, \mathbf{v} \in W^{8,\infty}(Q)$, we have

$$(3.8) \quad \tilde{\rho} \partial_t^2 \mathbf{u}_*(x, 0) = \mathbf{f}(x, 0), \quad \tilde{\rho} \partial_t^3 \mathbf{u}_*(x, 0) = \partial_t \mathbf{f}(x, 0).$$

For simplicity, for $x \in \bar{\Omega}$, we set

$$\mathbf{a} = \begin{pmatrix} -\frac{1}{\rho} \mathcal{L}_{\lambda,\mu}(x, \partial_x) \Phi \\ -\frac{1}{\rho} \mathcal{L}_{\lambda,\mu}(x, \partial_x) \Psi \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} \operatorname{div} \Phi \\ 0 \\ 0 \\ \operatorname{div} \Psi \\ 0 \\ 0 \end{pmatrix},$$

$$\mathbf{b}_2 = \begin{pmatrix} 0 \\ \operatorname{div} \Phi \\ 0 \\ 0 \\ \operatorname{div} \Psi \\ 0 \end{pmatrix}, \quad \mathbf{b}_3 = \begin{pmatrix} 0 \\ 0 \\ \operatorname{div} \Phi \\ 0 \\ 0 \\ \operatorname{div} \Psi \end{pmatrix},$$

$$(\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3) = \begin{pmatrix} \nabla \Phi + (\nabla \Phi)^T \\ \nabla \Psi + (\nabla \Psi)^T \end{pmatrix}, \quad \mathbf{G} = \begin{pmatrix} \tilde{\rho} \mathbf{z}_1(x) - (\lambda_* + \mu_*) \nabla(\operatorname{div} \Phi) - \mu_* \Delta \Phi \\ \tilde{\rho} \mathbf{z}_2(x) - (\lambda_* + \mu_*) \nabla(\operatorname{div} \Psi) - \mu_* \Delta \Psi \end{pmatrix}.$$

Then we can rewrite (3.8) as

$$(3.9) \quad \mathbf{a} \rho_* + \sum_{k=1}^3 \mathbf{b}_k \partial_k \lambda_*(x) = \mathbf{G} - \sum_{k=1}^3 \mathbf{d}_k \partial_k \mu_*(x).$$

In view of (3.9), we show the following.

LEMMA 3.2. *Let assumptions (\mathcal{H}_1) – (\mathcal{H}_2) hold true. Then there exists a constant $C > 0$ such that the following estimate holds:*

$$\sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \mu_*(x)|^2 + |\partial_x^\alpha \lambda_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \leq C(\tau^2 e^{2\tau(d_0-\eta)} + \mathcal{E} + e^{C\tau\mathcal{D}})$$

provided that τ is large.

Proof. By (3.9), for $j_1 \in \{1, 2, 3, 4, 5, 6\}$, we have

$$\{\mathbf{a}\}_{j_1 \rho_*} + \sum_{k=1}^3 \{\mathbf{b}_k\}_{j_1} \partial_k \lambda_*(x) = \{\mathbf{G}\}_{j_1} - \sum_{k=1}^3 \{\mathbf{d}_k\}_{j_1} \partial_k \mu_*(x) \quad \text{on } \bar{\Omega}.$$

For any fixed $x \in \bar{\Omega}$, we regard this as a system of five linear equations with respect to four unknowns ρ_* , $\partial_1 \lambda_*$, $\partial_2 \lambda_*$, $\partial_3 \lambda_*$, and so, for the existence of solutions, we need the consistency of the coefficients, that is,

$$\det_{j_1} \left(\mathbf{a}, \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \mathbf{G} - \sum_{k=1}^3 \mathbf{d}_k \partial_k \mu_* \right) = 0 \quad \text{on } \bar{\Omega}.$$

Hence

$$(3.10) \quad \sum_{k=1}^3 \det_{j_1} (\mathbf{a}, \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \mathbf{d}_k) \partial_k \mu_* = \det_{j_1} (\mathbf{a}, \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \mathbf{G}) := f \quad \text{on } \bar{\Omega}$$

by the linearity of the determinant with respect to the column vectors. Here we note that

$$(3.11) \quad \sum_{|\alpha| \leq 2} |\partial_x^\alpha f(x)| \leq C \left(\sum_{|\alpha| \leq 2} (|\partial_x^\alpha \mathbf{z}_1(x)| + |\partial_x^\alpha \mathbf{z}_2(x)|) + \sum_{|\alpha| \leq 2} (|\partial_x^\alpha \lambda_*(x)| + |\partial_x^\alpha \mu_*(x)|) \right).$$

By (3.10) the function $\mu_* \in C_0^2(\Omega)$ solves the following first-order partial differential equation:

$$\sum_{k=1}^3 a_k(x) \partial_k \mu_*(x) = f(x), \quad x \in \Omega,$$

where

$$a_k(x) = \det_{j_1} (\mathbf{a}, \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \mathbf{d}_k).$$

In view of (\mathcal{H}_1) , we can apply Lemma 2.4 to $v = \mu_*$, and we obtain

$$\tau^2 \sum_{|\alpha| \leq 2} \int_{\Omega} |\partial_x^\alpha \mu_*|^2 e^{2\tau\varphi_0(x)} dx \leq C \sum_{|\alpha| \leq 2} \int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha f(x)|^2 e^{2\tau\varphi_0(x)} dx,$$

where we have used $\mu_* \equiv 0$ in $\omega \supset \omega(3\epsilon_0)$. Hence (3.11) and Lemma 3.1 yield

$$\begin{aligned} & \tau^2 \sum_{|\alpha| \leq 2} \int_{\Omega} |\partial_x^\alpha \mu_*(x)|^2 e^{2\tau\varphi_0(x)} dx \leq C \sum_{j=1}^2 \sum_{|\alpha| \leq 2} \int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha \mathbf{z}_j(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & + C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \\ & \leq C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx + C(e^{C\tau\mathcal{D}} + \tau^4 e^{2\tau(d_0-\eta)} + \tau^2 \mathcal{E}) \end{aligned}$$

for all large $\tau > 0$. Similarly, by assumption (\mathcal{H}_2) we can argue for λ_* and obtain

$$\begin{aligned} & \tau^2 \sum_{|\alpha| \leq 2} \int_{\Omega} |\partial_x^\alpha \lambda_*(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & \leq C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx + C(e^{C\tau\mathcal{D}} + \tau^4 e^{2\tau(d_0-\eta)} + \tau^2 \mathcal{E}). \end{aligned}$$

Hence by adding the two inequalities, we have

$$\begin{aligned} & \tau^2 \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \mu_*(x)|^2 + |\partial_x^\alpha \lambda_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \\ & \leq C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \mu_*(x)|^2 + |\partial_x^\alpha \lambda_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \\ & \quad + C(\tau^4 e^{2\tau(d_0-\eta)} + \tau^2 \mathcal{E} + e^{C\tau} \mathcal{D}) \end{aligned}$$

for all large $\tau > 0$. By taking $\tau > 0$ large, we can absorb the first term on the right-hand side into the left-hand side, and the proof is complete. \square

3.3. Estimation for the density.

LEMMA 3.3. *Let assumptions (\mathcal{H}_1) – (\mathcal{H}_2) hold true. Then there exists a constant $C > 0$ such that*

$$\sum_{|\alpha| \leq 1} \int_{\Omega} |\partial_x^\alpha \rho_*(x)|^2 e^{2\tau\varphi_0(x)} dx \leq C(\tau^2 e^{2\tau(d_0-\eta)} + \mathcal{E} + e^{C\tau} \mathcal{D}),$$

provided that τ is large.

Proof. By (3.9), we have

$$\mathbf{a}\rho_*(x) = - \sum_{k=1}^3 \mathbf{b}_k \partial_k \lambda_*(x) + \mathbf{G} - \sum_{k=1}^3 \mathbf{d}_k \partial_k \mu_*(x), \quad x \in \Omega.$$

Moreover, by (\mathcal{H}_1) – (\mathcal{H}_2) , we see that $|\mathbf{a}(x)| > 0$ for $x \in \bar{\Omega}$, so that

$$|\rho_*(x)| \leq C |\mathbf{G}(x)| + C \sum_{|\alpha| \leq 1} (|\partial_x^\alpha \lambda_*(x)| + |\partial_x^\alpha \mu_*(x)|).$$

Similarly we have

$$|\partial_j \rho_*(x)| \leq C \left(|\mathbf{G}(x)| + |\nabla \mathbf{G}(x)| + \sum_{|\alpha| \leq 2} (|\partial_x^\alpha \lambda_*(x)| + |\partial_x^\alpha \mu_*(x)|) \right).$$

Hence

$$\begin{aligned} \sum_{|\alpha| \leq 1} \int_{\Omega} |\partial_x^\alpha \rho_*(x)|^2 e^{2\tau\varphi_0(x)} dx & \leq C \sum_{|\alpha| \leq 1} \int_{\Omega \setminus \omega(3\epsilon_0)} |\partial_x^\alpha \mathbf{G}(x)|^2 e^{2\tau\varphi_0(x)} dx \\ & \quad + C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx. \end{aligned}$$

Furthermore we see that

$$\sum_{|\alpha| \leq 1} |\partial_x^\alpha \mathbf{G}(x)|^2 \leq C \sum_{|\alpha| \leq 1} \left(|\partial_x^\alpha \mathbf{z}_1(x)|^2 + |\partial_x^\alpha \mathbf{z}_2(x)|^2 \right) + C \sum_{|\alpha| \leq 1} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right).$$

Therefore

$$\begin{aligned} \sum_{|\alpha| \leq 1} \int_{\Omega} |\partial_x^\alpha \rho_*(x)|^2 e^{2\tau\varphi_0(x)} dx & \leq C \sum_{|\alpha| \leq 1} \int_{\Omega \setminus \omega(3\epsilon_0)} \left(|\partial_x^\alpha \mathbf{z}_1(x)|^2 + |\partial_x^\alpha \mathbf{z}_2(x)|^2 \right) e^{2\tau\varphi_0(x)} dx \\ & \quad + C \sum_{|\alpha| \leq 2} \int_{\Omega} \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi_0(x)} dx. \end{aligned}$$

By Lemmas 3.1 and 3.2, we obtain the conclusion of the lemma. \square

3.4. Completion of the proof of the main result. In terms of Lemmas 3.2 and 3.3, we will now complete the proof of Theorem 1.0.

Since $\varphi(x, 0) > \varphi(x, t)$ for $t \neq 0$, by the Lebesgue theorem, we have

$$\begin{aligned} \sum_{|\alpha| \leq 1} \int_Q |\partial_x^\alpha \rho_*(x)|^2 e^{2\tau\varphi} dx dt &= \sum_{|\alpha| \leq 1} \int_\Omega |\partial_x^\alpha \rho_*(x)|^2 e^{2\tau\varphi(x,0)} \left(\int_{-T}^T e^{2\tau(\varphi(x,t) - \varphi(x,0))} dt \right) dx \\ &= o(1) \sum_{|\alpha| \leq 1} \int_\Omega |\partial_x^\alpha \rho_*(x)| e^{2\tau\varphi(x,0)} dx \end{aligned}$$

as $\tau \rightarrow \infty$. Similarly we obtain

$$\begin{aligned} &\sum_{|\alpha| \leq 2} \int_Q \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi} dx dt \\ &= o(1) \sum_{|\alpha| \leq 2} \int_\Omega \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) e^{2\tau\varphi(x,0)} dx, \end{aligned}$$

as $\tau \rightarrow \infty$. We set

$$\mathcal{E}_0 = \sum_{|\alpha| \leq 1} \|e^{\tau\varphi_0} \partial_x^\alpha \rho_*\|_{L^2(\Omega)}^2 + \sum_{|\alpha| \leq 2} (\|e^{\tau\varphi_0} \partial_x^\alpha \lambda_*\|_{L^2(\Omega)}^2 + \|e^{\tau\varphi_0} \partial_x^\alpha \mu_*\|_{L^2(\Omega)}^2).$$

Therefore, by Lemmas 3.2 and 3.3, we obtain

$$\mathcal{E}_0 \leq C(\tau^2 e^{2\tau(d_0 - \eta)} + e^{C\tau\mathcal{D}}) + o(1)\mathcal{E}_0$$

for all large $\tau > 0$.

By (2.6) and $\sup_{\tau > 0} (\tau^2 e^{-\tau\eta}) < \infty$, we have

$$\begin{aligned} \sum_{|\alpha| \leq 1} \int_\Omega |\partial_x^\alpha \rho_*(x)|^2 dx + \sum_{|\alpha| \leq 2} \int_\Omega \left(|\partial_x^\alpha \lambda_*(x)|^2 + |\partial_x^\alpha \mu_*(x)|^2 \right) dx &\leq C e^{-\tau\eta} + C e^{C\tau\mathcal{D}} \\ (3.12) \end{aligned}$$

for all large $\tau > \tau_0$. We replace C by $C e^{C\tau_0}$, and (3.12) holds for any $\tau > 0$. We may assume that $\mathcal{D} < 1$.

Now we choose $\tau > 0$ such that

$$e^{C\tau\mathcal{D}} = e^{-\tau\eta},$$

that is,

$$\tau = -\frac{1}{\eta + C} \log \mathcal{D}.$$

Therefore (3.12) implies

$$\|\rho_*\|_{H^1(\Omega)}^2 + \|\lambda_*\|_{H^2(\Omega)}^2 + \|\mu_*\|_{H^2(\Omega)}^2 \leq 2C\mathcal{D}^{\frac{\eta}{\eta+C}}$$

and

$$\|\rho_*\|_{H^1(\Omega)}^2 + \|\lambda_*\|_{H^2(\Omega)}^2 + \|\mu_*\|_{H^2(\Omega)}^2 \leq C \|\mathbf{u}_*\|_{H^6(\omega_T(\epsilon_0, 3\epsilon_0))}^{2\sigma},$$

with $\sigma = \frac{\eta}{\eta+C} \in (0, 1)$. Then, noting that $\|\mathbf{u}_*\|_{W^{s,\infty}(Q)} \leq M_1$ and the interpolation inequality

$$\|\mathbf{u}_*\|_{H^6(\omega_T(\epsilon_0, 3\epsilon_0))} \leq C \|\mathbf{u}_*\|_{H^2(\omega_T(\epsilon_0, 3\epsilon_0))}^{\frac{1}{3}} \|\mathbf{u}_*\|_{H^8(\omega_T(\epsilon_0, 3\epsilon_0))}^{\frac{2}{3}},$$

we apply Lemma 2.5 to $\mathbf{v} = \mathbf{u}_*$, so that

$$\|\rho_*\|_{H^1(\Omega)}^2 + \|\lambda_*\|_{H^2(\Omega)}^2 + \|\mu_*\|_{H^2(\Omega)}^2 \leq C \left[\log \left(2 + \frac{C}{\epsilon(\Sigma_1)} \right) \right]^{-2\kappa},$$

where $\kappa \in (0, 1)$. Thus the proof of Theorem 1.0 is complete.

4. Proof of Lemma 2.5. We will now prove Lemma 2.5. This will be done in terms of the FBI transformation. For $(\lambda, \mu, \rho) \in \Lambda$, let us recall that \mathbf{v} is a given solution to

$$(4.1) \quad \rho(x)\partial_t^2 \mathbf{v} - \mathcal{L}_{\lambda,\mu}(x, \partial_x)\mathbf{v} = R(x, t) \quad \text{in } Q \equiv \Omega \times (-T, T),$$

with the Dirichlet boundary condition

$$(4.2) \quad \mathbf{v}(x, t) = 0 \quad \text{on } \Sigma \equiv \Gamma \times (-T, T).$$

Here and henceforth we assume that

$$(4.3) \quad R(x, t) = 0, \quad (x, t) \in \omega \times (-T, T).$$

4.1. Preliminary and elliptic estimation. Since we can choose $T > 0$ sufficiently large, we may assume that $T > 2$. Denote for $r > 0$

$$\begin{aligned} \Omega_r &= \Omega \times (-r, r); & \omega_r(\epsilon_0, 3\epsilon_0) &= \omega(\epsilon_0, 3\epsilon_0) \times (-r, r); \\ \Sigma_r &= \Gamma \times (-r, r), & \Sigma_{1,r} &= \Gamma_1 \times (-r, r). \end{aligned}$$

We introduce $\theta \in C_0^\infty(\mathbb{R})$ to be a cutoff function defined by

$$\theta(t) = \begin{cases} 1, & |t| \leq (T-2), \\ 0, & |t| \geq (T-1). \end{cases}$$

Let $\gamma > 0$. We introduce the partial FBI transformation \mathcal{F}_γ . It is defined for $\mathbf{u} \in \{\mathcal{S}(\mathbb{R}^4)\}^3$, the space of rapidly decreasing functions, by

$$\mathbf{u}_{\gamma,t}(x, s) = \mathcal{F}_\gamma \mathbf{u}(x, z) = \sqrt{\frac{\gamma}{2\pi}} \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(z-y)^2} \theta(y) \mathbf{u}(x, y) dy, \quad z = t + is.$$

Then

$$|D_x^\alpha \mathcal{F}_\gamma \mathbf{u}(x, z)| \leq C \sqrt{\frac{\gamma}{2\pi}} e^{\gamma s^2} e^{-\frac{\gamma}{2}(\text{dist}(t, \text{supp}(\theta \mathbf{u}))^2)} \sup_{x \in \mathbb{R}^3} \|D_x^\alpha \mathbf{u}(x, \cdot)\|_{L^2(\mathbb{R})}^2$$

for any $\mathbf{u} \in C_0^\infty(\mathbb{R}^3 \times \mathbb{R})$ (see [45]).

Henceforth C_j, C denote generic constants which are independent of λ, γ, r, τ . Next we assume that T is sufficiently large, $s \in [-3r, 3r]$, and $t \in [-\frac{T}{2}, \frac{T}{2}]$. In particular we assume that $\frac{T}{2} > r$. We introduce a cutoff function χ_3 satisfying $0 \leq \chi_3 \leq 1$, $\chi_3 \in C_0^\infty(\mathbb{R}^3)$, and

$$\chi_3(x) = \begin{cases} 1, & \text{if } x \in \omega(6\epsilon_0), \\ 0, & \text{if } x \in \Omega \setminus \omega(7\epsilon_0). \end{cases}$$

Let $\mathbf{v}(x, t)$ satisfy (4.1). By setting $\mathbf{u}(x, t) = \chi_3(x)\mathbf{v}(x, t)$ and noting that $R(x, t)$ is zero in $\omega(7\epsilon_0)$, we obtain

$$(4.4) \quad \rho(x)\partial_t^2 \mathbf{u} - \mathcal{L}_{\lambda,\mu}(x, \partial_x)\mathbf{u} = -[\mathcal{L}_{\lambda,\mu}(x, \partial_x), \chi_3]\mathbf{v} \quad \text{in } Q$$

and

$$(4.5) \quad \mathbf{u}(x, t) = 0 \quad \text{on } \Sigma.$$

In connection with the operator $\rho(x)\partial_t^2 - \mathcal{L}_{\lambda,\mu}(x, \partial_x)$, we define an elliptic operator by

$$Q_{\rho,\lambda,\mu} = \rho(x)\partial_s^2 + \mathcal{L}_{\lambda,\mu}(x, \partial_x).$$

Noting that

$$\partial_s \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(is+t-y)^2} \theta(y)\mathbf{u}(x, y)dy = i \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(z-y)^2} \partial_y [\theta(y)\mathbf{u}(x, y)] dy,$$

by integration by parts, we have

$$(4.6) \quad \begin{aligned} Q_{\rho,\lambda,\mu}\mathbf{u}_{\gamma,t}(x, s) &= F_{\gamma,t}(x, s) + G_{\gamma,t}(x, s) := \mathbf{f}(x, s), \quad (x, s) \in \Omega_{3r}, \\ \mathbf{u}_{\gamma,t}(x, s) &= 0, \quad (x, s) \in \Sigma_{3r}, \end{aligned}$$

where

$$F_{\gamma,t}(x, s) = -\sqrt{\frac{\gamma}{2\pi}} \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(z-y)^2} (2\theta'(y)\partial_y \mathbf{u}(x, y) + \theta''(y)\mathbf{u}(x, y)) dy$$

and

$$G_{\gamma,t}(x, s) = \sqrt{\frac{\gamma}{2\pi}} \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(z-t)^2} \theta(y) [\mathcal{L}_{\lambda,\mu}(x, \partial_x), \chi] \mathbf{v}(x, y) dy.$$

Since $\theta' = \frac{d\theta}{dy}$ and $\theta'' = \frac{d^2\theta}{dy^2}$ are supported in $T - 2 \leq |y| \leq T - 1$, there exists $\eta > 0$, independent of T , such that

$$(4.7) \quad \|F_{\gamma,t}\|_{H^1(\Omega_{3r})} \leq Ce^{-\eta\gamma T} \|\mathbf{u}\|_{H^2(Q)} \quad \forall t \in \left[-\frac{T}{2}, \frac{T}{2}\right].$$

Moreover there exists $C_1 > 0$, independent of T , such that

$$(4.8) \quad \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})} \leq Ce^{C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} \quad \forall t \in \left[-\frac{T}{2}, \frac{T}{2}\right].$$

By the definition, we easily obtain

$$(4.9) \quad G_{\gamma,t}(x, s) = 0 \quad \forall x \in \omega(6\epsilon_0).$$

Let K be a compact set in $\bar{\Omega} \times (-3r, 3r)$ and let $\psi(x, s)$ be a C^∞ function satisfying $\nabla_{x,s}\psi(x, s) \neq 0$ on K . Let

$$\varphi(x, s) = e^{-\beta\psi(x,s)},$$

where $\beta > 0$ is sufficiently large.

Henceforth we set

$$\|u\|_{H^j_\tau(\Omega_{3r})}^2 = \sum_{|\alpha| \leq j} \tau^{2j-2|\alpha|} \|\partial^\alpha u\|_{L^2(\Omega_{3r})}^2$$

and

$$\|u\|_{H^j_\tau(\Sigma_{3r})}^2 = \sum_{|\alpha| \leq j} \tau^{2j-2|\alpha|} \|\partial^\alpha u\|_{L^2(\Sigma_{3r})}^2.$$

We note that the right-hand sides include not only tangential derivatives but also normal derivatives on Σ_{3r} .

Consider a scalar second-order elliptic operator

$$P(x, D) = \rho(x)\partial_s^2 + \sum_{j,k=1}^n a_{jk}(x)\partial_j\partial_k + \sum_{j=1}^n b_j(x)\partial_j + c.$$

If we assume that P has C^2 coefficients, then the following Carleman estimate holds true:

$$(4.10) \quad \begin{aligned} \frac{C}{\tau} \|e^{\tau\varphi} u\|_{H^2_\tau(\Omega_{3r})}^2 &\leq \|e^{\tau\varphi} Pu\|_{L^2(\Omega_{3r})}^2 + \|e^{\tau\varphi} u\|_{H^2_\tau(\Sigma_{3r})}^2, \\ C\tau \|e^{\tau\varphi} u\|_{H^1_\tau(\Omega_{3r})}^2 &\leq \|e^{\tau\varphi} Pu\|_{L^2(\Omega_{3r})}^2 + \tau \|e^{\tau\varphi} u\|_{H^1_\tau(\Sigma_{3r})}^2, \end{aligned}$$

whenever $u \in C_0^\infty(K)$ and $\tau > \tau_0$ (see, for example, [37] and [19], respectively, for the first and second Carleman estimates).

By (4.10) we can derive the following Carleman estimate for the elliptic system (4.6).

LEMMA 4.1. *There exist $\tau_0 > 0$ and $C > 0$ such that*

$$\frac{C}{\tau} \|e^{\tau\varphi} \mathbf{u}\|_{H^2_\tau(\Omega_{3r})}^2 \leq \|e^{\tau\varphi} Q_{\rho,\lambda,\mu} \mathbf{u}\|_{H^1(\Omega_{3r})}^2 + \|e^{\tau\varphi} \mathbf{u}\|_{H^2_\tau(\Sigma_{3r})}^2,$$

whenever $\mathbf{u} \in C_0^\infty(K)$ and all $\tau > \tau_0$.

Proof. In order to prove Lemma 4.1 we will extend system (4.6) for three unknown functions u_1, u_2, u_3 to a new one for four unknown functions by introducing $v = \text{div} \mathbf{u}$. We refer to Eller, Isakov, Nakamura, and Tataru [12] and Ikehata, Nakamura, and Yamamoto [15]. If \mathbf{u} solves

$$Q_{\rho,\lambda,\mu} \mathbf{u} = \rho(x)\partial_s^2 \mathbf{u} + \mathcal{L}_{\lambda,\mu}(x, \partial_x) \mathbf{u} = \mathbf{f},$$

then

$$\rho(x)\partial_s^2 v + (\lambda(x) + 2\mu(x))\Delta v + A_{1,1}(v, \mathbf{u}) = \text{div} \mathbf{f}$$

and

$$\rho(x)\partial_s^2 \mathbf{u} + \mu(x)\Delta \mathbf{u} + A_{1,2}(v, \mathbf{u}) = \mathbf{f},$$

where $A_{1,1}, A_{1,2}$ are (matrix) linear partial differential operators with measurable and bounded coefficients. By using the scalar Carleman estimate (4.10), we obtain

$$\frac{C}{\tau} \|e^{\tau\varphi} \mathbf{u}\|_{H^2_\tau(\Omega_{3r})}^2 \leq \|e^{\tau\varphi} \mathbf{f}\|_{L^2(\Omega_{3r})}^2 + \|e^{\tau\varphi} \mathbf{u}\|_{H^2_\tau(\Sigma_{3r})}^2 + \|e^{\tau\varphi} v\|_{H^1_\tau(\Omega_{3r})}^2$$

and

$$C\tau \|e^{\tau\varphi}v\|_{H^1_\tau(\Omega_{3r})}^2 \leq \|e^{\tau\varphi}\operatorname{div}\mathbf{f}\|_{L^2(\Omega_{3r})}^2 + \tau \|e^{\tau\varphi}v\|_{H^1_\tau(\Sigma_{3r})}^2 + \|e^{\tau\varphi}\mathbf{u}\|_{H^1_\tau(\Omega_{3r})}^2.$$

Therefore, by choosing τ large, we obtain

$$\begin{aligned} \frac{C}{\tau} \|e^{\tau\varphi}\mathbf{u}\|_{H^2_\tau(\Omega_{3r})}^2 &\leq \|e^{\tau\varphi}\mathbf{f}\|_{L^2(\Omega_{3r})}^2 + \frac{1}{\tau} \|e^{\tau\varphi}\nabla\mathbf{f}\|_{L^2(\Omega_{3r})}^2 \\ &\quad + \|e^{\tau\varphi}\mathbf{u}\|_{H^2_\tau(\Sigma_{3r})}^2 + \|e^{\tau\varphi}v\|_{H^1_\tau(\Sigma_{3r})}^2. \end{aligned}$$

This completes the proof of Lemma 4.1. \square

Now we can argue similarly to [6], with suitable modifications. We introduce a cutoff function χ_4 satisfying $0 \leq \chi_4 \leq 1$, $\chi_4 \in C_0^\infty(\mathbb{R})$, and

$$\chi_4(\eta) = \begin{cases} 0, & \text{if } \eta \leq \frac{1}{2}, \eta \geq 8, \\ 1, & \text{if } \eta \in [\frac{3}{4}, 7]. \end{cases}$$

Now we proceed to the estimation near Γ_1 .

4.2. Estimation near the boundary part Γ_1 . We shall estimate $\mathbf{u}_{\gamma,t}$ in a ball $B_1 = B(x^{(1)}, r) = \{x \in \mathbb{R}^3; |x - x^{(0)}| \leq r\}$ over a small interval $(-r, r)$ by the velocity trace (in the normal direction) in the given part $\Sigma_{1,3r} = \Gamma_1 \times (-3r, 3r) \subset \Sigma_{3r}$.

LEMMA 4.2. *Let $\mathbf{u}_{\gamma,t}$ be a solution to (4.6). Then there exist $\tilde{B}_1 \equiv B_1 \times (-r, r) \subset \Omega_r$ and $\nu_0 \in]0, 1[$ such that*

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \leq C \left(\|F_{\gamma,t}\|_{H^1(\Omega_{3r})} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})} \right)^{\nu_0} \left(\|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})} \right)^{1-\nu_0} \tag{4.11}$$

for some positive constant C .

Proof. Let us choose $\delta > 0$ and $x^{(0)} \in \mathbb{R}^3 \setminus \bar{\Omega}$ such that

$$\delta < \frac{\epsilon_0}{4}, \quad \overline{B(x^{(0)}, \delta)} \cap \bar{\Omega} = \emptyset, \quad B(x^{(0)}, 2\delta) \cap \Omega \neq \emptyset, \quad B(x^{(0)}, 4\delta) \cap \Gamma \subset \Gamma_1. \tag{4.12}$$

That is, $x^{(0)}$ is an outer point of $\bar{\Omega}$ and is near Γ_1 . We define the functions $\psi_0(x, s)$ and $\varphi_0(x, s)$ by

$$\psi_0(x, s) = \left| x - x^{(0)} \right|^2 + s^2, \quad \varphi_0(x, s) = e^{-\frac{\tilde{\beta}}{\delta^2} \psi_0(x, s)}.$$

Here we choose $\tilde{\beta} > 0$ sufficiently large. Denote

$$\mathbf{w}_{\gamma,t}(x, s) = \chi_4 \left(\frac{\psi_0}{\delta^2} \right) \mathbf{u}_{\gamma,t}(x, s).$$

By applying Lemma 4.1, we obtain

$$\frac{C}{\tau} \|e^{\tau\varphi_0}\mathbf{w}_{\gamma,t}\|_{H^2_\tau(\Omega_r)}^2 \leq \|e^{\tau\varphi_0}Q_{\rho,\lambda,\mu}\mathbf{w}_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 + \left\| e^{\tau\varphi_0}\chi_4 \left(\frac{\psi_0}{\delta^2} \right) \mathbf{u}_{\gamma,t} \right\|_{H^2_\tau(\Sigma_{3r})}^2$$

for $\tau > \tau_0$. Therefore by (4.9), we have

$$\begin{aligned} Q_{\rho,\lambda,\mu} \mathbf{w}_{\gamma,t}(x, s) &= \chi_4 \left(\frac{\psi_0}{\delta^2} \right) Q_{\rho,\lambda,\mu} \mathbf{u}_{\gamma,t}(x, s) + \left[Q_{\rho,\lambda,\mu}, \chi_4 \left(\frac{\psi_0}{\delta^2} \right) \right] \mathbf{u}_{\gamma,t}(x, s) \\ &= \chi_4 \left(\frac{\psi_0}{\delta^2} \right) (F_{\gamma,t}(x, s) + G_{\gamma,t}(x, s)) + \left[Q_{\rho,\lambda,\mu}, \chi_4 \left(\frac{\psi_0}{\delta^2} \right) \right] \mathbf{u}_{\gamma,t}(x, s) \\ &= \chi_4 \left(\frac{\psi_0}{\delta^2} \right) F_{\gamma,t}(x, s) + \left[Q_{\rho,\lambda,\mu}, \chi_4 \left(\frac{\psi_0}{\delta^2} \right) \right] \mathbf{u}_{\gamma,t}(x, s). \end{aligned}$$

Since $[Q_{\rho,\lambda,\mu}, \chi_4(\frac{\psi_0}{\delta^2})]$ is supported in

$$\left| x - x^{(0)} \right|^2 + s^2 \leq \frac{3}{4} \delta^2, \quad 7\delta^2 \leq \left| x - x^{(0)} \right|^2 + s^2 \leq 8\delta^2,$$

by taking (4.12) into account, we see that $|x - x^{(0)}| \geq \delta$ for all $x \in \bar{\Omega}$ and $\Omega \cap \{x; |x - x^{(0)}| \leq \frac{3}{4} \delta^2\} \neq \emptyset$, so that we obtain

$$\begin{aligned} \frac{C}{\tau} e^{2\tau e^{-4\tilde{\beta}}} \|\mathbf{u}_{\tau,t}\|_{H^2_\tau((\delta^2 \leq \psi_0 \leq 4\delta^2) \cap \Omega)}^2 &\leq \tau^2 e^{2\tau e^{-7\tilde{\beta}}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_0 \leq 8\delta^2)}^2 \\ &\quad + \tau^2 e^{2\tau e^{-\tilde{\beta}/2}} \|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 \\ &\quad + e^{2\tau e^{-\tilde{\beta}/2}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}^2. \end{aligned}$$

We can select $r > 0$ and $x^{(1)} \in \Omega$ such that

$$\text{dist}(x^{(1)}, \Gamma) \geq 4r, \quad \tilde{B}_1 = B(x^{(1)}, r) \times [-r, r] \subset \{\delta^2 \leq \psi_0(x, s) \leq 4\delta^2\}.$$

Then for $\tau > \tau_0$ we have

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)}^2 \leq C e^{C_1\tau} \left[\|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 + \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}^2 \right] + e^{-C_2\tau} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2.$$

The inequality holds for any $\tau > 0$ by replacing $C > 0$ by $C e^{C_1\tau_0}$. Now minimize the right-hand side with respect to τ , and with $\nu_0 = \frac{C_2}{C_1+C_2}$ we have

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)}^2 \leq C \left(\|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 + \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}^2 \right)^{\nu_0} \left(\|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2 \right)^{1-\nu_0}.$$

This completes the proof of the lemma. □

4.3. Estimation near the boundary. In this subsection we extend the estimation from \tilde{B}_1 to $\omega_r(\epsilon_0, 4\epsilon_0)$. In order to accomplish this, we use the techniques developed in [44]. This will be done by continuing estimates (4.11). Let $B(x^{(j)}, r)$, $2 \leq j \leq N$, be a covering of $\omega(\epsilon_0, 4\epsilon_0)$. We can assume that $x^{(j)}$ satisfies $\text{dist}(x^{(j)}, \Gamma) \geq 4r$. In what follows, we assume without any restriction in generality that

$$B(x^{(j+1)}, r) \subset B(x^{(j)}, 2r),$$

and we set

$$\tilde{B}_j = B(x^{(j)}, r) \times (-r, r); \quad 2 \leq j \leq N.$$

LEMMA 4.3. *Let $\mathbf{u}_{\gamma,t}$ be a solution to (4.6). Then there exist a constant $\nu \in (0, 1)$ and $C > 0$ such that*

$$(4.13) \quad \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_{k+1})} \leq C \left(\|F_{\gamma,t}\|_{H^1(\Omega_{3r})} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_k)} \right)^\nu \left(\|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})} \right)^{1-\nu}, \quad k \geq 1.$$

Proof. We define the functions $\psi_k(x, s)$ and $\varphi_k(x, s)$ by

$$\psi_k(x, s) = \left| x - x^{(k)} \right|^2 + s^2, \quad \varphi_k(x, s) = e^{-\frac{\tilde{\beta}}{r^2} \psi_k(x, s)}.$$

Moreover we set

$$\mathbf{w}_{\gamma,t}(x, s) = \chi_4 \left(\frac{\psi_k}{r^2} \right) \mathbf{u}_{\gamma,t}(x, s).$$

By applying Lemma 4.1 in an interior domain, we obtain

$$(4.14) \quad \frac{C}{\tau} \|e^{\tau\varphi_k} \mathbf{w}_{\gamma,t}\|_{H^2_\tau(\Omega_r)}^2 \leq \|e^{\tau\varphi_k} Q_{\rho,\lambda,\mu} \mathbf{w}_{\gamma,t}\|_{H^1(\Omega_{3r})}^2.$$

In the same way as the proof of Lemma 4.2, we have

$$Q_{\rho,\lambda,\mu} \mathbf{w}_{\lambda,t}(x, s) = \chi_4 \left(\frac{\psi_k}{r^2} \right) F_{\lambda,t}(x, s) + \left[Q_{\rho,\lambda,\mu}, \chi_4 \left(\frac{\psi_k}{r^2} \right) \right] \mathbf{u}_{\gamma,t}(x, s).$$

Since $[Q_{\rho,\lambda,\mu}, \chi_4(\frac{\psi_k}{r^2})]$ is supported in

$$\frac{r^2}{2} \leq \left| x - x^{(0)} \right|^2 + s^2 \leq r^2, \quad 7r^2 \leq \left| x - x^{(0)} \right|^2 + s^2 \leq 8r^2,$$

it follows from (4.14) that

$$\begin{aligned} \frac{C}{\tau} e^{2\tau e^{-5\tilde{\beta}}} \|\mathbf{u}_{\gamma,t}\|_{H^2_\tau(r^2 \leq \psi_k \leq 5r^2)}^2 &\leq \tau^2 e^{2\tau e^{-\tilde{\beta}/2}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq r^2)}^2 + \tau^2 e^{2\tau e^{-7\tilde{\beta}}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2 \\ &\quad + \tau^2 e^{2\tau e^{-\tilde{\beta}/2}} \|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2, \end{aligned}$$

and hence for τ large we obtain

$$\begin{aligned} C e^{2\tau e^{-5\tilde{\beta}}} \|\mathbf{u}_{\gamma,t}\|_{H^2_\tau(\psi_k \leq 5r^2)}^2 &\leq e^{2\tau e^{-\tilde{\beta}/3}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq r^2)}^2 + e^{2\tau e^{-6\tilde{\beta}}} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2 \\ &\quad + e^{2\tau e^{-\tilde{\beta}/3}} \|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2. \end{aligned}$$

Thus we obtain

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq 5r^2)}^2 \leq e^{C_1\tau} \left[\|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq r^2)}^2 + \|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 \right] + e^{-C_2\tau} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2.$$

Now minimize the right-hand side with respect to τ , and with $\nu = \frac{C_2}{C_1+C_2}$ we obtain

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq 5r^2)}^2 \leq C \left(\|F_{\gamma,t}\|_{H^1(\Omega_{3r})}^2 + \|\mathbf{u}_{\gamma,t}\|_{H^2(\psi_k \leq r^2)}^2 \right)^\nu \left(\|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}^2 \right)^{1-\nu}.$$

Since

$$\tilde{B}_{k+1} \subset \{ \psi_k(s, x) \leq 5r^2 \}, \quad \{ \psi_k(x, s) \leq r^2 \} \subset \tilde{B}_k,$$

we obtain (4.13). This completes the proof of the lemma. \square

LEMMA 4.4. *Let $\mathbf{u}_{\gamma,t}$ be a solution to (4.6). Then there exist a constant $C > 0$ and $\mu = \nu^n$ such that*

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq C \left(\|F_{\gamma,t}\|_{H^1(\Omega_{3r})} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \right)^\mu \left(\|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})} \right)^{1-\mu}, \quad n \geq 1.$$

Here $\nu \in (0,1)$ is the constant given in Lemma 4.3.

Proof. Put

$$a_k = \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_k)}, \quad A = \|F_{\gamma,t}\|_{H^1(\Omega_{3r})}, \quad B = \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})}.$$

By (4.13) we have

$$a_{k+1} \leq (C^{\frac{1}{1-\nu}} B)^{1-\nu} (a_k + A)^\nu.$$

By applying Lemma 4 in [38] (see also [37]), we obtain for all $\mu \in]0, \nu^n]$

$$a_n \leq 2^{\frac{1}{1-\nu}} C B^{1-\mu} (a_1 + A)^\mu.$$

This completes the proof of the lemma. \square

LEMMA 4.5. *Let $\mathbf{u}_{\lambda,t}$ be a solution to (4.6). Then there exist $C > 0$ and $C_1 > 0$ such that for all n there exist $C_n > 0$ and $T_n > 0$ such that*

$$(4.15) \quad C \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)}^2 \leq e^{-C_1\gamma} \|\mathbf{u}\|_{H^2(Q)}^2 + e^{C_n\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}^2$$

for all $t \in [-\frac{T}{2}, \frac{T}{2}]$, where $T > T_n$.

Proof. By Lemma 4.4 and the Young inequality, we easily obtain

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq \epsilon^p \|\mathbf{u}_{\gamma,t}\|_{H^2(\Omega_{3r})} + \epsilon^{-p'} \left[\|F_{\gamma,t}\|_{H^1(\Omega_{3r})} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \right]$$

for all $\epsilon > 0$. Here

$$p = \frac{1}{1-\mu}, \quad p' = \frac{1}{\mu}, \quad \text{and } \mu = \nu^n.$$

By using estimates (4.7) and (4.8), we have for all $t \in [-\frac{T}{2}, \frac{T}{2}]$

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq \epsilon^p e^{C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} + \epsilon^{-p'} \left[e^{-\eta T\gamma} \|\mathbf{u}\|_{H^2(Q)} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \right].$$

By selecting

$$\epsilon = e^{-\frac{2C_1}{p}\gamma},$$

we obtain

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq e^{-C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{-(\eta T - \frac{2C_1 p'}{p})\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{\frac{2C_1 p'}{p}\gamma} \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)}$$

for all $t \in [-\frac{T}{2}, \frac{T}{2}]$ and $\gamma > 0$. Take T sufficiently large such that

$$\eta T - \frac{2C_1 p'}{p} > C_1,$$

and we obtain

$$(4.16) \quad \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq e^{-C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{\kappa_1\gamma} \|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)},$$

where we set

$$\kappa_1 = \frac{2C_1 p'}{p}.$$

Similarly we obtain from Lemma 4.2 and the Young inequality

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \leq \epsilon^{p_0} e^{C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} + \epsilon^{-p'_0} \left[e^{-\eta T\gamma} \|\mathbf{u}\|_{H^2(Q)} + \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})} \right],$$

where

$$p_0 = \frac{1}{1-\nu_0}, \quad p'_0 = \frac{1}{\nu_0}.$$

Selecting $\epsilon = e^{-(\frac{2C_1+\kappa_1}{p_0})\gamma}$, we obtain for some positive constant κ_2

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \leq e^{-(C_1+\kappa_1)\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{-(\eta T - \frac{(2C_1+\kappa_1)p'_0}{p_0})\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{\kappa_2\gamma} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}.$$

Take T large such that

$$\left(\eta T - \frac{(2C_1 + \kappa_1)p'_0}{p_0} \right) > C_1 + \kappa_1.$$

Then we obtain

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_1)} \leq e^{-(C_1+\kappa_1)\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{\kappa_2\gamma} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}.$$

By inserting this into (4.16), we have

$$\|\mathbf{u}_{\gamma,t}\|_{H^2(\tilde{B}_n)} \leq e^{-C_1\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{C_n\gamma} \|\mathbf{u}_{\gamma,t}\|_{H^2(\Sigma_{1,3r})}$$

for some positive constant C_n . This completes the proof of (4.15). \square

4.4. End of the proof of Lemma 2.5. We shall complete the proof of Lemma 2.5 in this subsection.

We fix $T > \max_{1 \leq n \leq N} T_n$. Addition of inequalities (4.15) for $n \in \{1, \dots, N\}$ yields

$$(4.17) \quad \|\mathbf{u}_{\gamma,t}\|_{H^2(\omega_r(\epsilon_0, 3\epsilon_0))} \leq e^{-C_3\gamma} \|\mathbf{u}\|_{H^2(Q)} + e^{C_4\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}, \quad t \in \left[-\frac{T}{2}, \frac{T}{2} \right],$$

for some positive constants C_3 and C_4 . We set $\mathbf{u}_\gamma(x, t) = \mathbf{u}_{\gamma,t}(x, 0)$. Then we have

$$\mathbf{u}_\gamma(x, t) = \sqrt{\frac{\gamma}{2\pi}} \int_{\mathbb{R}} e^{-\frac{\gamma}{2}(t-y)^2} \theta(y) \mathbf{u}(x, y) dy = (K_\gamma * \theta \mathbf{u})(x, t),$$

where

$$K_\gamma(t) = \sqrt{\frac{\gamma}{2\pi}} e^{-\frac{\gamma}{2}t^2}.$$

LEMMA 4.6. *Let $T^1 = \frac{T}{2} - r$. Then we have*

$$\|\mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))}^2 \leq e^{-\mu\gamma} \|\mathbf{u}\|_{H^1(Q)}^2 + e^{\mu'\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}^2$$

for some positive constants μ and μ' .

Proof. By the Cauchy formula, for ϱ such that $0 < \varrho < r$, we obtain

$$\mathbf{u}_\gamma(x, a) = \frac{1}{2i\pi} \int_{|w-a|=\varrho} \frac{\mathbf{u}_\gamma(x, w)}{w-a} dw.$$

Thus, by using the polar coordinate, we obtain

$$|\mathbf{u}_\gamma(x, a)|^2 \leq C_5 \int_0^{2\pi} |\mathbf{u}_\gamma(x, a + \varrho e^{i\theta})|^2 d\theta.$$

We integrate $\varrho \in (0, r)$ and obtain

$$|\mathbf{u}_\gamma(x, a)|^2 \leq \frac{C_5}{r} \int_0^r \int_0^{2\pi} |\mathbf{u}_\gamma(x, a + \varrho e^{i\theta})|^2 d\theta d\varrho.$$

Therefore, for $x \in \omega(\epsilon_0, 3\epsilon_0)$ and $a \in [-\frac{T}{2} + r, \frac{T}{2} - r]$, we have

$$|\mathbf{u}_\gamma(x, a)|^2 \leq C_6 \int_{|s|\leq r, |t-a|\leq r} |\mathbf{u}_\gamma(x, t + is)|^2 ds dt = C_6 \int_{|s|\leq r, |t-a|\leq r} |\mathbf{u}_{\gamma,t}(x, s)|^2 ds dt$$

and

$$\begin{aligned} \|\mathbf{u}_\gamma(\cdot, a)\|_{L^2(\omega(\epsilon_0, 3\epsilon_0))}^2 &\leq C_6 \int_{-r}^r \int_{|t-a|\leq r} \|\mathbf{u}_{\gamma,t}\|_{L^2(\omega_r(\epsilon_0, 3\epsilon_0))}^2 dt ds \\ (4.18) \qquad \qquad \qquad &\leq C_7 \int_{-r}^r \int_{-T/2}^{T/2} \|\mathbf{u}_{\gamma,t}\|_{L^2(\omega_r(\epsilon_0, 3\epsilon_0))}^2 dt ds. \end{aligned}$$

We integrate $a \in [-T^1, T^1]$, and by using (4.17) we obtain

$$\|\mathbf{u}_\gamma\|_{L^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))}^2 \leq C_T (e^{-C_3\gamma} \|\mathbf{u}\|_{H^2(Q)}^2 + e^{C_4\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}^2).$$

By using the same argument to $\partial^\alpha \mathbf{u}_\gamma$, $|\alpha| \leq 2$, we complete the proof of Lemma 4.6. \square

LEMMA 4.7. *Let \mathbf{u} be a solution of (4.4). Then there exist $C_7 > 0$ and $C_8 > 0$ such that*

$$\|\mathbf{u}\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \leq \frac{C_7}{\sqrt{\gamma}} \|\mathbf{u}\|_{H^3(Q)} + e^{C_7\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}.$$

Proof. By $\widehat{\mathbf{u}}(x, \tau)$ we denote the Fourier transform of $\mathbf{u}(x, t)$ in t . We have

$$\widehat{\theta\mathbf{u}}(x, \tau) - \widehat{\mathbf{u}}_\gamma(x, \tau) = (1 - \widehat{K}_\gamma) \widehat{\theta\mathbf{u}}(x, \tau).$$

Furthermore we can directly verify that

$$\left| (1 - \widehat{K}_\gamma)(\tau) \right| \leq \frac{\tau^2}{\gamma},$$

so that we obtain for $T^1 = T/2 - r$

$$\|\mathbf{u} - \mathbf{u}_\gamma\|_{L^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \leq \frac{C_8}{\sqrt{\gamma}} \|\mathbf{u}\|_{H^1(Q)}.$$

Similarly we have

$$\|\mathbf{u} - \mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \leq \frac{C_9}{\sqrt{\gamma}} \|\mathbf{u}\|_{H^3(Q)}.$$

Hence

$$\begin{aligned} \|\mathbf{u}\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} &\leq C_{10} \left[\|\mathbf{u} - \mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} + \|\mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \right] \\ &\leq C_{11} \left[\frac{1}{\sqrt{\gamma}} \|\mathbf{u}\|_{H^3(Q)} + \|\mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \right]. \end{aligned}$$

On the other hand, by Lemma 4.6, we obtain

$$\|\mathbf{u}_\gamma\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))}^2 \leq e^{-\mu\gamma} \|\mathbf{u}\|_{H^2(Q)}^2 + e^{\mu'\gamma} \|\mathbf{u}\|_{H^2(\Sigma_1)}^2$$

for some positive constants μ and μ' . This complete the proof of the lemma. \square

We now turn to the proof of Lemma 2.6. By Lemma 4.7 and $\mathbf{u} = \chi_3 \mathbf{v}$, we obtain

$$\|\mathbf{v}\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \leq \frac{C_7}{\sqrt{\gamma}} \|\mathbf{v}\|_{H^3(Q)} + e^{C_7\gamma} \|\mathbf{v}\|_{H^2(\Sigma_1)}.$$

By (1.6), we obtain

$$\|\mathbf{u}\|_{H^2(\omega_{T^1}(\epsilon_0, 3\epsilon_0))} \leq \frac{M}{\sqrt{\gamma}} + e^{C_7\gamma} \epsilon(\Sigma_1).$$

By selecting

$$\gamma = \frac{1}{2C_7} \log \left(2 + \frac{M}{\epsilon(\Sigma_1)} \right),$$

the proof of Lemma 2.5 is complete. \square

Acknowledgment. The authors thank the anonymous referee for very useful comments.

REFERENCES

[1] C. BARDOS, G. LEBEAU, AND J. RAUCH, *Sharp sufficient conditions for the observation, control, and stabilization of waves from the boundary*, SIAM J. Control Optim., 30 (1992), pp. 1024–1065.
 [2] L. BAUDOIN AND J.-P. PUEL, *Uniqueness and stability in an inverse problem for the Schrödinger equation*, Inverse Problems, 18 (2002), pp. 1537–1554.
 [3] M. BELLASSOUED, *Unicité et contrôle pour le système de Lamé*, ESAIM Control Optim. Calc. Var., 6 (2001), pp. 561–592.
 [4] M. BELLASSOUED, *Global logarithmic stability in inverse hyperbolic problem by arbitrary boundary observation*, Inverse Problems, 20 (2004), pp. 1033–1052.
 [5] M. BELLASSOUED, *Uniqueness and stability in determining the speed of propagation of second-order hyperbolic equation with variable coefficients*, Appl. Anal., 83 (2004), pp. 983–1014.

- [6] M. BELLAÏOUE AND M. YAMAMOTO, *Logarithmic stability in determination of a coefficient in an acoustic equation by arbitrary boundary observation*, J. Math. Pures Appl. 85 (2006), pp. 193–224.
- [7] M. BELLAÏOUE AND M. YAMAMOTO, *Determination of a Coefficient in the Wave Equation with a Single Measurement*, preprint.
- [8] A. L. BUKHGEIM, *Introduction to the Theory of Inverse Problems*, VSP, Utrecht, 2000.
- [9] A. L. BUKHGEIM, J. CHENG, V. ISAKOV, AND M. YAMAMOTO, *Uniqueness in determining damping coefficients in hyperbolic equations*, in Analytic Extension Formulas and Their Applications, S. Saitoh et al., eds., Springer, New York, 2001, pp. 27–46.
- [10] A. L. BUKHGEIM AND M. V. KLIBANOV, *Global uniqueness of class of multidimensional inverse problems*, Soviet Math. Dokl., 24 (1981), pp. 244–247.
- [11] T. CARLEMAN, *Sur un problème d'unicité pour les systèmes d'équations aux dérivées partielles à deux variables indépendentes*, Ark. Mat. Astr. Fys. 2B (1939), pp. 1–9.
- [12] M. ELLER, V. ISAKOV, G. NAKAMURA, AND D. TATARU, *Uniqueness and stability in the Cauchy problem for Maxwell and elasticity systems*, in Nonlinear Partial Differential Equations and Their Applications, Collège de France Seminar, Vol. 14, North-Holland, Amsterdam, 2002, pp. 329–349.
- [13] A. V. FURSIKOV AND O. YU. IMANUVILOV, *Controllability of Evolution Equations*, Seoul National University, Seoul, 1996.
- [14] L. HÖRMANDER, *Linear Partial Differential Operators*, Springer, Berlin, 1963.
- [15] M. IKEHATA, G. NAKAMURA, AND M. YAMAMOTO, *Uniqueness in inverse problems for the isotropic Lamé system*, J. Math. Sci. Univ. Tokyo, 5 (1998), pp. 627–692.
- [16] O. YU. IMANUVILOV, *Controllability of parabolic equations*, Sb. Math., 186 (1995), pp. 879–900.
- [17] O. YU. IMANUVILOV, *On Carleman estimates for hyperbolic equations*, Asymptot. Anal., 32 (2002), pp. 185–220.
- [18] O. YU. IMANUVILOV, V. ISAKOV, AND M. YAMAMOTO, *An inverse problem for the dynamical Lamé system with two sets of boundary data*, Comm. Pure Appl. Math., 56 (2003), pp. 1366–1382.
- [19] O. YU. IMANUVILOV AND J.-P. PUEL, *Global Carleman estimates for weak solutions of elliptic nonhomogeneous Dirichlet problems*, Int. Math. Res. Not., 16 (2003), pp. 883–913.
- [20] O. YU. IMANUVILOV AND M. YAMAMOTO, *Lipshitz stability in inverse parabolic problems by Carleman estimate*, Inverse Problems, 14 (1998), pp. 1229–1249.
- [21] O. YU. IMANUVILOV AND M. YAMAMOTO, *Global Lipschitz stability in an inverse hyperbolic problem by interior observations*, Inverse Problems, 17 (2001), pp. 717–728.
- [22] O. YU. IMANUVILOV AND M. YAMAMOTO, *Determination of a coefficient in an acoustic equation with single measurement*, Inverse Problems, 19 (2003), pp. 157–171.
- [23] O. YU. IMANUVILOV AND M. YAMAMOTO, *Carleman estimates for the non-stationary Lamé system and the application to an inverse problem*, ESAIM Control Optim. Calc. Var., 11 (2005), pp. 1–56.
- [24] O. YU. IMANUVILOV AND M. YAMAMOTO, *Carleman estimates for the three-dimensional non-stationary Lamé system and applications to an inverse problem*, in Control Theory of Partial Differential Equations, Lect. Notes Pure Appl. Math. 242, O. Imanuvilov, G. Leugering, R. Triggiani, and B.-Y. Zhang, eds., Chapman & Hall/CRC, Boca Raton, FL, 2005, pp. 337–374.
- [25] V. ISAKOV, *A nonhyperbolic Cauchy problem for $\square_b \square_c$ and its applications to elasticity theory*, Comm. Pure Appl. Math., 39 (1986), pp. 747–767.
- [26] V. ISAKOV, *Inverse Problems for Partial Differential Equations*, Springer, Berlin, 2005.
- [27] V. ISAKOV AND M. YAMAMOTO, *Carleman estimates with the Neumann boundary condition and its applications to the observability inequality and inverse problems*, Contemp. Math., 268 (2000), pp. 191–225.
- [28] M. A. KAZEMI AND M. V. KLIBANOV, *Stability estimates for ill-posed Cauchy problems involving hyperbolic equations and inequality*, Appl. Anal., 50 (1993), pp. 93–102.
- [29] A. KHAÏDAROV, *On stability estimates in multidimensional inverse problems for differential equation*, Soviet Math. Dokl., 38 (1989), pp. 614–617.
- [30] M. V. KLIBANOV, *Inverse problems in the “large” and Carleman bounds*, Differ. Equ., 20 (1984), pp. 755–760.
- [31] M. V. KLIBANOV, *Inverse problems and Carleman estimates*, Inverse Problems, 8 (1992), pp. 575–596.
- [32] M. V. KLIBANOV AND J. MALINSKY, *Newton-Kantorovich method for 3-dimensional potential inverse scattering problem and stability of the hyperbolic Cauchy problem with time dependent data*, Inverse Problems, 7 (1991), pp. 577–595.
- [33] M. V. KLIBANOV AND A. TIMONOV, *Carleman Estimates for Coefficient Inverse Problems and*

- Numerical Applications*, VSP, Utrecht, 2004.
- [34] M. V. KLIBANOV AND M. YAMAMOTO, *Lipschitz stability of an inverse problem for an acoustic equation*, *Appl. Anal.*, 85 (2006), pp. 515–538.
 - [35] M. KUBO, *Uniqueness in inverse hyperbolic problems: Carleman estimate for boundary value problems*, *J. Math. Kyoto Univ.*, 40 (2000), pp. 451–473.
 - [36] M. M. LAVRENT'EV, V. G. ROMANOV, AND S. P. SHISHAT-SKIĬ, *Ill-posed Problems of Mathematics Physics and Analysis*, AMS, Providence, RI, 1986.
 - [37] G. LEBEAU AND L. ROBBIANO, *Contrôle exact de l'équation de la chaleur*, *Comm. Partial Differential Equations*, 20 (1995), pp. 335–356.
 - [38] G. LEBEAU AND L. ROBBIANO, *Stabilisation de l'équation des ondes par le bord*, *Duke Math. J.*, 86 (1997), pp. 465–491.
 - [39] S. LI, *An inverse problem for Maxwell's equations in bi-isotropic media*, *SIAM J. Math. Anal.*, 37 (2005), pp. 1027–1043.
 - [40] J. P. PUEL AND M. YAMAMOTO, *On a global estimate in a linear inverse hyperbolic problem*, *Inverse Problems*, 12 (1996), pp. 995–1002.
 - [41] J. P. PUEL AND M. YAMAMOTO, *Generic well ill-posedness in a multidimensional hyperbolic inverse problem*, *J. Inverse Ill-Posed Probl.*, 5 (1997), pp. 55–83.
 - [42] L. RACHELE, *An inverse problem in elastodynamics: Uniqueness of the wave speeds in the interior*, *J. Differential Equations*, 162 (2000), pp. 300–325.
 - [43] L. ROBBIANO, *Théorème d'unicité adapté au contrôle des solutions des problèmes hyperboliques*, *Comm. Partial Differential Equations*, 16 (1991), pp. 789–800.
 - [44] L. ROBBIANO, *Fonction de coût et contrôle des solutions des équations hyperboliques*, *Asymptot. Anal.*, 10 (1995), pp. 95–115.
 - [45] L. ROBBIANO AND C. ZUILY, *Uniqueness in the Cauchy problem for operators with partially holomorphic coefficients*, *Invent. Math.*, 131 (1998), pp. 493–539.
 - [46] D. TATARU, *Carleman estimates and unique continuation for solutions to boundary value problems*, *J. Math. Pures Appl.*, 75 (1996), pp. 367–408.
 - [47] R. TRIGGIANI AND P. F. YAO, *Carleman estimates with no lower-order terms for general Riemann wave equations. Global uniqueness and observability in one shot*, *Appl. Math. Optim.*, 46 (2002), pp. 331–375.
 - [48] V. G. YAKHNO, *Inverse Problems for Differential Equations of Elasticity*, Nauka, Novosibirsk, 1990.
 - [49] M. YAMAMOTO, *Uniqueness and stability in multidimensional hyperbolic inverse problems*, *J. Math. Pures Appl.*, 78 (1999), pp. 65–98.

THE HUNTER–SAXTON EQUATION: A GEOMETRIC APPROACH*

JONATAN LENELLS†

Abstract. We provide a rigorous foundation for the geometric interpretation of the Hunter–Saxton equation as the equation describing the geodesic flow of the \dot{H}^1 right-invariant metric on the quotient space $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$ of the infinite-dimensional Banach manifold $\mathcal{D}^k(\mathbb{S})$ of orientation-preserving H^k -diffeomorphisms of the unit circle \mathbb{S} modulo the subgroup of rotations $Rot(\mathbb{S})$. Once the underlying Riemannian structure has been established, the method of characteristics is used to derive explicit formulas for the geodesics corresponding to the \dot{H}^1 right-invariant metric, yielding, in particular, new explicit expressions for the spatially periodic solutions of the initial-value problem for the Hunter–Saxton equation.

Key words. diffeomorphism group, Hunter–Saxton equation, geodesic flow

AMS subject classifications. 35Q53, 53C21, 58D30

DOI. 10.1137/050647451

1. Introduction. The Hunter–Saxton equation

$$(1.1) \quad u_{txx} = -2u_x u_{xx} - uu_{xxx}, \quad t > 0, \quad x \in \mathbb{R},$$

models the propagation of weakly nonlinear orientation waves in a massive nematic liquid crystal director field, x being the space variable in a reference frame moving with the unperturbed wave speed and t being a slow time variable [6]. Equation (1.1) is a bivariational, completely integrable system with a bi-Hamiltonian structure [7], leading to the existence of an infinite family of commuting Hamiltonian flows together with associated conservation laws. Smooth solutions of (1.1) break down in finite time as the slope $u_x \rightarrow -\infty$. The existence of global weak solutions was considered in [8]. Weak dissipative solutions of (1.1) have been studied in [1], where also a discussion of different notions of weak solutions can be found.

It was first noticed in [9] that, for spatially periodic functions, (1.1) describes the geodesic flow on the homogeneous space $Rot(\mathbb{S}) \backslash \mathcal{D}(\mathbb{S})$ of the infinite-dimensional Lie group $\mathcal{D}(\mathbb{S})$ of orientation-preserving diffeomorphisms of the unit circle \mathbb{S} modulo the subgroup of rotations $Rot(\mathbb{S})$, endowed with the \dot{H}^1 right-invariant metric given at the identity by

$$\langle [u], [v] \rangle = \int_{\mathbb{S}} u_x v_x dx, \quad [u], [v] \in T_{[id]}(Rot(\mathbb{S}) \backslash \mathcal{D}(\mathbb{S})),$$

where $[id]$ denotes the equivalence class in $Rot(\mathbb{S}) \backslash \mathcal{D}(\mathbb{S})$ of the identity map $id \in \mathcal{D}(\mathbb{S})$. This makes (1.1) one of the equations (others include the well-known Euler, Burgers, Korteweg–de Vries, and Camassa–Holm equations [4, 3, 15, 13]) that arises as the Euler equation for the geodesic flow corresponding to a right-invariant metric.

In this paper we provide a rigorous foundation for this geometric interpretation of (1.1) within the periodic setting (for the case on the line further technical complications arise due to the need of considering weighted Sobolev spaces—see the discussion

*Received by the editors December 13, 2005; accepted for publication (in revised form) September 14, 2007; published electronically April 16, 2008.

<http://www.siam.org/journals/sima/40-1/64745.html>

†Department of Mathematics, University of California, Santa Barbara, CA 93106 (jonatan@math.ucsb.edu).

in [2]). Two similar but distinct approaches are possible. We may choose to consider the geodesic flow of the \dot{H}^1 right-invariant metric either (1) on the homogeneous space $Rot(\mathbb{S}) \setminus \mathcal{D}(\mathbb{S})$, where $\mathcal{D}(\mathbb{S})$ denotes the Fréchet Lie group of smooth diffeomorphisms of \mathbb{S} , or (2) on the Banach manifold $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ for $k > 3/2$, where $\mathcal{D}^k(\mathbb{S})$ incorporates all diffeomorphisms of \mathbb{S} of Sobolev class H^k . We will pursue the latter approach. Note that $\mathcal{D}^k(\mathbb{S})$ is a topological group but not a Lie group as the group operation $(\psi, \varphi) \mapsto \psi \circ \varphi$ for $\psi, \varphi \in \mathcal{D}^k(\mathbb{S})$ is continuous but not smooth due to derivative loss (cf. [5]). The advantage of working on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ is that the theory of Riemannian geometry on Banach manifolds is available—whereas nearly all results familiar from finite-dimensional Riemannian geometry immediately generalize to Banach manifolds (see [10]); a transition to Fréchet manifolds introduces several technical complications (see [5]). In particular, there are no general existence and uniqueness results for differential equations in Fréchet spaces, making it difficult to study geodesic flow and parallel translation. Also, on a Fréchet manifold, since the inverse mapping theorem does not hold, neither the Lie group exponential map nor the Riemannian exponential map is necessarily a local diffeomorphism at the identity. Another advantage of working with the wider class $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ is that when studying partial differential equations it is often preferable to work in Sobolev spaces rather than in the category of C^∞ -maps.

We will construct a smooth affine connection on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ compatible with the \dot{H}^1 right-invariant metric. Once this has been established, the general theory of Riemannian geometry on Banach manifolds immediately yields existence of a smooth curvature tensor, existence of normal neighborhoods, existence and uniqueness results for the geodesic flow and parallel translation, locally length-minimizing properties of the geodesics, etc. In this paper we focus on the geodesic flow and its relation to the Hunter–Saxton equation. In particular, we find that if $\varphi \in C^2([0, T]; \mathcal{D}^k(\mathbb{S}))$ is a C^2 -curve in $\mathcal{D}^k(\mathbb{S})$, then the quotient curve $t \mapsto [\varphi(t)]$ is a geodesic in $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ with respect to the \dot{H}^1 right-invariant metric if and only if $u = \varphi_t \circ \varphi^{-1} \in C([0, T]; H^k(\mathbb{S})) \cap C^1([0, T]; H^{k-1}(\mathbb{S}))$ satisfies (1.1). In a subsequent section the method of characteristics is used to derive explicit formulas for the geodesics corresponding to the \dot{H}^1 right-invariant metric. As a byproduct we also obtain explicit formulas for the spatially periodic solutions of the initial-value problem for (1.1). Note that this presents new solutions to (1.1) with respect to earlier presentations (cf. [6, 1])—the requirement that the solutions be spatially periodic introduces nontrivial constants of integration which alter the derivation and its outcome.

The quotient space $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$, the \dot{H}^1 right-invariant metric $\langle \cdot, \cdot \rangle$, and a Christoffel map Γ are introduced in section 2. In section 3 it is proved that the covariant derivative induced by Γ is the unique covariant derivative compatible with the \dot{H}^1 right-invariant metric. In section 4 the connection between the geodesic flow and the Hunter–Saxton equation is explained, while in section 5 explicit formulas for the geodesics are obtained by means of the method of characteristics.

2. The quotient space $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$. Let \mathbb{S} be the circle of length one and let D_x denote differentiation with respect to x . For $X = [0, 1]$ or $X = \mathbb{S}$ we let, for $n \geq 0$, $H^n(X)$ be the space of all functions on X of Sobolev class H^n . By restriction of a periodic function to the unit interval, $H^n(\mathbb{S})$ may be viewed as a closed linear subspace of $H^n[0, 1]$.

For an integer $k \geq 3$,¹ let $\mathcal{D}^k(\mathbb{S})$ denote the Banach manifold of orientation-

¹Even though $k > 3/2$ is sufficient for $\mathcal{D}^k(\mathbb{S})$ to be a topological group [4], we will assume $k \geq 3$ for simplicity; see [11] for an extension of the geometric approach which incorporates weak solutions.

preserving diffeomorphisms of \mathbb{S} of class H^k (cf. [14]). By $Rot(\mathbb{S}) \subset \mathcal{D}^k(\mathbb{S})$ we denote the subgroup of rotations $x \mapsto x + d$, $d \in \mathbb{R}$. Let $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$ be the space of right cosets $Rot(\mathbb{S}) \circ \varphi = \{\varphi(\cdot) + d \mid d \in \mathbb{R}\}$ for $\varphi \in \mathcal{D}^k(\mathbb{S})$. We will construct a global canonical chart on $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$.

Put $M^k = \{\varphi \in \mathcal{D}^k(\mathbb{S}) \mid \varphi(0) = 0\}$. Then the map

$$(2.1) \quad \varphi \mapsto (\varphi(0), \varphi(\cdot) - \varphi(0)) : \mathcal{D}^k(\mathbb{S}) \rightarrow \mathbb{S} \times M^k$$

is a diffeomorphism. Note that M^k can be characterized as

$$M^k = \{\varphi \in H^k[0, 1] \mid \varphi_x \in H^{k-1}(\mathbb{S}), \varphi_x > 0, \varphi(0) = 0, \varphi(1) = 1\},$$

or, equivalently,

$$(2.2) \quad M^k = \{u + id \mid u \in H^k(\mathbb{S}), u_x > -1, u(0) = 0\},$$

where $id \in \mathcal{D}^k(\mathbb{S})$ is the identity map $id(x) = x$ for $x \in \mathbb{S}$.

From (2.1) we obtain a natural identification $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S}) \simeq M^k$ given by

$$(2.3) \quad [\varphi] \mapsto \varphi - \varphi(0),$$

where $[\varphi]$ denotes the equivalence class of $\varphi \in \mathcal{D}^k(\mathbb{S})$. Let $\mathbf{E}^k \subset H^k(\mathbb{S})$ be the closed linear subspace

$$\mathbf{E}^k = \{u \in H^k(\mathbb{S}) \mid u(0) = 0\}$$

with topology induced from $H^k(\mathbb{S})$. The representation (2.2) shows that M^k is an open subset of the closed hyperplane $id + \mathbf{E}^k \subset H^k[0, 1]$. Hence M^k provides a global chart for $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$. Moreover, $T(Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})) \simeq TM^k \simeq M^k \times \mathbf{E}^k$, so that a vector field X on $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$ can be viewed as a map $M^k \rightarrow \mathbf{E}^k$.

The \dot{H}^1 right-invariant metric $\langle \cdot, \cdot \rangle$ on $Rot(\mathbb{S}) \backslash \mathcal{D}^k(\mathbb{S})$ is most easily defined in the global chart M^k as follows (see also [9]). Let $A = -D_x^2$ and introduce a positive definite symmetric bilinear form $\langle \cdot, \cdot \rangle_{id}$ on $T_{id}M^k \simeq \mathbf{E}^k$ by

$$\langle u, v \rangle_{id} = \int_{\mathbb{S}} uAv dx = \int_{\mathbb{S}} u_x v_x dx, \quad u, v \in \mathbf{E}^k.$$

and extend it to $T_{\varphi}M^k$ for any $\varphi \in M^k$ by right-invariance, so that, for $U, V \in T_{\varphi}M^k \simeq \mathbf{E}^k$,

$$\langle U, V \rangle_{\varphi} = \int_{\mathbb{S}} U \circ \varphi^{-1} A (V \circ \varphi^{-1}) dx.$$

We also introduce the closed linear subspace

$$\mathbf{F}^k = \left\{ f \in H^k(\mathbb{S}) \mid \int_{\mathbb{S}} f dx = 0 \right\}.$$

It is straightforward to check that $A = -D_x^2$ is an isomorphism $\mathbf{E}^k \rightarrow \mathbf{F}^{k-2}$. Denoting its inverse by $A^{-1} : \mathbf{F}^{k-2} \rightarrow \mathbf{E}^k$, we infer that, for $u \in H^{k-1}(\mathbb{S})$,

$$(2.4) \quad -(A^{-1}D_x(u))(x) = \int_0^x u(y)dy - x \int_{\mathbb{S}} u dx.$$

Now define $\Gamma : M^k \times \mathbf{E}^k \times \mathbf{E}^k \rightarrow \mathbf{E}^k$ by

$$(2.5) \quad \Gamma(\varphi, U, V) = -\frac{1}{2} \left(A^{-1} D_x ((U \circ \varphi^{-1})_x (V \circ \varphi^{-1})_x) \right) \circ \varphi,$$

and notice that Γ is right-invariant in the sense that

$$(2.6) \quad \Gamma(\psi, U, V) \circ \varphi = \Gamma(\psi \circ \varphi, U \circ \varphi, V \circ \varphi), \quad \varphi, \psi \in M^k, \quad U, V \in \mathbf{E}^k.$$

The motivation for the definition of Γ comes in section 3, where we will see that Γ is the Christoffel map (the infinite-dimensional analogue of the Christoffel symbols Γ^i_{jk} well known from finite-dimensional Riemannian geometry) for the affine connection on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ compatible with the \dot{H}^1 right-invariant metric.

Although $\mathcal{D}^k(\mathbb{S})$ is a smooth Banach manifold it is not a Lie group. Indeed, the group operation $(\psi, \varphi) \mapsto \psi \circ \varphi : \mathcal{D}^k(\mathbb{S}) \times \mathcal{D}^k(\mathbb{S}) \rightarrow \mathcal{D}^k(\mathbb{S})$ is continuous but not C^1 ; right multiplication $R_\varphi : \psi \mapsto \psi \circ \varphi$ is smooth, whereas left multiplication $L_\psi : \varphi \mapsto \psi \circ \varphi$ is continuous but not C^1 due to derivative loss (see [5]). Similarly, M^k is a Banach manifold and a topological group but not a Lie group. Therefore, it is not a priori clear that the \dot{H}^1 metric and the Christoffel map Γ are smooth objects (that right multiplication is smooth is not enough). The following two propositions deal with this technicality (see [12] for detailed proofs of similar results in the case of the Camassa–Holm equation).

PROPOSITION 2.1. *The map*

$$[\varphi] \mapsto \langle \cdot, \cdot \rangle_{[\varphi]} : Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S}) \rightarrow L^2_{sym}(T_{[\varphi]}(Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})); \mathbb{R})$$

is a smooth section of the bundle $L^2_{sym}(T(Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})); \mathbb{R})$.

PROPOSITION 2.2. *Let Γ be defined by (2.5). The map*

$$\varphi \mapsto \Gamma(\varphi, \cdot, \cdot) : M^k \rightarrow L^2_{sym}(\mathbf{E}^k; \mathbf{E}^k)$$

is smooth.

3. Covariant derivative. In this section the covariant derivative induced by the Christoffel map Γ is shown to be compatible with the \dot{H}^1 right-invariant metric.

We first recall the general definition of a covariant derivative. Let \mathcal{M} be a Banach manifold endowed with a Riemannian metric $\langle \cdot, \cdot \rangle$ and let $\mathfrak{X}(\mathcal{M})$ denote the space of smooth vector fields on \mathcal{M} . For $X, Y \in \mathfrak{X}(\mathcal{M})$ the Lie bracket $[X, Y]$ is defined locally by

$$[X, Y](m) = DY(m) \cdot X(m) - DX(m) \cdot Y(m).$$

DEFINITION 3.1. *An \mathbb{R} -bilinear operator $(X, Y) \mapsto \nabla_X Y : \mathfrak{X}(\mathcal{M}) \times \mathfrak{X}(\mathcal{M}) \rightarrow \mathfrak{X}(\mathcal{M})$ is a Riemannian covariant derivative if it satisfies*

- (a) $X(m) = 0$ implies $(\nabla_X Y)(m) = 0$ for $m \in \mathcal{M}$ and $X, Y \in \mathfrak{X}(\mathcal{M})$ (punctual dependence on X),
- (b) $\nabla_X Y - \nabla_Y X = [X, Y]$ for $X, Y \in \mathfrak{X}(\mathcal{M})$ (torsion-free),
- (c) $\nabla_X(fY) = (\mathcal{L}_X f)Y + f\nabla_X Y$ for $f \in C^\infty(\mathcal{M})$, $X, Y \in \mathfrak{X}(\mathcal{M})$ (derivation in Y),
- (d) $\mathcal{L}_X \langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X Z \rangle$ for $X, Y, Z \in \mathfrak{X}(\mathcal{M})$ (compatible with the metric).

Define the operator $\nabla : \mathfrak{X}(Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})) \times \mathfrak{X}(Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})) \rightarrow \mathfrak{X}(Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S}))$ in the global chart M^k by

$$(3.1) \quad (\nabla_X Y)(\varphi) = DY(\varphi) \cdot X(\varphi) - \Gamma(\varphi, Y(\varphi), X(\varphi)),$$

where $X, Y : M^k \rightarrow \mathbf{E}^k$ are representatives in the chart M^k of two vector fields on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$.

In the finite-dimensional case, given a Riemannian metric $\langle \cdot, \cdot \rangle$ on a manifold \mathcal{M} there automatically exists a Riemannian covariant derivative ∇ compatible with $\langle \cdot, \cdot \rangle$. For vector fields X, Y, Z on \mathcal{M} , $\nabla_X Y$ is defined as the unique vector field such that

$$(3.2) \quad 2\langle \nabla_X Y, Z \rangle = -\langle [Y, X], Z \rangle - \langle X, [Y, Z] \rangle - \langle Y, [X, Z] \rangle + \mathcal{L}_X \langle Y, Z \rangle + \mathcal{L}_Y \langle Z, X \rangle - \mathcal{L}_Z \langle X, Y \rangle.$$

Indeed, the bracket $\langle \cdot, \cdot \rangle$ establishes an isomorphism $T_m \mathcal{M} \rightarrow T_m^* \mathcal{M}$ for each $m \in \mathcal{M}$, so since the right-hand side is a continuous linear functional of $Z(m)$, existence of $(\nabla_X Y)(m)$ follows immediately.

This approach does not apply to $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ endowed with the \dot{H}^1 right-invariant metric. The right-hand side of (3.2) is a continuous linear functional of $Z(\varphi)$ for each $\varphi \in M^k$. But the topology of $T_\varphi M^k \simeq \mathbf{E}^k$ induced by the H^k inner product is much stronger than the topology defined by the \dot{H}^1 right-invariant metric $\langle \cdot, \cdot \rangle_\varphi$ —the \dot{H}^1 right-invariant metric is a *weak Riemannian metric* on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$. Therefore there are elements in $T_\varphi^* M^k$ that cannot be expressed as $\langle V, \cdot \rangle_\varphi$ for some $V \in T_\varphi M^k$; the spaces $T_\varphi M^k \simeq \mathbf{E}^k$ and $T_\varphi^* M^k \simeq \mathbf{E}^{k*}$ are in duality with respect to the H^k inner product, not with respect to $\langle \cdot, \cdot \rangle_\varphi$. The explicit formula for Γ will help us circumvent this difficulty.

However, even for weak Riemannian metrics *uniqueness* of the Riemannian covariant derivative can be deduced from (3.2). For if ∇ satisfies (a)–(d) of Definition 3.1, then, writing down property (d) for the cyclic permutations of $X, Y, Z \in \mathfrak{X}(M)$, we get

$$\begin{aligned} \mathcal{L}_X \langle Y, Z \rangle &= \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X Z \rangle, \\ \mathcal{L}_Y \langle Z, X \rangle &= \langle \nabla_Y Z, X \rangle + \langle Z, \nabla_Y X \rangle, \\ \mathcal{L}_Z \langle X, Y \rangle &= \langle \nabla_Z X, Y \rangle + \langle X, \nabla_Z Y \rangle. \end{aligned}$$

Adding the first two and subtracting the third of these relations, (3.2) drops out. Since $\langle \cdot, \cdot \rangle$ is nondegenerate, (3.2) shows the uniqueness of ∇ .

In the proof of Theorem 3.2 the identity

$$(3.3) \quad \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} U \circ (\varphi + \epsilon V)^{-1} = -(U \circ \varphi^{-1})_x V \circ \varphi^{-1}, \quad \varphi \in \mathcal{D}^k(\mathbb{S}), \quad U, V \in H^k(\mathbb{S}),$$

will be used. It is a consequence of

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} (\varphi + \epsilon V)^{-1} = -\frac{V \circ \varphi^{-1}}{\varphi_x \circ \varphi^{-1}} \quad \text{and} \quad \frac{U_x \circ \varphi^{-1}}{\varphi_x \circ \varphi^{-1}} = (U \circ \varphi^{-1})_x.$$

THEOREM 3.2. *The bilinear map ∇ given by (3.1) defines a unique Riemannian covariant derivative on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ compatible with the \dot{H}^1 right-invariant metric.*

Proof. Properties (a)–(c) are immediate from the local formula defining ∇ . To establish (d) we compute, for vector fields $X, Y, Z : M^k \rightarrow \mathbf{E}^k$,

$$\begin{aligned} & (\mathcal{L}_X \langle Y, Z \rangle)(\varphi) \\ &= \frac{d}{d\epsilon} \Big|_{\epsilon=0} \int_{\mathbb{S}} A(Y(\varphi + \epsilon X(\varphi)) \circ (\varphi + \epsilon X(\varphi))^{-1}) \\ & \quad \cdot Z(\varphi + \epsilon X(\varphi)) \circ (\varphi + \epsilon X(\varphi))^{-1} dx \\ &= \int_{\mathbb{S}} A \left((DY(\varphi) \cdot X(\varphi)) \circ \varphi^{-1} - (Y(\varphi) \circ \varphi^{-1})_x X(\varphi) \circ \varphi^{-1} \right) Z(\varphi) \circ \varphi^{-1} dx \\ & \quad + \int_{\mathbb{S}} A \left((DZ(\varphi) \cdot X(\varphi)) \circ \varphi^{-1} - (Z(\varphi) \circ \varphi^{-1})_x X(\varphi) \circ \varphi^{-1} \right) Y(\varphi) \circ \varphi^{-1} dx, \end{aligned}$$

where formula (3.3) was used to carry out the differentiation. Define $u, v, w \in \mathbf{E}^k$ by $u = X(\varphi) \circ \varphi^{-1}$, $v = Y(\varphi) \circ \varphi^{-1}$, and $w = Z(\varphi) \circ \varphi^{-1}$. We get

$$\begin{aligned} (3.4) \quad & (\mathcal{L}_X \langle Y, Z \rangle)(\varphi) = \int_{\mathbb{S}} A \left((DY(\varphi) \cdot X(\varphi)) \circ \varphi^{-1} \right) w dx \\ & + \int_{\mathbb{S}} A \left((DZ(\varphi) \cdot X(\varphi)) \circ \varphi^{-1} \right) v dx - \int_{\mathbb{S}} A(v_x u) w dx - \int_{\mathbb{S}} A(w_x u) v dx. \end{aligned}$$

On the other hand

$$\langle \nabla_X Y, Z \rangle_\varphi = \int_{\mathbb{S}} \left(DY(\varphi) \cdot X(\varphi) - \Gamma(\varphi, Y(\varphi), X(\varphi)) \right) \circ \varphi^{-1} A(Z(\varphi) \circ \varphi^{-1}) dx.$$

Since $\Gamma(\varphi, Y(\varphi), X(\varphi)) = -\frac{1}{2}(A^{-1}D_x(v_x u_x)) \circ \varphi$, we get

$$(3.5) \quad \langle \nabla_X Y, Z \rangle_\varphi = \int_{\mathbb{S}} (DY(\varphi) \cdot X(\varphi)) \circ \varphi^{-1} A(w) dx + \frac{1}{2} \int_{\mathbb{S}} (v_x u_x)_x w dx.$$

Now, recalling that $A = -D_x^2$, it is easy to check that

$$- \int_{\mathbb{S}} A(v_x u) w dx - \int_{\mathbb{S}} A(w_x u) v dx = \frac{1}{2} \int_{\mathbb{S}} (v_x u_x)_x w dx + \frac{1}{2} \int_{\mathbb{S}} (w_x u_x)_x v dx$$

so by (3.4) and (3.5) we obtain

$$(\mathcal{L}_X \langle Y, Z \rangle)(\varphi) = \langle \nabla_X Y, Z \rangle_\varphi + \langle Y, \nabla_X Z \rangle_\varphi.$$

This proves that ∇ also satisfies (d). \square

4. Geodesics and the Hunter–Saxton equation. Since the map Γ defined in (2.5) is a smooth Christoffel map (see Proposition 2.2), all the usual constructions for affine connections on Banach manifolds (cf. [10]) can easily be carried out on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$. Here we will be concerned with the geodesics on $Rot(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ and their relation to the Hunter–Saxton equation.

4.1. Parallel translation. Let $J \subset \mathbb{R}$ be an open interval and let $\varphi : J \rightarrow M^k$ be a C^2 -curve. A lift $V : J \rightarrow TM^k$ of φ is φ -parallel if

$$V_t = \Gamma(\varphi, V, \varphi_t), \quad t \in J,$$

which is equivalent to $\nabla_{\varphi_t} V \equiv 0$. Applying the general theory for Banach manifolds, we get the following result.

PROPOSITION 4.1. *Let $t_0 \in J$. Given $V_0 \in T_{\varphi(t_0)}M^k$, there exists a unique φ -parallel lift $t \mapsto V(t; V_0) : J \rightarrow TM^k$ such that $V(t_0; V_0) = V_0$.*

Define two continuous maps $u, v : J \rightarrow T_{id}M^k \simeq \mathbf{E}^k$ by

$$u(t) = T_{\varphi(t)}R_{\varphi(t)^{-1}}(\varphi_t(t)) = \varphi_t(t) \circ \varphi(t)^{-1}$$

and

$$v(t) = T_{\varphi(t)}R_{\varphi(t)^{-1}}(V(t)) = V(t) \circ \varphi(t)^{-1}.$$

Note that u, v are not C^1 -maps as

$$u_t = \varphi_{tt} \circ \varphi^{-1} + \frac{\varphi_{tx} \circ \varphi^{-1}}{\varphi_x \circ \varphi^{-1}}$$

is in general an element in \mathbf{E}^{k-1} but not in \mathbf{E}^k . Nevertheless, it is clear that

$$u, v \in C(J; \mathbf{E}^k) \cap C^1(J; \mathbf{E}^{k-1}).$$

THEOREM 4.2. *Let $\varphi : J \rightarrow M^k$ be a C^2 -curve and $V : J \rightarrow TM^k$ a lift of φ . Define $u, v : J \rightarrow \mathbf{E}^k$ by*

$$v(t) = V(t) \circ \varphi(t)^{-1}, \quad u(t) = \varphi_t(t) \circ \varphi(t)^{-1},$$

so that $u, v \in C(J; \mathbf{E}^k) \cap C^1(J; \mathbf{E}^{k-1})$. The following statements are equivalent:

- (a) V is φ -parallel.
- (b) u and v satisfy, for $t \in J$,

$$(4.1) \quad v_t = \Gamma(id, v, u) - v_x u \quad \text{in } \mathbf{E}^{k-1}.$$

- (c) u and v solve the equation

$$(4.2) \quad v_{txx} = -\frac{3}{2}v_{xx}u_x - \frac{1}{2}v_x u_{xx} - v_{xxx}u, \quad t \in J, x \in \mathbb{S}.$$

Proof. First note that

$$(4.3) \quad v_x u = (V \circ \varphi^{-1})_x \varphi_t \circ \varphi^{-1} = \frac{V_x \circ \varphi^{-1}}{\varphi_x \circ \varphi^{-1}} \varphi_t \circ \varphi^{-1} = -V_x \circ \varphi^{-1} \cdot (\varphi^{-1})_t.$$

Suppose V is φ -parallel. Using (4.3), we compute in \mathbf{E}^{k-1}

$$v_t = V_t \circ \varphi^{-1} + V_x \circ \varphi^{-1} \cdot (\varphi^{-1})_t = \Gamma(\varphi, V, \varphi_t) \circ \varphi^{-1} - v_x u = \Gamma(id, v, u) - v_x u,$$

where we used the right invariance (2.6) of Γ . Conversely, if (4.1) holds, then (4.3) yields

$$V_t \circ \varphi^{-1} = v_t - v_x u - V_x \circ \varphi^{-1} \cdot (\varphi^{-1})_t = \Gamma(id, v, u) - v_x u - V_x \circ \varphi^{-1} \cdot (\varphi^{-1})_t = \Gamma(\varphi, V, \varphi_t) \circ \varphi^{-1},$$

showing that V is φ -parallel. This establishes the equivalence of (a) and (b).

Suppose (b) holds. By the definition (2.5) of Γ we rewrite (4.1) as

$$(4.4) \quad v_t = -\frac{1}{2}A^{-1}D_x(v_x u_x) - v_x u.$$

Applying $A = -D_x^2$ to both sides of (4.4) gives (4.2). Conversely, suppose (c) holds. Since at each fixed time both sides of (4.2) belong to \mathbf{F}^{k-3} , we may apply the isomorphism $A^{-1} : \mathbf{F}^{k-3} \rightarrow \mathbf{E}^{k-1}$ to obtain (4.4). Hence (b) and (c) are equivalent. \square

4.2. Geodesics. A C^2 -map $\varphi : J \rightarrow M^k$ is a geodesic if

$$\varphi_{tt} = \Gamma(\varphi, \varphi_t, \varphi_t).$$

Just like for parallel translation we can express the geodesic equation as an equation for $u = \varphi_t \circ \varphi^{-1}$.

THEOREM 4.3. *Let $\varphi : J \rightarrow M^k$ be a C^2 -curve and define $u : J \rightarrow \mathbf{E}^k$ by $u(t) = \varphi_t(t) \circ \varphi(t)^{-1}$ so that $u \in C(J; \mathbf{E}^k) \cap C^1(J; \mathbf{E}^{k-1})$. Then φ is a geodesic if and only if u solves the Hunter-Saxton equation*

$$(4.5) \quad u_{txx} = -2u_x u_{xx} - uu_{xxx}, \quad t \in J, x \in \mathbb{S}.$$

Proof. φ is a geodesic if and only if φ_t is φ -parallel. The equivalence of (a) and (c) of Theorem 4.2 shows that φ_t is φ -parallel if and only if u satisfies (4.5). \square

The condition that $\varphi : J \rightarrow M^k \subset \mathcal{D}^k(\mathbb{S})$ imposes the condition $\varphi(t)|_{x=0} = 0$ for all t . We now remove this restriction.

Let $\rho : \mathcal{D}^k(\mathbb{S}) \rightarrow \text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ be the quotient map. Assume $\varphi : J \rightarrow \mathcal{D}^k(\mathbb{S})$ is a C^1 -curve and denote by $[\varphi] = \rho \circ \varphi : J \rightarrow \text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ the induced curve in $\text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$. A lift $V : J \rightarrow T\mathcal{D}^k(\mathbb{S})$ of φ gives rise to a lift $[V] : J \rightarrow T(\text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S}))$ of $[\varphi]$ given by $[V] = T\rho \circ V$. The next result, which is a fairly straightforward consequence of Theorem 4.2, characterizes the $[\varphi]$ -parallel lifts $[V]$ of $[\varphi]$.

PROPOSITION 4.4. *Let $\varphi : J \rightarrow \mathcal{D}^k(\mathbb{S})$ be a C^2 -curve and let V be a C^1 -lift of φ . Define $u, v \in C(J; H^k(\mathbb{S})) \cap C^1(J; H^{k-1}(\mathbb{S}))$ by*

$$u(t) = \varphi_t(t) \circ \varphi(t)^{-1}, \quad v(t) = V(t) \circ \varphi(t)^{-1}.$$

Then $[V] = T\rho \circ V$ is $[\varphi]$ -parallel if and only if u, v satisfy

$$(4.6) \quad v_{txx} = -\frac{3}{2}v_{xx}u_x - \frac{1}{2}v_x u_{xx} - v_{xxx}u, \quad t \in J, x \in \mathbb{S}.$$

The final theorem in this section says that $t \mapsto [\varphi(t)] : J \rightarrow \text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ is a geodesic if and only if $u = \varphi_t \circ \varphi^{-1}$ solves the Hunter-Saxton equation.

THEOREM 4.5. *Let $\varphi : J \rightarrow \mathcal{D}^k(\mathbb{S})$ be a C^2 -curve and define $u : J \rightarrow H^k(\mathbb{S})$ by $u(t) = \varphi_t(t) \circ \varphi(t)^{-1}$ so that $u \in C(J; H^k(\mathbb{S})) \cap C^1(J; H^{k-1}(\mathbb{S}))$. Then $[\varphi]$ is a geodesic in $\text{Rot}(\mathbb{S}) \setminus \mathcal{D}^k(\mathbb{S})$ if and only if u is a solution of the Hunter-Saxton equation*

$$(4.7) \quad u_{txx} = -2u_x u_{xx} - uu_{xxx}, \quad t \in J, x \in \mathbb{S}.$$

Proof. $[\varphi]$ is a geodesic if and only if $[\varphi]_t = [\varphi_t]$ is $[\varphi]$ -parallel. By Proposition 4.4 this occurs if and only if (4.6) holds with $u = v = \varphi_t \circ \varphi^{-1}$. \square

5. Explicit formulas. In this section the method of characteristics is adopted to obtain explicit formulas for the geodesics; as a byproduct new explicit solutions of the periodic Hunter–Saxton equation drop out.

DEFINITION 5.1. *By a solution of the Hunter–Saxton equation with initial data $u_0 \in H^k(\mathbb{S})$, $k \geq 3$, we mean a function $u(t, x)$ with $u \in C([0, T]; H^3(\mathbb{S})) \cap C^1([0, T]; H^2(\mathbb{S}))$ for some maximal time of existence $T > 0$ such that $u(0) = u_0$ and*

$$(5.1) \quad u_{txx} = -2u_x u_{xx} - uu_{xxx}, \quad t \in [0, T), \quad x \in \mathbb{S}.$$

The next lemma is straightforward to prove.

LEMMA 5.2. *Let $J \subset \mathbb{R}$ be an open interval and let $u \in C(J; H^3(\mathbb{S})) \cap C^1(J; H^2(\mathbb{S}))$. The map*

$$F : (t, \psi) \mapsto u(t) \circ \psi : J \times \mathcal{D}^2(\mathbb{S}) \rightarrow H^2(\mathbb{S})$$

is C^1 .

Now let $u_0 \in H^k(\mathbb{S})$ and suppose $u(t, x)$ is a solution of the Hunter–Saxton equation with initial data u_0 . Integration of (5.1) from 0 to x for each fixed time t gives

$$(5.2) \quad u_{tx} = -\frac{1}{2}u_x^2 - uu_{xx} + c(t), \quad t \in [0, T), \quad x \in \mathbb{S},$$

for some function $c : [0, T) \rightarrow \mathbb{R}$. Since $\int_{\mathbb{S}} u_{tx}(t, x) dx = 0$ for each $t \in [0, T)$, we get

$$\int_{\mathbb{S}} \left(-\frac{1}{2}u_x^2 - uu_{xx} + c(t) \right) dx = 0.$$

An integration by parts yields

$$c(t) = -\frac{1}{2} \int_{\mathbb{S}} u_x^2(t, x) dx.$$

Moreover, by (5.2),

$$\frac{1}{2} \frac{d}{dt} \int_{\mathbb{S}} u_x^2 dx = \int_{\mathbb{S}} u_x u_{tx} dx = - \int_{\mathbb{S}} \left(\frac{1}{2}u_x^3 + uu_x u_{xx} \right) dx = 0,$$

showing that $c(t) = -\frac{1}{2} \int_{\mathbb{S}} u_x^2(t, x) dx$ is a constant function. Assuming that u_0 is nontrivial, we may rescale u_0 so that $\int_{\mathbb{S}} u_{0x}^2 dx = 4$. Then (5.2) becomes

$$u_{tx} = -\frac{1}{2}u_x^2 - uu_{xx} - 2, \quad t \in [0, T), \quad x \in \mathbb{S}.$$

By Lemma 5.2 and the local existence and uniqueness theorem for differential equations in Banach spaces, there exists a unique map $\varphi \in C^1([0, T_1]; \mathcal{D}^2(\mathbb{S}))$ such that $\varphi(0) = id$ and

$$(5.3) \quad \varphi_t(t) = u(t) \circ \varphi(t), \quad t \in [0, T_1),$$

for some maximal existence time $T_1 > 0$. Since

$$(u_x \circ \varphi)_t = u_{tx} \circ \varphi + u_{xx} \circ \varphi \cdot \varphi_t = (u_{tx} + uu_{xx}) \circ \varphi,$$

this gives

$$(u_x \circ \varphi)_t = \left(-\frac{1}{2}u_x^2 - 2\right) \circ \varphi = -\frac{1}{2}(u_x \circ \varphi)^2 - 2, \quad t \in [0, T_1), x \in \mathbb{S}.$$

Therefore, for a fixed $x \in \mathbb{S}$, the function $t \mapsto (u_x \circ \varphi)(t, x)$ solves the differential equation

$$(5.4) \quad \dot{z}(t) = -\frac{1}{2}(z(t)^2 + 4).$$

The general solution of (5.4) is

$$z(t) = 2 \tan \left(\arctan \left(\frac{z(0)}{2} \right) - t \right).$$

As $(u_x \circ \varphi)(0, x) = u_{0x}(x)$ for $x \in \mathbb{S}$, we infer that

$$(5.5) \quad (u_x \circ \varphi)(t, x) = 2 \tan \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - t \right), \quad t \in [0, T_1), x \in \mathbb{S}.$$

Differentiation of (5.3) yields $\varphi_{tx} = u_x \circ \varphi \cdot \varphi_x$, so that, by (5.5),

$$\varphi_{tx}(t, x) = 2 \tan \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - t \right) \varphi_x(t, x).$$

Using the fact that $\varphi_x(0, x) = 1$ for $x \in \mathbb{S}$, we get

$$(5.6) \quad \varphi_x(t, x) = \exp \left(2 \int_0^t \tan \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - s \right) ds \right), \quad t \in [0, T_1), x \in \mathbb{S}.$$

Since

$$2 \int_0^t \tan \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - s \right) ds = \ln \left(\cos^2 \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - t \right) \right) - C(x),$$

where

$$C(x) = \ln \left(\cos^2 \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) \right) \right),$$

equation (5.6) yields

$$(5.7) \quad \varphi_x(t, x) = \frac{\cos^2 \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) - t \right)}{\cos^2 \left(\arctan \left(\frac{u_{0x}(x)}{2} \right) \right)}.$$

We rewrite this as

$$(5.8) \quad \varphi_x(t, x) = \left(\cos t + \frac{u_{0x}(x)}{2} \sin t \right)^2, \quad t \in [0, T_1), x \in \mathbb{S}.$$

Hence

$$\varphi(t, x) - \varphi(t, 0) = x \frac{1 + \cos 2t}{2} + \frac{1}{4} \int_0^x u_{0x}^2(y) dy \frac{1 - \cos 2t}{2} + \frac{u_0(x) - u_0(0)}{2} \sin 2t.$$

Since, by (2.4),

$$-(A^{-1}D_x(u_{0x}^2))(x) = \int_0^x u_{0x}^2(y)dy - x \int_{\mathbb{S}} u_{0x}^2 dx = \int_0^x u_{0x}^2(y)dy - 4x,$$

we infer that, for $t \in [0, T_1)$ and $x \in \mathbb{S}$,

$$(5.9) \quad \varphi(t, x) - \varphi(t, 0) = x - \frac{1}{8}(A^{-1}D_x(u_{0x}^2))(x)(1 - \cos 2t) + \frac{u_0(x) - u_0(0)}{2} \sin 2t.$$

The right-hand side of this equation is well-defined for all times $t < T^*$, where

$$T^* = \frac{\pi}{2} + \arctan\left(\frac{1}{2} \min_{x \in \mathbb{S}} u_{0x}(x)\right).$$

Indeed, it follows from (5.8) that T^* is the first time for which the right-hand side of (5.9) ceases to be a bijective map $\mathbb{S} \rightarrow \mathbb{S}$.

The invariance of (5.9) under the transformation $\varphi(t, x) \mapsto \varphi(t, x) + c(t)$ for an arbitrary function $c(t)$ corresponds to the invariance of the Hunter–Saxton equation under the transformation $u(t) \mapsto u(t, \cdot - c(t)) + c'(t)$. It follows that the geodesic $\varphi(t, x)$ can be extended to all times $t < T^*$ unless $\lim_{t \uparrow T} |\varphi(t, 0)| = \infty$ for some $T < T^*$. This artificial blow-up at some $T < T^*$ can be avoided by requiring that $u_0(0) = 0$ and that φ be a curve in M^k ; this fixes $\varphi(t, 0) = 0$ at all times. We summarize what we have obtained in a theorem.

THEOREM 5.3. *Suppose $u_0 \in H^k(\mathbb{S})$ satisfies $\langle u_0, u_0 \rangle_{id} = \int_{\mathbb{S}} u_{0x}^2 dx = 4$ and $u_0(0) = 0$. Let $\varphi : [0, T^*) \rightarrow M^k$ be the unique geodesic with $\varphi(0) = id$ and $\varphi_t(0) = u_0$ existing for some maximal time $T^* > 0$.*

Then

$$(5.10) \quad \varphi(t) = id - \frac{1}{8}(A^{-1}D_x(u_{0x}^2))(1 - \cos 2t) + \frac{u_0}{2} \sin 2t,$$

and T^ is the first time for which $\varphi(t) : \mathbb{S} \rightarrow \mathbb{S}$ ceases to be bijective given by*

$$T^* = \frac{\pi}{2} + \arctan\left(\frac{1}{2} \min_{x \in \mathbb{S}} u_{0x}(x)\right) < \frac{\pi}{2}.$$

Moreover, defining $u \in C([0, T^]; \mathbf{E}^k) \cap C^1([0, T^*]; \mathbf{E}^{k-1})$ by*

$$u(t) = \left(-\frac{1}{4}(A^{-1}D_x(u_{0x}^2)) \sin 2t + u_0 \cos 2t\right) \circ \varphi(t)^{-1},$$

the set of solutions of the Hunter–Saxton equation with initial data u_0 (see Definition 5.1) is exactly the set of maps

$$\{t \mapsto u(t, \cdot - c(t)) + c'(t)\} \subset C([0, T]; H^k(\mathbb{S})) \cap C^1([0, T]; H^{k-1}(\mathbb{S})),$$

where $T \leq T^$ is the maximal time of existence, $c : [0, T) \rightarrow \mathbb{R}$ is an arbitrary C^1 -function with $c(0) = c'(0) = 0$, and, if $T < T^*$, then $|c(t)| \rightarrow \infty$ as $t \uparrow T < T^*$.*

Remark. For nonconstant initial data u_0 the assumption $\int_{\mathbb{S}} u_{0x}^2 dx = 4$ is just a matter of scaling; it sets the speed of the geodesic to be 4. The case of a constant u_0 is degenerate as it means that the initial velocity of the geodesic vanishes. For the sake of completeness we note that if u_0 is constant the general solution of (1.1) is $u(t, x) \equiv \bar{u}(t)$, where $\bar{u} : [0, T) \rightarrow \mathbb{R}$ is a C^1 -function of time with $\bar{u}(0) = u_0$.

Acknowledgments. The research presented in this paper was carried out while the author was a visitor at the Mittag-Leffler Institute in Stockholm, Sweden. The author also thanks the two referees for helpful comments on a first version of the manuscript.

REFERENCES

- [1] A. BRESSAN AND A. CONSTANTIN, *Global solutions of the Hunter-Saxton equation*, SIAM J. Math. Anal., 37 (2005), pp. 996–1026.
- [2] A. CONSTANTIN, *Existence of permanent and breaking waves for a shallow water equation: A geometric approach*, Ann. Inst. Fourier (Grenoble), 50 (2000), pp. 321–362.
- [3] A. CONSTANTIN AND B. KOLEV, *On the geometric approach to the motion of inertial mechanical systems*, J. Phys. A, 35 (2002), pp. R51–R79.
- [4] D. EBIN AND J. E. MARSDEN, *Groups of diffeomorphisms and the motion of an incompressible fluid*, Ann. of Math. (2), 92 (1970), pp. 102–163.
- [5] R. HAMILTON, *The inverse function theorem of Nash and Moser*, Bull. Amer. Math. Soc., 7 (1982), pp. 65–222.
- [6] J. K. HUNTER AND R. SAXTON, *Dynamics of director fields*, SIAM J. Appl. Math., 51 (1991), pp. 1498–1521.
- [7] J. K. HUNTER AND R. SAXTON, *On a completely integrable nonlinear hyperbolic variational equation*, Phys. D, 79 (1994), pp. 361–386.
- [8] J. K. HUNTER AND Y. ZHENG, *On a nonlinear hyperbolic variational equation I. Global existence of weak solutions*, Arch. Rational Mech. Anal., 129 (1995), pp. 305–353.
- [9] B. KHESIN AND G. MISIOLEK, *Euler equations on homogeneous spaces and Virasoro orbits*, Adv. Math., 176 (2003), pp. 116–144.
- [10] S. LANG, *Differential and Riemannian Manifolds*, 3rd ed., Springer-Verlag, New York, 1995.
- [11] J. LENELLS, *Weak geodesic flow and global solutions of the Hunter-Saxton equation*, Discrete Contin. Dyn. Syst., 18 (2007), pp. 643–656.
- [12] J. LENELLS, *Riemannian geometry on the diffeomorphism group of the circle*, Ark. Mat., 45 (2007), pp. 297–325.
- [13] G. MISIOLEK, *A shallow water equation as a geodesic flow on the Bott-Virasoro group*, J. Geom. Phys., 24 (1998), pp. 203–208.
- [14] H. OMORI, *Infinite-Dimensional Lie Groups*, Transl. Math. Monogr. 158, AMS, Providence, RI, 1997.
- [15] V. OVSIENKO AND B. KHESIN, *The (super) KdV equation as an Euler equation*, Funct. Anal. Appl., 21(1987), pp. 81–82.

ON THE SCHRÖDINGER EQUATION WITH DISSIPATIVE NONLINEARITIES OF DERIVATIVE TYPE*

NAKAO HAYASHI[†], PAVEL I. NAUMKIN[‡], AND HIDEAKI SUNAGAWA[†]

Abstract. We consider the cubic nonlinear Schrödinger equations of derivative type with small initial data. We present a structural condition on the nonlinear terms under which the corresponding Cauchy problem has a dissipative nature.

Key words. Schrödinger equation, dissipative nonlinearity, derivative coupling

AMS subject classifications. 35Q55, 35B40

DOI. 10.1137/070689103

1. Introduction. We consider the initial value problem for the nonlinear Schrödinger equation of derivative type:

$$(1.1) \quad \begin{cases} i\partial_t u + \frac{1}{2}\partial_x^2 u = N(u, \partial_x u), & t > 0, x \in \mathbf{R}, \\ u(0, x) = u_0(x), & x \in \mathbf{R}, \end{cases}$$

where $i = \sqrt{-1}$, $\partial_t = \partial/\partial t$, $\partial_x = \partial/\partial x$, and u is a complex-valued unknown function. We will occasionally write u_x for $\partial_x u$, and \bar{u} denotes the complex conjugate of u . The nonlinear term $N(u, u_x)$ is a cubic homogeneous polynomial in $(u, \bar{u}, u_x, \bar{u}_x)$ with complex coefficients, and it satisfies so-called gauge invariance, that is,

$$(1.2) \quad N(e^{i\theta}v, e^{i\theta}q) = e^{i\theta}N(v, q) \quad \text{for } v, q \in \mathbf{C} \text{ and } \theta \in \mathbf{R}.$$

The aim of this paper is to present a structural condition on the nonlinear term N under which the corresponding forward Cauchy problem (1.1) has a dissipative nature. To explain the motivation, let us begin with the simplest case where N is independent of u_x , i.e., $N = \lambda|u|^2u$ with $\lambda \in \mathbf{C}$. Then it is easy to see that

$$(1.3) \quad \|u(t)\|_{L^2}^2 - 2\operatorname{Im} \lambda \int_0^t \|u(\tau)\|_{L^4}^4 d\tau = \|u_0\|_{L^2}^2,$$

which suggests a dissipative structure if $\operatorname{Im} \lambda < 0$. In fact, it is proved in [18] that the solution decays like $O((t \log t)^{-1/2})$ in L_x^∞ as $t \rightarrow +\infty$ if $\operatorname{Im} \lambda < 0$ and u_0 is small enough. Since the nontrivial free solution (i.e., the solution to (1.1) for $N \equiv 0$, $u_0 \neq 0$) only decays like $O(t^{-1/2})$, this gain of additional logarithmic time decay reflects a dissipative character. Now we turn our attention to the general gauge-invariant cubic nonlinear terms involving both u and u_x . Note that we cannot expect the conservation law like (1.3) anymore. However, as we shall show below, similar

*Received by the editors April 23, 2007; accepted for publication (in revised form) September 14, 2007; published electronically April 16, 2008.

<http://www.siam.org/journals/sima/40-1/68910.html>

[†]Department of Mathematics, Graduate School of Science, Osaka University, Toyonaka, Osaka 560-0043, Japan (nhayashi@math.wani.osaka-u.ac.jp, sunagawa@math.sci.osaka-u.ac.jp). The third author was partially supported by MEXT through a Grant-in-Aid for Young Scientists (B) (18740066).

[‡]Instituto de Matemáticas, Universidad Nacional Autónoma de México, Campus Morelia, AP 61-3 (Xangari), Morelia CP 58089, Michoacán, Mexico (pavelni@matmor.unam.mx). The research of this author is partially supported by CONACYT.

time decay is still valid if $\sup_{\xi \in \mathbf{R}} \operatorname{Im} N(1, i\xi) < 0$. It gives an answer to one of the questions left unsolved in the previous work [20].

Before stating our result, we introduce function spaces. For $s, \nu \geq 0$, let $H^{s, \nu}$ be the weighted Sobolev space given by $\{\phi \in L^2 : \|\phi\|_{H^{s, \nu}} = \|\langle x \rangle^\nu \langle i\partial_x \rangle^s \phi\|_{L^2} < \infty\}$, where $\langle x \rangle = (1 + x^2)^{1/2}$. In particular we set $H^s = H^{s, 0}$, which is typically an L^2 -based Sobolev space of order s . The main result is as follows.

THEOREM 1.1. *Suppose that N satisfies*

$$(1.4) \quad \operatorname{Im} N(1, i\xi) \leq 0 \quad \text{for } \xi \in \mathbf{R}.$$

Let $u_0 \in H^{2,1} \cap H^3$, and $\varepsilon = \|u_0\|_{H^{2,1}} + \|u_0\|_{H^3}$ is sufficiently small. Then (1.1) admits a unique global solution $u \in C([0, \infty); H^{2,1} \cap H^3)$. Moreover, the following asymptotic expression is valid as $t \rightarrow +\infty$ uniformly in $x \in \mathbf{R}$:

$$(1.5) \quad u(t, x) = \frac{a(\frac{x}{t}) \exp\left\{i\frac{x^2}{2t} - i|a(\frac{x}{t})|^2 \operatorname{Re} N(1, i\frac{x}{t}) \int_0^{\log t} \frac{d\sigma}{1 - 2\operatorname{Im} N(1, i\frac{x}{t})(|a(\frac{x}{t})|^2 \sigma + b(\frac{x}{t}))}\right\}}{\sqrt{t} \sqrt{1 - 2\operatorname{Im} N(1, i\frac{x}{t})(|a(\frac{x}{t})|^2 \log t + b(\frac{x}{t}))}} + O(t^{-3/4+\mu}),$$

where $\mu > 0$ is an arbitrarily small constant, and $a(\xi), b(\xi)$ are complex-valued continuous functions of $\xi \in \mathbf{R}$ which satisfy $|a(\xi)| \leq C\varepsilon \langle \xi \rangle^{-2}$, $|b(\xi)| \leq C\varepsilon^4 \langle \xi \rangle^{-4}$ with some positive constant C .

Remark 1.1. As an immediate consequence of (1.5), we can see that

$$\|u(t)\|_{L^\infty} = O((t \log t)^{-1/2}) \quad \text{as } t \rightarrow +\infty$$

if $\sup_{\xi \in \mathbf{R}} \operatorname{Im} N(1, i\xi) < 0$. Moreover, as we shall see in Remark 4.1 below, the L^2 norm of $u(t)$ also decays like $O((\log t)^{-1/2})$. Therefore, by interpolation, we obtain

$$\|u(t)\|_{L^p} = O(t^{-(1/2-1/p)} (\log t)^{-1/2}) \quad \text{as } t \rightarrow +\infty$$

for all $p \in [2, \infty]$.

Remark 1.2. Our result can be also viewed as an extension of [9] (see also [3], [2], [4], [10], [13], [15], [16], [17], [18], and [21]) because (1.5) is reduced to

$$u(t, x) = \frac{1}{\sqrt{t}} a(x/t) e^{i(x^2/(2t) - |a(x/t)|^2 \operatorname{Re} N(1, ix/t) \log t)} + O(t^{-3/4+\mu}) \quad \text{as } t \rightarrow +\infty$$

when $\operatorname{Im} N(1, i\xi) \equiv 0$. An example of the nonlinearity which was not considered previously but satisfies (1.4) is $N = \alpha \bar{u} u_x^2$ with $\operatorname{Im} \alpha > 0$.

Remark 1.3. When $\operatorname{Im} N(1, i\xi_0) > 0$ for some $\xi_0 \in \mathbf{R}$, the authors do not know any global existence or nonexistence results for (1.1). However, in view of the denominator of the leading term of (1.5), it is quite reasonable to expect that small amplitude solutions can blow up in finite time if condition (1.4) is violated. Concerning the lifespan T_ε of the solution for (1.1) with $u_0(x) = \varepsilon \varphi(x)$, the following explicit lower bound is obtained in [20]:

$$\liminf_{\varepsilon \downarrow 0} \varepsilon^2 \log T_\varepsilon \geq \frac{1}{\sup_{\xi \in \mathbf{R}} (2|\hat{\varphi}(\xi)|^2 \operatorname{Im} N(1, i\xi))},$$

where $\hat{\varphi}$ denotes the Fourier transform of φ . It may be an interesting open problem to consider whether the corresponding upper estimate holds or not.

Remark 1.4. When we put

$$A(s, \xi) = \frac{a(\xi) \exp\left\{-i|a(\xi)|^2 \operatorname{Re} N(1, i\xi) \int_0^s \frac{d\sigma}{1-2 \operatorname{Im} N(1, i\xi)(|a(\xi)|^2 \sigma + b(\xi))}\right\}}{\sqrt{1-2 \operatorname{Im} N(1, i\xi)(|a(\xi)|^2 s + b(\xi))}},$$

the asymptotic expression (1.5) can be interpreted as

$$u(t, x) \stackrel{t \rightarrow \infty}{\sim} \frac{e^{ix^2/(2t)}}{\sqrt{t}} A\left(\log t, \frac{x}{t}\right); \quad i\partial_s A = N(1, i\xi)|A|^2 A.$$

This tells us that asymptotic behavior of the solution for (1.1) is characterized by that of the solution for the simpler ordinary differential equation. At the level of the reduced equation, the dissipative nature is transparent via the identity

$$\partial_s (|A(s, \xi)|^2) = 2 \operatorname{Im} N(1, i\xi) |A(s, \xi)|^4.$$

Note that analogous results have been obtained in [19] for nonlinear Klein–Gordon equations and in [14] for nonlinear wave equations.

The rest of this paper is organized as follows: In the next section, we recall several useful identities and inequalities. Section 3 is devoted to getting an a priori estimate, from which global existence follows immediately. After that, we will derive the large time asymptotics of $u(t)$ in section 4 by applying a lemma on ordinary differential equations.

2. Preliminaries. We collect here several identities and inequalities which are useful in the proof of our main theorem. In what follows, we will denote several positive constants by the same letter, C , which is possibly different line by line.

First we put $\mathcal{J} = \mathcal{J}(t) = x + it\partial_x$. It is well known that

$$(2.1) \quad [\partial_x, \mathcal{J}] = 1 \quad \text{and} \quad \left[i\partial_t + \frac{1}{2}\partial_x^2, \mathcal{J} \right] = 0,$$

where $[\cdot, \cdot]$ denotes the commutator, i.e., $[\mathcal{A}, \mathcal{B}] = \mathcal{A}\mathcal{B} - \mathcal{B}\mathcal{A}$ for linear operators \mathcal{A} and \mathcal{B} . Also we have

$$(2.2) \quad \partial_x(\phi\bar{\psi}) = \frac{1}{it} \left\{ (\mathcal{J}\phi)\bar{\psi} - \phi(\overline{\mathcal{J}\psi}) \right\},$$

$$(2.3) \quad \mathcal{J}(\theta\phi\bar{\psi}) = (\mathcal{J}\theta)\phi\bar{\psi} + \theta(\mathcal{J}\phi)\bar{\psi} - \theta\phi(\overline{\mathcal{J}\psi})$$

for smooth functions θ , ϕ , and ψ . Next we denote by $\mathcal{U} = \mathcal{U}(t)$ the free Schrödinger evolution group defined by

$$(\mathcal{U}\phi)(x) = \frac{e^{-i\pi/4}}{\sqrt{2\pi t}} \int_{\mathbf{R}} e^{i(x-y)^2/(2t)} \phi(y) dy.$$

From the relation $\mathcal{J} = \mathcal{U}x\mathcal{U}^{-1}$, we see that

$$(2.4) \quad \|\mathcal{F}\mathcal{U}^{-1}\phi\|_{H_x^1} \leq C\|(1+|x|)\mathcal{U}^{-1}\phi\|_{L_x^2} \leq C(\|\phi\|_{L_x^2} + \|\mathcal{J}\phi\|_{L_x^2}),$$

where

$$(\mathcal{F}\phi)(\xi) = \hat{\phi}(\xi) = \frac{1}{\sqrt{2\pi}} \int_{\mathbf{R}} e^{-iy\xi} \phi(y) dy.$$

It is also well known that \mathcal{U} is decomposed into \mathcal{MDFM} , where

$$\begin{aligned} (\mathcal{M}\phi)(x) &= e^{ix^2/(2t)}\phi(x), \\ (\mathcal{D}\phi)(x) &= \frac{e^{-i\pi/4}}{\sqrt{t}}\phi\left(\frac{x}{t}\right). \end{aligned}$$

With this notation, we set $\mathcal{W} = \mathcal{FMF}^{-1}$ so that $\mathcal{U} = \mathcal{MDWF}$. It follows from the inequalities $\|\phi\|_{L^\infty} \leq C\|\phi\|_{L^2}^{1/2}\|\partial_x\phi\|_{L^2}^{1/2}$ and $|e^{ix^2/(2t)} - 1| \leq Ct^{-1/2}|x|$ that

$$(2.5) \quad \|(\mathcal{W} - 1)\phi\|_{L^\infty} \leq Ct^{-1/4}\|\phi\|_{H^1}, \quad \|(\mathcal{W}^{-1} - 1)\phi\|_{L^\infty} \leq Ct^{-1/4}\|\phi\|_{H^1}.$$

Consequently we have

$$\begin{aligned} (2.6) \quad \|\phi - \mathcal{MDFU}^{-1}\phi\|_{L^\infty} &= \|\mathcal{MD}(\mathcal{W} - 1)\mathcal{FU}^{-1}\phi\|_{L^\infty} \\ &= t^{-1/2}\|(\mathcal{W} - 1)\mathcal{FU}^{-1}\phi\|_{L^\infty} \\ &\leq Ct^{-3/4}\|\mathcal{FU}^{-1}\phi\|_{H^1} \\ &\leq Ct^{-3/4}(\|\phi\|_{L^2} + \|\mathcal{J}\phi\|_{L^2}) \end{aligned}$$

and

$$(2.7) \quad \begin{aligned} \|\phi\|_{L^\infty} &\leq \|\mathcal{MDFU}^{-1}\phi\|_{L^\infty} + \|\phi - \mathcal{MDFU}^{-1}\phi\|_{L^\infty} \\ &\leq t^{-1/2}\|\mathcal{FU}^{-1}\phi\|_{L^\infty} + Ct^{-3/4}(\|\phi\|_{L^2} + \|\mathcal{J}\phi\|_{L^2}). \end{aligned}$$

Also it follows from the inequality $\|\mathcal{W}^{-1}\phi\|_{L^\infty} \leq Ct^{1/2}\|\phi\|_{L^1}$ that

$$\begin{aligned} (2.8) \quad \|\mathcal{FU}^{-1}(\theta\phi\bar{\psi})\|_{L^\infty} &= \frac{1}{t}\left\|\mathcal{W}^{-1}\left\{(\mathcal{W}\mathcal{FU}^{-1}\theta)(\mathcal{W}\mathcal{FU}^{-1}\phi)\overline{(\mathcal{W}\mathcal{FU}^{-1}\psi)}\right\}\right\|_{L^\infty} \\ &\leq \frac{C}{t^{1/2}}\|(\mathcal{W}\mathcal{FU}^{-1}\theta)(\mathcal{W}\mathcal{FU}^{-1}\phi)(\mathcal{W}\mathcal{FU}^{-1}\psi)\|_{L^1} \\ &\leq \frac{C}{t^{1/2}}\|\theta\|_{L^2}\|\phi\|_{L^2}\|\mathcal{W}\mathcal{FU}^{-1}\psi\|_{L^\infty} \\ &\leq C\|\theta\|_{L^2}\|\phi\|_{L^2}\left\{\|\psi\|_{L^\infty} + t^{-3/4}(\|\psi\|_{L^2} + \|\mathcal{J}\psi\|_{L^2})\right\}. \end{aligned}$$

3. A priori estimates. This section is devoted to getting suitable a priori estimate for the solution of (1.1). Since the local existence is well known (see, e.g., [1], [11], [13]), we can deduce global existence from this estimate. From now on, let $u(t)$ be the solution of (1.1) for $t \in [0, T]$ with some $T > 0$ and define

$$E_\delta(T) = \sup_{0 \leq t \leq T} \left\{ (1+t)^{-\delta} \sum_{j=0}^1 \|\mathcal{J}^j u(t)\|_{H^{3-j}} + (1+t)^{1/2} \|u(t)\|_{W^{2,\infty}} \right\}$$

with $\delta \in (0, 1/8]$ fixed. Here and later on as well, $W^{k,p}$ denotes an L^p -based Sobolev space of order k . What we are going to show is the following.

LEMMA 3.1. *Let $K \geq 1$ and assume $E_\delta(T) \leq K\varepsilon$, where $\varepsilon = \|u_0\|_{H^{2,1}} + \|u_0\|_{H^3}$. Then we have*

$$E_\delta(T) \leq C_1 e^{C_2 K^2 \varepsilon^2} (1 + K^3 \varepsilon^2 + K^5 \varepsilon^4) \varepsilon,$$

where C_1, C_2 are positive constants independent of ε, K , and T , but possibly dependent on δ .

Once this lemma is established, global existence follows immediately. Indeed, when we put $\varepsilon_0 = K^{-2}$ and choose K so large that $C_1 e^{C_2 K^{-2}} (1 + K^{-1} + K^{-3}) \leq K/2$, it follows from the above lemma that $E_\delta(T) \leq K\varepsilon$ implies $E_\delta(T) \leq K\varepsilon/2$ for any $\varepsilon \in (0, \varepsilon_0]$. Then, by the continuity argument (see, e.g., [12]), we see that $E_\delta(T) \leq K\varepsilon$ must hold as long as the solution exists. Therefore the local solution can be extended to the global one.

Now, we turn to the proof of Lemma 3.1. It will be divided into two parts, i.e., L^∞ and L^2 estimates. The argument below is a refinement of section 2 of [9].

3.1. L^∞ estimates. In this part, we consider the estimate of $(1+t)^{1/2} \|u(t)\|_{W^{2,\infty}}$. We set

$$\alpha(t, \xi) = \mathcal{F}(\mathcal{U}^{-1}u(t))(\xi)$$

and $\alpha_k(t, \xi) = (i\xi)^k \alpha(t, \xi)$ so that $\partial_x^k u = \mathcal{M}\mathcal{D}\mathcal{W}\alpha_k$. Using the inequality (2.7) for $t \in [1, T]$ and the Sobolev imbedding for $t \in [0, 1]$, we can see that the problem is reduced to getting the bound of $\|\alpha_k(t)\|_{L^\infty}$ for $t \in [1, T]$, $0 \leq k \leq 2$, under the assumption $E_\delta(T) \leq K\varepsilon$. More precisely, our objective here is to obtain the following estimate:

$$\sup_{t \in [1, T]} \|\alpha_k(t)\|_{L^\infty} \leq C(\varepsilon + K^3 \varepsilon^3), \quad k = 0, 1, 2.$$

Let us first consider the case of $k = 0$. Applying the operator $\mathcal{F}\mathcal{U}^{-1}$ to the equation, we have

$$\begin{aligned} i\partial_t \alpha &= \mathcal{F}\mathcal{U}^{-1}N(u, u_x) \\ &= \mathcal{W}^{-1}\mathcal{D}^{-1}\mathcal{M}^{-1}N(\mathcal{M}\mathcal{D}\mathcal{W}\alpha, \mathcal{M}\mathcal{D}\mathcal{W}\alpha_1) \\ &= \frac{1}{t}\mathcal{W}^{-1}N(\mathcal{W}\alpha, \mathcal{W}\alpha_1) \\ &= \frac{1}{t}N(1, i\xi)|\alpha|^2\alpha + \rho_0, \end{aligned}$$

where

$$\rho_0 = \frac{1}{t} \left\{ \mathcal{W}^{-1}N(\mathcal{W}\alpha, \mathcal{W}\alpha_1) - N(\alpha, \alpha_1) \right\}.$$

Observe that (1.4) implies

$$\begin{aligned} \partial_t(|\alpha|^2) &= 2\operatorname{Im}(\bar{\alpha} i\partial_t \alpha) \\ &= \frac{2}{t} \operatorname{Im}N(1, i\xi)|\alpha|^4 + 2\operatorname{Im}(\bar{\alpha}\rho_0) \\ &\leq 2|\alpha||\rho_0| \end{aligned}$$

and that (2.5) yields

$$\begin{aligned} (3.1) \quad \|\rho_0\|_{L^\infty} &\leq t^{-1} \|(\mathcal{W}^{-1} - 1)N(\mathcal{W}\alpha, \mathcal{W}\alpha_1)\|_{L^\infty} \\ &\quad + t^{-1} \|N(\mathcal{W}\alpha, \mathcal{W}\alpha_1) - N(\alpha, \alpha_1)\|_{L^\infty} \\ &\leq CK^3 \varepsilon^3 t^{-5/4}. \end{aligned}$$

So we deduce that

$$\|\alpha(t)\|_{L^\infty} \leq \|\alpha(1)\|_{L^\infty} + \int_1^t \|\rho_0(\tau)\|_{L^\infty} d\tau \leq C\varepsilon + CK^3\varepsilon^3.$$

Next we turn to the case of $k = 1, 2$. We split $N(u, u_x)$ into $g_0(u, u_x)u + g_1(u, u_x)u_x$, where $g_0(u, q) = \lambda_1|u|^2 + \lambda_2u\bar{q} + \lambda_3|q|^2$, $g_1(u, q) = \lambda_4|u|^2 + \lambda_5q\bar{u} + \lambda_6|q|^2$ with complex constants $\lambda_1, \dots, \lambda_6$, and we set

$$h_k = \partial_x^k(N(u, u_x)) - \sum_{j=0}^1 g_j(u, u_x)\partial_x^{j+k}u$$

so that

$$i\partial_t\alpha_k = \mathcal{F}\mathcal{U}^{-1} \sum_{j=0}^1 g_j(u, u_x)\partial_x^{j+k}u + \mathcal{F}\mathcal{U}^{-1}h_k.$$

To find out the leading term of the right-hand side, we analyze the action of $\mathcal{F}\mathcal{U}^{-1}$ to $g_j(u, u_x)\partial_x^{j+k}u$ carefully. For $j = 0$, we have

$$\begin{aligned} \mathcal{F}\mathcal{U}^{-1}(g_0(u, u_x)\partial_x^k u) &= \mathcal{W}^{-1}\mathcal{D}^{-1}\mathcal{M}^{-1}(g_0(\mathcal{M}\mathcal{D}\mathcal{W}\alpha, \mathcal{M}\mathcal{D}\mathcal{W}\alpha_1)\mathcal{M}\mathcal{D}\mathcal{W}\alpha_k) \\ &= \frac{1}{t}\mathcal{W}^{-1}(g_0(\mathcal{W}\alpha, \mathcal{W}\alpha_1)\mathcal{W}\alpha_k) \\ &= \frac{1}{t}g_0(\alpha, \alpha_1)\alpha_k + \sigma_k, \end{aligned}$$

where

$$\sigma_k = \frac{1}{t}\left\{\mathcal{W}^{-1}(g_0(\mathcal{W}\alpha, \mathcal{W}\alpha_1)\mathcal{W}\alpha_k) - g_0(\alpha, \alpha_1)\alpha_k\right\}.$$

For $j = 1$, let us introduce

$$\tilde{h}_k = g_1(u, u_x)\partial_x^{1+k}u + (\lambda_4u\bar{u}_x + \lambda_5|u_x|^2 + \lambda_6u_x\bar{u}_{xx})\partial_x^k u.$$

Then, noting the relation $\lambda_4\alpha\bar{\alpha}_1 + \lambda_5|\alpha_1|^2 + \lambda_6\alpha_1\bar{\alpha}_2 = -i\xi g_1(1, i\xi)|\alpha|^2$, we see that

$$\begin{aligned} &\mathcal{F}\mathcal{U}^{-1}(g_1(u, u_x)\partial_x^{1+k}u) \\ &= -\mathcal{F}\mathcal{U}^{-1}\{(\lambda_4u\bar{u}_x + \lambda_5|u_x|^2 + \lambda_6u_x\bar{u}_{xx})\partial_x^k u\} + \mathcal{F}\mathcal{U}^{-1}\tilde{h}_k \\ &= -\frac{1}{t}\mathcal{W}^{-1}\left\{(\lambda_4(\mathcal{W}\alpha)(\overline{\mathcal{W}\alpha_1}) + \lambda_5|\mathcal{W}\alpha_1|^2 + \lambda_6(\mathcal{W}\alpha_1)(\overline{\mathcal{W}\alpha_2}))\mathcal{W}\alpha_k\right\} + \mathcal{F}\mathcal{U}^{-1}\tilde{h}_k \\ &= -\frac{1}{t}(\lambda_4\alpha\bar{\alpha}_1 + \lambda_5|\alpha_1|^2 + \lambda_6\alpha_1\bar{\alpha}_2)\alpha_k + \tilde{\sigma}_k + \mathcal{F}\mathcal{U}^{-1}\tilde{h}_k \\ &= \frac{1}{t}i\xi g_1(1, i\xi)|\alpha|^2\alpha_k + \tilde{\sigma}_k + \mathcal{F}\mathcal{U}^{-1}\tilde{h}_k, \end{aligned}$$

where

$$\begin{aligned} \tilde{\sigma}_k &= -\frac{\lambda_4}{t}\left\{\mathcal{W}^{-1}((\mathcal{W}\alpha)(\overline{\mathcal{W}\alpha_1})\mathcal{W}\alpha_k) - (\alpha\bar{\alpha}_1)\alpha_k\right\} \\ &\quad -\frac{\lambda_5}{t}\left\{\mathcal{W}^{-1}(|\mathcal{W}\alpha_1|^2\mathcal{W}\alpha_k) - |\alpha_1|^2\alpha_k\right\} \\ &\quad -\frac{\lambda_6}{t}\left\{\mathcal{W}^{-1}((\mathcal{W}\alpha_1)(\overline{\mathcal{W}\alpha_2})\mathcal{W}\alpha_k) - \alpha_1\bar{\alpha}_2\alpha_k\right\}. \end{aligned}$$

Summing up, we have

$$\begin{aligned} i\partial_t\alpha_k &= (g_0(1, i\xi) + i\xi g_1(1, i\xi))|\alpha|^2\alpha_k + \sigma_k + \tilde{\sigma}_k + \mathcal{F}\mathcal{U}^{-1}(h_k + \tilde{h}_k) \\ &= N(1, i\xi)|\alpha|^2\alpha_k + \rho_k, \end{aligned}$$

where $\rho_k = \sigma_k + \tilde{\sigma}_k + \mathcal{F}\mathcal{U}^{-1}(h_k + \tilde{h}_k)$. Thus we deduce as in the previous case that

$$\|\alpha_k(t)\|_{L^\infty} \leq \|\alpha_k(1)\|_{L^\infty} + \int_1^t \|\rho_k(\tau)\|_{L^\infty} d\tau \leq C\varepsilon + CK^3\varepsilon^3,$$

provided that the estimate $\|\rho_k\|_{L^\infty} \leq CK^3\varepsilon^3t^{-5/4}$ is verified. So the remaining task is to check this for $k = 1, 2$. In order to estimate $\|\mathcal{F}\mathcal{U}^{-1}(h_k + \tilde{h}_k)\|_{L^\infty}$, we rewrite h_k as

$$\sum_{j=0}^1 \sum_{l=1}^k \binom{k}{l} \partial_x^{l-1} \left(\partial_x (g_j(u, u_x)) \right) \partial_x^{j+k-l} u$$

and apply (2.2) to $\partial_x(g_j(u, u_x))$. Also we express \tilde{h}_k as

$$\lambda_4 u \partial_x (\bar{u} \partial_x^k u) + \lambda_5 u_x \partial_x (\bar{u} \partial_x^k u) + \lambda_6 u_x \partial_x (\bar{u}_x \partial_x^k u)$$

and use (2.2). Then we see that $h_k + \tilde{h}_k$ can be written as a linear combination of

$$\frac{1}{t} (\mathcal{J} \partial_x^{k_1} u) (\partial_x^{k_2} u) \overline{(\partial_x^{k_3} u)} \quad \text{or} \quad \frac{1}{t} \overline{(\mathcal{J} \partial_x^{k_1} u)} (\partial_x^{k_2} u) (\partial_x^{k_3} u)$$

with $k_1, k_2, k_3 \leq 2$. Hence it follows from (2.8) that

$$(3.2) \quad \|\mathcal{F}\mathcal{U}^{-1}(h_k + \tilde{h}_k)\|_{L^\infty} \leq CK^3\varepsilon^3t^{-3/2+2\delta}.$$

On the other hand, we deduce as the derivation of (3.1) that

$$(3.3) \quad \|\sigma_k + \tilde{\sigma}_k\|_{L^\infty} \leq CK^3\varepsilon^3t^{-5/4}.$$

From (3.2) and (3.3), it follows that

$$(3.4) \quad \|\rho_k\|_{L^\infty} \leq CK^3\varepsilon^3t^{-5/4} + CK^3\varepsilon^3t^{-3/2+2\delta} \leq CK^3\varepsilon^3t^{-5/4}$$

for $k = 1, 2$, as desired.

3.2. L^2 estimates. In the remainder of this section, we consider the bound of $(1+t)^{-\delta} \sum_{j=0}^1 \|\mathcal{J}^j u(t)\|_{H^{3-j}}$. It is enough to show that

$$(3.5) \quad \sum_{j=0}^1 \|\mathcal{J}^j u(t)\|_{L^2} \leq C\varepsilon + CK^3\varepsilon^3(1+t)^\delta$$

and

$$(3.6) \quad \sum_{j=0}^1 \|\partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} \leq Ce^{CK^2\varepsilon^2} (\varepsilon + K^3\varepsilon^3 + K^5\varepsilon^5)(1+t)^\delta$$

for $t \in [0, T]$ under the assumption $E_\delta(T) \leq K\varepsilon$.

First we consider the easier estimate (3.5). It follows from the standard energy method that

$$\frac{d}{dt} \|u(t)\|_{L^2} \leq C \|u\|_{W^{1,\infty}}^2 \|u\|_{H^1} \leq CK^3 \varepsilon^3 (1+t)^{-1+\delta}.$$

Also, when we remember (2.1) and (2.3), we see that

$$\begin{aligned} & \left(i\partial_t + \frac{1}{2}\partial_x^2\right) \mathcal{J}u \\ &= \frac{\partial N}{\partial u}(u, u_x) \mathcal{J}u - \frac{\partial N}{\partial \bar{u}}(u, u_x) \overline{\mathcal{J}u} + \frac{\partial N}{\partial q}(u, u_x) \mathcal{J}\partial_x u - \frac{\partial N}{\partial \bar{q}}(u, u_x) \overline{\mathcal{J}\partial_x u}, \end{aligned}$$

where the variable q is responsible for u_x , and $\partial/\partial u, \partial/\partial \bar{u}, \partial/\partial q, \partial/\partial \bar{q}$ are defined by

$$\begin{aligned} \frac{\partial}{\partial u} &= \frac{1}{2} \left(\frac{\partial}{\partial \operatorname{Re} u} - i \frac{\partial}{\partial \operatorname{Im} u} \right), & \frac{\partial}{\partial \bar{u}} &= \frac{1}{2} \left(\frac{\partial}{\partial \operatorname{Re} u} + i \frac{\partial}{\partial \operatorname{Im} u} \right), \\ \frac{\partial}{\partial q} &= \frac{1}{2} \left(\frac{\partial}{\partial \operatorname{Re} q} - i \frac{\partial}{\partial \operatorname{Im} q} \right), & \frac{\partial}{\partial \bar{q}} &= \frac{1}{2} \left(\frac{\partial}{\partial \operatorname{Re} q} + i \frac{\partial}{\partial \operatorname{Im} q} \right), \end{aligned}$$

respectively. Then the energy method again implies

$$\frac{d}{dt} \|\mathcal{J}u(t)\|_{L^2} \leq C \|u\|_{W^{1,\infty}}^2 (\|u\|_{H^1} + \|\mathcal{J}u\|_{H^1}) \leq CK^3 \varepsilon^3 (1+t)^{-1+\delta},$$

which gives us (3.5).

Next we consider (3.6). To obtain this estimate, we use the gauge transformation technique since the standard energy estimate may cause a derivative loss, as we have just seen. The following energy inequality is established in section 2 of [13] (see also [1] and [11]).

LEMMA 3.2. *Let $\psi(t, x)$ be a smooth function satisfying*

$$i\partial_t \psi + \frac{1}{2}\partial_x^2 \psi + b_1(t, x)\partial_x \psi + b_2(t, x)\partial_x \bar{\psi} = f(t, x)$$

for $(t, x) \in [0, T] \times \mathbf{R}$ with some smooth functions $b_1(t, x), b_2(t, x), f(t, x)$ and some $T > 0$. Then we have

$$\frac{d}{dt} \|e^{P(t,\cdot)} \psi(t, \cdot)\|_{L^2} \leq CB(t) \|e^{P(t,\cdot)} \psi(t, \cdot)\|_{L^2} + \|e^{P(t,\cdot)} f(t, \cdot)\|_{L^2}$$

for $t \in [0, T]$, where

$$(3.7) \quad P(t, x) = \int_{-\infty}^x \operatorname{Re} b_1(t, y) dy,$$

$$(3.8) \quad B(t) = \sum_{k=1}^2 (\|\partial_x b_k(t, \cdot)\|_{L^\infty} + \|b_k(t, \cdot)\|_{L^\infty}^2) + \sup_{x \in \mathbf{R}} \left| \int_{-\infty}^x \partial_t b_1(t, y) dy \right|.$$

Since $\partial_x^{3-j} \mathcal{J}^j u$ satisfies

$$\begin{aligned} & \left(i\partial_t + \frac{1}{2}\partial_x^2\right) (\partial_x^{3-j} \mathcal{J}^j u) \\ &= \frac{\partial N}{\partial q}(u, u_x) \partial_x (\partial_x^{3-j} \mathcal{J}^j u) + (-1)^j \frac{\partial N}{\partial \bar{q}}(u, u_x) \partial_x \overline{(\partial_x^{3-j} \mathcal{J}^j u)} + R_j \end{aligned}$$

with

$$(3.9) \quad \|R_j\|_{L^2} \leq CK^3\varepsilon^3(1+t)^{-1+\delta}$$

for $j = 0, 1$, we can apply Lemma 3.2 to obtain

$$(3.10) \quad \frac{d}{dt} \|e^{P(t)} \partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} \leq CB(t) \|e^{P(t)} \partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} + \|e^{P(t)} R_j(t)\|_{L^2},$$

where $P(t, x)$ and $B(t)$ are given by (3.7) and (3.8) with

$$b_1 = -\frac{\partial N}{\partial q}(u, u_x), \quad b_2 = (-1)^{j+1} \frac{\partial N}{\partial \bar{q}}(u, u_x).$$

Note that

$$\begin{aligned} \|u(t)\|_{H_x^1} &= \|\langle \xi \rangle \mathcal{F}U^{-1}u(t)\|_{L_\xi^2} \\ &= \|\langle \xi \rangle^{-1}(\alpha(t) + \xi^2\alpha(t))\|_{L_\xi^2} \\ &\leq \left(\int_{-\infty}^{\infty} \frac{d\xi}{1 + \xi^2} \right)^{1/2} (\|\alpha(t)\|_{L_\xi^\infty} + \|\alpha_2(t)\|_{L_\xi^\infty}) \\ &\leq CK\varepsilon, \end{aligned}$$

whence

$$(3.11) \quad e^{\pm P(t,x)} \leq e^{C\|u(t)\|_{H^1}^2} \leq e^{CK^2\varepsilon^2}.$$

As for the estimate of $B(t)$, we observe that the identity

$$\partial_t(\phi\bar{\psi}) = -\frac{i}{2}\partial_x(\phi\bar{\psi}_x - \phi_x\bar{\psi}) + i\phi \overline{\left(i\partial_t + \frac{1}{2}\partial_x^2\right)\psi} - i\bar{\psi} \left(i\partial_t + \frac{1}{2}\partial_x^2\right)\phi$$

adapted to $b_1 = \sum_{j,k=0}^1 \lambda_{jk}(\partial_x^j u)(\overline{\partial_x^k u})$. Then we have

$$\begin{aligned} (3.12) \quad B(t) &\leq C\|u\|_{W^{2,\infty}}^2 + C\|u\|_{W^{1,\infty}}^4 + C \sum_{j,k=0}^1 \int_{\mathbf{R}} |\partial_x^j u| |\partial_x^k N(u, u_x)| dx \\ &\leq C\|u\|_{W^{2,\infty}}^2 + C\|u\|_{W^{1,\infty}}^4 + C\|u\|_{W^{2,\infty}}^2 \|u\|_{H^1}^2 \\ &\leq C(K^2\varepsilon^2 + K^4\varepsilon^4)(1+t)^{-1}. \end{aligned}$$

By (3.9)–(3.12) we obtain

$$\begin{aligned} \frac{d}{dt} \|e^{P(t)} \partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} &\leq C(K^2\varepsilon^2 + K^4\varepsilon^4)(1+t)^{-1} e^{CK^2\varepsilon^2} \|\partial_x^{3-j} \mathcal{J}^j u(t, \cdot)\|_{L^2} \\ &\quad + Ce^{CK^2\varepsilon^2} K^3\varepsilon^3(1+t)^{-1+\delta} \\ &\leq Ce^{CK^2\varepsilon^2} (K^3\varepsilon^3 + K^5\varepsilon^5)(1+t)^{-1+\delta}, \end{aligned}$$

whence

$$\begin{aligned} \sum_{j=0}^1 \|\partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} &\leq e^{CK^2\varepsilon^2} \sum_{j=0}^1 \|e^{P(t)} \partial_x^{3-j} \mathcal{J}^j u(t)\|_{L^2} \\ &\leq e^{CK^2\varepsilon^2} \left\{ C\varepsilon + Ce^{CK^2\varepsilon^2} (K^3\varepsilon^3 + K^5\varepsilon^5) \int_0^t (1+\tau)^{-1+\delta} d\tau \right\} \\ &\leq Ce^{CK^2\varepsilon^2} (\varepsilon + K^3\varepsilon^3 + K^5\varepsilon^5)(1+t)^\delta. \end{aligned}$$

The proof of Lemma 3.1 is complete. \square

4. Large time asymptotics. Now we are in a position to show the asymptotic expression (1.5). In view of (2.6), it suffices to find the asymptotics of $\alpha(t, \xi)$. The following lemma is essentially due to [7], [8], [14], [19], though we shall state a slightly refined version here. For the convenience of the readers, the proof will be provided in the appendix.

LEMMA 4.1. *Let $\beta_0(\xi), \kappa(\xi)$ be bounded continuous functions and suppose that $\text{Im } \kappa(\xi) \leq 0$ for all $\xi \in \mathbf{R}$. Let $r(t, \xi)$ satisfy*

$$\sup_{\xi \in \mathbf{R}} |r(t, \xi)| = O(t^{-1-\lambda}) \quad (t \rightarrow +\infty)$$

with some constant $\lambda > 0$. If $\beta(t, \xi)$ solves the differential equation

$$(4.1) \quad i \frac{\partial \beta}{\partial t} = \varepsilon^2 \frac{\kappa(\xi)}{t} |\beta|^2 \beta + \varepsilon^2 r(t, \xi), \quad \beta(1, \xi) = \beta_0(\xi)$$

for sufficiently small ε , then there exist continuous functions $\beta_\infty(\xi)$ and $\gamma(\xi)$, which satisfy $|\beta_\infty(\xi)| \leq C$ and $|\gamma(\xi)| \leq C\varepsilon^2$, such that

$$\begin{aligned} \beta(t, \xi) &= \frac{\beta_\infty(\xi) \exp \left\{ -i\varepsilon^2 |\beta_\infty(\xi)|^2 \text{Re } \kappa(\xi) \int_0^{\log t} \frac{d\sigma}{1 - 2\varepsilon^2 \text{Im } \kappa(\xi) (|\beta_\infty(\xi)|^2 \sigma + \gamma(\xi))} \right\}}{\sqrt{1 - 2\varepsilon^2 \text{Im } \kappa(\xi) (|\beta_\infty(\xi)|^2 \log t + \gamma(\xi))}} \\ &\quad + O(t^{-\lambda+\mu}) \end{aligned}$$

as $t \rightarrow +\infty$ uniformly in $\xi \in \mathbf{R}$, where μ is an arbitrarily small number.

Now, we set $\beta(t, \xi) = \varepsilon^{-1} \langle \xi \rangle^2 \alpha(t, \xi)$, $\kappa(\xi) = \langle \xi \rangle^{-4} N(1, i\xi)$, and

$$r(t, \xi) = \frac{1}{\varepsilon^2} \left(i \frac{\partial \beta}{\partial t} - \varepsilon^2 \frac{\kappa(\xi)}{t} |\beta|^2 \beta \right).$$

Then, since

$$\begin{aligned} r(t, \xi) &= \frac{1}{\varepsilon^3} \left\{ \left(i \frac{\partial \alpha}{\partial t} - \frac{1}{t} N(1, i\xi) |\alpha|^2 \alpha \right) + \left(i \frac{\partial \alpha_2}{\partial t} - \frac{1}{t} N(1, i\xi) |\alpha|^2 \alpha_2 \right) \right\} \\ &= \frac{1}{\varepsilon^3} (\rho_0(t, \xi) + \rho_2(t, \xi)), \end{aligned}$$

we deduce from (3.1) and (3.4) that

$$\sup_{\xi \in \mathbf{R}} |r(t, \xi)| \leq \frac{1}{\varepsilon^3} (\|\rho_0(t)\|_{L^\infty} + \|\rho_2(t)\|_{L^\infty}) \leq Ct^{-5/4}.$$

Therefore we can apply Lemma 4.1 with $\lambda = 1/4$ to get the asymptotics of β as $t \rightarrow +\infty$. Putting $a(\xi) = e^{-i\pi/4}\varepsilon\langle\xi\rangle^{-2}\beta_\infty(\xi)$ and $b(\xi) = \varepsilon^2\langle\xi\rangle^{-4}\gamma(\xi)$, we have

$$\alpha(t, \xi) = \frac{a(\xi) \exp \left\{ i\pi/4 - i|a(\xi)|^2 \operatorname{Re} N(1, i\xi) \int_0^{\log t} \frac{d\sigma}{1 - 2 \operatorname{Im} N(1, i\xi)(|a(\xi)|^2 \sigma + b(\xi))} \right\}}{\sqrt{1 - 2 \operatorname{Im} N(1, i\xi)(|a(\xi)|^2 \log t + b(\xi))}} + O(t^{-1/4+\mu})$$

as $t \rightarrow +\infty$. Finally, using (2.6), we arrive at the desired asymptotic expression for $u(t, x)$. \square

Remark 4.1. We can also deduce the L^2 -decay of $u(t)$ because

$$\|u(t)\|_{L_x^2} = \varepsilon \|\langle\xi\rangle^{-2}\beta(t)\|_{L_\xi^2} \leq C\varepsilon \|\beta(t)\|_{L_\xi^\infty}$$

and

$$\|\beta(t)\|_{L_\xi^\infty} = O((\log t)^{-1/2}) \quad \text{as } t \rightarrow +\infty$$

if $\sup_{\xi \in \mathbf{R}} \operatorname{Im} N(1, i\xi) < 0$.

Remark 4.2. It is possible to weaken (1.2) to a certain extent, but impossible to remove it completely. In fact, when (1.2) is replaced by

$$(4.2) \quad N(e^{i\theta}, 0) = e^{i\theta} N(1, 0) \quad \text{for } \theta \in \mathbf{R},$$

we can modify the above argument combining the idea of [5], [6] (see also the appendix of [19]) to show that Theorem 1.1 is still valid if $N(1, i\xi)$ in the statement is replaced by

$$\frac{1}{2\pi} \int_0^{2\pi} N(e^{i\theta}, i\xi e^{i\theta}) e^{-i\theta} d\theta.$$

Note that (4.2) is just what excludes $u^3, \bar{u}^3, u\bar{u}^2$ from all possible cubic nonlinear terms, but it is not a technical assumption because for these three nonlinearities we can find a class of initial data for which the solution has another kind of asymptotic profile than (1.5) (see [7] and [8] for details).

Appendix. We give a proof of Lemma 4.1 following [19] and [14] with some modifications. Because of the uniqueness for (4.1), $\beta(t, \xi)$ admits the decomposition

$$\beta(t, \xi) = \frac{p(t, \xi)}{\sqrt{q(t, \xi)}},$$

where $p(t, \xi)$ and $q(t, \xi)$ satisfy

$$(A.1) \quad \begin{cases} \partial_t p(t, \xi) = -i\varepsilon^2 \frac{\operatorname{Re} \kappa(\xi)}{t} \frac{|p(t, \xi)|^2}{q(t, \xi)} p(t, \xi) - i\varepsilon^2 \sqrt{q(t, \xi)} r(t, \xi), \\ \partial_t q(t, \xi) = -2\varepsilon^2 \frac{\operatorname{Im} \kappa(\xi)}{t} |p(t, \xi)|^2, \\ p(1, \xi) = \beta_0(\xi), \quad q(1, \xi) = 1. \end{cases}$$

Note that $p(t, \xi)$ is complex-valued, while $q(t, \xi)$ is real and strictly positive. In order to obtain the desired conclusion, it is sufficient to get the asymptotics of $p(t, \xi)$ and $q(t, \xi)$ as $t \rightarrow +\infty$.

We first show that there exists a positive constant A such that

$$(A.2) \quad \sup_{(t,\xi) \in [1,\infty) \times \mathbf{R}} |p(t,\xi)| < A.$$

We shall argue by contradiction: Suppose that for any $A > \sup_{\xi \in \mathbf{R}} |\beta_0(\xi)|$ there exists a finite time $T_A \in (1, \infty)$ such that

$$\sup_{(t,\xi) \in [1,T_A) \times \mathbf{R}} |p(t,\xi)| \leq A \quad \text{and} \quad \sup_{\xi \in \mathbf{R}} |p(T_A,\xi)| = A.$$

Then, from the second equation of (A.1), we have

$$1 \leq q(t,\xi) \leq 1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) A^2 \log t$$

for $t \in [1, T_A]$. On the other hand, it follows from the first equation of (A.1) that

$$\begin{aligned} \partial_t \left(|p(t,\xi)|^2 \right) &= 2 \operatorname{Re} \left(\overline{p(t,\xi)} \partial_t p(t,\xi) \right) \\ &= 2\varepsilon^2 \operatorname{Re} \left(\overline{p(t,\xi)} \sqrt{q(t,\xi)} r(t,\xi) \right) \\ &\leq 2\varepsilon^2 |p(t,\xi)| \left| \sqrt{q(t,\xi)} r(t,\xi) \right|. \end{aligned}$$

So we have

$$\begin{aligned} |p(T_A,\xi)| &\leq |\beta_0(\xi)| + \int_1^{T_A} \varepsilon^2 \left| \sqrt{q(\tau,\xi)} r(\tau,\xi) \right| d\tau \\ &\leq C + \int_1^\infty C \varepsilon^2 (1 + A^2 \varepsilon^2 \log \tau)^{1/2} \tau^{-1-\lambda} d\tau \\ &\leq C(1 + A\varepsilon). \end{aligned}$$

When we choose $A = 4C$ and $\varepsilon_1 = 1/A$, we have

$$\sup_{\xi \in \mathbf{R}} |p(T_A,\xi)| \leq \frac{A}{2} < A$$

for $\varepsilon \in (0, \varepsilon_1]$, which is the desired contradiction. Hence (A.2) must hold for some A . Also we have

$$1 \leq q(t,\xi) \leq 1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) A^2 \log t$$

for any $t \geq 1$. Next we define

$$\Theta(t,\xi) = \int_1^t \frac{\operatorname{Re} \kappa(\xi) |p(\tau,\xi)|^2 d\tau}{q(\tau,\xi) \tau}$$

so that $\partial_t(p(t,\xi)e^{i\varepsilon^2\Theta(t,\xi)}) = -i\varepsilon^2 \sqrt{q(t,\xi)} r(t,\xi) e^{i\varepsilon^2\Theta(t,\xi)}$. Since

$$\sup_{\xi \in \mathbf{R}} \left| \sqrt{q(t,\xi)} r(t,\xi) \right| \leq C(1 + \varepsilon^2 \log t)^{1/2} t^{-1-\lambda} \leq C t^{-1-\lambda+\mu},$$

we obtain

$$\sup_{\xi \in \mathbf{R}} |p(t,\xi) - p_\infty(\xi) e^{-i\varepsilon^2\Theta(t,\xi)}| \leq C \varepsilon^2 t^{-\lambda+\mu},$$

where

$$p_\infty(\xi) = \beta_0(\xi) - i\varepsilon^2 \int_1^\infty \sqrt{q(\tau, \xi)} r(\tau, \xi) e^{i\varepsilon^2 \Theta(\tau, \xi)} d\tau.$$

We also set $q_\infty(t, \xi) = 1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) (|p_\infty(\xi)|^2 \log t + \gamma(\xi))$ with

$$\gamma(\xi) = \int_1^\infty (|p(\tau, \xi)|^2 - |p_\infty(\xi)|^2) \frac{d\tau}{\tau}.$$

Noting that

$$\begin{aligned} ||p(t, \xi)|^2 - |p_\infty(\xi)|^2| &\leq |p(t, \xi) - p_\infty(\xi) e^{-i\varepsilon^2 \Theta(t, \xi)}| (|p(t, \xi)| + |p_\infty(\xi)|) \\ &\leq C\varepsilon^2 t^{-\lambda+\mu}, \end{aligned}$$

we can see that

$$|q(t, \xi) - q_\infty(t, \xi)| \leq 2\varepsilon^2 |\operatorname{Im} \kappa(\xi)| \int_t^\infty ||p(\tau, \xi)|^2 - |p_\infty(\xi)|^2| \frac{d\tau}{\tau} \leq C\varepsilon^4 t^{-\lambda+\mu}.$$

Let us also introduce

$$\Phi(t, \xi) = \frac{\operatorname{Re} \kappa(\xi)}{t} \left(\frac{|p(t, \xi)|^2}{q(t, \xi)} - \frac{|p_\infty(\xi)|^2}{q_\infty(t, \xi)} \right).$$

Then we have

$$i\varepsilon^2 \Theta(t, \xi) = i\varepsilon^2 |p_\infty(\xi)|^2 \operatorname{Re} \kappa(\xi) \int_1^t \frac{d\tau}{q_\infty(\tau, \xi)} + i\varepsilon^2 \int_1^\infty \Phi(\tau, \xi) d\tau - i\varepsilon^2 \int_t^\infty \Phi(\tau, \xi) d\tau$$

and

$$\begin{aligned} \int_t^\infty |\Phi(\tau, \xi)| d\tau &\leq C \int_t^\infty \left(\frac{||p(\tau, \xi)|^2 - |p_\infty(\xi)|^2|}{q(\tau, \xi)} + \frac{|p_\infty(\xi)|^2 |q(\tau, \xi) - q_\infty(\tau, \xi)|}{q(\tau, \xi) q_\infty(\tau, \xi)} \right) \frac{d\tau}{\tau} \\ &\leq C\varepsilon^2 \int_t^\infty \frac{d\tau}{\tau^{1+\lambda-\mu}} \\ &\leq C\varepsilon^2 t^{-\lambda+\mu}. \end{aligned}$$

Therefore, putting $\beta_\infty(\xi) = p_\infty(\xi) \exp(-i\varepsilon^2 \int_1^\infty \Phi(\tau, \xi) d\tau)$, we obtain

$$\begin{aligned} p(t, \xi) &= p_\infty(\xi) e^{-i\varepsilon^2 \Theta(t, \xi)} + O(t^{-\lambda+\mu}) \\ &= \beta_\infty(\xi) \exp\left(-i\varepsilon^2 |p_\infty(\xi)|^2 \operatorname{Re} \kappa(\xi) \int_1^t \frac{d\tau}{q_\infty(\tau, \xi)}\right) + O(t^{-\lambda+\mu}) \\ &= \beta_\infty(\xi) \exp\left(-i\varepsilon^2 |\beta_\infty(\xi)|^2 \operatorname{Re} \kappa(\xi) \int_0^{\log t} \frac{d\sigma}{1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) (|\beta_\infty(\xi)|^2 \sigma + \gamma(\xi))}\right) \\ &\quad + O(t^{-\lambda+\mu}) \end{aligned}$$

as well as

$$\begin{aligned} \frac{1}{\sqrt{q(t, \xi)}} &= \frac{1}{\sqrt{q_\infty(t, \xi)}} + \frac{q_\infty(t, \xi) - q(t, \xi)}{\sqrt{q(t, \xi) q_\infty(t, \xi)} (\sqrt{q(t, \xi)} + \sqrt{q_\infty(t, \xi)})} \\ &= \frac{1}{\sqrt{1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) (|\beta_\infty(\xi)|^2 \log t + \gamma(\xi))}} + O(t^{-\lambda+\mu}), \end{aligned}$$

whence

$$\begin{aligned} \beta(t, \xi) &= \frac{p(t, \xi)}{\sqrt{q(t, \xi)}} \\ &= \frac{\beta_\infty(\xi) \exp\left(-i\varepsilon^2 |\beta_\infty(\xi)|^2 \operatorname{Re} \kappa(\xi) \int_0^{\log t} \frac{d\sigma}{1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) (|\beta_\infty(\xi)|^2 \sigma + \gamma(\xi))}\right)}{\sqrt{1 - 2\varepsilon^2 \operatorname{Im} \kappa(\xi) (|\beta_\infty(\xi)|^2 \log t + \gamma(\xi))}} \\ &\quad + O(t^{-\lambda+\mu}) \end{aligned}$$

as $t \rightarrow +\infty$. \square

REFERENCES

- [1] H. CHIHARA, *Local existence for the semilinear Schrödinger equations in one space dimension*, J. Math. Kyoto Univ., 34 (1994), pp. 353–367.
- [2] N. HAYASHI AND P. I. NAUMKIN, *Asymptotic behavior in time of solutions to the derivative nonlinear Schrödinger equation revisited*, Discrete Contin. Dynam. Systems, 3 (1997), pp. 383–400.
- [3] N. HAYASHI AND P. I. NAUMKIN, *Asymptotic behavior in time of solutions to the derivative nonlinear Schrödinger equation*, Ann. Inst. H. Poincaré Phys. Théor., 68 (1998), pp. 159–177.
- [4] N. HAYASHI AND P. I. NAUMKIN, *Asymptotics for large time of solutions to the nonlinear Schrödinger and Hartree equations*, Amer. J. Math., 120 (1998), pp. 369–389.
- [5] N. HAYASHI AND P. I. NAUMKIN, *Large time behavior of solutions for derivative cubic nonlinear Schrödinger equations without a self-conjugate property*, Funkcial. Ekvac., 42 (1999), pp. 311–324.
- [6] N. HAYASHI AND P. I. NAUMKIN, *Asymptotics of small solutions to nonlinear Schrödinger equations with cubic nonlinearities*, Int. J. Pure Appl. Math., 3 (2002), pp. 255–273.
- [7] N. HAYASHI AND P. I. NAUMKIN, *Large time behavior for the cubic nonlinear Schrödinger equation*, Canad. J. Math., 54 (2002), pp. 1065–1085.
- [8] N. HAYASHI AND P. I. NAUMKIN, *On the asymptotics for cubic nonlinear Schrödinger equations*, Complex Var. Theory Appl., 49 (2004), pp. 339–373.
- [9] N. HAYASHI, P. I. NAUMKIN, AND H. UCHIDA, *Large time behavior of solutions for derivative cubic nonlinear Schrödinger equations*, Publ. Res. Inst. Math. Sci., 35 (1999), pp. 501–513.
- [10] N. HAYASHI AND T. OZAWA, *Modified wave operators for the derivative nonlinear Schrödinger equation*, Math. Ann., 298 (1994), pp. 557–576.
- [11] N. HAYASHI AND T. OZAWA, *Remarks on nonlinear Schrödinger equations in one space dimension*, Differential Integral Equations, 7 (1994), pp. 453–461.
- [12] L. HÖRMANDER, *Lectures on Nonlinear Hyperbolic Differential Equations*, Math. Appl. (Berlin) 26, Springer-Verlag, Berlin, 1997.
- [13] S. KATAYAMA AND Y. TSUTSUMI, *Global existence of solutions for nonlinear Schrödinger equations in one space dimension*, Comm. Partial Differential Equations, 19 (1994), pp. 1971–1997.
- [14] H. KUBO, *Asymptotic behavior of solutions to semilinear wave equations with dissipative structure*, Discrete Contin. Dyn. Syst. (Suppl.), 2007 (2007), pp. 602–613.
- [15] H. LINDBLAD AND A. SOFFER, *Scattering and small data completeness for the critical nonlinear Schrödinger equation*, Nonlinearity, 19 (2006), pp. 345–353.
- [16] T. OZAWA, *Long range scattering for nonlinear Schrödinger equations in one space dimension*, Comm. Math. Phys., 139 (1991), pp. 479–493.
- [17] T. OZAWA, *On the nonlinear Schrödinger equations of derivative type*, Indiana Univ. Math. J., 45 (1996), pp. 137–163.
- [18] A. SHIMOMURA, *Asymptotic behavior of solutions for Schrödinger equation with dissipative nonlinearities*, Comm. Partial Differential Equations, 31 (2006), pp. 1407–1423.
- [19] H. SUNAGAWA, *Large time behavior of solutions to the Klein-Gordon equation with nonlinear dissipative terms*, J. Math. Soc. Japan, 58 (2006), pp. 379–400.
- [20] H. SUNAGAWA, *Lower bounds of the lifespan of small data solutions to the nonlinear Schrödinger equations*, Osaka J. Math., 43 (2006), pp. 771–789.
- [21] Y. TSUTSUMI, *The null gauge condition and the one dimensional nonlinear Schrödinger equation with cubic nonlinearity*, Indiana Univ. Math. J., 43 (1994), pp. 241–254.

QUASILINEAR PARABOLIC SYSTEMS WITH MIXED BOUNDARY CONDITIONS ON NONSMOOTH DOMAINS*

MATTHIAS HIEBER[†] AND JOACHIM REHBERG[‡]

Abstract. In this paper we investigate quasilinear systems of reaction-diffusion equations with mixed Dirichlet–Neumann boundary conditions on nonsmooth domains. Using techniques from maximal regularity and heat-kernel estimates we prove the existence of a unique solution to systems of this type.

Key words. quasilinear parabolic system, mixed Dirichlet–Neumann conditions, L^∞ -coefficients

AMS subject classifications. Primary, 35A05, 35B65; Secondary, 35K15/20

DOI. 10.1137/070683829

1. Introduction. The theory of quasilinear parabolic systems has many applications to evolution problems in natural sciences, see, e.g., [2], [1], [5], [6], [20], [10], [16], [30], and [41]. In this paper we investigate in particular systems of reaction-diffusion equations with *mixed* Dirichlet–Neumann boundary conditions on nonsmooth domains $\Omega \subset \mathbb{R}^n$ for $n = 2, 3$ of the form

$$(1.1) \quad \begin{aligned} u'_k - \operatorname{div}(G_k(v)\mu_k \nabla v_k) &= R_k(t, v), & t \in]T_0, T[, x \in \Omega, \\ u_k &= b_k F_k(v_k), & t \in [T_0, T[, x \in \Omega, \\ \nu \cdot \mu_k \nabla v_k &= 0, & t \in [T_0, T[, x \in \Gamma_N, \\ v_k &= \phi_k, & t \in [T_0, T[, x \in \Gamma_D, \\ v_k(T_0) &= v_{0k}, & x \in \Omega. \end{aligned}$$

Here $v = (v_1, \dots, v_m)$, $\mu_k \in L^\infty(\Omega, M_{n \times n})$ are diffusion coefficients, $b_k \in L^\infty(\Omega)$ are reference densities, and R_k, G_k, F_k denote the reaction, diffusion, and superposition terms for $k \in \{1, \dots, m\}$.

In many concrete problems which are described as a system of the form (1.1), the underlying domain is nonsmooth and the coefficient functions b_k and μ_k are discontinuous. We therefore aim for minimal smoothness assumptions on the boundary $\partial\Omega$ of Ω , the coefficient functions b_k and μ_k , as well as on the interface between the Neumann boundary part Γ_N of $\partial\Omega$ and the Dirichlet boundary part $\Gamma_D = \partial\Omega \setminus \Gamma_N$. More precisely, we generally assume that $\Omega \subset \mathbb{R}^n$ is a Lipschitz domain and $\Omega \cup \Gamma_N$ is regular in the sense of Gröger (see the definition below). Note that the situation where the boundary of Ω is smooth and consists of two *separated* parts, one with Dirichlet and the other with Neumann boundary conditions, has been studied before by Amann in [2].

Our setting includes even nonlocal diffusion terms as occurring in models describing the diffusion of bacteria (see [7], [8] and the references therein). In detail, the velocity at which the motion takes place is given by Fourier’s law and the (relative)

*Received by the editors February 27, 2007; accepted for publication (in revised form) September 14, 2007; published electronically April 23, 2008.

<http://www.siam.org/journals/sima/40-1/68382.html>

[†]Technische Universität Darmstadt, Fachbereich Mathematik, Schlossgartenstr. 7, D-64298 Darmstadt, Germany (hieber@mathematik.tu-darmstadt.de).

[‡]Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstr. 39, D-10117 Berlin, Germany (rehberg@wias-berlin.de).

diffusion coefficient $G = G_k$ depends on the solution v in the form

$$G(v) = \eta \left(\int_{\Omega} z(x)v(x) dx \right),$$

where η is a (strictly) positive, continuously differentiable function on \mathbb{R} .

Our approach includes reaction terms R_k depending discontinuously on time t , which is important in many examples (see [41], [26], [30]), in particular in the control theory of parabolic equations. Alternatively, the reader should think, e.g., of a manufacturing process for semiconductors, where at a certain moment light is switched on/off and, of course, parameters in the chemical process change abruptly.

An interesting example for a reaction term stems from the thermistor problem (see [38], [3] and the references therein, see also [6]) which describes the combined processes of heat conduction and electrical conduction in a body. This is of importance in the industrially important process of electrical welding. To be precise, the reaction term is of the form $\sigma(v)|\nabla\varphi|^2$, where φ is the solution of an auxiliary elliptic problem, into which v enters as a parameter (see details in section 4.2). Observe that the quadratic gradient term is a critical one and in general not easy to handle (see [38]), while in our context it does not cause additional difficulties.

Note that the original formulation of the evolution equation in terms of balance laws takes the form (see [39, Ch. 21], see also [5])

$$(1.2) \quad \frac{\partial}{\partial t} \int_{\Omega'} u_k dx + \int_{\partial\Omega'} \nu \cdot j_k d\sigma = \int_{\Omega'} R_k dx; \quad j_k = j_k(v) = G_k(v)\mu_k \nabla v_k,$$

where Ω' stands for any (Lipschitzian) subdomain of Ω . Within the variational theory of weak solutions, however, the indicator functions of the subdomains are not admissible test functions. Therefore the integral formulation (1.2) is equivalent to the above evolution equation only if the weak solutions have some additional regularity. It is the main advantage of the present concept that the divergence of the corresponding current $j_k(v)$ indeed is a function, not only a distribution. In a strict sense, only this justifies the application of Gauss' theorem to calculate the normal components of the currents over boundaries of suitable subdomains. Moreover, the fact $\operatorname{div} j_k \in L^p$ is also of importance for the numerical treatment of (1.1), as the formulation (1.2) is the basis of finite volume methods (see [18])—namely in the sense of local balances. Global existence results for (1.1) cannot be expected within such a general approach (see, e.g., [17] or [6] and the references therein).

In contrast to many papers where existence and uniqueness results for quasilinear parabolic systems are based on the construction of an appropriate evolution operator (see, e.g., [1]), our approach relies heavily on maximal L^p -estimates for the linear part of (1.1). In fact, after rewriting (1.1) as an abstract evolution equation in $L^p(\Omega)^m$ of the form

$$(1.3) \quad \begin{aligned} w' - H(t, w)(\operatorname{div}(\mu \nabla w)) &= S(t, w), \\ w(T_0) &= v_0 - \phi(T_0), \end{aligned}$$

our strategy to solve (1.3) follows the approach of Clément and Li [10] and Prüss [34]. The advantage of the given situation (1.1) is that subtle techniques from harmonic analysis as well as heat-kernel methods can be used to prove the central L^p -estimates of the linear part. In order to apply these methods in our situation, one needs embedding properties of certain interpolation spaces between the domain of the L^p -realization of

the underlying elliptic operators and $L^p(\Omega)$ into $W^{1,2p}(\Omega)$. This embedding property rests on the assumption that the operators formally defined by

$$-\nabla \cdot \mu_k \nabla + 1 : W_{\Gamma_N}^{1,q}(\Omega) \rightarrow W_{\Gamma_N}^{-1,q}(\Omega)$$

provide topological isomorphisms for some $q > n$. Note that this assumption restricts the physical dimension of the problems to two and three, because one knows that the solution for a mixed boundary value problem generically has singularities, precisely, one cannot expect that its gradient lies in L^4 ; see Shamir’s famous counterexample [37]. On the other hand, in two and three dimensions the assumption is fulfilled for many geometric constellations and (even discontinuous) coefficient functions; see section 4.

2. Preliminaries. Let $\Omega \subset \mathbb{R}^n$ be a bounded Lipschitz domain and assume that $n = 2$ or $n = 3$. (Concerning the notions “Lipschitz domain” and “domain with Lipschitz boundary” we follow [23, section 1.2.1].) Denote by $\Gamma_N \subset \partial\Omega$ an open subset of $\partial\Omega$. For $1 < q < \infty$ we define $W_{\Gamma_N}^{1,q}(\Omega)$ as the closure of

$$\{\psi|_{\Omega} : \psi \in C_c^\infty(\mathbb{R}^n), \text{supp } \psi \cap (\partial\Omega \setminus \Gamma_N) = \emptyset\}$$

in the Sobolev space $W^{1,q}(\Omega)$. If $q = 2$, we write $H^1(\Omega)$ or $H_{\Gamma_N}^1(\Omega)$ instead of $W^{1,2}(\Omega)$ or $W_{\Gamma_N}^{1,2}(\Omega)$. Of course, if $\Gamma_N = \emptyset$, then $W_{\Gamma_N}^{1,q}(\Omega) = W_0^{1,q}(\Omega)$. Moreover, throughout this work we always suppose that $\Omega \cup \Gamma_N$ is regular in the sense of Gröger (see [24]), this means that, for all $x \in \partial\Omega$ there exist open sets $U_x, V_x \subset \mathbb{R}^n$ and a bi-Lipschitz transform Ψ_x from U_x onto V_x such that $x \in U_x, \Psi_x(x) = 0$ and $\Psi_x(U_x \cap (\Omega \cup \Gamma_N))$ coincides with one of the sets

$$\begin{aligned} E_1 &:= \{x \in \mathbb{R}^n : \max_{l=1,\dots,n} |x_l| < 1, x_n < 0\}, \\ E_2 &:= \{x \in \mathbb{R}^n : \max_{l=1,\dots,n} |x_l| < 1, x_n \leq 0\}, \\ E_3 &:= \{x \in E_2 : x_n < 0 \text{ or } x_1 > 0\}. \end{aligned}$$

It is not hard to see that every Lipschitz domain, and also its closure, is regular in the sense of Gröger, the corresponding model sets are then E_1 or E_2 , respectively; see [23]. Moreover, if $\Omega \subset \mathbb{R}^2$ is a bounded Lipschitz domain and $\partial\Omega \setminus \Gamma_N$ is the finite union of (nondegenerate) closed arc pieces from the boundary, then $\Omega \cup \Gamma_N$ is regular in the sense of Gröger.

Finally, for $k \in \{1, \dots, m\}$, let $\mu_k \in L^\infty(\Omega, M_{n \times n})$, where $M_{n \times n}$ denotes the set of all real, symmetric $n \times n$ matrices. Suppose that additionally

$$(2.1) \quad \inf_{x \in \Omega} \inf_{|\zeta|=1} \mu_k(x) \zeta \cdot \zeta > 0.$$

For a closed subspace $V \subseteq H^1(\Omega)$ such that $H_0^1(\Omega) \subseteq V$ we define the form $a_k : V \times V \rightarrow \mathbb{R}$ by

$$a_k(u, v) := - \int_{\Omega} \mu_k \nabla u \cdot \nabla v \, dx, \quad u, v \in V.$$

The form induces a continuous mapping $\mathcal{A}_k : V \rightarrow V'$ such that

$$(2.2) \quad a_k(u, v) = (\mathcal{A}_k u | v), \quad u, v \in V.$$

Here, for $v \in L^2(\Omega)$, $f_v(u) := (v|u)_{L^2}$ defines an element $f_v \in V'$ and $v \mapsto f_v : L^2(\Omega) \rightarrow V'$ defines a continuous injection. In the following, we identify v with f_v .

We then define the operator A_k as

$$(2.3) \quad D(A_k) := \{u \in V : \exists f \in L^2(\Omega), a_k(u, \phi) = (f|\phi) \forall \phi \in V\},$$

$$(2.4) \quad A_k u := f.$$

It is well known that A_k generates an analytic semigroup on $L^2(\Omega)$ which is positivity preserving. Furthermore, this semigroup extends to a C_0 -semigroup of contractions on $L^p(\Omega)$ for all $1 < p < \infty$; see, [22, Thm. 4.9]. The realization of its generator in L^p is denoted by A_k^p .

3. Main result. We start this section by giving precise assumptions on the coefficients and functions being involved in problem (1.1). In order to do so, let $0 \leq T_0 < T_1$ and set $J :=]T_0, T_1[$. For $k \in \{1, \dots, m\}$ let $\mu_k \in L^\infty(\Omega, M_{n \times n})$ and assume that (2.1) is satisfied.

Moreover, let for every $k \in \{1, \dots, m\}$ the functions $b_k \in L^\infty(\Omega; \mathbb{R})$ be bounded from below by some positive constant.

We assume the following for all $k \in \{1, \dots, m\}$.

- (Op) There exists $p > \frac{n}{2}$ such that each $\mathcal{A}_k - Id$ is a topological isomorphism from $W_{\Gamma_N}^{1,2p}(\Omega)$ onto $W_{\Gamma_N}^{-1,2p}(\Omega)$. For all what follows we fix a number $r > \frac{4p}{2p-n}$.
- (Su) There exists $f_k \in C^2(\mathbb{R})$, positive, with strictly positive derivative, such that F_k is the superposition operator induced by f_k ; i.e., $F_k(v)(x) = (f_k \circ v)(x) = f_k(v(x))$, $x \in \Omega$.
- (Ga) The mapping $G_k : (W^{1,2p}(\Omega))^m \rightarrow W^{1,2p}(\Omega)$ is locally Lipschitz.
- (Gb) For any ball in $(W^{1,2p}(\Omega))^m$ there exists $\delta > 0$ such that $G_k(u) \geq \delta$ for all u from this ball.
- (Ra) The function $R_k : J \times (W^{1,2p}(\Omega))^m \rightarrow L^p(\Omega)$ is of Carathéodory type, i.e., $R_k(\cdot, u)$ is measurable for all $u \in (W^{1,2p}(\Omega))^m$ and $R_k(t, \cdot)$ is continuous for a.a. $t \in J$.
- (Rb) $R_k(\cdot, 0) \in L^r(J, L^p(\Omega))$ and for $\beta > 0$ there exists $g_\beta \in L^r(J)$ such that

$$\|R_k(t, u) - R_k(t, \tilde{u})\|_{L^p} \leq g_\beta(t) \|u - \tilde{u}\|_{W^{1,2p}}, \quad t \in J,$$

provided $\max(\|u\|_{W^{1,2p}}, \|\tilde{u}\|_{W^{1,2p}}) \leq \beta$.

- (BC) $\phi_k \in C(\bar{J}; W^{1,2p}(\Omega)) \cap W^{1,r}(J; L^p(\Omega))$ and $A_k \phi_k(t) = 0$ for all $t \in J$.
- (IC) $v_{0k} - \phi_k(T_0) \in (L^p(\Omega), D(A_k^p))_{1-\frac{1}{r}, r}$.

The assumptions imply that the system (1.1) may be (formally) rewritten as a quasilinear system of the form

$$(3.1) \quad \begin{aligned} w'_k - H_k(t, w) A_k w_k &= T_k(t, w), \quad k = 1, \dots, m, \\ w(T_0) &= v_0 - \phi(T_0), \end{aligned}$$

where

$$(3.2) \quad T_k(t, w) := (b_k f'_k(w_k + \phi_k(t)))^{-1} [\nabla G_k(w + \phi(t)) \cdot [\mu_k \nabla(w_k + \phi_k(t))] + Q_k(t, w) - \frac{\partial \phi_k}{\partial t}(t)]$$

with

$$(3.3) \quad H_k(t, z) := \frac{G_k(z + \phi(t))}{b_k f'_k(z_k + \phi_k(t))}, \quad t \in J, z \in (W^{1,2p}(\Omega))^m,$$

$$(3.4) \quad Q_k(t, z) := \frac{R_k(t, z + \phi(t))}{b_k f'_k(z_k + \phi_k(t))}, \quad t \in J, z \in (W^{1,2p}(\Omega))^m.$$

We are now in the position to state the main result of this paper.

THEOREM 3.1. *Let $1 < r, p < \infty$ such that $r > \frac{4p}{2p-n}$, where $n \in \{2, 3\}$. Assume that the assumptions (Op), (Su), (Ga), (Gb), (Ra), (Rb), (BC), and (IC) are satisfied. Then there exists a unique local solution $w = (w_1, \dots, w_m)$ for (3.1) on an interval $I =]T_0, T[$ satisfying*

$$(3.5) \quad w_k \in W^{1,r}(I; L^p(\Omega)) \cap L^r(I; D(A_k)), \quad k \in \{1, \dots, m\}.$$

COROLLARY 3.2. *Each w_k is Hölder-continuous simultaneously in space and time.*

Some remarks at this point are in order.

Remark 3.3.

- (a) We refer the reader to section 4 for precise geometric and smoothness conditions implying the validity of assumption (Op).
- (b) Besides the exponential, a typical example for a function f satisfying assumption (Su) is the Fermi–Dirac distribution function

$$f(t) := \frac{2}{\sqrt{\pi}} \int_0^\infty \frac{\sqrt{s}}{1 + e^{s-t}} ds.$$

- (c) Suppose that v_k coincides on Γ_D with a function $\phi \in C^1(J, W^{1,2p}(\Omega))$. Then there exists ϕ_k satisfying assumption (BC).
- (d) Note that condition (BC) implies $\nu \cdot \mu_k \nabla \phi_k = 0$ on Γ_N . This, together with the property (3.5), yields the Neumann boundary condition for v_k on Γ_N ; see [19, section II.2], [9, section 1.2].

4. Examples. Consider Ω and Γ_N , the subset of $\partial\Omega$ on which the Neumann boundary condition is prescribed. In this section we describe geometric configurations and coefficient functions for which Theorem 3.1 holds true. Furthermore, we present concrete examples of mappings G_k and reaction terms R_k fitting in our framework.

4.1. Geometric configurations and coefficients. In this subsection we will give a list of examples in order to show that (Op) is not an unjustified ad hoc assumption but fulfilled in many cases, even if the geometry is nonsmooth, the boundary conditions are mixed, and the coefficients may jump. Notice that for—realistic—nonsmooth situations a (strict) upper bound for the integrability index for the gradient of the solution is (at most) 4. That means in detail that, this limitation may be caused only by nonsmoothness of the domain (see [28, Thm. A]), only by the occurrence of mixed boundary conditions (despite smooth data, see [37]), or only by nonsmooth coefficients (see [33] or [14]).

We start with a result, due to Gröger [24], which completely covers the two-dimensional case.

PROPOSITION 4.1. *Assume that $\Omega \cup \Gamma_N$ is regular in the sense of Gröger and that the coefficient function satisfies the assumptions made in section 2. Then there exists $q > 2$ such that $\mathcal{A}_k - Id$ is a topological isomorphism from $W_{\Gamma_N}^{1,q}(\Omega)$ onto $W_{\Gamma_N}^{-1,q}(\Omega)$.*

In what follows we present three-dimensional settings, here always supposing that the domain Ω is a domain with Lipschitz boundary. We begin with the case where the coefficient function is continuous at least in a neighborhood of the boundary.

PROPOSITION 4.2. *Assume that $\Gamma_N = \emptyset$ or $\Gamma_N = \partial\Omega$ (pure Dirichlet or pure Neumann case). $\Omega_\circ \subset \Omega$ is another domain which is C^1 and which does not touch the boundary of Ω . $\mu_k|_{\Omega_\circ} \in BUC(\Omega_\circ)$ and $\mu_k|_{\Omega \setminus \bar{\Omega}_\circ} \in BUC(\Omega \setminus \bar{\Omega}_\circ)$. Then (Op) holds.*

Remark 4.3. The proposition is one of the main results from [15] and rests heavily on regularity results for the Dirichlet/Neumann Laplacian (see [28], [42]) and nontrivial estimates around points from $\partial\Omega_o$. In particular, the reader should carefully notice that the C^1 property of Ω_o is not dispensable without completely losing the result; see Elschner’s counterexample in [15], see also [14].

The next proposition describes two cases with pure Dirichlet conditions ($\Gamma_N = \emptyset$), where the discontinuities of the coefficient function are located also near the boundary.

PROPOSITION 4.4. *There exists $q > 3$ such that $\mathcal{A}_k - Id$ is a topological isomorphism from $W_{\Gamma_N}^{1,q}(\Omega)$ onto $W_{\Gamma_N}^{-1,q}(\Omega)$ if one of the following conditions is satisfied.*

- (i) Ω is a polyhedron. There are hyperplanes $\mathcal{H}_1, \dots, \mathcal{H}_n$ in \mathbb{R}^3 which meet at most in a vertex of the polyhedron such that the coefficient function μ_k is constantly a real, symmetric, positive definite 3×3 matrix on each of the connected components of $\Omega \setminus \cup_{l=1}^n \mathcal{H}_l$. Moreover, for every edge on the boundary, induced by a hyperplane \mathcal{H}_l , the angles between the outer boundary plane and the hyperplane \mathcal{H}_l do not exceed π .
- (ii) $\Omega_o \subset \Omega$ is a Lipschitz domain such that $\partial\Omega_o \cap \Omega$ is a C^1 surface. Moreover, $\partial\Omega$ and $\partial\Omega_o$ meet suitably, this means that for any point x from the boundary of $\partial\Omega \cap \partial\Omega_o$ within $\partial\Omega$ there is an open neighborhood \mathcal{U}_x of x in \mathbb{R}^3 and a C^1 diffeomorphism Φ_x from \mathcal{U}_x onto an open subset of \mathbb{R}^3 such that
 - $\Phi_x(\mathcal{U}_x \cap \Omega)$ equals a convex polyhedron \mathcal{K}_x ,
 - $\Phi_x(\mathcal{U}_x \cap \Omega \cap \partial\Omega_o) = \mathcal{K}_x \cap \mathcal{H}_x$, where \mathcal{H}_x is a plane which contains $\Phi_x(x)$ and an inner point of \mathcal{K}_x . $\mu_k|_{\Omega_o} \in BUC(\Omega_o)$ and $\mu_k|_{\Omega \setminus \bar{\Omega}_o} \in BUC(\Omega \setminus \bar{\Omega}_o)$.

The constellation (i) is treated in [14], the second in [15].

The last cases we present really affect the case of mixed boundary conditions.

PROPOSITION 4.5. *Assumption (Op) is also satisfied in any of the following two cases.*

- (i) Ω is a convex polyhedron, $\overline{\Gamma_N} \cap (\partial\Omega \setminus \Gamma_N)$ is a finite union of line segments, $\mu_k \equiv 1$.
- (ii) Ω is a three-dimensional prismatic domain with triangular basis. Γ is half of one of its upright sides. Ξ is a plane that intersects the (relative) boundary of Γ within $\partial\Omega$ in only finitely many points. The coefficient function is constant on both components of $\Omega \setminus \Xi$.

Remark 4.6. The assertion for (i) is shown in [11] (see Corollary 3.12), while the proof for (ii) is given in [25]. Notice that (ii) is by no means artificial: $\Omega \cup \Gamma$ may (alternatively) be taken as Gröger’s third local model set in the description of geometric settings including mixed boundary conditions; see section 2 or [24] for further details. Additionally, it is the first three-dimensional constellation treated in the literature—in view of (Op)—which includes mixed boundary conditions and a discontinuous coefficient function, which is necessary, e.g., in the modeling of heterogeneous semiconductors [36]. (See [25] for other settings and further details.)

Remark 4.7. The operator $\mathcal{A}_k - Id$ also fulfills (Op) if the following condition is satisfied: there is a covering of $\bar{\Omega}$ by open sets $\mathcal{U}_1, \dots, \mathcal{U}_l$ such that for every $j \in \{1, \dots, l\}$ the setting $\Omega_j := \Omega \cap \mathcal{U}_j$, $\Gamma_j := \Gamma_N \cap \mathcal{U}_j$, and the restriction of μ_k to $\Omega \cap \mathcal{U}_j$ satisfy one of the conditions of the foregoing propositions; compare [24, Lemma 2]. Additionally, one can show by perturbation arguments (as, e.g., carried out in [15]) that in many cases admissible constellations are preserved under C^1 diffeomorphisms.

In the following we illustrate two admissible three-dimensional settings. On the left-hand side of Figure 4.1 one assumes Neumann conditions on the top of the upper cuboid, otherwise Dirichlet conditions. On the right-hand side of the figure, the

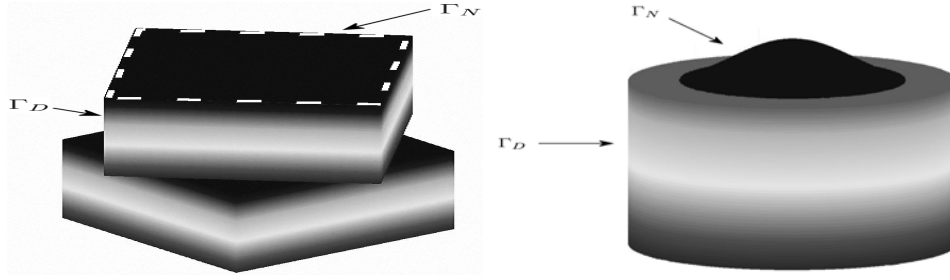


FIG. 4.1.

boundary of the cylinder is subject to Dirchlet conditions except for the upper “hat,” where Neumann conditions are prescribed.

4.2. Examples for the nonlinearities G_k and the reaction terms R_k .

Next we give two examples for the operators G_k .

Example 4.8. Let $g_k : \mathbb{R}^m \mapsto]0, \infty[$ be a twice continuously differentiable function and define $G_k(z)(x) = g_k(z(x))$ if $z \in (W^{1,2p})^m$ and $x \in \Omega$.

In many applications g_k depends only on one variable and is a multiple of the exponential function.

As the second example we present a nonlocal operator arising in the diffusion of bacteria; see [7], [8], and the references therein.

Example 4.9. Let η be a continuously differentiable function on \mathbb{R} which is bounded from above and below by positive constants. Assume $\varphi \in L^2(\Omega)$ and define

$$G_k(z) := \eta \left(\int_{\Omega} z_k \varphi dx \right), \quad z = (z_1, \dots, z_m) \in (W^{1,2p})^m.$$

Now we give two examples for mappings R_k .

Example 4.10. Assume that $[T_0, T_1[= \cup_{l=1}^j [t_l, t_{l+1}[$ is a (disjoint) decomposition of $[T_0, T_1[$ and let for $l \in \{1, \dots, j\}$

$$S_l : \mathbb{R}^m \times \mathbb{R}^{nm} \mapsto \mathbb{R}$$

be a function which satisfies the following condition: for any compact set $K \subset \mathbb{R}^m$ there is a constant L_K such that for any $a, \tilde{a} \in K$ and $b, \tilde{b} \in \mathbb{R}^{nm}$ the inequality

$$\begin{aligned} |S_l(a, b) - S_l(\tilde{a}, \tilde{b})| &\leq L_K |a - \tilde{a}|_{\mathbb{R}^m} (|b|_{\mathbb{R}^{nm}}^2 + |\tilde{b}|_{\mathbb{R}^{nm}}^2) \\ &\quad + L_K |b - \tilde{b}|_{\mathbb{R}^{nm}} (|b|_{\mathbb{R}^{nm}} + |\tilde{b}|_{\mathbb{R}^{nm}}) \end{aligned}$$

holds. We define a mapping $S : [T_0, T_1[\times \mathbb{R}^m \times \mathbb{R}^{nm} \mapsto \mathbb{R}$ by setting

$$S(t, a, b) := S_l(a, b) \quad \text{if } t \in [t_l, t_{l+1}[.$$

The function S defines a mapping R in the following way: if z is the restriction of an \mathbb{R}^m -valued, continuously differentiable function on \mathbb{R}^n to Ω , then we put

$$(4.1) \quad R(t, z, \nabla z)(x) = S(t, z(x), (\nabla z)(x)) \quad \text{for } x \in \Omega$$

and afterwards extend R by continuity to the whole set $[T_0, T_1[\times (W^{1,2p}(\Omega))^m$.

Example 4.11 (electrical and heat conduction). Assume that $\sigma : \mathbb{R} \mapsto]0, \infty[$ is a continuously differentiable function. Further, let $\mathcal{S} : W^{1,2p} \mapsto W^{1,2p}$ be the mapping which assigns to $z \in W^{1,2p}$ the solution φ of the elliptic problem

$$-\nabla \cdot \sigma(z) \nabla \varphi = 0$$

(including boundary conditions, see [38], [3], [6]). If one defines

$$R(z) = \sigma(z) |\nabla(\mathcal{S}(z))|^2,$$

then, under a reasonable supposition on the domain and the boundary conditions, the mapping R satisfies assumption (Ra).

5. Tools for the proof of Theorem 3.1. Let $1 < s < \infty$ and let B be a densely defined sectorial operator in a Banach space X . Let again $J =]T_0, T_1[$ for some $T_0, T_1 > 0$. We say that the linear evolution equation

$$(5.1) \quad \begin{aligned} u' + Bu &= f, \\ u(T_0) &= 0 \end{aligned}$$

admits maximal L^s -regularity on J if for any $f \in L^s(J; X)$ there exists a unique function $u \in W^{1,s}(J; X) \cap L^s(J; D(B))$ satisfying (5.1) in the L^s -sense. In that case, we write $B \in MR(s, X)$. Observe that

$$(5.2) \quad W^{1,s}(J; X) \cap L^s(J; D(B)) \hookrightarrow C(\bar{J}; X_s),$$

where X_s is the real interpolation space $(X, D(B))_{1-\frac{1}{s}, s}$. Consider now the quasi-linear problem

$$(5.3) \quad \begin{aligned} u'(t) + \mathcal{B}(t, u(t))u(t) &= F(t, u(t)), \quad t \in J, \\ u(T_0) &= u_0. \end{aligned}$$

Here $u_0 \in X_s$, $B := \mathcal{B}(T_0, u_0)$, and $\mathcal{B} : J \times X_s \rightarrow \mathcal{L}(D(B); X)$ is continuous. $F : J \times X_s \rightarrow X$ is a Carathéodory map. We assume the following Lipschitz conditions on \mathcal{B} and F .

(B) For each $R > 0$ there exists a constant $C_R > 0$, such that

$$(5.4) \quad \begin{aligned} &\|\mathcal{B}(t, u)v - \mathcal{B}(t, \tilde{u})v\|_X \\ &\leq C_R \|u - \tilde{u}\|_{X_s} \|v\|_{D(B)}, \quad t \in J, u, \tilde{u} \in X_s, \|u\|_s, \|\tilde{u}\|_s \leq R, v \in D(B). \end{aligned}$$

(F) $F(\cdot, 0) \in L^s(J; X)$ and for each $R > 0$ there is a function $\eta_R \in L^s(J)$ such that

$$(5.5) \quad \|F(t, u) - F(t, \tilde{u})\|_X \leq \eta_R(t) \|u - \tilde{u}\|_s \text{ a.a. } t \in J, u, \tilde{u} \in X_s, \|u\|_s, \|\tilde{u}\|_s \leq R.$$

Then the following existence and uniqueness result due to Clément and Li [10] and Prüss [34] holds true.

THEOREM 5.1. *Assume that (B) and (F) are satisfied and that $B := \mathcal{B}(T_0, u_0)$ has the property of maximal L^s -regularity. Then there exists $T \in]T_0, T_1[$ such that (5.3) admits a unique solution u on $I :=]T_0, T[$ satisfying*

$$u \in W^{1,s}(I; X) \cap L^s(I; D(B)).$$

In order to verify the crucial condition that $B = \mathcal{B}(T_0, u_0)$ has maximal L^s -regularity in our situation, we need the following results on traces, heat kernels, their multiplicative perturbations, and maximal L^s -regularity.

Consider a closed subspace V of $H^1(\Omega)$ which includes $H_0^1(\Omega)$. Let $\varrho \in L^\infty(\Omega, M_{n \times n})$ and assume that it is elliptic in the sense of (2.1). Define a bilinear form $a : V \times V \rightarrow \mathbb{R}$ on V by

$$a(u, v) = - \int_{\Omega} \varrho \nabla u \cdot \nabla v \, dx, \quad u, v \in V.$$

Let A be the operator associated with a in $L^2(\Omega)$ and let $(e^{tA})_{t \geq 0}$ be the semigroup on $L^2(\Omega)$ generated by A . The following result gives sufficient conditions on the subspace V such that $(e^{tA})_{t \geq 0}$ satisfies an upper Gaussian bound. More precisely, the following Proposition holds; see, [4, Ch. 4].

PROPOSITION 5.2. *Assume that V is a closed subspace of $H^1(\Omega)$ satisfying*

- (a) $H_0^1(\Omega) \subseteq V$,
- (b) V has the L^1 - H^1 extension property,
- (c) $u \in V$ implies $|u|, \inf(|u|, 1) \in V$,
- (d) $u \in V, v \in H^1(\Omega), |v| \leq u$ implies $v \in V$.

Then e^{tA} satisfies an upper Gaussian estimate, i.e.,

$$(e^{tA}f)(x) = \int_{\Omega} K_t(x, y)f(y)dy \quad \text{a.a. } x \in \Omega, f \in L^2(\Omega)$$

for some measurable function $K_t : \Omega \times \Omega \rightarrow \mathbb{R}_+$ and there exists constants $M, a > 0$ and $\omega \in \mathbb{R}$ such that

$$(5.6) \quad 0 \leq K_t(x, y) \leq \frac{M}{t^{\frac{n}{2}}} e^{-\frac{a|x-y|^2}{t}} e^{\omega t}, \quad t > 0, \text{ a.a. } x, y \in \Omega.$$

LEMMA 5.3. *Let $H_{\Gamma_N}^1(\Omega)$ be defined as above. Then $V := H_{\Gamma_N}^1(\Omega)$ satisfies the assumptions (a)–(d) of Proposition 5.2.*

Proof. We will not give a detailed proof—which is given in [4] in the case of domains with Lipschitz boundary. Let us only notice that assertion (a) is obvious and that (b) is assured by the fact that Lipschitz domains are extension domains for $H^{1,2}$; see [21, Thm. 3.10] or [32, section 1.1.16]. Inspecting the corresponding proofs (which are given via localization, Lipschitz diffeomorphism, and reflection), one recognizes that the extension mapping simultaneously extends continuously to L^1 . The assertions (c) and (d) may be proved essentially as in [4] (see Example 4.3). In fact, the continuity results for the mappings $u \rightarrow \inf(u, v)$ and $u \rightarrow |u|$ given in [31] are stated without any restriction on the regularity of the domain Ω . \square

Consider the semigroup e^{tA_k} on $L^2(\Omega)$ generated by A_k associated with the form a_k defined in (2.2) with $V = H_{\Gamma_N}^1(\Omega)$. It follows from Proposition 5.2 and Lemma 5.3 that e^{tA_k} is a positive semigroup on $L^2(\Omega)$ satisfying an upper Gaussian bound. Hence, $(e^{tA_k})_{t \geq 0}$ extends to a positive C_0 -semigroup of contractions on $L^q(\Omega)$ for all $1 \leq q < \infty$.

PROPOSITION 5.4. *Let $b \in L^\infty(\Omega, \mathbb{R})$ such that $\inf_{x \in \Omega} |b(x)| \geq \delta$ for some $\delta > 0$. Let $1 < s, q < \infty$. Then $bA_k \in MR(s, L^q(\Omega))$ for all $k \in \{1, \dots, m\}$.*

Proof. Let $k \in \{1, \dots, m\}$. By the above remark, e^{tA_k} is a positive contraction semigroup on $L^q(\Omega)$ satisfying an upper Gaussian bound. Hence, the kernel K_t of $e^{t(A_k - \alpha Id)}$ satisfies (5.6) with $\omega = 0$ for suitable $\alpha \in \mathbb{R}$. Moreover, $A_k - \alpha Id$ is self-adjoint in $L^2(\Omega)$. By a result due to Duong and Ouhabaz [13], the semigroup on

$L^2(\Omega)$ generated by $b(A_k - \alpha Id)$ satisfies an upper Gaussian bound with $\omega = 0$ as well. Thus $b(A_k - \alpha Id) \in MR(s, L^q(\Omega))$ by a result of Hieber and Prüss (see [27] or [12, Thm. 4.8]). Finally, $bA_k \in MR(s, L^q(\Omega))$ due to the lower order perturbation result of maximal regularity; see [12, Prop. 4.3]. \square

PROPOSITION 5.5. *Let $p > \frac{n}{2}$ be the number from Assumption (Op) and assume that $\theta \in (\frac{1}{2} + \frac{n}{4p}, 1]$. Then*

$$[L^p, D(A_k^p)]_\theta \hookrightarrow W_{\Gamma_N}^{1,2p}(\Omega).$$

A proof for the three-dimensional case is given in [35]; the two-dimensional case requires only obvious modifications. A complete, but technically more involved, proof for the two-dimensional case is contained in [29], see Thm. 5.2.

COROLLARY 5.6. *Let $r > \frac{4p}{2p-n}$. Then*

$$(L^p, D(A_k^p))_{1-\frac{1}{r}, r} \hookrightarrow W_{\Gamma_N}^{1,2p}(\Omega).$$

Proof. Let θ be any number from the interval $]\frac{1}{2} + \frac{n}{4p}, 1 - \frac{1}{r}[$. By interpolation

$$(L^p, D(A_k^p))_{1-\frac{1}{r}, r} \hookrightarrow (L^p, D(A_k^p))_{\theta, 1} \hookrightarrow [L^p, D(A_k^p)]_\theta.$$

Then the assertion follows from the embedding property of the complex interpolation space into $W_{\Gamma_N}^{1,2p}(\Omega)$ established in Proposition 5.5. \square

6. Proof of the main result. We first set $X := (L^p(\Omega))^m$, $\mathcal{D} := \times_{k=1}^m D(A_k^p)$, and $X_r := (X, \mathcal{D})_{1-\frac{1}{r}, r}$ for r as above. By assumption (IC), $w_0 \in X_r$. Further, for every pair $(t, z) \in [T_0, T_1] \times W^{1,2p}(\Omega)^m$ we define the mapping $H(t, z) : X \mapsto X$ via

$$(6.1) \quad \varphi := (\varphi_1, \dots, \varphi_m) \mapsto (H_1(t, z)\varphi_1, \dots, H_m(t, z)\varphi_m).$$

Since $H_k(t, z) \in L^\infty(\Omega)$ and since H_k possesses a strictly positive lower bound, it follows that

$$D(H_k(t, z)A_k^p) = D(A_k^p).$$

In particular, $D(H_k(T_0, w_0)A_k^p)$ is dense in $L^p(\Omega)$ (see [22, Thms. 4.5 and 4.7]).

Consider the mapping $\mathcal{B} : J \times X_r \rightarrow \mathcal{L}(\mathcal{D}; X)$ given by

$$\mathcal{B}(t, z)\varphi := H(t, z)(A_1^p\varphi_1, \dots, A_m^p\varphi_m), \quad \varphi = (\varphi_1, \dots, \varphi_m) \in \mathcal{D}.$$

By Corollary 5.6 and Morrey’s theorem we have

$$X_r \hookrightarrow (W_{\Gamma_N}^{1,2p}(\Omega))^m \hookrightarrow (C^\alpha(\Omega))^m$$

for some $\alpha > 0$. Thus, the assumed properties on F_k, G_k , and ϕ_k imply that

$$\mathcal{B} : J \times X_r \rightarrow \mathcal{L}(\mathcal{D}; X)$$

is continuous. Moreover, for $\beta > 0$ there exists $C_\beta > 0$ such that

$$\|H(t, z) - H(t, \tilde{z})\|_\infty \leq C_\beta \|z - \tilde{z}\|_{W^{1,2p}}$$

provided $t \in J$ and $\|z\|_{X_r}$ and $\|\tilde{z}\|_{X_r} \leq \beta$. Hence, (5.4) from assertion (B) is fulfilled.

Furthermore, (5.5) from assertion (F) holds due to the assumed properties of F_k, G_k, ϕ, R_k and Corollary 5.6. It remains to verify the key condition of Theorem 5.1, namely that $B := \mathcal{B}(T_0, w_0)$ has the property of maximal regularity. To this end, recall that $H(T_0, w_0) \in (L^\infty(\Omega))^m$ with a strictly positive lower bound in each component. Thus, $B \in MR(r, X)$ by Proposition 5.4. Finally, an application of Theorem 5.1 ends the proof of Theorem 3.1.

It remains to show that if w is a solution of (3.1), then $v := w + \phi$ provides a solution of (1.1). This will be done in the appendix. \square

We now give a proof of Corollary 3.2; in fact we prove the following sharper result.

LEMMA 6.1. *There exists $\beta > 0$ such that each component w_k of the solution w of (3.1) belongs to the space $C^\beta(]T_0, T[; W_\Gamma^{1,2p}(\Omega)) \hookrightarrow C^\beta(]T_0, T[; C^\alpha(\Omega))$.*

Proof. We write for short $D_k = D(A_k)$ and $I =]T_0, T[$. Then

$$W^{1,r}(I; L^p) \cap L^r(I; D_k) \hookrightarrow C(\bar{I}; (L^p, D_k)_{1-\frac{1}{r}, r}) \hookrightarrow C(\bar{I}; [L^p, D_k]_\theta)$$

if $\theta \in]0, 1 - \frac{1}{r}[$.

Moreover, we have the embedding

$$W^{1,r}(I; L^p) \hookrightarrow C^\delta(I; L^p) \quad \text{with} \quad \delta = 1 - \frac{1}{r}.$$

Fix $\theta \in]\frac{1}{2} + \frac{n}{4p}, 1 - \frac{1}{r}[$ and let $\lambda \in (0, 1)$ be given such that

$$\theta\lambda > \frac{1}{2} + \frac{n}{4p}.$$

In view of Proposition 5.5 and the reiteration theorem for complex interpolation (see [40, Ch. 1.9.3]) we obtain

$$\begin{aligned} \frac{\|w_k(t) - w_k(s)\|_{W^{1,2p}}}{|t - s|^{\delta(1-\lambda)}} &\leq c \frac{\|w_k(t) - w_k(s)\|_{[L^p, D_k]_{\theta\lambda}}}{|t - s|^{\delta(1-\lambda)}} \sim \frac{\|w_k(t) - w_k(s)\|_{[L^p, [L^p, D_k]_\theta]_\lambda}}{|t - s|^{\delta(1-\lambda)}} \\ &\leq \hat{c} \frac{\|w_k(t) - w_k(s)\|_{L^p}^{1-\lambda}}{|t - s|^{\delta(1-\lambda)}} \|w_k(t) - w_k(s)\|_{[L^p, D_k]_\theta}^\lambda \\ &= \hat{c} \left(\frac{\|w_k(t) - w_k(s)\|_{L^p}}{|t - s|^\delta} \right)^{1-\lambda} \left(2 \sup_{s \in I} \|w_k(s)\|_{[L^p, D_k]_\theta} \right)^\lambda. \quad \square \end{aligned}$$

7. Appendix. It remains to show that if w is a solution of (3.1), then $v := w + \phi$ provides a solution of (1.1). One easily recognizes that all the manipulations which transform (1.1) into (3.1) are straightforward to justify within the distributional calculus, except one. Therefore, we will give a strict justification of this point in the following lemma. Throughout this appendix, $f : \mathbb{R} \mapsto \mathbb{R}$ is always assumed to be twice continuously differentiable.

LEMMA 7.1. *Assume $p, r \in]1, \infty[$ and $v \in W^{1,r}(]T_0, T[; L^p) \cap C(]T_0, T[; C(\bar{\Omega}))$. Then the function $]T_0, T[\ni t \mapsto f(v(t))$ belongs to $W^{1,r}(]T_0, T[; L^p)$ and its distributional derivative is the function $]T_0, T[\ni t \mapsto f'(v(t))v'(t) \in L^r(]T_0, T[; L^p)$.*

Remark 7.2. We denote by $C^1(]T_0, T[; L^p)$ the space of all L^p -valued, continuously differentiable functions on $]T_0, T[$ with bounded derivatives on $]T_0, T[$.

In order to give a proof of Lemma 7.1 we use the following result.

LEMMA 7.3. *Let $]T_0, T[\ni t \mapsto \psi(t, \cdot)$ be a mapping belonging to $C(]T_0, T[; C(\bar{\Omega})) \cap C^1(]T_0, T[; L^p)$. Then the mapping*

$$(7.1) \quad]T_0, T[\ni t \mapsto f(\psi(t, \cdot))$$

takes its values in $C(\bar{\Omega}) \hookrightarrow L^p$. It is continuously differentiable when regarded as L^p -valued and its derivative in a point $s \in]T_0, T[$ is equal to the L^p -function $f'(\psi(s, \cdot))\psi'(s)$.

Proof. The first assertion is obvious. Concerning the second one, the set $\{\psi(t, x)/x \in \Omega, t \in [T_0, T]\}$ is bounded. Since f is twice continuously differentiable, for $s, t \in]T_0, T[$ and $x \in \Omega$ one may apply Taylor's formulae:

$$(7.2) \quad \frac{f(\psi(t, x)) - f(\psi(s, x))}{t - s} = f'(\psi(s, x)) \frac{[\psi(t, x) - \psi(s, x)]}{t - s} + \int_0^1 (1 - \tau) f''((1 - \tau)\psi(t, x) + \tau\psi(s, x)) d\tau \frac{[\psi(t, x) - \psi(s, x)]^2}{t - s}.$$

The family $\{f'(\psi(s, \cdot)) \frac{[\psi(t, \cdot) - \psi(s, \cdot)]}{t - s}\}_t$ converges by the supposition on the differentiability of the mapping $t \mapsto \psi(t, \cdot)$ in L^p to $f'(\psi(s, \cdot))\psi'(s)$ if t approaches s . It remains to show that the expression in (7.3) approaches zero in L^p . This follows easily from the uniform boundedness of the values $f''((1 - \tau)\psi(t, x) + \tau\psi(s, x))$, the boundedness of $\{\frac{[\psi(t, \cdot) - \psi(s, \cdot)]}{t - s}\}_t$ in L^p , and the convergence of $[\psi(t, \cdot) - \psi(s, \cdot)]$ to zero in $C(\bar{\Omega})$ for t approaching s . The continuity of the derivative follows from the continuity of ψ' and the continuity of the function $t \mapsto f'(\psi(t, \cdot))$ in $C(\bar{\Omega})$. \square

LEMMA 7.4. *Let $v \in W^{1,r}(]T_0, T[; L^p) \cap C(]T_0, T[; C(\bar{\Omega}))$. Then there is a sequence $\{\psi_l\}_l$ in $C(]T_0, T[; C(\bar{\Omega})) \cap C^1(]T_0, T[; L^p(\Omega))$ such that $\psi_l \mapsto v$ in $C(]T_0, T[; C(\bar{\Omega}))$ and $\psi'_l \mapsto v'$ in $L^r(]T_0, T[; L^p)$.*

Proof. Let us define a continuous extension \tilde{v} to all of \mathbb{R} which additionally has compact support as follows: we put

$$(7.4) \quad \hat{v}(t) := \begin{cases} v(T_0 + (T_0 - t)) & \text{if } t \in]T_0 - (T - T_0), T_0[, \\ v(t) & \text{if } t \in [T_0, T], \\ v(T - (t - T)) & \text{if } t \in]T, T + (T - T_0)[\end{cases}$$

(reflection at T_0, T , respectively). Afterwards we multiply \hat{v} by a real-valued, continuously differentiable function which is identical 1 on $[T_0, T]$ and which has its support in $]T_0 - (T - T_0)/2, T + (T - T_0)/2[$. We define this product as \tilde{v} and identify \tilde{v} with its extension by zero to whole \mathbb{R} . Obviously, $\tilde{v}|_{]T_0, T[} = v$; further one verifies the property $\tilde{v} \in W^{1,r}(\mathbb{R}; L^p) \cap C(\mathbb{R}; C(\bar{\Omega}))$. Let ϑ be the usual mollifier function

$$\vartheta(s) = \begin{cases} \frac{1}{\int e^{-\frac{1}{1-s^2}} ds} e^{-\frac{1}{1-s^2}} & \text{if } |s| < 1, \\ 0 & \text{else on } \mathbb{R} \end{cases}$$

and $\vartheta_l(s) := l\vartheta(ls)$. Now we put

$$(7.5) \quad \psi_l(t) := \begin{cases} \int_{T_0}^t (\tilde{v}' * \vartheta_l)(s) ds + (\tilde{v} * \vartheta_l)(T_0) & \text{if } t \geq T_0, \\ -\int_t^{T_0} (\tilde{v}' * \vartheta_l)(s) ds + (\tilde{v} * \vartheta_l)(T_0) & \text{if } t < T_0. \end{cases}$$

Then ψ_l is nothing else but $\tilde{v} * \vartheta_l$. This yields $\psi_l \mapsto v$ in $C(]T_0, T[; C(\bar{\Omega}))$. On the other hand, (7.5) immediately gives $\psi'_l = \tilde{v}' * \vartheta_l$. This means that $\psi'_l \mapsto \tilde{v}'$ in $L^r(\mathbb{R}; L^p)$, which implies $\psi'_l|_{]T_0, T[} \mapsto v'$ in $L^r(]T_0, T[; L^p)$. \square

We now turn to the proof of Lemma 7.1. Let $\{\psi_l\}_l$ be the sequence from the previous lemma and $\varphi \in C_0^\infty(]T_0, T[)$. Then, considering the function $]T_0, T[\ni t \mapsto f(v(t))$ as an L^p -valued distribution, one gets by the definition of the weak derivative

$$\begin{aligned} (f(v))'(\varphi) &= -f(v)(\varphi') = -\int_{T_0}^T f(v(s))\varphi'(s) ds = -\int_{T_0}^T \lim_{l \rightarrow \infty} f(\psi_l(s))\varphi'(s) ds \\ &= \lim_{l \rightarrow \infty} -\int_{T_0}^T f(\psi_l(s))\varphi'(s) ds. \end{aligned}$$

By Lemma 7.3, each $f(\psi_l)$ even has a strong (time) derivative which equals $f'(\psi_l)\psi'_l$. From this and integrating by parts one gets

$$-\int_{T_0}^T f(\psi_l(s))\varphi'(s) ds = \int_{T_0}^T f'(\psi_l(s))\psi'_l(s)\varphi(s)ds.$$

By construction, $\psi_l \mapsto v$ in $C([T_0, T]; C(\bar{\Omega}))$, $\psi'_l \mapsto v'$ in $L^r(]T_0, T[; L^p)$, which implies $f'(\psi_l(\cdot))\psi'_l\varphi \mapsto f'(v(\cdot))v'\varphi$ in $L^r(]T_0, T[; L^p)$. But the integral is a continuous mapping from $L^r(]T_0, T[; L^p)$ into L^p ; this finally gives

$$\begin{aligned} \int_{T_0}^T f'(v(s))v'(s)\varphi(s) ds &= \int_{T_0}^T \lim_{l \rightarrow \infty} f'(\psi_l(s))\psi'_l(s)\varphi(s)ds \\ &= \lim_{l \rightarrow \infty} \int_{T_0}^T f'(\psi_l(s))\psi'_l(s)\varphi(s)ds \\ &= \lim_{l \rightarrow \infty} - \int_{T_0}^T f(\psi_l(s))\varphi'(s) ds = (f(v))'(\varphi). \end{aligned}$$

Thus, Lemma 7.1 is proved.

Acknowledgments. The second author would like to thank K. Gröger and J. Griepentrog for stimulating discussions on the subject of the paper.

REFERENCES

- [1] H. AMANN, *Linear and Quasilinear Parabolic Problems*, Birkhäuser, Basel, Boston, Berlin, 1995.
- [2] H. AMANN, *Nonhomogeneous linear and quasilinear elliptic and parabolic boundary value problems*, in *Function Spaces, Differential Operators and Nonlinear Analysis*, Teubner-Texte Math. 133, H.-J. Schmeisser et al., eds., Stuttgart, 1993, pp. 9–126.
- [3] K. ANDREWS, P. SHI, M. SHILLOR, AND S. WRIGHT, *Thermoelastic contact with Barber's heat exchange condition*, *Appl. Math. Optim.*, 28 (1993), pp. 11–48.
- [4] W. ARENDT AND A. F. M. TER ELST, *Gaussian estimates for second order elliptic operators with boundary conditions*, *J. Operator Theory*, 38 (1997), pp. 87–130.
- [5] M. CHAPLAIN AND G. LOLAS, *Mathematical modelling of cancer cell invasion of tissue: The role of the urokinase plasminogen activation system*, *Math. Models Methods Appl. Sci.*, 15 (2005), pp. 1685–1734.
- [6] S. N. ANTONTSEV AND M. CHIPOT, *The thermistor problem: Existence, smoothness, uniqueness, blowup*, *SIAM J. Math. Anal.*, 25 (1994), pp. 1128–1156.
- [7] N. H. CHANG AND M. CHIPOT, *On some mixed boundary value problems with nonlocal diffusion*, *Adv. Math. Sci. Appl.*, 14 (2004), pp. 1–24.
- [8] M. CHIPOT AND B. LOVAT, *On the asymptotic behavior of some nonlocal problems*, *Positivity*, 3 (1999), pp. 65–81.
- [9] P. G. CIARLET, *The finite element method for elliptic problems*, *Stud. Math. Appl.*, North-Holland, Amsterdam, New York, Oxford, 1979.
- [10] P. CLÉMENT AND S. LI, *Abstract parabolic quasilinear equations and application to a groundwater flow problem*, *Adv. Math. Sci. Appl.*, 3 (1994), pp. 17–32.
- [11] M. DAUGE, *Neumann and mixed problems on curvilinear polyhedra*, *Integral Equations Operator Theory*, 15 (1992), pp. 227–261.
- [12] R. DENK, M. HIEBER, AND J. PRÜSS, *\mathcal{R} -boundedness, Fourier multipliers and problems of elliptic and parabolic type*, *Mem. Amer. Math. Soc.*, 166 (2003), pp. viii+114.
- [13] X. T. DUONG AND E. M. OUHAZ, *Complex multiplicative perturbations of elliptic operators: Heat kernel bounds and holomorphic functional calculus*, *Differential Integral Equations*, 12 (1999), pp. 395–418.
- [14] J. ELSCHNER, H. KAISER, J. REHBERG, AND G. SCHMIDT, *$W^{1,q}$ regularity results for elliptic transmission problems on heterogeneous polyhedra*, *Math. Models Methods Appl. Sci.*, 17 (2007), pp. 593–615.

- [15] J. ELSCHNER, J. REHBERG, AND G. SCHMIDT, *Optimal regularity for elliptic transmission problems including C^1 interfaces*, Interfaces Free Bound., 9 (2007), pp. 233–252.
- [16] J. ESCHER, *On quasilinear fully parabolic boundary value problems*, Differential Integral Equations, 7 (1994), pp. 1325–1343.
- [17] M. FILA AND H. MATANO, *Blow up in nonlinear heat equations from the dynamical systems point of view*, in Handbook of Dynamical Systems, Vol. 2, B. Fiedler, ed., North-Holland, Amsterdam, 2002.
- [18] J. FUHRMANN AND H. LANGMACH, *Stability and existence of solutions of time-implicit finite volume schemes for viscous nonlinear conservation laws*, Appl. Numer. Math., 37 (2001), pp. 201–230.
- [19] H. GAJEWSKI, K. GRÖGER, AND K. ZACHARIAS, *Nichtlineare Operatorgleichungen und Operatordifferentialgleichungen*, Akademie-Verlag, Berlin, 1974.
- [20] H. GAJEWSKI AND K. GRÖGER, *Reaction-diffusion processes of electrical charged species*, Math. Nachr., 177 (1996), pp. 109–130.
- [21] E. GIUSTI, *Metodi diretti nel calcolo delle variazioni*, Unione Matematica Italiana, Bologna, 1994.
- [22] J. A. GRIEPENTROG, H. C. KAISER, AND J. REHBERG, *Heat kernel and resolvent properties for second order elliptic differential operators with general boundary conditions on L^p* , Adv. Math. Sci. Appl., 11 (2001), pp. 87–112.
- [23] P. GRISVARD, *Elliptic Problems in Nonsmooth Domains*, Pitman, Boston, 1985.
- [24] K. GRÖGER, *A $W^{1,p}$ -estimate for solutions to mixed boundary value problems for second order elliptic differential equations*, Math. Ann., 283 (1989), pp. 679–687.
- [25] R. HALLER-DINTELMANN, H.-C. KAISER, AND J. REHBERG, *Elliptic Model Problems Including Mixed Boundary Conditions and Material Heterogeneities*, Preprint 1203, WIAS, Berlin 2007.
- [26] M. HEINKENSCHLOSS AND F. TROELTZSCH, *Analysis of the Lagrange-SQP-Newton method for the control of a phase field equation*, Control Cybernet., 28 (1999), pp. 177–211.
- [27] M. HIEBER AND J. PRÜSS, *Heat kernels and maximal L^p - L^q estimates for parabolic evolution equations*, Comm. Partial Differential Equations, 22 (1997), pp. 1647–1669.
- [28] D. JERISON AND C. KENIG, *The inhomogeneous Dirichlet problem in Lipschitz domains*, J. Funct. Anal., 130 (1995), pp. 161–219.
- [29] H. C. KAISER, J. REHBERG, AND H. NEIDHARDT, *Classical solutions of quasilinear parabolic systems on two dimensional domains*, NoDEA Nonlinear Differential Equation Appl., 13 (2006), pp. 287–310.
- [30] P. KREJCI, E. ROCCA, AND J. SPREKELS, *Nonlocal Temperature-Dependent Phase-Field Models for Non-Isothermal Phase Transitions*, Preprint 1006, WIAS, Berlin, 2007.
- [31] M. MARCUS AND V. MIZEL, *Every superposition operator mapping one Sobolev space into another is continuous*, J. Funct. Anal., 33 (1979), pp. 217–229.
- [32] V. MAZ'YA, *Sobolev Spaces*, Springer, Berlin, Heidelberg, New York, Tokyo, 1985.
- [33] N. G. MEYERS, *An L^p -estimate for the gradient of solutions of second order elliptic divergence equations*, Ann. Sc. Norm. Super. Pisa Cl. Sci. (4), 17 (1963), pp. 189–206.
- [34] J. PRÜSS, *Maximal regularity for evolution equations in L^p -spaces*, Conf. Semin. Math. Univ. Bari, 285 (2002), pp. 1–39.
- [35] J. REHBERG, *Quasilinear parabolic equations in L^p* , in Nonlinear Elliptic and Parabolic Problems. A special tribute to the work of Herbert Amann, Prog. Nonlinear Differential Equations Appl., M. Chipot and J. Escher, eds., 2005, pp. 413–419.
- [36] S. SELBERHERR, *Analysis and Simulation of Semiconductors*, Springer, Wien, 1984
- [37] E. SHAMIR, *Regularization of mixed second-order elliptic problems*, Israel J. Math. 6, (1968), pp. 150–168.
- [38] P. SHI, M. SHILLOR, AND X. XU, *Existence of a solution to the Stefan problem with Joule's heating*, J. Differential Equations, 105 (1993), pp. 239–263.
- [39] A. SOMMERFELD, *Thermodynamics and statistical mechanics*, Lectures on Theoretical Physics V, Academic Press, New York, 1956.
- [40] H. TRIEBEL, *Interpolation Theory, Function Spaces, Differential Operators*, North-Holland, Amsterdam, New York, Oxford, 1978.
- [41] A. UNGER AND F. TROELTZSCH, *Fast solutions of optimal control problems in the selective cooling of steel*, ZAMM Z. Angew. Math. Mech., 81 (2001), pp. 447–456.
- [42] D. ZANGER, *The inhomogeneous Neumann problem in Lipschitz domains*, Comm. Partial Differential Equations, 25 (2000), pp. 1771–1808.

A VARIATIONAL INEQUALITY ARISING FROM EUROPEAN INSTALLMENT CALL OPTIONS PRICING*

FAHUAI YI[†], ZHOU YANG[†], AND XIAOHUA WANG[‡]

Abstract. In this paper we consider a parabolic variational inequality arising from European continuous installment call options pricing and prove the existence and uniqueness of the solution to the problem. Moreover, we obtain C^∞ regularity and the bounds of the free boundary, as well as the limit of the free boundary as $\tau = T - t \rightarrow +\infty$. Eventually we show its numerical result by the binomial method.

Key words. free boundary, variational inequality, option pricing, European installment call options

AMS subject classification. 35R35

DOI. 10.1137/060670353

1. Introduction. In this paper we consider a parabolic variational inequality arising from the model of European continuous installment call options pricing. More precisely, we will find $C(S, t)$ satisfying

$$(1.1) \quad \begin{cases} \partial_t C + \frac{\sigma^2}{2} S^2 \partial_{SS} C + (r - q) S \partial_S C - rC = L^* & \text{if } C > 0 \text{ and } (S, t) \in (0, +\infty) \times (0, T], \\ \partial_t C + \frac{\sigma^2}{2} S^2 \partial_{SS} C + (r - q) S \partial_S C - rC \leq L^* & \text{if } C = 0 \text{ and } (S, t) \in (0, +\infty) \times (0, T], \\ C(S, T) = (S - K)^+, & S \in [0, +\infty), \end{cases}$$

where σ , r , L^* , and K are positive constants and q is a nonnegative constant.

In the appendix we present the financial and stochastic background of this problem.

If $L^* = 0$ in problem (1.1), it is standard European call option, which has an explicit analytic formula of solution and does not have free boundary at all (see [11]). We will find that the case $L^* > 0$ is more complicated than the case $L^* = 0$.

There are some papers in the field of install options, such as [7], [8], [1], [6], in which the authors developed the models and numerical analysis. Particularly, Alobaidi, Mallier, and Deakin showed the behavior of the free boundary $S_s(\tau)$ in (1.1) close to expire by Laplace transforms in [1], that is,

$$(1.2) \quad S_s(\tau) \sim K \exp\{-\sigma(-\tau \ln \tau)^{1/2}(1 + o(1))\} \quad \text{as } \tau = T - t \rightarrow 0^+.$$

*Received by the editors September 20, 2006; accepted for publication (in revised form) September 17, 2007; published electronically April 23, 2008. The project supported by National Natural Science Foundation of China (10671075), National Natural Science Foundation of Guangdong province (5005930), and University Special Research Fund for Ph.D. Program (20060574002).

<http://www.siam.org/journals/sima/40-1/67035.html>

[†]School of Mathematical Sciences, South China Normal University, Guangzhou 510631, China (fhyi@scnu.edu.cn, yangzhou1975@yahoo.com).

[‡]School of Mathematical Sciences, South China University of Technology, Guangzhou 510640, China.

Since (1.1) is a degenerate backward parabolic problem, we transform it into a familiar forward nondegenerate parabolic variational inequality problem. Thus, letting

$$(1.3) \quad V(x, \tau) = C(S, t)/K, \quad \tau = T - t, \quad x = \ln(S/K), \quad L = L^*/K,$$

then we have

$$(1.4) \quad \begin{cases} \partial_\tau V - \mathcal{L}_x V = -L & \text{if } V > 0 \text{ and } (x, \tau) \in \mathbb{R} \times (0, T], \\ \partial_\tau V - \mathcal{L}_x V \geq -L & \text{if } V = 0 \text{ and } (x, \tau) \in \mathbb{R} \times (0, T], \\ V(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}, \end{cases}$$

where

$$(1.5) \quad \mathcal{L}_x V = \frac{\sigma^2}{2} \partial_{xx} V + \left(r - q - \frac{\sigma^2}{2} \right) \partial_x V - rV.$$

We focus our attention on the monotonicity, regularity, and bounds of the free boundary as well as the limit of the free boundary as $\tau \rightarrow +\infty$ in problem (1.4).

As we know, the behavior $\partial_\tau V \geq 0$ is quite important for investigation of the regularity properties of the free boundary. Based on this behavior it can be deduced that $\partial_\tau V$ is continuous across the free boundary and C^∞ regularity of the free boundary follows (for the one-dimensional case, see [9], and for higher dimensions, see [4]). In the absence of condition $\partial_\tau V \geq 0$, the latest development is that $\partial_\tau V$ is continuous for almost all time (see [3]) and the free boundary possesses C^∞ regularity locally around some points which are energetically characterized (see [5]). Those results manifest that the analysis for regularity of the free boundary is not thoroughly solved.

In our case $V(x, 0) = e^x - 1$ if $x \geq 0$; combining Lemma 4.4 in this paper, we see that

$$\partial_\tau V(x, 0) = \mathcal{L}_x(e^x - 1) - L = r - qe^x - L \quad \text{if } x > 0.$$

It is clear that $\partial_\tau V(x, 0) < 0$ if $q > 0$ and x is large enough.

In the next section, we will construct a transformation (2.1) and deduce a new unknown function $v(y, \tau)$ satisfying the variational inequality (2.2), which is important to the property $\partial_\tau v \geq 0$; then we prove the existence and uniqueness of the $W_{p,loc}^{2,1}$ solution to the new parabolic variational inequality (2.2). The main work is in sections 3 and 4. In section 3 we prove that the new free boundary is monotonic and C^∞ -smooth based on the results in section 2. Moreover, we will show the starting point, the bounds of the free boundary, and the limit behavior of the free boundary as $\tau \rightarrow +\infty$. In section 4, we come back to consider the behavior of the free boundary of problem (1.4): it is C^∞ -smooth on $(0, +\infty)$, and its monotonicity depends on the relationship of parameters σ , r , q , and L in the problem. In section 5, we focus our attention on the monotonicity of free boundary with respect to the parameters r , q , and L . In the last section, we provide numerical result applying the binomial method.

2. Existence and uniqueness of $W_{p,loc}^{2,1}$ solution of problem (1.4). As mentioned in the previous section, applying the transformation

$$(2.1) \quad y = x + (r - q - p^*)\tau, \quad p^* = \max\{r, L\}, \quad v(y, \tau) = V(x, \tau),$$

problem (1.4) then becomes

$$(2.2) \quad \begin{cases} \partial_\tau v - \mathcal{L}_y v = -L & \text{if } v > 0 \text{ and } (y, \tau) \in \mathbb{R} \times (0, T], \\ \partial_\tau v - \mathcal{L}_y v \geq -L & \text{if } v = 0 \text{ and } (y, \tau) \in \mathbb{R} \times (0, T], \\ v(y, 0) = (e^y - 1)^+, & y \in \mathbb{R}, \end{cases}$$

where

$$(2.3) \quad \mathcal{L}_y v = \frac{\sigma^2}{2} \partial_{yy} v + \left(p^* - \frac{\sigma^2}{2} \right) \partial_y v - rv.$$

Since the problem lies in the unbounded domain $\Omega_T = \mathbb{R} \times (0, T]$, we first consider the problem in the bounded domain $\Omega_T^n = (-n, n) \times (0, T]$, $n \in \mathbb{Z}^+$:

$$(2.4) \quad \begin{cases} \partial_\tau v_n - \mathcal{L}_y v_n = -L & \text{if } v_n > 0 \text{ and } (y, \tau) \in \Omega_T^n, \\ \partial_\tau v_n - \mathcal{L}_y v_n \geq -L & \text{if } v_n = 0 \text{ and } (y, \tau) \in \Omega_T^n, \\ v_n(-n, \tau) = 0, \quad \partial_y v_n(n, \tau) = e^{n+(p^*-r)\tau}, & \tau \in [0, T], \\ v_n(y, 0) = (e^y - 1)^+, & y \in [-n, n]. \end{cases}$$

LEMMA 2.1. *For any fixed $n \in \mathbb{Z}^+$, there exists a unique solution $v_n \in C(\overline{\Omega_T^n}) \cap W_p^{2,1}(\Omega_T^n \setminus B_\delta(P_0))$ to problem (2.4), where $\forall 1 < p < +\infty$, and $\delta > 0$, $P_0 = (0, 0)$, $B_\delta(P_0) = \{(y, \tau) : y^2 + \tau^2 \leq \delta^2\}$. Moreover, if n is large enough, we have*

$$(2.5) \quad (e^y - 1)^+ \leq v_n \leq e^{y+(p^*-r)\tau},$$

$$(2.6) \quad \partial_\tau v_n \geq 0,$$

$$(2.7) \quad \partial_y v_n \geq 0.$$

Proof. As usual we define a penalty function $\beta_\varepsilon(t)$ (see Figure 1), which satisfies

$$\varepsilon > 0 \text{ and small enough, } \beta_\varepsilon(t) \in C^\infty(-\infty, +\infty),$$

$$\beta_\varepsilon(t) \leq 0, \quad 0 \leq \beta'_\varepsilon(t) \leq 2L/\varepsilon, \quad \beta''_\varepsilon \leq 0,$$

and

$$(2.8) \quad \beta_\varepsilon(t) = \begin{cases} 0, & t \geq 2\varepsilon, \\ \frac{2L}{\varepsilon}t - 3L - r\varepsilon, & t \leq 3\varepsilon/2. \end{cases}$$

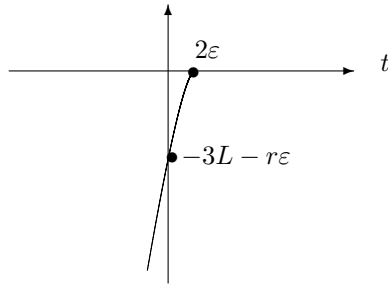


FIG. 1.

Then

$$(2.9) \quad \beta_\varepsilon(\varepsilon) = -L - r\varepsilon, \quad \beta_\varepsilon(0) = -3L - r\varepsilon.$$

Since $(e^y - 1)^+$ is not smooth enough, we need to smooth it. Define $\pi_\varepsilon(t)$ (see Figure 2):

$$(2.10) \quad \pi_\varepsilon(t) = \begin{cases} t, & t \geq \varepsilon, \\ 0, & t \leq -\varepsilon, \end{cases}$$

$$\pi_\varepsilon(t) \in C^\infty, \quad 0 \leq \pi'_\varepsilon(t) \leq 1, \quad \pi''_\varepsilon(t) \geq 0, \quad \lim_{\varepsilon \rightarrow 0^+} \pi_\varepsilon(t) = t^+.$$

Following the idea in [10], construct an approximation of problem (2.4):

$$(2.11) \quad \begin{cases} \partial_\tau v_{\varepsilon,n} - \mathcal{L}_y v_{\varepsilon,n} + \beta_\varepsilon(v_{\varepsilon,n}) = -L & \text{in } \Omega_T^n, \\ v_{\varepsilon,n}(-n, \tau) = 0, \quad \partial_y v_{\varepsilon,n}(n, \tau) = e^{n+(p^*-r)\tau}, \quad \tau \in [0, T], \\ v_{\varepsilon,n}(y, 0) = \pi_\varepsilon(e^y - 1), & y \in [-n, n]. \end{cases}$$

Applying Schauder's fixed point theorem, it is not difficult to get the existence of the $W_p^{2,1}$ solution to problem (2.11). The proof of uniqueness is standard as well, so we omit the details.

If we can prove that, as ε is small enough,

$$(2.12) \quad \pi_\varepsilon(e^y - 1) \leq v_{\varepsilon,n} \leq e^{y+(p^*-r)\tau},$$

then

$$-3L - r\varepsilon \leq \beta_\varepsilon(v_{\varepsilon,n}) \leq 0.$$

This means that $\beta_\varepsilon(v_{\varepsilon,n})$ is a bounded function and its bound is independent of ε . It is deduced that, by $W_p^{2,1}$ and C^α ($0 < \alpha < 1$) estimates of the parabolic problem,

$$\begin{aligned} |v_{\varepsilon,n}|_{W_p^{2,1}(\Omega_T^n \setminus B_\delta(P_0))} &\leq C, \\ |v_{\varepsilon,n}|_{C^\alpha(\overline{\Omega_T^n})} &\leq C, \end{aligned}$$

where C is independent of ε . It is not difficult to derive that, as $\varepsilon \rightarrow 0^+$,

$$v_{\varepsilon,n} \rightarrow v_n \text{ in } W_p^{2,1}(\Omega_T^n \setminus B_\delta(P_0)) \text{ weakly} \quad \text{and} \quad v_{\varepsilon,n} \rightarrow v_n \text{ in } C(\overline{\Omega_T^n}),$$

where v_n is the solution of problem (2.4).

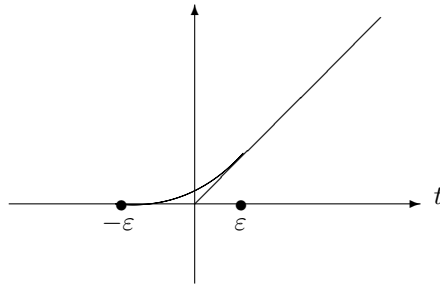


FIG. 2.

Next we prove (2.12). In the first, if ε is small enough, the properties of β_ε , π_ε , and (2.9) imply that

$$\begin{aligned} & \partial_\tau \pi_\varepsilon(e^y - 1) - \mathcal{L}_y \pi_\varepsilon(e^y - 1) + \beta_\varepsilon(\pi_\varepsilon(e^y - 1)) + L \\ &= -\frac{\sigma^2}{2} \pi_\varepsilon''(e^y - 1) e^{2y} - \max\{r, L\} \pi_\varepsilon'(e^y - 1) e^y + r \pi_\varepsilon(e^y - 1) + \beta_\varepsilon(\pi_\varepsilon(e^y - 1)) + L \\ &\leq \begin{cases} r \pi_\varepsilon(e^y - 1) + \beta_\varepsilon(\pi_\varepsilon(e^y - 1)) + L \leq r\varepsilon + \beta_\varepsilon(\varepsilon) + L \leq 0, & e^y - 1 \leq \varepsilon, \\ -\max\{r, L\} e^y + r(e^y - 1) + L \leq 0, & e^y - 1 > \varepsilon. \end{cases} \end{aligned}$$

Furthermore, from the initial and boundary conditions in (2.11), we deduce that if ε is small enough,

$$\begin{cases} \pi_\varepsilon(e^y - 1) = 0 = v_{\varepsilon, n}(y, \tau), & y = -n, \\ \partial_y \pi_\varepsilon(e^y - 1) = \partial_y(e^y - 1) = e^n \leq e^{n+(p^*-r)\tau} = \partial_y v_{\varepsilon, n}(y, \tau), & y = n, \\ \pi_\varepsilon(e^y - 1) = v_{\varepsilon, n}(y, 0), & \tau = 0. \end{cases}$$

From the comparison principle, we deduce $\pi_\varepsilon(e^y - 1) \leq v_{\varepsilon, n}$.

On the other hand, if $2\varepsilon < e^{-n}$, it is not difficult to check $\beta_\varepsilon(e^{y+(p^*-r)\tau}) \geq \beta_\varepsilon(e^{-n}) \geq \beta_\varepsilon(2\varepsilon) = 0$; hence, from (2.3),

$$\begin{aligned} & \partial_\tau e^{y+(p^*-r)\tau} - \mathcal{L}_y e^{y+(p^*-r)\tau} + \beta_\varepsilon(e^{y+(p^*-r)\tau}) + L \\ & \geq e^{y+(p^*-r)\tau} \left(p^* - r - \frac{\sigma^2}{2} - \left(p^* - \frac{\sigma^2}{2} \right) + r \right) + 0 + L = L \geq 0 \quad \forall (y, \tau) \in \Omega_T^n. \end{aligned}$$

Moreover, if $2\varepsilon < e^{-n}$, there holds

$$e^y = \pi_\varepsilon(e^y) \geq \pi_\varepsilon(e^y - 1).$$

From the initial and boundary conditions in (2.11), we see that if $2\varepsilon < e^{-n}$,

$$\begin{cases} e^{y+(p^*-r)\tau} \geq 0 = v_{\varepsilon, n}(y, \tau), & y = -n, \\ \partial_y e^{y+(p^*-r)\tau} = e^{n+(p^*-r)\tau} = \partial_y v_{\varepsilon, n}(y, \tau), & y = n, \\ e^{y+(p^*-r)\tau} = e^y \geq \pi_\varepsilon(e^y - 1) = v_{\varepsilon, n}(y, 0), & \tau = 0. \end{cases}$$

Then the comparison principle implies $e^{y+(p^*-r)\tau} \geq v_{\varepsilon, n}$; hence, we obtain (2.12), and (2.5) is a consequence of (2.12).

In the following, we prove (2.6). In fact, for any small $\delta > 0$, $v_{\varepsilon, n}(x, \tau + \delta)$ satisfies, by (2.11),

$$\begin{cases} \partial_\tau v_{\varepsilon, n}(y, \tau + \delta) - \mathcal{L}_y v_{\varepsilon, n}(y, \tau + \delta) + \beta_\varepsilon(v_{\varepsilon, n}(y, \tau + \delta)) = -L, & (y, \tau) \in (-n, n) \times (0, T - \delta], \\ v_{\varepsilon, n}(-n, \tau + \delta) = 0 = v_{\varepsilon, n}(-n, \tau), & \tau \in [0, T - \delta], \\ \partial_y v_{\varepsilon, n}(n, \tau + \delta) = e^{n+(p^*-r)(\tau+\delta)} \geq \partial_y v_{\varepsilon, n}(n, \tau), & \tau \in [0, T - \delta], \\ v_{\varepsilon, n}(y, 0 + \delta) \geq \pi_\varepsilon(e^y - 1) = v_{\varepsilon, n}(y, 0), & y \in [-n, n]. \end{cases}$$

Applying the comparison principle for solutions of PDEs with respect to initial and boundary values, we obtain

$$v_{\varepsilon, n}(y, \tau + \delta) \geq v_{\varepsilon, n}(y, \tau) \quad \forall (y, \tau) \in (-n, n) \times [0, T - \delta] \quad \text{and} \quad \partial_\tau v_{\varepsilon, n} \geq 0.$$

Take $\varepsilon \rightarrow 0^+$ in the above inequalities. Then (2.6) follows.

For proving (2.7), derive (2.11) with respect to y , and denote $W = \partial_y v_{\varepsilon, n}$; then

$$(2.13) \quad \begin{cases} \partial_\tau W - \mathcal{L}_y W + \beta'_\varepsilon(v_{\varepsilon, n})W = 0, & (y, \tau) \in \Omega_T^n, \\ W(-n, \tau) = \partial_y v_{\varepsilon, n}(-n, \tau), & \tau \in [0, T], \\ W(n, \tau) = e^{n+(p^*-r)\tau} \geq 0, & \tau \in [0, T], \\ W(y, 0) = \pi'_\varepsilon(e^y - 1)e^y \geq 0, & y \in [-n, n]. \end{cases}$$

Since $v_{\varepsilon, n}(x, \tau)$ achieves its minimum 0 at each point of $x = -n$, so $\partial_y v_{\varepsilon, n}(-n, \tau) \geq 0$, then the minimum principle implies

$$(2.14) \quad \partial_y v_{\varepsilon, n} = W \geq 0.$$

Take $\varepsilon \rightarrow 0^+$. Then we obtain (2.7).

At last, we prove uniqueness. Suppose v_1 and v_2 are two $W_{p, loc}^{2,1}(\Omega_T^n) \cap C(\overline{\Omega_T^n})$ solutions to problem (2.4), and denote

$$\mathcal{N} = \{(y, \tau) : v_1(y, \tau) < v_2(y, \tau), -n < y < n, 0 < \tau \leq T\}.$$

Suppose it is not empty; then if $(y, \tau) \in \mathcal{N}$,

$$v_2(y, \tau) > 0, \quad \partial_\tau v_2 - \mathcal{L}_y v_2 = -L.$$

Denote $W = v_2 - v_1$. Then W satisfies

$$\begin{cases} \partial_\tau W - \mathcal{L}_y W \leq 0, & (y, \tau) \in \mathcal{N}, \\ W(y, 0) = 0 & \text{on } \partial_p \mathcal{N} \setminus (\{n\} \times [0, T]), \\ \partial_y W(y, 0) = 0 & \text{on } \partial_p \mathcal{N} \cap (\{n\} \times [0, T]), \end{cases}$$

where $\partial_p \mathcal{N}$ is the parabolic boundary of the domain \mathcal{N} . Applying the A-B-P maximum principle (see [13]), we have $W \leq 0$ in \mathcal{N} , which contradicts the definition of \mathcal{N} . \square

THEOREM 2.2. *There exists a unique solution $v \in C(\overline{\Omega_T}) \cap W_p^{2,1}(\Omega_T^R \setminus B_\delta(P_0))$ to problem (2.2), where $\forall R > 0, \delta > 0, 1 < p < +\infty$. And $\partial_y v \in C(\Omega_T)$,*

$$(2.15) \quad (e^y - 1)^+ \leq v \leq e^{y+(p^*-r)\tau},$$

$$(2.16) \quad \partial_\tau v \geq 0,$$

$$(2.17) \quad \partial_y v \geq 0.$$

Proof. Applying

$$(\partial_\tau - \mathcal{L}_y)0 = 0,$$

we rewrite problem (2.4) as

$$(2.18) \quad \begin{cases} \partial_\tau v_n - \mathcal{L}_y v_n = f(y, \tau) & \text{in } \Omega_T^n, \\ v_n(-n, \tau) = 0, \quad \partial_y v_n(n, \tau) = e^{y+(p^*-r)\tau}, & \tau \in [0, T], \\ v_n(y, 0) = (e^y - 1)^+, & y \in [-n, n]. \end{cases}$$

$v_n \in W_{p, loc}^{2,1}(\Omega_T^n)$ implies $f(y, \tau) \in L_{loc}^p(\Omega_T^n)$ and

$$(2.19) \quad f(y, \tau) = I_{\{v_n > 0\}}(-L) \quad \text{a.e. in } \Omega_T^n,$$

where I_A denotes the indicator function of the set A .

Hence, for any fixed $R > \delta > 0$, if $n > R$, combining (2.5), we have the following $W_p^{2,1}$ uniform interior estimates in the domain $\Omega_T^R \setminus B_\delta(P_0)$:

$$\begin{aligned} & \|v_n\|_{W_p^{2,1}(\Omega_T^R \setminus B_\delta(P_0))} \\ & \leq C(\|v_n\|_{L^\infty(\Omega_T^R)} + \|(e^y - 1)^+\|_{C^2([-R, -\delta] \cup [\delta, R])} + \|f(y, \tau)\|_{L^\infty(\Omega_T^R)}) \leq C, \end{aligned}$$

where C depends on R but is independent of n . Let $n \rightarrow \infty$; then we have, possibly a subsequence,

$$v_n \rightharpoonup v_R \text{ in } W_p^{2,1}(\Omega_T^R \setminus B_\delta(P_0)) \text{ weakly as } n \rightarrow +\infty.$$

Moreover, the Sobolev imbedding theorem implies

$$\partial_y v_n \rightarrow \partial_y v_R \text{ in } C(\Omega_T^R \setminus B_\delta(P_0)) \text{ as } n \rightarrow +\infty.$$

Define $v = v_R$ if $y \in [-R, R]$. It is clear that v is reasonably defined and v is the solution of problem (2.2); moreover, $\partial_y v \in C(\Omega_T)$ and the C^α estimate implies $v \in C(\overline{\Omega_T})$.

Inequalities (2.15), (2.16), (2.17) are a consequence of (2.5), (2.6), (2.7), respectively. The proof of the uniqueness is the same as in the proof of Lemma 2.1. \square

From the transformation (2.1), it is not difficult to see the following.

THEOREM 2.3. *There exists a unique solution $V \in C(\overline{\Omega_T}) \cap W_p^{2,1}(\Omega_T^R \setminus B_\delta(P_0))$ to problem (1.4), where $\forall R > 0, \delta > 0, 1 < p < +\infty$. And $\partial_x V \in C(\Omega_T)$,*

$$(2.20) \quad (e^{x+(r-q-p^*)\tau} - 1)^+ \leq V \leq e^{x-q\tau},$$

$$(2.21) \quad \partial_x V \geq 0,$$

$$(2.22) \quad \partial_\tau V \geq (r - q - p^*) \partial_x V.$$

3. Characterizations of the free boundary of problem (2.2). In this section, we consider the behavior of the free boundary of problem (2.2). Denote

$$\begin{aligned} \mathbf{NR}_y &= \{(y, \tau) : v(y, \tau) > 0\} && \text{(nontransaction region),} \\ \mathbf{SR}_y &= \{(y, \tau) : v(y, \tau) = 0\} && \text{(stop region).} \end{aligned}$$

Applying (2.17), v is monotonic increasing with respect to y , so we can define the free boundary

$$y_s(\tau) = \sup\{y : v(y, \tau) = 0\}, \quad \tau > 0, \quad \text{between } \mathbf{NR}_y \text{ and } \mathbf{SR}_y.$$

THEOREM 3.1. $y_s(\tau) \in C[0, T] \cap C^\infty(0, T]$ and is strictly decreasing with $y_s(0) = 0$ (see Figure 3).

Proof. We divide the proof into four steps.

Step 1. From (2.16) and (2.17), we see that $y_s(\tau)$ is decreasing in $[0, T]$.

Step 2. Prove $y_s(\tau)$ is continuous in $[0, T]$ and $y_s(0) = 0$.

In the first we prove $y_s(\tau)$ is continuous in $[0, T]$. Otherwise, there exists a domain $(y_1, y_2) \times (\tau_0, T)$ ($y_1 < y_2, 0 \leq \tau_0 < T$) such that

$$\begin{cases} \partial_\tau v - \mathcal{L}_y v = -L, & (y, \tau) \in (y_1, y_2) \times (\tau_0, T), \\ v(y, \tau_0) = 0, & y \in (y_1, y_2). \end{cases}$$

Then we have $\partial_\tau v(y, \tau_0) = -L < 0$ for any $y_1 < y < y_2$, which contradicts (2.16).

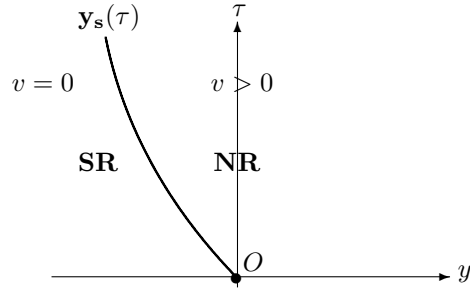


FIG. 3.

In the same way we can prove $y_s(0) \geq 0$; moreover, since $v(y, 0) = (e^y - 1)^+ > 0$ if $y > 0$, we have $y_s(0) = 0$.

Step 3. Prove $y_s(\tau)$ is strictly decreasing in $[0, T]$.

Otherwise there exists a domain $(y_0, 1) \times (\tau_1, \tau_2)$ ($y_0 \leq 0, 0 < \tau_1 < \tau_2 \leq T$) such that

$$(3.1) \quad \begin{cases} \partial_\tau v - \mathcal{L}_y v = -L, & (x, \tau) \in (y_0, 1) \times (\tau_1, \tau_2), \\ v(y_0, \tau) = 0, & \tau \in (\tau_1, \tau_2). \end{cases}$$

Then $W = \partial_\tau v$ satisfies

$$\begin{cases} \partial_\tau W - \mathcal{L}_y W = 0, & (x, \tau) \in (y_0, 1) \times (\tau_1, \tau_2), \\ W(y_0, \tau) = 0, & \tau \in (\tau_1, \tau_2). \end{cases}$$

Since $W = \partial_\tau v \geq 0$, W achieves its nonpositive minimum at $y = y_0$. Applying the maximum principle, we have $\partial_y W(y_0, \tau) > 0$ for any $\tau \in (\tau_1, \tau_2)$. On the other hand, we can deduce that $\partial_y v(y_0, \tau) = 0$ for any $\tau \in (\tau_1, \tau_2)$ by $\partial_y v \in C(\Omega_T)$. So, $\partial_y W(y_0, \tau) = \partial_{\tau y} v(y_0, \tau) = 0$ for any $\tau \in (\tau_1, \tau_2)$; thus we get a contradiction. Hence $y_s(\tau)$ is strictly decreasing in $(0, T]$.

Step 4. Since $\partial_\tau v \geq 0$ and 0 is the lower obstacle, it can be proved that $y_s(\tau) \in C^{0,1}(0, T]$ by the method developed by Friedman in [9]. Moreover $y_s(\tau) \in C^\infty(0, T]$ by a bootstrap argument. \square

THEOREM 3.2. *The free boundary $y_s(\tau)$ satisfies*

$$(3.2) \quad \ln \left[\frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + r\tau} \right) \right] \leq y_s(\tau) + (p^* - r)\tau \leq \ln \left[\frac{2L}{2r + \sigma^2} \left(1 + \frac{1}{L\tau} \right) \right].$$

Proof. We divide the proof into four steps.

Step 1. In the first, we give a transformation

$$(3.3) \quad z = y + (p^* - r)\tau, \quad v^*(z, \tau) = v(y, \tau) = V(x, \tau).$$

Then we have

$$(3.4) \quad \begin{cases} \partial_\tau v^* - \mathcal{L}_z v^* = -L & \text{if } v^* > 0 \text{ and } (z, \tau) \in \mathbb{R} \times (0, T], \\ \partial_\tau v^* - \mathcal{L}_z v^* \geq -L & \text{if } v^* = 0 \text{ and } (z, \tau) \in \mathbb{R} \times (0, T], \\ v^*(z, 0) = (e^z - 1)^+, & z \in \mathbb{R}, \end{cases}$$

where

$$(3.5) \quad \mathcal{L}_z v^* = \frac{\sigma^2}{2} \partial_{zz} v + \left(r - \frac{\sigma^2}{2} \right) \partial_z v - rv.$$

Denote the free boundary of problem (3.4) as $z_s(\tau)$, which is the counterpart of $y_s(\tau)$ and $z_s(\tau) = y_s(\tau) + (p^* - r)\tau$.

Step 2. We prove that for any $\tau > 0$, there holds

$$(3.6) \quad z_s(\tau) \geq z_1(\tau) = \ln \left[\frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + r\tau} \right) \right].$$

In fact, for any $T_0 > 0$, we will prove $z_s(T_0) \geq z_1(T_0)$.

Since $z_1(T_0) = \ln \left[\frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + rT_0} \right) \right]$, we define z_1 by

$$e^{z_1} = \frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + rT_0} \right).$$

It follows that

$$(3.7) \quad e^{z_1} = \frac{2L}{2r + \sigma^2} \frac{rT_0}{1 + rT_0} \leq \frac{LT_0}{1 + rT_0}.$$

We define

$$(3.8) \quad W(z, \tau) = \begin{cases} \frac{L}{1 + rT_0} (T_0 - \tau), & (z, \tau) \in (-\infty, z_1] \times [0, T_0], \\ e^z - \frac{1}{\alpha} e^{\alpha(z - z_1) + z_1} - \frac{L}{1 + rT_0} \tau, & (z, \tau) \in (z_1, +\infty) \times [0, T_0], \end{cases}$$

where $\alpha = -\frac{2r}{\sigma^2}$.

We claim that $W(x, \tau)$ possesses the following four properties:

- (I) $W \in W_{p,loc}^{2,1}(\mathbb{R} \times (0, T_0)) \cap C(\mathbb{R} \times [0, T_0])$.
- (II) $W \geq 0 \quad \forall (z, \tau) \in \mathbb{R} \times (0, T_0]$.
- (III) $W(z, 0) \geq (e^z - 1)^+$.
- (IV) $\partial_\tau W - \mathcal{L}_z W \geq -L, (z, \tau) \in \mathbb{R} \times [0, T_0]$.

Indeed, since

$$(3.9) \quad W(z_1 + 0, \tau) = e^{z_1} \left(1 - \frac{1}{\alpha} \right) - \frac{L}{1 + rT_0} \tau = \frac{L}{1 + rT_0} (T_0 - \tau),$$

$$(3.10) \quad \partial_z W(z_1 + 0, \tau) = e^{z_1} - e^{z_1} = 0.$$

Then property (I) follows. Next, we prove properties (II) and (III).

It is clear that $W(z, \tau) \geq 0$ for $(z, \tau) \in (-\infty, z_1] \times [0, T_0]$. Moreover, it can be seen that if $z > z_1$,

$$\partial_{zz} W = e^z - \alpha e^{\alpha(z - z_1) + z_1} \geq e^z > 0.$$

Combining (3.9) and (3.10), we know that for any $(z, \tau) \in (z_1, +\infty) \times [0, T_0]$, there hold

$$\partial_z W(z, \tau) > 0, \quad W(z, \tau) > \frac{L}{1 + rT_0} (T_0 - \tau) \geq 0.$$

Thus we get property (II). On the other hand, applying (3.7), we obtain

$$W(z, 0) = \begin{cases} \frac{LT_0}{1+rT_0} \geq e^{z_1} > (e^z - 1)^+, & z \leq z_1, \\ e^z - \frac{1}{\alpha} e^{\alpha(z-z_1)+z_1} > e^z > (e^z - 1)^+, & z > z_1. \end{cases}$$

Then we have property (III). In the following, we prove property (IV).

In fact, notice that 1 and α are the characteristic roots of the ODE $\mathcal{L}_z W = 0$. Then we see that for any $(z, \tau) \in (z_1, +\infty) \times (0, T_0]$,

$$\partial_\tau W - \mathcal{L}_z W = -\frac{L}{1+rT_0} - \frac{rL}{1+rT_0} \tau = -L \frac{1+r\tau}{1+rT_0} \geq -L.$$

On the other hand, for any $(z, \tau) \in (-\infty, z_1) \times (0, T_0]$, as above, we have

$$\partial_\tau W - \mathcal{L}_z W = -\frac{L}{1+rT_0} + \frac{rL}{1+rT_0} (T_0 - \tau) \geq -\frac{L}{1+rT_0} - \frac{rL}{1+rT_0} \tau \geq -L.$$

Hence we have property (IV).

Properties (I)–(IV) indicate that W is a supersolution of problem (3.4). Hence $W \geq v^*$ in $\mathbb{R} \times [0, T]$, especially

$$0 \leq v^*(z, T_0) \leq W(z, T_0) = 0, \quad z \leq z_1,$$

meaning that $z_s(T_0) \geq z_1(T_0)$. Then (3.6) follows.

Step 3. We prove that for any $\tau > 0$,

$$(3.11) \quad z_s(\tau) \leq z_2(\tau) = \ln \left[\frac{2L}{2r + \sigma^2} \left(1 + \frac{1}{L\tau} \right) \right].$$

In fact, for any $T_0 > 0$, denote

$$(3.12) \quad w(z, \tau) = \begin{cases} -1 + \frac{\tau}{T_0}, & (z, \tau) \in (-\infty, z_2] \times [0, T_0], \\ e^z - \frac{1}{\alpha} e^{\alpha(z-z_2)+z_2} - A + (-1 + \frac{\tau}{T_0}), & (z, \tau) \in (z_2, +\infty) \times [0, T_0], \end{cases}$$

where $\alpha = -\frac{2r}{\sigma^2}$ and

$$e^{z_2} = \frac{2L + 2/T_0}{2r + \sigma^2}, \quad A = \frac{\alpha - 1}{\alpha} e^{z_2} = \frac{L + 1/T_0}{r}.$$

We claim that $w(z, \tau)$ satisfies the following three properties:

(I) $w \in W_{p,loc}^{2,1}(\mathbb{R} \times (0, T)) \cap C(\mathbb{R} \times [0, T])$.

(II) $w(z, 0) \leq (e^z - 1)^+$.

(III) $\partial_\tau w - \mathcal{L}_z w \leq -L$ if $w > 0$ and $(z, \tau) \in \mathbb{R} \times (0, T_0]$.

Indeed, as in the previous step, property (I) is obvious. Next, we prove properties (II) and (III). Considering

$$w(z, 0) = \begin{cases} -1 < 0 < (e^z - 1)^+, & z \leq z_2, \\ (e^z - 1) - \left(\frac{1}{\alpha} e^{\alpha(z-z_2)+z_2} + A \right) < (e^z - 1) - \left(\frac{1}{\alpha} e^{z_2} + A \right) \\ = (e^z - 1) - e^{z_2} < (e^z - 1) \leq (e^z - 1)^+, & z > z_2, \end{cases}$$

property (II) follows.

In the following, we prove property (III). It is clear that $w(z, \tau) \leq 0$ for any $(z, \tau) \in (-\infty, z_2] \times [0, T_0]$. On the other hand, for any $(z, \tau) \in (z_2, +\infty) \times (0, T_0]$, by the method in Step 2, we see that

$$\partial_\tau w - \mathcal{L}_z w = \frac{1}{T_0} - rA + r \left(-1 + \frac{\tau}{T_0}\right) \leq \frac{1}{T_0} - rA = \frac{1}{T_0} - r \frac{L + 1/T_0}{r} = -L.$$

Hence we have property (III).

Based on properties (I)–(III), as in the proof of the uniqueness in Lemma 2.1, we can prove $w \leq v^*$ in $\mathbb{R} \times [0, T]$. Moreover, for $z > z_2$,

$$w(z, T_0) = (e^z - e^{z_2}) - \frac{1}{\alpha} e^{z_2} (e^{\alpha(z-z_2)} - 1) > 0.$$

Thus $v^*(z, T_0) > 0$ for any $z > z_2$, and (3.11) follows.

Step 4. From (3.6) and (3.11), we see that

$$\ln \left[\frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + r\tau} \right) \right] \leq z_s(\tau) \leq \ln \left[\frac{2L}{2r + \sigma^2} \left(1 + \frac{1}{L\tau} \right) \right].$$

Since $z_s(\tau) = y_s(\tau) + (p^* - r)\tau$, then (3.2) follows. \square

Remark. Let $\tau \rightarrow +\infty$ in (3.2). We can deduce that

$$y_s(\tau) + (p^* - r)\tau = z_s(\tau) \rightarrow \ln \left(\frac{2L}{2r + \sigma^2} \right) \quad \text{as} \quad \tau \rightarrow +\infty.$$

Moreover, combining $z_s(\tau) = y_s(\tau) + (p^* - r)\tau$ is continuous in $[0, T]$. We see that there exist two constants M_1 and M_2 such that $M_1 \leq z_s(\tau) \leq M_2$, which are independent of T .

4. Characterizations of the free boundary of problem (1.4). We have obtained some properties (Theorems 3.1 and 3.2) about the free boundary $y_s(\tau)$ of the new problem (2.2). In the following, we come back to consider the free boundary of the original problem (1.4).

Denote $x_s(\tau)$ as the free boundary between \mathbf{NR}_x and \mathbf{SR}_x , which are, respectively, the counterparts of \mathbf{NR}_y and \mathbf{SR}_y in the transformation (2.1). From the transformation (2.1), we have

$$V(x, \tau) = v(x + (r - q - p^*)\tau, \tau), \quad x_s(\tau) = y_s(\tau) + (q + p^* - r)\tau.$$

Theorems 3.1 and 3.2 imply Theorem 4.1

THEOREM 4.1. (1) $x_s(\tau) \in C[0, T] \cap C^\infty(0, T]$ with $x_s(0) = 0$ and satisfies

$$(4.1) \quad \ln \left[\frac{2L}{2r + \sigma^2} \left(1 - \frac{1}{1 + r\tau} \right) \right] \leq x_s(\tau) - q\tau \leq \ln \left[\frac{2L}{2r + \sigma^2} \left(1 + \frac{1}{L\tau} \right) \right].$$

(2) If $q = 0$ and $r \geq L$, $x_s(\tau) = y_s(\tau)$ is strictly decreasing (see Figure 4).

(3) If $q > 0$, $x_s(\tau) \rightarrow +\infty$ as $\tau \rightarrow +\infty$ (see Figure 5).

LEMMA 4.2. If $q = 0$ and $r < L \leq r + \sigma^2/2$ (see Figure 6), then

$$x_s(\tau) \leq 0.$$

Proof. Denote

$$(4.2) \quad w(x, \tau) = \begin{cases} 0, & (x, \tau) \in (-\infty, 0] \times [0, T], \\ C_1 e^x + C_2 e^{\alpha x} - \frac{L}{r}, & (x, \tau) \in (0, +\infty) \times [0, T], \end{cases}$$

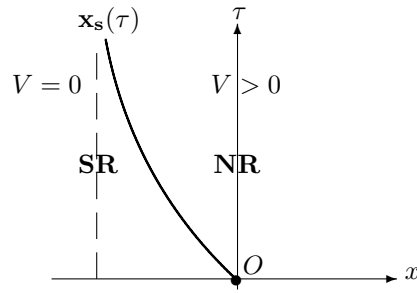


FIG. 4. $q = 0$ and $r \geq L$.

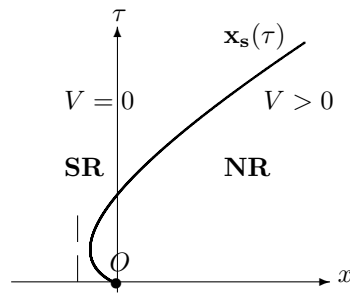


FIG. 5. $q > 0$.

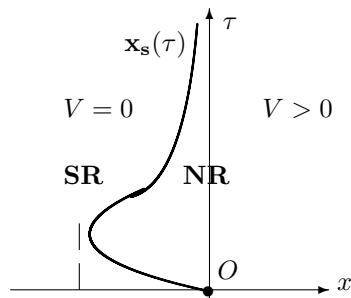


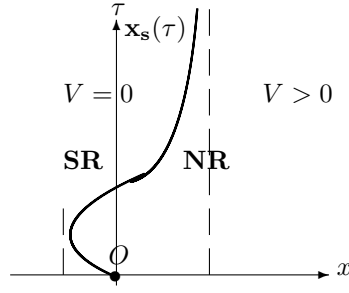
FIG. 6. $q = 0$ and $r < L \leq r + \sigma^2/2$.

where

$$(4.3) \quad \alpha = -2r/\sigma^2, \quad C_1 = \frac{-\alpha L}{r(1-\alpha)} = \frac{2L}{2r + \sigma^2} \leq 1, \quad C_2 = \frac{L}{r(1-\alpha)}.$$

We claim that $w(x, \tau)$ satisfies the following three properties:

$$(4.4) \quad \begin{cases} \text{(I)} & w \in W_{p,loc}^{2,1}(\mathbb{R} \times (0, T)) \cap C(\mathbb{R} \times [0, T]), \\ \text{(II)} & 0 \leq w(x, 0) \leq (e^x - 1)^+, \\ \text{(III)} & \partial_\tau w - \mathcal{L}_x w = -L \quad \text{if } w > 0 \text{ and } (x, \tau) \in \mathbb{R} \times (0, T]. \end{cases}$$

FIG. 7. $q = 0$ and $L > r + \sigma^2/2$.

Indeed, from (4.2) and (4.3), we have

$$(4.5) \quad w(0+0, \tau) = C_1 + C_2 - \frac{L}{r} = 0, \quad \partial_x w(0+0, \tau) = C_1 + \alpha C_2 = 0.$$

Then we get property (I). Next, we prove property (II). Indeed, from (4.2), we have

$$\partial_{xx} w(x, \tau) = C_1 e^x + C_2 \alpha^2 e^{\alpha x} > 0 \quad \forall x > 0.$$

This fact and (4.5) imply $\partial_x w(x, \tau) \geq 0$ and $w(x, \tau) \geq 0$ for any $x > 0$. Moreover, since as $x > 0$, there hold, by $C_1 \leq 1$,

$$\partial_x w(x, \tau) - \partial_x(e^x - 1) = C_1 e^x + C_2 \alpha e^{\alpha x} - e^x < 0 \quad \text{and} \quad w(0, 0) = 0,$$

we have $w(x, \tau) < e^x - 1$ for any $x > 0$. Then, combining this result with (4.2), we get conclusion (II).

In the following, we prove property (III). Notice that 1 and α are the characteristic roots of ODE $\mathcal{L}_x w = 0$; then we have, if $(x, \tau) \in (0, +\infty) \times [0, T_0]$,

$$\partial_\tau w - \mathcal{L}_x w = -r \frac{L}{r} = -L.$$

Based on properties (I)–(III), as in the proof of the uniqueness in Lemma 2.1, we can prove $w \leq V$. Since $w > 0$ for any $x > 0$, so $V > 0$ for any $x > 0$; therefore, $x_s(\tau) \leq 0$. \square

LEMMA 4.3. *If $q = 0$ and $L > r + \sigma^2/2$, then $x_s(\tau) \leq \ln \frac{2L}{2r + \sigma^2}$ (see Figure 7).*

Proof. Denote

$$(4.6) \quad w(x, \tau) = \begin{cases} 0, & (x, \tau) \in (-\infty, x_0] \times [0, T], \\ e^x - \frac{1}{\alpha} e^{\alpha(x-x_0)+x_0} - \frac{L}{r}, & (x, \tau) \in (x_0, +\infty) \times [0, T], \end{cases}$$

where

$$(4.7) \quad \alpha = \frac{-2r}{\sigma^2}, \quad e^{x_0} = \frac{\alpha L}{r(\alpha - 1)} = \frac{2L}{2r + \sigma^2} > 1.$$

The proof is analogous to that of Lemma 4.2. We claim that $w(x, \tau)$ still satisfies the three properties in (4.4). Indeed, from (4.6) and (4.7), we have

$$(4.8) \quad w(x_0+0, \tau) = \frac{\alpha - 1}{\alpha} e^{x_0} - \frac{L}{r} = 0, \quad \partial_x w(x_0+0, \tau) = e^{x_0} - \alpha \frac{e^{x_0}}{\alpha} = 0.$$

Then we get property (I).

Next, from (4.6), we have

$$\partial_{xx}w(x, \tau) = e^x - \alpha e^{\alpha(x-x_0)+x_0} > 0 \quad \forall (x, \tau) \in (x_0, +\infty) \times [0, T],$$

which, combined with (4.8), implies $\partial_x w(x, \tau) > 0$ and $w(x, \tau) > 0$ for any $x > x_0$. Moreover, since $x_0 > 0$ by $e^{x_0} > 1$, there hold, if $x > x_0$,

$$\partial_x w(x, \tau) - \partial_x(e^x - 1) = -e^{\alpha(x-x_0)+x_0} < 0 \quad \text{and} \quad w(x_0, 0) = 0 < e^{x_0} - 1;$$

hence, $w(x, \tau) < e^x - 1$ for any $x > x_0$. Combining this result with (4.6), we get conclusion (II).

Finally, as in the proof of Lemma 4.2, we can check

$$\partial_\tau w - \mathcal{L}_x w = -r \frac{L}{r} = -L, \quad \forall (x, \tau) \in (x_0, +\infty) \times [0, T_0];$$

hence $w \leq V$. Since $w > 0$ for any $x > x_0$, then $V > 0$ for any $x > x_0$; therefore $x_s(\tau) \leq x_0 = \ln \frac{2L}{2r+\sigma^2}$. \square

LEMMA 4.4. *There exists a $\tau_0 > 0$ such that $V(0, \tau) > 0$ for any $0 < \tau < \tau_0$ (see Figures 4-7), i.e.,*

$$\{x = 0, 0 < \tau < \tau_0\} \subset \mathbf{NR}, \quad x_s(\tau) < 0 \quad \text{for } 0 < \tau < \tau_0.$$

Proof. The solution $V(x, t)$ of (1.4) satisfies

$$(4.9) \quad \begin{cases} \partial_\tau V - \mathcal{L}_x V \geq -L, & (x, \tau) \in \mathbb{R} \times (0, T], \\ V(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}. \end{cases}$$

Consider the Cauchy problem

$$(4.10) \quad \begin{cases} \partial_\tau V - \mathcal{L}_x V = -L, & (x, \tau) \in \mathbb{R} \times (0, T], \\ V_1(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}. \end{cases}$$

Applying the comparison principle to problems (4.9) and (4.10), we see that $V(x, \tau) \geq V_1(x, \tau)$. If we can prove that there exists a $\tau_0 > 0$ such that $V_1(0, \tau) > 0$ for $0 < \tau < \tau_0$, then the result of the theorem is an immediate result. To do this we define

$$V_1(x, \tau) = V_2(x, \tau) + \frac{L}{r}(e^{-r\tau} - 1).$$

Then $V_2(x, \tau)$ satisfies

$$(4.11) \quad \begin{cases} \partial_\tau V_2 - \mathcal{L}_x V_2 = 0, & (x, \tau) \in \mathbb{R} \times (0, T], \\ V_2(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}. \end{cases}$$

This is a Cauchy problem which the price of the standard European call option satisfies. Its solution has an explicit formula (see [11]):

$$(4.12) \quad V_2(x, \tau) = e^{x-q\tau} N(\widehat{d}_1) - e^{-r\tau} N(\widehat{d}_2),$$

where

$$N(\widehat{d}_1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\widehat{d}_1} e^{-\eta^2/2} d\eta,$$

$$\widehat{d}_1 = \frac{x + (r - q + \sigma^2/2)\tau}{\sigma\sqrt{\tau}}, \quad \widehat{d}_2 = \frac{x + (r - q - \sigma^2/2)\tau}{\sigma\sqrt{\tau}}.$$

Now we check

$$\begin{aligned} \partial_\tau V_2(0, \tau) &= -qe^{-q\tau}N(d_1) + e^{-q\tau}n(d_1)\frac{r - q + \sigma^2/2}{2\sigma\sqrt{\tau}} + re^{-r\tau}N(d_2) \\ &\quad - e^{-r\tau}n(d_2)\frac{r - q - \sigma^2/2}{2\sigma\sqrt{\tau}} \geq -q + e^{-q\tau}n(d_1)\frac{\sigma}{2\sqrt{\tau}}, \end{aligned}$$

where

$$n(d_1) = e^{-d_1^2/2}/\sqrt{2\pi}, \quad d_1 = \frac{(r - q + \sigma^2/2)\sqrt{\tau}}{\sigma}, \quad d_2 = \frac{(r - q - \sigma^2/2)\sqrt{\tau}}{\sigma},$$

and we had utilized the equality $e^{-q\tau}n(d_1) = e^{-r\tau}n(d_2)$. As $\tau \rightarrow 0^+$, we have

$$d_1 \rightarrow 0, \quad n(d_1) \rightarrow 1, \quad \partial_\tau V_2(0, \tau) \rightarrow +\infty.$$

It follows that $\partial_\tau V_1(0, \tau) \rightarrow +\infty$ as $\tau \rightarrow 0^+$. Since $V_1(0, 0) = 0$ and $V_1 \in C(\mathbb{R} \times [0, T])$, then we see that there exists a $\tau_0 > 0$ such that $V_1(0, \tau) > 0$ for any $0 < \tau < \tau_0$. \square

From Theorem 4.1 and Lemmas 4.2–4.4, we can conclude the following theorem.

THEOREM 4.5. (1) *If $q > 0$, then $x_s(\tau)$ is not monotonic; moreover, $x_s(\tau)$ changes its sign (see Figure 5) and*

$$\lim_{\tau \rightarrow +\infty} [x_s(\tau) - q\tau] = \ln \frac{2L}{2r + \sigma^2}.$$

(2) *If $q = 0$ and $r < L \leq r + \sigma^2/2$, then $x_s(\tau) \leq 0$ (see Figure 6) and*

$$\lim_{\tau \rightarrow +\infty} x_s(\tau) = \ln \frac{2L}{2r + \sigma^2}.$$

(3) *If $q = 0$ and $L > r + \sigma^2/2$, then $x_s(\tau)$ is not monotonic; moreover, $x_s(\tau)$ changes its sign (see Figure 7) and $x_s(\tau) \leq \ln \frac{2L}{2r + \sigma^2}$ with*

$$\lim_{\tau \rightarrow +\infty} x_s(\tau) = \ln \frac{2L}{2r + \sigma^2}.$$

Proof. (1) From Theorem 4.1, we see that $x_s(0) = 0$ and $x_s(\tau) \rightarrow +\infty$ as $\tau \rightarrow +\infty$; moreover, by Lemma 4.4, there exists a $\tau_0 > 0$ such that $x_s(\tau) < 0$ for any $0 < \tau < \tau_0$, and hence $x_s(\tau)$ is not monotonic and changes sign. Inequality (4.1) implies

$$\lim_{\tau \rightarrow +\infty} (x_s(\tau) - q\tau) = \ln \frac{2L}{2r + \sigma^2}.$$

(2) and (3) are easily obtained as above, so we omit their proofs. \square

Remark on the monotonicity of $x_s(\tau)$ in (2) of Theorem 4.5. According to the result of asymptotic expansion of $S_s(\tau) = e^{x_s(\tau)}$ close to expire in (1.2), we know that

$$x_s(\tau) \sim -\sigma(-\tau \ln \tau)^{1/2}(1 + o(1)); \quad x'_s(\tau) \sim -\infty \quad \text{as } \tau = T - t \rightarrow 0^+.$$

So if τ is small enough, $x_s(\tau)$ is monotonic decreasing. Combining this result with (4.1), we guess that $x_s(\tau)$ should not be monotonic. But we cannot prove this conjecture.

5. The monotonicity of the free boundary $x_s(\cdot)$ with respect to the parameters r , q , and L . In this section, we have the following theorem.

THEOREM 5.1. $x_s(\tau)$ is decreasing with respect to r and increasing with respect to q and L .

Proof. We divide the proof into three steps.

Step 1. We prove

$$(5.1) \quad \partial_x V - V \geq 0.$$

For the proof of (5.1), we come back to the approximation problem (2.11) and denote

$$W^* = \partial_y v_{\varepsilon, n} - v_{\varepsilon, n} + 2\varepsilon.$$

We claim that $W^* \geq 0$; otherwise, we suppose W^* achieves its negative minimum at $(x_0, \tau_0) \in \overline{\Omega_T^n}$.

If n is large enough, from the definition of π_ε and the initial and boundary values in (2.11) and (2.13), we have

$$\begin{cases} W^*(-n, \tau) \geq 0 - 0 + 2\varepsilon > 0 & \text{(by (2.14)),} \\ W^*(n, \tau) \geq e^{n+(p^*-r)\tau} - e^{n+(p^*-r)\tau} + 2\varepsilon > 0 & \text{(by (2.12)),} \\ W^*(y, 0) = \partial_y \pi_\varepsilon(e^y - 1) - \pi_\varepsilon(e^y - 1) + \varepsilon \geq 0 - \varepsilon + 2\varepsilon > 0 & \text{if } e^y - 1 \leq \varepsilon, \\ W^*(y, 0) > e^y - (e^y - 1) + 2\varepsilon > 0 & \text{if } e^y - 1 > \varepsilon. \end{cases}$$

Hence, from (2.3), we deduce that

$$(5.2) \quad (x_0, \tau_0) \in \Omega_T^n \quad \text{and} \quad \partial_\tau W^*(x_0, \tau_0) - \mathcal{L}_y W^*(x_0, \tau_0) \leq 0.$$

On the other hand, (2.14) and $W^*(x_0, \tau_0) < 0$ imply $v_{\varepsilon, n}(x_0, \tau_0) > 2\varepsilon$. Then

$$\beta_\varepsilon(v_{\varepsilon, n}(x_0, \tau_0)) = 0, \quad \beta'_\varepsilon(v_{\varepsilon, n}(x_0, \tau_0)) = 0.$$

Combining (2.11) and (2.13), we see that at the point (x_0, τ_0)

$$\partial_\tau W^* - \mathcal{L}_y W^* = L + 2r\varepsilon > 0,$$

which contradicts (5.2). Hence,

$$W^* \geq 0 \quad \text{and} \quad \partial_y v_{\varepsilon, n} - v_{\varepsilon, n} \geq -2\varepsilon.$$

Taking $\varepsilon \rightarrow 0^+$ and $n \rightarrow \infty$, we obtain

$$\partial_y v - v \geq 0.$$

Eventually, the transformation (2.1) implies (5.1).

Step 2. Suppose that V_1 is the solution to (1.4), where r is r_1 , and that V_2 is the solution to (1.4), where r is r_2 and $r_1 > r_2$. By (5.1), we deduce $\partial_x V_1 - V_1 \geq 0$. Then V_1 satisfies

$$\begin{cases} \partial_\tau V_1 - \frac{\sigma^2}{2} \partial_{xx} V_1 - (r_2 - q - \frac{\sigma^2}{2}) \partial_x V_1 + r_2 V_1 = -L + (r_1 - r_2)(\partial_x V_1 - V_1) \geq -L & \text{if } V_1 > 0 \quad \text{and} \quad (x, \tau) \in \mathbb{R} \times (0, T], \\ \partial_\tau V_1 - \frac{\sigma^2}{2} \partial_{xx} V_1 - (r_2 - q - \frac{\sigma^2}{2}) \partial_x V_1 + r_2 V_1 \geq -L + (r_1 - r_2)(\partial_x V_1 - V_1) = -L & \text{if } V_1 = 0 \quad \text{and} \quad (x, \tau) \in \mathbb{R} \times (0, T], \\ V_1(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}. \end{cases}$$

By the comparison principle for the solution of variational inequalities with respect to nonhomogeneous terms, we have $V_1 \geq V_2$. If we denote $\mathbf{NR}_1, \mathbf{NR}_2$ to be the nontransaction region of V_1, V_2 , respectively, then $\mathbf{NR}_1 \supset \mathbf{NR}_2$. Hence $x_s(\tau)$ is decreasing with respect to r .

Step 3. As in Step 2, from (2.21), we can deduce that $x_s(\tau)$ is increasing with respect to q and L . \square

6. Numerical methods and results. Starting from problem (1.4), we have

$$(6.1) \quad \begin{cases} \min\{ \partial_\tau V - \frac{\sigma^2}{2} \partial_{xx} V - (r - q - \frac{\sigma^2}{2}) \partial_x V + rV + LV \} = 0, & x \in \mathbb{R}, \tau \in (0, T], \\ V(x, 0) = (e^x - 1)^+, & x \in \mathbb{R}. \end{cases}$$

Given mesh size $\Delta x, \Delta \tau > 0$, $V_j^n = V(j\Delta x, n\Delta \tau)$ represents the value of numerical approximation at $(j\Delta x, n\Delta \tau)$. Then the PDE is changed into the following difference equation:

$$(6.2) \quad \begin{cases} \min \left\{ \frac{V_j^n - V_j^{n-1}}{\Delta \tau} - \frac{\sigma^2}{2} \frac{V_{j+1}^{n-1} - 2V_j^{n-1} + V_{j-1}^{n-1}}{\Delta x^2} \right. \\ \quad \left. - (r - q - \frac{\sigma^2}{2}) \frac{V_{j+1}^{n-1} - V_{j-1}^{n-1}}{2\Delta x} + rV_j^n + L, V_j^n \right\} = 0, \\ V_j^0 = (e^{j\Delta x} - 1)^+. \end{cases}$$

This means

$$(6.3) \quad \begin{cases} V_j^n = \max \left\{ \frac{1}{1+r\Delta \tau} \left[\left(1 - \frac{\sigma^2 \Delta \tau}{\Delta x^2}\right) V_j^{n-1} + \frac{\sigma^2 + (r-q-\sigma^2/2)\Delta x}{2\Delta x^2} \Delta \tau V_{j+1}^{n-1} \right. \right. \\ \quad \left. \left. + \frac{\sigma^2 - (r-q-\sigma^2/2)\Delta x}{2\Delta x^2} \Delta \tau V_{j-1}^{n-1} - L\Delta \tau \right], 0 \right\}, \\ V_j^0 = \max\{0, e^{j\Delta x} - 1\}. \end{cases}$$

Choosing $\sigma^2 \Delta \tau (\Delta x)^{-2} = 1$, we have

$$(6.4) \quad \begin{cases} V_j^n = \max \left\{ \frac{1}{1+r\Delta \tau} \left[\left(\frac{1}{2} + \frac{r-q-\sigma^2/2}{2\sigma} \sqrt{\Delta \tau}\right) V_{j+1}^{n-1} \right. \right. \\ \quad \left. \left. + \left(\frac{1}{2} - \frac{r-q-\sigma^2/2}{2\sigma} \sqrt{\Delta \tau}\right) V_{j-1}^{n-1} - L\Delta \tau \right], 0 \right\}, \\ V_j^0 = \max\{0, e^{j\Delta x} - 1\}. \end{cases}$$

Denote $u = e^{\sigma \sqrt{\Delta \tau}}, d = u^{-1}, \rho = e^{r\Delta \tau}, p = (\rho e^{-q\Delta \tau} - d)(u - d)^{-1}$.

Then it is clear, as $\Delta \tau \rightarrow 0^+$, that the following hold:

$$\frac{1}{2} + \frac{r - q - \sigma^2/2}{2\sigma} \sqrt{\Delta \tau} = p + o(\Delta \tau), \quad \frac{1}{1 + r\Delta \tau} = \frac{1}{\rho} + O(\Delta \tau^2).$$

Neglecting a higher order of $\sqrt{\Delta \tau}$, we obtain

$$(6.5) \quad \begin{cases} V_j^n = \max \left\{ \frac{1}{\rho} [pV_{j+1}^{n-1} + (1 - p)V_{j-1}^{n-1} - L\Delta \tau], 0 \right\}, \\ V_j^0 = \max \{0, u^j - 1\}. \end{cases}$$

Consider the point $(x, \tau) = (j\Delta x, n\Delta \tau)$. Then

$$\begin{aligned} V_j^n &= V(x, \tau), \quad V_j^{n-1} = V(x, \tau - \Delta \tau), \\ V_{j+1}^{n-1} &= V(x + \Delta x, \tau - \Delta \tau), \quad V_{j-1}^{n-1} = V(x - \Delta x, \tau - \Delta \tau). \end{aligned}$$

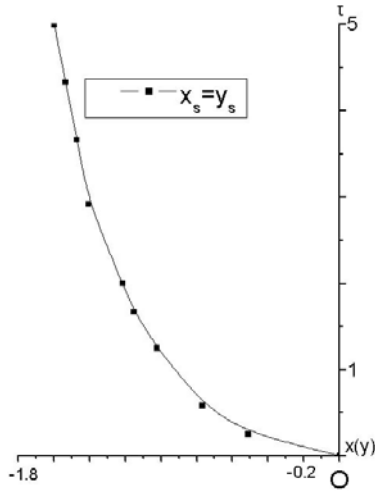


FIG. 8. $q = 0$ and $r \geq L$.

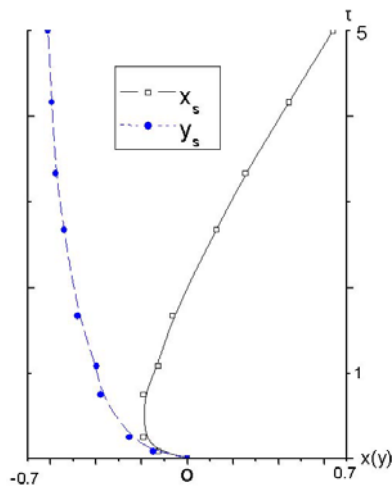


FIG. 9. $q > 0$.

Then we get Figures 8–11.

Figure 8 shows a plot of the optimal stop boundaries $x_s(\tau)$ ($y_s(\tau)$) as a function of time τ when $q = 0$, $r \geq L$. The parameter values used in the calculations are $r = 0.4$, $q = 0$, $\sigma = 0.7$, $L = 0.1$, $T = 5$, $n = 600$. $x_s(\tau)$, $y_s(\tau)$ are the free boundaries of problems (1.4) and (2.2), respectively. Observe that $x_s(0) = y_s(0) = 0$, $-1.587 \leq x_s(\tau) = y_s(\tau) \leq 0$, and $x_s(\tau), y_s(\tau)$ are strictly decreasing with respect to τ . The numerical result is consistent with that of our proof (see Figure 8).

Figure 9 shows a plot of the optimal stop boundaries $x_s(\tau)$ ($y_s(\tau)$) as a function of time τ when $q > 0$. The parameter values used in the calculations are $r = 0.25$, $q = 0.2$, $\sigma = 0.7$, $L = 0.3$, $T = 5$, $n = 600$. $x_s(\tau)$, $y_s(\tau)$ are the free boundaries of problems (1.4) and (2.2), respectively. Observe that $x_s(0) = y_s(0) = 0$, $x_s(\tau) \geq -0.1917$, and $x_s(\tau)$ is not monotonic and $y_s(\tau)$ is decreasing. The numerical result is consistent with that of our proof (see Figure 9).

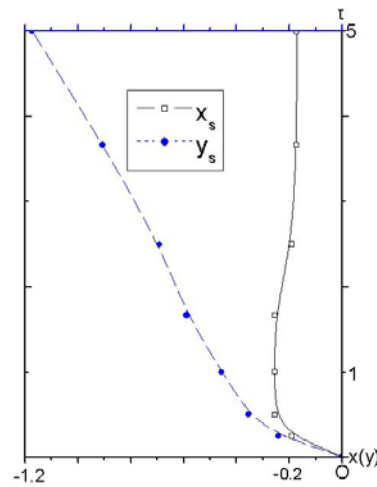
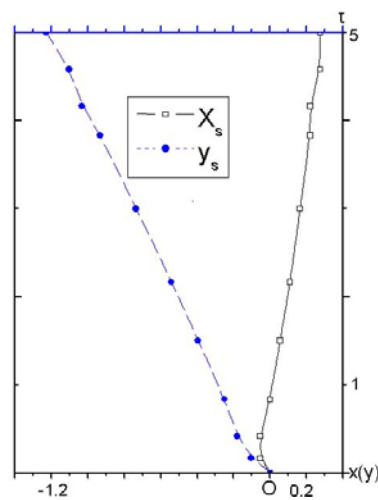
FIG. 10. $q = 0$ and $r < L \leq r + \sigma^2/2$.FIG. 11. $q = 0$ and $L > r + \sigma^2/2$.

Figure 10 shows a plot of the optimal stop boundaries $x_s(\tau)$ ($y_s(\tau)$) as a function of time τ when $q = 0$ and $r < L \leq r + \sigma^2/2$. The parameter values used in the calculations are $r = 0.1$, $q = 0$, $\sigma = 0.7$, $L = 0.3$, $T = 5$, $n = 600$. $x_s(\tau)$, $y_s(\tau)$ are the free boundaries of problems (1.4) and (2.2), respectively. Observe that $x_s(0) = y_s(0) = 0$, $-0.2556 \leq x_s(\tau) \leq -0.1278$, and $y_s(\tau)$ is decreasing. The numerical result is consistent with that of our proof (see Figure 10).

Figure 11 shows a plot of the optimal stop boundaries $x_s(\tau)$ ($y_s(\tau)$) as a function of time τ when $q = 0$ and $L > r + \sigma^2/2$. The parameter values used in the calculations are $r = 0.2$, $q = 0$, $\sigma = 0.6$, $L = 0.5$, $T = 5$, $n = 600$. $x_s(\tau)$, $y_s(\tau)$ are the free boundaries of problems (1.4) and (2.2), respectively. Observe that $x_s(0) = y_s(0) = 0$, $-0.0548 \leq x_s(\tau) \leq 0.2739$, and $x_s(\tau)$ is not monotonic and $y_s(\tau)$ is decreasing. The numerical result is consistent with that of our proof (see Figure 11).

Appendix. Formulation of the model. A European call option is a contract which gives the owner the right but not the obligation to buy an asset at a fixed price K at the expiry date T . Because the holder of the option stands to make a profit without risking a loss, the holder must pay some premium for the option.

In a conventional European call option contract (see [2]), the holder pays the premium entirely up front and acquires the right. In a continuous installment European call option contract, the holder pays a smaller up-front premium and then a constant stream of installments at a certain rate per unit time. However, the holder can choose at any time to stop making installment payments by stopping the option contract.

There are some papers about install options, such as [7], [8], [1], [6]. Particularly, there are a variational inequality model in [7] and some numerical results about the model. In the following, we deduce a parabolic variational inequality model using their idea.

We consider a standard model for perfect market, continuous trading, no-arbitrage opportunity, a constant interest rate $r > 0$, and an asset with constant continuous dividend yield $q \geq 0$ with price S following a geometric Brownian motion

$$(A.1) \quad dS = \mu S dt + \sigma S dB,$$

where B is a Wiener process on a risk neutral probability space, $\mu > 0$ is the expected return rate, and $\sigma > 0$ is the constant volatility.

Let $C(S, t)$ denote the value of a European continuous installment call option and let L^* be the continuous install rate. Applying Itô's formula to $C(S, t)$ and combining (A.1), we obtain the dynamics

$$(A.2) \quad dC = \left(\frac{\partial C}{\partial t} + \frac{\sigma^2}{2} S^2 \frac{\partial^2 C}{\partial S^2} + \mu S \frac{\partial C}{\partial S} \right) dt + \sigma S \frac{\partial C}{\partial S} dB.$$

We construct the Δ -hedging portfolio consisting of one continuous installment option and an amount $-\Delta$ asset. The value of this portfolio is

$$(A.3) \quad \Pi = C - \Delta S.$$

Then its dynamics is given by

$$(A.4) \quad d\Pi = dC - \Delta S q dt - \Delta dS - L^* dt \leq r \Pi dt = r(C - \Delta S) dt.$$

Combining (A.1) and (A.3), we get

$$(A.5) \quad \begin{aligned} 0 &\geq d\Pi - r(C - \Delta S) dt \\ &= \left(\frac{\partial C}{\partial t} + \frac{\sigma^2}{2} S^2 \frac{\partial^2 C}{\partial S^2} + \mu S \left(\frac{\partial C}{\partial S} - \Delta \right) + (r - q) \Delta S - rC - L^* \right) dt \\ &+ \sigma S \left(\frac{\partial C}{\partial S} - \Delta \right) dB. \end{aligned}$$

Setting $\Delta = \partial_S C$ the coefficient of dB vanishes. The portfolio is instantaneously riskless; then we see that $C(S, t)$ satisfies

$$(A.6) \quad \partial_t C + \frac{\sigma^2}{2} S^2 \partial_{SS} C + (r - q) S \partial_S C - rC \leq L^*.$$

Moreover, in the nontransaction region, $C(S, t)$ satisfies

$$(A.7) \quad \partial_t C + \frac{\sigma^2}{2} S^2 \partial_{SS} C + (r - q) S \partial_S C - rC = L^*.$$

At expiry date T , if $S > K$, the holder of the option will buy the asset with price S at K and can make a profit $S - K$. On the other hand, if $S \leq K$, he will not exercise the option and stand any loss. Hence, the value of the option is $(S - K)^+$ at expiry date T , that is,

$$(A.8) \quad C(S, T) = (S - K)^+.$$

Because the owner must keep paying premiums to keep the option alive, if S is small enough at some t , the present value of the expected pay-off may be less than the present value of the remaining installments; then the holder would allow the option to lapse and stop paying installment payments. Hence, the value of the option is 0 in the case in which (S, τ) lies in the stop region.

It is clear that $C(S, t) \geq 0$ in nontransaction region. Otherwise, he can choose to stop the option for increased profit; hence we have

$$(A.9) \quad C(S, t) \geq 0.$$

From the above deduction, we see that the following equality holds in the stop region and the nontransaction region:

$$(A.10) \quad \left[\partial_t C + \frac{\sigma^2}{2} S^2 \partial_{SS} C + (r - q) S \partial_S C - rC - L^* \right] C(S, \tau) = 0.$$

Then we see that the value of the European continuous installment call option $C(S, t)$ satisfies (A.6), (A.8)–(A.10), that is, (1.1). Moreover, from the smooth fit conditions [12], we know that $C, \partial_S C$ are continuous.

REFERENCES

- [1] G. ALOBAIDI, R. MALLIER, AND S. DEAKIN, *Laplace transforms and installment options*, Math. Models Methods Appl. Sci., 8 (2004), pp. 1167–1189.
- [2] F. BLACK AND M. SCHOLES, *The pricing of options and corporate liabilities*, J. Political Economy, 81 (1973), pp. 637–659.
- [3] A. BLANCHET, J. DOLBEAULT, AND R. MONNEAU, *On the continuity of the time derivative of the solution to the parabolic obstacle problem with variable coefficients*, J. Math. Pures Appl. (9), 85 (2006), pp. 371–414.
- [4] L. A. CAFFARELLI, *The regularity of free boundaries in higher dimensions*, J. Acta Math., 139 (1977), pp. 135–184.
- [5] L. A. CAFFARELLI, A. PETROSYAN, AND H. SHAHGHOLIAN, *Regularity of a free boundary in parabolic potential theory*, J. Amer. Math. Soc., 17 (2004), pp. 827–869.
- [6] P. CIURLIA AND I. ROKO, *Valuation of American continuous-installment option*, Comput. Economics, 25 (2005), pp. 143–165.
- [7] M. DAVIS, W. SCHACHERMAYER, AND R. TOMPKINS, *Pricing, no-arbitrage bounds and robust hedging of installment options*, Quantitative Finance, 6 (2001), pp. 597–610.
- [8] M. DAVIS, W. SCHACHERMAYER, AND R. TOMPKINS, *Installment options and static hedging*, J. Risk Fin., 3 (2002), pp. 46–52.
- [9] A. FRIEDMAN, *Parabolic variational inequalities in one space dimension and smoothness of the free boundary*, J. Funct. Anal., 18 (1975), pp. 151–176.
- [10] A. FRIEDMAN, *Variational Principle and Free Boundary Problems*, John Wiley & Sons, New York, 1982.
- [11] L. JIANG, *Mathematical Modeling and Methods of Option Pricing*, World Scientific, River Edge, NJ, 2005.
- [12] R. C. MERTON, *Theory of rational option pricing*, Bell J. Econom. and Management Sci., 4 (1973), pp. 141–184.
- [13] K. TAO, *On Aleksandrov, Bakel'man type maximum principle for second order parabolic equations*, Comm. Partial Differential Equations, 10 (1985) pp. 543–553.

DETERMINATION OF THE SPECTRAL GAP IN THE KAC MODEL FOR PHYSICAL MOMENTUM AND ENERGY-CONSERVING COLLISIONS*

ERIC A. CARLEN[†], JEFFREY S. GERONIMO[†], AND MICHAEL LOSS[†]

Abstract. The Kac model describes the local evolution of a gas of N particles with *three*-dimensional velocities by a random walk in which the steps correspond to binary collisions that conserve momentum as well as energy. The state space of this walk is a sphere of dimension $3N - 4$. The Kac conjecture concerns the spectral gap in the one-step transition operator Q for this walk. In this paper, we compute the exact spectral gap. As in previous work by Carlen, Carvalho, and Loss, where a lower bound on the spectral gap was proved, we use a method that relates the spectral properties of Q to the spectral properties of a simpler operator P , which is simply an average of certain noncommuting projections. The new feature is that we show how to use a knowledge of certain eigenfunctions and eigenvalues of P to determine spectral properties of Q , instead of simply using the spectral gap for P to bound the spectral gap for Q , inductively in N , as in previous work. The methods developed here can be applied to many other high-dimensional stochastic process, as we shall explain. We also use some deep results on Jacobi polynomials to obtain the required spectral information on P , and we show how the identity through which Jacobi polynomials enter our problem may be used to obtain new bounds on Jacobi polynomials.

Key words. Kac model, orthogonal polynomials, spectral gap

AMS subject classifications. 33C20, 82C40, 76P05

DOI. 10.1137/070695423

1. The Markov transition operator Q for the Kac walk. Let X_N be the N particle state space consisting of N -tuples $\vec{v} = (v_1, \dots, v_N)$ of vectors v_j in \mathbb{R}^3 with

$$\sum_{j=1}^N |v_j|^2 = 1 \quad \text{and} \quad \sum_{j=1}^N v_j = 0.$$

We think of a point \vec{v} as specifying the velocities of N particles and shall consider a random walk on X_N that was introduced by Kac [7]. At each step of this random walk, \vec{v} is updated due to the effect of a binary collision that conserves energy and momentum—hence the constraints defining X_N .

To specify this walk in more detail, we consider a collision in which particles i and j collide. Suppose that v_i^* and v_j^* are the postcollisional velocities, while v_i and v_j are the precollisional velocities. Then by momentum conservation, the center of mass velocity is conserved; i.e.,

$$v_i^* + v_j^* = v_i + v_j.$$

Furthermore, by energy conservation, i.e., $|v_i^*|^2 + |v_j^*|^2 = |v_i|^2 + |v_j|^2$, and the parallelogram law, it follows that

$$|v_i^* - v_j^*| = |v_i - v_j|.$$

*Received by the editors July 18, 2007; accepted for publication (in revised form) December 21, 2007; published electronically April 23, 2008.

<http://www.siam.org/journals/sima/40-1/69542.html>

[†]School of Mathematics, Georgia Institute of Technology, Atlanta, GA 30332 (carlen@math.gatech.edu, geronimo@math.gatech.edu, loss@math.gatech.edu). The first and third authors' work was partially supported by U.S. National Science Foundation grant DMS-060037. The second author's work was partially supported by U.S. National Science Foundation grant DMS-0500641.

This leads to a natural parameterization of all of the possible binary collision outcomes that conserve energy and momentum: The parameter σ is a unit vector in S^2 , and, when particles i and j collide, one updates $\vec{v} \rightarrow \vec{v}^* = R_{i,j,\sigma}(\vec{v})$, where

$$(1.1) \quad \begin{aligned} v_i^* &= \frac{v_i + v_j}{2} + \frac{|v_i - v_j|}{2} \sigma, \\ v_j^* &= \frac{v_i + v_j}{2} - \frac{|v_i - v_j|}{2} \sigma, \\ v_k^* &= v_k \quad \text{for } k \neq i, j. \end{aligned}$$

The *Kac walk* on X_N is a random walk in which the steps are such binary collisions between pairs of particles. At each step, one picks a pair (i, j) , $i < j$ uniformly at random, and also a unit vector σ in S^2 . One then makes the update described in (1.1). Of course it remains to specify the probabilistic rule according to which σ should be selected. In the physics being modeled here, the likelihood of selecting a particular σ will depend only on the resulting *scattering angle* θ , which is the angle between $v_i^* - v_j^*$ and $v_i - v_j$. In the parameterization above, this is the angle between σ and $v_i - v_j$. That is,

$$\cos(\theta) = \sigma \cdot \frac{v_i - v_j}{|v_i - v_j|}.$$

The *scattering rate function* b is a nonnegative integrable function on $[-1, 1]$ with

$$\frac{1}{2} \int_{-1}^1 b(u) du = 1.$$

Then for any $v_i \neq v_j$, and with $d\sigma$ being the uniform probability measure on S^2 ,

$$(1.2) \quad \int_{S^2} b\left(\sigma \cdot \frac{v_i - v_j}{|v_i - v_j|}\right) d\sigma = 1.$$

(If $v_i = v_j$, the collision has no effect and can be ignored.) One selects σ according to the probability density that is integrated in (1.2).

There are several choices of b of particular interest. One is the *uniform redirection model*, given by $b(x) = 1$ for all $-1 \leq x \leq 1$. In this case, the new direction of the relative velocity σ is chosen uniformly from S^2 .

Another is the *Morgenstern model* [10], [11], or the *uniform reflection model*: For any unit vector $\omega \in S^2$, let H_ω be the reflection given by

$$H_\omega(v) = v - 2(v \cdot \omega)\omega.$$

In the uniform reflection model, one updates the relative velocity according to

$$v_i - v_j \rightarrow H_\omega(v_i - v_j) = v_i^* - v_j^*,$$

with ω chosen uniformly. The relation between ω and σ is given by $\sigma = H_\omega((v_i - v_j)/|v_i - v_j|)$, and, by computing the Jacobian of the map $\omega \rightarrow \sigma$, one finds

$$b(x) = \frac{1}{\sqrt{2}\sqrt{1-x}}.$$

Both of these belong to the one-parameter family

$$(1.3) \quad b_\alpha(x) = (1 - \alpha)2^\alpha(1 - x)^{-\alpha}.$$

By leaving the particular choice of b open, this completes the specification of the steps in the Kac walk. For more detail and background, see [7] and [3].

The main object of study here is the spectrum of the one-step transition operator Q for this random walk, and the manner in which this spectrum depends on N . Q is defined as follows: Let \vec{v}_n be the state of the process after the n th step. The one-step Markov transition operator Q is given by taking the conditional expectation

$$Q\phi(\vec{v}) = \mathbf{E}(\phi(\vec{v}_{n+1}) \mid \vec{v}_n = \vec{v})$$

for any continuous function ϕ on X_N .

From the above description, one deduces the formula

$$(1.4) \quad Q\phi(\vec{v}) = \binom{N}{2}^{-1} \sum_{i < j} \int_{S^2} \phi(R_{i,j,\sigma}(\vec{v})) b\left(\sigma \cdot \frac{v_i - v_j}{|v_i - v_j|}\right) d\sigma.$$

Let σ_N denote the uniform probability measure on X_N , which is the normalized measure induced on X_N as a manifold embedded in \mathbb{R}^{3N} .

For any two unit vectors σ and ω , one sees from (1.1) that

$$R_{i,j,\sigma}(R_{i,j,\omega}\vec{v}) = R_{i,j,\sigma}\vec{v}.$$

From this and the fact that the measure $d\sigma_N$ is invariant under $\vec{v} \mapsto R_{i,j,\sigma}\vec{v}$, it follows that, for any continuous functions ϕ and ψ on X_N ,

$$\int_{X_N} \psi(\vec{v}) Q\phi(\vec{v}) d\sigma_N = \int_{X_N} \int_{S^2} \int_{S^2} \psi(R_{i,j,\omega}\vec{v}) \phi(R_{i,j,\sigma}\vec{v}) b(\omega \cdot \sigma) d\omega d\sigma d\sigma_N.$$

It follows that Q is a self-adjoint Markov operator on $L^2(X_N, \sigma_N)$. Moreover, it is clearly a Markov operator; that is, in addition to being self-adjoint, Q is positivity preserving and $Q1 = 1$.

The motivation for considering the spectral properties of Q stems from a theorem of Kac [7] that relates the continuous time version of the Kac walk to the nonlinear Boltzmann equation. For the details, see [7] or [3]. Let $\vec{v}(t)$ denote the random variable giving the state of the system at time t for the process run in continuous time with the jumps taking place in a Poisson stream with the mean time between jumps being $1/N$. Then the equation describing the evolution of the probability law of $\vec{v}(t)$ is called the *Kac master equation*: If the initial law on X_N has a density F_0 , then the law at time t has a density $F(\vec{v}, t)$ satisfying

$$\frac{\partial}{\partial t} F(\vec{v}, t) = N(Q - I)F(\vec{v}, t), \quad \text{with} \quad F(\vec{v}, 0) = F_0(\vec{v}).$$

The solution $F(\vec{v}, t)$ is of course given by

$$F(\vec{v}, t) = e^{t\mathcal{L}} F_0(\vec{v}),$$

where

$$\mathcal{L} = N(Q - I).$$

Since Q is a self-adjoint Markov operator, its spectrum lies in the interval $[-1, 1]$, and since $Q1 = 1$, the constant function is an eigenfunction with eigenvalue 1. It is easily seen that, as long as $b(x)$ is strictly positive on a neighborhood of $x = 1$, the eigenvalue 1 of Q has multiplicity one. It then follows that the spectrum of \mathcal{L} lies in $[-2N, 0]$ and that 0 is an eigenvalue of multiplicity one. We impose this assumption on b throughout what follows.

The Kac conjecture for this stochastic process pertains to the spectral gap

$$\Delta_N = \inf \left\{ - \int_{X_N} \phi(\vec{v}) \mathcal{L} \phi(\vec{v}) d\sigma_N \mid \int_{X_N} \phi^2(\vec{v}) d\sigma_N = 1, \int_{X_N} \phi(\vec{v}) d\sigma_N = 0 \right\}$$

and states that

$$\liminf_{N \rightarrow \infty} \Delta_N > 0.$$

This was proved by Carlen, Carvalho, and Loss [3], but without an explicit lower bound. Kac also made a similar conjecture for a simplified model with one-dimensional velocities and no conservation of momentum. For this model, the conjecture was first proved by Janvresse [6], though her approach provided no explicit lower bound. The sharp bound for the simplified model was first established in [2]. See Maslen [9] for a representation theoretic approach.

The main goal of the present paper is to compute *exactly* $\liminf_{N \rightarrow \infty} \Delta_N$. We shall be able to do this under an easily checked condition relating Δ_2 and the quantities

$$(1.5) \quad B_1 = \frac{1}{2} \int_{-1}^1 xb(x)dx \quad \text{and} \quad B_2 = \frac{1}{2} \int_{-1}^1 x^2b(x)dx.$$

The condition, given in (1.6) below, will turn out to be satisfied when b is given by b_α , as in (1.3), for all $0 \leq \alpha \leq 7/9$.

THEOREM 1.1. *Suppose that $B_2 > B_1$ and that*

$$(1.6) \quad \Delta_2 \geq \frac{20}{9}(1 - B_2).$$

Then for all $N \geq 3$,

$$(1.7) \quad \Delta_N = (1 - B_2) \frac{N}{(N - 1)}.$$

Moreover, the eigenspace is three-dimensional and is spanned by the functions

$$(1.8) \quad \phi(\vec{v}) = \sum_{j=1}^N |v_j|^2 v_j^\alpha, \quad \alpha = 1, 2, 3,$$

where v_j^α denotes the α th component of v_j .

As we shall see in the next section, for many choices of b , including the b_α , there is a simple monotonicity of the eigenvalues of Q for $N = 2$ which ensures that the eigenfunction providing the gap comes from a first degree polynomial and thus that

$$(1.9) \quad \Delta_2 = 2(1 - B_1).$$

When (1.9) is satisfied, the condition (1.6) reduces to $(1 - B_1)/(1 - B_2) > 10/9$.

Next, notice that the eigenfunctions listed in (1.8) are symmetric under permutation of the particle indices. Indeed, the operator Q commutes with such permutations, so that the subspace of functions with this symmetry is invariant. As explained in [7] and [3], it is the spectrum of Q on this subspace that is relevant for the study of the Boltzmann equation.

Moreover, notice that, in the collision rules (1.1), exchanging v_i^* and v_j^* has the same effect as changing σ to $-\sigma$. For this reason, if one's primary object of interest is the Boltzmann equation, one may freely assume that b is a symmetric function on $[-1, 1]$, since then replacing $b(x)$ by $(b(x) + b(-x))/2$ will have no effect on the spectrum of Q on the invariant subspace of symmetric functions or on the corresponding Boltzmann equation. (See the introduction of [4] for more discussion of this point in the context of the Boltzmann equation.) If B is symmetric, then $B_1 = 0$, and we do have $B_1 < B_2$.

However, it is interesting that the Kac conjecture holds without restriction to the symmetric subspace and that the methods developed here can be used to determine the spectral gap even when b is not symmetric, and the eigenfunctions corresponding to the gap eigenvalue are not symmetric.

When b is not symmetric, it may happen that $B_1 \leq B_2$. We shall give examples of this below. The next theorem gives the spectral gap and the eigenfunctions whenever $\Delta_2 = 2(1 - B_1)$, regardless of whether $B_1 < B_2$ or $B_2 < B_1$. However, it gives the exact value of Δ_N only for $N \geq 7$. Since we are interested in large values of N , this is fully satisfactory. Indeed, it is remarkable that the two theorems show that already at relatively small values of N , the behavior of the system is very close, qualitatively and quantitatively to the behavior in the large N limit.

THEOREM 1.2. *Suppose that $\Delta_2 = 2(1 - B_1)$. Then, for all $N \geq 7$,*

$$(1.10) \quad \Delta_N = \min\{ (1 - B_1), (1 - B_2) \} \frac{N}{(N - 1)}.$$

Moreover, if $B_2 > B_1$, the eigenspace is three-dimensional and is spanned by the functions

$$(1.11) \quad \phi(\vec{v}) = \sum_{j=1}^N |v_j|^2 v_j^\alpha, \quad \alpha = 1, 2, 3,$$

where v_j^α denotes the α th component of v_j .

On the other hand, if $B_2 < B_1$, the eigenspace is spanned by the functions of the form

$$(1.12) \quad |v_i|^2 - |v_j|^2 \quad \text{and} \quad v_i^\alpha - v_j^\alpha, \quad \alpha = 1, 2, 3,$$

for all $i < j$.

Finally, if $B_1 = B_2$, the eigenspace is spanned by both of the sets of functions listed in (1.8) and (1.12) together.

For the family of collision rates introduced so far, the b_α , one may apply Theorem 1.1, as we have indicated, but only for $\alpha \leq 7/9$. As we shall see in section 2, Theorem 1.2 applies for all $0 \leq \alpha < 1$ and in this case gives $\Delta_N = (N/N - 1)(1 - B_2)$ for $N \geq 7$. In order to illustrate the case in which Theorem 1.2 yields the gap $\Delta_N = (N/N - 1)(1 - B_1)$, we introduce

$$(1.13) \quad \tilde{b}_\alpha(x) = 2(\alpha + 1)1_{[0,1]}(x)x^\alpha, \quad \alpha \geq 0.$$

Since $x^2 < x$ on $(0, 1)$, it is clear that $B_2 < B_1$ for all α in this case. We show at the end of section 2 that, at least for $0 \leq \alpha \leq 1$, $\Delta_2 = 2(1 - B_1)$, so that Theorem 1.2 applies in these cases.

The method of proof is quite robust, and in section 10 we shall describe how it may be extended to determine the spectral gap of Q for still other choices of b that are not covered by Theorems 1.1 and 1.2.

The method of proof of these theorems relies on a basic strategy introduced in [3] but which is extended significantly here. The strategy consists of exploiting an inductive link between the spectral gap of Q and the one of an operator P , an average over projections introduced in section 3. In fact,

$$(1.14) \quad \Delta_N \geq \frac{N}{N-1}(1 - \mu_N)\Delta_{N-1},$$

where $1 - \mu_N$ is the gap of P . The eigenvalues of P are much easier to compute than the ones of Q since the range of P consists of sums of functions of single variables v_j .

In the case of the original model treated by Kac, one is in the happy circumstance that Q has a single gap eigenfunction ϕ which is also the gap eigenfunction of P for all N , and, when this is used as a trial function in the derivation of (1.14), one sees that (1.14) actually holds with equality, giving an identity relating Δ_N and Δ_{N-1} . Thus, starting at $N = 2$, where the gap can be easily calculated, the above formula yields a lower bound on Δ_N that turns out to be exact. The model treated in this paper does not have this simplifying feature, even when the gap eigenfunctions of Q are also the gap eigenfunctions of P . Nevertheless, the ideas that lead to (1.14) can be used in such a way that we can still calculate the gap of Q exactly. Very briefly, here is how.

Let $\mu_N^* < \mu_N$ be any number, and assume that there are finitely many eigenvalues $\mu_N^* \leq \mu_N^{(m)} \leq \dots \leq \mu_N^{(1)} \leq \mu_N$ of P . Denote the corresponding eigenspaces by E_j . Let V_j be the smallest invariant subspace of Q that contains E_j . Lemma 4.1 in section 4 provides us with the following dichotomy: *Either*

$$(1.15) \quad \Delta_N \geq \frac{N}{N-1}(1 - \mu_N^*)\Delta_{N-1}$$

or else

$$(1.16) \quad \text{the gap of } Q \text{ is the same as the gap of } Q \text{ restricted to } \bigoplus_{j=1}^m V_j.$$

If the threshold has been chosen so that the lower bound on Δ_N provided by (1.15) is at least as large as the upper bound on Δ_N provided by some trial function in $\bigoplus_{j=1}^m V_j$, then Δ_N is the gap of Q restricted to $\bigoplus_{j=1}^m V_j$. As we shall see, the V_j are finite-dimensional, so determining the gap of Q on $\bigoplus_{j=1}^m V_j$ is a tractable problem. In this case we have determined the exact value of Δ_N .

To proceed to the determination of Δ_N for all large N , one needs a strategy for choosing the threshold μ_N^* . The lower the value of μ_N^* that is chosen, the stronger the bound (1.15) will be, but also the higher the value of m will be. The basis for the choice of μ_N^* is a trial function calculation, providing a guess $\tilde{\Delta}_N$ for the value of Δ_N . Indeed, natural trial functions can often be chosen on the basis of physical considerations. (The spectrum of the linearized Boltzmann equation is the source in the case at hand.) To show that the guess is correct, so that $\tilde{\Delta}_N = \Delta_N$, we are led to choose μ_N^* so that

$$(1.17) \quad \tilde{\Delta}_N \leq \frac{N}{N-1}(1 - \mu_N^*)\tilde{\Delta}_{N-1}.$$

Since $\tilde{\Delta}_{N-1} \geq \Delta_{N-1}$, this forces us into the second alternative in the dichotomy discussed above, so that the gap eigenfunction for N particles lies in $\oplus_{j=1}^m V_j$. Indeed, if the physical intuition behind the guess was correct, the trial function leading to $\tilde{\Delta}_N$ will lie in $\oplus_{j=1}^m V_j$ and yield the gap.

Choosing μ_N^* small enough that (1.17) is satisfied might in principle lead to a value of m that depends on N . However, in the case at hand, we are fortunate and can work with a choice of μ_N^* that leads to a fixed and small value of m but for which (1.17) is satisfied for all sufficiently large values of N —hence the restriction to $N \geq 7$ in Theorem 1.2.

As will be clear from this summary of the strategy, the determination of the spectrum of P is the main technical step that must be accomplished. As we mentioned before, this is relatively simple, compared to the determination of the spectrum of Q , since the range of P consists of functions that are a sum of functions of a single variable.

For this reason, we can reduce the study of the spectrum of P to that of a much simpler Markov operator K acting on functions on the unit ball B in \mathbb{R}^3 . In the analysis of K , we shall draw on some deep results on Jacobi polynomials [8], [12]. In fact, it turns out that the connection between our eigenvalue problems and pointwise bounds on Jacobi polynomials is through a simple identity, and applications of this identity can be made in both directions: We not only use bounds on Jacobi polynomials to bound eigenvalues, we shall use simple eigenvalue estimates to sharpen certain best known bounds on Jacobi polynomials, as we briefly discuss in section 11.

First, however, we deal with a simpler technical problem, the computation of the spectral gap of Q for $N = 2$.

2. The spectral gap for $N = 2$. For $N = 2$, the state space X_2 consists of pairs $(v, -v)$, with $v \in \mathbb{R}^3$ satisfying $|v|^2 = 1/2$. Note that for $N = 2$ the collision rules (1.1) reduce to

$$v_1^* = \sigma/\sqrt{2} \quad \text{and} \quad v_2^* = -\sigma/\sqrt{2},$$

since $v_1 + v_2 = 0$.

The map $(v, -v) \mapsto \sqrt{2}v$ identifies X_2 with the unit sphere S^2 , and the measure $d\sigma_2$ on X_2 with $d\sigma$ on S^2 . Thus, we may think of Q as operating on functions on S^2 . In this representation, we have the formula

$$Q\phi(u) = \int_{S^2} \phi(\sigma)b(u \cdot \sigma)d\sigma.$$

Notice that if R is any rotation of \mathbb{R}^3

$$\begin{aligned} (Q\phi)(Ru) &= \int_{S^2} \phi(\sigma)b(Ru \cdot \sigma)d\sigma \\ &= \int_{S^2} \phi(R\sigma)b(Ru \cdot R\sigma)d\sigma \\ &= \int_{S^2} \phi(R\sigma)b(u \cdot \sigma)d\sigma = Q(\phi \circ R)(u). \end{aligned}$$

That is, $(Q\phi) \circ R = Q(\phi \circ R)$, and this means that for each n the space of spherical harmonics of degree n is an invariant subspace of Q , contained in an eigenspace of Q . In turn, this means that we can determine the spectrum of Q by computing its

action on the zonal spherical harmonics, i.e., those of the form $P_n(e \cdot u)$, where e is any fixed unit vector in \mathbb{R}^3 and P_n is the n th degree Legendre polynomial. Now, for any function $\phi(u)$ of the form $\phi(u) = f(e \cdot u)$,

$$Q\phi(u) = \int_{S^2} \phi(\sigma \cdot e) b(\sigma \cdot u) d\sigma.$$

We choose coordinates in which u and e span the x, z plane with

$$u = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \text{and} \quad e = \begin{bmatrix} \sqrt{1-t^2} \\ 0 \\ t \end{bmatrix},$$

so that $t = u \cdot e$. Then with

$$\sigma = \begin{bmatrix} \sin \theta \cos \psi \\ \sin \theta \sin \psi \\ \sin \theta \end{bmatrix},$$

$$Q\phi(u) = \mathcal{Q}f(e \cdot u),$$

where

$$\begin{aligned} \mathcal{Q}f(t) &= \frac{1}{4\pi} \int_0^\pi \int_0^{2\pi} f(t \cos \theta + \sqrt{1-t^2} \sin \theta \cos \varphi) b(\cos \theta) \sin \theta d\theta d\varphi \\ (2.1) \quad &= \frac{1}{4\pi} \int_0^\pi \int_{-1}^1 f(ts + \sqrt{1-t^2} \sqrt{1-s^2} \cos \varphi) b(s) ds d\varphi. \end{aligned}$$

Now, if f is any eigenfunction of \mathcal{Q} , with $\mathcal{Q}f = \lambda f$, then by evaluating both sides at $t = 1$, we have $\lambda f(1) = \frac{1}{2} \int_0^\pi \int_{-1}^1 f(s) b(s) ds$. Thus, the eigenvalue is given by

$$\lambda = \frac{1}{2} \int_{-1}^1 \frac{f(s)}{f(1)} b(s) ds.$$

As we have observed above, the eigenfunctions of \mathcal{Q} are the Legendre polynomials. Thus, if P_n is the Legendre polynomial of n th degree with the standard normalization $P_n(1) = 1$, and λ_n is the corresponding eigenvalue,

$$(2.2) \quad \lambda_n = \frac{1}{2} \int_{-1}^1 P_n(s) b(s) ds.$$

This explicit formula enables one to easily compute Δ_2 . For example, we can now easily prove the following.

LEMMA 2.1. *When $b(x) = b_\alpha(x)$, as in (1.3), then $1 - B_2 < 1 - B_1$ for all $\alpha < 1$, and moreover*

$$(2.3) \quad \Delta_2 = 2(1 - \lambda_1) = \frac{4(1 - \alpha)}{2 - \alpha} = (1 - B_2)(3 - \alpha),$$

so that (1.6) is satisfied for all α with $0 \leq \alpha \leq 7/9$.

Proof. By using Rodrigues's formula [13, equation 4.3.1]

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n$$

and integration by parts, one computes

$$\lambda_n = (1 - \alpha) \frac{(\alpha)_n}{(1 - \alpha)_{n+1}} = \frac{(\alpha)_n}{(2 - \alpha)_n},$$

where $(\alpha)_n = \alpha(\alpha + 1)(\alpha + 2) \dots (\alpha + n - 1)$. Notice that, for all $0 \leq \alpha < 1$, λ_n decreases as n increases, so with the collision rate given by b_α ,

$$(2.4) \quad \Delta_2 = 2(1 - \lambda_1) = \frac{4(1 - \alpha)}{2 - \alpha}.$$

Next, one computes

$$1 - B_1 = \frac{2(1 - \alpha)}{(2 - \alpha)} \quad \text{and} \quad 1 - B_2 = \frac{4(1 - \alpha)}{(2 - \alpha)(3 - \alpha)}.$$

Since $2 > 4/(3 - \alpha)$ for $\alpha < 1$, $1 - B_2 < 1 - B_1$ for all $\alpha < 1$. Moreover, from this computation, one readily obtains (2.3) and the statement concerning (1.6). \square

In particular, the condition (1.6) is satisfied in both the uniform redirection model ($\alpha = 0$) and the Morgenstern model ($\alpha = 1/2$). Thus in these cases we have the exact spectral gaps

$$(2.5) \quad \begin{aligned} \Delta_N &= \frac{2}{3} \frac{N}{N - 1} && \text{for the uniform redirection model,} \\ \Delta_N &= \frac{8}{15} \frac{N}{N - 1} && \text{for the Morgenstern model.} \end{aligned}$$

We close this section with a remark that may provide a useful perspective on what follows. In determining the spectral gap of Q for $N = 2$, general symmetry conditions told us right away what all of the eigenfunctions were. A less obvious, though still simple, argument then provided us with the explicit formula (2.2) for all of the eigenvalues. There is one last hurdle to cross: There are infinitely many eigenvalues given by (2.2), and for a general b , we cannot determine which is the second largest by computing them all explicitly. What was particularly nice about b_α is that in this case the eigenvalues of Q were monotone-decreasing:

$$\lambda_{n+1} \leq \lambda_n.$$

For other choices of b , this need not be the case. However, there are ways to use pointwise bounds on Legendre polynomials to reduce the problem of determining Δ_2 to the computation of a *finite* number of eigenvalues by using (2.2). For example, one has the classical bound (see [13, Theorem 7.3.3]):

$$(2.6) \quad |P_n(x)|^2 < \frac{2}{n\pi} \frac{1}{\sqrt{1 - x^2}}.$$

As long as $b(x)(1 - x^2)^{-1/4}$ is integrable, this gives an upper bound on λ_n that is proportional to $n^{-1/2}$: Define

$$\tilde{\lambda}_n = \left(\frac{1}{8\pi n} \right)^{1/2} \int_{-1}^1 b(x)(1 - x^2)^{-1/4} dx.$$

Then let n_0 be the least value of n such that $\tilde{\lambda}_n \leq \lambda_1$. Then the second largest eigenvalue of Q is

$$\max_{1 \leq n \leq n_0} \lambda_n.$$

We illustrate this by showing that, for the rate function \tilde{b}_α introduced in (1.13), $\Delta_2 = 2(1 - B_1)$ at least for $0 \leq \alpha \leq 1$. (Of course, the integrals in (2.2) can be computed exactly in this case; see [5, section 7.231, page 822]. However, we prefer to illustrate the use of (2.6).)

By (2.6) and (2.2),

$$(2.7) \quad \begin{aligned} |\lambda_n| &\leq (\alpha + 1) \left(\int_0^1 x^{2\alpha} dx \right)^{1/2} \left(\int_0^1 P_n(x)^2 dx \right)^{1/2} \\ &< \frac{\alpha + 1}{\sqrt{2\alpha + 1}} \frac{1}{\sqrt{n}}. \end{aligned}$$

Also, by (2.2), $\lambda_1 = B_1 = (\alpha + 1)/(\alpha + 2)$. Comparison of the formulas shows that, for $0 \leq \alpha \leq 1$, $\lambda_n < \lambda_1$ for all $n > 4$. Thus it suffices to check that $\lambda_j < \lambda_1$ for $j = 2, 3$, and 4 by direct computation with (2.2). By doing so, one finds that this is the case. Hence, Theorem 1.2 applies and yields $\Delta_N = (N/N - 1)(1 - B_1)$ for $N \geq 7$.

Further calculation would extend this result to higher values of α . Notice that, as α tends to infinity, $\tilde{b}_\alpha(x)$ is more and more concentrated at $x = 1$, which corresponds to $\theta = 0$. Thus, for large values of α , \tilde{b}_α represents a “grazing collision model.”

For $N > 2$, the operator Q is much more complicated, and direct determination of the spectrum is not feasible. Instead, we use an inductive procedure involving an auxiliary operator that we now introduce.

3. The average of projections operator P and its relation to Q . A simple convexity argument shows that, for each j ,

$$\sup\{|v_j|^2 : \vec{v} \in X_N\} = \frac{N - 1}{N}.$$

For each j , define $\pi_j(\vec{v})$ by

$$\pi_j(\vec{v}) = \sqrt{\frac{N}{N - 1}} v_j,$$

so that π_j maps X_N onto the unit ball B in \mathbb{R}^3 .

For any function ϕ in $L^2(X_N, d\sigma_N)$, and any j with $1 \leq j \leq N$, define $P_j(\phi)$ to be the orthogonal projection of ϕ onto the subspace of $L^2(X_N, d\sigma_N)$ consisting of square integrable functions that depend on \vec{v} through v_j alone. That is, $P_j(\phi)$ is the unique element of $L^2(X_N, d\sigma_N)$ of the form $f(\pi_j(\vec{v}))$ such that

$$\int_{X_N} \phi(\vec{v})g(\pi_j(\vec{v}))d\sigma_N = \int_{X_N} f(\pi_j(\vec{v}))g(\pi_j(\vec{v}))d\sigma_N$$

for all continuous functions g on B .

The *average of projections operator* P is then defined through

$$P = \frac{1}{N} \sum_{j=1}^N P_j.$$

If the individual projections P_j all commuted with one another, then the spectrum of P would be very simple: The eigenvalues of each P_j are 0 and 1. Moreover, $P_j\phi = \phi$ if and only if ϕ depends only on v_j so that it cannot then also satisfy $P_k\phi = \phi$ for

$k \neq j$, unless ϕ is constant. It would then follow that the eigenvalues of P would be 0, $1/N$, and 1, with the last having multiplicity one.

However, the individual projections P_j do not commute with one another, due to the nature of the constraints defining X_N .

We now define

$$(3.1) \quad \mu_N = \sup \left\{ \int_{X_N} \phi(\vec{v})P\phi(\vec{v})d\sigma_N \mid \int_{X_N} \phi^2(\vec{v})d\sigma_N = 1, \int_{X_N} \phi(\vec{v})d\sigma_N = 0 \right\}.$$

The P operator is simpler than the Q operator in that if ϕ is any eigenfunction of P with a nonzero eigenvalue, then clearly ϕ has the form

$$\phi = \sum_{j=1}^N f_j \circ \pi_j$$

for some functions f_1, \dots, f_N on B . For $N \geq 4$, most of the eigenfunctions of Q have a more complicated structure. Nonetheless, there is a close relation between the spectra of Q and P , as we now explain.

To do this, we need a more explicit formula for P , such as the formula (1.4) that we have for Q . The key to computing $P_j\phi$ is a factorization formula [3] for the measure σ_N . Define a map $T_N : X_{N-1} \times B \rightarrow X_N$ as follows:

$$(3.2) \quad T_N(\vec{y}, v) = \left(\alpha(v)y_1 - \frac{1}{\sqrt{N^2 - N}}v, \dots, \alpha(v)y_{N-1} - \frac{1}{\sqrt{N^2 - N}}v, \sqrt{\frac{N-1}{N}}v \right),$$

where

$$\alpha^2(v) = 1 - |v|^2.$$

This map induces coordinates (\vec{y}, v) on X_N , and, in terms of these coordinates, one has the integral factorization formula

$$\int_{X_N} \phi(\vec{v})d\sigma_N = \frac{|S^{3N-7}|}{|S^{3N-4}|} \int_B \left[\int_{X_{N-1}} \phi(T_N(\vec{y}, v))d\sigma_{N-1} \right] (1 - |v|^2)^{(3N-8)/2} dv.$$

It follows from this and the definition of P_N that

$$P_N\phi(\vec{v}) = f \circ \pi_N(\vec{v}),$$

where

$$f(v) = \int_{X_{N-1}} \phi(T_N(\vec{y}, v))d\sigma_{N-1}.$$

For $j < N$, one has analogous formulas for T_j and P_j , except the roles of v_N and v_j are interchanged.

Next, we make the definition for Q that is analogous to (3.1) for P : Define λ_N by

$$(3.3) \quad \lambda_N = \sup \left\{ \int_{X_N} \phi(\vec{v})Q\phi(\vec{v})d\sigma_N \mid \int_{X_N} \phi^2(\vec{v})d\sigma_N = 1, \int_{X_N} \phi(\vec{v})d\sigma_N = 0 \right\}.$$

With this explicit formula in hand, and the definitions of μ_N and λ_N , we come to the fundamental fact relating P and Q .

LEMMA 3.1. *For any square integrable function ϕ on X_N that is orthogonal to the constants,*

$$(3.4) \quad \langle \phi, Q\phi \rangle \leq \lambda_{N-1} \|\phi\|_2^2 + (1 - \lambda_{N-1}) \langle \phi, P\phi \rangle,$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product on $L^2(X_N, \sigma_N)$.

Proof. To bound $\langle \phi, Q\phi \rangle$ in terms of λ_{N-1} , define, for $1 \leq k \leq N$, the operator $Q^{(k)}$ by

$$Q^{(k)}\phi(\vec{v}) = \binom{N-1}{2}^{-1} \sum_{i < j, i \neq k, j \neq k} \int_{S^2} \phi(R_{i,j,\sigma}(\vec{v})) d\sigma.$$

That is, we leave out collisions involving the k th particle and average over the rest. Clearly,

$$Q = \frac{1}{N} \sum_{k=1}^N Q^{(k)}.$$

Therefore, for any ϕ in $L^2(X_N, \sigma_N)$,

$$\langle \phi, Q\phi \rangle = \frac{1}{N} \sum_{k=1}^N \langle \phi, Q^{(k)}\phi \rangle.$$

By using the coordinates (\vec{y}, v) induced by the map $T_k : X_{N-1} \times B \rightarrow X_N$, it is easy to see that, for $i \neq k, j \neq k$, $R_{i,j,\sigma}$ acts only on the \vec{y} variable. That is, for such i and j ,

$$R_{i,j,\sigma}(T_k(\vec{y}, v)) = T_k(R_{i,j,\sigma}(\vec{y}), v).$$

Thus, if we hold v fixed as a parameter, we can think of $(Q^{(k)}\phi)(T_k(\vec{y}, v))$ as resulting from applying the $N - 1$ dimensional version of Q to ϕ with v_k held fixed.

To estimate λ_N , we need to estimate $\langle \phi, Q\phi \rangle$ when ϕ is orthogonal to the constants. When ϕ is orthogonal to the constants, and we fix v , the function

$$\vec{y} \mapsto \phi(T_k(\vec{y}, v))$$

is not, in general, orthogonal to the constants on X_{N-1} . However, we can correct for that by adding and subtracting $P_k\phi$. Therefore

$$(3.5) \quad \begin{aligned} \langle (\phi - P_k\phi), Q^{(k)}(\phi - P_k\phi) \rangle &\leq \lambda_{N-1} \|\phi - P_k\phi\|_2^2 \\ &= \lambda_{N-1} (\|\phi\|_2^2 - \|P_k\phi\|_2^2) \\ &= \lambda_{N-1} (\|\phi\|_2^2 - \langle \phi, P_k\phi \rangle). \end{aligned}$$

Then since $Q^{(k)}P_k\phi = P_k\phi$ and since $P_k\phi$ is orthogonal to $\phi - P_k\phi$,

$$(3.6) \quad \begin{aligned} \langle \phi, Q^{(k)}\phi \rangle &= \langle (\phi - P_k\phi) + P_k\phi, Q^{(k)}((\phi - P_k\phi) + P_k\phi) \rangle \\ &= \langle (\phi - P_k\phi), Q^{(k)}(\phi - P_k\phi) \rangle + \langle P_k\phi, P_k\phi \rangle \\ &= \langle (\phi - P_k\phi), Q^{(k)}(\phi - P_k\phi) \rangle + \langle \phi, P_k\phi \rangle \\ &\leq \lambda_{N-1} (\|\phi\|_2^2 - \langle \phi, P_k\phi \rangle) + \langle \phi, P_k\phi \rangle. \end{aligned}$$

By averaging over k , we have (3.4). \square

Lemma 3.1 was used as follows in [3]: Any trial function ϕ for λ_N is a valid trial function for μ_N , so that

$$(3.7) \quad \lambda_N \leq \lambda_{N-1} + (1 - \lambda_{N-1})\mu_N.$$

Then since $\Delta_N = N(1 - \lambda_N)$, we have

$$(3.8) \quad \Delta_N \geq \frac{N}{N-1}(1 - \mu_N)\Delta_{N-1}.$$

Therefore, with $a_N = \frac{N}{N-1}(1 - \mu_N)$, for all $N \geq 3$,

$$\Delta_N \geq \left(\prod_{j=3}^N a_j \right) \Delta_2.$$

Thus, one route to proving a lower bound on Δ_N is to prove an upper bound on μ_N and hence a lower bound on a_N . This route led to a sharp lower bound for Δ_N —the exact value—for the one dimension Kac model investigated in [2]. However, it would not lead to a proof of Theorem 1.1. The reasons for this are worth pointing out before we proceed.

As we shall see below, the eigenspace of P with the eigenvalue μ_N —the gap eigenspace of P —is spanned by the functions specified in (1.8). Granted this, and granted Theorem 1.1, whenever condition (1.6) is satisfied:

$$\text{For } (1 - B_2) < (1 - B_1), \quad Q\phi = \lambda_N\phi \quad \Rightarrow \quad P\phi = \mu_N\phi,$$

while

$$\text{for } (1 - B_1) < (1 - B_2), \quad Q\phi = \lambda_N\phi \quad \Rightarrow \quad P\phi \neq \mu_N\phi.$$

In the second case $(1 - B_1) < (1 - B_2)$, the mismatch between the gap eigenspaces for Q and P means that equality cannot hold in (3.7), and hence the recursive relation (3.8) cannot possibly yield exact results in this case.

Moreover, even in the first case $(1 - B_2) < (1 - B_1)$, where there is a match between the gap eigenspaces of Q and P , there *still* will not be equality in (3.7). The reasons for this are more subtle: The inequality (3.7) comes from the key estimate (3.6). By considering (3.6), one sees that equality will hold there if and only if

$$Q^{(k)}(\phi - P_k\phi) = \lambda_{N-1}(\phi - P_k\phi)$$

for each k , where $(\phi - P_k\phi)$ is regarded as a function on X_{N-1} through the change of variables $T_k : (X_{N-1}, B) \rightarrow X_N$ that was introduced just before Lemma 3.1.

However, if ϕ is in the gap eigenspace for Q on X_N , Theorem 1.1 tells us that it is a linear combination of the three functions specified in (1.8), all of which are homogeneous of degree 3 in v . Because of the translation in (3.2), which is due to momentum conservation, $(\phi - P_k\phi)$ is regarded as a function on X_{N-1} that will *not* be homogeneous of degree 3—it will contain lower order terms. Hence $(\phi - P_k\phi)$ will not be in the gap eigenspace for $Q^{(k)}$.

The main result of the next section provides a way to use more detailed spectral information about P to sharpen the recursive estimate so that we do obtain the exact results announced in Theorem 1.1.

4. How to use more detailed spectral information on P to determine the spectral gap of Q . The following lemma is the key to using (3.4) to obtain sharp results for the model considered here.

LEMMA 4.1. *For any $N \geq 3$, let μ_N^* be a number with*

$$\mu_N^* < \mu_N$$

such that there are only finitely eigenvalues of P between μ_N^ and μ_N :*

$$\mu_N^* \leq \mu_N^{(m)} < \dots < \mu_N^{(1)} < \mu_N.$$

Let $\mu_N^{(0)}$ denote μ_N , and then, for $j = 0, \dots, m$, let E_j denote the eigenspace of P corresponding to $\mu_N^{(j)}$. Let V_j denote the smallest invariant subspace of Q that contains E_j . Let ν_j be the largest eigenvalue of Q on V_j .

Then either

$$(4.1) \quad \lambda_N = \max\{\nu_0, \dots, \nu_m\}$$

or else

$$(4.2) \quad \Delta_N \geq \frac{N}{N-1}(1 - \mu_N^*)\Delta_{N-1}.$$

If $\mu_N^ = \mu_N^{(m)}$, then we have the same alternative except with strict inequality in (4.2).*

Proof. If $\lambda_N > \max\{\nu_0, \dots, \nu_m\}$, then in the variational principle for λ_N we need only consider functions ϕ that are orthogonal to the constants and also in the orthogonal complement of each of the V_j . This means also that ϕ belongs to the orthogonal complement of each of the E_j . But then

$$\langle \phi, P\phi \rangle \leq \mu_N^* \|\phi\|_2^2.$$

By using this estimate in (3.4), we have (4.2). Moreover, if $\mu_N^* = \mu_N^{(m)}$, then strict inequality must hold in the last inequality. \square

Lemma 4.1 gives us the dichotomy between (1.15) and (1.16) that plays a key role in the strategy described in the introduction. To put this strategy into effect, we must first carry out a more detailed investigation of the spectrum of P . The main result of the next section reduces the investigation of the spectrum of P to the study of a simpler operator—the *correlation operator* K , which is a Markov operator on functions on the unit ball B in \mathbb{R}^3 .

5. The correlation operator K and its relation to P . While Q and P are both operators on spaces of functions of a large number of variables, the problem of computing the eigenvalues of P reduces to the problem of computing the eigenvalues of an operator on functions on B , the unit ball in \mathbb{R}^3 .

First, define the measure ν_N on B to be the “push forward” of σ_N under the map π_j . That is, for any continuous function f on B ,

$$\int_B f(v) d\nu_N = \int_{X_N} f(\pi_j(\vec{v})) d\sigma_N.$$

By the permutation invariance of σ_N , this definition does not depend on the choice of j . By direct calculation [3], one finds that

$$(5.1) \quad d\nu_N(v) = \frac{|S^{3N-7}|}{|S^{3N-4}|} (1 - |v|^2)^{(3N-8)/2} dv.$$

Now define the self-adjoint operator K on $L^2(B, d\nu_N)$ through the following quadratic form:

$$(5.2) \quad \langle f, Kf \rangle_{L^2(\nu)} = \int_{X_N} f(\pi_1(\vec{v}))f(\pi_2(\vec{v}))d\sigma_N$$

for all f in $L^2(B, d\nu_N)$. Equivalently,

$$(5.3) \quad (Kf) \circ \pi_1 = P_1(f \circ \pi_2).$$

Note that, by the permutation invariance of σ_N , one can replace the pair $(1, 2)$ of indices by any other pair of distinct indices without affecting the operator K defined by (5.3). This is the *correlation operator*.

To see the relation between the spectra of P and the spectra of K , suppose that ϕ is an eigenfunction of P that is symmetric under permutation of the particle indices. (These symmetric eigenfunctions are the ones that are significant in the physical application.) Then since any vector in the image of P has the form $\sum_{j=1}^N f_j \circ \pi_j$ for functions f_1, \dots, f_N on B , we must have, for ϕ symmetric,

$$(5.4) \quad \phi = \sum_{j=1}^N f \circ \pi_j.$$

Now we ask: For which choices of f will ϕ given by (5.4) be an eigenfunction of P ? To answer this, note that, by (5.3),

$$(5.5) \quad P_k\phi = f \circ \pi_k + \sum_{j=1, j \neq k}^N P_k(f \circ \pi_j).$$

Therefore, from (5.5) and the definition of K , $P_k\phi = f \circ \pi_k + (N - 1)(Kf) \circ \pi_k$. Thus, by averaging over k ,

$$(5.6) \quad P\phi = \frac{1}{N}\phi + \frac{N - 1}{N} \sum_{j=1}^N (Kf) \circ \pi_j.$$

In the case $Kf = \kappa f$, this reduces to

$$P\phi = \frac{1}{N}(1 + (N - 1)\kappa)\phi,$$

and thus eigenfunctions of K yield eigenfunctions of P . It turns out that all symmetric eigenfunctions arise in exactly this way and that all eigenfunctions, symmetric or not, arise in a similar way, specified in the next lemma.

LEMMA 5.1. *Let V be the orthogonal complement in $L^2(X_N, \sigma_N)$ of the kernel of P . There is a complete orthonormal basis of V consisting of eigenfunctions ϕ of P of one of the two forms:*

(i) *For some eigenfunction f of K , $\phi = \sum_{k=1}^N f \circ \pi_k$. In this case, if $Kf = \kappa f$, then $P\phi = \mu\phi$, where*

$$(5.7) \quad \mu = \frac{1}{N}(1 + (N - 1)\kappa).$$

(ii) For some eigenfunction f of K , and some pair of indices $i < j$, $\phi = f \circ \pi_i - f \circ \pi_j$. In this case, if $Kf = \kappa f$, then $P\phi = \mu\phi$, where

$$(5.8) \quad \mu = \frac{1 - \kappa}{N}.$$

Proof. Suppose that ϕ is an eigenfunction of P with nonzero eigenvalue μ , and ϕ is orthogonal to the constants. By the permutation invariance we may assume that either ϕ is invariant under permutations or that there is some pair permutation, which we may as well take to be $\sigma_{1,2}$, such that $\phi \circ \sigma_{1,2} = -\phi$. We will treat these two cases separately.

First, suppose that ϕ is symmetric. We have already observed that, in this case, the recipe $\phi = \sum_{j=1}^N f \circ \pi_j$, with f an eigenfunction of K , yields symmetric eigenfunctions of P . We now show that all symmetric eigenfunctions of P on V have this form.

First, simply because such a ϕ is in the image of P and is symmetric, ϕ must have the form (5.4). It remains to show that f must be an eigenfunction of K . Then by (5.6), $\mu\phi = P\phi$ becomes

$$\mu \sum_{k=1}^N f \circ \pi_k = \frac{1}{N} \sum_{k=1}^N (f + (N - 1)Kf) \circ \pi_k.$$

Apply P_1 to both sides to obtain

$$\frac{1}{N} ([f + (N - 1)Kf] + (N - 1)K[f + (N - 1)Kf]) = \mu(f + (N - 1)Kf),$$

which is

$$(5.9) \quad \frac{1}{N} (I + (N - 1)K)^2 f = \mu(I + (N - 1)K)f.$$

Since $\mu \neq 0$, f is not in the null space of either $I + (N - 1)K$ or $(I + (N - 1)K)^2$. It then follows from (5.9) that

$$\frac{1}{N} (I + (N - 1)K) f = \mu f.$$

Thus, when ϕ is symmetric, there is an eigenfunction f of K with eigenvalue κ such that $\phi = \sum_{k=1}^N f \circ \pi_k$ and

$$\mu = \frac{1}{N} (1 + (N - 1)\kappa).$$

We next consider the case in which

$$\phi \circ \sigma_{1,2} = -\phi.$$

Note that

$$P_k(\phi \circ \sigma_{1,2}) = P_k\phi = 0$$

whenever k is different from both 1 and 2. It follows that

$$\frac{1}{N} \sum_{k=1}^N P_k\phi = \frac{1}{N} (P_1\phi + P_2\phi).$$

The right-hand side is of the form $f(v_1) - f(v_2)$, and hence ϕ must have this form if it is an eigenvector. By taking $\phi = f \circ \pi_1 - f \circ \pi_2$ we have

$$\frac{1}{N} \sum_{k=1}^N P_k \phi = \frac{1}{N} ((f - Kf) \circ \pi_1 - (f - Kf) \circ \pi_2).$$

Hence when $P\phi = \mu\phi$ and ϕ is antisymmetric as above, there is an eigenvalue κ of K such that

$$\mu = \frac{1 - \kappa}{N}.$$

This proves the second part. \square

Lemma 5.1 reduces the computation of the spectrum of P to the computation of the spectrum of K . We undertake this in the next three sections.

6. Explicit form of the correlation operator K . For any two functions f and g on B that are square integrable with respect to ν_N , consider the bilinear form $\int_{X_N} f(\pi_1(\vec{v}))g(\pi_2(\vec{v}))d\sigma_N$. It is easily seen from (5.3) that

$$\langle f, Kg \rangle = \int_{X_N} f(\pi_1(\vec{v}))g(\pi_2(\vec{v}))d\sigma_N,$$

where $\langle \cdot, \cdot \rangle$ is the inner product on $L^2(B, \nu_N)$.

By computing the right-hand side using the factorization formula (3.2), but for T_1 instead of T_N , one finds, for $N > 3$:

$$Kg(v) = \frac{|S^{3N-10}|}{|S^{3N-7}|} \int_B g \left(\frac{\sqrt{N^2 - 2N}}{N - 1} \sqrt{1 - |v|^2}y - \frac{1}{N - 1}v \right) (1 - |y|^2)^{(3N-11)/2} dy.$$

The explicit form of K is slightly different for $N = 3$. We can see this different form as a limiting case, if we make the dimension a continuous fact. The following way of doing this will be convenient later on.

For $\alpha > -1$, define the constant C_α by

$$C_\alpha = \left(\int_B (1 - |y|^2)^\alpha dy \right)^{-1},$$

so that, for

$$\alpha = \frac{3N - 8}{2},$$

$$d\nu_N(v) = C_\alpha(1 - |y|^2)^\alpha dy,$$

and then

$$Kg(v) = C_{\alpha-3/2} \int_B g \left(\frac{\sqrt{N^2 - 2N}}{N - 1} \sqrt{1 - |v|^2}y - \frac{1}{N - 1}v \right) (1 - |y|^2)^{\alpha-3/2} dy.$$

Now, as N approaches 3, $\alpha - 3/2$ approaches -1 . Then the measure $C_\alpha(1 - |y|^2)^\alpha dy$ concentrates more and more on the boundary of the ball B , so that, in the limit, it becomes the uniform measure on S^2 . Understood in this way, the formula remains valid at $\alpha = 1/2$, i.e., at $N = 3$.

It is clear that K is a self-adjoint Markov operator on $L^2(B, \nu_N)$ and that 1 is an eigenvalue of multiplicity one. With more effort, there is much more that can be said; the spectrum of K can be completely determined.

7. The spectrum of K and ratios of Jacobi polynomials. In studying the spectrum of the correlation operator, it is in fact natural and useful to study a wider family of operators of this type. Fix any $\alpha > 1/2$ and any numbers a and b such that

$$a^2 + b^2 = 1.$$

Then define the generalized correlation operator, still simply denoted by K , through

$$(7.1) \quad Kg(v) = C_{\alpha-3/2} \int_B g \left(a\sqrt{1-|v|^2}y + bv \right) (1-|y|^2)^{\alpha-3/2} dy.$$

Notice that, as v and y range over B , the maximum of $|a\sqrt{1-|v|^2}y + bv|$ occurs when ay and bv are parallel. In that case,

$$|a\sqrt{1-|v|^2}y + bv| = |a||y|\sqrt{1-|v|^2} + |b||v| \leq (a^2 + b^2)^{1/2}((1-|v|^2)|y|^2 + |v|^2)^{1/2} \leq 1.$$

Thus, as v and y range over B , so does

$$(7.2) \quad u(y, v) = a\sqrt{1-|v|^2}y + bv,$$

and $g(a\sqrt{1-|v|^2}y + bv)$ is well defined for any function g on B . Thus, K is well defined.

Now when

$$(7.3) \quad a = \frac{\sqrt{N^2 - 2N}}{N - 1} \quad \text{and} \quad b = -\frac{1}{N - 1},$$

we know that K is self-adjoint because in that case it is defined in terms of a manifestly symmetric bilinear form. We shall show here that K is always self-adjoint for all $a^2 + b^2 = 1$ and that the eigenvalues of K are given by an explicit formula involving ratios of Jacobi polynomials.

To explain this, we fix some terminology and notation. For any numbers $\alpha > -1$ and $\beta > -1$, $P_n^{(\alpha, \beta)}$ denotes the n th degree polynomial in the sequence of orthogonal polynomials on $[-1, 1]$ for the measure

$$(1-x)^\alpha(1+x)^\beta dx$$

and is referred to as the n th degree Jacobi polynomial for (α, β) . As is well known, $\{P_n^{(\alpha, \beta)}\}_{n \geq 0}$ is a complete orthogonal basis for $L^2([-1, 1], (1-x)^\alpha(1+x)^\beta dx)$.

Of course, what we have said so far specifies $P_n^{(\alpha, \beta)}$ only up to a multiplicative constant. One common normalization is given by Rodrigues's formula

$$P_n^{(\alpha, \beta)}(x) = \frac{(-1)^n}{2^n n!} (1-x)^{-\alpha} (1+x)^{-\beta} \frac{d^n}{dx^n} ((1-x)^{\alpha+n} (1+x)^{\beta+n}).$$

For this normalization,

$$(7.4) \quad P_n^{(\alpha, \beta)}(1) = \binom{n + \alpha}{n} \quad \text{and} \quad P_n^{(\alpha, \beta)}(-1) = \binom{n + \beta}{n}.$$

LEMMA 7.1. *Fix any $\alpha > 1/2$ and any numbers a and b such that $a^2 + b^2 = 1$, and define K through the formula (7.1). Then K is a self-adjoint Markov operator, and*

the spectrum of K consists of eigenvalues $\kappa_{n,\ell}$ enumerated by nonnegative integers n and ℓ ; these eigenvalues are given by the explicit formula

$$(7.5) \quad \kappa_{n,\ell} = \frac{P_n^{(\alpha,\beta)}(-1 + 2b^2)}{P_n^{(\alpha,\beta)}(1)} b^\ell,$$

where $\beta = \ell + 1/2$ and α is the parameter α entering into the definition of K .

Proof. To see that K is self-adjoint, we write it as a bilinear form and change variables to reveal the symmetry. The change of variable that we make is naturally $(y, v) \rightarrow (u, v)$, with $u(y, v)$ given by (7.2). From (7.2), one computes $y = (u - bv)/(a\sqrt{1 - |v|^2})$, so that

$$(7.6) \quad \begin{aligned} 1 - |y|^2 &= \frac{a^2 - a^2|v|^2 - |u|^2 - b^2|v|^2 + 2bu \cdot v}{a^2(1 - |v|^2)} \\ &= \frac{b^2 - (|u|^2 + |v|^2) + 2bu \cdot v}{a^2(1 - |v|^2)}. \end{aligned}$$

The Jacobian is easy to work out, and one finds that $dudv = a^3(1 - |v|^2)^{3/2}dydv$, so that

$$(7.7) \quad \begin{aligned} &\int_B f(v)Kg(v)C_\alpha(1 - |v|^2)^\alpha dv \\ &= \int_B \int_B f(v)g(u)a^{-2\alpha} [a^2 - (|u|^2 + |v|^2) + 2bu \cdot v]_+^{\alpha-3/2} C_{\alpha-3/2} dudv. \end{aligned}$$

This shows that the operator K is self-adjoint on $L^2(B, C_\alpha(1 - |v|^2)^\alpha)$ for all $\alpha \geq 1/2$ and all a and b with $a^2 + b^2 = 1$.

Our next goal is to prove the eigenvalue formula (7.5). This shall follow from several simple properties of K .

First, K commutes with rotations in \mathbb{R}^3 . That is, if R is a rotation on \mathbb{R}^3 , it is evident that

$$K(g \circ R) = (Kg) \circ R.$$

Hence we may restrict our search for eigenfunctions g of K to functions of the form

$$g(v) = h(|v|)|v|^\ell \mathcal{Y}_{\ell,m}(v/|v|)$$

for some function h on $[0, \infty)$ and some spherical harmonic $\mathcal{Y}_{\ell,m}$.

Second, for each $n \geq 0$, K preserves the space of polynomials of degree n . To see this, notice that any monomial in $\sqrt{1 - |v|^2}y$ that is of odd degree is annihilated when integrated against $(1 - |y|^2)^{\alpha-3/2}dy$, and any even monomial in $\sqrt{1 - |v|^2}y$ is a polynomial in v .

By combining these two observations, we see that K has a complete basis of eigenfunctions of the form

$$g_{n,\ell,m}(v) = h_{n,\ell}(|v|^2)|v|^\ell \mathcal{Y}_{\ell,m}(v/|v|),$$

where $h_{n,\ell}$ is a polynomial of degree n .

To determine these polynomials, we use the fact that K is self-adjoint, so that the eigenfunctions $g_{n,\ell,m}$ can be taken to be orthogonal. In particular, for any two

distinct positive integers n and p , the eigenfunctions $g_{n,\ell,m}$ and $g_{p,\ell,m}$ are orthogonal in $L^2(B, C_\alpha(1 - |v|^2)^\alpha)$. Hence for each ℓ , and for $n \neq p$,

$$\int_{|v| \leq 1} h_{n,\ell}(|v|^2)h_{p,\ell}(|v|^2)(1 - |v|^2)^\alpha |v|^{2\ell} dv = 0.$$

By taking $r = |v|^2$ as a new variable, we have

$$\int_0^1 h_{n,\ell}(r)h_{p,\ell}(r)(1 - r)^\alpha r^{\ell+1/2} dr = 0.$$

This is the orthogonality relation for a family of Jacobi polynomials in one standard form, and this identifies the polynomials $h_{n,\ell}$. A more common standard form, and one that is used in the sources to which we shall refer, is obtained by the change of variable $t = 2r - 1$, so that the variable t ranges over the interval $[-1, 1]$. Then for $\alpha, \beta > -1$, $P_n^{(\alpha,\beta)}(t)$ is the n th degree orthogonal polynomial for the weight $(1 - t)^\alpha(1 + t)^\beta$. With the variables t and $|v|^2$ related as above, i.e.,

$$t = 2|v|^2 - 1,$$

$$h_{n,\ell}(|v|^2) = P_n^{(\alpha,\beta)}(t)$$

for

$$\beta = \ell + \frac{1}{2}.$$

Now that we have all of the eigenfunctions determined, a further observation gives us a simple formula for the eigenvalues. Consider any eigenfunction g with eigenvalue κ , so that $Kg(v) = \kappa g(v)$. Let \hat{e} be any unit vector in R^3 . Then since g is a polynomial and hence continuous,

$$(7.8) \quad \begin{aligned} \lim_{t \rightarrow 1} Kg(t\hat{e}) &= \lim_{t \rightarrow 1} \int_B g\left(a\sqrt{1 - t^2}y + bt\hat{e}\right) C_{\alpha-3/2}(1 - |y|)^{\alpha-3/2} dy \\ &= g(b\hat{e}), \end{aligned}$$

since $K1 = 1$. By combining this with $Kg(v) = \kappa g(v)$, we have

$$g(b\hat{e}) = \kappa g(\hat{e}).$$

Now consider any eigenfunction $g_{n,\ell,m}$ of the form given above, and let $\kappa_{n,\ell}$ be the corresponding eigenvalue, which will not depend on m . Then by taking any \hat{e} so that $\mathcal{Y}_{\ell,m}(\hat{e}) \neq 0$, we have

$$(7.9) \quad \kappa_{n,\ell} = \frac{h_{n,\ell}(b^2)}{h_{n,\ell}(1)} b^\ell.$$

By changing variables as above to express this as a ratio of Jacobi polynomials, we finally have proved (7.5). \square

One might expect the largest eigenvalues of K to correspond to eigenfunctions that are polynomials of low degree. After all, in a system of orthogonal polynomials, those with high degree will have many changes of sign, and one might expect considerable

cancellation when applying an averaging operator, such as K , to them. Therefore, let us compute the $\kappa_{n,\ell}$ for low values of n and ℓ . We find from (7.5), by using the value $b = -1/(N-1)$ from (7.3), that

$$(7.10) \quad \kappa_{0,1} = \kappa_{1,0} = \frac{-1}{N-1},$$

so that $\kappa_{n,\ell}$ is negative for $n + \ell = 1$. For $n + \ell = 2$, we find from (7.5) that

$$(7.11) \quad \begin{aligned} \kappa_{1,1}(N) &= \frac{5N-3}{3(N-1)^3}, \\ \kappa_{2,0}(N) &= \frac{(N-3)(15N^2-15N+4)}{3(3N-4)(N-1)^4}, \\ \kappa_{0,2}(N) &= \frac{1}{(N-1)^2}. \end{aligned}$$

Evidently, for large N ,

$$\kappa_{0,2}(N) = \frac{1}{N^2} + \mathcal{O}\left(\frac{1}{N^3}\right),$$

while

$$\kappa_{1,1}(N) = \frac{5}{3N^2} + \mathcal{O}\left(\frac{1}{N^3}\right) \quad \text{and} \quad \kappa_{0,2}(N) = \frac{5}{3N^2} + \mathcal{O}\left(\frac{1}{N^3}\right).$$

Thus, one might expect that, at least for large values of N , 1 , $\kappa_{1,1}(N)$, $\kappa_{2,0}(N)$, and $\kappa_{0,2}(N)$ are the four largest eigenvalues of K and that $\kappa_{0,1} = \kappa_{1,0}$ is the most negative, with all other eigenvalues of K lying strictly between these. We shall show in the next section that this is indeed the case for all $N \geq 4$ and that 1 and $\kappa_{1,1}$ are the two largest eigenvalues of K for all $N \geq 3$.

When we use Lemma 5.1 to convert this to spectral information on P , we find that $\kappa_{0,1}$, $\kappa_{1,0}$, and $\kappa_{0,2}$ all correspond to the same eigenvalues of P , namely,

$$\frac{1}{N} \left(1 + \frac{1}{N-1}\right) = \frac{1}{N} \left(1 + (N-1) \frac{1}{(N-1)^2}\right) = \frac{1}{N-1}.$$

This is the eigenvalue of P that shall play the role of $\mu_N^{(m)}$ in our application of Lemma 4.1.

Let us conclude this section by recording a number of useful calculations that can be made by using (7.5).

For $N = 3$, we have

$$(7.12) \quad \kappa_{1,1}(3) = \frac{1}{2} > \kappa_{2,2}(3) = \frac{13}{40} > \kappa_{0,2}(3) = \frac{1}{4} > \kappa_{2,0}(3) = 0.$$

For $N = 4$, we have

$$(7.13) \quad \kappa_{1,1}(4) = \frac{17}{81} > \kappa_{0,2}(4) = \frac{1}{9} > \kappa_{2,0}(4) = \frac{23}{243}.$$

For $N = 5$, we have

$$(7.14) \quad \kappa_{1,1}(5) = \frac{11}{96} > \kappa_{2,0}(5) = \frac{19}{264} > \kappa_{0,2}(5) = \frac{1}{16}.$$

In each case, the second largest eigenvalue after 1, among the ones listed, is $\kappa_{1,1}$. In the next section we shall see that the list is not misleading: $\kappa_{1,1}$ is the gap eigenvalue. However, note that the third largest eigenvalue comes from different values of n and ℓ for each of $N = 3$, $N = 4$, and $N = 5$. As we shall see, things do settle down for $N \geq 5$; the third largest eigenvalue does turn out to be $\kappa_{2,0}$ in all such cases.

LEMMA 7.2. *For all $N \geq 5$, $\kappa_{1,1}(N) > \kappa_{2,0}(N) > \kappa_{0,2}(N)$.*

Proof. From (7.11),

$$\kappa_{2,0}(N) - \kappa_{0,2}(N) = \frac{2N(3N^2 - 15N + 8)}{3(3N - 4)(N - 1)^4}.$$

A simple calculation shows that the roots of the polynomial in the numerator are less than 5, so that $\kappa_{2,0}(N) > \kappa_{0,2}(N)$ for $N \geq 5$. A similar argument applied to $\kappa_{1,1}(N) - \kappa_{2,0}(N)$ yields the conclusion of the lemma. \square

Our goal in the next section is to show that, for all $N \geq 4$, there are no eigenvalues $\kappa_{n,\ell}$ with $n + \ell > 2$ that are larger than the ones listed above and that, for $N = 3$, the three largest eigenvalues are $1 = \kappa_{0,0} > 1/2 = \kappa_{1,1} > 13/40 = \kappa_{2,2}$. However, since there is no simple monotonicity in $n + \ell$, this shall require some detailed estimate on ratios of Jacobi polynomials.

We shall also need to know that in all cases $\kappa_{0,1} = \kappa_{1,0} = -1/(N - 1)$ is the most negative eigenvalue. This will tell us the four largest eigenvalues of P for $N \geq 4$ and the three largest for $n = 3$, and this shall turn out to be enough to prove the main result, Theorem 1.1.

Finally, the value of $\kappa_{2,2}(N)$ will play an important role in the proof of Theorem 1.2, and so we record the expression here:

$$(7.15) \quad \kappa_{2,2}(N) = \frac{21N^3 - 60N^2 + 27N - 4}{(3N - 4)(N - 1)^6}.$$

8. The determination of the spectrum of K . The main result in this section is the following theorem.

THEOREM 8.1. *For $N \geq 5$ and all n and ℓ with $n + \ell > 2$,*

$$(8.1) \quad -\frac{1}{N - 1} \leq \kappa_{n,\ell}(N) < \kappa_{0,2}(N).$$

For $N = 4$ and all n and ℓ with $n + \ell > 2$,

$$(8.2) \quad -\frac{1}{N - 1} \leq \kappa_{n,\ell}(4) < \kappa_{2,0}(4).$$

For $N = 3$ and all n and ℓ with $n + \ell > 0$, except for $n = 1, \ell = 1$,

$$(8.3) \quad -\frac{1}{N - 1} \leq \kappa_{n,\ell}(3) \leq \kappa_{2,2}(3) = \frac{13}{40}.$$

We present the proof at the end of this section after a number of preparatory lemmas. These lemmas rest on two deep results about Jacobi polynomials. One is a formula due to Koornwinder [8] (see also [1, p. 31]) that was already applied in [3].

For all $-1 \leq x \leq 1$, all n , and all $\alpha > \beta$,

$$(8.4) \quad \frac{J_n^{(\alpha,\beta)}(x)}{J_n^{(\alpha,\beta)}(1)} = \int_0^\pi \int_0^1 \left[\frac{1 + x - (1 - x)r^2}{2} + i\sqrt{1 - x^2}r \cos(\theta) \right]^n dm_{\alpha,\beta}(r, \theta),$$

where

$$m_{\alpha,\beta}(r, \theta) = c_{\alpha,\beta}(1 - r^2)^{\alpha-\beta-1}r^{2\beta+1}(\sin \theta)^{2\beta} \, drd\theta,$$

and $c_{\alpha,\beta}$ is a normalizing constant that makes $dm_{\alpha,\beta}$ a probability measure.

Koornwinder’s bound is very useful for obtaining uniform control in n for given ℓ and N . But since in Lemma 7.1

$$(8.5) \quad \alpha = \frac{3N - 8}{2} \quad \text{and} \quad \beta = \ell + \frac{1}{2},$$

we can apply (8.4) only when

$$(8.6) \quad \ell < \ell^* = \frac{3N - 9}{2}.$$

As in [3], one may use this formula to show the following.

LEMMA 8.2. *For all ℓ with $2 \leq \ell < \ell^*$, and all $n > 0$ and $N \geq 3$,*

$$|\kappa_{n,\ell}(N)| < \frac{1}{(N - 1)^2} = \kappa_{0,2}(N).$$

Note that, while this lemma does not address the case $n = 0$, this is not a problem: We have the explicit formula

$$(8.7) \quad \kappa_{0,\ell} = \left(\frac{-1}{N - 1} \right)^\ell.$$

To handle large values of ℓ , we need another deep result, which is a uniform pointwise bound on the *orthonormal* Jacobi polynomials that was obtained by Nevai, Erdelyi, and Magnus [12]. Let $p_n^{\alpha,\beta}$ be the orthonormal Jacobi polynomial of degree n with a positive leading coefficient for the weight $w(x) = (1 - x)^\alpha(1 + x)^\beta$. It was shown in [12] that, for all $\alpha \geq -1/2$ and $\beta \geq -1/2$ and all nonnegative integers n ,

$$(8.8) \quad \max_{x \in [-1,1]} \sqrt{1 - x^2}w(x)p_n^{\alpha,\beta}(x)^2 \leq \frac{2e(2 + \sqrt{\alpha^2 + \beta^2})}{\pi}.$$

Of course, we could use the orthonormal Jacobi polynomials in the ratio formula (7.5), since any normalization factor would cancel out in the ratio. However, the exact formula (7.4) for the denominator in (7.5) is simplest in the other normalization. Hence we need the relation between $p_n^{\alpha,\beta}$ and $P_n^{\alpha,\beta}$, which is given by $p_n^{\alpha,\beta} = l_n P_n^{\alpha,\beta}$, where

$$l_n = \left(\frac{2n + \alpha + \beta + 1}{2^{\alpha+\beta+1}} \frac{\Gamma(n + 1)\Gamma(n + \alpha + \beta + 1)}{\Gamma(n + \alpha + 1)\Gamma(n + \beta + 1)} \right)^{1/2}.$$

Therefore

$$(8.9) \quad \frac{P_n^{\alpha,\beta}(x)^2}{P_n^{\alpha,\beta}(1)^2} \leq \frac{1}{l_n^2} \frac{2e\Gamma(n + 1)^2\Gamma(\alpha + 1)^2(2 + \sqrt{\alpha^2 + \beta^2})}{\sqrt{1 - x^2}w(x)\pi\Gamma(n + \alpha + 1)^2}.$$

At this point it is perhaps worth noting that, since the spectrum of K lies in $[-1, 1]$, any upper bound on its eigenvalues by a number larger than one is vacuous. This implies that for certain regions the identity (7.9) will provide a stronger bound than (8.9). We shall return to this point at the end of the paper.

Substituting $x = -1 + \frac{2}{(N-1)^2}$, $\beta = \ell + \frac{1}{2}$, and $\alpha = \frac{3}{2}N - 4$ in (8.9) and then multiplying by $\frac{1}{(N-1)^{2\ell}}$ yields

$$(8.10) \quad \kappa_{n,\ell}(N) \leq \tilde{\kappa}_{n,\ell}(N),$$

where

$$(8.11) \quad \tilde{\kappa}_{n,\ell}^2(N) = \frac{2e}{\pi} g_1(n, \ell, N) g_2(N) g_3(n, N) g_4(n, \ell, N),$$

with

$$(8.12) \quad \begin{aligned} g_1(n, \ell, N) &= \left(\frac{4 + \sqrt{9N^2 - 48N + 65 + 4\ell^2 + 4\ell}}{3N + 4n + 2\ell - 5} \right), \\ g_2(N) &= \left(\frac{(N-1)^2}{N(N-2)} \right)^{(3N-7)/2}, \\ g_3(n, N) &= \frac{\Gamma(n+1)\Gamma(\frac{3}{2}N-3)}{\Gamma(n+\frac{3}{2}N-3)}, \\ g_4(n+\ell, N) &= \frac{(N-1)^2\Gamma(n+\ell+\frac{3}{2})\Gamma(\frac{3}{2}N-3)}{\Gamma(n+\ell+\frac{3}{2}N-\frac{5}{2})}. \end{aligned}$$

Our goal now is to extract a reasonably tight upper bound for $\tilde{\kappa}_{n,\ell}(N)$ with as much monotonicity in n , ℓ , and N as possible. The next lemmas address this goal.

LEMMA 8.3. For $\ell \geq 0$, $N \geq 3$, and $n \geq 0$,

$$(8.13) \quad g_1(n, \ell, N) \leq \left(\frac{4}{3N + 4n + 2\ell - 5} + 1 \right),$$

where the right-hand side is clearly decreasing in n , ℓ , and N .

Proof. Note that, for $n \geq 0$, $\ell \geq 0$, and $N \geq 3$,

$$(8.14) \quad \frac{\sqrt{9N^2 - 48N + 65 + 4\ell^2 + 4\ell}}{3N + 4n + 2\ell - 5} \leq 1$$

since then

$$(8.15) \quad \begin{aligned} &(3N + 4n + 2\ell - 5)^2 - (9N^2 - 48N + 65 + 4\ell^2 + 4\ell) \\ &= (24N - 40)n + (12N - 24)\ell + 16n^2 + 16n\ell + 18N - 40 > 0. \quad \square \end{aligned}$$

LEMMA 8.4. For $N \geq 4$, $g_2(N)$ is a decreasing function of N .

Proof. Let $h(x) = (1 - 1/x^2)^{2-3x/2}$, so that $g_2(N) = h(N-1)$. By computing the derivative of $\ln(h(x))$, one finds that it is negative for $x \geq 3$. \square

LEMMA 8.5. For $n \geq 0$ and $N \geq 3$, $g_3(n, N)$ is a decreasing function of n and N .

Proof. For n a nonzero integer

$$(8.16) \quad \frac{\Gamma(n+1)\Gamma(\frac{3}{2}N-3)}{\Gamma(n+\frac{3}{2}N-3)} = \frac{n}{n+\frac{3}{2}N-4} \frac{n-1}{n+\frac{3}{2}N-5} \cdots \frac{1}{\frac{3}{2}N-3}.$$

Since each factor is less than 1 for $N \geq 3$ and is a decreasing function of N , the assertion follows. \square

LEMMA 8.6. For $N \geq 3$, $g_4(n + \ell, N)$ is a decreasing function of $n + \ell$, with

$$(8.17) \quad \lim_{n+\ell \rightarrow \infty} g_4(n + \ell, N) = 0.$$

Moreover, for $n + \ell \geq \ell^* = 3(N - 3)/2$,

$$(8.18) \quad g_4(n, \ell, N) \leq \frac{(N - 1)^2 \Gamma(\frac{3}{2}N - 3)^2}{\Gamma(3N - 7)} \leq f(N),$$

where

$$(8.19) \quad f(N) = \frac{(N - 1)^2 \sqrt{\pi} (\frac{3}{2}N - 4)}{2^{3N-8}}.$$

Finally, for $N \geq 5$, $(N - 1)^4 f(N)$ is a decreasing function of N .

Proof. Since

$$(8.20) \quad \frac{\Gamma(n + \ell + \frac{5}{2})}{\Gamma(n + \ell + \frac{3}{2}N - \frac{3}{2})} \frac{\Gamma(n + \ell + \frac{3}{2}N - \frac{5}{2})}{\Gamma(n + \ell + \frac{3}{2})} = \frac{(n + \ell + \frac{3}{2})}{n + \ell + \frac{3}{2}N - \frac{5}{2}} < 1$$

for $N \geq 3$, it follows that, for fixed nonnegative integers N ,

$$(8.21) \quad \frac{\Gamma(n + \ell + \frac{3}{2})}{\Gamma(n + \ell + \frac{3}{2}N - \frac{5}{2})}$$

is a decreasing function of $n + \ell$. Hence, for $n + \ell \geq 3(N - 3)/2$,

$$\frac{\Gamma(n + \ell + \frac{3}{2})}{\Gamma(n + \ell + \frac{3}{2}N - \frac{5}{2})} \leq \frac{\Gamma(\frac{3}{2}N - 3)}{\Gamma(3N - 7)}.$$

This together with the definition of g_4 proves the first inequality in (8.18). Use of the duplication formula for the Γ function yields

$$\frac{\Gamma(\frac{3}{2}N - 3)^2}{\Gamma(3N - 7)} = \frac{\sqrt{\pi} \Gamma(\frac{3}{2}N - 3)}{2^{3N-8} \Gamma(\frac{3}{2}N - \frac{7}{2})} = \frac{\sqrt{\pi} (\frac{3}{2}N - 4) \Gamma(\frac{3}{2}N - 4)}{2^{3N-8} \Gamma(\frac{3}{2}N - \frac{7}{2})} < \frac{\sqrt{\pi} (\frac{3}{2}N - 4)}{2^{3N-8}}.$$

This implies the second inequality in (8.18). A check of the logarithmic derivative of $(N - 1)^4 f(N)$ shows that it is negative for $N \leq 5$. \square

Now, by combining the results in the last four lemmas, we have, for $N \geq 3$ and $n + \ell \geq \ell^* = 3(N - 3)/2$,

$$(8.22) \quad \tilde{\kappa}_{n,\ell}^2(N) \leq \hat{\kappa}_{n,l}^2(N) \leq \kappa^2(N),$$

where

$$(8.23) \quad \hat{\kappa}_{n,l}^2(N) = \frac{2e}{\pi} \left(\frac{4}{3N + 4n + 2\ell - 5} + 1 \right) g_2(N) g_3(n, N) g_4(n, l, N)$$

and

$$(8.24) \quad \kappa^2(N) = \frac{2e}{\pi} \left(\frac{4}{6N - 14} + 1 \right) g_2(N) f(N),$$

where g_2 , g_3 , and f are given by (8.12) and (8.19).

We are now ready to prove the main theorem of this section.

Proof of Theorem 8.1. First, we take care of large values of N . By Lemmas 8.4 and 8.6, $(N - 1)^4 \kappa(N)$ is a decreasing function of N for $N \geq 5$. Direct computation shows that at $N = 12$ this quantity is less than one. Hence for $N \geq 12$, $\kappa(N) \leq (N - 1)^{-4} = \kappa_{0,2}^2$. For $\ell \geq \ell^*$, so that (8.22) is satisfied, this proves (8.1) for $N \geq 12$. On the other hand, if $2 \leq \ell < \ell^*$, we have this from Lemma 8.2 or (8.7). Thus, in any case, (8.1) is valid for $N \geq 12$.

For $4 \leq N \leq 11$, we again use Lemma 8.2 or (8.7) for $2 \leq \ell < \ell^*$ and computation of $\hat{\kappa}_{n,\ell}$. By (8.17), for each such N there is a finite value $k(N)$ so that we need only consider values of $n + \ell < k(N)$. By checking these cases, we obtain (8.1) and (8.2).

We finally turn to $N = 3$, which requires the greatest amount of computation. First for $n = 0$, we have from (8.7)

$$\kappa_{0,\ell}(3) = \left(\frac{-1}{2}\right)^\ell$$

so $\kappa_{0,1}(3) = -1/2$ and $|\kappa_{0,\ell}(3)| < 1/3$ for $\ell \geq 2$.

The exact forms of the eigenvalues are simple enough to be useful for $n = 1$ and 2 as well. We have

$$\kappa_{1,\ell}(N) = \frac{(-1)^{\ell+1} [2\ell N + 3(N - 1)]}{3(N - 1)^{\ell+2}}$$

and

$$\kappa_{2,\ell}(N) = \frac{(-1)^\ell ((4\ell^2 + 16\ell + 15)N^3 - (8\ell^2 + 44\ell + 60)N^2 + (49 + 16\ell)N - 12)}{3(3N - 4)(N - 1)^{\ell+4}}.$$

By specializing to $N = 3$,

$$\kappa_{1,\ell}(3) = (-1)^{\ell+1} \frac{\ell + 1}{2^{\ell+1}}$$

so that $|\kappa_{1,\ell}(3)| \leq 3/8$ for $\ell \geq 2$. Likewise, for $N = 3$,

$$\kappa_{2,\ell}(3) = \frac{\ell}{20} \frac{(7 + 3\ell)}{2^\ell} (-1)^\ell,$$

which implies that

$$(8.25) \quad |\kappa_{2,\ell}(3)| \leq |\kappa_{22}(3)| = \frac{13}{40}.$$

For higher values of n , we estimate $\kappa_{n,\ell}^2$ by means of $\hat{\kappa}_{n,\ell}^2$. Since $\ell^* = 0$ for $N = 3$, we may use Lemma 8.6 for all ℓ , and then by (8.17), for each fixed n , there is a maximal value $\ell(n)$ that needs to be considered and a maximum value of n that needs to be considered. Table 1 gives the values of n, ℓ , and $\hat{\kappa}_{n,\ell}^2(3)$ when $\tilde{\kappa}_{n,\ell}^2(3) < 1/4$. The monotonicity of $\kappa_{n,\ell}^2(3)$ in n and l shows that $\hat{\kappa}_{n,\ell}^2(3) \leq \hat{\kappa}_{n_0,\ell_0}^2(3)$ for $n \geq n_0$ and $\ell \geq \ell_0$ where (n_0, ℓ_0) is chosen from the table. The remaining values can be computed

TABLE 1

n	ℓ	$\hat{\kappa}_{n,\ell}^2$	n	ℓ	$\hat{\kappa}_{n,\ell}^2$	n	ℓ	$\hat{\kappa}_{n,\ell}^2$	n	ℓ	$\hat{\kappa}_{n,\ell}^2$
3	1253	0.10562	20	210	0.10547	37	90	0.10543	54	34	0.10514
4	989	0.10562	21	198	0.10559	38	86	0.10531	55	31	0.10538
5	817	0.10556	22	188	0.10546	39	82	0.105277	56	29	0.10506
6	694	0.10561	23	178	0.10552	40	78	0.10528	57	26	0.10538
7	604	0.10555	24	169	0.10549	41	74	0.10537	58	24	0.10511
8	533	0.10560	25	161	0.10537	42	70	0.10551	59	21	0.10551
9	477	0.10558	26	153	0.10540	43	67	0.10523	60	19	0.10529
10	431	0.10558	27	145	0.10558	44	63	0.10552	61	17	0.10509
11	393	0.10555	28	138	0.10561	45	60	0.10534	62	14	0.10560
12	360	0.10561	29	132	0.10542	46	57	0.10521	63	12	0.10545
13	333	0.10548	30	126	0.10534	47	54	0.10512	64	10	0.10534
14	308	0.10558	31	120	0.10540	48	50	0.10562	65	8	0.10523
15	287	0.10554	32	114	0.10554	49	48	0.10508	66	6	0.10514
16	268	0.10556	33	109	0.10543	50	45	0.10512	67	4	0.10509
17	251	0.10558	34	104	0.10540	51	42	0.10521	68	2	0.10506
18	236	0.10554	35	99	0.10546	52	39	0.10535	69	0	0.10503
19	222	0.10560	36	94	0.10562	53	36	0.10552			

from the exact formula for $\kappa_{n,\ell}(3)$ from (7.1), and the results are all consistent with (8.3). \square

9. The determination of the spectrum of P . For given values of N , n , and ℓ , let $\mu_{n,\ell}(N)$ be the eigenvalue of P corresponding to the eigenvalue $\kappa_{n,\ell}(N)$ of K through Theorem 8.1, where we use (5.7) if $\kappa_{n,\ell}(N) > 0$ and use (5.8) if $\kappa_{n,\ell}(N) < 0$. (This is the relevant choice, as we are concerned with the largest eigenvalues of P .)

By consulting the calculations in (7.10) for $n + \ell = 1$, and in (7.12), (7.13), and (7.14) for $n + \ell = 2$, and finally the bounds in Theorem 8.1 for $n + \ell > 2$, we see that for all $N \geq 3$ the largest eigenvalues of K is $\kappa_{1,1}$, and the least (most negative) is $\kappa_{0,1} = \kappa_{1,0}$. Thus, by turning to Lemma 5.1 and using the positive eigenvalue in (5.7) and the negative one in (5.8), we see that the positive one yields the greater value for each N . Thus, the gap eigenvalue of P , μ_N , is given by

$$(9.1) \quad \mu_N = \mu_{1,1}(N) = \frac{3N - 1}{3(N - 1)^2}.$$

Use of this result in (3.8) would yield a strictly positive lower bound on Δ_N , uniform in N , but, as we have said above, it would not yield the exact lower bound. To obtain this, we now carry out the strategy outlined in the introduction.

First, we combine Lemma 5.1 and Theorem 8.1 to produce the information necessary for the application of Lemma 4.1. We must now make a choice of the thresholds μ_N^* that appear in Lemma 4.1. The choice we shall make is based on trial function computations with Q that suggest that the gap eigenfunctions are the ones specified in Theorem 1.1.

Notice that in Theorem 1.1 the formula given for Δ_N is of the form $C \frac{N}{N-1}$ for some constant C . This value can be guessed by computing the eigenvalues of Q on the invariant subspace of polynomials of degree 4 or less in the v_j . If we are to prove this guess correct by using (4.2) of Lemma 4.1, we require a value of μ_N^* such that

$$(9.2) \quad \frac{N}{N-1}(1 - \mu_N^*) \frac{N-1}{N-2} \geq \frac{N}{N-1},$$

at least for $N \geq 4$. (The guess is valid only for $N-1 \geq 3$. For $N-1 = 2$, there is a different value of Δ_2 which has been determined already in section 2.)

The largest value of μ_N^* that will satisfy (9.2) is

$$(9.3) \quad \mu_N^* = \frac{1}{N-1} \quad \text{for} \quad N \geq 4.$$

This turns out to be an eigenvalue of P : Indeed, we have found in (7.11) that $\kappa_{0,2} = 1/(N-1)^2$. Furthermore, we have found in (7.10) that $\kappa_{0,1} = \kappa_{1,0} = -1/(N-1)$. By using the first of these results in (5.7) of Lemma 5.1 and the second in (5.8), we find that

$$\mu_{0,2} = \mu_{1,0} = \mu_{0,1} = \frac{1}{N-1}.$$

For $N = 3$ we need to make a different choice, as the spectrum of Q is quite different for $N = 2$ and for $N \geq 3$. The choice that will work is $\mu_3^* = \mu_{2,2}(3)$. Since $\kappa_{2,2}(3) = 13/40$, we have from Lemma 5.1 that $\mu_{2,2}(3) = \frac{1}{3}(1 + 2(13/40)) = \frac{11}{20}$. Thus,

$$(9.4) \quad \mu_3^* = \frac{11}{20}.$$

Now, to apply Lemma 4.1, we need the eigenspaces of P for the eigenvalues μ satisfying $1 > \mu > \mu_N^*$. By Theorem 8.1 and (7.13), for $N = 3$ and $N = 4$, there is just one such eigenvalue, namely, $\mu_{1,1}(4)$, the gap eigenvalue, and for $N \geq 5$, there are two: $\mu_{1,1}(N)$ and $\mu_{2,0}(N)$.

Let $E_{n,\ell}$ be the eigenspace of P corresponding to the eigenvalue $\mu_{n,\ell}(N)$. For all values of n and ℓ with $n + \ell \leq 2$, we have determined the corresponding eigenfunctions of K and thus, through Lemma 5.1, the corresponding eigenfunctions of P . Thus, we have the following explicit descriptions of the $E_{n,\ell}$ for $n + \ell \leq 2$.

First, for $n + \ell = 1$, the eigenvalues of K are negative, and so by Lemma 5.1, the eigenfunctions are antisymmetric. If we are concerned only with the spectrum of Q on the subspace of symmetric functions (which is all that is of significance for Kac's application to the Boltzmann equation), we can ignore these eigenspaces. However, they turn out to be very simple. The $n = 0, \ell = 1$ eigenfunctions of K are degree one spherical harmonics, and the $n = 1, \ell = 0$ eigenfunctions of K are degree one Jacobi polynomials in $|v|^2$. Hence

$$(9.5) \quad E_{0,1} \quad \text{is spanned by the functions} \quad v_i^\alpha - v_j^\alpha, \alpha = 1, 2, 3, \quad \text{and} \quad i < j,$$

while

$$(9.6) \quad E_{1,0} \quad \text{is spanned by the functions} \quad |v_i|^2 - |v_j|^2, i < j.$$

Next, for $n + \ell = 2$, the eigenvalues of K are positive, and so by Lemma 5.1, the eigenfunctions are symmetric. The $n = 0, \ell = 2$ eigenfunctions of K are degree two spherical harmonics and so have the form

$$f_{0,2}(v) = \sum_{\alpha,\beta=1}^3 A_{\alpha,\beta} v^\alpha v^\beta$$

for some traceless symmetric 3×3 matrix A . Hence, by Lemma 5.1,

$$(9.7) \quad E_{0,2} \text{ is spanned by the functions } \sum_{j=1}^N f_{0,2}(v_j),$$

with $f_{0,2}$ given as above.

For $n = 1, \ell = 1$, the eigenfunctions of K are the product of a degree one spherical harmonic and a degree one Jacobi polynomial in $|v|^2$. When we sum over the particles, the constant term in the Jacobi polynomial drops out due to the momentum constraint, and we see that

$$(9.8) \quad E_{1,1} \text{ is spanned by the functions } \sum_{j=1}^N f_{1,1}(v_j),$$

where

$$f_{1,1}(v) = |v|^2 v^\alpha, \quad \alpha = 1, 2, 3.$$

Finally, for $n = 2, \ell = 0$, the eigenfunction of K is a degree two Jacobi polynomial in $|v|^2$. After summing on the particles, the linear term can be absorbed into the constant by the energy constraint, and so we see that

$$(9.9) \quad E_{2,0} \text{ is spanned by the function } \sum_{j=1}^N f_{2,0}(v_j),$$

where

$$f_{2,0}(v) = |v|^4 - \int_B |v|^4 d\nu_N.$$

We close this section with a lemma that we shall need to prove Theorem 1.2. There we shall need to know the next largest eigenvalue of P below $\max_{n+\ell \leq 2} \mu_{n,\ell}(N)$. One might guess that this occurs for some values of n and ℓ with $n + \ell = 3$, but this is not the case: By (8.3) of Theorem 8.1, and (7.12), we see that for $N = 3$ the most negative eigenvalue of K is $-1/2$, and by Lemma 5.1, this corresponds to the eigenvalue $1/2$ of P . On the other hand, the largest eigenvalue of K apart from $\kappa_{1,1}(3)$ is $\kappa_{2,2}(3) = 13/40$. This corresponds to the eigenvalue $11/20$ of P . Since $11/20 > 1/2$, we do indeed have

$$\sup_{n+\ell > 2} \mu_{n,\ell}(3) = \mu_{2,2}(3) = \frac{11}{20}.$$

It seems likely, on the basis of computations that we have made, that in fact

$$(9.10) \quad \sup_{n+\ell > 2} \mu_{n,\ell}(N) = \mu_{2,2}(N)$$

for all $N \geq 3$. However, for the proof of Theorem 1.2, all that we require is the following.

LEMMA 9.1. *For $N = 3, 4, 5, 6$, and 7 , (9.10) is true.*

Proof. Note that the case $N = 3$ has already been proved in the remarks above. To deal with the other cases, we proceed essentially as in the proof of Theorem 8.1, using (8.23) to reduce the number of cases to be checked to a finite number and then checking these. We will therefore be brief in our remarks on the remaining cases.

Perhaps the most important point to recall is that (8.23) is valid for $n + \ell \geq 3(N - 3)/2$. Since the right-hand side evaluates to zero for $N = 3$, we could use it without restriction. For $N = 7$, though, $3(N - 3)/2$ evaluates to 6, and so we may use (8.23) only for $n + \ell \geq 6$. So these cases must be checked by direct computation of the eigenvalues by using 7.1 and then converting these to eigenvalues of P by using Lemma 5.1.

Then, by using (8.23) for $n + \ell > 6$, one finds that

$$\kappa_{n,\ell}^2(7) < \kappa_{2,2}^2(7)$$

unless $0 \leq n \leq 6$ and $0 \leq \ell \leq 27$. By computing the rest of the eigenvalues of P in this 6 by 27 rectangle, we find that the stated result is true for $N = 7$.

A similar analysis takes care of $N = 4$, $N = 5$, and $N = 6$. \square

We shall not need to know the corresponding eigenfunctions in our application of Lemma 9.1, since we will be concerned only with the eigenspaces of eigenvalues lying strictly above $\mu_{2,2}(N)$, and those have been determined already in this section.

10. The spectrum of Q on invariant subspaces containing eigenspaces of P . For each n and ℓ , let $V_{n,\ell}$ be the smallest invariant subspace of Q containing $E_{n,\ell}$. As we shall see, for $n + \ell \leq 2$, $V_{n,\ell} = E_{n,\ell}$ except for $n = 2, \ell = 0$, in which case $V_{2,0}$ is two-dimensional, while $E_{1,0}$ is one-dimensional. This is established in the next lemma, which also specifies the spectrum of Q on these invariant subspaces. The eigenvalues of course depend on the particular choice of b in the definition of Q , but in a very simple way: The dependence on b is only through the quantities $(1 - B_1)$ and $(1 - B_2)$, where B_j is the j th moment of b , as defined in (1.5).

LEMMA 10.1. *Every nonzero function in $E_{0,1}$ and in $E_{1,0}$ is an eigenfunction of Q with eigenvalue*

$$(10.1) \quad \lambda_{0,1}^Q = \lambda_{1,0}^Q = 1 - (1 - B_1) \frac{1}{N - 1},$$

so that $V_{0,1} = E_{0,1}$ and $V_{1,0} = E_{1,0}$.

Every nonzero function in $E_{1,1}$ is an eigenfunction of Q with the eigenvalue

$$(10.2) \quad \lambda_{1,1}^Q = 1 - (1 - B_2) \frac{1}{(N - 1)},$$

so that $V_{1,1} = E_{1,1}$.

Furthermore, every nonzero function in $E_{0,2}$ is an eigenfunction of Q with the eigenvalue

$$(10.3) \quad \lambda_{0,2}^Q = 1 - (1 - B_2) \frac{3}{2(N - 1)},$$

so that $V_{0,2} = E_{0,2}$.

Finally, while $V_{2,0}$ is larger than $E_{2,0}$, there are only two eigenvalues of Q in $V_{2,0}$. These are

$$(10.4) \quad 1 - (1 - B_2) \left(\frac{1}{N(N - 1)} \left[(2N - 1) \pm \sqrt{N^2 - 3N + 1} \right] \right).$$

For all $N \geq 3$, the largest of these eigenvalues is $\lambda_{1,1}^Q$.

Before beginning the proof, we note that, if $(1/2)b(x)dx$ is a Dirac mass at $x = 1$, the collisions are all trivial (zero scattering angle), and thus $Q = I$ in this case. But also in this case $(1 - B_1) = (1 - B_2) = 0$, so all of the eigenvalues $\lambda_{n,\ell}^Q$ listed above are 1—as they must be for $Q = I$.

Proof. We begin with the last case, $n = 2, \ell = 0$, which is the most involved. Consider the function

$$(10.5) \quad \phi = \sum_{i=1}^N |v_i|^4,$$

and note that $\phi - \int_{X_N} \phi d\sigma_N$ spans $E_{2,0}$, as we have noted above.

One simple way to calculate $Q\phi$ is to take advantage of the permutation symmetry of Q : Define the symmetrization operator \mathcal{S} by

$$\mathcal{S}f(v_1, \dots, v_N) = \frac{1}{N!} \sum_{\pi} f(v_{\pi(1)}, \dots, v_{\pi(N)}),$$

where the sum runs over all permutations π of $\{1, \dots, N\}$. Then it is easy to see that

$$\mathcal{S}(|v_1|^4) = \frac{1}{N} \phi.$$

Thus, since $\mathcal{S}Q = Q\mathcal{S}$,

$$Q\phi = N\mathcal{S}Q(|v_1|^4).$$

One now directly calculates $Q(|v_1|^4)$ and then symmetrizes. In carrying out the calculation, we make use of the following.

LEMMA 10.2. *Let c and d be any two vectors in \mathbb{R}^3 , and let e be any unit vector in \mathbb{R}^3 . Then with B_1 and B_2 defined as in (1.5), we have the following identities:*

$$\int_{S^2} (c \cdot \sigma)b(e \cdot \sigma)d\sigma = (c \cdot e)B_1$$

and

$$\int_{S^2} (c \cdot \sigma)(d \cdot \sigma)b(e \cdot \sigma)d\sigma = (c \cdot d)\frac{1 - B_2}{2} + (e \cdot c)(e \cdot d)\frac{3B_2 - 1}{2}.$$

Proof. We choose coordinates in which c and e span the x, z plane with

$$e = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \text{and} \quad c = \begin{bmatrix} c^1 \\ 0 \\ c^3 \end{bmatrix}.$$

Then with

$$\sigma = \begin{bmatrix} \sin \theta \cos \psi \\ \sin \theta \sin \psi \\ \sin \theta \end{bmatrix},$$

the computations are easily accomplished. \square

Now to compute $Q\phi$, go back to the definition of Q given in (1.4), and note first of all that with $\eta(\vec{v}) = |v_1|^4$, unless $i = 1$,

$$\eta(R_{i,j,\sigma}(\vec{v})) = \eta(\vec{v}).$$

Hence

$$Q\eta(\vec{v}) = \left(1 - \frac{2}{N}\right)\eta(\vec{v}) + \frac{2}{N(N-1)} \sum_{j=2}^N \int_{S^2} \eta(R_{1,j,\sigma}(\vec{v}))b\left(\sigma \cdot \frac{v_i - v_j}{|v_i - v_j|}\right) d\sigma.$$

Then from (1.1),

$$\begin{aligned} \eta(R_{1,j,\sigma}(\vec{v})) &= \left| \frac{v_1 + v_j}{2} + \frac{|v_1 - v_j|}{2} \sigma \right|^4 \\ &= \frac{1}{8} \left| |v_1|^2 + |v_j|^2 + |v_1 - v_j|(v_1 + v_j) \cdot \sigma \right|^2 \\ &= \frac{1}{8} \left((|v_1|^2 + |v_j|^2)^2 + 2(|v_1|^2 + |v_j|^2)|v_1 - v_j|(v_1 + v_j) \cdot \sigma \right. \\ &\quad \left. + |v_1 - v_j|^2((v_1 + v_j) \cdot \sigma)^2 \right). \end{aligned} \tag{10.6}$$

Integrating over S^2 using Lemma 10.2 yields

$$\begin{aligned} \int_{S^2} \eta(R_{1,j,\sigma}(\vec{v}))b\left(\sigma \cdot \frac{v_i - v_j}{|v_i - v_j|}\right) d\sigma &= \frac{1}{8}(|v_1|^2 + |v_j|^2)^2 \\ &\quad + \frac{1}{4}(|v_1|^2 + |v_j|^2)B_1(|v_1|^2 - |v_j|^2) \\ &\quad + |v_1 - v_j|^2|v_1 + v_j|^2 \frac{1 - B_2}{16} \\ &\quad + ((v_1 - v_j) \cdot (v_1 + v_j))^2 \frac{3B_2 - 1}{16}. \end{aligned} \tag{10.7}$$

The right-hand side simplifies to

$$\begin{aligned} &\frac{1}{8}(|v_1|^4 + |v_j|^4 + 2|v_1|^2|v_j|^2) + \frac{B_1}{4}(|v_1|^4 - |v_j|^4) \\ &\quad + (|v_1|^4 + |v_j|^4 + 2|v_1|^2|v_j|^2 - 4(v_1 \cdot v_j)^2) \frac{1 - B_2}{16} \\ &\quad + (|v_1|^4 + |v_j|^4 - 2|v_1|^2|v_j|^2) \frac{3B_2 - 1}{16}. \end{aligned} \tag{10.8}$$

It is now a simple matter to carry out the sum on $j \geq 2$. By using the identities

$$\sum_{j=2}^N |v_j|^4 = \phi(\vec{v}) - |v_1|^4 \quad \text{and} \quad \sum_{j=2}^N |v_j|^2 = \phi(\vec{v}) - |v_1|^2$$

and the symmetrizing, one finds that

$$Q\phi = \phi - \frac{1 - B_2}{N} \left[\frac{N + 1}{N - 1} \phi + \frac{1}{(N - 1)} \psi - \frac{2}{(N - 1)} \right], \tag{10.9}$$

where

$$\psi = \sum_{i \neq j} (v_i \cdot v_j)^2. \tag{10.10}$$

For $N \geq 4$, the two functions ϕ and ψ are linearly independent, although for $N = 3$ they are not. In fact for $N = 3$, one has the identity

$$(10.11) \quad \psi = 2\phi - \frac{1}{2}.$$

Evidently, for $N \geq 4$, $\phi - \int_{X_N} \phi d\sigma_N$ is not an eigenfunction of Q , so that $E_{0,2}$ is not an eigenspace of Q . We are required to compute $Q\psi$.

We again take advantage of the permutation symmetry and note that

$$\mathcal{S}((v_1 \cdot v_2)^2) = \frac{2}{N(N-1)}\psi \quad \text{and} \quad Q\psi = \frac{N(N-1)}{2}\mathcal{S}(Q(v_1 \cdot v_2)^2).$$

We carry out the calculation in the same way that we calculated $Q\phi$ and find that

$$(10.12) \quad Q\psi = \psi - \frac{1-B_2}{N} \left[3\psi + \frac{(N-3)}{(N-1)}\phi - \frac{1}{N} \right].$$

We see that the subspace spanned by $\phi - \int_{X_N} \phi d\sigma_N$ and $\psi - \int_{X_N} \phi d\sigma_N$ is invariant under Q . By using (10.12) and (10.9) we easily find that the two eigenvalues of $N(I - Q)$ on the two-dimensional space $V_{2,0}$ are the eigenvalues of

$$\frac{1-B_2}{N-1} \begin{bmatrix} N+1 & 1 \\ N-3 & 3N-3 \end{bmatrix},$$

which are

$$\frac{1-B_2}{N-1} \left[(2N-1) \pm \sqrt{N^2 - 3N + 1} \right].$$

The minus sign clearly gives the lesser of these and gives the gap for $N(I - Q)$ on $V_{2,0}$. From here, one easily deduces (10.4).

A further, much simpler calculation shows that the three functions

$$(10.13) \quad \psi_{1,1}^\alpha = \sum_{k=1}^N |v_k|^2 v_k^\alpha,$$

where α indexes the components, are also eigenfunctions of Q ; more precisely,

$$(10.14) \quad Q\psi_{1,1}^\alpha = \left(1 - \frac{1-B_2}{N-1} \right) \psi_{1,1}^\alpha.$$

Thus the unique eigenvalue of A on $V_{1,1}$ is

$$\lambda_{1,1}^Q = 1 - \frac{1-B_2}{N-1}.$$

For $E_{0,2}$, a simple computation shows that the functions

$$(10.15) \quad \psi_{0,2}^{\alpha,\beta} = \sum_k v_k^\alpha v_k^\beta,$$

where $\alpha \neq \beta$ are indices for the components, are also eigenfunctions for Q ; in fact,

$$(10.16) \quad Q\psi_{0,2}^{\alpha,\beta} = \left(1 - \frac{3(1-B_2)}{2(N-1)} \right) \psi_{0,2}^{\alpha,\beta}.$$

Thus, $V_{0,2} = E_{0,2}$, and the unique eigenvalue of Q on this subspace is

$$\lambda_{0,2}^Q = 1 - \frac{1}{N-1}.$$

Finally, we consider the spectrum on Q on the eigenspaces of P corresponding to $n + \ell = 1$. In this case, as noted above, the eigenfunctions are antisymmetric, so that if we are concerned only with the spectrum of Q on the subspace of symmetric functions (which is all that is of significance for Kac’s application to the Boltzmann equation), we can ignore these eigenspaces. However, if we define $\eta_{0,1}(\vec{v}) = v_1 - v_2$ and $\eta_{1,0}(\vec{v}) = |v_1|^2 - |v_2|^2$, we find, as above, that

$$(10.17) \quad Q\eta_{1,0} = \left(1 - (1 - B_1)\frac{1}{N-1}\right)\eta_{1,0} \quad \text{and}$$

$$(10.18) \quad Q\eta_{0,1} = \left(1 - (1 - B_1)\frac{1}{N-1}\right)\eta_{0,1}. \quad \square$$

Now that we have all of our eigenvalues, we need to order them. By a simple comparison, we determine that for all N the largest eigenvalue of Q on our three invariant subspaces with $n + \ell = 2$ is $\lambda_{1,1}^Q$. This is true for all choices of b , since the only dependence on b in these eigenvalues is a common factor of $(1 - B_2)$.

It is worth noting, however, that for large N

$$\lambda_{1,1}^Q = 1 - (1 - B_2)\frac{1}{N} + \mathcal{O}\left(\frac{1}{N^2}\right) \quad \text{and} \quad \lambda_{2,0}^Q = 1 - (1 - B_2)\frac{1}{N} + \mathcal{O}\left(\frac{1}{N^2}\right),$$

so that these eigenvalues merge as N tends to infinity. Still, for all finite N ,

$$\lambda_{1,1}^Q = 1 - (1 - B_2)\frac{1}{N-1}$$

is strictly larger.

Next, the invariant subspaces of Q with $n + \ell = 1$ are also eigenspaces of Q with the eigenvalue

$$\lambda_{0,1}^Q = \lambda_{1,0}^Q = 1 - (1 - B_1)\frac{1}{N-1}.$$

In summary, the largest eigenvalue of Q on the invariant subspaces $V_{n,\ell}$ in $L^2(X_N, d\sigma_N)$ with $n + \ell = 1, 2$ and $N \geq 3$ is either

$$1 - (1 - B_2)\frac{1}{N-1} \quad \text{or} \quad 1 - (1 - B_1)\frac{1}{N-1},$$

depending on which of these is larger. In particular,

$$(10.19) \quad \Delta_3 \leq \min\{(1 - B_2), (1 - B_1)\}\frac{3}{2}.$$

With the above arguments we have all of the ingredients needed to prove Theorem 1.1.

Proof of Theorem 1.1. First, we wish to apply Lemma 4.1 to estimate Δ_3 in terms of Δ_2 . In (9.4), we have set $\mu_3^* = 11/20$, and with this choice of the threshold, we have seen that there is just one eigenvalue of P between μ_3^* and 1, namely, the gap

eigenvalue $\mu_3 = \mu_{1,1}(3) = \mu_{0,1}(3) = \mu_{1,0}(3)$. Thus, from Lemma 4.1 and the eigenvalue computations in Lemma 10.1, either the gap eigenvalue of Q for $N = 3$ is

$$(10.20) \quad \max \left\{ 1 - (1 - B_2) \frac{1}{N-1}, 1 - (1 - B_1) \frac{1}{N-1} \right\}$$

or else

$$(10.21) \quad \Delta_3 \geq \frac{3}{2} \left(1 - \frac{11}{20} \right) \Delta_2 = \frac{27}{40} \Delta_2.$$

If (10.20) does give the gap eigenvalue of Q for $N = 3$, then

$$(10.22) \quad \Delta_3 = \min \{ (1 - B_1), (1 - B_2) \} \frac{3}{2}.$$

Since, according to Lemma 4.1, at least one of (10.21) and (10.22) is true, the condition

$$(10.23) \quad \frac{27}{40} \Delta_2 \geq \frac{3}{2} \min \{ (1 - B_1), (1 - B_2) \}$$

and (10.19) ensure that (10.22) is true and thus give us the gap eigenvalue for $N = 3$. Note that the condition (10.23) is equivalent to the condition (1.6) in Theorem 1.1.

Now we proceed by induction. For any $n \geq 4$, assume that

$$(10.24) \quad \Delta_{N-1} = \min \{ (1 - B_1), (1 - B_2) \} \frac{N-1}{N-2}.$$

In (9.3) we have set

$$\mu_N^* = \frac{1}{N-1}$$

for all $N \geq 4$, and we have seen that the only eigenvalues μ of P with $1 > \mu \geq \mu_N^*$ are the gap eigenvalue $\mu_N = \mu_{1,1}(N) = \mu_{0,1}(N) = \mu_{1,0}(N)$, and for $N > 4$, $\mu_{2,0}(N)$. Thus, by Lemma 4.1 and the eigenvalue computations in Lemma 10.1, either the gap eigenvalue of Q for N is

$$(10.25) \quad \max \left\{ 1 - (1 - B_2) \frac{1}{N-1}, 1 - (1 - B_1) \frac{1}{N-1} \right\}$$

or else

$$(10.26) \quad \Delta_N > \frac{N}{N-1} \left(1 - \frac{1}{N-1} \right) \Delta_{N-1}.$$

There is strict inequality in (10.26) since all remaining eigenvalues of P not taken into account in (10.25) are strictly less than μ_N^* . By the inductive hypothesis (10.24) yields

$$\Delta_N > \min \{ (1 - B_1), (1 - B_2) \} \frac{N}{N-1}.$$

This is impossible, as the trial functions leading to (10.25) yield the upper bound

$$(10.27) \quad \Delta_N \leq \min \{ (1 - B_1), (1 - B_2) \} \frac{N}{N-1}.$$

Thus equality holds in (10.27), which completes the proof of the inductive step. Because of the strict inequality in (10.26), the only eigenfunctions with the gap eigenvalue are found in the invariant subspaces considered here, i.e., in the $V_{n,\ell}$ with $0 < n + \ell \leq 2$. By the results of Lemma 10.1, this yields the statement in Theorem 1.1 concerning the gap eigenfunctions of Q . \square

Proof of Theorem 1.2. We proceed as in the previous proof except that, for low values of N , we shall use a different choice for the threshold μ_N^* , namely,

$$(10.28) \quad \mu_N^* = \mu_{2,2}(N).$$

We know from Lemma 9.1 that for all $N \leq 7$

$$\mu_{n,\ell}(N) \leq \mu_{2,2}(N) \quad \text{for all } n + \ell > 2.$$

Thus, at least for such N , the only eigenvalues μ of P with $\mu > \mu_N^* = \mu_{2,2}(N)$ are those with $n + \ell \leq 2$. We have already computed the gap for Q on the invariant subspaces containing these eigenvalues, and we have found that the gap in these subspaces is

$$\tilde{\Delta}_N = \min\{ (1 - B_1), (1 - B_2) \} \frac{N}{N - 1}.$$

If for any $N_0 \geq 3$ it turns out that $\tilde{\Delta}_{N_0} = \Delta_{N_0}$, the gap on the whole space, then we can switch from that point onwards to the use of $\mu_N^* = \mu_{0,2}(N)$ as in the proof of Theorem 1.1 to show that $\tilde{\Delta}_N = \Delta_N$ for all $N \geq N_0$ and that the eigenfunctions are exactly as claimed for all $N > N_0$.

We now show that it is always the case that $\tilde{\Delta}_{N_0} = \Delta_{N_0}$ for some $N_0 \leq 7$. To do this, pick any value $N_1 \geq 4$, and suppose that for $3 \leq j \leq N_1$ we have

$$(10.29) \quad \Delta_j < \min\{ (1 - B_1), (1 - B_2) \} \frac{j}{j - 1}.$$

Then by Lemmas 4.1 and 9.1, by using the value $\mu_j^* = \mu_{2,2}(j)$, for $3 \leq j \leq N_1$, we have

$$\Delta_{N_1} \geq \frac{N_1}{2} \prod_{j=3}^{N_1} (1 - \mu_{2,2}(j)) \Delta_2.$$

By using the hypothesis $\Delta_2 = 2(1 - B_1)$, we have

$$\Delta_{N_1} \geq \frac{N_1}{2} \prod_{j=3}^{N_1} (1 - \mu_{2,2}(j)) 2(1 - B_1).$$

Of course we can rewrite this as

$$(10.30) \quad \begin{aligned} \Delta_{N_1} &\geq \frac{N_1}{2} \left(\prod_{j=4}^{N_1} (1 - \mu_{0,2}(j)) \right) (1 - \mu_{2,2}(3)) \left(\prod_{j=4}^{N_1} \frac{(1 - \mu_{2,2}(j))}{(1 - \mu_{0,2}(j))} \right) 2(1 - B_1) \\ &= \frac{N_1}{N_1 - 1} \left(\prod_{j=4}^{N_1} \frac{(1 - \mu_{2,2}(j))}{(1 - \mu_{0,2}(j))} \right) \frac{9}{10} (1 - B_1), \end{aligned}$$

since, as in the last proof, $(1 - \mu_{2,2}(3)) = 9/20$, and

$$\frac{N_1}{2} \prod_{j=4}^{N_1} (1 - \mu_{0,2}(j)) = \frac{N_1}{N_1 - 1}.$$

Now, by direct computation, we find that

$$\prod_{j=4}^7 \frac{(1 - \mu_{2,2}(j))}{(1 - \mu_{0,2}(j))} = \frac{558018643}{495720000} > \frac{10}{9}.$$

For $N_1 \geq 7$, this would lead to

$$\Delta_{N_1} > \frac{N_1}{N_1 - 1} (1 - B_1),$$

and this is impossible, since we have a trial function showing that the gap cannot be so large. Hence it must be that (10.29) is false for some $j \leq 7$. By what we have said above, from this point onward, we can proceed as in the proof of Theorem 1.1, and we obtain Theorem 1.2. \square

While the results presented here cover a very wide range of models, it is possible to come up with choices of b for which $\Delta_2 \neq 2(1 - B_1)$. If one found a need to deal with such an example, one might have to go deeper into the spectrum of P . It is very likely that Lemma 9.1 holds for all $N \geq 3$, based on extensive computations. These computations also show that, as N increases, $\mu_{2,1}(N)$ comes very close to $\mu_{2,2}(N)$, so that to get much more leverage one would need to compute all of the eigenvalues of Q on the smallest invariant subspaces of Q that contain both of these eigenspaces of P . This could be done by using the methods illustrated above, but the computations would be considerably more involved than the ones we have presented in this section. Thus, having treated a wide range of models, we shall conclude our discussion of Q here. In the brief final section, we discuss a point we raised earlier concerning bounds on Jacobi polynomials.

11. Bounds on Jacobi polynomials. As alluded to in section 8, the identity (7.9), together with the trivial bound on the $|\kappa_{n,\ell}| \leq 1$, which comes from the fact that K is a Markov operator, will for certain regions provide a stronger bound than (8.8), the bound of Nevai, Erdelyi, and Magnus. We close this section by showing how (7.9) can be used to obtain better bounds.

To begin, write

$$(11.1) \quad b^{2\beta-1} \frac{P_n^{\alpha,\beta}(-1 + 2b^2)}{P_n^{\alpha,\beta}(1)} \leq \frac{2e}{\pi} \frac{\Gamma(n+1)}{b(1-b^2)^{\alpha+1/2}} \frac{2 + \sqrt{\alpha^2 + \beta^2}}{2n + \alpha + \beta + 1} \frac{\Gamma(n + \beta + 1)}{\Gamma(n + \alpha + \beta + 1)} \frac{\Gamma(\alpha + 1)^2}{\Gamma(n + \alpha + 1)},$$

where $\beta = l + 1/2$, with l an integer. In regions where the right-hand side of the above equation becomes larger than one, the simple bound

$$b^{2\beta-1} \frac{P_n^{\alpha,\beta}(-1 + 2b^2)}{P_n^{\alpha,\beta}(1)} \leq 1$$

becomes stronger. In the region $2n + 1 < \alpha < \beta$, we find $\frac{2 + \sqrt{\alpha^2 + \beta^2}}{2n + \alpha + \beta + 1} > \frac{1}{4}$. This plus Stirling's formula with the remainder yields

$$\begin{aligned}
& \frac{2e}{\pi} \frac{\Gamma(n+1)}{b(1-b^2)^{\alpha+1/2}} \frac{2+\sqrt{\alpha^2+\beta^2}}{2n+\alpha+\beta+1} \frac{\Gamma(n+\beta+1)}{\Gamma(n+\alpha+\beta+1)} \frac{\Gamma(\alpha+1)^2}{\Gamma(n+\alpha+1)} \\
& > \frac{e^n}{\sqrt{2\pi}} \frac{\Gamma(n+1)}{b(1-b^2)^{\alpha+1/2}} \frac{\alpha^{\alpha+1/2-n}\beta^{-\alpha}}{(1+\frac{\alpha}{\beta})^{n+\alpha+\beta+1}} \frac{(1+\frac{1}{\alpha})^{2\alpha+1}}{(1+\frac{n+1}{\alpha})^{n+\alpha+1/2}(1+\frac{n+1}{\alpha+\beta})^\beta} * r \\
& > \frac{e^n}{\sqrt{2\pi}} \frac{\Gamma(n+1)}{b(1-b^2)^{\alpha+1/2}} \frac{\alpha^{\alpha+1/2-n}\beta^{-\alpha}}{2^{2n+2\alpha+2\beta+3/2}} r,
\end{aligned}$$

where $r = (1 - \frac{1}{12(n+\alpha+\beta+1)})(1 - \frac{1}{12(n+\alpha+1)})$ and n is assumed to be fixed. Choosing $b(1-b^2)^{\alpha+1/2}$ so that the last inequality is greater than one provides a region where the combination of (7.9) and $|\kappa_{n,\ell}| \leq 1$ does better than (7.9). It would be interesting to obtain better bounds on $|\kappa_{n,\ell}|$ by direct analysis of K and to use these to sharpen the argument just made.

Acknowledgment. We thank Doron Lubinsky for valuable discussions concerning (8.8), the bound of Nevai, Erdelyi, and Magnus, and related results.

REFERENCES

- [1] R. ASKEY, *Orthogonal Polynomials and Special Functions*, CBMS-NSF Regional Conf. Ser. in Appl. Math. 21, SIAM, Philadelphia, 1975.
- [2] E. CARLEN, M. C. CARVALHO, AND M. LOSS, *Many-body aspects of approach to equilibrium*, J. EDP, 11 (2000), pp. 1–12.
- [3] E. CARLEN, M. CARVALHO, AND M. LOSS, *Determination of the spectral gap for Kac's master equation and related stochastic evolution*, Acta Math., 191 (2003), pp. 1–54.
- [4] E. CARLEN AND X. LU, *Fast and slow convergence to equilibrium Maxwellian molecules via Wild sums*, J. Statist. Phys., 112 (2003), pp. 59–134.
- [5] I. S. GRADSHTEYN AND I. M. RYZHIK, *Tables of Integrals Series and Products*, Academic Press, New York, 1965.
- [6] E. JANVRESSE, *Spectral Gap for Kac's model of Boltzmann Equation*, Ann. Probab., 29 (2001), pp. 288–304.
- [7] M. KAC, *Foundations of kinetic theory*, in Proceedings of the 3rd Berkeley Symposium on Mathematical Statistics and Probability, J. Neyman, ed., Vol. 3, University of California, Berkeley, 1956, pp. 171–197.
- [8] T. H. KOORNWINDER, *The addition formula for Jacobi polynomials. I, summary of results*, Indag. Math., 34 (1972), pp. 188–191.
- [9] D. MASLEN, *The eigenvalues of Kac's master equation*, Math. Z., 243 (2003), pp. 291–331.
- [10] D. MORGENSTERN, *General existence and uniqueness proof for spatially homogeneous solutions of the Maxwell-Boltzmann equation in the case of Maxwellian Molecules*, Proc. Natl. Acad. Sci. USA, 40 (1954), pp. 719–721.
- [11] D. MORGENSTERN, *Analytical studies related to the Maxwell-Boltzmann equation*, J. Ration. Mech. Anal., 4 (1955), pp. 154–183.
- [12] P. NEVAI, T. ERDÉLYI, AND A. MAGNUS, *Generalized Jacobi weights, Christoffel functions, and Jacobi polynomials*, SIAM J. Math. Anal., 25 (1994), pp. 602–614.
- [13] G. SZEGÖ, *Orthogonal Polynomials*, Amer. Math. Soc. Colloq. Publ. 23, AMS, Providence, RI, 1939.

ORBITAL STABILITY OF BOUND STATES OF SEMICLASSICAL NONLINEAR SCHRÖDINGER EQUATIONS WITH CRITICAL NONLINEARITY*

TAI-CHIA LIN[†] AND JUNCHENG WEI[‡]

Abstract. We consider the orbital stability of single-spike bound states of semiclassical nonlinear Schrödinger equations with critical nonlinearity and a trap potential. Due to the effect of the trap potential, we derive the asymptotic expansion formulas and obtain the necessary conditions for orbital stability and instability of single-spike bound states. Our argument is applied to two-component systems of nonlinear Schrödinger equations with a common trap potential, cubic nonlinearity in two spatial dimensions. The orbital stability of bound states with spikes of these systems is investigated. Our results show the existence of stable spikes in two-dimensional Bose–Einstein condensates.

Key words. orbital stability, spike, trap potential

AMS subject classifications. 35J50, 35Q55

DOI. 10.1137/070683842

1. Introduction. The nonlinear Schrödinger (NLS) equation with a trap potential is central to the understanding of many physical phenomena. For example, it has become a well-known model referred to as the Gross–Pitaevskii equation governing the evolution of Bose–Einstein condensates (BEC) given by

$$(1.1) \quad -i\hbar \frac{\partial \psi}{\partial t} = \frac{\hbar^2}{2m} \Delta \psi - V_{\text{trap}} \psi - \mu |\psi|^2 \psi$$

for $x \in \mathbb{R}^N$, $N \leq 3$, and $t > 0$, where $\psi = \psi(x, t) \in \mathbb{C}$ is the wavefunction of BEC and $V_{\text{trap}} = V_{\text{trap}}(x)$ is the trap potential. Also, \hbar is the Planck constant, m is the atom mass, and $\mu \sim 4\pi \frac{\hbar^2}{2m} a$, where a denotes the s-wave scattering length.

In BEC, spikes may occur when the s-wave scattering length is negative and large. Due to Feshbach resonance, the s-wave scattering length of a single condensate can be tuned over a very large range by adjusting the externally applied magnetic field. As the s-wave scattering length of a single condensate is negative and large enough, the interactions of atoms are strongly attractive and the associated condensate tends to increase its density at the center of the trap potential in order to lower the interaction energy (cf. [28]). Under the effect of trap potentials, spikes of BEC are observable by physical experiments (cf. [10]) so there must be stability to ensure spikes appearing in the condensate wavefunction (cf. [7]). In [24], stable bright solitons (spikes) of BEC can be observed by numerical simulations, provided that the strength of the trap potential exceeds a threshold value. Here we want to develop mathematical theorems to support the existence of stable spikes in BEC.

*Received by the editors February 27, 2007; accepted for publication (in revised form) August 15, 2007; published electronically April 25, 2008.

<http://www.siam.org/journals/sima/40-1/68384.html>

[†]Department of Mathematics, National Taiwan University, Taipei, 106 Taiwan (tclin@math.ntu.edu.tw). The research of this author was partially supported by a research grant from NSC of Taiwan.

[‡]Department of Mathematics, The Chinese University of Hong Kong, Shatin, Hong Kong (wei@math.cuhk.edu.hk). The research of this author was partially supported by an Earmarked Grant from RGC of Hong Kong.

To get spikes in BEC, we may assume the s-wave scattering length a , i.e., μ is negative and large. Setting $\hbar^2 = \hbar^2/(2m\mu)$, $V_{trap}(x) = \mu^{-1}V(x)$, and suitable time scale, (1.1) with negative and large μ can be equivalent to a semiclassical NLS given by

$$(1.2) \quad -i\hbar \frac{\partial \psi}{\partial t} = \hbar^2 \Delta \psi - V\psi + |\psi|^2 \psi, \quad x \in \mathbb{R}^N, \quad t > 0,$$

where $0 < \hbar \ll 1$ is a small parameter (semiclassical limit) and $V = V(x)$ is a smooth nonnegative function. We may generalize (1.2) to an NLS equation having the form

$$(1.3) \quad -i\hbar \frac{\partial \psi}{\partial t} = \hbar^2 \Delta \psi - V\psi + |\psi|^{p-1} \psi, \quad x \in \mathbb{R}^N, \quad t > 0,$$

with critical nonlinearity

$$(1.4) \quad p = 1 + \frac{4}{N}, \quad N \geq 1.$$

In particular, when $N = 2$, (1.3) is exactly the same as (1.2).

Bound states of (1.3) are of the form $\psi(x, t) = e^{i\lambda t/\hbar} u(x)$, where $\lambda > 0$ and u satisfies the following nonlinear elliptic equation:

$$(1.5) \quad \hbar^2 \Delta u - (V + \lambda)u + u^p = 0, \quad u \in H^1(\mathbb{R}^N), \quad u > 0 \quad \text{in } \mathbb{R}^N,$$

with zero Dirichlet boundary condition, i.e., $u(x) \rightarrow 0$ as $|x| \rightarrow \infty$.

In the case when $V(x) \equiv 0$, for any $\lambda > 0$, problem (1.5) admits a unique radially symmetric ground state, which is stable for any $\lambda > 0$ if $p < 1 + \frac{4}{N}$ and unstable for any $\lambda > 0$ if $p \geq 1 + \frac{4}{N}$ (see [5], [6], and [36]).

When $V(x) \not\equiv 0$, the existence of single- or multiple-spike solutions of (1.5) was first established by Floer and Weinstein [11] in the one-dimensional case, i.e., $N = 1$ and $1 < p < 5$, and later extended by Oh [25], [26] to the higher-dimensional case, i.e., $N \geq 2$ and $1 < p < \frac{N+2}{N-2}$ under the condition that the trap potential V has nondegenerate critical points. When the trap potential V becomes degenerate, there have been many works in recent years. The reader may refer to [1], [3], [17], [8], [9] [19], [29], [30], [32], [33], and the references therein.

The trap potential V may also play a crucial role in the orbital (dynamic) stability of single-spike bound states. As the trap potential V is switched off, it is well known that all bound states of (1.3) with the condition (1.4) are orbitally unstable if the dimension $N = 2$ (cf. [37]). To stabilize bound states, one has to turn on the trap potential. However, in general, some nonzero trap potentials may still cause dynamic instability in BEC. For instance, one may find bending-wave instability of vortex ring dynamics under some nonzero trap potentials (cf. [18]). Consequently, to get the dynamic stability of single-spike bound states, we have to choose trap potentials properly. For suitable trap potentials, Oh [26] and Grillakis, Shatah, and Strauss [14] proved that when $N = 1$, the single-spike bound state (concentrating at local nondegenerate minimum of the trap potential V) is stable if $1 < p < 1 + \frac{4}{N}$ and unstable if $p > 1 + \frac{4}{N}$. Generically, the case of $p = 1 + \frac{4}{N}$ is left open and referred to as a critical case in the literature. In this paper, we give an affirmative answer for such a case by studying the orbital stability and instability of single-spike bound states when the trap potential V has nondegenerate critical points.

In [25] and [26], a single-spike bound state solution u_h of (1.5) can be obtained, provided the trap potential V is of class $(V)_a$ and fulfills other conditions. Hereafter,

we set u_h as a single-spike bound state constructed in [25] and [26] and satisfying (1.5). Of course, the trap potential V is also of class $(V)_a$ and fulfills other conditions in [25] and [26]. Hence $\psi_h(x, t) = e^{i\lambda t/h}u_h(x)$ may form an orbit of (1.3). From [14], the orbital stability of ψ_h 's is defined as follows: For all $\epsilon > 0$, there exists $\delta > 0$ such that if $\|\psi_0 - u_h\|_{H^1} < \delta$ and ψ is a solution of (1.3) in some interval $[0, t_0]$ with $\psi|_{t=0} = \psi_0$, then $\psi(t, \cdot)$ can be extended to a solution in $0 \leq t < \infty$ and

$$\sup_{0 < t < \infty} \inf_{s \in \mathbb{R}} \|\psi(\cdot, t) - \psi_h(\cdot, s)\|_{H^1} < \epsilon.$$

Otherwise, the orbit ψ_h is called orbitally unstable. To check the orbital stability of ψ_h , we use the linearized operator defined by

$$(1.6) \quad L_h = h^2\Delta - (V + \lambda) + pu_h^{p-1}, \quad p = 1 + \frac{4}{N}.$$

Observe that u_h may depend on the variable λ . Moreover, we assume u_h to be nondegenerate due to [16]. Let $n(L_h)$ be the number of positive eigenvalues of L_h and

$$(1.7) \quad d(\lambda) = \int_{\mathbb{R}^N} \left[\frac{h^2}{2} |\nabla u_h|^2 + \frac{1}{2}(V + \lambda)u_h^2 - \frac{1}{p+1}u_h^{p+1} \right].$$

Assume that d is nondegenerate, i.e., $d'' \neq 0$. Let $p(d'') = 1$ if $d'' > 0$ and $p(d'') = 0$ if $d'' < 0$. According to the general theory of orbital stability of bound states (cf. [14], [15]), u_h is orbitally stable if $n(L_h) = p(d'')$ and orbitally unstable if $n(L_h) - p(d'')$ is odd (see page 309 of [15]).

It is worth noting that if $V \equiv C$ and $p = 1 + \frac{4}{N}$, then $d''(\lambda) = 0$, i.e., the function d becomes degenerate, where C is a positive constant. Consequently, we may assume that the trap potential V has nondegenerate critical points in order to derive the asymptotic expansion formulas for the operator L_h and the function d as the parameter h goes to zero. These formulas are crucial to obtaining the orbital stability and instability of single-spike bound states as follows.

THEOREM 1.1. *Let N be a positive integer and $p = 1 + \frac{4}{N}$. For $0 < h < 1$, let u_h be a bound state of (1.3) concentrated at a nondegenerate critical point x_0 of the potential V such that $\Delta V(x_0) \neq 0$. Let m denote the number of negative eigenvalues of the matrix $(\nabla^2 V(x_0))$. Suppose the parameter h is sufficiently small. Then u_h is orbitally stable if x_0 is a nondegenerate local minimum point of the potential V . Furthermore, u_h is orbitally unstable if $m - \frac{1}{2}(1 + \frac{\Delta V(x_0)}{|\Delta V(x_0)|})$ is even.*

Remark. In [12], Fukuizumi considered the orbital stability of the standing wave solution to (1.3) in the critical case $p = 1 + \frac{4}{N}$ under the following conditions:

$$(1.8) \quad h = 1, \quad V(x) = V(|x|), u(x, t) = u(|x|, t).$$

He studied the stability of the standing wave solution $e^{i\lambda t}u(x)$ for λ large (in the radially symmetric class). After suitable scaling, the standing wave solution satisfies

$$(1.9) \quad \epsilon^2 \Delta u - (1 + \epsilon^2 V(|x|))u + u^p = 0, u > 0, \quad \text{where } \epsilon^2 = \frac{1}{\lambda}.$$

There is a slightly subtle difference between (1.5) and (1.9). However, the main difference is that we consider the full functional space here. The method of proving Theorem 1.1 can also be applied to obtain a more general result to Fukuizumi's problem for general $V(x)$.

Another motivation of studying (1.3) in the critical case may come from two-component systems of NLS equations which describe a double condensate, i.e., a binary mixture of BEC (cf. [28]). To get stable spikes of a double condensate with two spatial dimensions, we study orbitally stable bound states with spikes of a two-component system of NLS equations given by

$$(1.10) \quad \begin{cases} -ih \frac{\partial \Phi}{\partial t} = h^2 \Delta \Phi - V \Phi + |\Phi|^2 \Phi + \beta |\Psi|^2 \Phi, \\ -ih \frac{\partial \Psi}{\partial t} = h^2 \Delta \Psi - V \Psi + |\Psi|^2 \Psi + \beta |\Phi|^2 \Psi \end{cases}$$

for $x \in \mathbb{R}^2$ and $t > 0$, where $V = V(x)$ is a smooth nonnegative function, $\beta \in \mathbb{R}$ is a nonzero constant, and $0 < h \ll 1$ is a small parameter. Bound states of (1.10) are of the form $\Phi(x, t) = e^{i\lambda t/h} u(x)$ and $\Psi(x, t) = e^{i\lambda t/h} v(x)$, where (u, v) satisfies the following nonlinear elliptic system:

$$(1.11) \quad \begin{cases} h^2 \Delta u - (V + \lambda)u + u^3 + \beta uv^2 = 0, & x \in \mathbb{R}^2, \\ h^2 \Delta v - (V + \lambda)v + v^3 + \beta u^2 v = 0, & x \in \mathbb{R}^2, \\ u(x), v(x) > 0, & u, v \in H^1(\mathbb{R}^2). \end{cases}$$

Note that in \mathbb{R}^2 , the nonlinearity u^3, v^3 are a critical nonlinearity by the simple algebra $p = 3 = 1 + \frac{4}{N}$ with $N = 2$. The system (1.11) admits a bound state solution $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} v_h)$, where $\beta > -1$ and u_h satisfies (1.5). Generically, such a solution may be neither a unique positive solution nor a ground state solution of the system (1.11). Thus the stability problem is nontrivial. Here we want to get the orbital stability of such a solution using suitable trap potentials V 's. To study the orbital stability of such a bound state solution, we set the linearized operator of (1.10) around $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} v_h)$ given by

$$(1.12) \quad \mathbb{L}_h \begin{pmatrix} \phi \\ \psi \end{pmatrix} = \begin{pmatrix} h^2 \Delta \phi - (V + \lambda)\phi + \frac{3+\beta}{1+\beta} u_h^2 \phi + \frac{2\beta}{1+\beta} u_h^2 \psi \\ h^2 \Delta \psi - (V + \lambda)\psi + \frac{3+\beta}{1+\beta} u_h^2 \psi + \frac{2\beta}{1+\beta} u_h^2 \phi \end{pmatrix}.$$

Furthermore, we also need a function defined as follows:

$$(1.13) \quad \begin{aligned} d(\lambda_1, \lambda_2) &= \int_{\mathbb{R}^2} \frac{h^2}{2} |\nabla u_{h,\lambda_1,\lambda_2}|^2 + \frac{V(x) + \lambda + \lambda_1}{2} u_{h,\lambda_1,\lambda_2}^2 - \frac{1}{4} \int_{\mathbb{R}^2} u_{h,\lambda_1,\lambda_2}^4 \\ &\quad + \int_{\mathbb{R}^2} \frac{h^2}{2} |\nabla v_{h,\lambda_1,\lambda_2}|^2 + \frac{V(x) + \lambda + \lambda_2}{2} v_{h,\lambda_1,\lambda_2}^2 - \frac{1}{4} \int_{\mathbb{R}^2} v_{h,\lambda_1,\lambda_2}^4 \\ &\quad - \frac{\beta}{2} \int_{\mathbb{R}^2} u_{h,\lambda_1,\lambda_2}^2 v_{h,\lambda_1,\lambda_2}^2, \end{aligned}$$

where $(u_{h,\lambda_1,\lambda_2}, v_{h,\lambda_1,\lambda_2})$ is the solution of

$$(1.14) \quad \begin{cases} h^2 \Delta u - (V + \lambda + \lambda_1)u + u^3 + \beta uv^2 = 0 & \text{in } \mathbb{R}^2, \\ h^2 \Delta v - (V + \lambda + \lambda_2)v + v^3 + \beta u^2 v = 0 & \text{in } \mathbb{R}^2, \end{cases}$$

such that $(u_{h,\lambda_1,\lambda_2}, v_{h,\lambda_1,\lambda_2}) \rightarrow (\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} v_h)$ as $|\lambda_1| + |\lambda_2| \rightarrow 0$.

Suppose the solution $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$ is nondegenerate, i.e., the operator \mathbb{L}_h has no zero eigenvalue. Let $n(\mathbb{L}_h)$ denote the positive eigenvalues of \mathbb{L}_h , and set p as the number of positive eigenvalues of the Hessian matrix $(\nabla^2 d(0, 0))$. From [14] and [15], we know that the solution $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$ is orbitally stable if $n(\mathbb{L}_h) = p$ and orbitally unstable if $n(\mathbb{L}_h) - p$ is odd. The parameter β may affect the orbital stability of the solution $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$. Now we state our result as follows.

THEOREM 1.2. *For $0 < h < 1$, let u_h be a single-spike solution concentrated at a local minimum point of the function V . Suppose the parameter h is sufficiently small. Then $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$ is an orbitally stable solution to (1.10) if $0 < \beta \neq 1$.*

Remark. The orbital instability of $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$ for $-1 < \beta < 0$ can also be investigated. However, the condition is quite complicated so we may omit it here. On the other hand, as $\beta = 1$, the system (1.11) may have infinitely many solutions with the form $(u, v) = (w, \eta w)$ for $\eta \neq 0$, where w is the solution of $h^2 \Delta w - (V + \lambda)w + (1 + \eta^2)w^3 = 0$ in \mathbb{R}^2 . This may provide a reason to ignore the case $\beta = 1$ in Theorem 1.2.

For the existence of other bound states to the system (1.11), the reader may refer to [2], [4], [13], [20], [21], [23], [31], and the references therein. Our result here seems to be the first in studying the orbital stability of (1.11) with a trapping potential.

The argument of Theorem 1.2 is applicable to studying another two-component system of NLS equations having symbiotic bright solitons (cf. [22] and [27]) given by

$$(1.15) \quad \begin{cases} -ih \frac{\partial \Phi}{\partial t} = h^2 \Delta \Phi - V\Phi - |\Phi|^2 \Phi + \beta |\Psi|^2 \Phi, \\ -ih \frac{\partial \Psi}{\partial t} = h^2 \Delta \Psi - V\Psi - |\Psi|^2 \Psi + \beta |\Phi|^2 \Psi \end{cases}$$

for $x \in \mathbb{R}^2$ and $t > 0$, where $V = V(x)$ is a smooth nonnegative function, $\beta \in \mathbb{R}$ is a nonzero constant, and $0 < h \ll 1$ is a small parameter. It is remarkable that the coefficients of the terms $|\Phi|^2 \Phi$ and $|\Psi|^2 \Psi$ of the system (1.15) have opposite sign to those of the system (1.10). As for the system (1.11), bound states of (1.15) are of the form $\Phi(x, t) = e^{i\lambda t/h} u(x)$ and $\Psi(x, t) = e^{i\lambda t/h} v(x)$, where (u, v) satisfies the following nonlinear elliptic system:

$$(1.16) \quad \begin{cases} h^2 \Delta u - (V + \lambda)u - u^3 + \beta uv^2 = 0, & x \in \mathbb{R}^2, \\ h^2 \Delta v - (V + \lambda)v - v^3 + \beta u^2 v = 0, & x \in \mathbb{R}^2, \\ u(x), v(x) > 0, & u, v \in H^1(\mathbb{R}^2). \end{cases}$$

It is easy to check that the system (1.16) has a solution $(\frac{1}{\sqrt{\beta-1}}u_h, \frac{1}{\sqrt{\beta-1}}u_h)$ for $\beta > 1$. As for Theorem 1.2, we may have the following corollary.

COROLLARY 1.3. *For $0 < h < 1$, let u_h be a single-spike solution concentrated at a local minimum point of the function V . Suppose the parameter h is sufficiently small. Then $(\frac{1}{\sqrt{\beta-1}}u_h, \frac{1}{\sqrt{\beta-1}}u_h)$ is an orbitally stable solution to (1.15) if $\beta > 1$.*

The proof of Corollary 1.3 is quite similar to that of Theorem 1.2 so we may neglect the detailed proof here.

The rest of this paper is organized as follows. In section 2, we figure out the properties of u_h and the spectrum of the linearized operator L_h as the parameter h goes to zero. Then we state the proof of Theorem 1.1 in section 3. Finally, we provide the proof of Theorem 1.2 in section 4.

2. Properties of u_h . In this section, we study the properties of u_h , which is a single-spike solution concentrated at a nondegenerate critical point x_0 of $V(x)$. Let x_h be the unique local maximum point of u_h . So $x_h \rightarrow x_0$. Let us recall the following results of Grossi [16].

LEMMA 2.1.

- (1) $x_h = x_0 + o(h)$;
- (2) u_h is unique and nondegenerate, i.e., L_h has no zero eigenvalue.

Proof. (1) follows from Lemma 5.4 of [16], while (2) follows from Theorem 1.1 of [16]. \square

We need the following two lemmas. The first one is an asymptotic behavior of u_h .

LEMMA 2.2.

$$(2.1) \quad u_h(x_h + hy) = (V(x_h) + \lambda)^{\frac{1}{p-1}} w(\sqrt{V(x_h) + \lambda y}) + h^2 \phi_0 + o(h^2),$$

where w is the unique positive solution of

$$(2.2) \quad \Delta w - w + w^p = 0, w(0) = \max_{y \in \mathbb{R}^N} w(y), w > 0 \text{ in } \mathbb{R}^N, w \rightarrow 0 \text{ as } |y| \rightarrow +\infty,$$

ϕ_0 satisfies

$$(2.3) \quad \Delta \phi_0 - (V(x_h) + \lambda)\phi_0 + pw_{x_h}^{p-1}\phi_0 - \frac{1}{2} \sum_{i,j} V_{ij}(x_0)y_i y_j w_{x_h} = 0$$

with

$$(2.4) \quad w_{x_h}(y) := (V(x_h) + \lambda)^{\frac{1}{p-1}} w(\sqrt{V(x_h) + \lambda y}).$$

Proof. Note that for fixed s , $w_s(y)$ satisfies

$$(2.5) \quad \Delta w_s - (V(s) + \lambda)w_s + w_s^p = 0.$$

Let $\phi_h(y) = u_h(x_h + hy) - w_{x_h}(y)$. Then $|\phi_h| \rightarrow 0$ uniformly and ϕ_h satisfies

$$(2.6) \quad \Delta \phi_h - (V(x_h + hy) + \lambda)\phi_h + pw_{x_h}^{p-1}\phi_h + N(\phi_h) - (V(x_h + hy) - V(x_h))w_{x_h} = 0,$$

where $N(\phi_h) = (w_{x_h} + \phi_h)^p - w_{x_h}^p - pw_{x_h}^{p-1}\phi_h$. Note that $\nabla \phi_h(0) = 0$ and

$$(2.7) \quad \begin{aligned} (V(x_h + hy) - V(x_h)) &= (\nabla V(x_h))hy + \frac{1}{2} \sum_{i,j} V_{ij}(x_h)h^2 y_i y_j + O(h^3|y|^3) \\ &= o(h^2)|y| + \frac{1}{2} \sum_{i,j} V_{ij}(x_0)h^2 y_i y_j + o(h^2|y|^2). \end{aligned}$$

Here we have used Lemma 2.1(1).

Now we claim that $|\phi_h| \leq ch^2$. In fact, suppose not. We may assume that $|\phi_h|_{L^\infty} h^{-2} \rightarrow \infty$. Let $\tilde{\phi}_h = \frac{\phi_h}{|\phi_h|_{L^\infty}}$. Then $\tilde{\phi}_h$ satisfies

$$(2.8) \quad \Delta \tilde{\phi}_h - (V + \lambda)\tilde{\phi}_h + pw_{x_h}^{p-1}\tilde{\phi}_h + \frac{N(\phi_h)}{|\phi_h|_{L^\infty}} - \frac{(V(x_h + hy) - V(x_h))w_{x_h}}{|\phi_h|_{L^\infty}} = 0.$$

Note that by (2.7),

$$(2.9) \quad \frac{|V(x_h + hy) - V(x_h)||w_{x_h}|}{|\phi_h|_{L^\infty}} \leq \frac{h^2|y|^2|w_{x_h}|}{|\phi_h|_{L^\infty}} \leq o(1)|y|^2|w_{x_h}|$$

for $|y| \geq 1$. Let y_h be the global maximum point of $\tilde{\phi}_h$, i.e., $\tilde{\phi}_h(y_h) = \max_y \frac{\phi_h(y)}{|\phi_h|_{L^\infty}} = 1$. Then by (2.8) and (2.9) and the maximum principle, we have $|y_h| \leq C$. Here we have used the fact that $V \geq 0$ and $\lambda > 0$.

By the usual elliptic regularity theory, we may take a subsequence $\tilde{\phi}_h \rightarrow \bar{\phi}_0$, where $\bar{\phi}_0$ satisfies

$$(2.10) \quad \Delta \bar{\phi}_0 - (V(x_0) + \lambda)\bar{\phi}_0 + pw_{x_0}^{p-1}\bar{\phi}_0 = 0, \quad \nabla \bar{\phi}_0(0) = 0.$$

Since $\nabla \bar{\phi}_0(0) = 0$, we see that $\bar{\phi}_0 = \sum_{j=1}^N c_j \frac{\partial w_{x_0}}{\partial y_j}$, and hence $c_j = 0$. Consequently, $\bar{\phi}_0 \equiv 0$. This may contradict the fact that $1 = \tilde{\phi}_h(y_h) \rightarrow \bar{\phi}_0(y_0)$ for some y_0 . Therefore $|\phi_h| \leq ch^2$. Now we let $\phi_h = h^2\phi_0 + h^2\bar{\phi}_h$. Then, as for the previous argument, we may have $\bar{\phi}_h = o(1)$ and complete the proof of Lemma 2.2. \square

As in Proposition 3.6 of [19], one may get two lemmas, as follows.

LEMMA 2.3. *For each $s \in \mathbb{R}^N$, the map*

$$(2.11) \quad L_s \phi := \Delta \phi - (V(s) + \lambda)\phi + pw_s^{p-1}\phi$$

is invertible from K_s^\perp to C_s^\perp , where

$$K_s^\perp = \left\{ \phi \in H^2(\mathbb{R}^N) \mid \int_{\mathbb{R}^N} \phi \frac{\partial w_s}{\partial y_j}(y) dy = 0, j = 1, \dots, N \right\} \subset H^2(\mathbb{R}^N),$$

$$C_s^\perp = \left\{ \phi \in L^2(\mathbb{R}^N) \mid \int_{\mathbb{R}^N} \phi \frac{\partial w_s}{\partial y_j}(y) dy = 0, j = 1, \dots, N \right\} \subset L^2(\mathbb{R}^N).$$

LEMMA 2.4. *The map*

$$(2.12) \quad L_0 \phi := \Delta \phi - (V(x_0) + \lambda)\phi + pw_{x_0}^{p-1}\phi$$

admits the following eigenvalues:

$$\lambda_1 > 0, \lambda_2 = \dots = \lambda_{N+1} = 0, \lambda_{N+2} < 0,$$

where the kernel of L_0 is spanned by $\frac{\partial w_{x_0}}{\partial y_j}$, $j = 1, \dots, N$.

Our main result in this section is the following.

THEOREM 2.5. *The eigenvalue problem*

$$(2.13) \quad L_h \psi_h = \lambda_h \psi_h$$

admits eigenvalues

$$(2.14) \quad \lambda_{h,1} > \lambda_{h,2} > \dots > \lambda_{h,N+1} > \lambda_{h,N+2},$$

satisfying, as $h \rightarrow 0$, $\lambda_{h,1} \rightarrow \lambda_1 > 0$, $\lambda_{h,N+2} \rightarrow \lambda_{N+2} < 0$, and

$$(2.15) \quad \frac{\lambda_{h,j}}{h^2} \rightarrow c_0 \nu_{j-1}, \quad j = 2, \dots, N + 1,$$

where c_0 is a negative constant and ν_j 's are eigenvalues of the Hessian matrix $(\nabla^2 V(x_0))$.

Proof. We follow the proofs given in section 5 of [35]. Assume that $\|\psi_h\|_{L^2} = 1$. It is easy to see that for eigenvalues $\lambda_h \in [\frac{1}{2}\lambda_{N+2}, \frac{1}{2}\lambda_1]$, as $h \downarrow 0$, $\lambda_h \rightarrow \lambda_j$ for some j , where λ_j 's are given in Lemma 2.4. Now we focus on the case $\lambda_{h,j} \rightarrow 0$, i.e., $\lambda_h \rightarrow 0$ as $h \downarrow 0$. Then the corresponding eigenfunctions can be written as

$$(2.16) \quad \psi_h(x_h + hy) = \sum_{j=1}^N c_j \frac{\partial w_{x_h}}{\partial y_j}(y) + \psi_h^\perp(y),$$

where $\int_{\mathbb{R}^N} \frac{\partial w_{x_h}}{\partial y_j} \psi_h^\perp(y) dy = 0$, $j = 1, 2, \dots, N$. Hence by (2.13) and (2.16), ψ_h^\perp may satisfy

$$(2.17) \quad \begin{aligned} \Delta \psi_h^\perp - (V(x_h + hy) + \lambda) \psi_h^\perp + p w_{x_h}^{p-1}(y) \psi_h^\perp + p(u_h^{p-1}(x_h + hy) - w_{x_h}^{p-1}(y)) \psi_h^\perp \\ + \sum_j c_j L_h \frac{\partial w_{x_h}}{\partial y_j} = \lambda_h \left(\sum_j c_j \frac{\partial w_{x_h}}{\partial y_j} + \psi_h^\perp \right). \end{aligned}$$

Using (2.1) and (2.7) of Lemma 2.2, we have

$$(2.18) \quad \begin{aligned} L_h \frac{\partial w_{x_h}}{\partial y_j} &= \Delta \left(\frac{\partial w_{x_h}}{\partial y_j} \right) - (V(x_h + hy) + \lambda) \frac{\partial w_{x_h}}{\partial y_j} + p u_h^{p-1}(x_h + hy) \frac{\partial w_{x_h}}{\partial y_j} \\ &= (V(x_h) - V(x_h + hy)) \frac{\partial w_{x_h}}{\partial y_j} + p(u_h^{p-1}(x_h + hy) - w_{x_h}^{p-1}(y)) \frac{\partial w_{x_h}}{\partial y_j} \\ &= O(h^2). \end{aligned}$$

From Lemma 2.3, the map $L_{x_h} = \Delta - (V(x_h) + \lambda) + p w_{x_h}^{p-1}$ is invertible in the space $K_{x_h}^\perp$. Thus by (2.1), (2.7), and (2.18), (2.17) may give

$$(2.19) \quad \|\psi_h^\perp\|_{H^2} \leq c(h^2 + |\lambda_h|) \sum_j |c_j|.$$

Now we set $z_j(y) = \frac{\partial w_{x_h}}{\partial y_j}(y)$ for $j = 1, \dots, N$. Then multiplying (2.17) by z_k and integrating over \mathbb{R}^N , it is obvious that

$$(2.20) \quad \int_{\mathbb{R}^N} (L_h \psi_h^\perp) z_k dy + \sum_j c_j \int_{\mathbb{R}^N} \left(L_h \frac{\partial w_{x_h}}{\partial y_j} \right) z_k dy = \lambda_h \left(\sum_j c_j \int_{\mathbb{R}^N} z_j z_k \right) dy.$$

Here we have used the fact that $\psi_h^\perp \in K_{x_h}^\perp$. Using (2.18), (2.19), $\lambda_h = o(1)$, and integration by parts, we obtain

$$(2.21) \quad \int_{\mathbb{R}^N} (L_h \psi_h^\perp) z_k = \int_{\mathbb{R}^N} \psi_h^\perp L_h z_k = o(h^2)$$

and

$$\begin{aligned} \int_{\mathbb{R}^N} (L_h z_j) z_k &= \int_{\mathbb{R}^N} (V(x_h) - V(x_h + hy)) z_j z_k + p \int_{\mathbb{R}^N} (u_h^{p-1} - w_{x_h}^{p-1}) z_j z_k \\ (2.22) \qquad \qquad \qquad &:= I_1 + I_2, \end{aligned}$$

where

$$\begin{aligned} I_1 &= \int_{\mathbb{R}^N} (V(x_h) - V(x_h + hy)) z_j z_k \\ &= o(h^2) - \frac{h^2}{2} \sum_{l,m} V_{lm}(x_h) \int_{\mathbb{R}^N} y_l y_m z_j z_k \\ (2.23) \qquad \qquad \qquad &= -\frac{h^2}{2} V_{jk}(x_h) \int_{\mathbb{R}^N} y_j z_j y_k z_k. \end{aligned}$$

Here we have used (2.7) to get (2.23). For I_2 , we use Lemma 2.2:

$$\begin{aligned} I_2 &= p(p-1)h^2 \int_{\mathbb{R}^N} \phi_0 w_{x_h}^{p-2} z_j z_k + o(h^2) \\ &= -h^2 \int_{\mathbb{R}^N} L_0 \left(\frac{\partial^2 w}{\partial y_j \partial y_k} \right) \phi_0 + o(h^2) \\ &= -h^2 \int_{\mathbb{R}^N} (L_0 \phi_0) \frac{\partial^2 w}{\partial y_j \partial y_k} + o(h^2) \\ &= -\frac{h^2}{2} \sum_{l,m} V_{lm}(x_h) \int_{\mathbb{R}^N} y_l y_m w_{x_h} \frac{\partial^2 w}{\partial y_j \partial y_k} + o(h^2) \\ &= \frac{h^2}{2} \sum_{l,m} V_{lm}(x_h) \int_{\mathbb{R}^N} \frac{\partial}{\partial y_j} (y_l y_m w_{x_h}) z_k + o(h^2) \\ (2.24) \qquad \qquad \qquad &= \frac{h^2}{2} V_{jk}(x_h) \int_{\mathbb{R}^N} y_j y_k z_j z_k - \frac{h^2}{2} V_{jk}(x_h) \int_{\mathbb{R}^N} w_{x_h}^2 + o(h^2). \end{aligned}$$

Here we have used the following identity:

$$\begin{aligned} \int_{\mathbb{R}^N} \frac{\partial}{\partial y_j} (y_l y_m w_{x_h}) z_k &= \int_{\mathbb{R}^N} \delta_{jl} y_m w_{x_h} z_k + \int_{\mathbb{R}^N} \delta_{jm} y_l w_{x_h} z_k + \int_{\mathbb{R}^N} y_l y_m z_j z_k \\ &= \frac{1}{2} \int_{\mathbb{R}^N} \delta_{jl} y_m \frac{\partial}{\partial y_k} (w_{x_h}^2) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^N} \delta_{jm} y_l \frac{\partial}{\partial y_k} (w_{x_h}^2) + \int_{\mathbb{R}^N} y_l y_m z_j z_k \\ &= -(\delta_{jl} \delta_{km} + \delta_{jm} \delta_{lk}) \frac{1}{2} \int_{\mathbb{R}^N} w_{x_h}^2 + \int_{\mathbb{R}^N} y_l y_m z_j z_k. \end{aligned}$$

Combining (2.23) and (2.24), we have

$$(2.25) \quad I_1 + I_2 = -\frac{h^2}{2} V_{jk}(x_h) \int_{\mathbb{R}^N} w_{x_0}^2 + o(h^2).$$

Substituting (2.21) and (2.25) into (2.20), we may obtain $\lambda_j/h^2 \rightarrow c_0\nu_j$ for $j = 1, \dots, N$, where $c_0 = -\frac{\int_{\mathbb{R}^N} w_{x_0}^2 dy}{\int_{\mathbb{R}^N} z_j^2 dy}$ is a negative constant. The rest of the proof follows from a perturbation result, similar to pages 1473–1474 of [35]. We may omit the details here. \square

From Theorem 2.5, we may deduce the following.

THEOREM 2.6. *u_h is smooth in λ . Moreover, let $R_h = \frac{\partial u_h}{\partial \lambda}$. Then R_h satisfies*

$$(2.26) \quad L_h R_h - u_h = 0$$

and

$$(2.27) \quad R_h = \sum_{j=1}^N c_j^h z_j + R_0 + o(1),$$

where $R_0 = \frac{\partial}{\partial \lambda} w_{x_h} = (V(x_h) + \lambda)^{-1} (\frac{1}{p-1} w_{x_h} + \frac{1}{2} y \cdot \nabla w_{x_h})$ and $|c_j^h| = O(1)$ for $j = 1, \dots, N$.

Proof. Since u_h is unique and L_h is invertible, it is easy to see that u_h is smooth in λ and R_h satisfies (2.26). Now we decompose R_h as

$$R_h = \sum_{j=1}^N c_j^h z_j + R_0 + \bar{R}_h,$$

where $\int_{\mathbb{R}^N} z_j \bar{R}_h = 0$, $j = 1, \dots, N$. Then \bar{R}_h satisfies

$$(2.28) \quad L_h \bar{R}_h - u_h + L_h R_0 + \sum_{j=1}^N c_j^h L_h z_j = 0.$$

As for the proof of Theorem 2.5, we have

$$(2.29) \quad \|\bar{R}_h\|_{H^2} \leq c (|c_j^h| h^2 + \|L_h R_0 - u_h\|_{L^2}).$$

From (2.4) and (2.5), it is easy to check that $R_0 = \frac{\partial}{\partial \lambda} w_{x_h} = (V(x_h) + \lambda)^{-1} (\frac{1}{p-1} w_{x_h} + \frac{1}{2} y \cdot \nabla w_{x_h})$ and $L_{x_h} R_0 = w_{x_h}$ by differentiating (2.5) with respect to λ . Hence

$$L_h R_0 - u_h = p(w_h^{p-1}(x_h + hy) - w_{x_h}^{p-1}(y))R_0 - (V(x_h + hy) - V(x_h))R_0 + w_{x_h} - u_h.$$

Consequently, by Lemma 2.2 and (2.7), we obtain

$$(2.30) \quad \|L_h R_0 - u_h\|_{L^2} = O(h^2),$$

and then by (2.29),

$$(2.31) \quad \|\bar{R}_h\|_{H^2} = (1 + |c_j^h|)O(h^2).$$

To estimate c_j^h 's, we may multiply (2.28) by z_k and integrate over \mathbb{R}^N . Then

$$(2.32) \quad \sum_{j=1}^N c_j^h \int_{\mathbb{R}^N} (L_h z_j) z_k + \int_{\mathbb{R}^N} (L_h R_0 - u_h) z_k + \int_{\mathbb{R}^N} (L_h \bar{R}_h) z_k = 0.$$

Hence by (2.22) and (2.25), (2.32) may imply

$$(2.33) \quad |c_j^h| \leq \frac{C}{h^2} \left(\left| \int_{\mathbb{R}^N} (L_h R_0 - u_h) z_k \right| + \left| \int_{\mathbb{R}^N} (L_h \bar{R}_h) z_k \right| \right).$$

Using integration by parts and (2.18), we have

$$(2.34) \quad \int_{\mathbb{R}^N} (L_h \bar{R}_h) z_k = \int_{\mathbb{R}^N} (L_h z_k) \bar{R}_h = \|\bar{R}_h\|_{L^2} O(h^2).$$

Therefore by (2.30), (2.31), (2.33), and (2.34), we may obtain $|c_j^h| = O(1)$ and complete the proof. \square

3. Proof of Theorem 1.1. Let $p = 1 + \frac{4}{N}$. By Theorem 2.5, L_h has $m + 1$ positive eigenvalues and no zero eigenvalue, where m is the number of negative eigenvalues of the matrix $(\nabla^2 V(x_0))$. Let us now compute $d''(\lambda)$.

From (1.7), it is easy to get

$$d'(\lambda) = \frac{1}{2} \int_{\mathbb{R}^N} u_h^2$$

and hence

$$d''(\lambda) = \int_{\mathbb{R}^N} \frac{\partial u_h}{\partial \lambda} u_h = \int_{\mathbb{R}^N} R_h u_h.$$

By direct computations,

$$(3.1) \quad L_h \left(\frac{1}{p-1} u_h + \frac{1}{2} h y \cdot \nabla u_h \right) = \frac{1}{2} h y \cdot \nabla V(x_h + h y) u_h + (V(x_h + h y) + \lambda) u_h.$$

Consequently,

$$(3.2) \quad \begin{aligned} (V(x_h) + \lambda) \int_{\mathbb{R}^N} R_h u_h &= \int_{\mathbb{R}^N} R_h (V(x_h) - V(x_h + h y)) u_h \\ &\quad + \int_{\mathbb{R}^N} R_h (V(x_h + h y) + \lambda) u_h \\ &= \int_{\mathbb{R}^N} R_h \left(V(x_h) - V(x_h + h y) - \frac{1}{2} h y \cdot \nabla V(x_h + h y) \right) u_h \\ &\quad + \int_{\mathbb{R}^N} R_h L_h \left(\frac{1}{p-1} u_h + \frac{1}{2} h y \cdot \nabla u_h \right). \end{aligned}$$

Using integration by parts and (2.26), we may obtain

$$(3.3) \quad \begin{aligned} \int_{\mathbb{R}^N} R_h L_h \left(\frac{1}{p-1} u_h + \frac{1}{2} h y \cdot \nabla u_h \right) &= \int_{\mathbb{R}^N} (L_h R_h) \left(\frac{1}{p-1} u_h + \frac{1}{2} h y \cdot \nabla u_h \right) \\ &= \int_{\mathbb{R}^N} u_h \left(\frac{1}{p-1} u_h + \frac{1}{2} h y \cdot \nabla u_h \right) \\ &= \left(\frac{1}{p-1} - \frac{N}{4} \right) \int_{\mathbb{R}^N} u_h^2 = 0, \end{aligned}$$

since $p = 1 + \frac{4}{N}$. So by (2.7), (3.2), (3.3), Lemma 2.2, and Theorem 2.6, we have

$$\begin{aligned}
 d''(\lambda) &= \frac{1}{V(x_h) + \lambda} \int_{\mathbb{R}^N} R_h \left(V(x_h) - V(x_h + hy) - \frac{1}{2} hy \cdot \nabla V(x_h + hy) \right) u_h \\
 &= \frac{h^2}{V(x_h) + \lambda} \int_{\mathbb{R}^N} R_h \left[\sum_{i,j} -V_{ij}(x_h) y_i y_j \right] u_h + o(h^2) \\
 &= \frac{h^2}{V(x_h) + \lambda} \int_{\mathbb{R}^N} \left(\sum_l c_l^h z_l + R_0 + o(1) \right) \left(\sum_{i,j} -V_{ij}(x_h) y_i y_j \right) \\
 &\quad (w_{x_h} + O(h^2)) + o(h^2) \\
 &= -\frac{h^2}{V(x_h) + \lambda} \sum_i V_{ii}(x_h) \int_{\mathbb{R}^N} R_0 y_i^2 w_{x_h} + o(h^2) \\
 &= -\frac{h^2}{(V(x_h) + \lambda)^2} \sum_i V_{ii}(x_h) \int_{\mathbb{R}^N} \left(\frac{1}{p-1} w_{x_h} + \frac{1}{2} y \cdot \nabla w_{x_h} \right) y_i^2 w_{x_h} + o(h^2) \\
 &= -\frac{h^2}{(V(x_h) + \lambda)^2} \left(\frac{1}{p-1} - \frac{N+2}{4} \right) \sum_i V_{ii}(x_h) \int_{\mathbb{R}^N} y_i^2 w_{x_h}^2 + o(h^2) \\
 (3.4) \quad &= \frac{h^2}{2(V(x_0) + \lambda)^2} \Delta V(x_0) \int_{\mathbb{R}^N} y_i^2 w_{x_0}^2 + o(h^2) \quad \left(\text{because } p = 1 + \frac{4}{N} \right).
 \end{aligned}$$

If x_0 is a local minimum point, then $m = 0$ and $n(L_h) = 1$. Since the Hessian matrix $(\nabla^2 V(x_0))$ is positive definite, then

$$(3.5) \quad d''(\lambda) > 0, \quad p(d''(\lambda)) = 1,$$

which implies that u_h is orbitally unstable by the orbital stability criteria of [14] and [15].

If x_0 is not a local minimum, then $m \geq 1$ and $n(L_h) \geq 2$. In this case, by the formula (3.4), $p(d''(\lambda)) = \frac{1}{2}(1 + \frac{\Delta V(x_0)}{|\Delta V(x_0)|})$. By the instability criteria of [15], we conclude that u_h is orbitally unstable if $m - \frac{1}{2}(1 + \frac{\Delta V(x_0)}{|\Delta V(x_0)|})$ is even. This completes the proof of Theorem 1.1.

4. Proof of Theorem 1.2. Let $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$ be a solution of (1.11). The linearized operator of (1.10) around $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$ is

$$(4.1) \quad \mathbb{L}_h \begin{pmatrix} \phi \\ \psi \end{pmatrix} = \begin{pmatrix} h^2 \Delta \phi - (V(x) + \lambda) \phi + \frac{3+\beta}{1+\beta} u_h^2 \phi + \frac{2\beta}{1+\beta} u_h^2 \psi \\ h^2 \Delta \psi - (V(x) + \lambda) \psi + \frac{3+\beta}{1+\beta} u_h^2 \psi + \frac{2\beta}{1+\beta} u_h^2 \phi \end{pmatrix}.$$

We first define a sequence of numbers $\beta_j \in (-1, 0)$. By Lemma 4.2 of [34], the eigenvalue problem

$$(4.2) \quad \Delta \psi - (V(x_0) + \lambda) \psi + \mu w_{x_0}^2 \psi = 0$$

admits eigenvalues

$$(4.3) \quad \mu_1 = 1, \mu_2 = \dots = \mu_{N+1} = 3, \mu_{N+2} > 3.$$

We then define β_j by

$$(4.4) \quad \beta_j = \frac{3 - \mu_j}{1 + \mu_j}, \quad j = 1, 2, \dots$$

The following lemma shows the nondegeneracy of \mathbb{L} .

LEMMA 4.1. \mathbb{L}_h has no zero eigenvalue if $\beta \neq \beta_j, j = 1, 2, \dots$

Proof. Let $\mathbb{L}_h \begin{pmatrix} \phi \\ \psi \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$. Then by an orthonormal transformation it is equivalent to

$$(4.5) \quad L_{h,1} \tilde{\phi} = 0,$$

$$(4.6) \quad L_{h,2} \tilde{\psi} = 0,$$

where $L_{h,1} = L_h, L_{h,2} = h^2 \Delta - (V(x_h + hy) + \lambda) + \frac{3-\beta}{1+\beta} u_h^2$. By Theorem 2.5, we may conclude that $\tilde{\phi} = 0$. It remains to consider (4.6). As $h \rightarrow 0$, (4.6) may tend to the limiting equation given by

$$(4.7) \quad \Delta \psi - (V(x_0) + \lambda) \psi + \frac{3 - \beta}{1 + \beta} w_{x_0}^2 \psi = 0.$$

Since $\beta \neq \beta_j$, i.e., $\frac{3-\beta}{1+\beta} \neq \mu_j$, then by Lemma 4.2 of [34], (4.7) has only a trivial solution. Therefore, $\psi = 0$ and we may complete the proof. \square

The next lemma computes the Morse index of $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$. Here the Morse index is defined to be the number of positive eigenvalues of \mathbb{L}_h , which is just $n(\mathbb{L}_h)$.

LEMMA 4.2. If $-1 < \beta < 0$ and $\beta \notin \{\beta_2, \dots, \beta_j, \dots\}$, then the Morse index of $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$ is at least $N + 2$. If $0 < \beta < 1$, then the Morse index of $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$ is two. If $\beta > 1$, then the Morse index of $(\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} u_h)$ is one.

Proof. The eigenvalue problem $\mathbb{L}_h \begin{pmatrix} \phi \\ \psi \end{pmatrix} = \bar{\lambda} \begin{pmatrix} \phi \\ \psi \end{pmatrix}$ can be decomposed to

$$(4.8) \quad L_{h,1} \tilde{\phi} = \bar{\lambda} \tilde{\phi},$$

$$(4.9) \quad L_{h,2} \tilde{\psi} = \bar{\lambda} \tilde{\psi}.$$

By Theorem 2.5, $L_{h,1}$ has only one positive eigenvalue. It remains to consider the spectrum of $L_{h,2}$. If $\beta < 0$, then the eigenvalue problem

$$(4.10) \quad \Delta \psi - (V(x_0) + \lambda) \psi + \frac{3 - \beta}{1 + \beta} w_{x_0}^2 \psi = \bar{\lambda} \psi$$

has at least $N + 1$ positive eigenvalues. We may define a space of functions by

$$\mathbf{V} = \text{span} \left\{ w_{x_0}, \frac{\partial w_{x_0}}{\partial y_j}, j = 1, \dots, N, \right\}.$$

Since $-1 < \beta < 0$, we have $\frac{3-\beta}{1+\beta} > 3$. Hence

$$(4.11) \quad \int_{\mathbb{R}^N} \left[|\nabla \phi|^2 + (V(x_0) + \lambda) \phi^2 - \frac{3 - \beta}{1 + \beta} w_{x_0}^2 \phi^2 \right] < 0 \quad \forall \phi \in \mathbf{V}.$$

Thus by the variational characterization of the eigenvalues of (4.10), we see that $\lambda_{N+1} > 0$. Moreover, by the perturbation argument, (4.9) has at least $N + 1$ positive eigenvalues.

So when $\beta < 0$, \mathbb{L}_h has at least $N + 2$ positive eigenvalues.

When $0 < \beta < 1$, $1 < \frac{3-\beta}{1+\beta} < 3$, (4.10) has only one positive eigenvalue. So the Morse index is two.

When $\beta > 1$, (4.10) has no positive eigenvalue. So the Morse index is one. \square

Since \mathbb{L}_h is invertible, $(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h)$ is nondegenerate. Thus the system

$$(4.12) \quad \begin{cases} h^2 \Delta u - (V(x) + \lambda + \lambda_1)u + u^3 + \beta uv^2 = 0 & \text{in } \mathbb{R}^N, \\ h^2 \Delta v - (V(x) + \lambda + \lambda_2)v + v^3 + \beta u^2 v = 0 & \text{in } \mathbb{R}^N \end{cases}$$

has a solution $(u_{h,\lambda_1,\lambda_2}, v_{h,\lambda_1,\lambda_2})$ satisfying

$$(4.13) \quad (u_{h,\lambda_1,\lambda_2}, v_{h,\lambda_1,\lambda_2}) = \left(\frac{1}{\sqrt{1+\beta}}u_h, \frac{1}{\sqrt{1+\beta}}u_h \right) + O((|\lambda_1| + |\lambda_2|)h^{-2})$$

as $|\lambda_1| + |\lambda_2| \ll 1$.

Let us define

$$(4.14) \quad \begin{aligned} d(\lambda_1, \lambda_2) &= \int_{\mathbb{R}^N} \frac{h^2}{2} |\nabla u_{h,\lambda_1,\lambda_2}|^2 + \frac{V(x) + \lambda + \lambda_1}{2} u_{h,\lambda_1,\lambda_2}^2 - \frac{1}{4} \int_{\mathbb{R}^N} u_{h,\lambda_1,\lambda_2}^4 \\ &\quad + \int_{\mathbb{R}^N} \frac{h^2}{2} |\nabla v_{h,\lambda_1,\lambda_2}|^2 + \frac{V(x) + \lambda + \lambda_2}{2} v_{h,\lambda_1,\lambda_2}^2 - \frac{1}{4} \int_{\mathbb{R}^N} v_{h,\lambda_1,\lambda_2}^4 \\ &\quad - \frac{\beta}{2} \int_{\mathbb{R}^N} u_{h,\lambda_1,\lambda_2}^2 v_{h,\lambda_1,\lambda_2}^2. \end{aligned}$$

It is easy to see that

$$\begin{aligned} \frac{\partial d}{\partial \lambda_1} &= \frac{1}{2} \int_{\mathbb{R}^N} u_{h,\lambda_1,\lambda_2}^2, & \frac{\partial^2 d}{\partial \lambda_1^2} &= \int_{\mathbb{R}^N} u_{h,\lambda_1,\lambda_2} \frac{\partial u_{h,\lambda_1,\lambda_2}}{\partial \lambda_1}, \\ \frac{\partial d}{\partial \lambda_2} &= \frac{1}{2} \int_{\mathbb{R}^N} v_{h,\lambda_1,\lambda_2}^2, & \frac{\partial^2 d}{\partial \lambda_2^2} &= \int_{\mathbb{R}^N} v_{h,\lambda_1,\lambda_2} \frac{\partial v_{h,\lambda_1,\lambda_2}}{\partial \lambda_2}, \\ \frac{\partial^2 d}{\partial \lambda_1 \partial \lambda_2} &= \int_{\mathbb{R}^N} u_{h,\lambda_1,\lambda_2} \frac{\partial u_{h,\lambda_1,\lambda_2}}{\partial \lambda_2}. \end{aligned}$$

Now we may define functions as

$$\Phi_1 = \left. \frac{\partial u_{h,\lambda_1,\lambda_2}}{\partial \lambda_1} \right|_{(\lambda_1,\lambda_2)=(0,0)} \quad \text{and} \quad \Psi_1 = \left. \frac{\partial u_{h,\lambda_1,\lambda_2}}{\partial \lambda_2} \right|_{(\lambda_1,\lambda_2)=(0,0)}.$$

Then by (4.12), (Φ_1, Ψ_1) satisfies

$$(4.15) \quad \mathbb{L}_h \begin{pmatrix} \Phi_1 \\ \Psi_1 \end{pmatrix} = \begin{pmatrix} u_{h,0,0} \\ 0 \end{pmatrix}.$$

Similarly, if we set $\Phi_2 = \frac{\partial v_{h,\lambda_1,\lambda_2}}{\partial \lambda_1}|_{(\lambda_1,\lambda_2)=(0,0)}$ and $\Psi_2 = \frac{\partial v_{h,\lambda_1,\lambda_2}}{\partial \lambda_2}|_{(\lambda_1,\lambda_2)=(0,0)}$, then by (4.12), we have

$$(4.16) \quad \Psi_2 = \Phi_1, \quad \Phi_2 = \Psi_1.$$

Let $\mathbb{B} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$. Then (4.15) is equivalent to

$$(4.17) \quad \mathbb{B} \mathbb{L}_h \begin{pmatrix} \Phi_1 \\ \Psi_1 \end{pmatrix} = \begin{pmatrix} L_{h,1}(\Phi_1 + \Psi_1) \\ L_{h,2}(\Phi_1 - \Psi_1) \end{pmatrix} = \begin{pmatrix} u_{h,0,0} \\ u_{h,0,0} \end{pmatrix}.$$

So

$$(4.18) \quad \Phi_1 + \Psi_1 = R_{h,1} \quad \text{and} \quad \Phi_1 - \Psi_1 = R_{h,2},$$

where

$$(4.19) \quad R_{h,1} = \frac{1}{\sqrt{1+\beta}} R_h \quad \text{and} \quad R_{h,2} = L_{h,2}^{-1} \left(\frac{1}{\sqrt{1+\beta}} u_h \right).$$

Now we compute the Hessian matrix

$$\begin{aligned} (\nabla^2 d)|_{(\lambda_1,\lambda_2)=(0,0)} &= \begin{pmatrix} \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h \Phi_1 & \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h \Psi_1 \\ \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h \Psi_1 & \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h \Phi_1 \end{pmatrix}, \\ \mathbb{B} (\nabla^2 d)|_{(\lambda_1,\lambda_2)=(0,0)} \mathbb{B}^T &= \begin{pmatrix} \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,1} & \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,1} \\ \frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,2} & -\frac{1}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,2} \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \\ &= \begin{pmatrix} \frac{2}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,1} & 0 \\ 0 & \frac{2}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,2} \end{pmatrix}. \end{aligned}$$

By the results in section 3, $\frac{2}{\sqrt{1+\beta}} \int_{\mathbb{R}^N} u_h R_{h,1} = \frac{2}{1+\beta} \int_{\mathbb{R}^N} u_h R_h > 0$. It is enough to compute $\sqrt{1+\beta} \int_{\mathbb{R}^N} u_h R_{h,2} = \int_{\mathbb{R}^N} u_h L_{h,2}^{-1}(u_h)$. Note that as $h \rightarrow 0^+$,

$$\int_{\mathbb{R}^N} u_h L_{h,2}^{-1}(u_h) \rightarrow \int_{\mathbb{R}^N} w_{x_0} L_{\mu}^{-1}(w_{x_0}),$$

where

$$(4.20) \quad L_{\mu} \phi = \Delta \phi - (V(x_0) + \lambda) \phi + \mu w_{x_0}^2 \phi$$

with $\mu = \frac{3-\beta}{1+\beta}$.

Let $\rho(\mu) = \int_{\mathbb{R}^N} w_{x_0} L_{\mu}^{-1}(w_{x_0})$ and ϕ_{μ} the unique solution of $\Delta \phi_{\mu} - (V(x_0) + \lambda) \phi_{\mu} + \mu w_{x_0}^2 \phi_{\mu} = w_{x_0}$, i.e., $L_{\mu} \phi_{\mu} = w_{x_0}$ for $\mu \neq \mu_j, j = 1, 2, \dots$. Then $\frac{\partial \phi_{\mu}}{\partial \mu}$ satisfies

$$L_{\mu} \left(\frac{\partial \phi_{\mu}}{\partial \mu} \right) = -w_{x_0}^2 \phi_{\mu}, \quad \text{i.e.,} \quad \frac{\partial \phi_{\mu}}{\partial \mu} = -L_{\mu}^{-1}(w_{x_0}^2 \phi_{\mu}).$$

Hence

$$\begin{aligned} \rho'(\mu) &= \int_{\mathbb{R}^N} w_{x_0} \frac{\partial \phi_{\mu}}{\partial \mu} = - \int_{\mathbb{R}^N} w_{x_0} L_{\mu}^{-1}(w_{x_0}^2 \phi_{\mu}) \\ &= - \int_{\mathbb{R}^N} (L_{\mu}^{-1} w_{x_0})(w_{x_0}^2 \phi_{\mu}) \\ &= - \int_{\mathbb{R}^N} w_{x_0}^2 \phi_{\mu}^2 < 0 \quad \text{for} \quad \mu \neq \mu_j, j = 1, 2, \dots, \end{aligned}$$

i.e.,

$$(4.21) \quad \rho'(\mu) < 0 \quad \text{for } \mu \neq \mu_j, j = 1, 2, \dots$$

Due to (4.3), ρ is smooth on $(-\infty, 1) \cup (1, 3) \cup (3, \infty) \setminus \{\mu_j : j = N + 2, N + 3, \dots\}$. On the other hand, as $\mu \rightarrow 3$,

$$(4.22) \quad \phi_\mu \rightarrow L_0^{-1} w_{x_0} = \frac{1}{2} w_{x_0} + \frac{1}{2} y \cdot \nabla w_{x_0},$$

$$(4.23) \quad \rho(\mu) \rightarrow \int_{\mathbb{R}^2} w_{x_0} \left(\frac{1}{2} w_{x_0} + \frac{1}{2} y \cdot \nabla w_{x_0} \right) = 0.$$

Here we have used the fact that $N = 2$ and $p = 3$. Thus for $1 < \mu < 3$, $\rho(\mu) > 0$. This implies that for $0 < \beta < 1$, $\int_{\mathbb{R}^2} u_h R_{h,2} > 0$ and thus $(\nabla^2 d(0, 0))$ has *two* positive eigenvalues.

Now we consider $\mu \in (-\infty, 1)$, i.e., $\beta > 1$. By the standard maximal principle, $\phi_\mu < 0$ in \mathbb{R}^2 for $\mu < 0$. Consequently, $\rho(\mu) < 0$ for $\mu < 0$. Hence by (4.21), $\rho(\mu) < 0$ for $\mu \in (-\infty, 1)$, i.e., $\beta > 1$. This implies that $\int_{\mathbb{R}^2} u_h R_{h,2} < 0$ and thus $(\nabla^2 d(0, 0))$ has only *one* positive eigenvalue.

In conclusion, we see that the matrix $(\nabla^2 d(0, 0))$ has two positive eigenvalues when $0 < \beta < 1$ and one positive eigenvalue when $\beta > 1$. That is $p = 2$ when $0 < \beta < 1$ and $p = 1$ when $\beta > 1$. By Lemmas 4.1 and 4.2, we also deduce that $n(\mathbb{L}_h) = 2$ when $0 < \beta < 1$ and $n(\mathbb{L}_h) = 1$ when $\beta > 1$. Hence for $\beta > 0, \beta \neq 1$, we have $n(\mathbb{L}_h) = p$. Therefore, we conclude that $(u_{h,0,0}, v_{h,0,0}) = (\frac{1}{\sqrt{1+\beta}} u_h, \frac{1}{\sqrt{1+\beta}} v_h)$ is orbitally stable if $0 < \beta, \beta \neq 1$. This completes the proof of Theorem 1.2.

Acknowledgment. We thank the referees for bringing our attention to [12] and [13] and for many helpful suggestions.

REFERENCES

- [1] A. AMBROSETTI, M. BADIÀLE, AND S. CINGOLANI, *Semiclassical states of nonlinear Schrödinger equations*, Arch. Ration. Mech. Anal., 140 (1997), pp. 285–300.
- [2] A. AMBROSETTI AND E. COLORADO, *Bound and ground states of coupled nonlinear Schrödinger equations*, C. R. Math. Acad. Sci. Paris, 342 (2006), pp. 453–458.
- [3] A. AMBROSETTI, A. MALCHIODI, AND W.-M. NI, *Singularly perturbed elliptic equations with symmetry: Existence of solutions concentrating on spheres, Part I*, Comm. Math. Phys., 235 (2003), pp. 427–466.
- [4] T. BARTSCH, Z.-Q. WANG, AND J. WEI, *Bound states for a coupled Schrödinger system*, J. Fixed Point Theory Appl., 2 (2007), pp. 353–367.
- [5] H. BERESTYCKI AND T. CAZENAVE, *Instabilité des états stationnaires dans les équations de Schrödinger et de Klein-Gordon non linéaires*, C. R. Acad. Sci. Paris Sér. I Math., 293 (1981), pp. 489–492.
- [6] T. CAZENAVE AND P. L. LIONS, *Orbital stability of standing waves for some nonlinear Schrödinger equations*, Comm. Math. Phys., 85 (1982), pp. 549–561.
- [7] S. L. CORNISH, S. T. THOMPSON, AND C. E. WIEMAN, *Formation of bright matter-wave solitons during the collapse of Bose–Einstein condensates*, Phys. Rev. Lett., 96 (2006), article 170401.
- [8] M. DEL PINO AND P. FELMER, *Semiclassical states for nonlinear Schrödinger equation*, J. Funct. Anal., 149 (1997), pp. 245–265.
- [9] M. DEL PINO AND P. FELMER, *Semiclassical states of nonlinear Schrödinger equations: A variational reduction method*, Math. Ann., 324 (2002), pp. 1–32.
- [10] E. A. DONLEY, N. R. CLAUSSEN, S. L. CORNISH, J. L. ROBERTS, E. A. CORNELL, AND C. E. WIEMAN, *Dynamics of collapsing and exploding Bose–Einstein condensates*, Nature, 19 (2001), pp. 295–299.

- [11] A. FLOER AND A. WEINSTEIN, *Nonspreading wave packets for the cubic Schrödinger equation with a bounded potential*, J. Funct. Anal., 69 (1986), pp. 397–408.
- [12] R. FUKUIZUMI, *Stability of standing waves for nonlinear Schrödinger equations with critical power nonlinearity and potentials*, Adv. Differential Equations, 10 (2005), pp. 259–276.
- [13] R. FUKUIZUMI AND L. JEANJEAN, *Stability of standing waves for a nonlinear Schrödinger equation with a repulsive Dirac delta potential*, Discrete Contin. Dyn. Syst. to appear.
- [14] M. GRILLAKIS, J. SHATAH, AND W. STRAUSS, *Stability theory of solitary waves in the presence of symmetry I*, J. Funct. Anal., 74 (1987), pp. 160–197.
- [15] M. GRILLAKIS, J. SHATAH, AND W. STRAUSS, *Stability theory of solitary waves in the presence of symmetry II*, J. Funct. Anal., 94 (1990), pp. 308–348.
- [16] M. GROSSI, *On the number of single-peaked solutions of the nonlinear Schrödinger equations*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 19 (2002), pp. 261–280.
- [17] C. GUI, *Existence of multi-bump solutions for nonlinear Schrödinger equations via variational method*, Comm. Partial Differential Equations, 21 (1996), pp. 787–820.
- [18] T. L. HORNG, S. C. GOU, AND T. C. LIN, *Bending-wave instability of a vortex ring in a trapped Bose–Einstein condensate*, Phys. Rev. A, 74 (2006), article 041603.
- [19] X. KANG AND J. WEI, *On interacting bumps of semiclassical states of nonlinear Schrödinger equations*, Adv. Differential Equations, 5 (2000), pp. 899–928.
- [20] T.-C. LIN AND J.-C. WEI, *Ground state of N coupled Nonlinear Schrödinger Equations in \mathbb{R}^n , $n \leq 3$* , Comm. Math. Phys., 255 (2005), pp. 629–653.
- [21] T. C. LIN AND J. WEI, *Spikes in two-component systems of nonlinear Schrödinger equations with trapping potentials*, J. Differential Equations, 229 (2006), pp. 538–569.
- [22] T. C. LIN AND J. WEI, *Symbiotic bright solitary wave solutions of coupled nonlinear Schrödinger equations*, Nonlinearity, 19 (2006), pp. 2755–2773.
- [23] L. A. MAIA, E. MONTEFUSCO, AND B. PELLACCI, *Positive solutions for a weakly coupled nonlinear Schrödinger system*, J. Differential Equations, 229 (2006), pp. 743–767.
- [24] D. MIHALACHE, D. MAZILU, F. LEDERER, B. A. MALOMED, L.-C. CRASOVAN, Y. V. KARTASHOV, AND L. TORNER, *Stable three-dimensional solitons in attractive Bose–Einstein condensates loaded in an optical lattice*, Phys. Rev. A, 72 (2005), article 021601.
- [25] Y.-G. OH, *On positive multi-bump states of nonlinear Schrödinger equation under multiple well potentials*, Comm. Math. Phys., 131 (1990), pp. 223–253.
- [26] Y.-G. OH, *Stability of semiclassical bound state of nonlinear Schrödinger equations with potentials*, Comm. Math. Phys., 121 (1989), pp. 11–33.
- [27] V. M. PEREZ-GARCIA AND J. BELMONTE BEITIA, *Symbiotic solitons in heteronuclear multi-component Bose–Einstein condensates*, Phys. Rev. A, 72 (2005), article 033620.
- [28] L. PITAEVSKII AND S. STRINGARI, *Bose–Einstein Condensation*, Oxford University Press, Oxford, 2003.
- [29] P. RABINOWITZ, *On a class of nonlinear Schrödinger equations*, Z. Angew. Math. Phys., 43 (1992), pp. 270–291.
- [30] B. SIRAKOV, *Standing wave solutions of the nonlinear Schrödinger equation in \mathbb{R}^N* , Ann. Mat. Pura Appl. (4), 4 (2002), pp. 73–83.
- [31] B. SIRAKOV, *Least energy solitary waves for a system of nonlinear Schrödinger equations*, Comm. Math. Phys., to appear.
- [32] X. WANG, *On concentration of positive bound states of nonlinear Schrödinger equations*, Comm. Math. Phys., 153 (1993), pp. 229–243.
- [33] Z.-Q. WANG, *Existence and symmetry of multi-bump solutions for nonlinear Schrödinger equations*, J. Differential Equations, 159 (1999), pp. 102–137.
- [34] J. WEI, *On the construction of single-peaked solutions to a singularly perturbed elliptic Dirichlet problem*, J. Differential Equations, 129 (1996), pp. 315–333.
- [35] J. WEI, *On the interior spike layer solutions for some singular perturbation problems*, Proc. Roy. Soc. Edinburgh Sect. A, 128 (1998), pp. 849–874.
- [36] M. WEINSTEIN, *Nonlinear Schrödinger equations and sharp interpolation estimates*, Comm. Math. Phys., 87 (1983), pp. 567–576.
- [37] M. I. WEINSTEIN, *The Connection between Finite and Infinite-Dimensional Dynamical Systems*, AMS, Providence, RI, 1989.

FORCED VIBRATIONS OF A NONHOMOGENEOUS STRING*

PIETRO BALDI[†] AND MASSIMILIANO BERTI[‡]

Abstract. We prove existence of vibrations of a nonhomogeneous string under a nonlinear time periodic forcing term in the case in which the forcing frequency avoids resonances with the vibration modes of the string (nonresonant case). The proof relies on a Lyapunov–Schmidt reduction and a Nash–Moser iteration scheme.

Key words. nonlinear wave equations, periodic solutions, small divisors, Nash–Moser iteration scheme

AMS subject classifications. 35B10, 35L70, 58C15

DOI. 10.1137/060665038

1. Introduction. In this paper we study forced vibrations of a nonhomogeneous string,

$$(1) \quad \begin{cases} \rho(x)u_{tt} - (p(x)u_x)_x = \mu f(x, \omega t, u), \\ u(0, t) = u(\pi, t) = 0, \end{cases}$$

where $\rho(x) > 0$ is the mass per unit length, $p(x) > 0$ is the modulus of elasticity multiplied by the cross-sectional area (see [15, p. 291]), $\mu > 0$ is a parameter, and the nonlinear forcing term $f(x, \omega t, u)$ is $(2\pi/\omega)$ -periodic in time (i.e., $f(x, \cdot, u)$ is 2π -periodic).

Equation (1) is a nonlinear model also for propagation of waves in nonisotropic media describing seismic phenomena; see, e.g., [2].

We look for $(2\pi/\omega)$ -time periodic solutions $u(x, t)$ of (1).

This problem has received wide attention since the pioneering paper of Rabinowitz [26] dealing with the weakly nonlinear homogeneous string with $\rho(x) = p(x) = 1$, μ small, and $\omega = 1$. In this case the forcing frequency ω enters in resonance with the proper eigenfrequencies $\omega_j = j \in \mathbb{N}$ of the string.

For functions 2π -periodic in time and satisfying spatial Dirichlet boundary conditions, the spectrum $\{-l^2 + j^2, l \in \mathbb{Z}, j \geq 1\}$ of the D’Alembertian operator $\partial_{tt} - \partial_{xx}$ possesses the zero eigenvalue with infinite multiplicity (for $|l| = j$), but the other eigenvalues are well separated. The corresponding infinite-dimensional bifurcation problem is solved in [26] for nonlinearities f which are monotone in u ; see [7] for nonmonotone f .

Subsequently many other results, both of bifurcation and of a global nature ($\mu = 1$), have been obtained, still for rational forcing frequencies $\omega \in \mathbb{Q}$, relying on the separation properties of the spectrum; see, e.g., [27, 28, 14, 31, 4].

When the forcing frequency $\omega \in \mathbb{R} \setminus \mathbb{Q}$ is irrational (nonresonant case) the situation is completely different. Indeed the wave operator $\omega^2 \partial_{tt} - \partial_{xx}$ does not possess the

*Received by the editors July 13, 2006; accepted for publication (in revised form) September 14, 2007; published electronically April 25, 2008. This work was supported by MURST under the national project “Variational Methods and Nonlinear Differential Equations.”

<http://www.siam.org/journals/sima/40-1/66503.html>

[†]SISSA, via Beirut 2-4, 34014 Trieste, Italy (baldi@sissa.it).

[‡]Dipartimento di Matematica e Applicazioni “R. Caccioppoli,” Università di Napoli “Federico II,” via Cintia, 80126 Napoli, Italy (m.berti@unina.it).

zero eigenvalue, but its spectrum $\{-\omega^2 l^2 + j^2, l \in \mathbb{Z}, j \geq 1\}$ accumulates to zero for almost every ω . This is a “small divisors problem.”

We underline that this “small divisors” phenomenon arises naturally for more realistic model equations like (1) where the density $\rho(x)$ and the modulus of elasticity $p(x)$ are not constant. Indeed in this case the eigenfrequencies ω_j of the string are no longer integer numbers, having the asymptotic expansion

$$(2) \quad \omega_j^2 = \frac{j^2}{c^2} + b + O\left(\frac{1}{j}\right)$$

with suitable constants c, b depending on ρ, p ; see (66).

If $\omega = m/n \in \mathbb{Q}$, good separation properties of the spectrum can be recovered when $p(x) = \rho(x)$ (so $c = 1$) and assuming the extra condition $b \neq 0$; see [3, 29]. Indeed in this case the linear spectrum

$$-\omega^2 l^2 + \omega_j^2 = -\omega^2 l^2 + j^2 + b + O\left(\frac{1}{j}\right)$$

possesses at most finitely many zero eigenvalues and the remaining part of the spectrum is far away from zero. On the other hand, if $b = 0$, the eigenvalues with $(l, j) \in (n, m)\mathbb{N}$ tend to zero (also in the case $\omega \in \mathbb{Q}$!).

Existence of weak solutions in the nonresonant case was proved by Acquistapace [1] for $\rho = 1, \mu$ small, and for a zero measure set of forcing frequencies ω for which the eigenvalues $-\omega^2 l^2 + \omega_j^2$ are far away from zero. These frequencies are essentially the numbers whose continued fraction expansion is bounded; see [30].

For a similar zero measure set of frequencies, McKenna [23] has obtained some result when $\mu = 1, \rho = p = 1$, and $f(x, t, u) = g(u) + h(x, t)$ with g uniformly Lipschitz, via a fixed point argument; see also [5]; for related results using variational methods see [18, 10].

Existence of classical solutions of (1) for a positive measure set of frequencies was proved by Plotnikov and Yungerman [24] for the homogeneous string $\rho = p = 1, \mu$ small, and f monotone in u . This monotonicity condition allows one to control the constant coefficient in the asymptotic expansion of the eigenvalues (like b in (2)) of some perturbed linearized operator.

Recently Fokam [19] has proved existence of classical periodic solutions for large frequencies ω in a set of asymptotically full measure for the homogeneous string $\rho = p = 1$ plus a potential, $\mu = 1$ and $f = u^3 + h(x, t)$ with h a trigonometric polynomial odd in time and space.

In the present paper we prove existence of classical solutions of the nonhomogeneous string (1) for every $\rho(x), p(x) > 0$ for general nonlinear terms $f(x, \omega t, u)$, and for (μ, ω) belonging to a large measure Cantor set B_γ , when the ratio μ/ω is small. Our Theorem 1 covers both the case $\mu \rightarrow 0$ (weak forcing) and the case $\omega \rightarrow +\infty$ (rapid forcing).

In the limit $\mu/\omega \rightarrow 0$ the solution we find tends to a static equilibrium $v(x)$ with smaller, zero average oscillations $w(x, t)$ of amplitude $O(\mu/\omega)$; see (13), (14), and Figure 1. The nonlinearity f selects such v through the infinite-dimensional bifurcation equation (10), which possesses nondegenerate solutions under natural assumptions on f ; see Hypothesis (V). This problem is not present in [19], where, thanks to the symmetry assumptions on f , there is no bifurcation equation.

Considering the structure of the expected solution, it is natural to attack the problem via a Lyapunov–Schmidt decomposition.

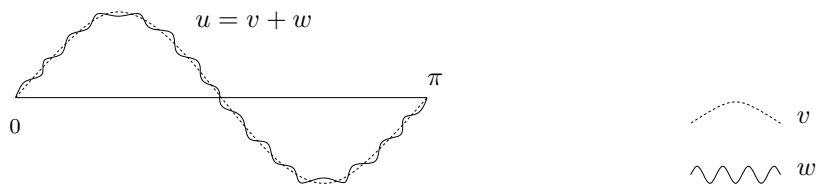


FIG. 1. The solution $u(x,t) = v(x) + w(x,t)$ of (1).

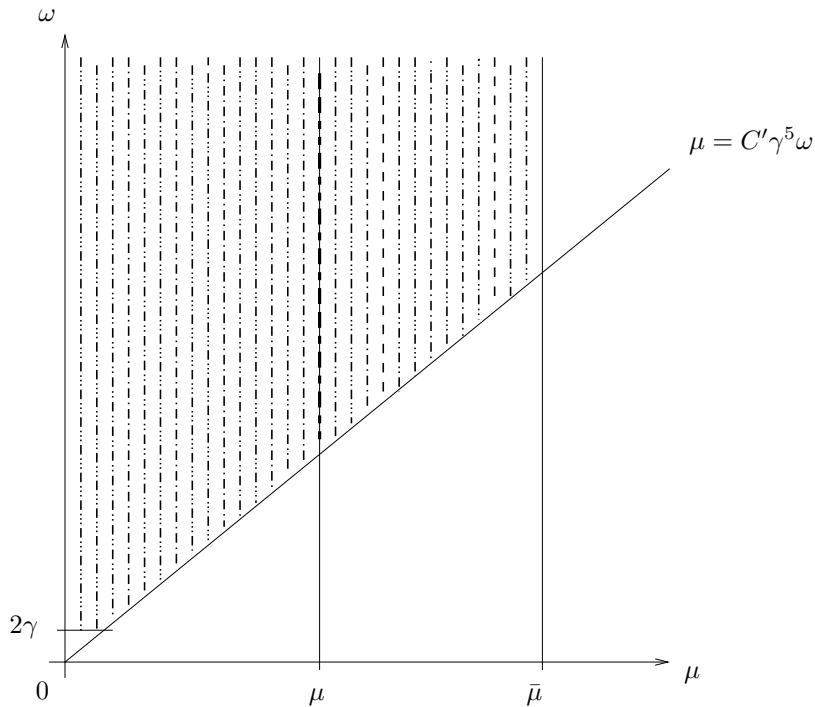


FIG. 2. The Cantor set B_γ .

In the range equation (to find w) a small divisors problem arises, and we solve it with a Nash–Moser-type iterative scheme. The inversion of the “linearized operators”—which is the core of any Nash–Moser scheme—is obtained adapting the techniques of [8] to the present time-dependent case (section 6). This method is also reminiscent of the approach of Kuksin (unpublished) explained by Bourgain in [13, pp. 90–94]. See also the works of Craig and Wayne [16, 17] and Bourgain [11, 12] for related techniques.

It is in the solution of the range equation where the interaction between the forcing frequency ω and the normal modes of oscillations of the string linearized at different positions (approximating better and better the final string configuration) appears.

The set B_γ of “nonresonant” parameters (μ, ω) for which we find a solution of the range equation (and then of (1)) is constructed avoiding these primary resonances. In particular the forcing frequency ω must not enter in resonance with the normal frequencies of oscillations of the string linearized at the limit solution; see (15). At the end of the construction we obtain a large measure Cantor set B_γ which looks like Figure 2. Outside this set the effect of resonance phenomena shall in general destroy the existence of periodic solutions like those found in Theorem 1.

Finally we recall that related existence results of periodic and quasi-periodic solutions for autonomous Hamiltonian PDEs have been obtained via KAM-type techniques since the pioneering works of Kuksin [21] and Wayne [32]; see also [22] and the references therein.

We now present rigorously our results.

1.1. Main result. After a time rescaling we look for 2π -periodic solutions of

$$(3) \quad \begin{cases} \omega^2 \rho(x) u_{tt} - (p(x) u_x)_x = \mu f(x, t, u), \\ u(0, t) = u(\pi, t) = 0, \end{cases}$$

where $\mu \in [0, \bar{\mu}]$ for some given $\bar{\mu} > 0$, under the 2π -periodic forcing term

$$(4) \quad f(x, t, u) = \sum_{l \in \mathbb{Z}} f_l(x, u) e^{ilt} = f_0(x, u) + \bar{f}(x, t, u),$$

where

$$\bar{f}(x, t, u) := \sum_{l \neq 0} f_l(x, u) e^{ilt}.$$

We suppose that f is analytic in (t, u) ; more precisely,

$$f(x, t, u) = \sum_{l \in \mathbb{Z}, k \in \mathbb{N}} f_{lk}(x) u^k e^{ilt},$$

where $f_{lk}(x) \in H^1((0, \pi); \mathbb{C})$, $f_{-l,k} = f_{lk}^*$ (the symbol z^* denotes the complex conjugate of $z \in \mathbb{C}$), and we assume the following hypothesis on the decay of the coefficients $\|f_{lk}\|_{H^1}$.

HYPOTHESIS (F). *There exist $2\sigma_0 > 0, r > 0$ such that*

$$\sum_{l \in \mathbb{Z}} \|f_{lk}\|_{H^1}^2 (1 + l^2) e^{(2\sigma_0)2|l|} := C_k^2(f) < \infty \quad \text{and} \quad \sum_{k=0}^{+\infty} C_k(f) r^k < \infty.$$

For example, trigonometric polynomials in t and polynomials in u , namely,

$$(5) \quad f(x, t, u) = \sum_{|l| \leq L, 0 \leq k \leq K} f_{lk}(x) u^k e^{ilt}$$

for some $L, K \in \mathbb{N}$, satisfy Hypothesis (F) for every σ_0, r .

Remark 1. We notice that if $f(x, t, 0) = \sum_{l \in \mathbb{Z}} f_{l0}(x) e^{ilt} \neq 0$, then (3) does not possess the trivial solution $u = 0$.

We look for periodic solutions of (3) in the Hilbert space

$$X_{\sigma,s} := \left\{ u : \mathbb{T} \rightarrow H_0^1((0, \pi); \mathbb{R}), u(x, t) = \sum_{l \in \mathbb{Z}} u_l(x) e^{ilt}, \quad u_l \in H_0^1((0, \pi); \mathbb{C}), \right.$$

$$\left. u_{-l} = u_l^*, \quad \|u\|_{\sigma,s}^2 := \sum_{l \in \mathbb{Z}} \|u_l\|_{H^1}^2 (1 + l^{2s}) e^{2\sigma|l|} < \infty \right\}$$

of 2π -periodic-in-time functions valued in $H^1((0, \pi); \mathbb{R})$ which have a bounded analytic extension on the complex strip $|\text{Im } t| < \sigma$ with trace function on $|\text{Im } t| = \sigma$ belonging to $H^s(\mathbb{T}; H^1((0, \pi); \mathbb{C}))$.

For $s > 1/2$, $X_{\sigma,s}$ is a multiplicative Banach algebra:

$$(6) \quad \|uv\|_{\sigma,s} \leq c_s \|u\|_{\sigma,s} \|v\|_{\sigma,s} \quad \forall u, v \in X_{\sigma,s}$$

with

$$(7) \quad c_s := 2^s \left(\sum_{n \in \mathbb{Z}} \frac{1}{1+n^{2s}} \right)^{1/2};$$

see, e.g., [6]. We shall use the notation X_σ , resp., $\|\cdot\|_\sigma$, for $X_{\sigma,1}$, resp., $\|\cdot\|_{\sigma,1}$.

1.2. The Lyapunov–Schmidt reduction. To find solutions of (3) we implement the Lyapunov–Schmidt reduction according to the decomposition

$$X_{\sigma,s} = V \oplus (W \cap X_{\sigma,s}),$$

where

$$V := H_0^1(0, \pi), \quad W := \left\{ w = \sum_{l \neq 0} w_l(x) e^{ilt} \in X_{0,s} \right\},$$

writing every $u \in X_{\sigma,s}$ as $u = u_0(x) + \sum_{l \neq 0} u_l(x) e^{ilt}$.

Projecting (3) with

$$u = v + w, \quad v \in V, w \in W,$$

yields

$$(8) \quad \begin{cases} -(pv')' = \mu \Pi_V f(v+w) & \text{(bifurcation equation),} \\ L_\omega w = \mu \Pi_W f(v+w) & \text{(range equation),} \end{cases}$$

where Π_V, Π_W denote the projectors, $f(u)(x, t) := f(x, t, u(x, t))$, and

$$L_\omega u := \omega^2 \rho(x) u_{tt} - (p(x) u_x)_x.$$

We shall find solutions of (8) when μ/ω is small. In this limit w tends to 0 and the bifurcation equation reduces to the time-independent equation

$$(9) \quad -(pv')' = \mu f_0(v)$$

because, by (4), for $w = 0$

$$\Pi_V f(v) = \Pi_V f_0(x, v(x)) + \Pi_V \bar{f}(x, t, v(x)) = f_0(v).$$

The infinite-dimensional “0th order bifurcation equation” (9) is a second order ODE, which, under natural conditions on f_0 , possesses nondegenerate solutions satisfying the boundary conditions $v(0) = v(\pi) = 0$.

HYPOTHESIS (V). *The problem*

$$(10) \quad \begin{cases} -(p(x)v'(x))' = \mu f_0(x, v(x)), \\ v(0) = v(\pi) = 0 \end{cases}$$

admits a solution $\bar{v} \in H_0^1(0, \pi)$ which is nondegenerate, namely, the linearized equation

$$-(ph')' = \mu f'_0(\bar{v})h$$

possesses in $H_0^1(0, \pi)$ only the trivial solution $h = 0$.

We note that, for $\mu = 0$, the trivial solution $\bar{v} = 0$ is nondegenerate, so, by the implicit function theorem, Hypothesis (V) is automatically satisfied for μ small. We deal also with the case μ not small; see, for example, Lemmas 2 and 3.

By the implicit function theorem, Hypothesis (V) implies the existence of a smooth map

$$(\mu, w) \mapsto v(\mu, w) \in V$$

such that $v(\mu, w)$ solves the bifurcation equation in (8); see Lemma 4.

Remark 2. For a discussion about the difficulties caused by a degenerate solution, we refer to [9].

Let λ_j denote the eigenvalues of the Sturm–Liouville problem

$$(11) \quad \begin{cases} -(p(x)y'(x))' = \lambda\rho(x)y(x), \\ y(0) = y(\pi) = 0 \end{cases}$$

and $\omega_j := \sqrt{\lambda_j}$. These are the frequencies of the free vibrations of the string (note that all the eigenvalues λ_j are positive). Physically, it is the sequence of the fundamental tone ω_1 and all its overharmonics $\omega_2, \omega_3, \dots$ which compose the musical note of the string.

For $\gamma \in (0, 1)$ we define

$$(12) \quad A_\gamma := \left\{ (\mu, \omega) \in (\mu_1, \mu_2) \times (\gamma, +\infty) : \frac{\mu}{\omega} < C'\gamma^5, \quad |\omega l - \omega_j| > \frac{\gamma}{l^\tau} \right. \\ \left. \forall l = 1, \dots, N_0, \quad j \geq 1 \right\},$$

where ω_j are given by (11), and (μ_1, μ_2) , $N_0 \in \mathbb{N}$, $C' > 0$ shall be fixed in the next theorem.

THEOREM 1 (existence). *Suppose $p(x), \rho(x) > 0$ are of class $H^3(0, \pi)$, f satisfies Hypothesis (F), and Hypothesis (V) holds for some $\mu_0 \in [0, \bar{\mu}]$.*

Fix $\tau \in (1, 2)$, $\gamma \in (0, 1)$. There exist a neighborhood (μ_1, μ_2) of μ_0 , $N_0 \in \mathbb{N}$, positive constants C, C' (depending on $\rho, p, f, \bar{\mu}, \bar{v}, \tau$), a map

$$\tilde{w} \in C^\infty(A_\gamma, X_{\sigma_0/2} \cap W),$$

and a Cantor set $B_\gamma \subset A_\gamma$ of positive measure such that, for all $(\mu, \omega) \in B_\gamma$,

$$(13) \quad \tilde{u}(\mu, \omega) := v(\mu, \tilde{w}(\mu, \omega)) + \tilde{w}(\mu, \omega) \in V \oplus (W \cap X_{\sigma_0/2})$$

is a classical solution of (3) and satisfies

$$\tilde{u}(\cdot, t) \in H^3(0, \pi) \cap H_0^1(0, \pi) \quad \forall t \in \mathbb{R}.$$

The Cantor set B_γ is defined in (15) and satisfies the measure estimate (56).

Furthermore, for all $(\mu, \omega) \in A_\gamma$ the following estimates hold:

$$(14) \quad \|\tilde{w}(\mu, \omega)\|_{\sigma_0/2} \leq C \frac{\mu}{\gamma\omega}, \quad \|v(\mu, \tilde{w}(\mu, \omega)) - v(\mu, 0)\|_{H^1} \leq C \frac{\mu}{\gamma\omega},$$

and $\|v(\mu, 0) - \bar{v}\|_{H^1} \leq C|\mu - \mu_0|$.

The neighborhood (μ_1, μ_2) of μ_0 is fixed in Lemma 4, the integer N_0 is fixed in Lemma 9, and the constant C' is fixed in Lemma 13.

Estimate (14) shows how close the solution \tilde{u} is to the static configuration $v(\mu, 0)$; see Figure 1.

Remark 3. We underline that the function $\tilde{w}(\mu, \omega)$, as well as $\tilde{u}(\mu, \omega)$, is defined for all the values of the parameters $(\mu, \omega) \in A_\gamma$ and not only on the Cantor set B_γ ($\tilde{w}(\mu, \omega)$ is introduced in Lemma 11). What is true is that if $(\mu, \omega) \in B_\gamma$, then $\tilde{w}(\mu, \omega)$ solves the range equation; see Theorem 3. As a consequence, if $(\mu, \omega) \in B_\gamma$, then $\tilde{u}(\mu, \omega)$ solves (3).

The Cantor set B_γ is explicitly defined by

$$B_\gamma := \left\{ (\mu, \omega) \in (\mu_1, \mu_2) \times (2\gamma, +\infty) : |\omega l - \omega_j| > \frac{2\gamma}{l^\tau} \quad \forall l = 1, \dots, N_0, \quad j \geq 1, \right. \\ \left. (15) \quad \frac{\mu}{\omega} < C'\gamma^5, \quad \left| \omega l - \frac{j}{c} \right| > \frac{2\gamma}{l^\tau}, \quad |\omega l - \tilde{\omega}_j(\mu, \omega)| > \frac{2\gamma}{l^\tau} \quad \forall l, j \geq 1 \right\},$$

where

$$(16) \quad c := \frac{1}{\pi} \int_0^\pi \left(\frac{\rho(x)}{p(x)} \right)^{1/2} dx$$

and $\tilde{\lambda}_j(\mu, \omega) := \tilde{\omega}_j^2(\mu, \omega)$ denote the eigenvalues of the Sturm–Liouville problem

$$(17) \quad \begin{cases} -(py')' - \mu \Pi_V f'(\tilde{u}(\mu, \omega)) y = \lambda \rho y, \\ y(0) = y(\pi) = 0. \end{cases}$$

Note that B_γ is constructed by means of the function $\tilde{u}(\mu, \omega)$, which is defined for all $(\mu, \omega) \in A_\gamma$; see Remark 3.

Remark 4. If some $\tilde{\lambda}_j(\mu, \omega)$ is negative, then $\tilde{\omega}_j(\mu, \omega) = i\sqrt{|\tilde{\lambda}_j(\mu, \omega)|}$ is a purely imaginary complex number and the nonresonance conditions in (15) are trivially satisfied.

The Cantor set B_γ is large in a measure theoretical sense; see section 4.3. In particular, for all $\mu \in (\mu_1, \mu_2)$,

$$S(\mu) := \{ \omega : (\mu, \omega) \in \cup_{\gamma \in (0,1)} B_\gamma \}$$

has asymptotically full measure at $\omega \rightarrow +\infty$, i.e.,

$$(18) \quad \lim_{\omega \rightarrow +\infty} |S(\mu) \cap (\omega, \omega + 1)| = 1.$$

Analogously,

$$S(\omega) := \{ \mu : (\mu, \omega) \in \cup_{\gamma \in (0,1)} B_\gamma \}$$

satisfies, for all $\omega' > 0$ and for all $\gamma' \in (0, 1)$,

$$(19) \quad \lim_{\mu \rightarrow 0} \left| \left\{ \omega \in (\omega', \omega' + 1) : \frac{|S(\omega) \cap (0, \mu)|}{\mu} \geq 1 - \gamma' \right\} \right| = 1.$$

Finally we discuss the regularity of the solution $\tilde{u}(x, t)$ found in Theorem 1 with respect to x (by construction \tilde{u} is analytic with respect to t).

THEOREM 2 (regularity). *Assume the hypotheses of Theorem 1. In addition, suppose that, for some $m \geq 3$,*

$$(20) \quad \rho(x) \in H^m(0, \pi), \quad p(x) \in H^{m+1}(0, \pi), \quad f_{lk}(x) \in H^m(0, \pi) \quad \forall l, k$$

and, for some $r_m > 0$,

$$(21) \quad \sum_{l \in \mathbb{Z}, k \geq 0} \|f_{lk}\|_{H^m r_m^k} < +\infty.$$

If $\|\tilde{u}(\cdot, t)\|_{H^1 r_m^{-1}}$ is small enough, then

$$(22) \quad \tilde{u}(\cdot, t) \in H^{m+2}(0, \pi) \cap H_0^1(0, \pi).$$

This conclusion holds true, for example, when $f_0(x, 0) = d_u f_0(x, 0) = 0$ and $\mu/\gamma\omega$ is small enough.

Note that the regularity (22) requires no skewsymmetry assumptions on f and requires just a smallness condition for the H^1 norm of $\tilde{u}(\cdot, t)$.

Remark 5. If $f(x, t, u)$ is a trigonometric polynomial in t and a polynomial in u as in (5), then the series in (21) is a finite sum. Therefore the conclusion (22) is true without smallness conditions for \tilde{u} .

In particular, if $\rho(x), p(x), f_{lk}(x)$ are C^∞ (for example, $f = \cos x \cos t(1 + u^2)$), then the solution \tilde{u} is C^∞ also in the variable x (the above f does not satisfy the skewsymmetry assumption (24)).

The subtle problem to prove Theorem 2 is that, because of Dirichlet boundary conditions, the Sobolev regularity of a function with respect to x is *not* characterized by the rapid decaying properties of the Fourier coefficients (unless we assume skewsymmetry assumptions on the nonlinearity and restrict solutions to $u(x, t)$ odd in x ; see Remark 6). Theorem 2 is proved in section 5 via bootstrap arguments.

1.3. Outline of the proof. In section 2 we prove that, under assumption (F), the composition operator induced by the nonlinearity f on $X_{\sigma,s}$ is an analytic map.

In section 3, we find a solution $v(\mu, w)$ of the infinite-dimensional bifurcation equation in (8). Thanks to assumption (V) (which is verified on several examples in Lemmas 2 and 3), $v(\mu, w)$ is obtained in Lemma 4 by a standard implicit function theorem.

In section 4 we solve the range equation by means of an iterative Nash–Moser implicit function theorem. The final theorem, Theorem 3, is proved in several steps.

In section 4.1 we find inductively a sequence of approximate solutions $w_n(\mu, \omega)$ defined on smaller and smaller subsets A_n of the parameters (μ, ω) (see (39)). The reason for these “excisions” is to avoid resonance phenomena in order to prove the invertibility of the linearized operators obtained at each step of the iteration; see conditions (33)–(34) in Lemma 7.

In section 4.2 we extend these approximate solutions $w_n(\mu, \omega)$ to C^∞ -functions $\tilde{w}_n(\mu, \omega)$ defined for all the values of the parameters (μ, ω) and converging (superexponentially fast) to a C^∞ map \tilde{w} defined for all (μ, ω) ; see Lemma 11. It is in proving the regularity of w_n with respect to the parameters (μ, ω) that we find it convenient to define the approximate solutions w_n as exact solutions of (41) (with $k = n$); see Remark 9.

In Lemma 12, we prove that the Cantor set B_γ , defined in (15) by means of \tilde{w} , is contained in A_n (which depends on w_{n-1}) for each n . This is a consequence of the superexponentially fast convergence of \tilde{w}_n to \tilde{w} ; see (52).

In section 4.3 we prove that B_γ is a large set in a measure theoretical sense.

In all the previous steps we have to assume smallness conditions for μ/ω . The most restrictive one is $\mu/\omega < C'\gamma^5$ in the definition (12) of A_γ .

In section 5 we conclude the proof of the existence Theorem 1, and we prove the regularity Theorem 2.

In section 6 we study the key step for the inversion of the linearized operators. Lemma 7 is obtained by a variant of the techniques developed in [8]. In particular, the key estimate on the small divisors of Lemma 18 is reminiscent of the method of Kuksin explained in [13, pp. 90–94].

Notation. The symbols K, K_i, K'_i shall denote positive constants depending only on $\rho, p, f, \bar{\mu}, \bar{\nu}, \tau$.

2. Regularity of the composition operator. We first prove the analyticity of the composition operator

$$u(x, t) \mapsto f(x, t, u(x, t))$$

induced by f on $X_{\sigma,s}$.

By the Banach algebra property (6) of $X_{\sigma,s}$ the composition operator

$$u \mapsto u^k \quad \forall k \in \mathbb{N}$$

is an analytic map from $X_{\sigma,s}$ into itself. Thanks to the rapid decay of the coefficients $\|f_{lk}\|_{H^1}$ assumed in Hypothesis (F), this property holds true also for the composition operator $f(u)$.

LEMMA 1. *Let f satisfy assumption (F). For every $\sigma \in [0, \sigma_0], s > 1/2$, the composition operator f is analytic on the ball $\{u \in X_{\sigma,s} : \|u\|_{\sigma,s} < r/c_s\}$, where c_s is defined in (7).*

Proof. First note that

$$\sum_{l \in \mathbb{Z}} \|u_l\|_{\infty} \leq \sqrt{\frac{\pi}{2}} \sum_{l \in \mathbb{Z}} \|u_l\|_{H^1} \leq \sqrt{\frac{\pi}{2}} \left(\sum_{l \in \mathbb{Z}} \|u_l\|_{H^1}^2 (1 + l^{2s}) \right)^{1/2} \left(\sum_{l \in \mathbb{Z}} \frac{1}{1 + l^{2s}} \right)^{1/2}$$

so $\|u\|_{\infty} \leq c_s \|u\|_{\sigma,s}$ for all $u \in X_{\sigma,s}, \sigma \geq 0, s > 1/2$, and $f(x, t, u(x, t))$ is well-defined.

By definition of the norm $\|\cdot\|_{\sigma,s}$, there exists $C := C(\sigma_0, s) > 0$ such that for all $\sigma \in [0, \sigma_0]$ and for all $k \in \mathbb{N}$,

$$\left\| \sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right\|_{\sigma,s} \leq C \left\| \sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right\|_{2\sigma_0, 1}.$$

Next

$$\left\| \sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right\|_{2\sigma_0, 1}^2 = \sum_{l \in \mathbb{Z}} \|f_{lk}\|_{H^1}^2 (1 + l^2) e^{(2\sigma_0)2|l|} =: C_k^2(f) < +\infty$$

by assumption (F). Therefore $\sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \in X_{\sigma,s}$ and

$$(23) \quad \left\| \sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right\|_{\sigma,s} \leq C C_k(f).$$

Using the algebra property (6) of $X_{\sigma,s}$ and (23)

$$\begin{aligned} \|f(u)\|_{\sigma,s} &\leq \sum_{k=0}^{\infty} \left\| \left(\sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right) u^k \right\|_{\sigma,s} \\ &\leq \sum_{k=0}^{\infty} c_s \left\| \sum_{l \in \mathbb{Z}} f_{lk}(x) e^{ilt} \right\|_{\sigma,s} \|u^k\|_{\sigma,s} \\ &\leq C \sum_{k=0}^{\infty} C_k(f) (c_s \|u\|_{\sigma,s})^k < C \sum_{k=0}^{\infty} C_k(f) r^k < +\infty \end{aligned}$$

for $c_s \|u\|_{\sigma,s} < r$, by (F) again.

The analyticity of the composition operator f with respect to $\|\cdot\|_{\sigma,s}$ follows from the properties of the power series as explained in [25, Appendix A]. \square

We emphasize that the analyticity of f as a map in $X_{\sigma,s}$ is not an assumption but follows from (F).

Remark 6. If $f(x, t, u)$ admits an analytic extension, which is 2π -periodic in x and skewsymmetric, namely,

$$(24) \quad f(-x, t, -u) = -f(x, t, u),$$

then the Dirichlet problem on $[0, \pi]$ is equivalent to the 2π -periodic problem within the space of all functions odd in x . In this case also the spatial regularity is characterized by the decay properties of the Fourier coefficients. Therefore we could look for analytic solutions of (3),

$$u(x, t) = \sum_{l \in \mathbb{Z}} u_l(x) e^{ilt},$$

which are periodic and odd in x , more precisely with

$$u_l(x) \in Y := \left\{ y(x) = \sum_{j \geq 1} y_j \sin(jx) : \sum_{j \geq 1} |y_j|^2 j^{2b} e^{2aj} < +\infty \right\}$$

for some $a \geq 0, b > 1/2$. Without the oddness assumption (24) the composition operator f does not map this subspace into itself. It is for this reason that we consider the space $X_{\sigma,s}$ of functions valued in $H_0^1(0, \pi)$: also without (24), f sends $X_{\sigma,s}$ into itself (Lemma 1).

Throughout this paper we shall use spaces $X_{\sigma,s}$ with $\sigma \in [\sigma_0/2, \sigma_0]$ and $s \in \mathcal{S} := \{1, 1 - \frac{\tau-1}{2}, 1 + \frac{(\tau-1)\tau}{2-\tau}\}$. So we choose $\bar{c} := \max_{s \in \mathcal{S}} c_s$ as a multiplicative algebra constant for all spaces $X_{\sigma,s}$.

By Lemma 1, f is analytic in the ball

$$\left\{ u \in X_{\sigma,s} : \|u\|_{\sigma,s} < R_0 := \frac{r}{\bar{c}} \right\}$$

and f, f', f'', \dots are bounded, uniformly in σ, s .

3. The bifurcation equation. Now we give some examples in which Hypothesis (V) holds.

LEMMA 2. *Suppose $f_0(x, u) = u^m$ for $m \geq 3$ odd and $p(x) \equiv 1$. Then, for all μ , there exists an unbounded sequence of nondegenerate solutions v_n of (10).*

Proof. All the solutions of the autonomous equation $-v'' = \mu v^m$ are periodic and can be parametrized by their energy

$$E = \frac{1}{2} v'^2 + \frac{\mu}{m+1} v^{m+1}.$$

Let T_E denote the period of the solution v_E . We can suppose $v_E(0) = 0$, so $v'_E(0) = \sqrt{2E}$. The other boundary condition $v_E(\pi) = 0$ is satisfied iff

$$(25) \quad k \frac{T_E}{2} = \pi \quad \text{for some } k \in \mathbb{N}.$$

By symmetry and energy conservation $v_E(T_E/4) = [(m + 1)E/\mu]^{\frac{1}{m+1}}$. So

$$\begin{aligned} T_E &= 4 \int_0^{\left[\frac{(m+1)E}{\mu}\right]^{\frac{1}{m+1}}} \left[2\left(E - \frac{\mu x^{m+1}}{m+1}\right)\right]^{-1/2} dx \\ &= \frac{4(m+1/\mu)^{\frac{1}{m+1}}}{E^{\frac{1}{2} - \frac{1}{m+1}}} \int_0^1 \frac{dy}{\sqrt{2(1-y^{m+1})}} =: \frac{C(m, \mu)}{E^{\frac{1}{2} - \frac{1}{m+1}}} \end{aligned}$$

by the change of variable $y = x [E(m + 1)/\mu]^{-\frac{1}{m+1}}$, and (25) is satisfied at infinitely many energy levels. Let $\bar{E} > 0$ such that $T_{\bar{E}} = 2\pi/k$ and denote the solution $\bar{v} := v_{\bar{E}}$.

Let us prove that \bar{v} is nondegenerate. Any solution h of the linearized equation at \bar{v} ,

$$(26) \quad -h''(x) = \mu m \bar{v}^{m-1}(x) h(x),$$

can be written as $h = A\bar{v}' + B\beta$, $A, B \in \mathbb{R}$, because $\bar{v}'(x)$ and $\beta(x) := (\partial_E v_E)|_{E=\bar{E}}(x)$ are solutions of (26); they are independent because $\bar{v}'(0) \neq 0$ while $\beta(0) = 0$. If $h(0) = 0$, then $A = 0$. We claim that $\beta(\pi) \neq 0$; as a consequence, if $h(\pi) = 0$, then $B = 0$, and so $h = 0$, i.e., \bar{v} is nondegenerate. To prove that $\beta(\pi) \neq 0$, we differentiate at \bar{E} the identity $v_E(kT_E/2) = 0$,

$$\beta(\pi) + \bar{v}'(\pi)(\partial_E T_E)|_{E=\bar{E}} = 0.$$

Since $\bar{v}'(\pi) = (-1)^k \sqrt{2\bar{E}} \neq 0$ and $\partial_E T_E \neq 0$, we get $\beta(\pi) \neq 0$. \square

LEMMA 3. *If $f_0(x, 0) = d_u f_0(x, 0) = 0$, then $\bar{v} = 0$ is a nondegenerate solution of (10) for every μ .*

Proof. The linearized equation $-(ph')' = 0$, $h(0) = h(\pi) = 0$ has only the trivial solution. \square

When Hypothesis (V) holds at some (μ_0, \bar{v}) , we solve first the bifurcation equation in (8) using the standard implicit function theorem. We find, for every w small enough and μ in a neighborhood of μ_0 , a unique solution $v(\mu, w)$ of the bifurcation equation.

LEMMA 4 (solution of the bifurcation equation). *There exist $0 < R < R_0$, a neighborhood $[\mu_1, \mu_2]$ of μ_0 , and a C^∞ map*

$$[\mu_1, \mu_2] \times \{w \in W \cap X_{\sigma,s} : \|w\|_{\sigma,s} < R\} \rightarrow V, \quad (\mu, w) \mapsto v(\mu, w)$$

such that $v(\mu, w)$ solves the bifurcation equation in (8).

Proof. The linear operator

$$h \mapsto -(ph')' - \mu_0 d_v \Pi_V f(v)[h] = -(ph')' - \mu_0 f'_0(v) h$$

is invertible on $H_0^1(0, \pi)$ by Hypothesis (V). Then we apply the implicit function theorem. \square

Remark 7. The solutions of the 0th order bifurcation equation (10) found in Lemmas 2 and 3 are nondegenerate for every μ , so, in that case, we can continue $v(\mu, w)$ for all $[\mu_1, \mu_2] = [0, \bar{\mu}]$.

We denote by $\lambda_j(\mu, w) := \omega_j^2(\mu, w)$ the eigenvalues of the Sturm–Liouville problem

$$(27) \quad \begin{cases} -(py')' - \mu \Pi_V f'(v(\mu, w) + w) y = \lambda \rho y, \\ y(0) = y(\pi) = 0. \end{cases}$$

LEMMA 5. *The eigenvalues of (27) satisfy the continuity property*

$$(28) \quad |\lambda_j(\mu, w) - \lambda_j(\mu', w')| \leq K(|\mu - \mu'| + \|w - w'\|_{\sigma,s})$$

for some constant $K > 0$ independent of j .

Proof. For the proof of the lemma, see the appendix. \square

The nondegeneracy of $\bar{v} = v(\mu_0, 0)$ means that $\lambda_j(\mu_0, 0) \neq 0$ for all j . By (28),

$$(29) \quad \delta_0 := \inf \left\{ |\lambda_j(\mu, w)| : j \geq 1, \mu \in [\mu_1, \mu_2], \|w\|_{\sigma_0/2} \leq R \right\} > 0,$$

taking, if necessary, $|\mu_2 - \mu_1|$ and R smaller in Lemma 4.

Note also that the index j_0 of the smallest positive eigenvalue is constant, independently of (μ, w) .

4. Solution of the range equation. It remains to solve the range equation

$$(30) \quad L_\omega w = \mu \Pi_W \mathcal{F}(\mu, w),$$

where

$$\mathcal{F}(\mu, w) := f(v(\mu, w) + w).$$

By Lemmas 1 and 4, \mathcal{F} is C^∞ and bounded, together with its derivatives, on $[\mu_1, \mu_2] \times B_R$, where $B_R := \{w \in W \cap X_{\sigma,s} : \|w\|_{\sigma,s} < R\}$.

4.1. The Nash–Moser recursive scheme. We define the sequence of finite-dimensional subspaces

$$W^{(n)} := \left\{ w = \sum_{1 \leq |l| \leq N_n} w_l(x) e^{ilt} \right\} \subset W,$$

where

$$N_n := N_0 2^n$$

and $N_0 \in \mathbb{N}$ will be fixed in Lemma 9. We also set

$$W^{(n)\perp} := \left\{ w = \sum_{|l| > N_n} w_l(x) e^{ilt} \in W \right\}$$

and denote by P_n , resp., P_n^\perp , the projection on $W^{(n)}$, resp., $W^{(n)\perp}$. Note that $P_n \circ \Pi_W = P_n$.

LEMMA 6 (smoothing estimate). *For $w \in W^{(n)\perp}$, if $0 < \sigma'' < \sigma'$,*

$$(31) \quad \|w\|_{\sigma'',s} \leq \exp[-(\sigma' - \sigma'')N_n] \|w\|_{\sigma',s}.$$

Proof. It follows from the definition of the norms $\|\cdot\|_{\sigma,s}$ and $W^{(n)\perp}$; see, e.g., [16, 8]. \square

The key property for the construction of the iterative sequence is the invertibility of the linear operator

$$(32) \quad \begin{aligned} \mathcal{L}_n(w)h &:= -L_\omega h + \mu P_n[d_w \mathcal{F}(\mu, w)h] \\ &= -L_\omega h + \mu P_n[f'(v(\mu, w) + w)(h + d_w v(\mu, w)[h])] \quad \forall h \in W^{(n)}. \end{aligned}$$

LEMMA 7 (inversion of the linear problem). *Let $\omega > 0$, $\tau \in (1, 2)$, $\gamma \in (0, 1)$, $\gamma < \omega$, and $\sigma \in (0, \sigma_0]$. Assume the ‘‘Melnikov’’ nonresonance conditions*

$$(33) \quad \left| \omega l - \frac{j}{c} \right| > \frac{\gamma}{l^\tau} \quad \forall l = 1, 2, \dots, N_n, \quad \forall j \geq 1,$$

where c is defined in (16), and

$$(34) \quad |\omega^2 l^2 - \lambda_j(\mu, w)| > \frac{\gamma \omega}{l^{\tau-1}} \quad \forall l = 1, 2, \dots, N_n, \quad j \geq 1,$$

where $\lambda_j(\mu, w)$ are the eigenvalues of (27).

Let $u := v(\mu, w) + w$. There exist K_1, K'_1 such that if

$$(35) \quad \frac{\mu}{\gamma^3 \omega} \|\Pi_W f'(u)\|_{\sigma, 1 + \frac{\tau(\tau-1)}{2-\tau}} < K'_1,$$

then $\mathcal{L}_n(w)$ is invertible and

$$(36) \quad \|\mathcal{L}_n(w)^{-1} h\|_\sigma \leq \frac{K_1 N_n^{\tau-1}}{\gamma \omega} \|h\|_\sigma \quad \forall h \in W^{(n)}.$$

Proof. For the proof of the lemma, see section 6. \square

Remark 8. The condition $\omega > 0$ means that (1) is nonautonomous. Indeed, if $\omega = 0$, the nonlinearity $f(x, \omega t, u) = f(x, 0, u)$ is independent of t .

For $\vartheta := 3\sigma_0/\pi^2$ we define the sequence

$$(37) \quad \sigma_{n+1} := \sigma_n - \frac{\vartheta}{(n+1)^2}, \quad \sigma_0 > \sigma_1 > \sigma_2 > \dots > \frac{\sigma_0}{2}.$$

Let A_0 denote the open set

$$A_0 := \left\{ (\mu, \omega) \in (\mu_1, \mu_2) \times (\gamma, +\infty) : |\omega l - \omega_j| > \frac{\gamma}{l^\tau} \quad \forall l = 1, \dots, N_0, \quad j \geq 1 \right\},$$

where ω_j are defined by (11).

LEMMA 8 (approximate solution). *There exist K_2, K'_2 such that if $(\mu, \omega) \in A_0$ and $\mu N_0^{\tau-1}/\gamma \omega < K'_2$, then there exists a solution $w_0 := w_0(\mu, \omega) \in W^{(0)}$ of*

$$L_\omega w_0 = \mu P_0 \mathcal{F}(\mu, w_0)$$

satisfying $\|w_0\|_{\sigma_0} \leq \mu K_2 N_0^{\tau-1}/\gamma \omega$.

Proof. By definition of A_0 , the eigenvalues of $(1/\rho)L_\omega$ satisfy

$$|\omega^2 l^2 - \lambda_j| > \frac{\gamma \omega}{l^{\tau-1}} \quad \forall l = 1, 2, \dots, N_0, \quad j \geq 1,$$

so L_ω is invertible on $W^{(0)}$ and, for some K ,

$$(38) \quad \|L_\omega^{-1} h\|_{\sigma_0} \leq \frac{K N_0^{\tau-1}}{\gamma \omega} \|h\|_{\sigma_0} \quad \forall h \in W^{(0)}.$$

Then we look for a solution $w_0 \in W^{(0)}$ of $w_0 = \mu L_\omega^{-1} P_0 \mathcal{F}(\mu, w_0)$. The right-hand side term is a contraction in $\{\|w_0\|_{\sigma_0} < R\}$ if $\mu N_0^{\tau-1}/\gamma \omega$ is sufficiently small. \square

Given $w_n \in W^{(n)}$, $\|w_n\|_{\sigma_n} < R$, and $A_n \subseteq A_0$, we define

$$(39) \quad A_{n+1} := \left\{ (\mu, \omega) \in A_n : \left| \omega l - \omega_j(\mu, w_n) \right| > \frac{\gamma}{l^\tau}, \quad \left| \omega l - \frac{j}{c} \right| > \frac{\gamma}{l^\tau} \right. \\ \left. \forall l = 1, 2, \dots, N_{n+1}, \quad j \geq 1 \right\} \subseteq A_n,$$

where $\lambda_j(\mu, w_n) = \omega_j^2(\mu, w_n)$ are defined in (27) with $w = w_n$.

In Lemma 8 we have constructed $h_0 := w_0$ for $(\mu, \omega) \in A_0$. Next, we proceed by induction. By means of w_0 we define the set A_1 as above, and we find $w_1 := h_0 + h_1 \in W^{(1)}$ for every $(\mu, \omega) \in A_1$ by Lemma 9 below. Then we define A_2 , we find $w_2 \in W^{(2)}$, and so on. The main goal of the construction is to prove that, at the end of the recurrence, the set of parameters $(\mu, \omega) \in \cap_n A_n$ is actually a large set (see Lemmas 12 and 13).

LEMMA 9 (inductive step). *Fix $\chi := 3/2$. There exist $N_0 \in \mathbb{N}$ (depending only on $\rho, p, f, \bar{\mu}, \bar{\nu}, \tau$) and $K'_3 \leq K'_2/N_0^{\tau-1}$ with the following property.*

Suppose that $h_i \in W^{(i)}$ for all $i = 0, \dots, n$ satisfy

$$(40) \quad \|h_i\|_{\sigma_i} < \frac{\mu K_3 N_0^{\tau-1}}{\gamma \omega} \exp(-\chi^i),$$

where $K_3 := eK_2$ and K_2 is the constant in Lemma 8; for all $k = 0, \dots, n$, let $w_k := h_0 + \dots + h_k$ satisfy $\|w_k\|_{\sigma_k} < R$ and

$$(41) \quad L_\omega w_k = \mu P_k \mathcal{F}(\mu, w_k)$$

and suppose that $(\mu, \omega) \in A_n$, where A_{i+1} is constructed by means of w_i as shown above.

If $(\mu, \omega) \in A_{n+1}$ and $\mu/\gamma^3 \omega < K'_3$, then there exists $h_{n+1} \in W^{(n+1)}$ satisfying

$$(42) \quad \|h_{n+1}\|_{\sigma_{n+1}} < \frac{\mu K_3 N_0^{\tau-1}}{\gamma \omega} \exp(-\chi^{n+1})$$

such that $w_{n+1} = w_n + h_{n+1}$ verifies $\|w_{n+1}\|_{\sigma_{n+1}} < R$ and

$$(43) \quad L_\omega w_{n+1} = \mu P_{n+1} \mathcal{F}(\mu, w_{n+1}).$$

Proof. In short $\mathcal{F}(w) := \mathcal{F}(\mu, w)$ and $D\mathcal{F}(w) := d_w \mathcal{F}(\mu, w)$. Equation (43) for $w_{n+1} = w_n + h_{n+1}$ is $L_\omega[w_n + h_{n+1}] = \mu P_{n+1} \mathcal{F}(w_n + h_{n+1})$.

By assumption, w_n satisfies (41) for $k = n$, namely, $L_\omega w_n = \mu P_n \mathcal{F}(w_n)$, so the equation for h_{n+1} can be written as

$$(44) \quad \mathcal{L}_{n+1}(w_n)h_{n+1} + \mu(P_{n+1} - P_n)\mathcal{F}(w_n) + \mu P_{n+1}Q = 0,$$

where, as defined in (32), $\mathcal{L}_{n+1}(w_n)h_{n+1} := -L_\omega h_{n+1} + \mu P_{n+1} D\mathcal{F}(w_n)h_{n+1}$, and Q denotes the quadratic remainder

$$Q = Q(w_n, h_{n+1}) := \mathcal{F}(w_{n+1}) - \mathcal{F}(w_n) - D\mathcal{F}(w_n)h_{n+1}.$$

Step 1. Inversion of $\mathcal{L}_{n+1}(w_n)$. We verify the assumptions of Lemma 7. By definition of A_{n+1} , ω satisfies (33). If $\lambda_j(\mu, w_n) < 0$, then $|\omega^2 l^2 - \lambda_j(\mu, w_n)| \geq \omega^2 l^2 > \gamma \omega / l^{\tau-1}$ because $\omega > \gamma$. If $\lambda_j(\mu, w_n) > 0$, we have

$$|\omega^2 l^2 - \lambda_j(\mu, w_n)| \geq |\omega l - \omega_j(\mu, w_n)| \omega l > \frac{\gamma \omega}{l^{\tau-1}} \quad \forall l = 1, \dots, N_{n+1}$$

because $(\mu, \omega) \in A_{n+1}$. In both cases the nonresonance condition (34) holds.

To verify (35) we need an estimate for w_n . Let $\eta := \tau(\tau - 1)/(2 - \tau)$ and $\alpha > 0$. Using the elementary inequality

$$\frac{1 + l^{2(1+\eta)}}{1 + l^2} \cdot \frac{e^{2(\sigma-\alpha)|l|}}{e^{2\sigma|l|}} \leq \frac{2l^{2\eta}}{e^{2\alpha|l|}} \leq 2 \max_{y>0} (y^{2\eta} e^{-2\alpha y}) = 2 \left(\frac{\eta}{\alpha e}\right)^{2\eta} \quad \forall l \neq 0,$$

we deduce

$$\|h_i\|_{\sigma_{n+1}, 1+\eta} \leq \frac{C_\eta}{(\sigma_i - \sigma_{n+1})^\eta} \|h_i\|_{\sigma_i},$$

where $C_\eta := \sqrt{2}(\eta/e)^\eta$. Since $\sigma_i - \sigma_{n+1} \geq \sigma_i - \sigma_{i+1}$ for every $i \leq n$,

$$\|w_n\|_{\sigma_{n+1}, 1+\eta} \leq \sum_{i=0}^n \|h_i\|_{\sigma_{n+1}, 1+\eta} \leq C_\eta \sum_{i=0}^n \frac{\|h_i\|_{\sigma_i}}{(\sigma_i - \sigma_{i+1})^\eta} \leq S_\eta \frac{\mu K_3 N_0^{\tau-1}}{\gamma\omega},$$

using (40) where $S_\eta := (C_\eta/\vartheta^\eta) \sum_{i=0}^{+\infty} (i+1)^{2\eta} \exp(-\chi^i) < +\infty$. If

$$\frac{S_\eta \mu K_3 N_0^{\tau-1}}{\gamma\omega} < R,$$

then

$$\|f'(u_n)\|_{\sigma_{n+1}, 1+\eta} \leq K$$

for some K , where $u_n := v(\mu, w_n) + w_n$. Hence Hypothesis (35) is verified for $\mu/\gamma^3\omega$ sufficiently small.

Analogously we get $\|w_n\|_{\sigma_n} < R$ if $\mu N_0^{\tau-1}/\gamma\omega$ is small enough.

By Lemma 7 the operator $\mathcal{L}_{n+1}(w_n)$ is invertible on $W^{(n+1)}$ and

$$(45) \quad \|\mathcal{L}_{n+1}(w_n)^{-1}h\|_{\sigma_{n+1}} \leq \frac{K_1 N_{n+1}^{\tau-1}}{\gamma\omega} \|h\|_{\sigma_{n+1}} \quad \forall h \in W^{(n+1)}.$$

Equation (44) amounts to the fixed point problem

$$h_{n+1} = -\mu \mathcal{L}_{n+1}(w_n)^{-1} [(P_{n+1} - P_n)\mathcal{F}(w_n) + P_{n+1}Q] := \mathcal{G}(h_{n+1})$$

for $h_{n+1} \in W^{(n+1)}$.

Step 2. \mathcal{G} is a contraction. We prove that \mathcal{G} is a contraction on the ball $B_{n+1} := \{\|h\|_{\sigma_{n+1}} < r_{n+1}\}$, where $r_{n+1} := (\mu K_3 N_0^{\tau-1}/\gamma\omega) \exp(-\chi^{n+1})$, implying (42). By (31)

$$\|(P_{n+1} - P_n)\mathcal{F}(w_n)\|_{\sigma_{n+1}} \leq \|\mathcal{F}(w_n)\|_{\sigma_n} \exp[-(\sigma_n - \sigma_{n+1})N_n].$$

Since $\|w_n\|_{\sigma_n} < R$, we have $\|Q\|_{\sigma_{n+1}} \leq K \|h_{n+1}\|_{\sigma_{n+1}}^2$. Hence, by (45),

$$\|\mathcal{G}(h_{n+1})\|_{\sigma_{n+1}} \leq K \frac{\mu N_{n+1}^{\tau-1}}{\gamma\omega} \left(\exp[-(\sigma_n - \sigma_{n+1})N_n] + \|h_{n+1}\|_{\sigma_{n+1}}^2 \right).$$

Therefore $\mathcal{G}(B_{n+1}) \subseteq B_{n+1}$ if

$$(46) \quad \frac{\mu K N_{n+1}^{\tau-1}}{\gamma\omega} \exp[-(\sigma_n - \sigma_{n+1})N_n] < \frac{r_{n+1}}{2}, \quad \frac{\mu K N_{n+1}^{\tau-1}}{\gamma\omega} r_{n+1}^2 < \frac{r_{n+1}}{2}.$$

By the definition of σ_n in (37) and $N_n := N_0 2^n$, the first inequality is verified for every $n \geq 0$ if $\sigma_0 N_0$ is greater than a constant depending only on χ, K, K_3 . The second inequality is verified for every $n \geq 0$ if $\mu N_0^{\tau-1} / \gamma \omega$ is small enough.

The estimate for $\|\mathcal{G}h - \mathcal{G}k\|$, $h, k \in B_{n+1}$ is similar. The lemma now follows from the contraction mapping theorem. \square

Remark 9. In the previous scheme h_{n+1} is found as an exact solution of (44). We find this convenient to prove the regularity of h_{n+1} with respect to the parameters (μ, ω) in Lemma 10. However, other schemes are possible. For example, we could define h_{n+1} as a solution of the linearized equation $\mathcal{L}_{n+1}(w_n)h + \mu(P_{n+1} - P_n)\mathcal{F}(w_n) = 0$.

COROLLARY 1 (existence). *Suppose $A_\infty := \bigcap_{n \geq 0} A_n \neq \emptyset$. If $(\mu, \omega) \in A_\infty$ and $\mu/\gamma^3 \omega < K'_3$, then*

$$w_\infty(\mu, \omega) := \sum_{n \geq 0} h_n(\mu, \omega) \in W \cap X_{\sigma_0/2}$$

is a solution of the range equation (30) satisfying $\|w_\infty\|_{\sigma_0/2} \leq K_\infty \mu/\gamma \omega$ for some K_∞ .

Proof. Since w_n solves (41) for $k = n$,

$$-L_\omega w_n + \mu \Pi_W f(u_n) = \mu P_n^\perp f(u_n) \in W^{(n)\perp},$$

where $u_n := v(\mu, w_n) + w_n$. By (31)

$$\lim_{n \rightarrow +\infty} \|-L_\omega w_n + \mu f(u_n)\|_{\sigma_0/2} \leq \lim_{n \rightarrow +\infty} K \exp[-(\sigma_n - \sigma_0/2)N_n] = 0.$$

Since $w_n \rightarrow w_\infty$ in $\|\cdot\|_{\sigma_0/2}$ also $f(u_n) \rightarrow f(u_\infty)$ in the same norm, while $L_\omega w_n \rightarrow L_\omega w_\infty$ in the sense of distributions. So w_∞ is a weak solution of the range equation (30). \square

Remark 10. We shall prove, as a consequence of Lemma 12 and section 4.3, that A_∞ is actually a positive measure set. One possible way to prove it uses the Whitney extension of w_∞ of section 4.2.

4.2. Whitney C^∞ extension. The functions h_n constructed in Lemmas 8 and 9 depend smoothly on the parameters (μ, ω) .

LEMMA 10. *There exist K_4 and $K'_4 \leq K'_3$ such that the maps*

$$h_i : A_i \cap \{\mu/\gamma^3 \omega < K'_4\} \rightarrow W^{(i)}$$

are C^∞ and

$$\|\partial_\omega h_i(\mu, \omega)\|_{\sigma_i} \leq \frac{K_4 \mu}{\gamma^2 \omega} \exp(-\chi_0^i), \quad \|\partial_\mu h_i(\mu, \omega)\|_{\sigma_i} \leq \frac{K_4}{\gamma \omega} \exp(-\chi_0^i),$$

where $\chi_0 := (1 + \chi)/2 = 5/4$.

Proof. Since $w_0 = \mu L_\omega^{-1} P_0 \mathcal{F}(\mu, w_0)$, by the implicit function theorem the map w_0 is C^∞ . Differentiating the identity $L_\omega(L_\omega^{-1} h) = h$ with respect to ω , by (38) we get $\|\partial_\omega L_\omega^{-1} h\|_{\sigma_0} \leq (K/\gamma^2 \omega) \|h\|_{\sigma_0}$. For $\mu/\gamma \omega$ small,

$$\|\partial_\omega w_0\|_{\sigma_0} \leq \frac{K \mu}{\gamma^2 \omega}.$$

Differentiating with respect to μ we get $\|\partial_\mu w_0\|_{\sigma_0} \leq K'/\gamma \omega$ for some K' .

By induction, suppose that h_i depends smoothly on $(\mu, \omega) \in A_i$ for every $i = 0, \dots, n$. For $(\mu, \omega) \in A_{n+1}$, by (43), h_{n+1} is a solution of

$$(47) \quad -L_\omega h_{n+1} + \mu P_{n+1}[\mathcal{F}(w_n + h_{n+1}) - \mathcal{F}(w_n)] + \mu(P_{n+1} - P_n)\mathcal{F}(w_n) = 0.$$

By the implicit function theorem $h_{n+1} \in C^\infty$ once we prove that

$$\mathcal{L}_{n+1}(w_{n+1})[z] := -L_\omega z + \mu P_{n+1} D\mathcal{F}(w_n + h_{n+1})[z]$$

is invertible. By (45), $\mathcal{L}_{n+1}(w_n)$ is invertible. Hence it is sufficient that

$$\left\| \mathcal{L}_{n+1}^{-1}(w_n)(\mathcal{L}_{n+1}(w_{n+1}) - \mathcal{L}_{n+1}(w_n)) \right\|_{\sigma_{n+1}} < \frac{1}{2},$$

which holds true for $\mu^2/\gamma\omega$ small enough; indeed, by (42),

$$\|\mathcal{L}_{n+1}(w_{n+1}) - \mathcal{L}_{n+1}(w_n)\|_{\sigma_{n+1}} \leq K\mu\|h_{n+1}\|_{\sigma_{n+1}} \leq \frac{\mu^2 K' N_0^{\tau-1}}{\gamma\omega} \exp(-\chi^{n+1}).$$

Finally (45) implies

$$(48) \quad \|\mathcal{L}_{n+1}(w_{n+1})^{-1}\|_{\sigma_{n+1}} \leq \frac{2K_1 N_{n+1}^{\tau-1}}{\gamma\omega}.$$

Differentiating (47) with respect to ω

$$(49) \quad \begin{aligned} \mathcal{L}_{n+1}(w_{n+1})[\partial_\omega h_{n+1}] &= 2\omega\rho(x)(h_{n+1})_{tt} - \mu(P_{n+1} - P_n)D\mathcal{F}(w_n)\partial_\omega w_n \\ &\quad - \mu P_{n+1}[D\mathcal{F}(w_n + h_{n+1}) - D\mathcal{F}(w_n)]\partial_\omega w_n \end{aligned}$$

and, using (48) and (31),

$$\begin{aligned} \|\partial_\omega h_{n+1}\|_{\sigma_{n+1}} &\leq \frac{KN_{n+1}^{\tau-1}}{\gamma\omega} \left(\omega N_{n+1}^2 \|h_{n+1}\|_{\sigma_{n+1}} + \frac{\mu\|\partial_\omega w_n\|_{\sigma_n}}{\exp[(\sigma_n - \sigma_{n+1})N_n]} \right. \\ &\quad \left. + \mu\|h_{n+1}\|_{\sigma_{n+1}} \|\partial_\omega w_n\|_{\sigma_n} \right). \end{aligned}$$

We note that $\|\partial_\omega w_n\|_{\sigma_n} \leq \sum_{i=0}^n \|\partial_\omega h_i\|_{\sigma_i}$. Using (46) the sequence $a_n := \|\partial_\omega h_n\|_{\sigma_n}$ satisfies

$$\begin{aligned} a_{n+1} &\leq \frac{KN_{n+1}^{\tau-1}}{\gamma\omega} \left(\omega N_{n+1}^2 r_{n+1} + \frac{\omega\gamma r_{n+1}}{N_{n+1}^{\tau-1}} \sum_{i=0}^n a_i + \mu r_{n+1} \sum_{i=0}^n a_i \right) \\ &\leq b_{n+1} \left(1 + \sum_{i=0}^n a_i \right), \quad \text{where } b_{n+1} := \frac{K\mu}{\gamma^2\omega} N_{n+1}^{\tau+1} \exp(-\chi^{n+1}), \end{aligned}$$

recalling that $r_{n+1} = (\mu K/\gamma\omega) \exp(-\chi^{n+1})$. By induction, for $K\mu/\omega\gamma^2 < 1$, we have $a_n \leq 2b_n$ and

$$\|\partial_\omega h_{n+1}\|_{\sigma_{n+1}} \leq \frac{K\mu}{\gamma^2\omega} N_{n+1}^{\tau+1} \exp(-\chi^{n+1}) \leq \frac{K'\mu}{\gamma^2\omega} \exp(-\chi_0^{n+1}),$$

where $\chi_0 := (1 + \chi)/2$. It follows that $\|\partial_\omega w_{n+1}\|_{\sigma_{n+1}} \leq K\mu/\gamma^2\omega$.

Differentiating (47) with respect to μ we obtain the estimate for $\partial_\mu h_{n+1}$. \square
 Define, for $\nu_0 > 0$,

$$(50) \quad A_n^* := \left\{ (\mu, \omega) \in A_n : \text{dist}((\mu, \omega), \partial A_n) > \frac{\nu_0 \gamma^4}{N_n^3} \right\},$$

$$\tilde{A}_n := \left\{ (\mu, \omega) \in A_n : \text{dist}((\mu, \omega), \partial A_n) > \frac{2\nu_0 \gamma^4}{N_n^3} \right\} \subset A_n^*.$$

LEMMA 11 (Whitney extension). *There exists a C^∞ map*

$$\tilde{w} : A_0 \cap \left\{ (\mu, \omega) : \frac{\mu}{\gamma^3 \omega} < K_4' \right\} \rightarrow W \cap X_{\sigma_0/2}$$

satisfying

$$(51) \quad \|\tilde{w}(\mu, \omega)\|_{\sigma_0/2} \leq \frac{K_5 \mu}{\gamma \omega},$$

$$\|\partial_\omega \tilde{w}(\mu, \omega)\|_{\sigma_0/2} \leq \frac{C(\nu_0) \mu}{\gamma^5 \omega}, \quad \|\partial_\mu \tilde{w}(\mu, \omega)\|_{\sigma_0/2} \leq \frac{C(\nu_0)}{\gamma^5 \omega}$$

for some K_5 and for some $C(\nu_0) > 0$, such that, for $(\mu, \omega) \in \tilde{A}_\infty := \bigcap_{n \geq 0} \tilde{A}_n$, $\tilde{w}(\mu, \omega)$ solves the range equation (30).

Moreover there exists a sequence of C^∞ maps

$$\tilde{w}_n : A_0 \cap \left\{ (\mu, \omega) : \frac{\mu}{\gamma^3 \omega} < K_4' \right\} \rightarrow W^{(n)}$$

such that $\tilde{w}_n(\mu, \omega) = w_n(\mu, \omega)$ for $(\mu, \omega) \in \tilde{A}_n$, and

$$(52) \quad \|\tilde{w}(\mu, \omega) - \tilde{w}_n(\mu, \omega)\|_{\sigma_0/2} \leq \frac{K_5 \mu}{\gamma \omega} \exp(-\chi^n).$$

Proof. Let $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^+$ be a C^∞ -function supported in the open ball $B(0, 1)$ of center 0 and radius 1 and with $\int_{\mathbb{R}^2} \varphi = 1$. Let $\varphi_n : \mathbb{R}^2 \rightarrow \mathbb{R}^+$ be the mollifier

$$\varphi_n(x) := \frac{N_n^6}{\nu_0^2 \gamma^8} \varphi\left(\frac{N_n^3}{\nu_0 \gamma^4} x\right).$$

Supp $(\varphi_n) \subset B(0, \nu_0 \gamma^4 / N_n^3)$ and $\int_{\mathbb{R}^2} \varphi_n = 1$. We define $\psi_n : \mathbb{R}^2 \rightarrow \mathbb{R}$ as

$$\psi_n(x) := (\varphi_n * \chi_{A_n^*})(x) = \int_{\mathbb{R}^2} \varphi_n(y - x) \chi_{A_n^*}(y) dy,$$

where $\chi_{A_n^*}$ is the characteristic function of the set A_n^* . ψ_n is C^∞ ,

$$(53) \quad |D\psi_n(x)| \leq \int_{\mathbb{R}^2} |D\varphi_n(x - y)| \chi_{A_n^*}(y) dy \leq \frac{N_n^3}{\nu_0 \gamma^4} C,$$

where $C := \int_{\mathbb{R}^2} |D\varphi| dy$,

$$0 \leq \psi_n(x) \leq 1, \quad \text{supp}(\psi_n) \subset A_n, \quad \psi_n(x) = 1 \quad \forall x \in \tilde{A}_n.$$

We define, for $(\mu, \omega) \in A_0$, the C^∞ -functions

$$\tilde{h}_n(\mu, \omega) := \begin{cases} \psi_n(\mu, \omega)h_n(\mu, \omega) & \text{if } (\mu, \omega) \in A_n, \\ 0 & \text{if } (\mu, \omega) \notin A_n, \end{cases}$$

and

$$\tilde{w}_n(\mu, \omega) := \sum_{i=0}^n \tilde{h}_i, \quad \tilde{w}(\mu, \omega) := \sum_{i \geq 0} \tilde{h}_i,$$

which is a series if $(\mu, \omega) \in A_\infty := \bigcap_{n \geq 0} A_n$.

The estimate for $\|\tilde{w}\|_{\sigma_0/2}$ follows by $\|\tilde{h}_i\|_{\sigma_i} \leq \|h_i\|_{\sigma_i}$ (because $0 \leq \psi_i \leq 1$) and (40). The estimates for the derivatives in (51) follow by differentiating the product $\tilde{h}_i = \psi_i h_i$ and using (53), (40), and Lemma 10. Similarly it follows that \tilde{w} is in C^∞ ; see [8] for details.

For $(\mu, \omega) \in \tilde{A}_n$, $\psi_n(\mu, \omega) = 1$, implying $\tilde{w}_n = w_n$. As a consequence, for $(\mu, \omega) \in \tilde{A}_\infty := \bigcap_{n \geq 0} \tilde{A}_n$, by Corollary 1, $\tilde{w} = w_\infty$ solves (30).

Finally, using (40),

$$\|\tilde{w} - \tilde{w}_n\|_{\sigma_0/2} \leq \sum_{i \geq n+1} \|\tilde{h}_i\|_{\sigma_i} \leq \sum_{i \geq n+1} \frac{K\mu}{\gamma\omega} \exp(-\chi^i) \leq \frac{K'\mu}{\gamma\omega} \exp(-\chi^n). \quad \square$$

In the next lemma we fix the constant ν_0 introduced in (50).

LEMMA 12. *There exist $\nu_0 > 0$ and $K'_5 \leq K'_4$ such that if $\mu/\gamma^3\omega < K'_5$, then*

$$B_\gamma \subseteq \tilde{A}_n \subset A_n \quad \forall n \geq 0,$$

where B_γ is defined in (15) taking $C' \leq K'_5$.

Proof. The proof is by induction. Let $(\mu, \omega) \in B_\gamma$. Then $(\mu, \omega) \in \tilde{A}_0$ if A_0 contains the closed ball of center (μ, ω) and radius $2\nu_0\gamma^4/N_0^3$. Let (ω', μ') belong to such a ball. Then, for all $l = 1, \dots, N_0$,

$$|\omega'l - \omega_j| \geq |\omega l - \omega_j| - |\omega - \omega'|l > \frac{2\gamma}{l^\tau} - \frac{2\nu_0\gamma^4}{N_0^3}l \geq \frac{\gamma}{l^\tau}$$

if $\nu_0 \leq 1/2$.

Suppose now that $B_\gamma \subseteq \tilde{A}_n$ and let $(\mu, \omega) \in B_\gamma$. To prove that $(\mu, \omega) \in \tilde{A}_{n+1}$, we have to show that the closed ball of center (μ, ω) and radius $2\nu_0\gamma^4/N_{n+1}^3$ is contained in A_{n+1} . Let (μ', ω') belong to such a ball. The nonresonance condition on $|\omega'l - j/c|$ is verified, as above, for $\nu_0 \leq 1/2$. For the other condition, we denote in short $\omega_j^n(\mu', \omega') := \omega_j(\mu', w_n(\mu', \omega'))$ (see (27) for the definition of $\omega_j(\mu, w)$). It results, for all $l = 1, \dots, N_{n+1}$, in

$$\begin{aligned} |\omega'l - \omega_j^n(\mu', \omega')| &\geq |\omega l - \tilde{\omega}_j(\mu, \omega)| - |\omega - \omega'|l - |\omega_j^n(\mu', \omega') - \tilde{\omega}_j(\mu, \omega)| \\ &> \frac{2\gamma}{l^\tau} - \frac{2\nu_0\gamma^4 l}{N_{n+1}^3} - |\omega_j^n(\mu', \omega') - \tilde{\omega}_j(\mu, \omega)| \\ (54) \quad &> \frac{3\gamma}{2l^\tau} - |\omega_j^n(\mu', \omega') - \tilde{\omega}_j(\mu, \omega)| \end{aligned}$$

if $\nu_0 \leq 1/4$. Now we estimate the last term

$$|\omega_j^n(\mu', \omega') - \tilde{\omega}_j(\mu, \omega)| = \frac{|\lambda_j^n(\mu', \omega') - \tilde{\lambda}_j(\mu, \omega)|}{|\tilde{\omega}_j(\mu, \omega)| + |\omega_j^n(\mu', \omega')|} \leq \frac{|\lambda_j^n(\mu', \omega') - \tilde{\lambda}_j(\mu, \omega)|}{\sqrt{\delta_0}}$$

by (29), both for $j < j_0$ and for $j \geq j_0$. By the comparison principle (28)

$$\delta_0^{-1/2} |\lambda_j^n(\mu', \omega') - \tilde{\lambda}_j(\mu, \omega)| \leq K|\mu - \mu'| + K\|w_n(\mu', \omega') - \tilde{w}(\mu, \omega)\|_{\sigma_0/2}.$$

By Lemma 10, $\|\partial_\omega w_n\|_{\sigma_0/2}, \|\partial_\mu w_n\|_{\sigma_0/2} \leq K/\gamma^2 \omega$ for some other K , and being $\omega, \omega' > \gamma$,

$$K\|w_n(\mu', \omega') - w_n(\mu, \omega)\|_{\sigma_0/2} \leq \frac{K'}{\gamma^3} \frac{\nu_0 \gamma^4}{N_{n+1}^3} < \frac{\gamma}{8l^\tau} \quad \forall l = 1, \dots, N_{n+1}$$

if ν_0 is small enough ($1 < \tau < 2$). On the other hand, since $(\mu, \omega) \in \tilde{A}_n$ we have $w_n(\mu, \omega) = \tilde{w}_n(\mu, \omega)$ (Lemma 11) and, by (52),

$$K\|w_n(\mu, \omega) - \tilde{w}(\mu, \omega)\|_{\sigma_0/2} \leq \frac{K'\mu}{\gamma\omega} \exp(-\chi^n) < \frac{\gamma}{8l^\tau} \quad \forall l = 1, \dots, N_{n+1}$$

for $\mu/\gamma^2 \omega$ sufficiently small. By (54), collecting the previous estimates,

$$|\omega^l - \omega_j^n(\mu', \omega')| > \frac{\gamma}{l^\tau} \quad \forall l = 1, \dots, N_{n+1}$$

and (μ', ω') belongs to A_{n+1} . \square

4.3. Measure of the Cantor set B_γ . In the following $R := (\mu', \mu'') \times (\omega', \omega'')$ denotes a rectangle contained in the region $\{(\mu, \omega) \in [\mu_1, \mu_2] \times (2\gamma, +\infty) : \mu < K'_6 \gamma^5 \omega\}$. Furthermore we consider $\omega'' - \omega'$ as a fixed quantity (“of order 1”).

LEMMA 13. *There exist K_6 and $K'_6 \leq K'_5$ such that, taking $C' \leq K'_6$ in the definition (15) of B_γ , for all $\mu \in (\mu_1, \mu_2)$ the section*

$$S_\gamma(\mu) := \{\omega : (\mu, \omega) \in B_\gamma\}$$

satisfies the measure estimate

$$(55) \quad |S_\gamma(\mu) \cap (\omega', \omega'')| \geq (1 - K_6 \gamma)(\omega'' - \omega').$$

As a consequence, for every $R := (\mu', \mu'') \times (\omega', \omega'')$

$$(56) \quad |B_\gamma \cap R| \geq |R| (1 - K_6 \gamma).$$

Proof. We consider the inequalities $|\omega l - \tilde{\omega}_j(\mu, \omega)| > 2\gamma/l^\tau$ in the definition of B_γ . The analogous inequalities for $|\omega l - \omega_j|$ and $|\omega l - j/c|$ are simpler because j/c and ω_j do not depend on (μ, ω) .

The complementary set we have to estimate is

$$C := \bigcup_{l, j \geq 1} \mathcal{R}_{lj},$$

where $\mathcal{R}_{lj} := \{\omega \in (\omega', \omega'') : |\omega l - \tilde{\omega}_j(\mu, \omega)| \leq 2\gamma/l^\tau\}$.

We claim that

$$(57) \quad |\partial_\omega \tilde{\omega}_j(\mu, \omega)| \leq \frac{K\mu}{\gamma^5 \omega}.$$

Indeed, by the same arguments as in the proof of Lemma 12 and the comparison principle (28), we have

$$|\tilde{\omega}_j(\mu, \omega) - \tilde{\omega}_j(\mu, \omega')| \leq K \|\tilde{w}(\mu, \omega) - \tilde{w}(\mu, \omega')\|_{\sigma_0/2} \leq \frac{K\mu}{\gamma^5 \omega} |\omega - \omega'|$$

using (51). As a consequence of (57)

$$\partial_\omega (l\omega - \tilde{\omega}_j(\mu, \omega)) \geq l - \frac{K\mu}{\gamma^5 \omega} \geq \frac{l}{2} \quad \forall l \geq 1$$

for $\mu/\gamma^5 \omega$ small enough; we deduce $|\mathcal{R}_{lj}| \leq 4\gamma/l^{\tau+1}$.

Furthermore the set \mathcal{R}_{lj} is nonempty only if

$$\omega' l - \frac{2\gamma}{l^\tau} < \tilde{\omega}_j(\mu, \omega) < \omega'' l + \frac{2\gamma}{l^\tau}.$$

So, for every fixed l , the number of indices j such that $\mathcal{R}_{lj} \neq \emptyset$ is

$$\#\{j\} \leq \frac{1}{\delta} \left(l(\omega'' - \omega') + \frac{4\gamma}{l^\tau} \right) + 1 \leq Kl(\omega'' - \omega'),$$

where

$$\delta := \inf \left\{ |\tilde{\omega}_{j+1}(\mu, \omega) - \tilde{\omega}_j(\mu, \omega)| : j \geq 1, (\mu, \omega) \in B_\gamma \right\}.$$

For $\|\tilde{w}\|_{\sigma_0/2} \leq K'\mu/\gamma\omega < R$ we have $\delta \geq \delta_1$, where

$$(58) \quad \delta_1 := \inf \left\{ |\omega_{j+1}(\mu, \omega) - \omega_j(\mu, \omega)| : j \geq 1, \mu \in [\mu_1, \mu_2], \|\omega\|_{\sigma_0/2} \leq R \right\} > 0,$$

as proved in the appendix.

In conclusion, the measure of the complementary set is

$$|\mathcal{C}| \leq \sum_{l=1}^{+\infty} \frac{4\gamma}{l^{\tau+1}} Kl(\omega'' - \omega') \leq K'(\omega'' - \omega')\gamma$$

and (55) is proved. Integrating on (μ', μ'') we obtain (56). \square

By Fubini's theorem also the section $S_\gamma(\omega)$ is large for ω in a large set.

LEMMA 14. *Let*

$$S_\gamma(\omega) := \{\mu : (\mu, \omega) \in B_\gamma\}.$$

For every $R := (\mu', \mu'') \times (\omega', \omega'')$, $\gamma' \in (0, 1)$ we obtain

$$(59) \quad \left| \left\{ \omega \in (\omega', \omega'') : \frac{|S_\gamma(\omega) \cap (\mu', \mu'')|}{\mu'' - \mu'} \geq 1 - \gamma' \right\} \right| \geq (\omega'' - \omega') \left(1 - K_6 \frac{\gamma}{\gamma'} \right).$$

Proof. Consider

$$\Omega^+ := \{ \omega \in (\omega', \omega'') : |S_\gamma(\omega) \cap (\mu', \mu'')| \geq (\mu'' - \mu')(1 - \gamma') \},$$

$$\Omega^- := \{ \omega \in (\omega', \omega'') : |S_\gamma(\omega) \cap (\mu', \mu'')| < (\mu'' - \mu')(1 - \gamma') \}.$$

Using Fubini's theorem

$$\begin{aligned}
 |B_\gamma \cap R| &= \int_{\omega'}^{\omega''} |S_\gamma(\omega) \cap (\mu', \mu'')| d\omega \\
 &= \int_{\Omega^+} |S_\gamma(\omega) \cap (\mu', \mu'')| d\omega + \int_{\Omega^-} |S_\gamma(\omega) \cap (\mu', \mu'')| d\omega \\
 (60) \qquad &\leq (\mu'' - \mu')|\Omega^+| + (\mu'' - \mu')(1 - \gamma')|\Omega^-|.
 \end{aligned}$$

By (56), $|B_\gamma \cap R| \geq (\omega'' - \omega')(\mu'' - \mu')(1 - K_6\gamma)$ and therefore, by (60),

$$(61) \qquad (\omega'' - \omega')(1 - K_6\gamma) \leq |\Omega^+| + (1 - \gamma')|\Omega^-| = (\omega'' - \omega') - \gamma'|\Omega^-|$$

because $|\Omega^+| + |\Omega^-| = \omega'' - \omega'$. Then

$$|\Omega^-| \leq (\omega'' - \omega')K_6 \frac{\gamma}{\gamma'}$$

and, by the first inequality in (61), $|\Omega^+| \geq (\omega'' - \omega')(1 - K_6\gamma/\gamma')$, which is (59). \square

Inequalities (55) and (59) imply the measure estimates (18)–(19).

The main conclusions of this section are summarized in the following theorem, which follows by Lemmas 11, 12, and 13.

THEOREM 3 (solution of the range equation). *There exist $\tilde{w} \in C^\infty(A_\gamma, W \cap X_{\sigma_0/2})$ satisfying (51) and the large (see (56)) Cantor set B_γ defined in (15) such that, for every $(\mu, \omega) \in B_\gamma$, the function $\tilde{w}(\mu, \omega)$ solves the range equation (30).*

5. Proof of Theorems 1 and 2.

Proof of Theorem 1. By Theorem 3 for all $(\mu, \omega) \in B_\gamma$ the function $\tilde{w}(\mu, \omega) \in X_{\sigma_0/2}$ solves the range equation (30). By Lemma 4, $v(\mu, \tilde{w}(\mu, \omega))$ solves the bifurcation equation in (8), and therefore

$$\tilde{u} := v(\mu, \tilde{w}(\mu, \omega)) + \tilde{w}(\mu, \omega) \in X_{\sigma_0/2}$$

is a solution of (3). Estimates (14) follow by (51).

Since \tilde{u} solves

$$(62) \qquad -(p(x)\tilde{u}_x)_x = \mu f(x, t, \tilde{u}) - \omega^2 \rho(x)\tilde{u}_{tt}$$

we deduce

$$-(p(x)\tilde{u}_x(t, x))_x \in H^1(0, \pi) \quad \forall t \in \mathbb{R}.$$

This implies, since $p(x) \in H^3(0, \pi)$, that

$$\tilde{u}(t, x) \in H^3(0, \pi) \cap H_0^1(0, \pi) \subset C^2(0, \pi) \quad \forall t \in \mathbb{R}. \quad \square$$

Proof of Theorem 2. For every fixed t , by the algebra property of H^m

$$\|f(x, t, u(x, t))\|_{H^m} \leq \sum_{l,k} \|f_{lk}(x)u^k(x)\|_{H^m} \leq K \sum_{l,k} \|f_{lk}\|_{H^m} \|u^k\|_{H^m}$$

for some $K > 0$.

Using the Gagliardo–Nirenberg-type inequality

$$\|u^k\|_{H^m} \leq (C_m \|u\|_{H^1})^{k-1} \|u\|_{H^m}$$

valid for every $u \in H_0^1 \cap H^m$ (see, e.g., [26, 20]), we get

$$(63) \quad \|f(x, t, u(x, t))\|_{H^m} \leq K \|u\|_{H^m} \sum_{l, k} \|f_{lk}\|_{H^m} (C_m \|u\|_{H^1})^{k-1},$$

which is convergent for $\|u\|_{H^1} < r_m/C_m$ by (21).

The solution \tilde{u} satisfies (62) and $\tilde{u}(\cdot, t) \in H^3(0, \pi)$ for all t .

By assumption $\|\tilde{u}\|_{H^1} < r_m/C_m$. By induction, assume $\tilde{u}(\cdot, t) \in H^k$ for $k = 3, \dots, m$. Hence $\tilde{u}_{tt}(\cdot, t) \in H^k$ and $\rho(x)\tilde{u}_{tt}(\cdot, t) \in H^k$ because $\rho \in H^m$. Furthermore, by (63), $f(x, t, \tilde{u}) \in H^k$. We deduce, by (62), that $p(x)\tilde{u}_x \in H^{k+1}$. Finally $\tilde{u} \in H^{k+2}$ because $p \in H^{m+1}$.

If $f_0(x, 0) = d_u f_0(x, 0) = 0$, then, by Lemma 3, we can take $v(\mu, 0) = 0$ for all μ . Therefore, by (14),

$$\|\tilde{u}(t, \cdot)\|_{H^1} \leq \|\tilde{u}\|_{\sigma_0/2} \leq \frac{2C\mu}{\gamma\omega} \quad \forall t,$$

and, for $\mu/\gamma\omega$ small enough, we deduce the regularity in (22). \square

6. Inversion of the linearized problem. Here we prove Lemma 7. Decomposing in Fourier series

$$f'(u) = \sum_{k \in \mathbb{Z}} a_k(x) e^{ikt}$$

we write, for all $h = \sum_{1 \leq |l| \leq N_n} h_l(x) e^{ilt} \in W^{(n)}$,

$$\begin{aligned} -L_\omega h + \mu P_n[f'(u)h] &= \sum_{1 \leq |l| \leq N_n} [\omega^2 l^2 \rho h_l + \partial_x(p \partial_x h_l)] e^{ilt} \\ &\quad + \mu P_n \left[\left(\sum_{k \in \mathbb{Z}} a_k e^{ikt} \right) \left(\sum_{1 \leq |l| \leq N_n} h_l e^{ilt} \right) \right] \\ &= \sum_{1 \leq |l| \leq N_n} [\omega^2 l^2 \rho h_l + \partial_x(p \partial_x h_l) + \mu a_0 h_l] e^{ilt} \\ &\quad + \mu \sum_{|l|, |k+l| \in \{1, \dots, N_n\}, k \neq 0} a_k h_l e^{i(k+l)t}. \end{aligned}$$

Hence $\mathcal{L}_n(w)$ defined in (32) can be decomposed as

$$(64) \quad \mathcal{L}_n(w)h = \rho(Dh + M_1h + M_2h),$$

where

$$\begin{aligned} Dh &:= \frac{1}{\rho} \sum_{|l|=1}^{N_n} [\omega^2 l^2 \rho h_l + (p h_l)' + \mu a_0 h_l] e^{ilt}, \\ (65) \quad M_1h &:= \frac{\mu}{\rho} \sum_{|l|, |k| \in \{1, \dots, N_n\}, l \neq k} a_{k-l} h_l e^{ikt}, \\ M_2h &:= \frac{\mu}{\rho} P_n[f'(u) d_w v(\mu, w)[h]]. \end{aligned}$$

Note that D is a diagonal operator in time Fourier basis. To study the eigenvalues of D , we use Sturm–Liouville-type techniques.

LEMMA 15 (Sturm–Liouville). *The eigenvalues $\lambda_j(\mu, w)$ of the Sturm–Liouville problem (27) form a strictly increasing sequence which tends to $+\infty$. Every $\lambda_j(\mu, w)$ is simple and the following asymptotic formula holds:*

$$(66) \quad \lambda_j(\mu, w) = \frac{j^2}{c^2} + b + M(\mu, w) + r_j(\mu, w), \quad |r_j(\mu, w)| \leq \frac{K}{j}$$

for all $j \geq 1$, $(\mu, w) \in [\mu_1, \mu_2] \times B_R$, where

$$c := \frac{1}{\pi} \int_0^\pi \left(\frac{\rho}{p}\right)^{1/2} dx, \quad b := \frac{1}{4\pi c} \int_0^\pi \left[\frac{(\rho p)'}{\rho \sqrt[4]{\rho p}}\right]' \frac{1}{\sqrt[4]{\rho p}} dx,$$

$$M(\mu, w) := -\frac{\mu}{c\pi} \int_0^\pi \frac{\Pi_V f'(v(\mu, w) + w)}{\sqrt{\rho p}} dx.$$

The eigenfunctions $\varphi_j(\mu, w)$ of (27) form an orthonormal basis of $L^2(0, \pi)$ with respect to the scalar product $(y, z)_{L^2_\rho} := c^{-1} \int_0^\pi yz\rho dx$. For K big enough

$$(y, z)_{\mu, w} := \frac{1}{c} \int_0^\pi p y' z' + [K\rho - \mu \Pi_V f'(v(\mu, w) + w)] yz dx$$

defines an equivalent scalar product on $H_0^1(0, \pi)$ and

$$(67) \quad K' \|y\|_{H^1} \leq \|y\|_{\mu, w} \leq K'' \|y\|_{H^1} \quad \forall y \in H_0^1.$$

$\varphi_j(\mu, w)$ is also an orthogonal basis of $H_0^1(0, \pi)$ with respect to the scalar product $(\cdot, \cdot)_{\mu, w}$ and, for $y = \sum_{j \geq 1} \hat{y}_j \varphi_j(\mu, w)$,

$$(68) \quad \|y\|_{L^2_\rho}^2 = \sum_{j \geq 1} \hat{y}_j^2, \quad \|y\|_{\mu, w}^2 = \sum_{j \geq 1} \hat{y}_j^2 (\lambda_j(\mu, w) + K).$$

Proof. For the proof of the lemma, see the appendix. \square

We develop

$$Dh = \sum_{1 \leq |l| \leq N_n} D_l h_l e^{ilt},$$

where

$$D_l z := \frac{1}{\rho} [\omega^2 l^2 \rho z + (p z')' + \mu a_0 z] \quad \forall z \in H_0^1(0, \pi)$$

and $a_0 = \Pi_V f(v(\mu, w) + w)$.

By Lemma 15 each D_l is diagonal with respect to the basis $\varphi_j(\mu, w)$:

$$z = \sum_{j=1}^{+\infty} \hat{z}_j \varphi_j(\mu, w) \in H_0^1(0, \pi) \Rightarrow D_l z = \sum_{j=1}^{+\infty} (\omega^2 l^2 - \lambda_j(\mu, w)) \hat{z}_j \varphi_j(\mu, w).$$

LEMMA 16. *Suppose all the eigenvalues $\omega^2 l^2 - \lambda_j(\mu, w)$ are not zero. Then*

$$|D_l|^{-1/2} z := \sum_{j=1}^{+\infty} \frac{\hat{z}_j \varphi_j(\mu, w)}{\sqrt{|\omega^2 l^2 - \lambda_j(\mu, w)|}}$$

satisfies

$$(69) \quad \left\| |D_l|^{-1/2} z \right\|_{H^1} \leq \frac{K}{\sqrt{\alpha_l}} \|z\|_{H^1} \quad \forall z \in H_0^1(0, \pi),$$

where $\alpha_l := \min_{j \geq 1} |\omega^2 l^2 - \lambda_j(\mu, w)| > 0$.

Proof. By (68) $\| |D_l|^{-1/2} z \|_{\mu, w}^2 \leq (1/\alpha_l) \|z\|_{\mu, w}^2$. Hence (69) follows by the equivalence of the norms (67). \square

LEMMA 17 (inversion of D). *Assume the nonresonance condition (34). Then $|D|^{-1/2} : W^{(n)} \rightarrow W^{(n)}$ defined by*

$$|D|^{-1/2} h := \sum_{1 \leq |l| \leq N_n} |D_l|^{-1/2} h_l e^{ilt}$$

satisfies

$$\| |D|^{-1/2} h \|_{\sigma, s} \leq \frac{K}{\sqrt{\gamma\omega}} \|h\|_{\sigma, s + \frac{\tau-1}{2}} \leq \frac{KN_n^{\frac{\tau-1}{2}}}{\sqrt{\gamma\omega}} \|h\|_{\sigma, s} \quad \forall h \in W^{(n)}.$$

Proof. By (69) and $\alpha_{-l} = \alpha_l \geq \gamma\omega/|l|^{\tau-1}$

$$\begin{aligned} \| |D|^{-1/2} h \|_{\sigma, s}^2 &= \sum_{1 \leq |l| \leq N_n} \| |D_l|^{-1/2} h_l \|_{H^1}^2 (1 + l^{2s}) e^{2\sigma|l|} \\ &\leq \sum_{1 \leq |l| \leq N_n} \frac{K^2 |l|^{\tau-1}}{\gamma\omega} \|h_l\|_{H^1}^2 (1 + l^{2s}) e^{2\sigma|l|} \\ &\leq \frac{K'}{\gamma\omega} \|h\|_{\sigma, s + \frac{\tau-1}{2}}^2 \end{aligned}$$

because $|l|^{\tau-1}(1 + l^{2s}) < 2(1 + |l|^{2s+\tau-1})$ for all $|l| \geq 1$. \square

To prove the invertibility of $\mathcal{L}_n(w)$ we write (64) as

$$(70) \quad \mathcal{L}_n(w) = \rho |D|^{1/2} (U + T_1 + T_2) |D|^{1/2},$$

where

$$(71) \quad \begin{cases} U := |D|^{-1/2} D |D|^{-1/2}, \\ T_i := |D|^{-1/2} M_i |D|^{-1/2}, \quad i = 1, 2. \end{cases}$$

With respect to the basis $\varphi_j(\mu, w) e^{ilt}$ the operator U is diagonal and its (l, j) th eigenvalue is $\text{sign}(\omega^2 l^2 - \lambda_j(\mu, w)) \in \{\pm 1\}$, implying that the operator norm is

$$(72) \quad \|U\|_{\sigma} := \sup_{\|h\|_{\sigma} \leq 1} \|Uh\|_{\sigma} = 1.$$

The smallness of T_1 requires an analysis of the small divisors. Formula (66) implies, by Taylor expansion, the asymptotic dispersion relation

$$(73) \quad \left| \omega_j(\mu, w) - \frac{j}{c} \right| \leq \frac{K}{j},$$

and there exists K such that, for every $x \geq 0$,

$$(74) \quad |x^2 - \lambda_{j^*}(\mu, w)| = \min_{j \geq 1} |x^2 - \lambda_j(\mu, w)| \Rightarrow j^* \geq Kx.$$

LEMMA 18 (analysis of the small divisors). *Assume the nonresonance conditions (33)–(34) and $\omega > \gamma$. Then for all $|k|, |l| \in \{1, \dots, N_n\}$, $k \neq l$,*

$$\alpha_l \alpha_k \geq \left(\frac{K\gamma^3\omega}{|k-l|^{\frac{\tau(\tau-1)}{2-\tau}}} \right)^2,$$

where $\alpha_l := \min_{j \geq 1} |\omega^2 l^2 - \lambda_j(\mu, w)|$.

Proof. Since $\alpha_{-l} = \alpha_l$ for all l , we can suppose $l, k \geq 1$.

We distinguish two cases, if k, l are close to or far from each other. Let $\beta := (2 - \tau)/\tau \in (0, 1)$.

Case 1. Let $2|k - l| > (\max\{k, l\})^\beta$. By (34)

$$\alpha_k \alpha_l \geq \frac{(\gamma\omega)^2}{(kl)^{\tau-1}} \geq \frac{(\gamma\omega)^2}{(\max\{k, l\})^{2(\tau-1)}} \geq \frac{C(\gamma\omega)^2}{|k-l|^{\frac{2(\tau-1)}{\beta}}}.$$

Case 2. Let $0 < 2|k - l| \leq (\max\{k, l\})^\beta$. In this case $2k \geq l \geq k/2$. Indeed, if $k > l$, then $2(k - l) \leq k^\beta$, so $2l \geq 2k - k^\beta \geq k$ because $\beta \in (0, 1)$ —analogously if $l > k$.

Let i , resp., j , be an integer which realizes the minimum α_k , resp., α_l , and write in short $\lambda_j(\mu) := \lambda_j(\mu, w)$, $\omega_j(\mu) := \omega_j(\mu, w)$.

If both $\lambda_i(\mu), \lambda_j(\mu) \leq 0$, then $\alpha_l \geq \omega^2 l^2$, $\alpha_k \geq \omega^2 k^2$, $\alpha_l \alpha_k \geq \omega^4 > \gamma^2 \omega^2$.

If only $\lambda_j(\mu) \leq 0$, then $\alpha_l \alpha_k \geq \gamma \omega^3 l^2 / k^{\tau-1} \geq 2^{1-\tau} \gamma \omega^3 \geq 2^{1-\tau} \gamma^2 \omega^2$.

The really resonant cases happen if $\lambda_i(\mu), \lambda_j(\mu) > 0$. Suppose, for example, that $\max\{k, l\} = k$. By (73), $|\omega_j(\mu) - (j/c)| \leq K/j$, and, by (74), $i \geq K\omega k$, $j \geq K\omega l$. Hence, using also (33),

$$\begin{aligned} |(\omega k - \omega_i(\mu)) - (\omega l - \omega_j(\mu))| &= |\omega(k - l) - (\omega_i(\mu) - \omega_j(\mu))| \\ &\geq \left| \omega(k - l) - \frac{i - j}{c} \right| - \frac{K}{\omega l} - \frac{K}{\omega k} \\ &\geq \frac{\gamma}{(k - l)^\tau} - \frac{3K}{\omega k} \geq \frac{2^\tau \gamma}{k^{\beta\tau}} - \frac{3K}{\omega k} \end{aligned}$$

because $2(k - l) \leq k^\beta$, $2l \geq k$. Since $\beta\tau < 1$ and $k \leq 2l$,

$$|(\omega k - \omega_i(\mu)) - (\omega l - \omega_j(\mu))| \geq \frac{1}{2} \left(\frac{\gamma}{k^{\beta\tau}} + \frac{\gamma}{l^{\beta\tau}} \right) \quad \forall k \geq \left(\frac{K}{\omega\gamma} \right)^{\frac{1}{1-\beta\tau}} =: k^*.$$

We reach the same conclusion if $\max\{k, l\} = l$. It follows that, for $\max\{k, l\} \geq k^*$, there holds $|\omega k - \omega_i(\mu)| \geq \gamma/2k^{\beta\tau}$ or $|\omega l - \omega_j(\mu)| \geq \gamma/2l^{\beta\tau}$. Suppose $|\omega k - \omega_i(\mu)| \geq \gamma/2k^{\beta\tau}$. Then

$$\alpha_k = |\omega^2 k^2 - \omega_i^2(\mu)| \geq |\omega k - \omega_i(\mu)| \omega k \geq \frac{\gamma\omega}{2} k^{1-\beta\tau}.$$

Since $l \leq 2k$, for α_l we can use (34),

$$\alpha_k \alpha_l \geq \frac{\gamma\omega k^{1-\beta\tau}}{2} \frac{\gamma\omega}{l^{\tau-1}} \geq \frac{\gamma^2\omega^2}{2^\tau} k^{2-\tau-\beta\tau} = \frac{\gamma^2\omega^2}{2^\tau}$$

because $2 - \tau - \beta\tau = 0$.

On the other hand, if $\max\{k, l\} < k^* = (K/\omega\gamma)^{1/(\tau-1)}$, we can use (34) for both k, l :

$$\alpha_k \alpha_l \geq \frac{(\gamma\omega)^2}{(kl)^{\tau-1}} > \frac{(\gamma\omega)^2}{(k^*)^{2(\tau-1)}} = (\gamma\omega)^2 \left(\frac{\omega\gamma}{K}\right)^{\frac{1}{\tau-1} 2(\tau-1)} > \frac{\gamma^6 \omega^2}{K^2}$$

(using $\omega > \gamma$). Since $\gamma < 1$, taking the minimum for all these cases concludes the proof. \square

LEMMA 19 (estimate of T_1). *Assume the nonresonance conditions (33)–(34), $\omega > \gamma$, and $\Pi_W f'(u) = \sum_{l \neq 0} a_l(x) e^{ilt} \in X_{\sigma, 1 + \frac{\tau(\tau-1)}{2-\tau}}$. There exists K such that*

$$\|T_1 h\|_\sigma \leq \frac{K\mu}{\gamma^3 \omega} \|\Pi_W f'(u)\|_{\sigma, 1 + \frac{\tau(\tau-1)}{2-\tau}} \|h\|_\sigma \quad \forall h \in W^{(n)}.$$

Proof. For all $h \in W^{(n)}$, $T_1 h = \sum_{1 \leq |k| \leq N_n} (T_1 h)_k e^{ikt}$, where

$$\begin{aligned} (T_1 h)_k &= |D_k|^{-1/2} (M_1 |D|^{-1/2} h)_k \\ &= |D_k|^{-1/2} \left[\sum_{1 \leq |l| \leq N_n, l \neq k} \mu \frac{a_{k-l}}{\rho} |D_l|^{-1/2} h_l \right]. \end{aligned}$$

Setting $A_m := \|a_m/\rho\|_{H^1}$, using (69) and Lemma 18, we obtain

$$(75) \quad \|(T_1 h)_k\|_{H^1} \leq K\mu \sum_{1 \leq |l| \leq N_n, l \neq k} \frac{A_{k-l}}{\sqrt{\alpha_k} \sqrt{\alpha_l}} \|h_l\|_{H^1} \leq \frac{K\mu}{\gamma^3 \omega} S_k,$$

where

$$S_k := \sum_{|l| \leq N_n, l \neq k} A_{k-l} |k-l|^{\frac{\tau(\tau-1)}{2-\tau}} \|h_l\|_{H^1}.$$

By (75) we get, defining $S(t) := \sum_{|k|=1}^{N_n} S_k e^{ikt}$,

$$\begin{aligned} \|T_1 h\|_\sigma^2 &= \sum_{|k|=1}^{N_n} \|(T_1 h)_k\|_{H^1}^2 (1+k^2) e^{2\sigma|k|} \\ &\leq \left(\frac{K\mu}{\gamma^3 \omega}\right)^2 \sum_{|k|=1}^{N_n} S_k^2 (1+k^2) e^{2\sigma|k|} = \left(\frac{K\mu}{\gamma^3 \omega}\right)^2 \|S\|_\sigma^2. \end{aligned}$$

Since $S = P_n(\varphi\psi)$ with $\varphi(t) := \sum_{l \in \mathbb{Z}} A_l |l|^{\frac{\tau(\tau-1)}{2-\tau}} e^{ilt}$ and $\psi(t) := \sum_{|l|=1}^{N_n} \|h_l\|_{H^1} e^{ilt}$

$$\|T_1 h\|_\sigma \leq \frac{K\mu}{\gamma^3 \omega} \|\varphi\|_\sigma \|\psi\|_\sigma \leq \frac{K\mu}{\gamma^3 \omega} \left\| \Pi_W f'(u) \right\|_{\sigma, 1 + \frac{\tau(\tau-1)}{2-\tau}} \|h\|_\sigma$$

because $\|\varphi\|_\sigma \leq 2 \|\Pi_W f'(u)\|_{\sigma, 1 + \frac{\tau(\tau-1)}{2-\tau}}$ and $\|\psi\|_\sigma = \|h\|_\sigma$. \square

LEMMA 20 (estimate of T_2). *Suppose that $\Pi_W f'(u) \in X_{\sigma, 1 + \frac{\tau-1}{2}}$. Then*

$$\|T_2 h\|_\sigma \leq \frac{K\mu}{\gamma\omega} \|\Pi_W f'(u)\|_{\sigma, 1 + \frac{\tau-1}{2}} \|h\|_\sigma \quad \forall h \in W^{(n)}$$

for some K .

Proof. By the definitions (71) and (65) and by Lemma 17,

$$\begin{aligned} \|T_2 h\|_\sigma &\leq \frac{K}{\sqrt{\gamma\omega}} \|M_2 |D|^{-1/2} h\|_{\sigma, 1+\frac{\tau-1}{2}} \\ &\leq \frac{K'\mu}{\sqrt{\gamma\omega}} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau-1}{2}} \|d_w v(\mu, w)[|D|^{-1/2} h]\|_{\sigma, 1+\frac{\tau-1}{2}} \\ &= \frac{K'\mu}{\sqrt{\gamma\omega}} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau-1}{2}} \|d_w v(\mu, w)[|D|^{-1/2} h]\|_{H^1} \end{aligned}$$

because $d_w v(\mu, w)[|D|^{-1/2} h] \in V$. By Lemmas 4 and 17

$$\|d_w v(\mu, w)[|D|^{-1/2} h]\|_{H^1} \leq K \| |D|^{-1/2} h \|_{\sigma, 1-\frac{\tau-1}{2}} \leq \frac{K}{\sqrt{\gamma\omega}} \|h\|_{\sigma, 1},$$

implying the thesis. \square

Proof of Lemma 7. By (72), $\|U\|_\sigma = 1$. If

$$(76) \quad \|T_1 + T_2\|_\sigma < \frac{1}{2},$$

then by Neumann series $U + T_1 + T_2$ is invertible in $(W^{(n)}, \|\cdot\|_\sigma)$ and

$$\|(U + T_1 + T_2)^{-1}\|_\sigma < 2.$$

By Lemmas 19 and 20, condition (76) is verified if

$$(77) \quad \|T_1\|_\sigma \leq \frac{K\mu}{\gamma^3\omega} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau(\tau-1)}{2-\tau}} < \frac{1}{4}$$

and

$$(78) \quad \|T_2\|_\sigma \leq \frac{K\mu}{\gamma\omega} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau-1}{2}} \leq \frac{K\mu}{\gamma^3\omega} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau(\tau-1)}{2-\tau}} < \frac{1}{4}$$

(we recall that $\gamma \in (0, 1)$ and $(\tau - 1)/2 < \tau(\tau - 1)/(2 - \tau)$ because $\tau > 1$). Both conditions (77) and (78) are satisfied if

$$\frac{\mu}{\gamma^3\omega} \|\Pi_W f'(u)\|_{\sigma, 1+\frac{\tau(\tau-1)}{2-\tau}} < \frac{1}{4K} =: K'_1,$$

which is condition (35). Hence, inverting (70)

$$\mathcal{L}_n(w)^{-1} h = |D|^{-1/2} (U + T_1 + T_2)^{-1} |D|^{-1/2} \left(\frac{h}{\rho}\right),$$

which, using Lemma 17, yields (36). \square

7. Appendix.

Proof of Lemma 15. Let $a(x) \in L^2(0, \pi)$. Under the ‘‘Liouville change of variable’’

$$(79) \quad x = \psi(\xi) \Leftrightarrow \xi = g(x), \quad g(x) := \frac{1}{c} \int_0^x \left(\frac{\rho(s)}{p(s)}\right)^{1/2} ds,$$

we have that $(\lambda, y(x))$ satisfies

$$(80) \quad \begin{cases} -(p(x)y'(x))' + a(x)y(x) = \lambda\rho(x)y(x), \\ y(0) = y(\pi) = 0 \end{cases}$$

iff $(\nu, z(\xi))$ satisfies

$$(81) \quad \begin{cases} -z''(\xi) + [q(\xi) + \alpha(\xi)]z(\xi) = \nu z(\xi), \\ z(0) = z(\pi) = 0, \end{cases}$$

where

$$\begin{aligned} \nu &= c^2\lambda, & r(x) &= \sqrt[4]{p(x)\rho(x)}, & z(\xi) &= y(\psi(\xi))r(\psi(\xi)), \\ \alpha(\xi) &= c^2 \frac{a(\psi(\xi))}{\rho(\psi(\xi))}, & q(\xi) &= c^2 Q(\psi(\xi)), & Q &= \frac{p}{\rho} \frac{r''}{r} + \frac{1}{2} \left(\frac{p}{\rho} \right)' \frac{r'}{r}. \end{aligned}$$

By [25, Theorem 4 in Chapter 2, p. 35], the eigenvalues of (81) form an increasing sequence ν_j satisfying the asymptotic expansion

$$\nu_j = j^2 + \frac{1}{\pi} \int_0^\pi (q + \alpha) d\xi - \frac{1}{\pi} \int_0^\pi \cos(2j\xi)(q(\xi) + \alpha(\xi)) d\xi + r_j, \quad |r_j| \leq \frac{C}{j},$$

where $C := C(\|q + \alpha\|_{L^2})$ is a positive constant. Moreover every ν_j is simple [25, Theorem 2, p. 30].

Since p, ρ are positive and belong to H^3 , if $a \in H^1$, then $q, \alpha \in H^1$. Integrating by parts, $|\int_0^\pi \cos(2j\xi)(q + \alpha) d\xi| \leq \|q + \alpha\|_{H^1}/j$ and so

$$\nu_j = j^2 + \frac{1}{\pi} \int_0^\pi (q + \alpha) d\xi + r'_j, \quad |r'_j| \leq \frac{C'}{j}$$

for some $C' := C'(\|q + \alpha\|_{H^1})$. Dividing by c^2 and using the inverse Liouville change of variable, we obtain the formula for the eigenvalues $\lambda_j(a)$ of (80),

$$(82) \quad \lambda_j(a) = \frac{j^2}{c^2} + \frac{1}{\pi c} \int_0^\pi \frac{Q\sqrt{\rho}}{\sqrt{p}} dx + \frac{1}{\pi c} \int_0^\pi \frac{a}{\sqrt{\rho p}} dx + r_j(a), \quad |r_j(a)| \leq \frac{C}{j}$$

for some $C(\rho, p, \|a\|_{H^1}) > 0$. Formula (66) follows for $a(x) = -\mu \Pi_V f'(v(\mu, w) + w)(x)$ and some algebra.

By [25, Theorem 7, p. 43], the eigenfunctions of (81) form an orthonormal basis for L^2 . Applying the Liouville change of variable (79) in the integrals, the eigenfunctions $\varphi_j(a)$ of (80) form an orthonormal basis for L^2 with respect to the scalar product $(\cdot, \cdot)_{L^2_\rho}$.

Finally, since $\varphi_j := \varphi_j(a)$ solves

$$-(p\varphi_j')' + (K\rho + a)\varphi_j = (\lambda_j(a) + K)\rho\varphi_j,$$

multiplying by φ_i and integrating by parts gives

$$(\varphi_j, \varphi_i)_{\mu, w} = \delta_{i, j}(\lambda_j(a) + K),$$

and (68) follows (note that $\lambda_j(a) + K > 0$ for all j and for K large enough). \square

Proof of Lemma 5. Let $a, b \in H^1(0, \pi)$ and consider $\alpha := c^2 a(\psi)/\rho(\psi)$, $\beta := c^2 b(\psi)/\rho(\psi)$ constructed as above via the Liouville change of variable (79). By [25, p. 34], for every j

$$(83) \quad |\lambda_j(a) - \lambda_j(b)| = \frac{1}{c^2} |\nu_j(\alpha) - \nu_j(\beta)| \leq \frac{1}{c^2} \|\alpha - \beta\|_\infty \leq K \|a - b\|_{H^1},$$

and (28) follows by the mean value theorem because $\mu \Pi_V f(v(\mu, w) + w)$ has bounded derivatives on bounded sets. \square

Proof of (58). By the asymptotic formula (73)

$$\min_{j \geq 1} |\omega_{j+1}(\mu, w) - \omega_j(\mu, w)| \geq \frac{1}{c} - \frac{2K}{j} > \frac{1}{2c}$$

if $j > 4Kc$, uniformly in $\mu \in [\mu_1, \mu_2]$, $w \in B_R$. For $1 \leq j \leq 4Kc$ the minimum

$$m_j := \min_{(\mu, w) \in [\mu_1, \mu_2] \times B_R} |\omega_{j+1}(\mu, w) - \omega_j(\mu, w)|$$

is attained because $a \mapsto \lambda_j(a)$ is a compact function on H^1 by the compact embedding $H^1(0, \pi) \hookrightarrow L^\infty(0, \pi)$ and by (83) (see also [25, Theorem 3, pp. 31 and 34]). Each $m_j > 0$ because all the eigenvalues λ_j are simple. \square

REFERENCES

- [1] P. ACQUISTAPACE, *Soluzioni periodiche di un'equazione iperbolica non lineare*, Boll. Un. Mat. Ital. B (5), 13 (1976), pp. 760–777.
- [2] A. BAMBERGER, G. CHAVENT, AND P. LAILLY, *About the stability of the inverse problem in 1D wave equations—Applications to the interpretation of seismic profiles*, Appl. Math. Optim., 5 (1979), pp. 1–47.
- [3] V. BARBU AND N. H. PAVEL, *Periodic solutions to nonlinear one dimensional wave equation with x -dependent coefficients*, Trans. Amer. Math. Soc., 349 (1997), pp. 2035–2048.
- [4] T. BARTSCH, Y. H. DING, AND C. LEE, *Periodic solutions of a wave equation with concave and convex nonlinearities*, J. Differential Equations, 153 (1999), pp. 121–141.
- [5] J. BERKOVITS AND J. MAWHIN, *Diophantine approximation, Bessel functions and radially symmetric periodic solutions of semilinear wave equations in a ball*, Trans. Amer. Math. Soc., 353 (2001), pp. 5041–5055.
- [6] M. BERTI, *Nonlinear Oscillations of Hamiltonian PDEs*, Progr. Nonlinear Differential Equations Appl. 74, H. Brézis, ed., Birkhäuser, Boston, 2008.
- [7] M. BERTI AND L. BIASCO, *Forced vibrations of wave equations with non-monotone nonlinearities*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 23 (2006), pp. 437–474.
- [8] M. BERTI AND P. BOLLE, *Cantor families of periodic solutions for completely resonant nonlinear wave equations*, Duke Math. J., 134 (2006), pp. 359–419.
- [9] M. BERTI AND P. BOLLE, *Cantor families of periodic solutions for wave equations via a variational principle*, Adv. Math., 217 (2008), pp. 1671–1727.
- [10] M. BERTI AND M. PROCESI, *Quasi-periodic solutions of completely resonant forced wave equations*, Commun. Partial Differential Equations, 31 (2006), pp. 959–985.
- [11] J. BOURGAIN, *Construction of quasi-periodic solutions for Hamiltonian perturbations of linear equations and applications to nonlinear PDE*, Int. Math. Res. Not., 11 (1994), pp. 475–497.
- [12] J. BOURGAIN, *Periodic solutions of nonlinear wave equations*, in Harmonic Analysis and Partial Differential Equations, Chicago Lectures in Math., University of Chicago Press, Chicago, IL, 1999, pp. 69–97.
- [13] J. BOURGAIN, *Nonlinear Schrödinger equations*, in Hyperbolic Equations and Frequency Interactions (Park City, UT, 1995), IAS/Park City Math. Ser. 5, AMS, Providence, RI, 1999, pp. 3–157.
- [14] H. BRÉZIS AND L. NIRENBERG, *Forced vibrations for a nonlinear wave equation*, Comm. Pure Appl. Math., 31 (1978), pp. 1–30.
- [15] R. COURANT AND D. HILBERT, *Methods of Mathematical Physics*, Vol. I, Interscience Publishers, New York, 1953.

- [16] W. CRAIG, *Problèmes de Petits Diviseurs dans les Équations aux Dérivées Partielles*, Panor Synthèses 9, Société Mathématique de France, Paris, 2000.
- [17] W. CRAIG AND E. WAYNE, *Newton's method and periodic solutions of nonlinear wave equations*, Comm. Pure Appl. Math., 46 (1993), pp. 1409–1498.
- [18] R. DE LA LLAVE, *Variational methods for quasi-periodic solutions of partial differential equations*, in Hamiltonian Systems and Celestial Mechanics (Pátzcuaro, 1998), World Sci. Monogr. Ser. Math. 6, World Scientific, River Edge, NJ, 2000, pp. 214–228.
- [19] J.-M. FOKAM, *Forced Vibrations via Nash-Moser Iteration*, Ph.D. thesis, University of Texas, Austin, TX.
- [20] A. FRIEDMAN, *Partial Differential Equations*, Robert E. Krieger, Huntington, NY, 1976.
- [21] S. KUKSIN, *Hamiltonian perturbations of infinite-dimensional linear systems with imaginary spectrum*, Funktsional. Anal. i Prilozhen., 21 (1987), pp. 22–37.
- [22] S. KUKSIN, *Analysis of Hamiltonian PDEs*, Oxford Lecture Ser. Math. Appl. 19, Oxford University Press, New York, 2000.
- [23] P. J. MCKENNA, *On solutions of a nonlinear wave question when the ratio of the period to the length of the intervals is irrational*, Proc. Amer. Math. Soc., 93 (1985), pp. 59–64.
- [24] P. I. PLOTNIKOV AND L. N. YUNGERMAN, *Periodic solutions of a weakly nonlinear wave equation with an irrational relation of period to interval length*, Differential Equations, 24 (1988), pp. 1059–1065.
- [25] J. PÖSCHEL AND E. TRUBOWITZ, *Inverse Spectral Theory*, Academic Press, Orlando, FL, 1987.
- [26] P. RABINOWITZ, *Periodic solutions of nonlinear hyperbolic partial differential equations*, Comm. Pure Appl. Math., 20 (1967), pp. 145–205.
- [27] P. RABINOWITZ, *Time periodic solutions of nonlinear wave equations*, Manuscripta Math., 5 (1971), pp. 165–194.
- [28] P. RABINOWITZ, *Free vibration of a semilinear wave equation*, Comm. Pure Appl. Math, 31 (1978), pp. 31–68.
- [29] I. A. RUDAKOV, *Periodic solutions of a nonlinear wave equation with nonconstant coefficients*, Math. Notes, 76 (2004), pp. 395–406.
- [30] W. M. SCHMIDT, *Diophantine Approximation*, Lecture Notes in Math. 785, Springer-Verlag, Berlin, 1980.
- [31] K. TANAKA, *Infinitely many periodic solutions for the equation: $u_{tt} - u_{xx} \pm |u|^{p-1}u = f(t, x)$* , Trans. Amer. Math. Soc., 307 (1988), pp. 615–645.
- [32] E. WAYNE, *Periodic and quasi-periodic solutions of nonlinear wave equations via KAM theory*, Comm. Math. Phys., 127 (1990), pp. 479–528.

NORMAL FORMS, QUASI-INVARIANT MANIFOLDS, AND BIFURCATIONS OF NONLINEAR DIFFERENCE-ALGEBRAIC EQUATIONS*

R. BEARDMORE[†] AND K. WEBSTER[†]

Abstract. We study the existence of quasi-invariant manifolds in a neighborhood of a fixed point of the *difference*-algebraic equation (Δ AE) $F(z_n, z_{n+1}) = 0$, where $F : \mathbb{R}^{2m} \rightarrow \mathbb{R}^m$ is a smooth map satisfying $F(0, 0) = 0$. We demonstrate the existence of quasi-invariant manifolds on which one can define forward and backward orbits of the Δ AE under mild assumptions on its linearization at the fixed point $z = 0$. Indeed, by assuming this linearization to be a regular matrix pencil, one obtains a functional equation satisfied by invariant manifolds which can be solved using an extension of the contraction mapping to spaces that satisfy an interpolation property. If the Δ AE under study is permitted to depend smoothly on a parameter, we then obtain a Neimark–Sacker bifurcation theorem as a corollary that can be deduced from the existence of a normal form for nonlinear Δ AEs.

Key words. invariant manifolds, bifurcations, difference-algebraic equations

AMS subject classifications. 39B72, 39A11, 37G99

DOI. 10.1137/050638618

1. Introduction. The purpose of this paper is to provide an analysis of the invariant manifolds and bifurcations found in a class of *difference*-algebraic equations (Δ AEs) of the form

$$(1.1) \quad F(z_n, z_{n+1}) = 0;$$

the nomenclature and chosen acronym for (1.1) have been taken from [4, 22]. We assume that $F(= F(z, \bar{z})) : \mathbb{R}^{2m} \rightarrow \mathbb{R}^m$ is a smooth map satisfying $F(0, 0) = 0$ and say that (1.1) is *singular* because the partial derivative $d_{\bar{z}}F(0, 0)$ is not an isomorphism from \mathbb{R}^m to itself. The purpose of the first part of this paper is to provide conditions under which (1.1) has suitably defined invariant manifolds that contain the fixed point $z = 0$, where the main difficulty to overcome in this analysis is the fact that forward orbits of (1.1) are not necessarily uniquely defined in a neighborhood of the fixed point.

The second part of the paper utilizes the existence of the aforementioned invariant manifolds to investigate the presence of bifurcations in Δ AEs in the sense that by extending F to be a C^k -mapping of the form $F : \mathbb{R}^{2m} \times \mathbb{R} \rightarrow \mathbb{R}^{2m}$, where $k \geq 5$, we examine the structure of invariant sets in the family of Δ AEs

$$(1.2) \quad F(z_n, z_{n+1}, \mu) = 0.$$

Our rationale is taken from bifurcation theory for maps which leads to the following question. If $F(0, 0, \mu) = 0$ for all μ in some interval and the one-parameter family of matrix pencils

$$\mathcal{P}(\mu) := (A(\mu), B(\mu)) := (d_{\bar{z}}F(0, 0, \mu), d_zF(0, 0, \mu)),$$

*Received by the editors August 22, 2005; accepted for publication (in revised form) August 28, 2007; published electronically May 7, 2008.

<http://www.siam.org/journals/sima/40-1/63861.html>

[†]Department of Mathematics, Imperial College, 180 Queen's Gate, London, SW7 2AZ, United Kingdom (r.beardmore@ic.ac.uk, k.webster@ic.ac.uk). The first author was supported by Nuffield Foundation grant NAL/00511/G. The second author was supported by EPSRC grant GR/S/17215/01.

where F has (z, \bar{z}, μ) as its argument, is such that $\mathcal{P}(\mu_0)$ has a finite eigenvalue of unit modulus in the complex plane, does an invariant set of (1.2) bifurcate from the fixed point $z = 0$ at $\mu = \mu_0$?

Due to the lack of forward uniqueness we modify what we mean by the term *invariant*, which we do by using the prefixed *quasi-invariant*, and say that a set $\mathcal{Q} \subset F^{-1}\{0\}$ is quasi-invariant for (1.1) if, for $(z, \bar{z}) \in \mathcal{Q}$, there is a subsequent iterate $(\bar{z}, \bar{\bar{z}}) \in \mathcal{Q}$ for some $\bar{\bar{z}} \in \mathbb{R}^m$.

The paper is organized in the following way: The remainder of section 1 briefly covers the linear prerequisites for (1.1). Section 2 then provides some motivating applications. Section 3 presents the basic definitions of how (1.1) defines a local dynamical system and gives the invariant manifold equation of fixed points of (1.1). Section 4 provides a reformulation of the invariant manifold equation from section 3 as a nonlinear fixed-point problem in suitable Banach spaces which then is shown to have a solution in section 4.3. Section 5 gives a normal form for (1.1) with and without the presence of a bifurcation parameter. This section concludes with theorems that can be deduced using these normal forms, giving bifurcation results for (1.1) when that parameter is included. Finally, section 6 finishes the paper with a series of examples.

1.1. The linear case: Kronecker normal form. As a precursor to the analysis of the nonlinear problem (1.1), consider the linear case

$$(1.3) \quad Bz_n + Az_{n+1} = 0,$$

where $A, B : \mathbb{R}^m \rightarrow \mathbb{R}^m$ are linear maps and A is singular. In order to discuss the behavior of (1.3) we first introduce a normal form for matrix pencils.

When A is singular, the matrix pencil (A, B) is said to be *regular* if there is an $\omega \in \mathbb{C}$ such that $\det(\omega A + B) \neq 0$. The following result is well known for regular matrix pencils (see [7, 3]): There are complementary subspaces $K_1 \simeq \mathbb{R}^p, K_2 \simeq \mathbb{R}^q \subset \mathbb{R}^m$ such that $p + q = m$ and nonsingular linear mappings P, Q on \mathbb{R}^m , $L : K_1 \rightarrow K_1$, and $N : K_2 \rightarrow K_2$ such that

$$(1.4) \quad PAQ = \begin{pmatrix} I_p & 0 \\ 0 & N \end{pmatrix}, \quad PBQ = \begin{pmatrix} L & 0 \\ 0 & I_q \end{pmatrix};$$

I_p and I_q are identities on K_1 and K_2 , respectively. Moreover, there is a $\nu \geq 1$ such that $N^\nu = 0$, and ν is said to be the Kronecker index of (A, B) .

The Kronecker normal form (KNF) in (1.4) can be used to rewrite (1.3) as a coupled system of difference equations

$$(1.5) \quad Lu_n + u_{n+1} = 0, \quad v_n + Nv_{n+1} = 0,$$

which has the solution $u_n = (-L)^n u_0$ and $v_n \equiv 0$ for all n , and thus (1.3) has a quasi-invariant subspace that arises from the quasi-invariant space $\{(u, v) : v = 0\}$ associated with (1.5). It is the presence of the former that we shall exploit in the remainder of the paper to study nonlinear perturbations of (1.5) that arise from a consideration of problems of the form (1.1).

1.2. Notation. If we define the spectrum of a matrix pencil to be

$$\sigma(A, B) = \{\lambda \in \mathbb{C} : \det(\lambda A + B) = 0\},$$

then $\sigma(A, B) = -\sigma(L)$ (note the minus sign), and p as defined within the KNF above coincides with the number of finite eigenvalues of (A, B) , where eigenvalues are

counted according to their algebraic multiplicity. The matrix pencil (A, B) is said to be *hyperbolic* if $\sigma(A, B)$ is nonempty and contains no elements of unit modulus; otherwise, it is said to be *elliptic*. We shall also write $\rho(A, B) = \sup\{|\lambda| : \lambda \in \sigma(A, B)\}$ and denote the spectral radius of any linear mapping L by $\rho(L)$. Throughout we shall use $\#$ to denote the cardinality of a set of eigenvalues, counted according to algebraic multiplicity.

We shall use $BL(X, Y)$ to denote the space of continuous linear maps from one normed linear space X to another Y , even when X and Y are finite-dimensional. We shall use $B_\epsilon(x)$ for the open ball of radius ϵ about x , and $B_\epsilon(x; X)$ will specify that this ball is contained in the space X . If $L \in BL(X, Y)$, we shall denote the usual operator norm by $\|L\|_{BL(X, Y)}$, which is given by $\sup\{\|Lx\|_Y : x \in X, \|x\|_X = 1\}$. If the context is clear, we shall simply write $\|L\|$, and $BL(X)$ is also used for $BL(X, X)$. Throughout, if $F : X \rightarrow Y$ is a nonlinear mapping, then $dF(x) \in BL(X, Y)$ shall denote the Fréchet derivative, and when acting on $h \in X$ it will be written with square brackets, as in $dF(x)[h]$. Similarly, $d^2F(x)[h, k]$ denotes the second derivative, and this is bilinear in $[h, k]$.

If n is a positive integer, we shall use $\mathcal{O}_n(x)$ on occasion to denote any mapping, H , say, with the property that $\lim_{x \rightarrow 0} \|H(x)\|/\|x\|^n$ exists.

2. Motivation. There are several problems from control theory and numerical analysis that lead to discrete systems where the relationship between the current and future states of a system are not explicit; see [12, 10, 6, 14] for examples.

2.1. Discretized differential-algebraic equations. In [11] the authors apply a Runge–Kutta method to solve a differential-algebraic boundary-value problem arising from an optimal control problem, yielding a nonlinear *difference-algebraic equation* where the control plays the role of an implicit variable. For example, using a forward-Euler method to discretize the differential-algebraic equation (DAE)

$$\dot{x} = f(x, y), 0 = g(x, y) \quad ((x(0), y(0)) \text{ given})$$

yields the Δ AE

$$x_{n+1} = x_n + hf(x_n, y_n), 0 = g(x_n, y_n) \quad ((x_0, y_0) \text{ given}),$$

where h is a small parameter. A singularity in this context occurs when the partial derivative $d_y g(x, y)$ is singular on some subset of $g^{-1}\{0\}$.

Over the past decade a great deal of attention has been devoted to singular DAEs

$$(2.1) \quad F(z, \dot{z}) = 0,$$

where $d_z F(z, \dot{z})$ changes rank on some set, and our study of (1.1) can be viewed as an extension of the work undertaken on (2.1) to the discrete-time case.

It is well known that (2.1) supports a range of singular and regular behavior, including *impasse points* and *pseudoequilibria* [16, 19, 17, 20]. However, not a great deal of the current literature is devoted to the study of *bifurcations* of DAEs nor to the unfolding of singularities in DAEs, such as the image and kernel singularities defined in [23]. One reason for this is the difficulty of proving a suitable center manifold theorem that can cope with the kind of singularities peculiar to DAE. We do note, however, that a Hopf bifurcation theorem is presented in [9] for systems of the form $\dot{x} = f(x, \dot{x}, \alpha)$, where α is a bifurcation parameter, $x \in \mathbb{R}^m$, and f is *nonexpansive* with respect to \dot{x} , a case that does include certain DAE singularities; a Hopf bifurcation

theorem for regular DAEs can be found in [15]. We also note that there are results in the DAE literature that yield the existence of an invariant manifold containing an equilibrium point; for regular DAEs see [18], and for singular DAEs see [24, 2].

2.2. Output-nulling control. Take a discrete dynamical system of the form

$$(2.2) \quad x_{n+1} = f(x_n, u_n), \quad y_n = g(x_n, u_n),$$

where (x_n) is a sequence of states, (u_n) are controls, and (y_n) is a sequence of outputs. One may ask whether there is an admissible control that nullifies or fixes the outputs: Given (y_n) , does there exist a sequence pair $((x_n), (u_n))$ satisfying (2.2)? The resulting equation is an infinite-dimensional system of equations that has a structure reminiscent of the semiexplicit, index-1 DAE (for terminology, see [3]).

2.3. Optimal control. The preprint [14] is relevant to the present work as it presents an invariant manifold result which leads to the existence of a control for the following variational problem with an infinite horizon:

$$\min_{(u_k)} \left\{ \sum_{k=0}^{\infty} \ell(x_k, u_k) : x_{k+1} = f(x_k, u_k), x_0 \in \mathbb{R}^m, u_k \in B_\epsilon(0; \mathbb{R}^p) \right\}$$

such that $f(0, 0) = 0$ and $\ell(0, 0) = 0$. An optimal orbit satisfies the first-order optimality conditions given by the quasi-linear, *implicit* difference equation

$$(2.3) \quad x_{k+1} = f(x_k, u_k),$$

$$(2.4) \quad \frac{\partial H}{\partial x}(x_k, u_k, \lambda_{k+1}) = \lambda_k,$$

$$(2.5) \quad 0 = \frac{\partial H}{\partial u}(x_k, u_k, \lambda_{k+1}),$$

where H is the Hamiltonian $H(x, u, \lambda) = \lambda^T f(x, u) + \ell(x, u)$. In [14] the author demonstrates the existence of a *stable manifold* associated with (2.3)–(2.5) which has a dimension that coincides with the number of eigenvalues of the linearization about its fixed point, and the existence of this stable manifold then provides the necessary optimal control. The obstacle treated in [14] is the existence of a zero closed-loop eigenvalue which is analogous to the type of singularity treated in this paper. One can see the resemblance of (2.3)–(2.5) to (1.5) in that state variables (x_k) propagate forwards in time in (2.3)–(2.5), whereas adjoint variables (λ_k) propagate backwards, a property shared by (1.5).

3. A functional equation for quasi-invariant manifolds. Let us now define in what sense we expect (1.1) to induce a dynamical system. An element $z \in \mathbb{R}^m$ is said to be a *fixed point* of (1.1) if $F(z, z) = 0$. If z denotes the first argument of F and \bar{z} the second, as in $F(z, \bar{z})$, then we define the following conditions:

(A1) $z = 0$ is a fixed point of (1.1): $F(0, 0) = 0$,

(A2) $\det(d_{\bar{z}}F(0, 0)) = 0$, and

(A3) there is a $\xi \in \mathbb{C}$ such that $\det(d_z F(0, 0) + \xi d_{\bar{z}} F(0, 0)) \neq 0$.

Throughout we make use of the matrix pencil (A, B) , where

$$(3.1) \quad A := d_{\bar{z}}F(0, 0) \quad \text{and} \quad B := d_z F(0, 0),$$

and (A3) is the condition that (A, B) is regular. We shall assume that $F \in C^k(\mathbb{R}^{2m}, \mathbb{R}^m)$ for $k > 3$ and seek *local* and *global orbits* of the Δ AE (1.1) in the following sense.

DEFINITION 1. A sequence $(z_n)_{j=0}^J$ is said to be a J -orbit of (1.1) for $J \in \mathbb{N}$ if

$$F(z_j, z_{j+1}) = 0 \text{ for } 0 \leq j \leq J - 1 \text{ and } 2 \leq J < \infty,$$

and a local orbit is a J -orbit for some $J \geq 2$. If there is a $z_2 \in \mathbb{R}^m$ such that $(z_0, z_1; z_2)$ is a local 2-orbit, then we say that the initial condition (z_0, z_1) supports this orbit. A sequence $(z_n)_{j=0}^\infty$ is said to be a global orbit of (1.1) as it is a J -orbit for each $J \geq 2$.

Following the terminology used for DAEs, a pair (z_0, z_1) such that $F(z_0, z_1) = 0$ is said to be *consistent*, and if this pair supports some orbit, then it is said to be a consistent initial condition. We could also have analogously defined backward orbits for $J \leq -2$, but we omit this for brevity. Note that initial conditions lie in \mathbb{R}^{2m} and not \mathbb{R}^m , a property that is analogous to DAEs whereby initial positions *and* certain initial derivatives must be provided in order to obtain the existence of solutions.

We now give the definition which stipulates how we expect (1.1) to induce a dynamical system.

DEFINITION 2. (1.1) induces a local dynamical system on a manifold $\mathcal{M} \subset F^{-1}\{0\} \subset \mathbb{R}^{2m}$ which contains the origin of \mathbb{R}^{2m} if there is a ball $\mathcal{M}_r := \mathcal{M} \cap B_r(0; \mathbb{R}^{2m})$ such that for each $(z, \bar{z}) \in \mathcal{M}_r$ there is a unique $(\bar{z}, \bar{\bar{z}}) \in \mathcal{M}$. If these conditions hold, \mathcal{M} is said to be a solution manifold of (1.1).

This definition ensures that every point $(z, \bar{z}) \in \mathcal{M}_r$ supports the nontrivial 3-orbit $(z, \bar{z}; \bar{\bar{z}})$ and that the point $\bar{\bar{z}} \in \mathbb{R}^m$ is uniquely determined if we are to impose the requirement that $(\bar{z}, \bar{\bar{z}}) \in \mathcal{M}$.

DEFINITION 3. A set $\mathcal{Q} \subset F^{-1}\{0\} \subset \mathbb{R}^{2m}$ is said to be *quasi-invariant* if, for each $(z, \bar{z}) \in \mathcal{Q}$, there exists a $\bar{\bar{z}} \in \mathbb{R}^m$ such that $(\bar{z}, \bar{\bar{z}}) \in \mathcal{Q}$.

As an aside, note that (1.1) induces a trivial dynamical system on the quasi-invariant set $\{(0, 0)\}$ by virtue of (A1), even if assumption (A3) fails. Note also that a solution manifold \mathcal{M} is not necessarily unique; it is the local orbit within \mathcal{M} that must be uniquely determined. Indeed, there may well be many possible choices for $\bar{\bar{z}}$ in order to keep the orbit on $F^{-1}\{0\}$, many of which may not be elements of \mathcal{M} .

3.1. The functional equation. Our strategy for locating quasi-invariant manifolds of (1.1) is to study a functional equation obtained in an analogous manner to the center-manifold equation from the theory of invariant manifolds for maps. The solution of this equation then provides the manifold \mathcal{M} needed to form a local dynamical system for (1.1). We shall show in Theorem 1 that one can find a linear space $K_1 \subset \mathbb{R}^{2m}$ with an associated locally defined, differentiable map $\varphi : K_1 \rightarrow K_1$, a manifold $\mathcal{M} \subset F^{-1}\{0\}$, and a local diffeomorphism $\theta : K_1 \rightarrow \mathcal{M}$ such that

$$F(\theta(u)) = 0 \implies F(\theta(\varphi(u))) = 0.$$

As a result, we will be able to ensure that (1.1) induces a local dynamical system on \mathcal{M} essentially by iterating the map φ . This simply means that if $(z, \bar{z}) = \theta(u)$ is a consistent initial condition, then $(\bar{z}, \bar{\bar{z}}) = \theta(\varphi(u))$ and $(\bar{\bar{z}}, \bar{\bar{\bar{z}}}) = \theta(\varphi(\varphi(u)))$ provide subsequent iterates of (1.1).

Returning to (1.1), let us change the form of the problem by setting $w_n = z_{n+1}$, so that along an orbit of (1.1) we have

$$(3.2) \quad z_{n+1} = w_n,$$

$$(3.3) \quad 0 = Bz_n + Aw_n + \Phi(w_n, z_n),$$

where Φ is the C^k function which satisfies $\Phi(0,0) = 0, d\Phi(0,0) = 0$ and which is defined by

$$F(z, w) - Bz - Aw := \Phi(w, z).$$

The problem of finding an initial condition which is consistent, (z_0, w_0) , say, is of an algebraic nature, whereas the problem of finding an orbit which is supported by this initial condition is a *dynamic* problem. This means that the problem of finding a manifold of orbits of a Δ AE will lead not to an algebraic equation that one could tackle using an elementary version of the implicit function theorem but instead to a functional equation.

Let us now obtain this functional equation. By applying the condition that (A, B) is a regular matrix pencil (condition (A3)), it follows that

$$(3.4) \quad (\mathcal{A}, \mathcal{B}) := \left(\begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & I \\ B & A \end{pmatrix} \right)$$

is also a regular matrix pencil. If we define the vector $W_n = (z_n, w_n) \in \mathbb{R}^{2m}$, using (3.2)–(3.3), (1.1) can be written in the semilinear form

$$\mathcal{A}W_{n+1} = \mathcal{B}W_n + \Psi(W_n),$$

where Ψ is the C^k -mapping $\Psi(W) := (0, \Phi(W))$, so that $\Psi(0) = 0$ and $d\Psi(0) = 0$. There are mappings P and Q that put $(\mathcal{A}, \mathcal{B})$ in Kronecker normal form, and, by setting $W_n = QX_n$, we may write (3.2)–(3.3) in the form

$$(3.5) \quad [PAQ]X_{n+1} = [PBQ]X_n + P\Psi(QX_n),$$

where the terms in square brackets are in normal form:

$$\begin{bmatrix} I & 0 \\ 0 & N \end{bmatrix} X_{n+1} = \begin{bmatrix} C & 0 \\ 0 & I \end{bmatrix} X_n + P\Psi(QX_n).$$

Consequently, there are linear spaces $K_1 \simeq \mathbb{R}^p$ and $K_2 \simeq \mathbb{R}^q$ such that $K_1 \oplus K_2 \simeq \mathbb{R}^{2m}$ and $X_n = (u_n, v_n) \in K_1 \oplus K_2$, where (u_n, v_n) satisfies the difference equation in normal form

$$(NF) \quad \begin{cases} u_{n+1} = Cu_n + f(u_n, v_n), \\ Nv_{n+1} = v_n - g(u_n, v_n). \end{cases}$$

We now ask that there is a manifold given by the graph of some function h on which one can solve (NF) uniquely in a neighborhood of the fixed point $(u, v) = (0, 0)$ in the sense that $v_n = h(u_n)$ holds along orbits. This imposes the two conditions

$$u_{n+1} = Cu_n + f(u_n, h(u_n)) \quad \text{and} \quad Nh(u_{n+1}) = h(u_n) - g(u_n, h(u_n))$$

on h , and it follows that the local orbit (u_n, v_n) of (NF) can be found if h satisfies the functional equation

$$(3.6) \quad h(u) = Nh(Cu + f(u, h(u))) + g(u, h(u)), \quad h(0) = 0, \quad dh(0) = 0,$$

for all u in some neighborhood of the origin in \mathbb{R}^p . The boundary conditions in (3.6) ask first that the fixed point $(u, v) = (0, 0)$ of (NF) lies on the graph of h and then that this graph is tangent to the quasi-invariant subspace obtained on setting $f = 0$ and $g = 0$ in (NF).

3.2. Further preliminaries.

3.2.1. Perturbation of eigenvalues. For completeness we have included the following two preliminary results regarding the spectra of one-parameter families of matrix pencils, which are mappings of the form

$$\mathcal{P} : (-1, 1) \rightarrow BL(\mathbb{R}^m) \times BL(\mathbb{R}^m); \mu \mapsto (A(\mu), B(\mu)).$$

If we define the family of analytic functions

$$f_\mu(\omega) = \det(\omega A(\mu) + B(\mu)),$$

the multiplicity of an eigenvalue of a matrix pencil is then the multiplicity of the corresponding zero of $f_\mu(\cdot)$, which is at most m . The identity

$$(3.7) \quad \frac{d}{d\omega} f_\mu(\omega) = f_\mu(\omega) \operatorname{tr}[(\omega A(\mu) + B(\mu))^{-1} A(\mu)],$$

whenever this inverse is defined, can be used to obtain the following two lemmas.

LEMMA 1 (C^1 -dependence of eigenvalues). *Suppose that $\mathcal{P}(\mu) := (A(\mu), B(\mu))$ is a C^1 -parameterized family of real matrix pencils, with $\mu \in (-1, 1)$, such that $\mathcal{P}(0)$ is a regular matrix pencil. An element $\lambda_0 \in \sigma(\mathcal{P}(0))$ is said to be an algebraically simple eigenvalue of $\mathcal{P}(0)$ if*

$$\ker(\lambda_0 A(0) + B(0)) = \langle x_0 \rangle \quad \text{and} \quad x_0 \notin \operatorname{ran}(\lambda_0 A(0) + B(0)).$$

If λ_0 is an algebraically simple eigenvalue of $\mathcal{P}(0)$, then there is a C^1 -parameter family of algebraically simple eigenvalues $\lambda(\mu) \in \mathbb{C}$ of $\mathcal{P}(\mu)$ such that $\lambda(0) = \lambda_0$, with a corresponding C^1 family of unit eigenvectors $x(\mu)$, with $x(0) = x_0$.

Proof. This follows from the implicit function theorem applied to the system $F(\lambda, x, \mu) = (0, 0)$, where $F(\lambda, x, \mu) := [(\lambda A(\mu) + B(\mu))x, \|x\|_2^2 - 1]$. \square

LEMMA 2 (C^0 -dependence of eigenvalues). *Suppose that $\mathcal{P}(\mu) := (A(\mu), B(\mu))$ is a C^0 -parameterized family of real matrix pencils, with $\mu \in (-1, 1)$, such that $\mathcal{P}(0)$ is regular. If λ_0 is an eigenvalue of $\mathcal{P}(0)$ of algebraic multiplicity l , then it is isolated in the complex plane, and for each $\epsilon > 0$ there is a $\delta > 0$ such that if $|\mu| < \delta$, then $\mathcal{P}(\mu)$ has l eigenvalues (counted according to algebraic multiplicity) in the disk $D(\lambda_0, \epsilon)$.*

Proof. As $f_0(\cdot)$ does not vanish identically because $\mathcal{P}(0)$ is regular by assumption, neither can $f_\mu(\cdot)$ for sufficiently small μ . The *isolatedness* of eigenvalues of $\mathcal{P}(\mu)$ is a consequence of the fact that analytic functions have isolated zeros. Now by using (3.7) we integrate around a closed circle in the complex plane with center $\omega = \lambda_0$ and radius ϵ , from where

$$\#\{\sigma(\mathcal{P}(\mu)) \cap D(\lambda_0, \epsilon)\} = \frac{1}{2\pi i} \oint_{\partial D(\lambda_0, \epsilon)} \operatorname{tr}[(\omega A(\mu) + B(\mu))^{-1} A(\mu)] d\omega,$$

where $D(\lambda_0, \epsilon)$ is an open disk of radius ϵ about λ_0 in the complex plane. This quantity is integer-valued and depends continuously on μ , and the result now follows. \square

3.2.2. Notation. From this point we shall identify the linear space K_1 from the KNF with \mathbb{R}^p and K_2 with \mathbb{R}^q ; now let $|\cdot|_p$ and $|\cdot|_q$ denote norms on \mathbb{R}^p and \mathbb{R}^q , and let $\Omega_\delta = \{u \in \mathbb{R}^p : |u|_p < \delta\}$. We also assume that the unit sphere $\partial\Omega_1$ is a C^∞ manifold.

Let $C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ be the Banach space of continuous maps on $\overline{\Omega}_\delta$ with norm $\|h\|_{C^0} = \sup_{u \in \overline{\Omega}_\delta} |h(u)|_q$. Similarly, let $C^j(\overline{\Omega}_\delta, \mathbb{R}^q)$ be the space of all j -times continuously differentiable functions on $\overline{\Omega}_\delta$ with norm $\|h\|_{C^j} = \max_{0 \leq i \leq j} \sup_{u \in \overline{\Omega}_\delta} \|d^i h(u)\|_{C^0}$, where d^i denotes the i th (Fréchet) derivative, so that $d^j h(u)$ is a j -linear form which we denote $[k_1, \dots, k_j] \rightarrow d^j h(u)[k_1, \dots, k_j]$. Consequently, we have the norm of a higher derivative given by the formula

$$(3.8) \quad \|d^j h\|_{C^0} = \sup_{u \in \overline{\Omega}_\delta} \sup_{|k_i|_p \leq 1} |d^j h(u)[k_1, \dots, k_j]|_q.$$

If M is any multilinear form on a linear space Z and $z \in Z$, then $M[z]^{(k)}$ is shorthand for $M[z, z, \dots, z]$.

From the smoothness of the unit sphere $\partial\Omega_1$, it follows that the embedding of $C^{j+1}(\overline{\Omega}_\delta)$ into $C^j(\overline{\Omega}_\delta)$ is compact, so that if $(h_n) \subset C^{j+1}(\overline{\Omega}_\delta)$ is bounded in the norm of the latter space, there is a subsequence (h_{n_k}) which converges in $C^j(\overline{\Omega}_\delta)$ to some element of $C^j(\overline{\Omega}_\delta)$. We shall also make limited use of the Hölder spaces, which we denote by $C^{j+\alpha}(\overline{\Omega}_\delta)$ whenever j is an integer and $0 < \alpha < 1$, recalling the compact embedding $C^{j+\alpha}(\overline{\Omega}_\delta) \subset C^{j+\beta}(\overline{\Omega}_\delta)$ if $\alpha > \beta$.

It can be somewhat notationally cumbersome to include all of the references to the underlying spaces in all of the norms that we use, so we shall limit their use and expect that the precise meaning can be taken from context.

4. Solving the fixed-point problem (3.6). It is not (NF) that we shall seek to solve directly, but we make the substitution

$$u = \epsilon \tilde{u}, v = \epsilon \tilde{v}$$

in (NF) to give (after removal of the tildes for clarity)

$$(NF)_\epsilon \quad \begin{cases} u_{n+1} = Cu_n + \epsilon^{-1} f(\epsilon u_n, \epsilon v_n), \\ Nv_{n+1} = v_n - \epsilon^{-1} g(\epsilon u_n, \epsilon v_n). \end{cases}$$

As the functions f and g are higher than linear order at the origin, $(NF)_\epsilon$ is in fact smooth with respect to variations in ϵ .

Let us define the one-parameter family of C^k functions \mathbf{f}_ϵ and \mathbf{g}_ϵ (with C^{k-1} dependence on ϵ) by

$$\mathbf{f}_\epsilon(u, v) = \epsilon^{-1} f(\epsilon u, \epsilon v) \quad \text{and} \quad \mathbf{g}_\epsilon(u, v) = \epsilon^{-1} g(\epsilon u, \epsilon v),$$

respectively. For $j \in \mathbb{N}$ we also have

$$(4.1) \quad d^j \mathbf{f}_\epsilon(u, v) = \epsilon^{j-1} d^j f(\epsilon u, \epsilon v) \quad \text{and} \quad d^j \mathbf{g}_\epsilon(u, v) = \epsilon^{j-1} d^j g(\epsilon u, \epsilon v),$$

whenever these derivatives are defined.

Now seek an invariant manifold \mathcal{M} of $(NF)_\epsilon$ given by a graph on which

$$v_n = h(u_n),$$

and then \mathcal{M} can be realized as such a graph if there is a solution of the nonlinear functional equation

$$(4.2) \quad \begin{cases} h(u) = Nh(Cu + \mathbf{f}_\epsilon(u, h(u))) + \mathbf{g}_\epsilon(u, h(u)), \\ h(0) = 0, dh(0) = 0. \end{cases}$$

4.1. Preliminary estimates. The following are simple but essential estimates on the derivatives of f and g . By the mean-value inequality and the fact that the mapping $(u, v) \mapsto (f(u, v), g(u, v))$ and its derivative vanish at $(u, v) = (0, 0)$, there exists an $\ell (= \ell(\delta, r)) > 0$ such that

$$(4.3) \quad |f(u, v)|_p \leq \ell \| (u, v) \|^2, \quad |g(u, v)|_q \leq \ell \| (u, v) \|^2,$$

and

$$(4.4) \quad \|df(u, v)\| \leq \ell \| (u, v) \|, \quad \|dg(u, v)\| \leq \ell \| (u, v) \|,$$

whenever $|u|_p \leq \delta, |v|_q \leq r$, where here and throughout we use the norm

$$\| (u, v) \| = \max(|u|_p, |v|_q) \quad (\forall (u, v) \in \mathbb{R}^p \times \mathbb{R}^q).$$

Here $\|df(u, v)\|$ and $\|dg(u, v)\|$ both refer to induced operator norms, treating $df(u, v)$ and $dg(u, v)$ as linear mappings. By using the mean-value inequality we obtain

$$\|d^2g(u, v) - d^2g(0, 0)\| \leq \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d^3g(u, v)\| \| (u, v) \|,$$

and the triangle inequality gives

$$\|d^2g(u, v)\| \leq \|d^2g(0, 0)\| + \bar{\ell} \| (u, v) \| \quad (|u|_p \leq \delta, |v|_q \leq r),$$

where $\bar{\ell} (= \bar{\ell}(\delta, r)) = \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d^3g(u, v)\|$, whence

$$(4.5) \quad \|d^2g_\epsilon(u, v)\| \leq \epsilon (\|d^2g(0, 0)\| + \bar{\ell} \| (u, v) \|) \quad (|u|_p \leq \delta, |v|_q \leq r).$$

An analogous inequality holds for f and f_ϵ :

$$(4.6) \quad \|d^2f_\epsilon(u, v)\| \leq \epsilon (\|d^2f(0, 0)\| + \bar{\ell} \| (u, v) \|) \quad (|u|_p \leq \delta, |v|_q \leq r).$$

It is the $O(\epsilon)$ size of these quantities that will be important later.

4.2. Introducing a cutoff function. It is not (4.2) that we shall seek to solve directly, but we must employ a cutoff function to rewrite (4.2) in a fixed-point form that is amenable to a Picard iteration. This is not the case at present because if we were to define a nonlinear operator acting on h by the right-hand side of (4.2), there is no reason for it or its iterates to be well-defined on a suitable function space.

For any $\delta > 0$, there is a cutoff function $\psi \in C^\infty(\mathbb{R}^p)$ such that

$$\psi(u) = \begin{cases} u & \text{if } |u|_p \leq \delta/2, \\ 0 & \text{if } |u|_p \geq 3\delta/2 \end{cases}$$

and such that $|\psi(u)|_p \leq \delta$. By using this cutoff we define a Nemitskii operator π as follows:

$$(4.7) \quad \pi(h)(u) = \psi(Cu + f_\epsilon(u, h(u))) \quad (\forall u \in \bar{\Omega}_\delta, h : \bar{\Omega}_\delta \subset \mathbb{R}^p \rightarrow \mathbb{R}^q),$$

so that

$$\pi : C^0(\bar{\Omega}_\delta, \mathbb{R}^q) \rightarrow C^0(\bar{\Omega}_\delta, \mathbb{R}^p),$$

and the inequality $|\pi(h)(u)|_p \leq \delta$ holds pointwise. Moreover, by the C^k -regularity of f , it follows that π is itself a C^k -mapping for each $\epsilon > 0$ fixed, with Fréchet derivative

$$(4.8) \quad d\pi(h)[k](\cdot) = d\psi(C \cdot + f_\epsilon(\cdot, h))[d_v f_\epsilon(\cdot, h)[k]] \quad (\forall h, k \in C^0(\overline{\Omega}_\delta, \mathbb{R}^q)).$$

Also note the following estimate, which will be important later.

LEMMA 3. *Suppose that $h \in C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ satisfies $\|h\|_{C^0} \leq r$, then*

$$(4.9) \quad \|d\pi(h)\|_{BL(C^0)} \leq \epsilon \cdot \|\psi\|_{C^1} \ell \max(\delta, r),$$

and there are constants $\kappa_1, \kappa_2 > 0$, depending on ℓ, δ , and r , but not on ϵ , such that

$$(4.10) \quad \|d^2\pi(h)\|_{BL(C^0) \times BL(C^0)} \leq \epsilon \|\psi\|_{C^2} (\kappa_1(\ell, \delta, r) + \epsilon \kappa_2(\ell, \delta, r)).$$

Proof. By using (4.8) we obtain

$$\begin{aligned} \sup_{k \in C^0, \|k\|_{C^0}=1} \|d\pi(h)[k]\|_{C^0} &= \sup_{k \in C^0, \|k\|_{C^0}=1} \|d\psi(C \cdot + f_\epsilon(\cdot, h))[d_v f_\epsilon(\cdot, h)[k]]\|_{C^0} \\ &\leq \|\psi\|_{C^1} \cdot \sup_{k \in C^0, \|k\|_{C^0}=1} \|d_v f_\epsilon(\cdot, h)[k]\|_{C^0} \\ &\leq \|\psi\|_{C^1} \cdot \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d_v f_\epsilon(u, v)\| \\ &\leq \|\psi\|_{C^1} \cdot \ell \cdot \sup_{|u|_p \leq \delta, |v|_q \leq r} \|(\epsilon u, \epsilon v)\| \quad (\text{by (4.1) and (4.3)}), \end{aligned}$$

and the first part follows. The second part follows from

$$d^2\pi(h)[k_1, k_2] = d^2\psi(C \cdot + f_\epsilon(\cdot, h))[d_v f_\epsilon[k_1], d_v f_\epsilon[k_2]] + d\psi[d_{vv}^2 f_\epsilon(\cdot, h)[k_1, k_2]],$$

so that, for $i = 1, 2$,

$$\begin{aligned} \sup_{k_i \in C^0, \|k_i\|_{C^0}=1} \|d^2\pi(h)[k_1, k_2]\|_{C^0} &\leq \sup_{u \in \mathbb{R}^p} \|d\psi(u)\| \cdot \|d_{vv}^2 f_\epsilon(\cdot, h)\| \\ &\quad + \sup_{u \in \mathbb{R}^p} \|d^2\psi(u)\| \cdot \|d_v f_\epsilon(\cdot, h)\|^2 \\ &\leq \|\psi\|_{C^2} \left(\sup_{|u|_p \leq \delta, |v|_q \leq r} \|d_{vv}^2 f_\epsilon(u, v)\| + \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d_v f_\epsilon(u, v)\|^2 \right) \\ &\leq \|\psi\|_{C^2} \left(\epsilon \|d^2 f(0, 0)\| + \epsilon^2 \bar{\ell} \max(\delta, r) + \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d_v f(\epsilon u, \epsilon v)\|^2 \right) \quad (\text{by (4.5)}) \\ &\leq \|\psi\|_{C^2} \left(\epsilon \|d^2 f(0, 0)\| + \epsilon^2 \bar{\ell} \max(\delta, r) + \sup_{|u|_p \leq \delta, |v|_q \leq r} \|(\epsilon u, \epsilon v)\|^2 \ell^2 \right) \quad (\text{by (4.3)}), \end{aligned}$$

and the result follows directly from here. \square

In order to solve (4.2), we now tackle the following nonlinear fixed-point problem:

$$(4.11) \quad h = Nh(\pi(h)) + G_\epsilon(h),$$

where $G_\epsilon : C^0(\overline{\Omega}_\delta, \mathbb{R}^q) \rightarrow C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ is the Nemitskii operator

$$G_\epsilon(h)(u) := \mathbf{g}_\epsilon(u, h(u)).$$

Notice that the cutoff ψ has been used in (4.11), but, because ψ coincides with the identity on some balls around the origin, solutions of (4.11) will satisfy (4.2) on this ball.

This construction ensures that $h(\pi(h)) \in C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ whenever $h \in C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$, and we may, as a result, define the operator

$$T : C^0(\overline{\Omega}_\delta, \mathbb{R}^q) \rightarrow C^0(\overline{\Omega}_\delta, \mathbb{R}^q); \quad h \mapsto h(\pi(h)),$$

noting that the operators $T : C^1(\overline{\Omega}_\delta, \mathbb{R}^q) \rightarrow C^1(\overline{\Omega}_\delta, \mathbb{R}^q)$ and $T : C^2(\overline{\Omega}_\delta, \mathbb{R}^q) \rightarrow C^2(\overline{\Omega}_\delta, \mathbb{R}^q)$ are also well-defined as the restrictions of T to various subspaces of $C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ as f and g are C^3 functions.

However, it is not (4.11) that we shall solve, but we exploit the nilpotency of N to bring the *functional part* of (4.2) and (4.11), that is, $Nh(\pi(h))$, into a higher-order contribution to the problem. However, because $g(u, v)$ is a second-, or possibly higher-order function, the operator G_ϵ contains no linear terms, and this will help us to obtain a contractive sequence by iterating T .

By way of example, let us suppose that $N \neq 0$ but $N^2 = 0$. If there exists a solution of $h = Nh(\pi(h)) + G_\epsilon(h)$, there results $Nh = NG_\epsilon(h)$, and therefore $Nh(\pi(h)) = NG_\epsilon(h(\pi(h)))$. In this case we find that h must satisfy

$$(4.12) \quad h = NG_\epsilon(h(\pi(h))) + G_\epsilon(h),$$

and one observes that the *functional part* of the equation (that is, $h(\pi(h))$) now sits inside a higher-order term (and not a linear one as in (4.11)). Conversely, if h satisfies (4.12), then $Nh = NG_\epsilon(h)$, so that $NG_\epsilon(h(\pi(h))) = Nh(\pi(h))$ because $N^2 = 0$, and $h = Nh(\pi(h)) + G_\epsilon(h)$ follows.

When the nilpotency index of the map N is arbitrary we extend this idea in the following lemma, where here and in the remainder we shall write

$$\mathcal{G}_\epsilon(h) = \sum_{j=0}^{\nu-1} N^j G_\epsilon(T^j(h)),$$

and

$$T^{j+1} = T(T^j), \quad \text{where} \quad T^0 = I,$$

and the latter denotes the identity on $C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$.

LEMMA 4. *Suppose that $N^\nu = 0$ but $N^{\nu-1} \neq 0$; then h is a solution of (4.11) if it is a solution of the fixed-point problem*

$$(4.13) \quad h = \mathcal{G}_\epsilon(h).$$

Proof. Let us suppose that h is a solution of (4.13); then

$$Nh = N\mathcal{G}_\epsilon(h) = N \left[\sum_{j=0}^{\nu-1} N^j G_\epsilon(T^j(h)) \right] = \sum_{j=0}^{\nu-1} N^{j+1} G_\epsilon(T^j(h)),$$

and so, because $N^\nu = 0$, we obtain

$$\begin{aligned} Nh(\pi(h)) &= \sum_{j=0}^{\nu-1} N^{j+1} G_\epsilon(T^j(h(\pi(h)))) \\ &= \sum_{j=0}^{\nu-1} N^{j+1} G_\epsilon(T^{j+1}(h)) = \sum_{j=1}^{\nu-1} N^j G_\epsilon(T^j(h)) = \mathcal{G}_\epsilon(h) - G_\epsilon(h) = h - G_\epsilon(h), \end{aligned}$$

which therefore provides a solution of (4.11) as required. \square

4.3. The main result. Our strategy for solving (4.13), and hence (4.2), is to show that \mathcal{G}_ϵ satisfies a *refined* Banach contraction theorem of the type given in [26, p. 286]. The idea that we employ several times is encapsulated in the following idea. Consider Banach spaces $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$ such that $x_0 \in Y$ and $Y \subset X$, and moreover \mathcal{T} is a mapping satisfying $\mathcal{T} : Y \rightarrow Y$ and $\mathcal{T} : X \rightarrow X$. Now, if there is an $r > 0$ and a $\kappa \in [0, 1)$ such that

1. $\mathcal{T} : \overline{B}_r(x_0; X) \rightarrow \overline{B}_r(x_0; X)$,
 2. $\|\mathcal{T}(y) - \mathcal{T}(y')\|_X \leq \kappa \|y - y'\|_X$ for all $y, y' \in Y \cap \overline{B}_r(x_0; X)$,
- then \mathcal{T} has a fixed point $y^* \in \overline{B}_r(x_0; X)$. In addition, if
3. $\mathcal{T} : \overline{B}_\rho(x_0; Y) \rightarrow \overline{B}_\rho(x_0; Y)$ and
 4. there is an interpolating Banach space Z such that $Y \subset Z$ with compact embedding and $Z \subset X$ with continuous embedding,
- then $y^* \in Z$.

The point here is that we cannot ensure that $y^* \in Y$, although one still obtains a fixed point in some space from the standard iteration scheme. By using this idea one can prove an existence and regularity result for (4.2), where we have in mind $x_0 = 0, \mathcal{T} = \mathcal{G}_\epsilon, Y = C^{k+1}, Z = C^{k+\alpha}$, and $X = C^k$, where $\alpha \in (0, 1)$. We begin by providing the details to cover the cases $k = 0$ and $k = 1$.

The following theorem is the main result of this paper from which the invariant manifold and bifurcation theorems are deduced.

THEOREM 1. *Let $\alpha \in [0, 1)$. There exists an $\epsilon_0 > 0$ such that, for each $\epsilon \in (0, \epsilon_0)$, (4.13) has a solution $h \in C^{1+\alpha}(\overline{\Omega}_\delta, \mathbb{R}^q)$; moreover $h(0) = 0$ and $dh(0) = 0$.*

Proof. Let X_r be the C^0 -closed ball of radius r about zero in $C^0(\overline{\Omega}_\delta, \mathbb{R}^q)$ and Y_r the C^1 -closed ball of radius r about zero in $C^1(\overline{\Omega}_\delta, \mathbb{R}^q)$. Let $h_0 \in Y_r$, and define a sequence

$$h_{n+1} := \mathcal{G}_\epsilon(h_n).$$

We shall show that we can choose ϵ such that (h_n) is well-defined, contractive, and hence Cauchy in X_r , and it therefore converges. Throughout the remainder of the proof we shall use the positive constant

$$n_* := \sum_{j=0}^{\nu-1} \|N\|_{BL(\mathbb{R}^q)}^j.$$

We now give a proof of Theorem 1 in four short steps, each placing a stronger restriction on ϵ relative to the fixed choice of δ and r to ensure that \mathcal{G}_ϵ contracts when acting on C^1 functions, measured in the C^0 norm.

Step 1. If $\epsilon < r(\ell n_* \max(\delta, r)^2)^{-1} =: \epsilon_1$, then $\mathcal{G}_\epsilon : X_r \rightarrow X_r$.

Proof. Suppose that $h \in X_r$, and then $T(h) = h(\pi(h))$ is continuous as h is. Moreover if $\|h\|_{C^0} \leq r$, then $\|T(h)\|_{C^0} = \|h(\pi(h))\|_{C^0} \leq r$; similarly $\|T^j(h)\|_{C^0} \leq r$ for all $0 \leq j \leq \nu - 1$. By definition,

$$\begin{aligned} \|\mathcal{G}_\epsilon(h)\|_{C^0} &\leq \sum_{j=0}^{\nu-1} \|N\|^j \|\mathcal{G}_\epsilon(T^j(h))\|_{C^0} \leq \sum_{j=0}^{\nu-1} \|N\|^j \sup_{\|\bar{h}\|_{C^0} \leq r} \|\mathcal{G}_\epsilon(\bar{h})\|_{C^0} \\ &\leq n_* \sup_{|u|_p \leq \delta, |v|_q \leq r} |\mathbf{g}_\epsilon(u, v)|_q = n_* \sup_{|u|_p \leq \delta, |v|_q \leq r} \epsilon^{-1} |g(\epsilon u, \epsilon v)|_q. \end{aligned}$$

From (4.3) it follows that

$$\|\mathcal{G}_\epsilon(h)\|_{C^0} \leq \epsilon \ell n_* \sup_{|u|_p \leq \delta, |v|_q \leq r} \|(u, v)\|^2 = \epsilon \ell n_* \max(\delta, r)^2.$$

By assumption, $\epsilon \ell n_* \max(\delta, r)^2 < r$, and Step 1 is complete. \square

Step 2. If $\epsilon < \min\{\epsilon_1, [r\ell \max(\delta, r)\|\psi\|_{C^1}]^{-1}, r[n_*(1+r)\ell \max(\delta, r)]^{-1}\} =: \epsilon_2$, then $\mathcal{G}_\epsilon : Y_r \rightarrow Y_r$.

Proof. Suppose that $h \in Y_r$ so that $h \in X_r$, then $T^j(h)$ is differentiable on $\overline{\Omega}_\delta$ as h is; moreover throughout the remainder of the proof of Step 2 we shall write $H := \mathcal{G}_\epsilon(h)$ for brevity. It then follows that

$$(4.14) \quad dH(u) = \sum_{j=0}^{\nu-1} N^j [d_u \mathbf{g}_\epsilon(u, T^j(h)(u)) + d_v \mathbf{g}_\epsilon(u, T^j(h)(u)) [d_u(T^j(h))(u)]],$$

and we now need to estimate $\|dH(u)\|_{C^0}$. By using the fact that $\|h\|_{C^0} \leq r$, since $T : X_r \rightarrow X_r$ from Step 1, we obtain $|T^j(h)(u)|_q = |h(\pi(h(\dots)))|_q \leq r$, and therefore

$$\begin{aligned} \sup_{|u|_p \leq \delta} \|dH(u)\|_{BL(\mathbb{R}^p, \mathbb{R}^q)} &\leq \sum_{j=0}^{\nu-1} \|N\|^j (\|d_u \mathbf{g}_\epsilon(u, T^j(h)(u))\| \\ &\quad + \|d_v \mathbf{g}_\epsilon(u, T^j(h)(u))\| \|d_u(T^j(h))(u)\|) \\ &\leq n_* \sup_{0 \leq j \leq \nu-1} (\|d_u \mathbf{g}_\epsilon(u, T^j(h)(u))\| + \|d_v \mathbf{g}_\epsilon(u, T^j(h)(u))\| \|d_u(T^j(h))(u)\|) \\ &\leq n_* \sup_{\substack{0 \leq j \leq \nu-1 \\ |u|_p \leq \delta, |v|_q \leq r}} (\|d_u \mathbf{g}_\epsilon(u, v)\| + \|d_v \mathbf{g}_\epsilon(u, v)\| \|d_u(T^j(h))(u)\|) \\ &\leq \epsilon n_* \ell \max(r, \delta) \left(1 + \sup_{0 \leq j \leq \nu-1, |u|_p \leq \delta} \|d_u(T^j(h))(u)\| \right). \end{aligned}$$

Now we estimate the final bracketed term in the latter expression. The linear mapping obtained from differentiating $T^j(h)(u)$ with respect to u is

$$d_u(T^j(h))(u) = d_u(h(\pi(h(\dots \pi(h) \dots))),$$

which can be written as the recurrence

$$(4.15) \quad d_u(T^j(h))(u) = dh(\pi(T^{j-1}(h))(u)) d\pi(T^{j-1}(h)(u)) \cdot d_u(T^{j-1}(h)(u)),$$

where, by definition, $d_u(T^0(h))(u) = d_u(h)(u) = dh(u)$. By taking C^0 norms and setting

$$\xi_j := \|d_u(T^j(h))(u)\|_{BL(\mathbb{R}^p, \mathbb{R}^q)},$$

we obtain the relation

$$\xi_j \leq \|dh\|_{C^0} \sup_{\|\bar{h}\|_{C^0} \leq r} \|d\pi(\bar{h})\|_{BL(C^0)} \cdot \xi_{j-1},$$

where $\xi_0 = \|dh\|_{C^0} \leq \|h\|_{C^1} \leq r$. From (4.9) of Lemma 3, we find that

$$\xi_j \leq \epsilon \ell r \max(\delta, r) \|\psi\|_{C^1} \cdot \xi_{j-1} \leq (\epsilon \ell r \max(\delta, r) \|\psi\|_{C^1})^j \xi_0 \leq r$$

because $\epsilon \ell r \max(\delta, r) \|\psi\|_{C^1} < 1$ by assumption. As a result, the inequality

$$\begin{aligned} \sup_{|u|_p \leq \delta} \|dH(u)\|_{BL(\mathbb{R}^p, \mathbb{R}^q)} &\leq \epsilon n_* \ell \max(r, \delta) \left(1 + \sup_{0 \leq j \leq \nu-1} \xi_j \right) \\ &\leq \epsilon n_* \ell \max(r, \delta) (1 + r) \leq r \end{aligned}$$

also now follows from the assumption of the claim, and we have proven that $h \in Y_r$, as required. \square

Step 3. If $\epsilon < \epsilon_2$ and we define the $O(\epsilon)$ quantity $\kappa(\epsilon)$ by

$$\kappa(\epsilon) := \epsilon \cdot \ell \max(\delta, r) \sum_{i=0}^{\nu-1} (r\epsilon \cdot \|\psi\|_{C^1} \ell \max(\delta, r))^i,$$

then $\|h_{n+1} - h_n\|_{C^0} \leq \kappa(\epsilon) \|h_n - h_{n-1}\|_{C^0}$. As a result, (h_n) is Cauchy in X_r if ϵ is further restricted so that $\kappa(\epsilon) < 1$, and so (h_n) converges in X_r .

Proof. For brevity, let us write h in place of h_{n+1} and k for h_n after setting $h_0 = 0 \in Y_r$ and $h_{n+1} = \mathcal{G}_\epsilon(h_n)$. Now

$$\begin{aligned} \|\mathcal{G}_\epsilon(h) - \mathcal{G}_\epsilon(k)\|_{C^0} &\leq \sum_{j=0}^{\nu-1} \|N\|^j \|\mathbf{G}_\epsilon(T^j(h)) - \mathbf{G}_\epsilon(T^j(k))\|_{C^0} \\ &\leq n_* \sup_{0 \leq j \leq \nu-1} \|\mathbf{g}_\epsilon(u, T^j(h)(u)) - \mathbf{g}_\epsilon(u, T^j(k)(u))\|_{C^0}. \end{aligned}$$

The mean-value inequality yields

$$\begin{aligned} &|\mathbf{g}_\epsilon(u, T^j(h)(u)) - \mathbf{g}_\epsilon(u, T^j(k)(u))|_q \\ &\leq \sup_{z \in [T^j(h)(u), T^j(k)(u)]} \|d_v \mathbf{g}_\epsilon(u, z)\| \|T^j(h) - T^j(k)\|_{C^0} \\ &\leq \sup_{|u|_p \leq \delta, |v|_q \leq r} \|d_v \mathbf{g}_\epsilon(u, v)\| \|T^j(h) - T^j(k)\|_{C^0} \quad (\text{using Step 1}) \\ (4.16) \quad &\leq \epsilon \ell \max(r, \delta) \cdot \|T^j(h) - T^j(k)\|_{C^0} \quad (\text{by (4.3)}), \end{aligned}$$

where, for any $z_1, z_2 \in \mathbb{R}^p$, the generalized interval from z_1 to z_2 is given by

$$[z_1, z_2] := \{\lambda z_1 + (1 - \lambda)z_2 : 0 \leq \lambda \leq 1\}.$$

So let us define $\chi_j := \|T^j(h) - T^j(k)\|_{C^0}$, and we estimate χ_j as follows:

$$\begin{aligned} |T^j(h)(u) - T^j(k)(u)|_q &= |h(\pi(T^{j-1}(h))) - k(\pi(T^{j-1}(k)))|_q \\ &\leq |h(\pi(T^{j-1}(h))) - k(\pi(T^{j-1}(h)))|_q \\ &\quad + |k(\pi(T^{j-1}(h))) - k(\pi(T^{j-1}(k)))|_q \\ \implies \chi_j &\leq \|h - k\|_{C^0} + \|k(\pi(T^{j-1}(h))) - k(\pi(T^{j-1}(k)))\|_{C^0}. \end{aligned}$$

However, from Step 2 we have $\|k\|_{C^1} \leq r$, and therefore

$$|k(u) - k(u')|_q \leq r|u - u'|_p \quad \forall u, u' \in \bar{\Omega}_\delta.$$

Using the fact that π is a Fréchet differentiable mapping on $C^0(\bar{\Omega}_\delta, \mathbb{R}^q)$ with norm bounded according to (4.9), an application of the mean-value inequality gives

$$\chi_j \leq \|h - k\|_{C^0} + r \sup_{\|\bar{h}\|_{C^1} \leq r} \|d\pi(\bar{h})\|_{BL(C^0)} \cdot \chi_{j-1},$$

so that $\chi_j \leq \|h - k\|_{C^0} + \epsilon \cdot r \|\psi\|_{C^1} \ell \max(\delta, r) \chi_{j-1}$, where $\chi_0 = \|h - k\|_{C^0}$. The discrete Gronwall inequality now gives

$$\chi_j \leq \|h - k\|_{C^0} \sum_{i=0}^{\nu-1} (r\epsilon \|\psi\|_{C^1} \ell \max(\delta, r))^i,$$

and with (4.16) we have the desired inequality

$$\|\mathcal{G}_\epsilon(h) - \mathcal{G}_\epsilon(k)\|_{C^0} \leq \kappa(\epsilon) \cdot \|h - k\|_{C^0}.$$

The standard contraction argument now shows that $(h_n) \subset Y_r$ is Cauchy in X_r as claimed, and there therefore exists an $h \in X_r$ such that $h_n \xrightarrow{C^0} h$ as $n \rightarrow \infty$. \square

Since $(h_n) \subset Y_r$, it follows that there is a subsequence (n_j) such that $h_{n_j} \xrightarrow{C^\alpha} \bar{h}$ as $j \rightarrow \infty$ for some $\bar{h} \in X_r \cap C^\alpha(\bar{\Omega}_\delta)$, and we deduce that $\bar{h} = h \in C^\alpha(\bar{\Omega}_\delta)$. By restricting ϵ further we can actually ensure that h is differentiable, as follows.

Step 4. There is an $\epsilon_3 > 0$ such that $h \in C^{1+\alpha}(\bar{\Omega}_\delta)$ whenever $\epsilon < \min(\epsilon_2, \epsilon_3)$ and $\kappa(\epsilon) < 1$.

Proof. Let $h \in C^2(\bar{\Omega}_\delta)$ satisfy $\|h\|_{C^2} \leq r$, and recall that $H := \mathcal{G}_\epsilon(h)$. From (4.14) we obtain

$$\begin{aligned} d^2H(u) &= \sum_{j=0}^{\nu-1} N^j \{ d_{uu}^2 \mathbf{g}_\epsilon(u, T^j(h)(u)) + 2d_{uv}^2 \mathbf{g}_\epsilon(u, T^j(h)(u)) [I, d_u(T^j(h)(u))] \\ &\quad + d_{vv}^2 \mathbf{g}_\epsilon(u, T^j(h)(u)) [d_u(T^j(h)(u)), d_u(T^j(h)(u))] \\ &\quad + d_v \mathbf{g}_\epsilon(u, T^j(h)(u)) [d_{uu}^2(T^j(h)(u))] \}. \end{aligned}$$

In seeking a bound on $\|H\|_{C^2}$, we now examine the term $d_{uu}^2(T^j(h)(u))$ more closely as bounds on the remaining elements of d^2H can be obtained from Steps 1 and 2. By applying the chain rule (4.15) we obtain the recurrence

$$\begin{aligned} d_{uu}^2(T^j(h)(u)) &= d_u \{ dh(\pi(T^{j-1}(h))(u)) d\pi(T^{j-1}(h)(u)) \cdot d_u(T^{j-1}(h)(u)) \} \\ &= d^2h(\pi(T^{j-1}(h))(u)) [d\pi(T^{j-1}(h)(u)) \cdot d_u(T^{j-1}(h)(u))]^{(2)} \\ &\quad + dh(\pi(T^{j-1}(h))(u)) [d^2\pi(T^{j-1}(h)(u)) [d_u(T^{j-1}(h)(u))]^{(2)}] \\ &\quad + dh(\pi(T^{j-1}(h))(u)) [d\pi(T^{j-1}(h)(u)) [d_{uu}^2(T^{j-1}(h)(u))]], \end{aligned}$$

and taking norms gives

$$\begin{aligned} \|d_{uu}^2(T^j(h)(u))\|_{C^0} &\leq \|d^2h\|_{C^0} \|d\pi(h)\|^2 \|d_u(T^j(h)(u))\|_{C^0} \\ &\quad + \|dh\|_{C^0} \|d^2\pi(h)\| \|d_u(T^j(h)(u))\|_{C^0}^2 \\ &\quad + \|dh\|_{C^0} \|d\pi(h)\| \|d_{uu}^2(T^j(h)(u))\|_{C^0}. \end{aligned}$$

If we write

$$\eta_j := \|d_{uu}^2(T^j(h)(u))\|_{C^0},$$

and use the fact that $\xi_j = \|d_u(T^j(h))(u)\|_{BL(\mathbb{R}^p, \mathbb{R}^q)} \leq r$ for all $0 \leq j \leq \nu - 1$, which was established in Step 2, we obtain the difference inequality

$$\eta_j \leq r^3 \|d\pi(h)\|_{C^0}^2 + r^2 \|d^2\pi(h)\|_{C^0} + r \|d\pi\|_{C^0} \cdot \eta_{j-1}$$

such that $\eta_0 \leq r$ by definition. There results, for $j \geq 1$,

$$\eta_j \leq (r \|d\pi(h)\|_{C^0})^j \eta_0 + (r^3 \|d\pi(h)\|_{C^0}^2 + r^2 \|d^2\pi(h)\|_{C^0}) \sum_{i=0}^{j-1} (r \|d\pi(h)\|_{C^0})^i,$$

and the bounds (4.9) and (4.10) show that $\eta := \max_{0 \leq j \leq \nu-1} \eta_j$ has an $O(1)$ dependence on ϵ in the sense that there is an $M > 0$ such that $\eta \leq M$ whenever $0 \leq \epsilon \leq 1$. (In fact, one can choose M to be less than r if ϵ is sufficiently small.)

We now obtain

$$\begin{aligned}
\|d^2H\|_{C^0} &\leq n_* \sup_{0 \leq j \leq \nu-1} \{ \|d_{uu}^2 \mathbf{g}_\epsilon(u, T^j(h)(u))\|_{C^0} + 2\xi_j \|d_{uv}^2 \mathbf{g}_\epsilon(u, T^j(h)(u))\|_{C^0} \\
&\quad + \xi_j^2 \|d_{vv}^2 \mathbf{g}_\epsilon(u, T^j(h)(u))\|_{C^0} + \eta_j \|d_v \mathbf{g}_\epsilon(u, T^j(h)(u))\|_{C^0} \} \\
&\leq n_* \sup_{\|\bar{h}\|_{C^0} \leq r, 0 \leq j \leq \nu-1} \{ \|d_{uu}^2 \mathbf{g}_\epsilon(u, \bar{h})\|_{C^0} + 2r \|d_{uv}^2 \mathbf{g}_\epsilon(u, \bar{h})\|_{C^0} \\
(4.17) \quad &\quad + r^2 \|d_{vv}^2 \mathbf{g}_\epsilon(u, \bar{h})\|_{C^0} + \eta \|d_v \mathbf{g}_\epsilon(u, \bar{h})\|_{C^0} \},
\end{aligned}$$

where η has been used to bound the last term in (4.17). By using (4.5) to estimate the second derivative terms we find that

$$\begin{aligned}
\|d^2H\|_{C^0} &\leq \epsilon \cdot n_* \sup_{0 \leq j \leq \nu-1} \{ \|d^2g(0, 0)\| + \epsilon \bar{\ell} \max(r, \delta) \\
&\quad + 2r(\|d^2g(0, 0)\| + \epsilon r \bar{\ell} \max(r, \delta)) \\
(4.18) \quad &\quad + r^2(\|d^2g(0, 0)\| + \epsilon \bar{\ell} \max(r, \delta)) + \eta \ell \max(r, \delta) \}.
\end{aligned}$$

It is immediate from (4.18) and Steps 1 and 2 that a suitably small choice of ϵ ensures that $\|H\|_{C^2} \leq r$ whenever $\|h\|_{C^2} \leq r$. As a result, if we impose the following restriction on the initial guess for a fixed point of \mathcal{G}_ϵ :

$$h_0 \in Z_r := \{h \in C^2(\bar{\Omega}_\delta, \mathbb{R}^q) : \|h\|_{C^2} \leq r\},$$

then the C^0 -convergent sequence from Step 3 also satisfies $(h_n) \subset Z_r$. This means that a $C^{1+\alpha}$ -convergent subsequence can now be extracted from (h_n) , so that the convergence of h_n to h actually occurs in $C^{1+\alpha}$ and h therefore lies in this smoother space. \square

We have shown that there is a differentiable solution of (4.2) on a sufficiently small ball around the origin, but there remains to prove the last part of Theorem 1 regarding the behavior of h at the origin. So let h be a C^1 solution of (4.2) on some ball $\bar{\Omega}_\delta$, and put $\zeta = h(0)$. It follows that ζ is a solution of the algebraic equation

$$(4.19) \quad -\zeta + Nh(f_\epsilon(0, \zeta)) + \mathbf{g}_\epsilon(0, \zeta) = 0,$$

and (4.19) has solution $\zeta = 0$. As the linearization of the left-hand side of (4.19) at $\zeta = 0$ is a multiple of the identity, the inverse function theorem ensures that $\zeta = 0$ is the only solution of (4.19) in some neighborhood of zero, and this ensures that $h(0) = 0$. Since any solution of (4.13) provides one of (4.2), we can differentiate (4.2) with respect to u and set $u = 0$; this gives

$$dh(0) = Ndh(0)[C],$$

but then we can continue in an inductive manner to deduce that

$$Ndh(0)[C] = N^2dh(0)[C^2] = \dots = N^\nu dh(0)[C^\nu] = 0,$$

as N is nilpotent. We find that $dh(0) = 0$, and this concludes the proof of Theorem 1. \square

The question of maximal regularity of a solution of (4.2), or equivalently (4.13), is not addressed, although the method of proof used in Theorem 1 can be continued by restricting ϵ further as required to show that \mathcal{G}_ϵ maps the ball of radius r in $C^j(\bar{\Omega}_\delta, \mathbb{R}^p)$

into itself. This ensures that the sequence (h_n) constructed in the proof of Theorem 1 converges to h in as strong a C^j norm as we like, provided that f and g are sufficiently smooth. We cannot, however, be sure that the resulting solution h lies in $C^\infty(\bar{\Omega}_\delta, \mathbb{R}^p)$ because the interval in which ϵ must reside so as to obtain a C^j fixed point of \mathcal{G}_ϵ could shrink indefinitely as j grows.

Theorem 1 does not ensure the existence of a continuous fixed point of (4.2) which is not C^α for some $0 < \alpha < 1$. If we take a continuous initial guess for a solution, $h_0 \in C^0(\bar{\Omega}_\delta, \mathbb{R}^q)$, say, then, although there is a sequence of continuous functions defined by $h_{n+1} = \mathcal{G}_\epsilon(h_n)$, there is no reason to suspect that the sequence of iterates (h_n) will satisfy the property of being a C^0 -contractive sequence. On the other hand, a suitably small C^1 initial guess will lead to C^α convergence of the resulting iterates for any $\alpha \in [0, 1)$.

4.4. Stability implies uniqueness. There are some simple cases where uniqueness and smoothness can be easily established. The most obvious is where $N = 0$, and then Theorem 1 can be proven using the elementary implicit function theorem. Another occurs when the matrix denoted C that arises from the Kronecker normal form of $(\mathcal{A}, \mathcal{B})$ in (3.5) satisfies $\|C\| < 1$ in the norm induced by $|\cdot|_p$. In this case the cutoff function ψ used above is not needed in order to obtain a well-defined operator π . If we define the Nemitskii operator

$$\pi(h)(u) = Cu + f_\epsilon(u, h(u)),$$

and $h \in C^0(\bar{\Omega}_\delta, \mathbb{R}^q)$ satisfies $\|h\|_{C^0} \leq r$, then $|Cu + f_\epsilon(u, h(u))|_p \leq \|C\| \|u\|_p + |f_\epsilon(u, h(u))|_p \leq \|C\| \delta + \epsilon \cdot \ell \max(\delta, r)^2$. In order for $h(\pi(h))$ to be well-defined, we now need only to choose ϵ such that $\|C\| \delta + \epsilon \cdot \ell \max(\delta, r)^2 \leq \delta$, which can be done. The proof of Theorem 1 then goes through with this minor modification, and the resulting fixed point of \mathcal{G}_ϵ is unique in the space of continuously differentiable functions.

4.5. Polynomial approximation. Let us remove the dependence of (4.2) on ϵ for clarity and return to the fixed-point problem (3.6) directly, which we recall defines a nonlinear operator G via

$$(4.20) \quad \begin{aligned} h(u) &= Nh(Cu + f(u, h(u))) + g(u, h(u)) =: (Gh)(u), \\ &\text{where } h(0) = 0, dh(0) = 0. \end{aligned}$$

PROPOSITION 1. *Suppose that $h \in C^k(\bar{\Omega}_\delta, \mathbb{R}^q)$ is a solution of (4.20), and suppose that $\bar{h} \in C^k(\bar{\Omega}_\delta, \mathbb{R}^q)$ satisfies*

$$\bar{h}(u) - G(\bar{h})(u) = e(u), \quad \bar{h}(0) = 0,$$

where e is a given C^k function that satisfies $e(0) = 0, de(0) = 0, \dots, d^k e(0) = 0$. Then

$$|h(u) - \bar{h}(u)|_q = o(|u|_p^k) \text{ as } u \rightarrow 0.$$

Proof. If we define the function $\Delta := h - \bar{h}$, then we have to prove that

$$\Delta(0) = 0, d\Delta(0) = 0, \dots, d^k \Delta(0) = 0,$$

and the result then follows from the basic properties of the derivative. Clearly $\Delta(0) = 0$, and differentiating (4.20) with respect to u gives

$$\begin{aligned} dh(u) &= Ndh(Cu + f(u, h(u)))[C + d_u f(u, h(u)) + d_v f(u, h(u))[dh(u)]] \\ &\quad + d_u g(u, h(u)) + d_v g(u, h(u))[dh(u)], \end{aligned}$$

and similarly,

$$\begin{aligned} d\bar{h}(u) &= Nd\bar{h}(Cu + f(u, \bar{h}(u)))[C + d_u f(u, \bar{h}(u)) + d_v f(u, \bar{h}(u))[d\bar{h}(u)]] \\ &\quad + d_u g(u, \bar{h}(u)) + d_v g(u, \bar{h}(u))[d\bar{h}(u)] + de(u). \end{aligned}$$

As a result, because $h(0) = \bar{h}(0) = 0$, we find that

$$d\Delta(0) = Nd\Delta(0)[C] + de(0) = Nd\Delta(0)[C],$$

by the assumption that $de(0) = 0$, so that

$$d\Delta(0) = Nd\Delta(0)[C] = N^2 d\Delta(0)[C^2] = \dots = N^\nu d\Delta(0)[C^\nu] = 0.$$

Continuing in a similar vein, we obtain

$$\begin{aligned} d^2 h(u) &= Nd^2 h(Cu + f(u, h(u)))[C + d_u f(u, h(u)) \\ &\quad + d_v f(u, h(u))[dh(u)]^{(2)} + d_{uu}^2 g(u, h(u)) + 2d_{uv}^2 g(u, h(u))[I, dh(u)] \\ &\quad + d_{vv}^2 g(u, h(u))[dh(u), dh(u)] + d_v g(u, h(u))[d^2 h(u)], \end{aligned}$$

with a similar expression for $d^2 \bar{h}(u)$, with the additional presence of the term $d^2 e(u)$. We find that

$$(4.21) \quad d^2 h(0) = Nd^2 h(0)[C, C] + d_{uu}^2 g(0, 0),$$

and, because $d^2 e(0) = 0$, (4.21) also holds with $h(0)$ replaced by $\bar{h}(0)$. We deduce that the bilinear form $d^2 \Delta(0)$ satisfies

$$d^2 \Delta(0)[X, Y] = Nd^2 \Delta(0)[CX, CY] \quad (\forall X, Y \in \mathbb{R}^p),$$

so that

$$d^2 \Delta(0)[X, Y] = N^\nu d^2 \Delta(0)[C^\nu X, C^\nu Y] = 0 \quad (\forall X, Y \in \mathbb{R}^p).$$

We omit the details, but by continuing inductively and assuming $d^j e(0) = 0$, one obtains the result that the j -linear form $d^j \Delta(0)$ satisfies

$$d^j \Delta(0)[X_1, X_2, \dots, X_j] = Nd^j \Delta(0)[CX_1, CX_2, \dots, CX_j]$$

for all $X_1, \dots, X_j \in \mathbb{R}^p$. The nilpotency of N now ensures that the latter quantity is zero. \square

A simple corollary to Proposition 1 is that if (4.20) has two infinitely differentiable solutions h and \bar{h} defined on some neighborhood of zero, then they agree beyond all orders at zero:

$$\lim_{u \rightarrow 0} \frac{|h(u) - \bar{h}(u)|_q}{|u|_p^k} = 0 \quad (\forall k \geq 1).$$

5. Applications.

5.1. Nonlinear normal form. The first application of Theorem 1 is the following result, which says that (1.1) induces a local dynamical system on a manifold in the sense of Definition 2. Recall the definition of the matrix pencil (A, B) via

$$(A, B) := (d_{\bar{z}}F(0, 0), d_zF(0, 0)),$$

where F has (z, \bar{z}) as its argument, and note that the following matrix pencil defined on \mathbb{R}^{2m} :

$$(\mathcal{A}, \mathcal{B}) := \left(\left(\begin{array}{cc} I & 0 \\ 0 & 0 \end{array} \right), \left(\begin{array}{cc} 0 & I \\ B & A \end{array} \right) \right),$$

satisfies $\sigma(\mathcal{A}, \mathcal{B}) = -\sigma(A, B)$.

THEOREM 2 (nonlinear Kronecker normal form). *Let $\alpha \in [0, 1)$, and suppose that (A1)–(A3) hold; then there is a linear space K_1 and a $C^{1+\alpha}$ -manifold \mathcal{M} modeled on K_1 such that $\dim(K_1) = \#\sigma(A, B)$ and*

1. *for all $(z, w) \in \mathcal{M}$ there results $F(z, w) = 0$,*
2. *there is an r' such that for each $(z, w) \in \mathcal{M}$, with $\|(z, w)\| < r'$, there is a unique $(\bar{z}, \bar{w}) \in \mathcal{M}$ such that $w = \bar{z}$, and*
3. *there is a $C^{1+\alpha}$ -diffeomorphism $\theta : K_1 \rightarrow \mathcal{M}$ such that if $(z, w) = \theta(u)$ for some $u \in K_1$, then $(\bar{z}, \bar{w}) = \theta(\varphi(u))$, where*

$$\varphi(u) = Cu + p(u)$$

for some linear map $C : K_1 \rightarrow K_1$ that satisfies $\sigma(C) = \sigma(A, B)$. Moreover, $p : B_\delta(0; K_1) \rightarrow K_1$ satisfies $p(0) = 0$ and $d_u p(0) = 0$.

Proof. From Theorem 1, first identify the linear space K_1 from the KNF of $(\mathcal{A}, \mathcal{B})$ with \mathbb{R}^p and K_2 with \mathbb{R}^q . Now define the $C^{1+\alpha}$ -graph

$$\hat{\mathcal{M}} := \{(u, h(u)) \in \mathbb{R}^p \oplus \mathbb{R}^q : u \in \Omega_\delta\},$$

and let the r -neighborhood of zero in $\hat{\mathcal{M}}$ be $\hat{\mathcal{M}}_r := \{(u, h(u)) \in \mathbb{R}^p \oplus \mathbb{R}^q : |u|_p < r\}$ whenever $r < \delta$.

As a consequence of Theorem 1, there is an $r > 0$ such that for each $(u, v) = (u, h(u)) \in \hat{\mathcal{M}}_r$ we can find a pair (\bar{u}, \bar{v}) given by $(\bar{u}, h(\bar{u})) \in \hat{\mathcal{M}}$ such that (NF) is satisfied:

$$\bar{u} = Cu + f(u, v), \quad N\bar{v} = v - g(u, v),$$

where C is obtained from the KNF of $(\mathcal{A}, \mathcal{B})$ so that $\sigma(C) = -\sigma(\mathcal{A}, \mathcal{B}) = \sigma(A, B)$.

Let $\varphi(u) := Cu + f(u, h(u))$, and for each $u \in \mathbb{R}^p$ of sufficiently small norm set

$$\theta(u) := Q(u, h(u)) \quad \text{and} \quad \mathcal{M} \equiv Q(\hat{\mathcal{M}}), \quad \mathcal{M}_r \equiv Q(\hat{\mathcal{M}}_r),$$

where the linear map Q is taken from the KNF of $(\mathcal{A}, \mathcal{B})$ from the discussion that immediately follows (3.4). The map θ then provides a local diffeomorphism between Ω_δ and \mathcal{M} ; moreover both

$$F(\theta(u)) = 0 \quad \text{and} \quad F(\theta(\varphi(u))) = 0$$

follow by the construction of h . The following diagram illustrates how a map is induced on \mathcal{M} in this way:

$$\begin{array}{ccc} \mathcal{M}_r & \longrightarrow & \mathcal{M} \\ \theta^{-1} \downarrow & & \uparrow \theta \\ \Omega_r & \xrightarrow{\varphi} & \Omega_\delta \end{array} ,$$

where $\theta^{-1}(Q(u, h(u))) = u$, and this concludes the proof. \square

The following immediate corollaries of Theorem 1 provide some information regarding the stable and unstable behavior of orbits in a neighborhood of a fixed point of (1.1).

COROLLARY 1. *Suppose that F satisfies (A1)–(A3) and (A, B) is a regular matrix pencil with $\rho(A, B) < 1$; then there is a $C^{1+\alpha}$ -solution manifold \mathcal{M} of (1.1) containing 0 such that each $(z, \bar{z}) \in \mathcal{M}$ supports a global orbit $(z_n)_{n=0}^\infty$ with $(z_n, z_{n+1}) \in \mathcal{M}$, $z_0 = z$, $z_1 = \bar{z}$, and $\lim_{n \rightarrow \infty} z_n = 0$.*

Proof. The orbit is constructed by iterating the map $\varphi(u) = Cu + p(u)$ given in part 3 of Theorem 2: Because $\rho(A, B) < 1$ we have $\rho(C) < 1$ so that φ is a contraction near the origin in some norm, and the result follows. \square

The following are the natural definitions of stable and unstable sets associated with fixed points of (1.1); note that they are subsets of \mathbb{R}^{2m} and not \mathbb{R}^m .

DEFINITION 4. *The set*

$$W_{\text{loc}}^s(0) := \{(z, \bar{z}) \in \mathbb{R}^{2m} : \exists \text{ global orbit } (z_n)_{n=0}^\infty, z_0 = z, z_1 = \bar{z}, \lim_{n \rightarrow \infty} z_n = 0\}$$

is the local stable set associated with the zero fixed point of (1.1), and

$$W_{\text{loc}}^u(0) := \{(z, \bar{z}) \in \mathbb{R}^{2m} : \exists \text{ global orbit } (z_n)_{n=0}^\infty, z_{-1} = z, z_0 = \bar{z}, \lim_{n \rightarrow -\infty} z_n = 0\}$$

is the local unstable set.

In case $\sigma(A, B)$ contains elements outside the unit disk, one can apply the stable manifold theorem for maps to φ in Theorem 2 to give the following result.

COROLLARY 2. *Suppose that (A1)–(A3) are satisfied and (A, B) possesses n_s eigenvalues in the open unit disk; then (1.1) possesses a subset of the local stable set which is a differentiable manifold of dimension n_s .*

There is an analogous corollary to show that the unstable set is nonempty and contains a manifold of dimension n_u , where n_u is the number of elements of $\sigma(A, B)$ lying outside the closed unit disk. This result is obtained by applying Corollary 2 to (1.1) but with time running backwards.

COROLLARY 3. *Suppose that (A1)–(A3) are satisfied and (B, A) possesses n_u eigenvalues in the open unit disk; then (1.1) possesses a subset of the local unstable set which is a differentiable manifold of dimension n_u .*

Proof. Let us rewrite (1.1) in the form

$$(5.1) \quad F(z_{n-1}, z_n) = 0$$

to emphasize the fact that we are seeking an orbit that propagates backwards in time, with (z_{-1}, z_0) given. The linearization of (5.1) is of the form

$$Bz_{n-1} + Az_n$$

and if $\det B \neq 0$, then we may locally solve (5.1) for $z_{n-1} = f(z_n)$, and then one can apply the stable manifold theorem for maps to this.

On the other hand if $\det B = 0$, then conditions (A1)–(A3), appropriately modified by exchanging the roles of z and \bar{z} because time is flowing backwards, still apply to (5.1) because (B, A) is a regular matrix pencil due to the fact that (A, B) is regular. The result then follows from Corollary 2. \square

As A is a singular mapping it follows that the finite spectrum of (B, A) satisfies

$$\sigma(B, A) = (\sigma(A, B) \setminus \{0\})^{-1} \cup \{0\},$$

whether or not B is singular, and hence $\sigma(B, A)$ contains zero so that $n_u \geq 1$. An unstable manifold therefore always exists for (1.1) under conditions (A1)–(A3).

5.2. Bifurcation theorems. We now consider a C^k -mapping $F : \mathbb{R}^{2m} \times \mathbb{R} \rightarrow \mathbb{R}^{2m}$, where $k \geq 5$, and examine the family of difference equations

$$(5.2) \quad F(z_n, z_{n+1}, \mu) = 0.$$

Define the one-parameter family of matrix pencils

$$\mathcal{P}(\mu) := (A(\mu), B(\mu)) := (d_{\bar{z}}F(0, 0, \mu), d_zF(0, 0, \mu)),$$

where F has (z, \bar{z}, μ) as its argument.

THEOREM 3 (parameterized nonlinear KNF). *Suppose that $(0, 0)$ is a fixed point of (5.2) for all $\mu \in \mathbb{R}$ and that $\mathcal{P}(0)$ is a regular matrix pencil. Then there is a linear space K_1 and a $C^{1+\alpha}$ -parameter family of $C^{1+\alpha}$ -manifolds \mathcal{M}_μ modeled on K_1 such that $\dim(K_1) = \#\sigma(\mathcal{P}(0))$ and*

1. for all $(z, w) \in \mathcal{M}_\mu$ there results $F(z, w, \mu) = 0$,
2. there is an r' (independent of μ) such that for each $(z, w) \in \mathcal{M}_\mu$, with $\|(z, w)\| < r'$, there is a unique $(\bar{z}, \bar{w}) \in \mathcal{M}_\mu$ such that $w = \bar{z}$, and
3. there is a $C^{1+\alpha}$ -parameter family of $C^{1+\alpha}$ -diffeomorphisms $\theta_\mu : K_1 \rightarrow \mathcal{M}_\mu$ such that if $(z, w) = \theta_\mu(u)$ for some $u \in K_1$, then $(\bar{z}, \bar{w}) = \theta_\mu(\varphi(u, \mu))$, where

$$\varphi(u, \mu) = C(\mu)u + p(u, \mu).$$

Moreover, $C(\cdot)$ is a $C^{1+\alpha}$ -parameter family of maps in $BL(K_1)$ and, for some $\delta > 0$, $p : B_\delta(0; K_1) \times B_\delta(0; \mathbb{R}) \rightarrow K_1$ satisfies

$$p(0, \mu) \equiv 0, \quad d_u p(0, \mu) \equiv 0.$$

4. If $\lambda : [-\delta, \delta] \rightarrow \mathbb{C}$ is a continuous (and so bounded) curve, then $\lambda(\mu) \in \sigma(\mathcal{P}(\mu))$ for all $\mu \in [-\delta, \delta]$ if and only if $\lambda(\mu) \in \sigma(C(\mu))$ for all $\mu \in [-\delta, \delta]$.

Proof. Consider the suspended difference equation

$$(S) \quad \begin{cases} \mu_{n+1} = \mu_n, \\ z_{n+1} = w_n, \\ 0 = F(z_n, w_n, \mu_n). \end{cases}$$

Let us write

$$F(z, w, \mu) = A(\mu)w + B(\mu)z + \mathcal{F}(z, w, \mu),$$

where $F(0, 0, \mu) = 0$, $d_z F(0, 0, \mu) = B(\mu)$, and $d_w F(0, 0, \mu) = A(\mu)$, and then consider the following matrix pencil on \mathbb{R}^{2m} :

$$(\mathcal{A}, \mathcal{B}(\mu)) := \left(\begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & I \\ B(\mu) & A(\mu) \end{pmatrix} \right).$$

This satisfies $\sigma(\mathcal{A}, \mathcal{B}(\mu)) = -\sigma(\mathcal{P}(\mu))$ and is a regular matrix pencil when $\mu = 0$, and we can exploit this fact using the resulting Kronecker normal form to put (S) into a normal form. If we set $\mathbf{z} = (z, w)$, then we may write (S) as

$$\begin{aligned} \bar{\mu} &= \mu, \\ \mathcal{A}\bar{\mathbf{z}} &= \mathcal{B}(\mu)\mathbf{z} + \mathcal{F}(\mathbf{z}, \mu), \end{aligned}$$

where an overbar is used to denote a forward iterate. Now there is a matrix pair (P, Q) such that $PAQ = \begin{pmatrix} I_{K_1} & 0 \\ 0 & N \end{pmatrix}$ and $P\mathcal{B}(0)Q = \begin{pmatrix} C & 0 \\ 0 & I_{K_2} \end{pmatrix}$, where $Q : \mathbb{R}^{2m} \rightarrow K_1 \oplus K_2 = \mathbb{R}^{p+q}$ and N is nilpotent. With $(u, v) := \mathbf{w} = Q\mathbf{z} \in K_1 \oplus K_2$, we obtain

$$\begin{aligned} (5.3) \quad & \bar{\mu} = \mu, \\ (5.4) \quad & \bar{u} = \alpha(\mu)u + \beta(\mu)v + \mathcal{G}_1(u, v, \mu), \\ (5.5) \quad & N\bar{v} = \gamma(\mu)u + \delta(\mu)v + \mathcal{G}_2(u, v, \mu), \end{aligned}$$

where

$$\alpha : K_1 \rightarrow K_1, \quad \beta : K_2 \rightarrow K_1, \quad \gamma : K_1 \rightarrow K_2, \quad \text{and} \quad \delta : K_2 \rightarrow K_2$$

are differentiable linear maps in μ , and

$$\alpha(0) = C, \quad \beta(0) = 0, \quad \gamma(0) = 0, \quad \text{and} \quad \delta(0) = I_{K_2},$$

where $C \in BL(K_1)$ is provided by the KNF of $\mathcal{P}(0)$ and $\sigma(C) = -\sigma(\mathcal{A}, \mathcal{B}(0)) = \sigma(\mathcal{P}(0))$. Moreover, \mathcal{G}_1 and \mathcal{G}_2 represent $\mathcal{O}_2(u, v)$ -functions parameterized by μ .

Seeking an invariant manifold on which $v = h(u, \mu)$, we put (5.3)–(5.5) into the form

$$\begin{aligned} (5.6) \quad & \bar{\mu} = \mu, \\ (5.7) \quad & \bar{u} = Cu + \mathcal{O}_2(u, v, \mu), \\ (5.8) \quad & N\bar{v} = v + \mathcal{O}_2(u, v, \mu), \end{aligned}$$

where $\mathcal{O}_2(u, v, \mu)$ denotes a function of (u, v, μ) which vanishes to second or higher order at the origin. From Theorem 1 we obtain a local invariant manifold of (5.6)–(5.7) on which $v = h(u, \mu)$. Moreover $h(0, 0) = 0, dh(0, 0) = 0$, and we may assume that h is C^1 .

It follows that $h(u, \mu)$ satisfies the functional equation

$$Nh(\alpha u + \beta h(u, \mu) + \mathcal{G}_1(u, h(u, \mu), \mu)) = \gamma(\mu)u + \delta h(u, \mu) + \mathcal{G}_2(u, h(u, \mu), \mu),$$

and if we set $x = h(0, \mu)$, then x satisfies the equation

$$Nh(\beta x + \mathcal{G}_1(0, x, \mu), \mu) = \delta(\mu)x + \mathcal{G}_2(0, x, \mu).$$

The latter is an algebraic equation for x which we denote $a(x, \mu) = 0$; moreover $a(0, \mu) = 0$ holds for all μ near 0, whence $h(0, \mu) \equiv 0$. In addition, a short calculation shows that

$$d_x a(0, \mu) = \delta(\mu) - Nd_u h(0, \mu)[\beta(\mu)],$$

which is an identity mapping when $\mu = 0$. The implicit function theorem now ensures that $x = 0$ is the only solution of $a(x, \mu) = 0$ for all μ near 0.

The functional equation satisfied by $d_u h(u, \mu)$ is then

$$-Nd_u h(\alpha u + \beta h + \mathcal{G}_1(u, h, \mu))[\alpha + \beta d_u h + d_u \mathcal{G}_1 + d_v \mathcal{G}_1 \cdot d_u h] + \gamma + \delta d_u h + d_u \mathcal{G}_2 + d_v \mathcal{G}_2 \cdot d_u h = 0,$$

where various arguments have been omitted for brevity. If we write τ for the linear map $d_u h(0, \mu)$, then

$$\gamma(\mu) + \delta(\mu)\tau = N\tau \cdot [\alpha(\mu) + \beta(\mu)\tau].$$

This is a Riccati equation for τ that can be solved near $\mu = 0$ for τ as a function of μ using the implicit function theorem and the properties of α, β, γ , and δ . The result that $\tau(0) = 0$ then follows because N is nilpotent and $\tau(0) = N\tau(0)[C]$.

We are now in a position to define the one-parameter family of matrices, denoted by $C(\mu)$ in the statement of the theorem, namely,

$$C(\mu) := \alpha(\mu) + \beta(\mu)\tau(\mu),$$

so that $C(0) = C$. If $C(\mu)$ has an eigenvalue λ , say, then

$$(\alpha + \beta\tau)w = \lambda w \implies \gamma + \delta\tau = N\tau \cdot [\lambda w],$$

whence

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \begin{bmatrix} w \\ \tau w \end{bmatrix} = \lambda \begin{bmatrix} w \\ N \cdot \tau w \end{bmatrix},$$

and therefore

$$-\lambda \in \sigma \left(\begin{bmatrix} I & 0 \\ 0 & N \end{bmatrix}, \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \right) = \sigma(\mathcal{A}, \mathcal{B}(\mu)) = -\sigma(\mathcal{P}(\mu)).$$

We have deduced that

$$\sigma(\alpha(\mu) + \beta(\mu)\tau(\mu)) \subseteq \sigma(\mathcal{P}(\mu)),$$

but the left-hand side of this inclusion has $\dim(K_1)$ elements counted according to algebraic multiplicity, whereas the right-hand side may have more unless, that is, $\mu = 0$, in which case the inclusion is replaced by an equality because $\sigma(\alpha(0)) = \sigma(C) = \sigma(\mathcal{P}(0))$.

As a result, if $\lambda(\mu) \in C([-\delta, \delta], \mathbb{C})$ is a continuous path of eigenvalues of $\mathcal{P}(\mu)$, then by virtue of the fact that $\lambda(0) \in \sigma(\alpha(\mu) + \beta(\mu)\tau(\mu)|_{\mu=0})$, it follows by the continuous dependence of eigenvalues on μ and by counting their location in the complex plane that $\lambda(\mu) \in \sigma(\alpha(\mu) + \beta(\mu)\tau(\mu)|_{\mu \neq 0})$. \square

Conclusion 4 of Theorem 3 is not equivalent to saying that a locus of eigenvalues of $\mathcal{P}(\mu)$ is necessarily a locus of eigenvalues of $C(\mu)$. This is because $\mathcal{P}(\mu)$ may have other eigenvalues for small μ which become unbounded as μ tends to zero, and such curves cannot correspond to eigenvalues of $C(\mu)$. This happens, for instance, in the singularity-induced bifurcation theorem from [1] because an eigenvalue has a pole with respect to the bifurcation parameter.

5.2.1. Existence of bifurcations for (1.1). One can easily prove a period-doubling bifurcation theorem for (1.1) without having recourse to the nonlinear normal form given in Theorem 2, but we now include this result for completeness.

THEOREM 4. *Suppose that (5.2) satisfies the following hypotheses:*

1. $1 \notin \sigma(\mathcal{P}(0))$ and $F(0, 0, \mu) \equiv 0$,
2. $\ker(-A(0) + B(0)) = \langle k \rangle$, so that $-1 \in \sigma(\mathcal{P}(0))$, and
3. $B'(0)k \notin \text{im}(-A(0) + B(0))$.

Then $\mu = 0$ is a period-doubling bifurcation point for (5.2) from the trivial solution $z = 0$.

Proof. Consider the algebraic equation $G(z, w, \mu) = 0$, where $G : \mathbb{R}^{2m+1} \rightarrow \mathbb{R}^{2m}$ is given by

$$G(z, w, \mu) := \begin{bmatrix} F(z, w, \mu) \\ F(w, z, \mu) \end{bmatrix},$$

and moreover G has the trivial solution branch. Also define $\overline{G}(z, \mu) := F(z, z, \mu)$, so that

$$d_{(z,w)}G(0, 0, \mu) = \begin{bmatrix} A(\mu) & B(\mu) \\ B(\mu) & A(\mu) \end{bmatrix},$$

and $d_z\overline{G}(0, \mu) = A(\mu) + B(\mu)$. By assumption, $A(0) + B(0)$ is invertible, and therefore $d_z\overline{G}(0, \mu)$ is an invertible map for small $|\mu|$, so that if $G(z, w, \mu) = 0$ near $\mu = 0$, then $z \neq w$ unless $z = w = 0$. The theorem now follows from the simple eigenvalue bifurcation theorem applied to G at $\mu = 0$, noting that the kernel of $d_{(z,w)}G(0, 0, 0)$ is $(k, -k)^T$. \square

One can of course formulate a similarly straightforward fold bifurcation for (1.1) in an entirely analogous fashion. However, the following theorem relies on Theorem 3 in a nontrivial way.

THEOREM 5 (Neimark–Sacker bifurcation). *Suppose that (5.2) has the fixed point $z = 0$ for all $\mu \in \mathbb{R}$ and that $\lambda(\mu) \in \sigma(\mathcal{P}(\mu))$ is a curve which satisfies the following:*

1. $\mathcal{P}(0)$ is a regular matrix pencil;
2. $|\lambda(0)| = 1$, and $\lambda(0)$ is an algebraically simple eigenvalue of $\mathcal{P}(0)$;
3. $\lambda(0)^n \neq 1$ for $n \in \{1, 2, 3, 4\}$;
4. $\frac{d}{d\mu}|\lambda(\mu)|\big|_{\mu=0} \neq 0$.

Then modulo a further nonresonance condition¹ there is a half-interval $J \subset \mathbb{R}$ containing 0 in its closure such that (5.2) possesses a quasi-invariant circle $\Gamma_\mu \subset \mathbb{R}^{2m}$ for all $\mu \in J$. Moreover, if $\text{diam}(\Gamma_\mu) = \sup\{\|z - w\| : z, w \in \Gamma_\mu\}$, then $\lim_{\mu \rightarrow 0} \text{diam}(\Gamma_\mu) = 0$.

Proof. Theorem 5 follows immediately from the Neimark–Sacker bifurcation for maps applied to $\varphi(u, \mu)$ from Theorem 3 (part 3). \square

6. Examples. *Example 1* (output-nulling control problem). The results in this paper give sufficient conditions for a positive answer to the following question:

- (Q) *Given $f(= f(x, u)) : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^n, g(= g(x)) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that $f(0, 0) = 0, g(0) = 0$, does there exist a sequence of states (x_n) given by the iterates of f and controls (u_n) such that $g(x_n) \equiv 0$ and $x_n \rightarrow 0$ as $n \rightarrow \infty$?*

Thus, we seek a global orbit of the ΔAE

$$x_{n+1} = f(x_n, u_n), \quad g(x_n) = 0.$$

¹See [21] or [25, p. 376], where the open condition “ $a \neq 0$ ” is given, and this requires the computation of third-order terms in the normal form for this bifurcation.

A necessary condition for the existence of such a solution can be obtained by substituting the dynamic part of the problem into the constraint, to give the hidden constraint

$$g(f(x_n, u_n)) = 0 \quad (\forall n \geq 1).$$

Hence, provided the function $g(f(x, u))$ has an invertible partial u -derivative at $(x, u) = (0, 0)$, by the stable manifold theorem the response to **(Q)** is affirmative if the spectrum of the x -derivative of $f(x, U(x))$, also evaluated at $(x, u) = (0, 0)$, contains an element of the open unit disk. Here $U(x)$ denotes the solution of the equation $g(f(x, U)) = 0$ given locally by the implicit function theorem. Note for a moment that if the stated u -derivative $dg(0)d_u f(0, 0)$ is invertible, it follows that the matrix pencil

$$(A, B) := \left(\left(\begin{array}{cc} I_x & 0 \\ 0 & 0 \end{array} \right), \left(\begin{array}{cc} d_x f & d_u f \\ dg & 0 \end{array} \right) \right) \Big|_{(x,u)=(0,0)}$$

is regular, and thus (A1)–(A3) hold for this problem.

However, by using Theorem 1 one can dispense with the condition that $dg(0)d_u f(0, 0)$ is invertible. In fact, let us assume that $K := \ker(dg(0)d_u f(0, 0)) \neq \{0\}$. In this case, the matrix $\lambda A + B$ is invertible if $-\lambda \notin \sigma(d_x f(0, 0))$ and $dg(\lambda I_x + d_x f)^{-1}d_u f$ is also invertible at $(x, u) = (0, 0)$. However, for λ large we have

$$\begin{aligned} \lambda^2 dg(\lambda I_x + d_x f)^{-1}d_u f &= \lambda dg(I_x + \lambda^{-1}d_x f)^{-1}d_u f \\ &= \lambda dg(I_x - \lambda^{-1}d_x f + O(\lambda^{-2}))d_u f \\ &= \lambda dg \cdot d_u f - dg \cdot d_x f \cdot d_u f + O(\lambda^{-1}), \end{aligned}$$

evaluating all of the stated derivatives at $(x, u) = (0, 0)$. As a result, if the weaker condition holds that the pencil $(dg \cdot d_u f, dg \cdot d_x f \cdot d_u f)$ is regular, then (A, B) is regular and (A1)–(A3) still apply. The response to **(Q)** is again affirmative if $\sigma(A, B)$ contains at least one member of the open unit disk.

In fact one can show that output-nulling control problems are well-posed in sequence spaces as follows. First consider the linear problem

$$(6.1) \quad Az_{n+1} + Bz_n = \Gamma_n,$$

where $n \in \mathbb{N}$ and (Γ_n) is a given sequence in

$$\ell_{\mathbb{N}}^{\infty}(\mathbb{R}^m) = \left\{ (z_n)_{n \in \mathbb{N}} : z_n \in \mathbb{R}^m, \sup_{n \in \mathbb{N}} \|z_n\|_{\mathbb{R}^m} < \infty \right\},$$

but $\det A = 0$. If (A, B) is a regular matrix pencil with index ν , the KNF allows us to write (6.1) in the form

$$(6.2) \quad u_{n+1} = Cu_n + \alpha_n,$$

$$(6.3) \quad Nv_{n+1} = v_n + \beta_n,$$

where $(u_n, v_n) \in \mathbb{R}^{p+q}$. In order to solve (6.1), let us consider the linear operator $I - N\sigma$ on a space of sequences $\ell_{\mathbb{N}}^{\infty}(\mathbb{R}^q)$. We take linear maps $T \in BL(\mathbb{R}^q)$ to act pointwise on $\ell_{\mathbb{N}}^{\infty}(\mathbb{R}^q)$, so

$$T(w_n)_{n \in \mathbb{N}} = (Tw_n)_{n \in \mathbb{N}} \quad (\forall (w_n)_{n \in \mathbb{N}} \in \ell_{\mathbb{N}}^{\infty}(\mathbb{R}^q)),$$

and we define σ the forward-shift map by

$$\sigma(w_n)_{n \in \mathbb{N}} = (w_{n+1})_{n \in \mathbb{N}} \quad (\forall (w_n)_{n \in \mathbb{N}} \in \ell_{\mathbb{N}}^{\infty}(\mathbb{R}^q)).$$

Any such T will commute with σ , and one can see by a direct multiplication that

$$(I - N\sigma) \sum_{i=0}^{\nu-1} N^i \sigma^i = I,$$

where I denotes the identity on $\ell_{\mathbb{N}}^{\infty}(\mathbb{R}^q)$.

One can solve (6.2) in a suitably weighted sequence space if no restrictions are to be placed on the spectrum of (A, B) . Equation (6.3) can also be solved:

$$\mathbf{v} = -(I - N\sigma)^{-1} \beta = - \sum_{i=0}^{\nu-1} N^i \sigma^i \beta,$$

where $\mathbf{v} = (v_i)_{i \in \mathbb{N}}$ and $\beta = (\beta_i)_{i \in \mathbb{N}}$, whence

$$v_n = - \sum_{i=0}^{\nu-1} N^i \beta_{n+i}.$$

So from a temporal point of view the current values of the state depend on future values of the input, but (6.1) is still well-posed in a sequence space as \mathbf{v} depends continuously on β .

This means that a second- or higher-order nonlinear perturbation of (6.1),

$$(6.4) \quad Az_{n+1} + Bz_n + \mathcal{F}(z_n) = \Gamma_n,$$

say, where $\mathcal{F}(0) = 0, d\mathcal{F}(0) = 0$, can be written as an infinite-dimensional problem

$$(6.5) \quad (A\sigma + B)\mathbf{z} + \mathcal{F}(\mathbf{z}) = \mathbf{\Gamma}$$

in $\ell_{\mathbb{N}}^{\infty}(\mathbb{R}^m)$, and one can apply the inverse function theorem to solve locally for small-norm solutions of the form $\mathbf{z} = \mathbf{z}(\mathbf{\Gamma})$, where $\mathbf{z}(\mathbf{0}) = \mathbf{0}$.

This solution can be found via the Picard iteration $\mathbf{z}(\mathbf{\Gamma}) = \lim_{n \rightarrow \infty} \mathbf{y}^{(n)}$, where $\mathbf{y}^{(0)} = \mathbf{0}$ and

$$\mathbf{y}^{(n+1)} = -(A\sigma + B)^{-1} [\mathcal{F}(\mathbf{y}^{(n)}) - \mathbf{\Gamma}].$$

As a result, writing the solution sequence $\mathbf{z}(\mathbf{\Gamma})$ as $(z_n(\mathbf{\Gamma}))_{n \in \mathbb{N}}$, it is clear that the nonlinear perturbation will have the effect of making each z_n depend on infinitely many elements of the sequence $\mathbf{\Gamma}$, unless ν happens to equal 1.

This effect has been observed before in the literature in the context of delay DAEs [5, 13], where it is noted in the former reference that linear systems of delay DAEs can act like advanced systems when their index is two or higher. The problem (6.5) is displaying exactly this behavior.

Example 2. This example serves to illustrate how we can use Theorem 1 to deduce qualitative similarities between a DAE and its discrete counterpart. Take the DAE

$$(6.6) \quad \dot{x} = f(x, y),$$

$$(6.7) \quad g(x, y) = 0,$$

subject to $x \in \mathbb{R}^n, y \in \mathbb{R}^m$ with an equilibrium at the origin, so that $f(0, 0) = 0, g(0, 0) = 0$. Let us impose a singularity of the form

$$(6.8) \quad \ker(d_y g(0, 0)) = \langle k \rangle \text{ such that } d_y f(0, 0)d_x g(0, 0)k \notin \text{im}(d_y g(0, 0)),$$

and the equilibrium solution is isolated:

$$\det \begin{pmatrix} d_x f(0, 0) & d_y f(0, 0) \\ d_x g(0, 0) & d_y g(0, 0) \end{pmatrix} \neq 0.$$

From [2] it is known that (6.6)–(6.7) has an invariant manifold W of dimension $n - 1$ that contains the origin and intersects the singularity in an $n - 2$ -dimensional manifold of pseudoequilibria.

Now consider the forward-Euler method in state-space form [8, p. 375] applied to (6.6)–(6.7), resulting in the difference equation

$$(6.9) \quad x_{i+1} - x_i = hf(x_i, y_i),$$

$$(6.10) \quad g(x_{i+1}, y_{i+1}) = 0.$$

Using Theorem 1, in order to show that (6.9)–(6.10) possesses a quasi-invariant manifold of solutions W_h that contains the origin and has dimension $n - 1$, we need only show that (A1)–(A3) hold, which entails showing that the derivative at the origin of (6.9)–(6.10) is a regular matrix pencil. Hence we seek a $\xi \in \mathbb{C}$ such that

$$\det \left(\xi \begin{bmatrix} I & 0 \\ d_x g & d_y g \end{bmatrix} + \begin{bmatrix} I + hd_x f & hd_y f \\ 0 & 0 \end{bmatrix} \right) \Big|_{(x,y)=(0,0)} \neq 0.$$

However, the conditions in (6.8) ensure the existence of such a ξ for any fixed $h \neq 0$, and the existence of W_h follows. The dimension of W_h comes from counting the number of finite eigenvalues of the linearization of (6.9)–(6.10) which is given in [1] as $n - 1$.

Example 3. Consider again (6.6)–(6.7) but now with a parameter α

$$(6.11) \quad \dot{x} = f(x, y, \alpha), g(x, y, \alpha) = 0,$$

and suppose that $(x, y) = (0, 0)$ is an equilibrium locus for all $\alpha \in \mathbb{R}$. Now suppose that the conditions for a Hopf bifurcation are formally satisfied by (6.11) at $\alpha = \alpha_0$: $\omega(\alpha) \in \sigma(M, -L(\alpha))$ is a locus of algebraically simple eigenvalues, where

$$M = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}, L(\alpha) = \begin{bmatrix} d_x f(0, 0, \alpha) & d_y f(0, 0, \alpha) \\ d_x g(0, 0, \alpha) & d_y g(0, 0, \alpha) \end{bmatrix};$$

and $\omega(\alpha_0) = i\omega_0$ is an eigenvalue of $(M, -L(\alpha))$ such that $\frac{d}{d\alpha} \Re(\omega(\alpha)) \Big|_{\alpha=\alpha_0}$ is non-zero. Also, no other eigenvalues of the regular matrix pencil $(M, -L(\alpha_0))$ have a zero real part.

Then, modulo an open condition on the third-order derivatives of the smooth functions f and g , we can show that the backward-Euler method

$$(6.12) \quad x_{i+1} - x_i = hf(x_{i+1}, y_{i+1}),$$

$$(6.13) \quad -h \cdot g(x_{i+1}, y_{i+1}) = 0,$$

satisfies the conditions of the Neimark–Sacker bifurcation theorem for all sufficiently small $h > 0$. Note that $(M - hL(\alpha), -M)$ is the linearization of (6.12)–(6.13) about the zero fixed point and is a regular matrix pencil for $h > 0$ and $\alpha \approx \alpha_0$ such that

$$(1 - h\omega(\alpha))^{-1} \in \sigma(M - hL(\alpha), -M) \quad (\forall \alpha).$$

In order to apply Theorem 5 to (6.12)–(6.13) we first note that the eigenvalue locus $(1 - h\omega(\alpha))^{-1}$ has unit length if and only if $1 - h\omega(\alpha)$ has unit length. If we define functions $R(\alpha)$ and $I(\alpha)$ by $\omega(\alpha) = R(\alpha) + iI(\alpha)$, then $|1 - h\omega(\alpha)| = 1$ if and only if $(1 - hR(\alpha))^2 + h^2I(\alpha)^2 = 1$, which holds when

$$(6.14) \quad h \left(-R(\alpha) + \frac{h}{2}(R(\alpha)^2 + I(\alpha)^2) \right) = 0.$$

As a result, we define the function $b(\alpha, h) := -R(\alpha) + \frac{h}{2}(R(\alpha)^2 + I(\alpha)^2)$ and note that $b(\alpha, h) = 0$ for $h > 0$ ensures that the linearization of (6.12)–(6.13) at $(0, 0)$ has an eigenvalue of unit modulus. Now $b(\alpha_0, 0) = 0$ and $\frac{\partial b}{\partial \alpha}(\alpha_0, 0) = -R'(\alpha_0) \neq 0$ by assumption, and, as a result, one may solve $b(\alpha, h) = 0$ locally using the implicit function theorem for $\alpha = \alpha(h)$ such that $\alpha(0) = \alpha_0$.

From this calculation one can show that the numerical scheme (6.12)–(6.13) has a quasi-invariant circle for α in some half-neighborhood of $\alpha(h) = \alpha_0 + \frac{h\omega_0^2}{R'(\alpha_0)} + O(h^2)$ provided $h > 0$ is sufficiently small.

Note that no assumption is made regarding the invertibility of $d_y g(0, 0, 0)$. If this mapping were invertible in addition to the formal conditions given above for Hopf bifurcation, then the existence of a locus of periodic solutions of (6.11) could be deduced. However, without the invertibility of $d_y g(0, 0, 0)$, it is not known whether a Hopf bifurcation occurs in (6.11), but (6.12)–(6.13) has a locus of invariant circles nevertheless.

6.1. Concluding remark. The results in this paper can be used to show that second-order problems of the form

$$(6.15) \quad F(z_n, z_{n+1}, z_{n+2}) = 0$$

have quasi-invariant manifolds provided that the appropriate matrix pencil is regular simply by rewriting (6.15) as a first-order problem. However, the methods of this paper do not easily extend to the study of invariant manifolds of the system that one would like to study if (1.1) had a period-2 orbit (z, w, z, w, \dots) , namely, the system

$$(6.16) \quad F(z_n, z_{n+1}) = 0,$$

$$(6.17) \quad F(z_{n+1}, z_{n+2}) = 0,$$

where $F(z, w) = 0$ and $F(w, z) = 0$, with $w \neq z$. Another approach is needed to answer the question of whether overdetermined systems of this type have any invariant manifolds associated with them.

REFERENCES

- [1] R. E. BEARDMORE, *The singularity-induced bifurcation and its Kronecker normal form*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 126–137.
- [2] R. E. BEARDMORE AND R. LAISTER, *The flow of a DAE near a singular equilibrium*, SIAM J. Matrix Anal. Appl., 24 (2002), pp. 106–120.

- [3] K. E. BRENNAN, S. L. CAMPBELL, AND L. R. PETZOLD, *Numerical Solution of Initial-Value Problems in Differential-Algebraic Equations*, North-Holland, Amsterdam, 1989.
- [4] R. BRU, C. COLL, AND E. SANCHEZ, *Structural properties of positive linear time-invariant difference-algebraic equations*, *Linear Algebra Appl.*, 349 (2002), pp. 1–10.
- [5] S. L. CAMPBELL, *Singular linear systems of differential equations with delays*, *Appl. Anal.*, 11 (1980), pp. 129–136.
- [6] T. FLIEGNER, Ü. KOTTA, AND H. NIJMEIJER, *Solvability and right-inversion of implicit nonlinear discrete-time systems*, *SIAM J. Control Optim.*, 34 (1996), pp. 2092–2115.
- [7] F. R. GANTMACHER, *Theory of Matrices*, Vol. 2, Chelsea, New York, 1977.
- [8] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations*, 2nd revised ed., Springer Ser. Comput. Math. II, Springer, New York, 2002.
- [9] T. KACZYNSKI AND W. KRAWCEWICZ, *A local Hopf bifurcation theorem for a certain class of implicit differential equations*, *Canad. Math. Bull.*, 36 (1993), pp. 183–189.
- [10] P. E. KLOEDEN AND P. MARÍN-RUBIO, *Weak pullback attractors of non-autonomous difference inclusions*, *J. Difference Equ. Appl.*, 9 (2003), pp. 489–502.
- [11] J. LAURENT-VARIN, F. BONNANS, N. BEREND, C. TALBOT, AND M. HADDOU, *On the refinement of discretization for optimal control problems*, in *Proceedings of the 16th IFAC Symposium on Automatic Control in Aerospace*, St. Petersburg, 2004.
- [12] J.-Y. LIN AND Z.-H. YANG, *A discrete optimal control for descriptor systems*, *IEEE Trans. Automat. Control*, 34 (1989), pp. 177–181.
- [13] H. LOGEMAN, *Destabilizing effects of small time delays on feedback-controlled descriptor systems*, *Linear Algebra Appl.*, (1998), pp. 131–153.
- [14] C. NAVASCA, *Local stable manifold for the bidirectional discrete-time dynamics*, *SIAM J. Control Optim.*, submitted.
- [15] P. J. RABIER, *The Hopf bifurcation theorem for quasilinear differential-algebraic equations*, *Comput. Methods Appl. Mech. Engrg.*, 170 (1999), pp. 355–371.
- [16] P. J. RABIER AND W. C. RHEINBOLDT, *On impasse points of quasilinear differential-algebraic equations*, *Math. Anal. Appl.*, 181 (1994), pp. 429–454.
- [17] P. J. RABIER AND W. C. RHEINBOLDT, *Theoretical and Numerical Analysis of Differential-Algebraic Equations*, *Handbook of Numerical Analysis*, Vol. VIII, Elsevier Science, New York, 2003, Part 4, pp. 183–540.
- [18] S. REICH, *On the qualitative behaviour of DAEs*, *Circuits Systems Signal Proc.*, 14 (1995), pp. 427–443.
- [19] G. REIßIG AND H. BOCHE, *On singularities of autonomous implicit ordinary differential equations*, *IEEE CAS I*, 50 (2003), pp. 922–931.
- [20] R. RIAZA AND P. J. ZUFIRIA, *Stability of singular equilibria in quasilinear implicit differential equations*, *J. Differential Equations*, 171 (2001), pp. 24–53.
- [21] R. J. SACKER, *On Invariant Surfaces and Bifurcation of Periodic Solutions of Ordinary Differential Equations*, Ph.D. thesis, New York University, Courant Institute of Mathematical Sciences, 1964.
- [22] G. SÖDERLIND, *Remarks on the stability of high-index DAEs with respect to parametric perturbations*, *Computing*, 49 (1992), pp. 303–314.
- [23] J. SOTOMAYOR AND M. ZHITOMIRSKII, *Impasse singularities of differential systems of the form $A(x)\dot{x} = f(x)$* , *J. Differential Equations*, 169 (2001), pp. 567–587.
- [24] V. VENKATASUBRAMANIAN, *Singularity induced bifurcation in the Van Der Pol oscillator*, *IEEE Trans. Circuits Syst. I Fund. Theory Appl.*, 41 (1994), pp. 765–769.
- [25] S. WIGGINS, *Introduction to Applied Nonlinear Dynamical Systems and Chaos*, *Texts Appl. Math.*, Springer, New York, 1990.
- [26] E. ZEIDLER, *Nonlinear Analysis: Linear Monotone Operators*, Vol. II, Springer, New York, 1990.

ON THE WELL-POSEDNESS FOR THE VISCOUS SHALLOW WATER EQUATIONS*

QIONGLEI CHEN[†], CHANGXING MIAO[†], AND ZHIFEI ZHANG[‡]

Abstract. In this paper, we prove the existence and uniqueness of the solutions for the two-dimensional viscous shallow water equations with low regularity assumptions on the initial data as well as the initial height bounded away from zero.

Key words. shallow water equations, well-posedness, Bony’s paraproduct decomposition, weight Besov space

AMS subject classifications. 35Q35, 76D

DOI. 10.1137/060660552

1. Introduction. In this paper, we study the two-dimensional (2D) viscous shallow water equations with a more general diffusion,

$$(1.1) \quad \begin{cases} h_t + \operatorname{div}(hu) = 0, \\ h(u_t + u \cdot \nabla u) - \nu \nabla \cdot (hD(u)) - \nu \nabla (h \operatorname{div}(u)) + h \nabla h = 0, \\ u(0, \cdot) = u_0, \quad h(0, \cdot) = h_0, \end{cases}$$

where $h(t, x)$ is the height of fluid surface, $u(t, x) = (u_1(t, x), u_2(t, x))$ is the horizontal velocity vector field, $D(u) = \frac{1}{2}(\nabla u + \nabla u^t)$ is the deformation tensor, and $\nu > 0$ is the viscous coefficient.

Recently, the viscous shallow water equations have been widely studied by mathematicians; see the review paper [4]. Bui [5] proved the local existence and uniqueness of classical solutions to the Cauchy–Dirichlet problem for shallow water equations with initial data h_0, u_0 in Hölder spaces as well as h_0 bounded away from a vacuum. Kloeden [17] and Sundbye [20] independently proved global existence and uniqueness of classical solutions to the Cauchy–Dirichlet problem in Sobolev spaces. Later, Sundbye [21] also proved global existence and uniqueness of classical solutions to the Cauchy problem. However, for all of the above results (except those of [5]), the authors only consider the case when the initial data h_0 is a small perturbation of some positive constant \bar{h}_0 and u_0 is small in some sense. Very recently, Wang and Xu [23] proved the local well-posedness of the Cauchy problem in Sobolev spaces for the large data u_0 and h_0 closing to \bar{h}_0 . More precisely, they obtained the following result.

THEOREM 1.1 (see [23]). *Let \bar{h}_0 be a strictly positive constant and $s > 2$. Assume that*

- (i) $(u_0, h_0 - \bar{h}_0) \in H^s(\mathbb{R}^2) \otimes H^s(\mathbb{R}^2)$;
- (ii) $\|h_0 - \bar{h}_0\|_{H^s} \ll \bar{h}_0$.

*Received by the editors May 23, 2006; accepted for publication (in revised form) September 17, 2007; published electronically May 16, 2008.

<http://www.siam.org/journals/sima/40-2/66055.html>

[†]Institute of Applied Physics and Computational Mathematics, P.O. Box 8009, Beijing 100088, People’s Republic of China (chen.qionglei@iapcm.ac.cn, miao_changxing@iapcm.ac.cn). The research of these authors was partially supported by the NSF of China (grants 10701012 and 10725102).

[‡]School of Mathematical Science, Peking University, Beijing 100871, People’s Republic of China (zffzhang@math.pku.edu.cn). This author’s research was supported by the NSF of China (grant 10601002).

Then there exist a positive time T and a unique solution (u, h) of (1.1) such that

$$(1.2) \quad u, h - \bar{h}_0 \in L^\infty([0, T], H^s), \quad \nabla u \in L^2([0, T]; H^s).$$

Moreover, there exists a strictly positive constant c such that if

$$(1.3) \quad \|u_0\|_{H^s} + \|h_0 - \bar{h}_0\|_{H^s} \leq c,$$

then we can choose $T = +\infty$.

One purpose of this paper is to study the well-posedness of (1.1) for the initial data with the minimal regularity. For the incompressible Navier–Stokes equations, such research has been initiated by Fujita and Kato [16]; see also [6, 7, 18] for other relevant results. They proved local well-posedness for the incompressible Navier–Stokes equations in the scaling invariant space. The scaling invariance means that if (u, p) is a solution of the incompressible Navier–Stokes equations with initial data $u_0(x)$, then

$$(1.4) \quad u_\lambda(t, x) \triangleq \lambda u(\lambda^2 t, \lambda x), \quad p_\lambda(t, x) \triangleq \lambda^2 p(\lambda^2 t, \lambda x)$$

is also a solution of the incompressible Navier–Stokes equations with $u_{0,\lambda} \triangleq \lambda u_0(\lambda x)$. Obviously, $\dot{H}^{\frac{d}{2}-1}(\mathbb{R}^d)$ is a scaling invariant space under the scaling of (1.4), i.e.,

$$\|u_\lambda\|_{\dot{H}^{\frac{d}{2}-1}} = \|u\|_{\dot{H}^{\frac{d}{2}-1}}.$$

Equations (1.1) have no scaling invariance like the incompressible Navier–Stokes equations. However, due to the similarity of the structure between (1.1) and the incompressible Navier–Stokes equations, we still solve (1.1) for initial data whose regularity fits the scaling of (1.4). It should be pointed out that Danchin was the first to consider the same problem for the compressible Navier–Stokes equations, and some ideas of this paper are motivated by [11].

The second purpose of this paper is to prove the local well-posedness of (1.1) under the more natural assumption that the initial height is bounded away from zero. For the initial data with slightly higher regularity, this can be easily obtained by modifying the argument of Danchin [13]. However, for the initial data with low regularity, his method is not applicable anymore, since the proof of [13] relies on the fact that some profits can be gained from the inclusion map $B^s \hookrightarrow L^\infty$ in the case of $s > \frac{d}{2}$. For this reason, we have to introduce some kind of weighted Besov space E_T^s (see section 3), which is crucial to eliminating the condition that the initial height h_0 is close to \bar{h}_0 . One important observation is that the E_T^s norm of the solution is small for small time T .

Before stating our main result, let us first introduce some notations and definitions. Choose a radial function $\varphi \in \mathcal{S}(\mathbb{R}^d)$ such that

$$\text{supp } \varphi \subset \left\{ \xi \in \mathbb{R}^d; \frac{5}{6} \leq |\xi| \leq \frac{12}{5} \right\}, \quad \sum_{k \in \mathbb{Z}} \varphi(2^{-k}\xi) = 1, \quad \xi \in \mathbb{R}^d \setminus \{0\}.$$

Here $\varphi_k(\xi) = \varphi(2^{-k}\xi)$, $k \in \mathbb{Z}$.

DEFINITION 1.1. *Let $k \in \mathbb{Z}$. The Littlewood–Paley projection operators Δ_k and S_k are defined as follows:*

$$\Delta_k f = \varphi(2^{-k}D)f, \quad S_k f = \sum_{j \leq k-1} \Delta_j f \quad \text{for } f \in \mathcal{S}'(\mathbb{R}^d).$$

We denote the space $\mathcal{Z}'(\mathbb{R}^d)$ by the dual space of $\mathcal{Z}(\mathbb{R}^d) = \{f \in \mathcal{S}(\mathbb{R}^d); D^\alpha \hat{f}(0) = 0; \forall \alpha \in \mathbb{N}^d \text{ multi-index}\}$; it also can be identified by the quotient space of $\mathcal{S}'(\mathbb{R}^d)/\mathcal{P}$ with the polynomials space \mathcal{P} . The formal equality

$$f = \sum_{k \in \mathbb{Z}} \Delta_k f$$

holds true for $f \in \mathcal{Z}'(\mathbb{R}^d)$ and is called the homogeneous Littlewood–Paley decomposition. It has nice properties of quasi-orthogonality: with our choice of φ ,

$$(1.5) \quad \Delta_j \Delta_k f = 0 \quad \text{if } |j - k| \geq 2 \quad \text{and} \quad \Delta_j (S_{k-1} \Delta_k f) = 0 \quad \text{if } |j - k| \geq 4.$$

DEFINITION 1.2. *Let $s \in \mathbb{R}$, $1 \leq p, r \leq +\infty$. The homogeneous Besov space $\dot{B}_{p,r}^s$ is defined by*

$$\dot{B}_{p,r}^s = \{f \in \mathcal{Z}'(\mathbb{R}^d) : \|f\|_{\dot{B}_{p,r}^s} < +\infty\},$$

where

$$\|f\|_{\dot{B}_{p,r}^s} = \begin{cases} \left(\sum_{k \in \mathbb{Z}} 2^{ksr} \|\Delta_k f\|_p^r \right)^{\frac{1}{r}} & \text{for } r < +\infty, \\ \sup_{k \in \mathbb{Z}} 2^{ks} \|\Delta_k f\|_p & \text{for } r = +\infty. \end{cases}$$

If $p = r = 2$, $\dot{B}_{2,2}^s = \dot{H}^s$, and if $d = 2$, we have $\dot{B}_{2,1}^1 \hookrightarrow L^\infty$ and

$$\|f\|_\infty \leq C \|f\|_{\dot{B}_{2,1}^1}.$$

We refer to [8, 22] for more details.

In addition to the general time-space space such as $L^\rho(0, T; \dot{B}_{p,r}^s)$, we introduce a useful mixed time-space homogeneous Besov space $\tilde{L}_T^\rho(\dot{B}_{p,r}^s)$, which was initiated in [10] and will be used in the proof of the uniqueness.

DEFINITION 1.3. *Let $s \in \mathbb{R}$, $1 \leq p, r, \rho \leq +\infty$, $0 < T \leq +\infty$. The mixed time-space homogeneous Besov space $\tilde{L}_T^\rho(\dot{B}_{p,r}^s)$ is defined by*

$$\tilde{L}_T^\rho(\dot{B}_{p,r}^s) = \{f \in \mathcal{Z}'(\mathbb{R}^{d+1}) : \|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,r}^s)} < +\infty\},$$

where

$$\|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,r}^s)} = \left\| 2^{ks} \left(\int_0^T \|\Delta_k f(t)\|_p^\rho dt \right)^{\frac{1}{\rho}} \right\|_{\ell^r}.$$

Using the Minkowski inequality, it is easy to verify that

$$L_T^\rho(\dot{B}_{p,r}^s) \subseteq \tilde{L}_T^\rho(\dot{B}_{p,r}^s) \quad \text{if } \rho \leq r \quad \text{and} \quad \tilde{L}_T^\rho(\dot{B}_{p,r}^s) \subseteq L_T^\rho(\dot{B}_{p,r}^s) \quad \text{if } \rho \geq r.$$

Next, we introduce a hybrid-index Besov space which plays an important role in the study of compressible fluids and was initiated in [11, 12].

DEFINITION 1.4. *Let $s, \sigma \in \mathbb{R}$, and set*

$$\|f\|_{\tilde{B}_2^{s,\sigma}} \triangleq \sum_{k \leq 0} 2^{ks} \|\Delta_k f\|_2 + \sum_{k > 0} 2^{k\sigma} \|\Delta_k f\|_2.$$

Let $m = -[\frac{d}{2} + 1 - s]$; we define

$$\begin{aligned} \tilde{B}_2^{s,\sigma}(\mathbb{R}^d) &= \{f \in \mathcal{S}'(\mathbb{R}^d) : \|f\|_{\tilde{B}_2^{s,\sigma}} < +\infty\} \quad \text{if } m < 0, \\ \tilde{B}_2^{s,\sigma}(\mathbb{R}^d) &= \{f \in \mathcal{S}'(\mathbb{R}^d)/\mathcal{P}_m : \|f\|_{\tilde{B}_2^{s,\sigma}} < +\infty\} \quad \text{if } m \geq 0, \end{aligned}$$

where \mathcal{P}_m denotes the set of polynomials of degree $\leq m$.

Throughout this paper, we will denote $\tilde{B}_{2,1}^s$ by B^s and $\tilde{B}_2^{s,\sigma}$ by $\tilde{B}^{s,\sigma}$. The following facts can be easily verified by using the definition of $\tilde{B}^{s,\sigma}$:

- (i) $\tilde{B}^{s,s} = \dot{B}_{2,1}^s$.
- (ii) If $s \leq \sigma$, then $\tilde{B}^{s,\sigma} = \dot{B}_{2,1}^s \cap \dot{B}_{2,1}^\sigma$. Otherwise, $\tilde{B}^{s,\sigma} = \dot{B}_{2,1}^s + \dot{B}_{2,1}^\sigma$.

Now we state our main result as follows.

THEOREM 1.2. *Let \bar{h}_0 be a positive constant. Assume that*

- (i) $(u_0, h_0 - \bar{h}_0) \in B^0(\mathbb{R}^2) \otimes \tilde{B}^{0,1}(\mathbb{R}^2)$;
- (ii) $h_0 \geq \bar{h}_0$.

Then there exist a positive time T and a unique solution (u, h) of (1.1) such that

$$(1.6) \quad \begin{aligned} u &\in C([0, T]; B^0) \cap L^1(0, T; B^2), \\ h - \bar{h}_0 &\in C([0, T]; \tilde{B}^{0,1}) \cap L^1(0, T; \tilde{B}^{2,1}), \quad h \geq \frac{1}{2}\bar{h}_0. \end{aligned}$$

Moreover, there exists a strictly positive constant c such that if

$$(1.7) \quad \|u_0\|_{B^0} + \|h_0 - \bar{h}_0\|_{\tilde{B}^{0,1}} \leq c,$$

then we can choose $T = +\infty$.

The structure of this paper is as follows. In section 2, we recall some useful multilinear estimates in the Besov spaces. In section 3, we prove the existence of the solution. In section 4, we prove the uniqueness of the solution. Finally, in the appendix, we prove some multilinear estimates in the weighted Besov spaces.

Throughout the paper, C denotes various ‘‘harmless’’ large finite constants, and c denotes various ‘‘harmless’’ small constants. We shall sometimes use $X \lesssim Y$ to denote the estimate $X \leq CY$ for some constant C . We denote $\|\cdot\|_p$ by the L^p norm of a function.

2. Multilinear estimates in the Besov spaces. Let us first recall Bony’s paraproduct decomposition.

DEFINITION 2.1. *We shall use the following Bony paraproduct decomposition (see [1, 3]):*

$$(2.1) \quad fg = T_f g + T_g f + R(f, g),$$

with

$$(2.2) \quad T_f g = \sum_{k \in \mathbb{Z}} S_{k-1} f \Delta_k g \quad \text{and} \quad R(f, g) = \sum_{k \in \mathbb{Z}} \sum_{|k'-k| \leq 1} \Delta_k f \Delta_{k'} g.$$

Next, let us recall some useful lemmas and multilinear estimates in the Besov spaces.

LEMMA 2.1 (Bernstein’s inequality). *Let $1 \leq p \leq q \leq +\infty$. Assume that $f \in \mathcal{S}'(\mathbb{R}^d)$; then for any $\gamma \in \mathbb{Z}^d$, there exist constants C_1, C_2 independent of f, j*

such that

$$\begin{aligned} \text{supp} \hat{f} \subseteq \{|\xi| \leq A_0 2^j\} &\Rightarrow \|\partial^\gamma f\|_q \leq C_1 2^{j|\gamma| + jd(\frac{1}{p} - \frac{1}{q})} \|f\|_p, \\ \text{supp} \hat{f} \subseteq \{A_1 2^j \leq |\xi| \leq A_2 2^j\} &\Rightarrow \|f\|_p \leq C_2 2^{-j|\gamma|} \sup_{|\beta|=|\gamma|} \|\partial^\beta f\|_p. \end{aligned}$$

The proof can be found in [8].

PROPOSITION 2.2. *If $s > 0$, $f, g \in B^s \cap L^\infty$, then $fg \in B^s \cap L^\infty$ and*

$$(2.3) \quad \|fg\|_{B^s} \leq C(\|f\|_\infty \|g\|_{B^s} + \|g\|_\infty \|f\|_{B^s}).$$

If $s_1, s_2 \leq \frac{d}{2}$ such that $s_1 + s_2 > 0$, $f \in B^{s_1}$, and $g \in B^{s_2}$, then $fg \in B^{s_1+s_2-\frac{d}{2}}$ and

$$(2.4) \quad \|fg\|_{B^{s_1+s_2-\frac{d}{2}}} \leq C\|f\|_{B^{s_1}} \|g\|_{B^{s_2}}.$$

If $|s| < \frac{d}{2}$, $1 \leq r \leq +\infty$, $f \in \dot{B}_{2,r}^s$, and $g \in B^{\frac{d}{2}}$, then $fg \in \dot{B}_{2,r}^s$ and

$$(2.5) \quad \|fg\|_{\dot{B}_{2,r}^s} \leq C\|f\|_{\dot{B}_{2,r}^s} \|g\|_{B^{\frac{d}{2}}}.$$

If $s \in (-\frac{d}{2}, \frac{d}{2}]$, $f \in B^s$, and $g \in \dot{B}_{2,\infty}^{-s}$, then $fg \in \dot{B}_{2,\infty}^{-\frac{d}{2}}$ and

$$(2.6) \quad \|fg\|_{\dot{B}_{2,\infty}^{-\frac{d}{2}}} \leq C\|f\|_{B^s} \|g\|_{\dot{B}_{2,\infty}^{-s}}.$$

If $1 \leq \rho_1, \rho_2, \rho \leq \infty$, $s \in (-\frac{d}{2}, \frac{d}{2}]$, $f \in \tilde{L}_T^{\rho_1}(B^s)$, and $g \in \tilde{L}_T^{\rho_2}(\dot{B}_{2,\infty}^{-s})$, then there holds

$$(2.7) \quad \|fg\|_{\tilde{L}_T^\rho(\dot{B}_{2,\infty}^{-\frac{d}{2}})} \leq C\|f\|_{\tilde{L}_T^{\rho_1}(B^s)} \|g\|_{\tilde{L}_T^{\rho_2}(\dot{B}_{2,\infty}^{-s})},$$

where $\frac{1}{\rho_1} + \frac{1}{\rho_2} = \frac{1}{\rho}$.

Proof. For the sake of simplicity, we present only the proof of (2.4); the others can be deduced in the same way (see also [14, 19]). By Bony's paraproduct decomposition and the property of quasi-orthogonality (1.5), for fixed $j \in \mathbb{Z}$, we write

$$\begin{aligned} \Delta_j(fg) &= \sum_{|k-j| \leq 3} \Delta_j(S_{k-1} f \Delta_k g) + \sum_{|k-j| \leq 3} \Delta_j(S_{k-1} g \Delta_k f) \\ &\quad + \sum_{k \geq j-2} \sum_{|k-k'| \leq 1} \Delta_j(\Delta_k f \Delta_{k'} g) \\ &\triangleq I + II + III. \end{aligned}$$

Thanks to the definition of Besov space B^s , we have

$$(2.8) \quad \|fg\|_{B^{s_1+s_2-\frac{d}{2}}} \leq \left(\sum_{j \in \mathbb{Z}} 2^{(s_1+s_2-\frac{d}{2})j} \|I\|_2 \right) + \dots + \left(\sum_{j \in \mathbb{Z}} 2^{(s_1+s_2-\frac{d}{2})j} \|III\|_2 \right) \triangleq I' + II' + III'.$$

It suffices to estimate the above three terms separately. Using Young's inequality and Lemma 2.1, we have

$$\begin{aligned} \|\Delta_j(S_{k-1}f\Delta_k g)\|_2 &\lesssim \|S_{k-1}f\|_\infty \|\Delta_k g\|_2 \lesssim \sum_{k' \leq k-2} \|\Delta_{k'} f\|_\infty \|\Delta_k g\|_2 \\ &\lesssim \sum_{k' \leq k-2} 2^{k's_1} \|\Delta_{k'} f\|_2 2^{k'(\frac{d}{2}-s_1)} \|\Delta_k g\|_2 \\ &\lesssim \|f\|_{B^{s_1}} \|\Delta_k g\|_2 2^{k(\frac{d}{2}-s_1)}, \end{aligned}$$

where we have used the fact $s_1 \leq \frac{d}{2}$ in the last inequality. Hence, we get

$$\begin{aligned} I' &\lesssim \|f\|_{B^{s_1}} \sum_{j \in \mathbb{Z}} 2^{(s_1+s_2-\frac{d}{2})j} \sum_{|k-j| \leq 3} 2^{k(\frac{d}{2}-s_1)} \|\Delta_k g\|_2 \\ (2.9) \quad &\lesssim \|f\|_{B^{s_1}} \sum_{|\ell| \leq 3} 2^{-(s_1+s_2-\frac{d}{2})\ell} \sum_{j \in \mathbb{Z}} 2^{s_2(j+\ell)} \|\Delta_{j+\ell} g\|_2 \lesssim \|f\|_{B^{s_1}} \|g\|_{B^{s_2}}. \end{aligned}$$

Similarly, using the fact $s_2 \leq \frac{d}{2}$, we can obtain

$$(2.10) \quad II' \lesssim \|f\|_{B^{s_1}} \|g\|_{B^{s_2}}.$$

Now we turn to estimate III' . From Lemma 2.1 and the Hölder inequality, it follows that

$$\|\Delta_j(\Delta_k f \Delta_{k'} g)\|_2 \lesssim 2^{j\frac{d}{2}} \|\Delta_k f \Delta_{k'} g\|_1 \lesssim 2^{j\frac{d}{2}} \|\Delta_k f\|_2 \|\Delta_{k'} g\|_2.$$

So, we get by the Minkowski inequality that for $s_1 + s_2 > 0$

$$\begin{aligned} III' &\lesssim \sum_{j \in \mathbb{Z}} 2^{(s_1+s_2-\frac{d}{2})j} 2^{j\frac{d}{2}} \left(\sum_{k \geq j-2} \sum_{|k-k'| \leq 1} \|\Delta_k f\|_2 \|\Delta_{k'} g\|_2 \right) \\ (2.11) \quad &\lesssim \sum_{\ell \geq -2} 2^{-(s_1+s_2)\ell} \sum_{j \in \mathbb{Z}} 2^{s_1(j+\ell)} \|\Delta_{j+\ell} f\|_2 \|g\|_{B^{s_2}} \lesssim \|f\|_{B^{s_1}} \|g\|_{B^{s_2}}. \end{aligned}$$

Summing up (2.8)–(2.11), we get the desired inequality (2.4). \square

PROPOSITION 2.3. (1) *Let $s > 0$. Assume that $F \in W_{loc}^{[s]+2, \infty}(\mathbb{R}^d)$ such that $F(0) = 0$. Then there exists a constant $C(s, d, F)$ such that if $u \in B^s \cap L^\infty$, it holds that*

$$(2.12) \quad \|F(u)\|_{B^s} \leq C(1 + \|u\|_\infty)^{[s]+1} \|u\|_{B^s};$$

and if $u \in \dot{B}_{2, \infty}^s \cap L^\infty$, it holds that

$$(2.13) \quad \|F(u)\|_{\dot{B}_{2, \infty}^s} \leq C(1 + \|u\|_\infty)^{[s]+1} \|u\|_{\dot{B}_{2, \infty}^s}.$$

(2) *Assume that $G \in W_{loc}^{[\frac{d}{2}]+3, \infty}(\mathbb{R}^d)$ such that $G'(0) = 0$. Then there exists a functions $C(s, d, G)$ such that if $-\frac{d}{2} < s \leq \frac{d}{2}$, $u, v \in B^{\frac{d}{2}} \cap L^\infty$ and $u - v \in B^s$, it holds that*

$$(2.14) \quad \|G(u) - G(v)\|_{B^s} \leq C(\|u\|_\infty, \|v\|_\infty) (\|u\|_{B^{\frac{d}{2}}} + \|v\|_{B^{\frac{d}{2}}}) \|u - v\|_{B^s};$$

and if $|s| < \frac{d}{2}$, $u, v \in B^{\frac{d}{2}} \cap L^\infty$ and $u - v \in \dot{B}_{2,\infty}^s$, it holds that

$$(2.15) \quad \|G(u) - G(v)\|_{\dot{B}_{2,\infty}^s} \leq C(\|u\|_\infty, \|v\|_\infty)(\|u\|_{B^{\frac{d}{2}}} + \|v\|_{B^{\frac{d}{2}}})\|u - v\|_{\dot{B}_{2,\infty}^s}.$$

Proof. We can refer to [2, 19] for the proof of (1). For (2), we refer to [11, 15]. For example, we write

$$G(u) - G(v) = (u - v) \int_0^1 G'(v + \tau(u - v)) d\tau.$$

Then it follows from (2.5) that for $|s| < \frac{d}{2}$

$$\|G(u) - G(v)\|_{\dot{B}_{2,\infty}^s} \leq C\|u - v\|_{\dot{B}_{2,\infty}^s} \|G'(v + \tau(u - v))\|_{B^{\frac{d}{2}}},$$

which together with (2.12) implies (2.15). \square

PROPOSITION 2.4. *Let A be a homogeneous smooth function of degree m . Assume that $-\frac{d}{2} < s_1, t_1, s_2, t_2 \leq 1 + \frac{d}{2}$. Then it holds that if $k \geq 1$,*

$$(2.16) \quad \begin{aligned} & |(A(D)\Delta_k(v \cdot \nabla f), A(D)\Delta_k f)| \\ & \lesssim \alpha_k 2^{-k(s_2 - m)} \|v\|_{B^{\frac{d}{2}+1}} \|f\|_{\tilde{B}^{s_1, s_2}} \|A(D)\Delta_k f\|_2; \end{aligned}$$

if $k \leq 0$,

$$(2.17) \quad \begin{aligned} & |(A(D)\Delta_k(v \cdot \nabla f), A(D)\Delta_k f)| \\ & \lesssim \alpha_k 2^{-k(s_1 - m)} \|v\|_{B^{\frac{d}{2}+1}} \|f\|_{\tilde{B}^{s_1, s_2}} \|A(D)\Delta_k f\|_2; \end{aligned}$$

if $k \geq 1$,

$$(2.18) \quad \begin{aligned} & |(A(D)\Delta_k(v \cdot \nabla f), \Delta_k g) + (\Delta_k(v \cdot \nabla g), A(D)\Delta_k f)| \\ & \lesssim \alpha_k \|v\|_{B^{\frac{d}{2}+1}} (2^{-kt_2} \|g\|_{\tilde{B}^{t_1, t_2}} \|A(D)\Delta_k f\|_2 + 2^{-k(s_2 - m)} \|f\|_{\tilde{B}^{s_1, s_2}} \|\Delta_k g\|_2); \end{aligned}$$

and if $k \leq 0$,

$$(2.19) \quad \begin{aligned} & |(A(D)\Delta_k(v \cdot \nabla f), \Delta_k g) + (\Delta_k(v \cdot \nabla g), A(D)\Delta_k f)| \\ & \lesssim \alpha_k \|v\|_{B^{\frac{d}{2}+1}} (2^{-kt_1} \|g\|_{\tilde{B}^{t_1, t_2}} \|A(D)\Delta_k f\|_2 + 2^{-k(s_1 - m)} \|f\|_{\tilde{B}^{s_1, s_2}} \|\Delta_k g\|_2), \end{aligned}$$

where $\sum_{k \in \mathbb{Z}} \alpha_k \leq 1$.

For the proof we refer to [12].

3. Existence. In this section, we prove the existence of the solution for the 2D viscous shallow water equations. Without loss of generality, we assume that $\bar{h}_0 = 1$ and $\nu = 1$. Replacing h by $h + 1$ in (1.1), we rewrite (1.1) as

$$(3.1) \quad \begin{cases} h_t + \operatorname{div} u + \operatorname{div}(hu) = 0, \\ u_t - (\nabla \cdot D(u) + \nabla \operatorname{div} u) + u \cdot \nabla u - \frac{\nabla h}{1+h}(D(u) + \operatorname{div} u) + \nabla h = 0, \\ u(0, \cdot) = u_0, \quad h(0, \cdot) = h_0. \end{cases}$$

3.1. The linearized system. In this subsection, we consider the linearized system of (3.1):

$$(3.2) \quad \begin{cases} h_t + v \cdot \nabla h + \operatorname{div} u = \mathcal{H}, \\ u_t - (\nabla \cdot D(u) + \nabla \operatorname{div} u) + v \cdot \nabla u + \nabla h = \mathcal{G}, \\ u(0, \cdot) = u_0, \quad h(0, \cdot) = h_0. \end{cases}$$

Let us first introduce some definitions. Set

$$e_k^r(t) \triangleq (1 - e^{-cr2^{2k}t})^{\frac{1}{r}}, \quad \omega_k(t) = \sum_{\tilde{k} \geq k} 2^{-(\tilde{k}-k)} (e_{\tilde{k}}^1(t) + e_{\tilde{k}}^2(t)),$$

where c is a positive constant which will be determined later. We remark that

$$\omega_k(t) \leq C \quad \text{for any } k \in \mathbb{Z},$$

which will be constantly used in the following.

DEFINITION 3.1. Let $s \in \mathbb{R}$ and $T > 0$. The function space E_T^s is defined by

$$E_T^s = \{f \in \mathcal{Z}'((0, T) \times \mathbb{R}^d) : \|f\|_{E_T^s} < +\infty\},$$

where

$$\|f\|_{E_T^s} \triangleq \sum_{k \in \mathbb{Z}} 2^{ks} \omega_k(T) \|\Delta_k f\|_{L_T^\infty(L^2)}.$$

DEFINITION 3.2. Let $s_1, s_2 \in \mathbb{R}$ and $T > 0$. The function space $\tilde{E}_T^{s_1, s_2}$ is defined by

$$\tilde{E}_T^{s_1, s_2} = \{f \in \mathcal{Z}'((0, T) \times \mathbb{R}^d) : \|f\|_{\tilde{E}_T^{s_1, s_2}} < +\infty\},$$

where

$$\|f\|_{\tilde{E}_T^{s_1, s_2}} \triangleq \sum_{k \leq 0} 2^{ks_1} \omega_k(T) \|\Delta_k f\|_{L_T^\infty(L^2)} + \sum_{k \geq 1} 2^{ks_2} \omega_k(T) \|\Delta_k f\|_{L_T^\infty(L^2)}.$$

Remark 3.1. If $s_1 \leq s_2$, then $\tilde{E}_T^{s_1, s_2} = E_T^{s_1} \cap E_T^{s_2}$. Otherwise, $\tilde{E}_T^{s_1, s_2} = E_T^{s_1} + E_T^{s_2}$.

Let (u, h) be a smooth solution of (3.2). We want to establish the following a priori estimates for (h, u) :

$$(3.3) \quad \begin{aligned} & \|u\|_{L_T^1(B^2)} + \|u\|_{L_T^2(B^1)} + \|h\|_{\tilde{E}_T^{0,1}} \\ & \leq C \sum_{k \in \mathbb{Z}} \omega_k(T) E_k(0) + C \sum_{k \in \mathbb{Z}} \omega_k(T) \|\Delta_k \mathcal{G}(t)\|_{L_T^1(L^2)} \\ & \quad + C \sum_{k \geq 1} \omega_k(T) \|\nabla \Delta_k \mathcal{H}(t)\|_{L_T^1(L^2)} + C \sum_{k < 1} \omega_k(T) \|\Delta_k \mathcal{H}(t)\|_{L_T^1(L^2)} \\ & \quad + C \|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + C \|h\|_{\tilde{E}_T^{0,1}} \|v\|_{L_T^1(B^2)} \end{aligned}$$

and

$$(3.4) \quad \begin{aligned} & \|u\|_{L_T^\infty(B^0)} + \|h\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h\|_{L_T^1(\tilde{B}^{2,1})} \\ & \leq E_0 + C \left(\|\mathcal{H}\|_{L_T^1(\tilde{B}^{0,1})} + \|\mathcal{G}\|_{L_T^1(B^0)} + \int_0^T V'(t) (\|u(t)\|_{B^0} + \|h(t)\|_{\tilde{B}^{0,1}}) dt \right), \end{aligned}$$

where $V(t) = \|v(t')\|_{L^1_t(B^2)}$ and

$$E_0 = \sum_{k \in \mathbb{Z}} E_k(0), \quad E_k(t) = \begin{cases} E_{hk}(t), & k \geq 1, \\ E_{lk}(t), & k < 1, \end{cases}$$

with

$$E_{hk}^2(t) = \frac{1}{2} \|u_k(t)\|_2^2 + \|\nabla h_k(t)\|_2^2 + (u_k(t), \nabla h_k(t)),$$

$$E_{lk}^2(t) = \frac{1}{2} \|u_k(t)\|_2^2 + \frac{1}{2} \|h_k(t)\|_2^2 + \frac{1}{8} (u_k(t), \nabla h_k(t)).$$

Let us begin with the proof of (3.3) and (3.4). Set

$$u_k = \Delta_k u, \quad h_k = \Delta_k h, \quad \mathcal{H}^k = \Delta_k \mathcal{H}, \quad \mathcal{G}^k = \Delta_k \mathcal{G}.$$

Then we get by applying the operator Δ_k to (3.2) that

$$(3.5) \quad \begin{cases} \partial_t h_k + \Delta_k(v \cdot \nabla h) + \operatorname{div} u_k = \mathcal{H}_k, \\ \partial_t u_k - (\nabla \cdot D(u_k) + \nabla \operatorname{div} u_k) + \Delta_k(v \cdot \nabla u) + \nabla h_k = \mathcal{G}_k, \\ u_k(0, \cdot) = \Delta_k u_0, \quad h_k(0, \cdot) = \Delta_k h_0. \end{cases}$$

Multiplying the second equation of (3.5) by u_k , and integrating the resulting equation over \mathbb{R}^2 , we obtain

$$(3.6) \quad \frac{1}{2} \frac{d}{dt} \|u_k\|_2^2 + \frac{1}{2} \|\nabla u_k\|_2^2 + \frac{3}{2} \|\operatorname{div} u_k\|_2^2 + (\nabla h_k, u_k) = (\mathcal{G}_k, u_k) - (\Delta_k(v \cdot \nabla u), u_k).$$

In the following, we will deal with the high frequency and the low frequency of h in a different manner.

High frequencies: $k \geq 1$. First, applying ∇ to the first equation of (3.5) and multiplying it by ∇h_k , then integrating the resulting equation over \mathbb{R}^2 , we obtain

$$(3.7) \quad \frac{1}{2} \frac{d}{dt} \|\nabla h_k\|_2^2 + (\nabla \operatorname{div} u_k, \nabla h_k) = (\nabla \mathcal{H}_k, \nabla h_k) - (\nabla \Delta_k(v \cdot \nabla h), \nabla h_k).$$

Second, applying the operator ∇ to the first equation of (3.5) and taking the L^2 product of the resulting equation with u_k , then taking the L^2 product of the second equation of (3.5) with ∇h_k , we get by summing them up that

$$(3.8) \quad \begin{aligned} & \frac{d}{dt} (u_k, \nabla h_k) - \|\operatorname{div} u_k\|_2^2 - 2(\nabla \operatorname{div} u_k, \nabla h_k) + \|\nabla h_k\|_2^2 \\ & = (\nabla \mathcal{H}_k, u_k) + (\mathcal{G}_k, \nabla h_k) - (\nabla \Delta_k(v \cdot \nabla h), u_k) - (\Delta_k(v \cdot \nabla u), \nabla h_k), \end{aligned}$$

where we used the fact that

$$(\nabla \cdot D(u_k) + \nabla \operatorname{div} u_k, \nabla h_k) = 2(\nabla \operatorname{div} u_k, \nabla h_k).$$

Then we get by summing up (3.6), (3.7)×2, and (3.8) that

$$\begin{aligned}
& \frac{d}{dt} \left[\frac{1}{2} \|u_k\|_2^2 + \|\nabla h_k\|_2^2 + (u_k, \nabla h_k) \right] \\
& \quad + \left[\|\nabla h_k\|_2^2 + \frac{1}{2} \|\nabla u_k\|_2^2 + \frac{1}{2} \|\operatorname{div} u_k\|_2^2 + (\nabla h_k, u_k) \right] \\
& = \left[(\nabla \mathcal{H}_k, u_k) + 2(\nabla \mathcal{H}_k, \nabla h_k) + (\mathcal{G}_k, u_k) + (\mathcal{G}_k, \nabla h_k) \right] \\
& \quad - (\Delta_k(v \cdot \nabla u), u_k) - 2(\nabla \Delta_k(v \cdot \nabla h), \nabla h_k) \\
& \quad - \left[(\nabla \Delta_k(v \cdot \nabla h), u_k) + (\Delta_k(v \cdot \nabla u), \nabla h_k) \right] \\
(3.9) \quad & \triangleq I + II + III + IV.
\end{aligned}$$

Note that

$$(u_k, \nabla h_k) \leq \frac{1}{3} \|u_k\|_2^2 + \frac{3}{4} \|\nabla h_k\|_2^2;$$

hence, we get by the definition of E_{hk} that

$$(3.10) \quad \frac{1}{6} (\|u_k\|_2^2 + \|\nabla h_k\|_2^2) \leq E_{hk}^2 \leq 2(\|u_k\|_2^2 + \|\nabla h_k\|_2^2).$$

Similarly, using the fact that $\frac{5}{6}2^k \geq \frac{5}{3}$ and (3.10), we have

$$(3.11) \quad \|\nabla h_k\|_2^2 + \frac{1}{2} \|\nabla u_k\|_2^2 + \frac{1}{2} \|\operatorname{div} u_k\|_2^2 + (\nabla h_k, u_k) \geq \frac{1}{8} E_{hk}^2.$$

By summing up (3.9)–(3.11), we obtain

$$(3.12) \quad \frac{d}{dt} E_{hk}^2 + cE_{hk}^2 \leq C|I + II + III + IV|.$$

In order to obtain (3.3), we use Lemma 5.1 to deal with the right-hand terms of (3.12). First, we get by using the Cauchy–Schwarz inequality and (3.10) that

$$(3.13) \quad |I| \leq C(\|\nabla \mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2) E_{hk}.$$

From Lemma 5.1 and (3.10), it follows that

$$(3.14) \quad |II + III + IV| \leq C(\|\mathcal{F}_k^1(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2) E_{hk}.$$

By summing up (3.12), and (3.13)–(3.14), we obtain

$$(3.15) \quad \frac{d}{dt} E_{hk} + cE_{hk} \leq C \left(\|\nabla \mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \|\mathcal{F}_k^1(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2 \right),$$

which implies that

$$\begin{aligned}
(3.16) \quad & \|E_{hk}(t)\|_{L_T^\infty} \leq E_{hk}(0) + C \left(\|\nabla \mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \right. \\
& \quad \left. + \|\mathcal{F}_k^1(t)\|_{L_T^1(L^2)} + \|\tilde{\mathcal{F}}_k^0(t)\|_{L_T^1(L^2)} \right).
\end{aligned}$$

Furthermore, by (5.4) and (5.5), it holds that

$$\begin{aligned} & \sum_{k \in \mathbb{Z}} \omega_k(T) \left(\|\mathcal{F}_k^1(t)\|_{L_T^1(L^2)} + \|\tilde{\mathcal{F}}_k^0(t)\|_{L_T^1(L^2)} \right) \\ & \leq C \left(\|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + \|h\|_{E_T^1} \|v\|_{L_T^1(B^2)} \right). \end{aligned}$$

Multiplying $\omega_k(T)$ on both sides of (3.16), then summing up the resulting equation over $k \geq 1$, we obtain

$$\begin{aligned} (3.17) \quad & \sum_{k \geq 1} \omega_k(T) \|E_{hk}(t)\|_{L_T^\infty} \leq \sum_{k \geq 1} \omega_k(T) E_{hk}(0) \\ & + C \sum_{k \geq 1} \omega_k(T) \left(\|\nabla \mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \right) \\ & + C \left(\|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + \|h\|_{E_T^1} \|v\|_{L_T^1(B^2)} \right). \end{aligned}$$

Next, we use the decay effect of the parabolic operators to estimate $\|u\|_{L_T^2(B^1) \cap L_T^1(B^2)}$. It follows from (3.6) and Lemma 5.1 that

$$\frac{d}{dt} \|u_k\|_2 + c2^{2k} \|u_k\|_2 \leq C (\|\nabla h_k(t)\|_2 + \|\mathcal{G}_k(t)\|_{L^2} + \|\tilde{\mathcal{F}}_k^0(t)\|_{L^2}),$$

which implies that

$$\|u_k\|_2 \leq e^{-ct2^{2k}} \|u_k(0)\|_2 + Ce^{-ct2^{2k}} *_t \left(\|\nabla h_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2 \right),$$

where the sign $*$ denotes the convolution of functions defined in \mathbb{R}^+ ; more precisely,

$$e^{-ct2^{2k}} *_t f \triangleq \int_0^t e^{-c(t-\tau)2^{2k}} f(\tau) d\tau.$$

Taking the L^r norm for $r = 1, 2$ with respect to t , we get by using Young's inequality that

$$\|u_k\|_{L_T^r(L^2)} \leq C2^{-2k/r} e_k^r(T) \left(\|u_k(0)\|_2 + \|\nabla h_k\|_{L_T^1(L^2)} + \|\mathcal{G}_k\|_{L_T^1(L^2)} + \|\tilde{\mathcal{F}}_k^0\|_{L_T^1(L^2)} \right),$$

which together with (5.5) implies that

$$\begin{aligned} (3.18) \quad & \sum_{k \geq 1} \left(2^{2k} \|u_k\|_{L_T^1(L^2)} + 2^k \|u_k\|_{L_T^2(L^2)} \right) \leq C \sum_{k \geq 1} \omega_k(T) \|u_k(0)\|_2 \\ & + C \sum_{k \geq 1} \omega_k(T) \left(\|\nabla h_k\|_{L_T^1(L^2)} + \|\mathcal{G}_k\|_{L_T^1(L^2)} \right) + C \|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)}, \end{aligned}$$

where we used the fact that

$$e_k^1(T) + e_k^2(T) \leq \omega_k(T).$$

On the other hand, it follows from (3.15) that

$$\|E_{hk}\|_2 \leq e^{-ct} E_{hk}(0) + Ce^{-ct} *_t \left(\|\nabla \mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \|\mathcal{F}_k^1(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2 \right).$$

Taking the L^1 norm with respect to t , we get by using Young's inequality that

$$(3.19) \quad \begin{aligned} \|E_{hk}\|_{L_T^1} &\leq C(1 - e^{-cT})E_{hk}(0) + C(1 - e^{-cT})\left(\|\nabla\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)}\right) \\ &\quad + \|\mathcal{F}_k^1(t)\|_{L_T^1(L^2)} + \|\tilde{\mathcal{F}}_k^0(t)\|_{L_T^1(L^2)}. \end{aligned}$$

Note that for $k \geq 1$

$$1 - e^{-ct} \leq 1 - e^{-ct2^{2k}} \leq \omega_k(t),$$

which together with (3.19) and Lemma 5.1 gives

$$(3.20) \quad \begin{aligned} \sum_{k \geq 1} \|E_{hk}\|_{L_T^1} &\leq C \sum_{k \geq 1} \omega_k(T)E_{hk}(0) + C \sum_{k \geq 1} \omega_k(T)\left(\|\nabla\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)}\right) \\ &\quad + C\left(\|u\|_{L_T^2(B^1)}\|v\|_{L_T^2(B^1)} + \|h\|_{E_T^1}\|v\|_{L_T^1(B^2)}\right). \end{aligned}$$

Plugging (3.20) into (3.18), we obtain

$$(3.21) \quad \begin{aligned} &\sum_{k \geq 1} \left(2^{2k}\|u_k\|_{L_T^1(L^2)} + 2^k\|u_k\|_{L_T^2(L^2)}\right) \\ &\leq C \sum_{k \geq 1} \omega_k(T)E_{hk}(0) + C \sum_{k \geq 1} \omega_k(T)\left(\|\nabla\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)}\right) \\ &\quad + C\left(\|u\|_{L_T^2(B^1)}\|v\|_{L_T^2(B^1)} + \|h\|_{E_T^1}\|v\|_{L_T^1(B^2)}\right). \end{aligned}$$

On the other hand, in order to obtain (3.4), we use Proposition 2.4 to deal with the right-hand terms of (3.12). Applying (2.16) with $s_1 = s_2 = 0$ to *II*, (2.16) with $s_1 = 0, s_2 = 1$ to *III*, and (2.18) with $t_1 = t_2 = 0, s_1 = 0, s_2 = 1$ to *IV*, we obtain

$$(3.22) \quad |II + III + IV| \leq CE_{hk}\alpha_k V'(t)(\|u\|_{B^0} + \|h\|_{\bar{B}^{0,1}}),$$

with $\sum_{k \in \mathbb{Z}} \alpha_k \leq 1$ and $V(t) = \|v(t')\|_{L_t^1(B^2)}$. From (3.13) and (3.22), it follows that

$$\frac{d}{dt}E_{hk} + cE_{hk} \leq C\left(\|\nabla\mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \alpha_k V'(t)(\|u\|_{B^0} + \|h\|_{\bar{B}^{0,1}})\right),$$

from which a similar proof of (3.21) ensures that

$$(3.23) \quad \begin{aligned} &\sum_{k \geq 1} \left(\|E_{hk}\|_{L_T^1} + \|E_{hk}\|_{L_T^\infty}\right) \leq C \sum_{k \geq 1} E_{hk}(0) \\ &\quad + C\left(\|\mathcal{H}\|_{L_T^1(\bar{B}^{0,1})} + \|\mathcal{G}\|_{L_T^1(B^0)} + \int_0^T V'(t)(\|u(t)\|_{B^0} + \|h(t)\|_{\bar{B}^{0,1}})dt\right). \end{aligned}$$

Low frequencies: $k < 1$. Multiplying the first equation of (3.5) by h_k , we get by integrating the resulting equation over \mathbb{R}^2 that

$$(3.24) \quad \frac{1}{2} \frac{d}{dt} \|h_k\|_2^2 + (\operatorname{div} u_k, h_k) = (\mathcal{H}_k, h_k) - (\Delta_k(v \cdot \nabla h), h_k).$$

Summing up (3.6), (3.8) $\times \frac{1}{8}$, and (3.24), we obtain

$$\begin{aligned}
 & \frac{d}{dt} \left[\frac{1}{2} \|u_k\|_2^2 + \frac{1}{2} \|h_k\|_2^2 + \frac{1}{8} (u_k, \nabla h_k) \right] \\
 & + \left[\frac{1}{8} \|\nabla h_k\|_2^2 + \frac{1}{2} \|\nabla u_k\|_2^2 + \frac{11}{8} \|\operatorname{div} u_k\|_2^2 - \frac{1}{4} (\nabla \operatorname{div} u_k, \nabla h_k) \right] \\
 & = \left[\frac{1}{8} (\nabla \mathcal{H}_k, u_k) + (\mathcal{H}_k, h_k) + (\mathcal{G}_k, u_k) + \frac{1}{8} (\mathcal{G}_k, \nabla h_k) \right] \\
 & - (\Delta_k(v \cdot \nabla u), u_k) - (\Delta_k(v \cdot \nabla h), h_k) \\
 & - \frac{1}{8} \left[(\nabla \Delta_k(v \cdot \nabla h), u_k) + (\Delta_k(v \cdot \nabla u), \nabla h_k) \right] \\
 (3.25) \quad & \triangleq I + II + III + IV.
 \end{aligned}$$

Note that $2^k \leq 1$. We get by the Cauchy–Schwarz inequality that

$$\frac{1}{8} (u_k, \nabla h_k) \leq \frac{3}{10} \|u_k\|_2 \|h_k\|_2 \leq \frac{1}{4} \|u_k\|_2^2 + \frac{1}{4} \|h_k\|_2^2;$$

hence, we get by the definition of E_{lk} that

$$(3.26) \quad \frac{1}{4} (\|u_k\|_2^2 + \|h_k\|_2^2) \leq E_{lk}^2 \leq 2(\|u_k\|_2^2 + \|h_k\|_2^2).$$

Similarly, we can prove

$$\frac{1}{4} (\nabla \operatorname{div} u_k, \nabla h_k) \leq \frac{3}{5} \|\operatorname{div} u_k\|_2 \|\nabla h_k\|_2 \leq \frac{9}{10} \|\nabla u_k\|_2^2 + \frac{1}{10} \|\nabla h_k\|_2^2,$$

which together with (3.26) implies that

$$\begin{aligned}
 & \frac{1}{8} \|\nabla h_k\|_2^2 + \frac{1}{2} \|\nabla u_k\|_2^2 + \frac{11}{8} \|\operatorname{div} u_k\|_2^2 - \frac{1}{4} (\nabla \operatorname{div} u_k, \nabla h_k) \\
 (3.27) \quad & \geq \frac{1}{160} 2^{2k} (\|u_k\|_2^2 + \|h_k\|_2^2) \geq \frac{1}{320} 2^{2k} E_{lk}^2.
 \end{aligned}$$

By summing up (3.25)–(3.27), we obtain

$$(3.28) \quad \frac{d}{dt} E_{lk}^2 + c 2^{2k} E_{lk}^2 \leq C |I + II + III + IV|.$$

In order to obtain (3.3), we use Lemma 5.1 to estimate the right-hand terms of (3.28). Using the fact that $2^k \leq 1$, we get by the Cauchy–Schwarz inequality and (3.26) that

$$(3.29) \quad |I| \leq C (\|\mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2) E_{lk}.$$

Using Lemma 5.1 and (3.26), we have

$$(3.30) \quad |II + III + IV| \leq C (\|\mathcal{F}_k^1(t)\|_2 + \|\mathcal{F}_k^0(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2) E_{lk}.$$

By summing up (3.28)–(3.30), we obtain

$$\frac{d}{dt} E_{lk} + c2^{2k} E_{lk} \leq C \left(\|\mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \|\mathcal{F}_k^1(t)\|_2 + \|\mathcal{F}_k^0(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2 \right),$$

which implies that

$$\begin{aligned} E_{lk} &\leq e^{-c2^{2k}t} E_{lk}(0) \\ &\quad + C e^{-c2^{2k}t} *_t \left(\|\mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \|\mathcal{F}_k^1(t)\|_2 + \|\mathcal{F}_k^0(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2 \right). \end{aligned}$$

Taking the L^r norm with respect to t , we get by using Young's inequality that

$$\begin{aligned} \|E_{lk}\|_{L_T^r} &\leq C 2^{-2k/r} e_k^r(T) \left(E_{lk}(0) + \|\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \right. \\ &\quad \left. + \|\mathcal{F}_k^1(t)\|_{L_T^1(L^2)} + \|\mathcal{F}_k^0(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_{L_T^1(L^2)} \right), \end{aligned}$$

from which, together with Lemma 5.1, it follows that

$$\begin{aligned} (3.31) \quad &\sum_{k < 1} \omega_k(T) \|E_{lk}\|_{L_T^\infty} \\ &\leq C \sum_{k < 1} \omega_k(T) E_{lk}(0) + \sum_{k < 1} \omega_k(T) \left(\|\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \right) \\ &\quad + C \left(\|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + \|h\|_{\tilde{E}_T^{0,1}} \|v\|_{L_T^1(B^2)} \right) \end{aligned}$$

and

$$\begin{aligned} &\sum_{k < 1} (2^{2k} \|E_{lk}\|_{L_T^1} + 2^k \|E_{lk}\|_{L_T^2}) \\ &\leq C \sum_{k < 1} \omega_k(T) E_{lk}(0) + \sum_{k < 1} \omega_k(T) \left(\|\mathcal{H}_k(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \right) \\ (3.32) \quad &\quad + C \left(\|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + \|h\|_{\tilde{E}_T^{0,1}} \|v\|_{L_T^1(B^2)} \right). \end{aligned}$$

On the other hand, in order to obtain (3.4), we use Proposition 2.4 to deal with the right-hand terms of (3.28). Applying (2.17) with $s_1 = s_2 = 0$ to II , (2.17) with $s_1 = 0, s_2 = 1$ to III , and (2.19) with $t_1 = t_2 = 0, s_1 = 0, s_2 = 1$ to IV , we obtain

$$(3.33) \quad |II + III + IV| \leq C E_{lk} \alpha_k V'(t) (\|u\|_{B^0} + \|h\|_{\tilde{B}^{0,1}}),$$

with $\sum_{k \in \mathbb{Z}} \alpha_k \leq 1$ and $V(t) = \|v(t')\|_{L_t^1(B^2)}$. From (3.32) and (3.36), it follows that

$$\frac{d}{dt} E_{lk} + c2^{2k} E_{lk} \leq C \left(\|\mathcal{H}_k(t)\|_2 + \|\mathcal{G}_k(t)\|_2 + \alpha_k V'(t) (\|u\|_{B^0} + \|h\|_{\tilde{B}^{0,1}}) \right),$$

which, together with a similar proof of (3.21), ensures that

$$\begin{aligned} &\sum_{k < 1} (2^{2k} \|E_{lk}\|_{L_T^1} + \|E_{lk}\|_{L_T^\infty}) \leq \sum_{k < 1} E_{hk}(0) \\ (3.34) \quad &\quad + C \left(\|\mathcal{H}\|_{L_T^1(\tilde{B}^{0,1})} + \|\mathcal{G}\|_{L_T^1(B^0)} + \int_0^T V'(t) (\|u(t)\|_{B^0} + \|h(t)\|_{\tilde{B}^{0,1}}) dt \right). \end{aligned}$$

The completion of the a priori estimates. First, adding up (3.17), (3.21), (3.31), and (3.32) yields that

$$\begin{aligned}
 & \|u\|_{L_T^1(B^2)} + \|u\|_{L_T^2(B^1)} + \|h\|_{\tilde{E}_T^{0,1}} \\
 & \leq C \sum_{k \in \mathbb{Z}} \omega_k(T) E_k(0) + C \sum_{k \in \mathbb{Z}} \omega_k(T) \|\mathcal{G}_k(t)\|_{L_T^1(L^2)} \\
 & \quad + C \sum_{k \geq 1} \omega_k(T) \|\nabla \mathcal{H}_k(t)\|_{L_T^1(L^2)} + C \sum_{k < 1} \omega_k(T) \|\mathcal{H}_k(t)\|_{L_T^1(L^2)} \\
 (3.35) \quad & + C \|u\|_{L_T^2(B^1)} \|v\|_{L_T^2(B^1)} + C \|h\|_{\tilde{E}_T^{0,1}} \|v\|_{L_T^1(B^2)},
 \end{aligned}$$

where we used the fact that

$$\|h\|_{E_T^1} \leq C \|h\|_{\tilde{E}_T^{0,1}}.$$

On the other hand, adding up (3.23) and (3.34) gives rise to

$$\begin{aligned}
 & \|u\|_{L_T^\infty(B^0)} + \|h\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h\|_{L_T^1(\tilde{B}^{2,1})} \\
 (3.36) \quad & \leq E_0 + C \left(\|\mathcal{H}\|_{L_T^1(\tilde{B}^{0,1})} + \|\mathcal{G}\|_{L_T^1(B^0)} + \int_0^T V'(t) (\|u\|_{B^0} + \|h\|_{\tilde{B}^{0,1}}) dt \right),
 \end{aligned}$$

which, together with the Gronwall inequality, implies that

$$\begin{aligned}
 & \|u\|_{L_T^\infty(B^0)} + \|h\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h\|_{L_T^1(\tilde{B}^{2,1})} \\
 (3.37) \quad & \leq C e^{C \|v\|_{L_T^1(B^2)}} \left(E_0 + \|\mathcal{H}\|_{L_T^1(\tilde{B}^{0,1})} + \|\mathcal{G}\|_{L_T^1(B^0)} \right).
 \end{aligned}$$

Finally, let us remark that

$$E_0 \approx (\|h_0\|_{\tilde{B}^{0,1}} + \|u_0\|_{B^0}).$$

3.2. The uniform estimate of the approximate sequence of solutions. In this subsection, we will construct the approximate solutions of (3.1) and present the uniform estimate of the approximate solutions. Let us first define the approximate sequence $(h^n, u^n)_{n \in \mathbb{N}}$ of (3.1) by the following system:

$$(3.38) \quad \begin{cases} \partial_t h^{n+1} + u^n \cdot \nabla h^{n+1} + \operatorname{div} u^{n+1} = \mathcal{H}^n, \\ \partial_t u^{n+1} - (\nabla \cdot D(u^{n+1}) + \nabla \operatorname{div} u^{n+1}) + u^n \cdot \nabla u^{n+1} + \nabla h^{n+1} = \mathcal{G}^n, \\ (h^{n+1}, u^{n+1})|_{t=0} = \sum_{|k| \leq n+N} \Delta_k(h_0, u_0), \end{cases}$$

where

$$\mathcal{H}^n \triangleq -h^n \operatorname{div} u^n, \quad \mathcal{G}^n \triangleq \frac{\nabla h^n}{1+h^n} \tilde{\nabla} u^n, \quad \text{with } \tilde{\nabla} u^n = D(u^n) + \operatorname{div} u^n,$$

and N is a fixed large integer such that

$$1 + h^n(0) \geq \frac{3}{4} \quad \text{for } n \geq 1.$$

Set $(h^0, u^0) = (0, 0)$ and solve the linear system. We can define $(h^n, u^n)_{n \in \mathbb{N}_0}$ by the induction. Next, we are going to prove by the induction that there exist positive constants η , K , and T such that the following bounds hold for all $n \in \mathbb{N}_0$:

$$(3.39) \quad 1 + h^n \geq \frac{1}{2},$$

$$(3.40) \quad \|u^n\|_{L_T^1(B^2) \cap L_T^2(B^1)} + \|h^n\|_{\tilde{E}_T^{0,1}} \leq \eta,$$

$$(3.41) \quad \|u^n\|_{L_T^\infty(B^0)} + \|h^n\|_{L_T^\infty(\tilde{B}^{0,1}) \cap L_T^1(\tilde{B}^{2,1})} \leq KE_0.$$

Assume that (3.39)–(3.41) hold for (h^n, u^n) . We need to prove that (3.39)–(3.41) also hold for (h^{n+1}, u^{n+1}) . Applying the a priori estimates (3.35) and (3.37) to (h^{n+1}, u^{n+1}) , we obtain

$$\begin{aligned} & \|u^{n+1}\|_{L_T^1(B^2)} + \|u^{n+1}\|_{L_T^2(B^1)} + \|h^{n+1}\|_{\tilde{E}_T^{0,1}} \\ & \leq C\mathcal{Q}_0(T) + C \sum_{k \in \mathbb{Z}} \omega_k(T) \|\mathcal{G}_k^n(t)\|_{L_T^1(L^2)} + C \sum_{k \geq 1} \omega_k(T) \|\nabla \mathcal{H}_k^n(t)\|_{L_T^1(L^2)} \\ & \quad + C \sum_{k < 1} \omega_k(T) \|\mathcal{H}_k^n(t)\|_{L_T^1(L^2)} + C \|u^{n+1}\|_{L_T^2(B^1)} \|u^n\|_{L_T^2(B^1)} \\ (3.42) \quad & + C \|h^{n+1}\|_{\tilde{E}_T^{0,1}} \|u^n\|_{L_T^1(B^2)} \end{aligned}$$

and

$$\begin{aligned} & \|u^{n+1}\|_{L_T^\infty(B^0)} + \|h^{n+1}\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h^{n+1}\|_{L_T^1(\tilde{B}^{2,1})} \\ (3.43) \quad & \leq C e^{C \|u^n\|_{L_T^1(B^2)}} \left(E_0 + \|\mathcal{H}^n\|_{L_T^1(\tilde{B}^{0,1})} + \|\mathcal{G}^n\|_{L_T^1(B^0)} \right), \end{aligned}$$

with

$$\mathcal{Q}_0(T) \triangleq \sum_{k \in \mathbb{Z}} \omega_k(T) E_k(0).$$

Thanks to (2.4), we have

$$\|\mathcal{H}^n\|_{B^0} \leq C \|h^n\|_{B^0} \|u^n\|_{B^2} \quad \text{and} \quad \|\mathcal{H}^n\|_{B^1} \leq C \|h^n\|_{B^1} \|u^n\|_{B^2},$$

which, together with the fact that $\tilde{B}^{0,1} = B^0 \cap B^1$, yields

$$(3.44) \quad \|\mathcal{H}^n\|_{L_T^1(\tilde{B}^{0,1})} \leq C \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})} \|u^n\|_{L_T^1(B^2)} \leq CKE_0\eta.$$

We rewrite \mathcal{G}^n as

$$\frac{\nabla h^n}{1+h^n} \tilde{\nabla} u^n = (1+h^n) \nabla \left(\frac{h^n}{1+h^n} \right) \tilde{\nabla} u^n.$$

Using (2.4) and (2.12), we get

$$\begin{aligned} \|\mathcal{G}^n\|_{L_T^1(B^0)} & \leq C \left\| \nabla \left(\frac{h^n}{1+h^n} \right) \right\|_{L_T^\infty(B^0)} \|(1+h^n) \tilde{\nabla} u^n\|_{L_T^1(B^1)} \\ & \leq C(1 + \|h^n\|_{L_T^\infty(L^\infty)})^2 \|h^n\|_{L_T^\infty(B^1)} (1 + \|h^n\|_{L_T^\infty(B^1)}) \|u^n\|_{L_T^1(B^2)} \\ & \leq C(1 + \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})})^3 \|h^n\|_{L_T^\infty(B^1)} \|u^n\|_{L_T^1(B^2)} \\ (3.45) \quad & \leq CKE_0(1 + KE_0)^3 \eta. \end{aligned}$$

Plugging (3.44) and (3.45) into (3.43) yields

$$(3.46) \quad \begin{aligned} & \|u^{n+1}\|_{L_T^\infty(B^0)} + \|h^{n+1}\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h^{n+1}\|_{L_T^1(\tilde{B}^{2,1})} \\ & \leq C e^{C\eta} \left(E_0 + K E_0 (1 + K E_0)^3 \eta \right). \end{aligned}$$

We take $\eta > 0$ small enough and $K = 4C$ such that

$$(3.1) \quad e^{C\eta} \leq 2, \quad K(1 + K E_0)^3 \eta \leq 1,$$

from which, together with (3.46), it follows that

$$\|u^{n+1}\|_{L_T^\infty(B^0)} + \|h^{n+1}\|_{L_T^\infty(\tilde{B}^{0,1})} + \|h^{n+1}\|_{L_T^1(\tilde{B}^{2,1})} \leq K E_0.$$

This proves (3.41) for (u^{n+1}, h^{n+1}) .

Next, we prove (3.40) for (u^{n+1}, h^{n+1}) . Applying Lemma 5.2 with $s_1 = 0$ and $s_2 = 1$, (2.4) with $s_1 = s_2 = 1$, and Lemma 5.4 with $s = 1$, we obtain

$$(3.47) \quad \begin{aligned} \sum_{k \in \mathbb{Z}} \omega_k(T) \|\mathcal{G}_k^n(t)\|_{L_T^1(L^2)} & \leq C \left\| \nabla \left(\frac{h^n}{1+h^n} \right) \right\|_{E_T^0} \|(1+h^n)\tilde{\nabla} u^n\|_{L_T^1(B^1)} \\ & \leq C(1 + \|h^n\|_{L_T^\infty(L^\infty)})^3 \|h^n\|_{E_T^1} (1 + \|h^n\|_{L_T^\infty(B^1)}) \|u^n\|_{L_T^1(B^2)} \\ & \leq C(1 + \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})})^4 \|h^n\|_{\tilde{E}_T^{0,1}} \|u^n\|_{L_T^1(B^2)} \\ & \leq C(1 + K E_0)^4 \eta^2. \end{aligned}$$

On the other hand, we apply Lemma 5.2 with $s_1 = 0, s_2 = 1$ to get

$$\begin{aligned} & \sum_{k \geq 1} \omega_k(T) \|\nabla \mathcal{H}_k^n(t)\|_{L_T^1(L^2)} + \sum_{k < 1} \omega_k(T) \|\mathcal{H}_k^n(t)\|_{L_T^1(L^2)} \\ & \leq C \sum_{k \in \mathbb{Z}} \omega_k(T) (\|\nabla h_k^n\|_{L_T^\infty(L^2)} + \|h_k^n\|_{L_T^\infty(L^2)}) \|\operatorname{div} u^n\|_{L_T^1(B^1)} \\ & \quad + C \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{2k} \|u_k^n\|_{L_T^1(L^2)} \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})} \\ & \triangleq I + II. \end{aligned}$$

Obviously, we have

$$(3.48) \quad I \leq C \|h^n\|_{\tilde{E}_T^{0,1}} \|u^n\|_{L_T^1(B^2)} \leq C \eta^2.$$

In order to estimate II , we first fix $k_0 \geq 1$ such that

$$(3.49) \quad \sum_{k \geq k_0} \|u_k(0)\|_2 \leq \frac{\eta}{16CKE_0}.$$

Then we write

$$\begin{aligned} II & = C \sum_{k \geq k_0} \omega_k(T) 2^{2k} \|u_k^n\|_{L_T^1(L^2)} \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})} + C \sum_{k \leq k_0} \omega_k(T) 2^{2k} \|u_k^n\|_{L_T^1(L^2)} \|h^n\|_{L_T^\infty(\tilde{B}^{0,1})} \\ & \triangleq II_1 + II_2. \end{aligned}$$

Using (3.18), (3.47), and (3.49), we obtain

$$\begin{aligned}
II_1 &\leq CK E_0 \left[\sum_{k \geq k_0} \omega_k(T) \|u_k(0)\|_2 + \sum_{k \geq k_0} \omega_k(T) \left(\|\nabla h_k^n\|_{L_T^1(L^2)} + \|\mathcal{G}_k^{n-1}\|_{L_T^1(L^2)} \right) \right. \\
&\quad \left. + \|u^n\|_{L_T^2(B^1)} \|u^{n-1}\|_{L_T^2(B^1)} \right] \\
(3.50) \quad &\leq CK E_0 \left[\frac{\eta}{16CK E_0} + \sum_{k \geq k_0} \omega_k(T) \|\nabla h_k^n\|_{L_T^1(L^2)} + (1 + KE_0)^4 \eta^2 \right].
\end{aligned}$$

On the other hand, thanks to (3.19) and Lemma 5.1, we have

$$\begin{aligned}
\sum_{k \geq k_0} \omega_k(T) \|\nabla h_k^n\|_{L_T^1(L^2)} &\leq C(1 - e^{-cT}) \sum_{k \geq k_0} \omega_k(T) E_{hk}(0) \\
&\quad + C(1 - e^{-cT}) \sum_{k \geq k_0} \omega_k(T) \\
&\quad \left(\|\nabla \mathcal{H}_k^{n-1}(t)\|_{L_T^1(L^2)} + \|\mathcal{G}_k^{n-1}(t)\|_{L_T^1(L^2)} \right) \\
&\quad + C \|u^n\|_{L_T^2(B^1)} \|u^{n-1}\|_{L_T^2(B^1)} + C \|h^n\|_{E_T^1} \|u^{n-1}\|_{L_T^2(B^1)} \\
&\leq C(1 - e^{-cT}) E_0 + C(1 - e^{-cT}) KE_0 \eta + C(1 + KE_0)^4 \eta^2,
\end{aligned}$$

where we used (3.44) and (3.47) in the second inequality. Plugging the above inequality into (3.50) yields

$$(3.51) \quad II_1 \leq CK E_0 \left[\frac{\eta}{16CK E_0} + (1 - e^{-cT})(E_0 + KE_0 \eta) + (1 + KE_0)^4 \eta^2 \right].$$

Note that for $k \leq k_0$, we can choose $T > 0$ small enough such that

$$(3.52) \quad \omega_k(T) \leq \frac{1}{16CK E_0},$$

so we get

$$(3.52) \quad |II_2| \leq \frac{\eta}{16}.$$

Plugging (3.47), (3.48), (3.51), and (3.52) into (3.42), we get

$$\begin{aligned}
&\|u^{n+1}\|_{L_T^1(B^2)} + \|u^{n+1}\|_{L_T^2(B^1)} + \|h^{n+1}\|_{\tilde{E}_T^{0,1}} \\
&\leq C \mathcal{Q}_0(T) + \frac{\eta}{8} + C(1 + KE_0)^5 \eta^2 + CK E_0 (1 - e^{-cT})(E_0 + KE_0 \eta) \\
(3.53) \quad &+ C \eta (\|u^{n+1}\|_{L_T^2(B^1)} + \|h^{n+1}\|_{\tilde{E}_T^{0,1}}).
\end{aligned}$$

Note that $\mathcal{Q}_0(0) = 0$. We can take T, η small enough such that

$$\begin{aligned}
&C \eta \leq \frac{1}{2}, \quad C \mathcal{Q}_0(T) \leq \frac{\eta}{8}, \quad C(1 + KE_0)^5 \eta < \frac{1}{8}, \quad \text{and} \\
(3.53) \quad &CK E_0 (1 - e^{-cT})(E_0 + KE_0 \eta) \leq \frac{\eta}{8},
\end{aligned}$$

which, together with (3.53), gives

$$\|u^{n+1}\|_{L_T^1(B^2)} + \|u^{n+1}\|_{L_T^2(B^1)} + \|h^{n+1}\|_{\tilde{E}_T^{0,1}} \leq \eta.$$

Finally, let us prove (3.39) for h^{n+1} . We rewrite the first equation of (3.38) as

$$\partial_t(1 + h^{n+1}) + u^n \cdot \nabla(1 + h^{n+1}) + \operatorname{div} u^{n+1} - \mathcal{H}^n = 0.$$

Then $1 + h^{n+1}$ can be represented as

$$(3.54) \quad \begin{aligned} (1 + h^{n+1})(t, x) &= (1 + h_0^{n+1})((\psi^n)_t^{-1}(x)) + \int_0^t \operatorname{div} u^{n+1}(\tau, \psi_\tau^n((\psi^n)_t^{-1}(x))) d\tau \\ &+ \int_0^t \mathcal{H}^n(\tau, \psi_\tau^n((\psi^n)_t^{-1}(x))) d\tau, \end{aligned}$$

where the flow map ψ_t^n is defined by

$$\begin{cases} \partial_t \psi_t^n(x) = u^n(t, \psi_t^n(x)), \\ \psi_t^n|_{t=0} = x. \end{cases}$$

Thanks to the inclusion map $B^1 \hookrightarrow L^\infty$ and (2.4), we get

$$\begin{aligned} \int_0^t \|\operatorname{div} u^{n+1}(\tau, \psi_\tau^n((\psi^n)_t^{-1}(x)))\|_\infty d\tau &\leq \|u^{n+1}\|_{L_t^1(B^2)} \leq \eta, \\ \int_0^t \|\mathcal{H}^n(\tau, \psi_\tau^n((\psi^n)_t^{-1}(x)))\|_\infty d\tau &\leq \|h^n \operatorname{div} u^n\|_{L_t^1(B^1)} \\ &\leq C \|h^n\|_{L_t^\infty(\tilde{B}^{0,1})} \|u^n\|_{L_t^1(B^2)} \leq CK E_0 \eta, \end{aligned}$$

from which, together with (3.54), it follows that

$$(3.55) \quad 1 + h^{n+1} \geq \frac{3}{4} - (1 + CK E_0) \eta.$$

We take η small enough such that

$$(3.56) \quad (1 + CK E_0) \eta \leq \frac{1}{4},$$

which, together with (3.55), ensures that

$$1 + h^{n+1} \geq \frac{1}{2}.$$

So far, we have shown that T, η can be chosen small enough such that the assumptions (\mathfrak{R}_1) – (\mathfrak{R}_4) hold, under which the approximate solutions $(u^n, h^n)_{n \in \mathbb{N}_0}$ are uniformly bounded in

$$\mathcal{E}_T \triangleq \left(L_T^\infty(B^0) \cap L_T^1(B^2) \right) \times \left(L_T^\infty(\tilde{B}^{0,1}) \cap L_T^1(\tilde{B}^{2,1}) \right).$$

It should be pointed out that if $\|u_0\|_{B^0} + \|h_0\|_{\tilde{B}^{0,1}}$ is small enough, we can take $T = +\infty$ such that the assumptions (\mathfrak{R}_1) – (\mathfrak{R}_4) hold.

3.3. The existence of the solution. Now let us turn to prove the existence of the solution, and the standard compact arguments will be used. We should point out that in order to obtain the convergence of the terms $u^n \cdot \nabla u^{n+1}$ and $u^n \cdot \nabla h^{n+1}$ in (3.38), we need to show that both (u^n, h^n) and (u^{n+1}, h^{n+1}) have the same limit. Indeed, following the argument in the proof of uniqueness (section 4), it can be proved that $(u^n, h^n)_n$ is a Cauchy sequence in some weak topology which particularly ensures the uniqueness of the limit for any subsequence of $(u^n, h^n)_n$. In section 3.2, we have showed that the approximate solutions $(h^n, u^n)_{n \in \mathbb{N}}$ satisfy (3.39)–(3.41), and without loss of generality, we can assume the following:

$$(3.56) \quad 1 + h^n \geq \frac{1}{2},$$

$$(3.57) \quad \|u^n\|_{L^\infty(B^0) \cap L^1_T(B^2)} + \|h^n\|_{L^\infty(\tilde{B}^{0,1}) \cap L^1_T(\tilde{B}^{2,1})} \leq KE_0.$$

Using the interpolation and the fact that $B^0 \cap B^1 = \tilde{B}^{0,1}$, we have

$$\begin{aligned} \|h^n\|_{L^2_T(B^1)} &\lesssim \|h^n\|_{L^\infty_T(\tilde{B}^{0,1})}^{\frac{1}{2}} \|h^n\|_{L^1_T(\tilde{B}^{2,1})}^{\frac{1}{2}}, & \|u^n\|_{L^2_T(B^1)} &\lesssim \|u^n\|_{L^\infty_T(B^0)}^{\frac{1}{2}} \|u^n\|_{L^1_T(B^2)}^{\frac{1}{2}}, \\ \|h^n\|_{L^4_T(B^{\frac{1}{2}})} &\lesssim \|h^n\|_{L^\infty_T(\tilde{B}^{0,1})}^{\frac{1}{2}} \|h^n\|_{L^2_T(B^1)}^{\frac{1}{2}}, & \|u^n\|_{L^{\frac{4}{3}}_T(B^{\frac{3}{2}})} &\lesssim \|u^n\|_{L^\infty_T(B^0)}^{\frac{1}{4}} \|u^n\|_{L^1_T(B^2)}^{\frac{3}{4}}, \end{aligned}$$

from which, together with (3.57), it follows that

$$(3.58) \quad \|h^n\|_{L^2_T(B^1)} + \|u^n\|_{L^2_T(B^1)} + \|h^n\|_{L^4_T(B^{\frac{1}{2}})} + \|u^n\|_{L^{\frac{4}{3}}_T(B^{\frac{3}{2}})} \lesssim KE_0.$$

Now, we show that (h^n, u^n) is uniformly bounded in $C^{\frac{1}{2}}_{loc}(B^0) \times C^{\frac{1}{4}}_{loc}(B^{-\frac{1}{2}})$. Using (2.4), (3.57), and (3.58), it is easy to verify that

$$\begin{aligned} \|u^n \cdot \nabla h^{n+1}\|_{L^2_T(B^0)} &\lesssim \|u^n\|_{L^2_T(B^1)} \|h^{n+1}\|_{L^\infty_T(\tilde{B}^{0,1})} \lesssim (KE_0)^2, \\ \|h^n \operatorname{div} u^n\|_{L^2_T(B^0)} &\lesssim \|u^n\|_{L^2_T(B^1)} \|h^n\|_{L^\infty_T(\tilde{B}^{0,1})} \lesssim (KE_0)^2, \end{aligned}$$

from which, together with the first equation of (3.38), it follows that $\partial_t h^n$ is uniformly bounded in $L^2_T(B^0)$, which implies that h^n is uniformly bounded in $C^{\frac{1}{2}}_{loc}(B^0)$. On the other hand, thanks to (2.4), (3.58), and (2.12), we have

$$\begin{aligned} \|u^n \cdot \nabla u^{n+1}\|_{L^{\frac{4}{3}}_T(B^{-\frac{1}{2}})} &\lesssim \|u^n\|_{L^\infty_T(B^0)} \|u^{n+1}\|_{L^{\frac{4}{3}}_T(B^{\frac{3}{2}})} \lesssim (KE_0)^2, \\ \left\| \frac{\nabla h^n}{1+h^n} \tilde{\nabla} u^n \right\|_{L^{\frac{4}{3}}_T(B^{-\frac{1}{2}})} &\lesssim C(1 + \|h^n\|_{L^\infty_T(B^1)})^3 \|u^n\|_{L^{\frac{4}{3}}_T(B^{\frac{3}{2}})} \lesssim C(1 + KE_0)^3 KE_0, \end{aligned}$$

from which, together with the second equation of (3.38), it follows that $\partial_t u^n$ is uniformly bounded in $L^{\frac{4}{3}}_T(B^{-\frac{1}{2}})$, which implies that u^n is uniformly bounded in $C^{\frac{1}{4}}_{loc}(B^{-\frac{1}{2}})$.

Next, we claim that the inclusions $B^0 \cap B^1 \hookrightarrow L^2$ and $B^{-\frac{1}{2}} \cap B^0 \hookrightarrow \dot{H}^{-\frac{1}{2}}$ are locally compact. Indeed, these can be proved by noting that for $s' < s$, $\dot{H}^{s'} \cap \dot{H}^s \hookrightarrow \dot{H}^{s'}$ is locally compact and for $s \in \mathbb{R}$, $B^s \hookrightarrow \dot{H}^s$. Then, by the Arzelà–Ascoli theorem and

Cantor’s diagonal process, there exist a subsequence (u^{n_k}, h^{n_k}) and a function (u, h) such that

$$(3.59) \quad (u^{n_k}, h^{n_k}) \rightarrow (u, h) \quad \text{in } C_{loc}(\dot{H}_{loc}^{-\frac{1}{2}}) \times C_{loc}(L^2_{loc}),$$

as $n_k \rightarrow \infty$. On the other hand, (u^{n_k}, h^{n_k}) is uniformly bounded in \mathcal{E}_T . Then there exists a subsequence (which is still denoted by (u^{n_k}, h^{n_k})) such that

$$(u^{n_k}, h^{n_k}) \rightharpoonup (u, h) \quad \text{in } \mathcal{E}_T,$$

where “ \rightharpoonup ” denotes weak* convergence.

Finally, let us prove that (u, h) solves (1.1) in the sense of distribution. We need only prove that the nonlinear terms such as $u^n \cdot \nabla h^n$, $\frac{\nabla h^n}{1+h^n} \tilde{\nabla} u^n$, etc., tend to the corresponding nonlinear terms in the sense of distribution. This can be done by using the uniform estimates of (u^n, h^n) , (u, h) in \mathcal{E}_T and the convergence result (3.59). Here, we show only the case of the term $Y(h^n) \tilde{\nabla} u^n$ (where $Y(z) \triangleq \nabla z / (1+z)$); the other terms can be treated in the same way. For any test function $\theta \in C_0^\infty([0, T^*) \times \mathbb{R}^2)$, we write

$$\begin{aligned} & \langle Y(h^n) \tilde{\nabla} u^n - Y(h) \tilde{\nabla} u, \theta \rangle \\ &= \left\langle (1+h^n) \nabla \left(\frac{h^n}{1+h^n} - \frac{h}{1+h} \right) \tilde{\nabla} u^n, \theta \right\rangle \\ & \quad + \left\langle (h^n - h) \nabla \left(\frac{h}{1+h} \right) \tilde{\nabla} u^n, \theta \right\rangle + \left\langle (1+h) \nabla \left(\frac{h}{1+h} \right) \tilde{\nabla} (u^n - u), \theta \right\rangle \\ & \triangleq I_1 + I_2 + I_3. \end{aligned}$$

Thanks to (2.4) and (3.56), we have

$$\begin{aligned} I_1 &\leq \left\| \frac{\psi(h^n - h)}{(1+h^n)(1+h)} \right\|_2 \|\nabla((1+h^n) \tilde{\nabla} u^n \theta)\|_2 \lesssim \|\psi(h^n - h)\|_2 \|(1+h^n) \tilde{\nabla} u^n\|_{B^1} \\ &\lesssim \|\psi(h^n - h)\|_2 (1 + \|h^n\|_{\tilde{B}^{0,1}}) \|u^n\|_{B^2}, \end{aligned}$$

where $\psi \in C_0^\infty([0, T^*) \times \mathbb{R}^2)$, and $\psi = 1$ on $\text{supp } \theta$. For I_2 , we have

$$\begin{aligned} I_2 &\leq \|\theta(h^n - h)\|_2 \left\| \nabla \left(\frac{h}{1+h} \right) \tilde{\nabla} u^n \right\|_2 \lesssim \|\theta(h^n - h)\|_2 \|\nabla h\|_2 \|\nabla u^n\|_{L^\infty} \\ &\lesssim \|\theta(h^n - h)\|_2 \|h\|_{\tilde{B}^{0,1}} \|u^n\|_{B^2}. \end{aligned}$$

Using (3.56) and the interpolation, we get

$$\begin{aligned} I_3 &\leq \left\| (1+h) \nabla \left(\frac{h}{1+h} \right) \right\|_2 \|\tilde{\nabla} (u^n - u) \theta\|_2 \lesssim (1 + \|h\|_\infty) \|\nabla h\|_2 \|(u^n - u) \theta\|_{\dot{H}^1} \\ &\lesssim (1 + \|h\|_{\tilde{B}^{0,1}}) \|h\|_{B^1} \|u^n - u\|_{H^2}^{\frac{3}{5}} \|(u^n - u) \theta\|_{\dot{H}^{-\frac{1}{2}}}^{\frac{2}{5}}. \end{aligned}$$

Thus, by (3.59), we get as $n \rightarrow \infty$

$$\langle Y(h^n) \tilde{\nabla} u^n - Y(h) \tilde{\nabla} u, \theta \rangle \rightarrow 0.$$

Following the argument in [11], we can also prove that (u, h) is continuous in time with values in $B^0 \times \tilde{B}^{0,1}$.

4. Uniqueness. In this section, we will prove the uniqueness of the solution. First, let us recall some known results.

LEMMA 4.1 (Osgood’s lemma). *Let ρ be a measurable positive function and γ a positive locally integrable function, each defined on the domain $[t_0, t_1]$. Let $\mu : [0, \infty) \rightarrow [0, \infty)$ be a continuous nondecreasing function, with $\mu(0) = 0$. Let $a \geq 0$, and assume that for all t in $[t_0, t_1]$,*

$$\rho(t) \leq a + \int_{t_0}^t \gamma(\tau)\mu(\rho(\tau))d\tau.$$

If $a > 0$, then

$$-\mathcal{M}(\rho(t)) + \mathcal{M}(a) \leq \int_{t_0}^t \gamma(\tau)d\tau, \quad \text{where } \mathcal{M}(x) = \int_x^1 \frac{d\tau}{\mu(\tau)}.$$

If $a = 0$ and $\mathcal{M} = \infty$, then $\rho \equiv 0$.

This lemma can be understood as a generalization of the classical Gronwall lemma and can be found in [8].

PROPOSITION 4.2. *Let $s \in (-\frac{d}{p}, 1 + \frac{d}{p})$ and $1 \leq p, r \leq +\infty$. Let v be a vector field such that $\nabla v \in L_T^1(\dot{B}_{p,r}^{\frac{d}{p}} \cap L^\infty)$. Assume that $f_0 \in \dot{B}_{p,r}^s$, $g \in L_T^1(\dot{B}_{p,r}^s)$ and $f \in L_T^\infty(\dot{B}_{p,r}^s) \cap C([0, T]; \mathcal{S}')$ is the solution of*

$$\begin{cases} \partial_t f + v \cdot \nabla f = g, \\ f(0, x) = f_0. \end{cases}$$

Then there exists a constant $C(s, p, d)$ such that for $t \in [0, T]$

$$(4.1) \quad \|f\|_{\tilde{L}_t^\infty(\dot{B}_{p,r}^s)} \leq C e^{CV(t)} \left(\|f_0\|_{\dot{B}_{p,r}^s} + \int_0^t e^{-CV(\tau)} \|g(\tau)\|_{\dot{B}_{p,r}^s} d\tau \right),$$

where $V(t) \triangleq \int_0^t \|\nabla v(\tau)\|_{\dot{B}_{p,r}^{\frac{d}{p}} \cap L^\infty} d\tau$. If $r < +\infty$, then f belongs to $C([0, T]; \dot{B}_{p,r}^s)$.

The proof can be found in [15].

PROPOSITION 4.3. *Let $T > 0$, $s \in \mathbb{R}$, and $1 \leq q, r \leq +\infty$. Assume that $u_0 \in \dot{B}_{2,q}^s$, $g \in \tilde{L}_T^1(\dot{B}_{2,q}^s)$, and u is the solution of*

$$\begin{cases} \partial_t u - \nu \tilde{\Delta} u = g, \\ u(0, x) = u_0, \end{cases}$$

where $\tilde{\Delta} u = \nabla \cdot D(u) + \nabla \operatorname{div} u$. Then there exists a constant $C(s, d, \nu)$ such that

$$(4.2) \quad \begin{aligned} (r\nu)^{\frac{1}{r}} \|u\|_{\tilde{L}_T^r(\dot{B}_{2,q}^{s+\frac{2}{r}})} &\leq \left(\sum_{k \in \mathbb{Z}} (1 - e^{-r\nu 2^{2k}T})^{\frac{q}{r}} 2^{qks} \|\Delta_k u_0\|_2^q \right)^{\frac{1}{q}} \\ &+ C \left(\sum_{k \in \mathbb{Z}} (1 - e^{-r\nu 2^{2k}T})^{\frac{q}{r}} 2^{qks} \|\Delta_k g\|_{L_T^1(L^2)}^q \right)^{\frac{1}{q}}. \end{aligned}$$

If $q < +\infty$, then u belongs to $C([0, T]; \dot{B}_{2,q}^s)$.

The proof is similar to the case when the diffusion term $\tilde{\Delta} u$ is replaced by Δu . We refer the reader to [9] for details.

Now we introduce the logarithmic interpolation inequality (see [14]).

PROPOSITION 4.4. *For any $1 \leq p, \rho \leq +\infty$, $s \in \mathbb{R}$, and $0 < \epsilon \leq 1$, we have*

$$(4.3) \quad \|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,1}^s)} \leq C \frac{\|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,\infty}^s)}}{\epsilon} \log \left(e + \frac{\|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,\infty}^{s-\epsilon})} + \|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,\infty}^{s+\epsilon})}}{\|f\|_{\tilde{L}_T^\rho(\dot{B}_{p,\infty}^s)}} \right).$$

Now, let us prove the uniqueness of the solution of (3.1). Let $(u_1, h_1), (u_2, h_2) \in (L_T^\infty(B^0) \cap L_T^1(B^2)) \times L_T^\infty(\tilde{B}^{0,1})$ be two solutions of (3.1) with the same initial data. The difference $\vartheta \triangleq h_2 - h_1$, $w \triangleq u_2 - u_1$ satisfies the following system:

$$(4.4) \quad \begin{cases} \partial_t \vartheta + u_2 \cdot \nabla \vartheta = -\operatorname{div} w - w \nabla h_1 - \vartheta \operatorname{div} u_2 - h_1 \operatorname{div} w, \\ \partial_t w - \nu \tilde{\Delta} w = -\nabla \vartheta - u_2 \cdot \nabla w - w \cdot \nabla u_1 + \nu(1 + h_1) \nabla \left(\frac{h_1}{1 + h_1} \right) \tilde{\nabla} w \\ \quad + \nu(1 + h_1) \nabla \left(\frac{h_2}{1 + h_2} - \frac{h_1}{1 + h_1} \right) \tilde{\nabla} u_2 + \nu \vartheta \nabla \left(\frac{h_2}{1 + h_2} \right) \tilde{\nabla} u_2, \\ \vartheta(0, x) = 0, \quad w(0, x) = 0. \end{cases}$$

Without loss of generality, we assume that there holds for sufficiently small T

$$(4.5) \quad 1 + h_1 \geq \frac{1}{2},$$

$$(4.6) \quad \|h_1\|_{\tilde{E}_T^{0,1}} \leq \epsilon,$$

where $\epsilon > 0$ is small enough. Applying Proposition 4.2 to the first equation of (4.4) yields

$$(4.7) \quad \|\vartheta(t)\|_{\dot{B}_{2,\infty}^0} \lesssim \int_0^t e^{C(V_2(t)-V_2(\tau))} \|w \cdot \nabla h_1 + \vartheta \operatorname{div} u_2 + h_1 \operatorname{div} w + \operatorname{div} w\|_{\dot{B}_{2,\infty}^0} d\tau,$$

with $V_2(t) \triangleq \int_0^t \|\nabla u_2\|_{\dot{B}_{2,\infty}^1 \cap L^\infty} d\tau$. It follows from (2.5) with $s = 0$ that

$$\begin{aligned} \|w \cdot \nabla h_1\|_{\dot{B}_{2,\infty}^0} &\lesssim \|\nabla h_1\|_{\dot{B}_{2,\infty}^0} \|w\|_{B^1} \lesssim \|w\|_{B^1} \|h_1\|_{B^1}, \\ \|\vartheta \operatorname{div} u_2\|_{\dot{B}_{2,\infty}^0} &\lesssim \|\vartheta\|_{\dot{B}_{2,\infty}^0} \|u_2\|_{B^2}, \\ \|h_1 \operatorname{div} w\|_{\dot{B}_{2,\infty}^0} &\lesssim \|\operatorname{div} w\|_{\dot{B}_{2,\infty}^0} \|h_1\|_{B^1} \lesssim \|w\|_{B^1} \|h_1\|_{B^1}, \end{aligned}$$

where we have used $B^1 \hookrightarrow \dot{B}_{2,\infty}^1$. Plugging the above estimates into (4.7), we get

$$(4.8) \quad \|\vartheta(t)\|_{\dot{B}_{2,\infty}^0} \lesssim \int_0^t e^{C(V_2(t)-V_2(\tau))} \left[\|w\|_{B^1} (1 + \|h_1\|_{B^1}) + \|\vartheta\|_{\dot{B}_{2,\infty}^0} \|u_2\|_{B^2} \right] d\tau.$$

Recalling that $u^i \in L_T^1(B^2)$, we can take a $T \in (0, \infty)$ small enough so that

$$C \|u_2\|_{L_T^1(B^2)} \leq \frac{1}{4},$$

which together with (4.8) implies that for $t \leq T$

$$(4.9) \quad \|\vartheta\|_{L_t^\infty(\dot{B}_{2,\infty}^0)} \lesssim \|w\|_{L_t^1(B^1)} (1 + \|h_1\|_{L_t^\infty(B^1)}).$$

Applying (4.3) to the term $\|w\|_{L_t^1(B^1)}$ yields

$$(4.10) \quad \|\vartheta\|_{L_t^\infty(\dot{B}_{2,\infty}^0)} \lesssim \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)} \log \left(e + \frac{\|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^0)} + \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^2)}}{\|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)}} \right) (1 + \|h_1\|_{L_t^\infty(B^1)}).$$

Thanks to $B^s \hookrightarrow \dot{B}_{2,\infty}^s$ and $\tilde{B}^{0,1} \hookrightarrow B^1$, we have

$$(4.11) \quad \|\vartheta\|_{L_t^\infty(\dot{B}_{2,\infty}^0)} \lesssim \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)} \log \left(e + \frac{W(t)}{\|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)}} \right)$$

with

$$W(t) \triangleq \sum_{i=1}^2 \|u^i\|_{\tilde{L}_t^1(B^0)} + \|u^i\|_{\tilde{L}_t^1(B^2)},$$

and for finite t , $W(t) < +\infty$.

Next, we deal with the second equation of (4.4). We get by applying (2.7) with $s = 1$, $s = 0$, respectively, that

$$(4.12) \quad \|u_2 \cdot \nabla w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})} \lesssim \|u_2\|_{L_t^2(B^1)} \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)},$$

$$(4.13) \quad \|w \cdot \nabla u_1\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})} \lesssim \|u_1\|_{L_t^2(B^1)} \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)}.$$

We can deduce $h_i \in C(0, T; \mathbb{R}^2)$ ($i = 1, 2$) from the fact $B^1 \hookrightarrow C$. Moreover, due to (4.5), we can assume $h_1(t, x) + 1 \geq \frac{1}{2}$ for all $t \leq T$, $x \in \mathbb{R}^2$. Since h_1, h_2 have the same initial data, from the continuity of h_2 , there exists a $\tilde{T} \leq T$ such that

$$h_2(x, t) + 1 \geq \frac{1}{4} \quad \text{for all } t \in [0, \tilde{T}], \quad x \in \mathbb{R}^2.$$

It follows from (2.6) with $s = 1$, (2.15), and $B^1 \hookrightarrow \dot{B}_{2,\infty}^1 \cap L^\infty$ that

$$\begin{aligned} & \left\| (1 + h_1) \nabla \left(\frac{h_2}{1 + h_2} - \frac{h_1}{1 + h_1} \right) \tilde{\nabla} u_2 \right\|_{\dot{B}_{2,\infty}^{-1}} \\ & \lesssim \left\| (1 + h_1) \nabla \left(\frac{h_2}{1 + h_2} - \frac{h_1}{1 + h_1} \right) \right\|_{\dot{B}_{2,\infty}^{-1}} \|\tilde{\nabla} u_2\|_{B^1} \\ & \lesssim (1 + \|h_1\|_{B^1}) \left\| \frac{h_2}{1 + h_2} - \frac{h_1}{1 + h_1} \right\|_{\dot{B}_{2,\infty}^0} \|u_2\|_{B^2} \\ & \lesssim (1 + \|h_1\|_{B^1}) (\|h_1\|_{B^1} + \|h_2\|_{B^1}) \|\vartheta\|_{\dot{B}_{2,\infty}^0} \|u_2\|_{B^2}, \end{aligned}$$

which, together with $L_t^1(\dot{B}_{2,\infty}^{-1}) \subset \tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})$, yields

$$(4.14) \quad \begin{aligned} & \left\| (1 + h_1) \nabla \left(\frac{h_2}{1 + h_2} - \frac{h_1}{1 + h_1} \right) \tilde{\nabla} u_2 \right\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})} \\ & \lesssim \int_0^t (1 + \|h_1\|_{B^1}) (\|h_1\|_{B^1} + \|h_2\|_{B^1}) \|\vartheta\|_{\dot{B}_{2,\infty}^0} \|u_2\|_{B^2} d\tau. \end{aligned}$$

Thanks to (2.6), (2.12), and $L_t^1(\dot{B}_{2,\infty}^{-1}) \subset \tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})$, we get

$$(4.15) \quad \left\| \vartheta \nabla \left(\frac{h_2}{1+h_2} \right) \tilde{\nabla} u_2 \right\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^{-1})} \lesssim \int_0^t \|\vartheta\|_{\dot{B}_{2,\infty}^0} \|h_2\|_{B^1} \|u_2\|_{B^2} d\tau.$$

Thanks to Lemma 5.3 with $s_1 = s_2 = 0$, Lemma 5.4 with $s = 1$, and (2.5) with $s = 0$, we have

$$(4.16) \quad \begin{aligned} & \sup_{k \in \mathbb{Z}} \omega_k(t) 2^{-k} \left\| \Delta_k \left((1+h_1) \nabla \left(\frac{h_1}{1+h_1} \right) \tilde{\nabla} w \right) \right\|_{L_t^1(L^2)} \\ & \lesssim \left\| \nabla \left(\frac{h_1}{1+h_1} \right) \right\|_{E_t^0} \|(1+h_1) \tilde{\nabla} w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^0)} \\ & \lesssim \left\| \frac{h_1}{1+h_1} \right\|_{E_t^1} (1 + \|h_1\|_{L_t^\infty(B^1)}) \|\nabla w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^0)} \\ & \lesssim \|h_1\|_{\tilde{E}_t^{0,1}} (1 + \|h_1\|_{L_t^\infty(\tilde{B}^{0,1})})^4 \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)}. \end{aligned}$$

In terms of Proposition 4.3, (4.12)–(4.16), and $\tilde{B}^{0,1} \hookrightarrow B^1$, we finally obtain

$$(4.17) \quad \begin{aligned} & \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)} + \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)} \\ & \lesssim \|u_2\|_{L_t^2(B^1)} \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)} + \|u_1\|_{L_t^2(B^1)} \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)} \\ & \quad + \|h_1\|_{\tilde{E}_t^{0,1}} (1 + \|h_1\|_{L_t^\infty(\tilde{B}^{0,1})})^4 \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)} \\ & \quad + \int_0^t (1 + \|h_1\|_{\tilde{B}^{0,1}}) (1 + \|h_1\|_{\tilde{B}^{0,1}} + \|h_2\|_{\tilde{B}^{0,1}}) (1 + \|u_2\|_{B^2}) \|\vartheta\|_{\dot{B}_{2,\infty}^0} d\tau. \end{aligned}$$

Let us define

$$Z(t) \triangleq \|w\|_{\tilde{L}_t^1(\dot{B}_{2,\infty}^1)} + \|w\|_{\tilde{L}_t^2(\dot{B}_{2,\infty}^0)}.$$

Due to (4.6), if T is chosen small enough, then the first three terms on the right side of (4.17) can be absorbed by the left side $Z(t)$. Noting that $r \log(e + \frac{W(T)}{r})$ is increasing, from (4.11) and (4.17), it follows that

$$(4.18) \quad \begin{aligned} Z(t) & \lesssim \int_0^t (1 + W'(\tau)) Z(\tau) \log \left(e + \frac{W(\tau)}{Z(\tau)} \right) d\tau \\ & \lesssim \int_0^t (1 + W'(\tau)) Z(\tau) \log \left(e + \frac{W(T)}{Z(\tau)} \right) d\tau. \end{aligned}$$

It is easy to verify that

$$1 + W'(\tau) \in L_{loc}^1(\mathbb{R}^+) \quad \text{and} \quad \int_0^1 \frac{dr}{r \log(e + \frac{W(T)}{r})} = +\infty.$$

Hence by the Osgood lemma, we have $Z \equiv 0$ on $[0, \tilde{T}]$, i.e., $w \equiv 0$; then from (4.9), $\vartheta = h_2 - h_1 \equiv 0$. Then a standard continuous argument gives the uniqueness.

5. Appendix. In this appendix, we prove some multilinear estimates in the weighted Besov space.

LEMMA 5.1. *Let A be a homogeneous smooth function of degree m . Assume that $-\frac{d}{2} < \rho \leq \frac{d}{2}$. Then it holds that*

$$(5.1) \quad \left| (A(D)\Delta_k(v \cdot \nabla h), A(D)\Delta_k h) \right| \leq C \|\mathcal{F}_k^m(t)\|_2 \|A(D)\Delta_k h\|_2,$$

$$(5.2) \quad \left| (A(D)\Delta_k(v \cdot \nabla u), A(D)\Delta_k u) \right| \leq C \|\tilde{\mathcal{F}}_k^m(t)\|_2 \|A(D)\Delta_k u\|_2,$$

and

$$(5.3) \quad \left| (A(D)\Delta_k(v \cdot \nabla h), \Delta_k u) + (\Delta_k(v \cdot \nabla u), A(D)\Delta_k h) \right| \leq C (\|\mathcal{F}_k^m(t)\|_2 + \|\tilde{\mathcal{F}}_k^0(t)\|_2) (\|\Delta_k u\|_2 + \|A(D)\Delta_k h\|_2),$$

where $\mathcal{F}_k^m(t)$ and $\tilde{\mathcal{F}}_k^m(t)$ satisfy

$$(5.4) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(\rho-m)} \|\mathcal{F}_k^m(t)\|_{L_T^1(L^2)} \leq C \|h\|_{E_T^\rho} \|v\|_{L_T^1(B^{\frac{d}{2}+1})},$$

$$(5.5) \quad \sum_{k \in \mathbb{Z}} 2^{k(\rho-m)} \|\tilde{\mathcal{F}}_k^m(t)\|_{L_T^1(L^2)} \leq C \|u\|_{L_T^2(B^{\rho+1})} \|v\|_{L_T^2(B^{\frac{d}{2}})}.$$

Proof. Let us first prove (5.1). Using Bony's paraproduct decomposition, we write

$$(5.6) \quad (A(D)\Delta_k(v \cdot \nabla h), A(D)\Delta_k h) = (A(D)\Delta_k(T'_{\partial_j h} v^j), A(D)\Delta_k h) + J_k,$$

where

$$\begin{aligned} T'_f g &= T_f g + R(f, g), \quad \text{and} \\ J_k &= \sum_{|k'-k| \leq 3} ([A(D)\Delta_k, S_{k'-1} v^j] \Delta_{k'} \partial_j h, A(D)\Delta_k h) \\ &\quad + \sum_{|k'-k| \leq 3} ((S_{k'-1} - S_{k-1}) v^j A(D)\Delta_k \Delta_{k'} \partial_j h, A(D)\Delta_k h) \\ &\quad + (S_{k-1} v^j A(D)\Delta_k \partial_j h, A(D)\Delta_k h). \end{aligned}$$

We get by integration by parts that

$$(S_{k-1} v^j A(D)\Delta_k \partial_j h, A(D)\Delta_k h) = -\frac{1}{2} (S_{k-1} \operatorname{div} v A(D)\Delta_k h, A(D)\Delta_k h).$$

Let us set

$$\begin{aligned} \mathcal{F}_{k,0}^m(t) &= A(D)\Delta_k(T'_{\partial_j h} v^j), \\ \mathcal{F}_{k,1}^m(t) &= \sum_{|k'-k| \leq 3} [A(D)\Delta_k, S_{k'-1} v^j] \Delta_{k'} \partial_j h, \\ \mathcal{F}_{k,2}^m(t) &= \sum_{|k'-k| \leq 3} (S_{k'-1} - S_{k-1}) v^j A(D)\Delta_k \Delta_{k'} \partial_j h, \\ \mathcal{F}_{k,3}^m(t) &= -\frac{1}{2} S_{k-1} \operatorname{div} v A(D)\Delta_k h. \end{aligned}$$

By the Cauchy–Schwarz inequality, we get

$$(A(D)\Delta_k(v \cdot \nabla h), A(D)\Delta_k h) \leq \|\mathcal{F}_k^m(t)\|_2 \|A(D)\Delta_k h\|_2,$$

with $\mathcal{F}_k^m(t) = \sum_{i=0}^3 \mathcal{F}_{k,i}^m(t)$. So, it remains to prove that $\mathcal{F}_k^m(t)$ satisfies (5.4). For the simplicity, we set

$$\tilde{\Delta}_k = \sum_{|k'-k| \leq 1} \Delta_{k'}, \quad \tilde{\tilde{\Delta}}_k = \sum_{|k'-k| \leq 3} \Delta_{k'}.$$

Thanks to the definition of $\mathcal{F}_{k,0}^m(t)$ and Lemma 2.1, we have

$$\begin{aligned} \|\mathcal{F}_{k,0}^m(t)\|_{L_T^1(L^2)} &\leq \sum_{|k'-k| \leq 3} 2^{km} \|S_{k'-1} \partial_j h\|_{L_T^\infty(L^\infty)} \|\Delta_{k'} v^j\|_{L_T^1(L^2)} \\ &\quad + \sum_{k' \geq k-2} 2^{(m+\frac{d}{2})k} \|\Delta_k(\Delta_{k'} \partial_j h \tilde{\tilde{\Delta}}_{k'} v^j)\|_{L_T^1(L^1)} \\ &\triangleq I + II. \end{aligned}$$

Thanks to Lemma 2.1, we have

$$\begin{aligned} 2^{k(\rho-m)} I &\lesssim 2^{k\rho} \sum_{k' \leq k+1} 2^{k'(1+\frac{d}{2})} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \|\tilde{\tilde{\Delta}}_k v\|_{L_T^1(L^2)} \\ &\lesssim \sum_{k' \leq k+1} 2^{(k'-k)(1+\frac{d}{2}-\rho)} 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} 2^{k(1+\frac{d}{2})} \|\tilde{\tilde{\Delta}}_k v\|_{L_T^1(L^2)}, \end{aligned}$$

from which, together with the definition of $\omega_k(T)$, it follows that

$$\begin{aligned} &\sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(\rho-m)} I \\ &\lesssim \sum_{k' \in \mathbb{Z}} 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \sum_{k \geq k'-1} \omega_k(T) 2^{(k'-k)(1+\frac{d}{2}-\rho)} 2^{k(1+\frac{d}{2})} \|\tilde{\tilde{\Delta}}_k v\|_{L_T^1(L^2)} \\ &\lesssim \sum_{k' \in \mathbb{Z}} \omega_{k'}(T) 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \sum_{k \geq k'-1} 2^{(k'-k)(\frac{d}{2}-\rho)} 2^{k(1+\frac{d}{2})} \|\tilde{\tilde{\Delta}}_k v\|_{L_T^1(L^2)} \\ (5.7) \quad &\lesssim \|h\|_{E_T^\rho} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}, \end{aligned}$$

where we used the assumption $\rho \leq \frac{d}{2}$ in the last inequality. Set $e_k(T) = e_k^1(T) + e_k^2(T)$. Using Lemma 2.1, we also have

$$\begin{aligned} \omega_k(T) 2^{k(\rho-m)} II &\lesssim \omega_k(T) 2^{k(\rho+\frac{d}{2})} \sum_{k' \geq k-2} 2^{k'} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \|\tilde{\tilde{\Delta}}_{k'} v\|_{L_T^1(L^2)} \\ &\lesssim 2^{k(\rho+\frac{d}{2})} \sum_{k' \geq k-2} 2^{k'} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \|\tilde{\tilde{\Delta}}_{k'} v\|_{L_T^1(L^2)} \sum_{k' \geq \tilde{k} \geq k} 2^{-(\tilde{k}-k)} e_{\tilde{k}}(T) \\ &\quad + 2^{k(\rho+\frac{d}{2})} \sum_{k' \geq k-2} 2^{k'} \|\Delta_{k'} h\|_{L_T^\infty(L^2)} \|\tilde{\tilde{\Delta}}_{k'} v\|_{L_T^1(L^2)} \\ &\quad \sum_{\tilde{k} \geq k, \tilde{k} \geq k'} 2^{-(\tilde{k}-k)} e_{\tilde{k}}(T) \\ &\triangleq II_1 + II_2. \end{aligned}$$

Note that for $\tilde{k} \leq k'$

$$e_{\tilde{k}}^{\gamma}(T) \leq e_{k'}(T) \leq \omega_{k'}(T),$$

from which, together with $\rho > -\frac{d}{2}$, we deduce that

$$\begin{aligned} \sum_{k \in \mathbb{Z}} II_1 &\lesssim \sum_{k' \in \mathbb{Z}} \omega_{k'}(T) 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^{\infty}(L^2)} 2^{k'(\frac{d}{2}+1)} \|\tilde{\Delta}_{k'} v\|_{L_T^1(L^2)} \sum_{k \leq k'+2} 2^{(k-k')(\rho+\frac{d}{2})} \\ (5.8) \quad &\lesssim \|h\|_{E_T^{\rho}} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}. \end{aligned}$$

Similarly, we can obtain

$$\begin{aligned} \sum_{k \in \mathbb{Z}} II_2 &\lesssim \sum_{k' \in \mathbb{Z}} 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^{\infty}(L^2)} \sum_{k \leq k'+2} 2^{(k-k')(\frac{d}{2}+\rho)} \sum_{\tilde{k} \geq k'} 2^{-(\tilde{k}-k)} e_{\tilde{k}}^{\gamma}(T) \|v\|_{L_T^1(B^{\frac{d}{2}+1})} \\ &\lesssim \sum_{k' \in \mathbb{Z}} \omega_{k'}(T) 2^{k'\rho} \|\Delta_{k'} h\|_{L_T^{\infty}(L^2)} \sum_{k \leq k'+2} 2^{(k-k')(\frac{d}{2}+\rho+1)} \|v\|_{L_T^1(B^{\frac{d}{2}+1})} \\ (5.9) \quad &\lesssim \|h\|_{E_T^{\rho}} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}. \end{aligned}$$

By summing up (5.7)–(5.9), we obtain

$$(5.10) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(\rho-m)} \|\mathcal{F}_{k,0}^m(t)\|_{L_T^1(L^2)} \lesssim \|h\|_{E_T^{\rho}} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}.$$

Note that $A(D)\Delta_k = 2^{km}\tilde{\varphi}(2^{-k}D)$ with $\tilde{\varphi}(\xi) = A(\xi)\varphi(\xi)$. Setting $\tilde{\theta} = \mathcal{F}^{-1}\tilde{\varphi}$, we get by using Taylor's formula that

$$\begin{aligned} \mathcal{F}_{k,1}^m(t) &= \sum_{|k'-k| \leq 3} 2^{k(m-1)} \int_{\mathbb{R}^d} \int_0^1 \tilde{\theta}(y) \\ &\quad (y \cdot S_{k'-1} \nabla v^j(x - 2^{-k}\tau y)) \Delta_{k'} \partial_j h(x - 2^{-k}y) d\tau dy, \end{aligned}$$

from which, together with Lemma 2.1, it follows that

$$\begin{aligned} \|\mathcal{F}_{k,1}^m(t)\|_{L_T^1(L^2)} &\lesssim 2^{k(m-1)} \sum_{|k'-k| \leq 3} \|S_{k'-1} \nabla v^j\|_{L_T^1(L^{\infty})} \|\Delta_{k'} \partial_j h\|_{L_T^{\infty}(L^2)} \\ &\lesssim 2^{km} \sum_{|k'-k| \leq 3} \|\Delta_{k'} h\|_{L_T^{\infty}(L^2)} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}; \end{aligned}$$

thus, we get

$$(5.11) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(\rho-m)} \|\mathcal{F}_{k,1}^m(t)\|_{L_T^1(L^2)} \lesssim \|h\|_{E_T^{\rho}} \|v\|_{L_T^1(B^{\frac{d}{2}+1})}.$$

Thanks to the fact $|k' - k| \leq 3$ and Lemma 2.1, we have

$$\|(S_{k'-1} - S_{k-1})v^j A(D)\Delta_k \Delta_{k'} \partial_j h\|_{L_T^1(L^2)} \lesssim 2^{km} \|\Delta_k h\|_{L_T^{\infty}(L^2)} \|v\|_{L_T^1(B^{\frac{d}{2}+1})},$$

from which it follows that

$$(5.12) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(\rho-m)} (\|\mathcal{F}_{k,2}^m(t)\|_{L_T^1(L^2)} + \|\mathcal{F}_{k,3}^m(t)\|_{L_T^1(L^2)}) \lesssim \|h\|_{E_T^\rho} \|v\|_{L_T^1(B^{\frac{d}{2}+1})},$$

which, together with (5.10) and (5.11), yields (5.4).

Using the decomposition (5.6), with h instead of u , and Lemma 2.1, (5.2) can be easily proved. We omit the proof here. In order to prove (5.3), we use the decomposition

$$(A(D)\Delta_k(v \cdot \nabla h), \Delta_k u) + (\Delta_k(v \cdot \nabla u), A(D)\Delta_k h) = I_k + J_k,$$

with

$$\begin{aligned} I_k &= (A(D)\Delta_k(T'_{\partial_j h} v^j), \Delta_k u) + (\Delta_k(T'_{\partial_j u} v^j), A(D)\Delta_k h) \\ &\triangleq (\mathcal{F}_{k,0}^m(t), \Delta_k u) + (\tilde{\mathcal{F}}_{k,0}^0(t), A(D)\Delta_k h), \\ J_k &= \sum_{|k'-k| \leq 3} \left([A(D)\Delta_k, S_{k'-1} v^j] \Delta_{k'} \partial_j h, \Delta_k u \right) \\ &\quad + \left((S_{k'-1} - S_{k-1}) v^j A(D)\Delta_k \Delta_{k'} \partial_j h, \Delta_k u \right) \\ &\quad + \sum_{|k'-k| \leq 3} \left([\Delta_k, S_{k'-1} v^j] \Delta_{k'} \partial_j u, A(D)\Delta_k h \right) \\ &\quad + \left((S_{k'-1} - S_{k-1}) v^j \Delta_k \Delta_{k'} \partial_j u, A(D)\Delta_k h \right) \\ &\quad - \left(S_{k-1} \operatorname{div} v A(D)\Delta_k h, \Delta_k u \right) \\ &\triangleq (\mathcal{F}_{k,1}^m(t), \Delta_k u) + (\mathcal{F}_{k,2}^m(t), \Delta_k u) + (\tilde{\mathcal{F}}_{k,1}^0(t), A(D)\Delta_k h) \\ &\quad + (\tilde{\mathcal{F}}_{k,2}^0, A(D)\Delta_k h) + (\mathcal{F}_{k,3}^m(t), \Delta_k u), \end{aligned}$$

from which, using a similar proof of (5.4), we obtain (5.3). This completes the proof of Lemma 5.1. \square

LEMMA 5.2. *Let $s_1 \leq \frac{d}{2} - 1$, $s_2 \leq \frac{d}{2}$, and $s_1 + s_2 > 0$. Then it holds that*

$$(5.13) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(s_1+s_2-\frac{d}{2})} \|\Delta_k(fg)\|_{L_T^1(L^2)} \leq C \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks_1} \|\Delta_k f\|_{L_T^{r_1}(L^2)} \|g\|_{L_T^{r_2}(B^{s_2})},$$

where $1 \leq r_1, r_2 \leq \infty$ and $\frac{1}{r_1} + \frac{1}{r_2} = 1$.

Proof. Using Bony's paraproduct decomposition, we write

$$\begin{aligned} \Delta_k(fg) &= \sum_{|k'-k| \leq 3} \Delta_k(S_{k'-1} f \Delta_{k'} g) + \sum_{|k'-k| \leq 3} \Delta_k(S_{k'-1} g \Delta_{k'} f) \\ &\quad + \sum_{k' \geq k-2} \Delta_k(\Delta_{k'} f \tilde{\Delta}_{k'} g) \triangleq I + II + III. \end{aligned}$$

A similar proof of (5.7) ensures that for $s_1 \leq \frac{d}{2} - 1$

$$\sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(s_1+s_2-\frac{d}{2})} \|I\|_{L_T^1(L^2)} \lesssim \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks_1} \|\Delta_k f\|_{L_T^{r_1}(L^2)} \|g\|_{L_T^{r_2}(B^{s_2})},$$

while II can be directly deduced for $s_2 \leq \frac{d}{2}$. On the other hand, a similar proof of (5.8) and (5.9) gives for $s_1 + s_2 > 0$

$$\sum_{k \in \mathbb{Z}} \omega_k(T) 2^{k(s_1+s_2-\frac{d}{2})} \|III\|_{L_T^1(L^2)} \lesssim \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks_1} \|\Delta_k f\|_{L_T^{r_1}(L^2)} \|g\|_{L_T^{r_2}(B^{s_2})}.$$

This completes the proof of Lemma 5.2. \square

Similarly, we can also prove the following lemma.

LEMMA 5.3. *Let $s_1 \leq \frac{d}{2} - 1$, $s_2 < \frac{d}{2}$, and $s_1 + s_2 \geq 0$. Then it holds that*

$$(5.14) \quad \sup_{k \in \mathbb{Z}} \omega_k(T) 2^{k(s_1+s_2-\frac{d}{2})} \|\Delta_k(fg)\|_{L_T^1(L^2)} \leq C \|f\|_{E_T^{s_1}} \|g\|_{\tilde{L}_T^1(\dot{B}_2^{s_2, \infty})}.$$

LEMMA 5.4. *Let $s > 0$. Assume that $F \in W_{loc}^{[s]+3, \infty}(\mathbb{R}^d)$ with $F(0) = 0$. Then it holds that*

$$(5.15) \quad \|F(f)\|_{E_T^s} \leq C(1 + \|f\|_{L_T^\infty(L^\infty)})^{[s]+2} \|f\|_{E_T^s}.$$

Proof. We decompose $F(f)$ as

$$\begin{aligned} F(f) &= \sum_{k' \in \mathbb{Z}} F(S_{k'+1}f) - F(S_{k'}f) = \sum_{k' \in \mathbb{Z}} \Delta_{k'} f \int_0^1 F'(S_{k'}f + \tau \Delta_{k'} f) d\tau \\ &\triangleq \sum_{k' \in \mathbb{Z}} \Delta_{k'} f m_{k'}, \end{aligned}$$

where $m_{k'} = \int_0^1 F'(S_{k'}f + \tau \Delta_{k'} f) d\tau$. Furthermore, we write

$$\Delta_k F(f) = \sum_{k' < k} \Delta_k(\Delta_{k'} f m_{k'}) + \sum_{k' \geq k} \Delta_k(\Delta_{k'} f m_{k'}) \triangleq I + II.$$

By Lemma 2.1, we have

$$(5.16) \quad \begin{aligned} \|I\|_{L_T^\infty(L^2)} &\leq \sum_{k' < k} \|\Delta_k(\Delta_{k'} f m_{k'})\|_{L_T^\infty(L^2)} \\ &\leq \sum_{k' < k} 2^{-k|\alpha|} \sup_{|\gamma|=|\alpha|} \|D^\gamma \Delta_k(\Delta_{k'} f m_{k'})\|_{L_T^\infty(L^2)}, \end{aligned}$$

with α to be determined later. Note that for $|\gamma| \geq 0$, we have

$$\|D^\gamma m_{k'}\|_\infty \lesssim 2^{k'|\gamma|} (1 + \|f\|_\infty)^{|\gamma|} \|F'\|_{W^{|\gamma|, \infty}},$$

from which, together with (5.16), it follows that

$$2^{ks} \|I\|_{L_T^\infty(L^2)} \lesssim 2^{k(s-|\alpha|)} \sum_{k' < k} 2^{k'|\alpha|} \|\Delta_{k'} f\|_{L_T^\infty(L^2)} (1 + \|f\|_{L_T^\infty(L^\infty)})^{|\alpha|} \|F'\|_{W^{|\alpha|, \infty}}.$$

Thus, if we take $|\alpha| = [s] + 2$, we get

$$\begin{aligned}
 & \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks} \|I\|_{L_T^\infty(L^2)} \\
 & \lesssim \sum_{k' \in \mathbb{Z}} 2^{k's} \omega_{k'}(T) \|\Delta_{k'} f\|_{L_T^\infty(L^2)} \sum_{k > k'} 2^{(k-k')(s-|\alpha|+1)} (1 + \|f\|_{L_T^\infty(L^\infty)})^{|\alpha|} \|F'\|_{W^{|\alpha|, \infty}} \\
 (5.17) \quad & \lesssim (1 + \|f\|_{L_T^\infty(L^\infty)})^{[s]+2} \|F'\|_{W^{[s]+2, \infty}} \|f\|_{E_T^s}.
 \end{aligned}$$

Next, let us turn to the proof of II . We get by using Lemma 2.1 that

$$\|II\|_{L_T^\infty(L^2)} \lesssim \sum_{k' \geq k} \|\Delta_{k'} f\|_{L_T^\infty(L^2)}.$$

Then we write

$$\begin{aligned}
 \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks} \|II\|_{L_T^\infty(L^2)} & \lesssim \sum_{k \in \mathbb{Z}} 2^{ks} \sum_{k' \geq k} \|\Delta_{k'} f\|_{L_T^\infty(L^2)} \sum_{k' \geq \tilde{k} \geq k} 2^{-(\tilde{k}-k)} e_{\tilde{k}}(T) \\
 & + \sum_{k \in \mathbb{Z}} 2^{ks} \sum_{k' \geq k} \|\Delta_{k'} f\|_{L_T^\infty(L^2)} \sum_{\tilde{k} \geq k', \tilde{k} \geq k} 2^{-(\tilde{k}-k)} e_{\tilde{k}}(T),
 \end{aligned}$$

from which a similar proof of (5.8) and (5.9) ensures that

$$(5.18) \quad \sum_{k \in \mathbb{Z}} \omega_k(T) 2^{ks} \|II\|_{L_T^\infty(L^2)} \lesssim \|f\|_{E_T^s}.$$

By summing up (5.17) and (5.18), we deduce the inequality (5.15). This completes the proof of Lemma 5.4. \square

Acknowledgment. The authors thank the referees for their invaluable comments and suggestions which helped improve the paper greatly.

REFERENCES

- [1] H. BAHOURI AND J.-Y. CHEMIN, *Equations de transport relatives à des champs des vecteurs non-lipschitziens et mécanique des fluides*, Arch. Rational Mech. Anal., 127 (1994), pp. 159–199.
- [2] H. BAHOURI AND J.-Y. CHEMIN, *Équations d'ondes quasilinéaires et estimations de Strichartz*, Amer. J. Math., 121 (1999), pp. 1337–1377.
- [3] J.-M. BONY, *Calcul symbolique et propagation des singularités pour les équations aux dérivées partielles non linéaires*, Ann. Sci. École Norm. Sup. (4), 14 (1981), pp. 209–246.
- [4] D. BRESCH, B. DESJARDINS, AND G. MÉTIVIER, *Recent Mathematical Results and Open Problems about Shallow Water Equations*, preprint.
- [5] A. T. BUI, *Existence and uniqueness of a classical solution of an initial boundary value problem of the theory of shallow waters*, SIAM J. Math. Anal., 12 (1981), pp. 229–241.
- [6] M. CANNONE, *Ondelettes, paraproducts et Navier-Stokes*, Nouveaux essais, Diderot éditeurs, Paris, 1995.
- [7] M. CANNONE, *Harmonic analysis tools for solving the incompressible Navier–Stokes equations*, in Handbook of Mathematical Fluid Dynamics, Vol. III, North-Holland, Amsterdam, 2004, pp. 161–244.
- [8] J.-Y. CHEMIN, *Perfect Incompressible Fluids*, Oxford University Press, New York, 1998.
- [9] J.-Y. CHEMIN, *Théorèmes d'unicité pour le système de Navier-Stokes tridimensionnel*, J. Anal. Math., 77 (1999), pp. 27–50.

- [10] J.-Y. CHEMIN AND N. LERNER, *Flot de champs de vecteurs non lipschitziens et équations de Navier-Stokes*, J. Differential Equations, 121 (1992), pp. 314–328.
- [11] R. DANCHIN, *Global existence in critical spaces for compressible Navier-Stokes equations*, Invent. Math., 141 (2000), pp. 579–614.
- [12] R. DANCHIN, *Global existence in critical spaces for flows of compressible viscous and heat-conductive gases*, Arch. Rational Mech. Anal., 160 (2001), pp. 1–39.
- [13] R. DANCHIN, *Local theory in critical spaces for compressible viscous and heat-conductive gases*, Comm. Partial Differential Equations, 26 (2001), pp. 1183–1233.
- [14] R. DANCHIN, *Density-dependent incompressible viscous fluids in critical spaces*, Proc. Roy. Soc. Edinburgh Sect. A, 133 (2003), pp. 1311–1334.
- [15] R. DANCHIN, *On the uniqueness in critical spaces for compressible Navier-Stokes equations*, Nonlinear Differential Equations Appl., 12 (2005), pp. 111–128.
- [16] H. FUJITA AND T. KATO, *On the Navier-Stokes initial value problem I*, Arch. Rational Mech. Anal., 16 (1964), pp. 269–315.
- [17] P. E. KLOEDEN, *Global existence of classical solutions in the dissipative shallow water equations*, SIAM J. Math. Anal., 16 (1985), pp. 301–315.
- [18] Y. MEYER, *Wavelets, Paraproducts and Navier-Stokes Equations. Current Developments in Mathematics*, International Press, Boston, MA, 1996.
- [19] T. RUNST AND W. SICKEL, *Sobolev spaces of fractional order, Nemytskij operators, and nonlinear partial differential equations*, de Gruyter Ser. Nonlinear Anal. Appl. 3, Walter de Gruyter, Berlin, 1996.
- [20] L. SUNDBYE, *Global existence for the Dirichlet problem for the viscous shallow water equations*, J. Math. Anal. Appl., 202 (1996), pp. 236–258.
- [21] L. SUNDBYE, *Global existence for the Cauchy problem for the viscous shallow water equations*, Rocky Mountain J. Math., 28 (1998), pp. 1135–1152.
- [22] H. TRIEBEL, *Theory of Function Spaces*, Monogr. Math. 78, Birkhäuser Verlag, Basel, Boston, Stuttgart, 1983.
- [23] W.-K. WANG AND C.-J. XU, *The Cauchy problem for viscous shallow water equations*, Rev. Mat. Iberoamericana, 21 (2005), pp. 1–24.

BLOW-UP AND DECAY OF THE SOLUTION OF THE WEAKLY DISSIPATIVE DEGASPERIS–PROCESI EQUATION*

SHUYIN WU[†] AND ZHAOYANG YIN[†]

Abstract. In this paper, we mainly study several problems on the weakly dissipative Degasperis–Procesi equation. We first establish the local well-posedness of the equation, derive a precise blow-up scenario, and present two blow-up criteria for strong solutions to the equation. We then give the precise blow-up rate of blow-up solutions to the equation. We finally prove that the equation has global strong solutions and these global solutions decay to zero as time goes to infinity provided the potentials associated with their initial data are of one sign.

Key words. the weakly dissipative Degasperis–Procesi equation, blow-up, blow-up rate, global strong solutions, decay of solutions

AMS subject classifications. 35G25, 35L05

DOI. 10.1137/07070855X

1. Introduction. Recently, Degasperis and Procesi [20] studied the following family of third-order dispersive conservation laws:

$$(1.1) \quad u_t + c_0 u_x + \gamma u_{xxx} - \alpha^2 u_{txx} = (c_1 u^2 + c_2 u_x^2 + c_3 u u_{xx})_x,$$

where α , c_0 , c_1 , c_2 , and c_3 are real constants. They found [20] that there are only three equations that satisfy the asymptotic integrability condition within this family: the KdV equation, the Camassa–Holm equation, and the Degasperis–Procesi equation.

If $\alpha = c_2 = c_3 = 0$, then (1.1) becomes the well-known KdV equation which describes the unidirectional propagation of waves at the free surface of shallow water under the influence of gravity. In this model $u(t, x)$ represents the wave’s height above a flat bottom, x is proportional to the distance in the direction of propagation, and t is proportional to the elapsed time. The KdV equation is completely integrable, and its solitary waves are solitons [21, 39]. The Cauchy problem of the KdV equation has been the subject of a number of studies, and a satisfactory local or global (in time) existence theory is now in hand (for example, see [31, 43]). It is shown that the KdV equation is globally well-posed for $u_0 \in L^2(\mathbb{R})$; cf. [43]. It is observed that the KdV equation does not accommodate wave breaking (by wave breaking we mean that the wave remains bounded, but its slope becomes unbounded in finite time [45]).

For $c_1 = -\frac{3}{2}c_3/\alpha^2$ and $c_2 = c_3/2$, (1.1) becomes the Camassa–Holm equation, modeling the unidirectional propagation of shallow water waves over a flat bottom. Again $u(t, x)$ stands for the fluid velocity at time t in the spatial x direction, and c_0 is a nonnegative parameter related to the critical shallow water speed [3, 22, 30]. The Camassa–Holm equation is also a model for the propagation of axially symmetric waves in hyperelastic rods [15, 18]. It has a bi-Hamiltonian structure [33, 27] and is completely integrable [3, 8]. Its solitary waves are smooth if $c_0 > 0$ and peaked in the limiting case $c_0 = 0$; see [4]. The orbital stability of the peaked solitons is proved

*Received by the editors November 19, 2007; accepted for publication (in revised form) January 30, 2008; published electronically May 16, 2008.

<http://www.siam.org/journals/sima/40-2/70855.html>

[†]Department of Computer Science, Sun Yet-sen University, 510275 Guangzhou, China (isswsy@mail.sysu.edu.cn, mcsyzy@mail.sysu.edu.cn). The second author’s research was partially supported by the AvH Foundation, the NNSF of China (10531040), and the NSF of Guangdong Province.

in [16] and that of the smooth solitons in [17]. The explicit interaction of the peaked solitons is given in [1].

The Cauchy problem of the Camassa–Holm equation has been studied extensively. It has been shown that this equation is locally well-posed [9, 34, 42] for initial data $u_0 \in H^s(\mathbb{R})$, $s > \frac{3}{2}$. More interestingly, it has global strong solutions [7, 9] and also finite time blow-up solutions [7, 9, 10, 13]. On the other hand, it has global weak solutions in $H^1(\mathbb{R})$ [2, 11, 14, 47]. It is also known that if u is the solution of the Camassa–Holm equation with the initial data u_0 in $H^1(\mathbb{R})$, then we have the following a priori estimate:

$$\|u(t, \cdot)\|_{L^\infty(\mathbb{R})} \leq \sqrt{2}\|u(t, \cdot)\|_{H^1(\mathbb{R})} \leq \sqrt{2}\|u_0(\cdot)\|_{H^1(\mathbb{R})}$$

for all $t > 0$. The advantage of the Camassa–Holm equation in comparison with the KdV equation lies in the fact that the Camassa–Holm equation has peaked solitons and models wave breaking [4].

If $c_1 = -2c_3/\alpha^2$ and $c_2 = c_3$ in (1.1), then after rescaling, shifting the dependent variable, and applying a Galilean boost [19], we find the Degasperis–Procesi equation of the form

$$(1.2) \quad u_t - u_{txx} + 4uu_x = 3u_xu_{xx} + uu_{xxx}, \quad t > 0, \quad x \in \mathbb{R}.$$

The integrability of (1.2) was proved in [19] by constructing a Lax pair. It was also shown in [19] that (1.2) has a bi-Hamiltonian structure and admits exact peakon solutions which are analogous to the Camassa–Holm peakons.

The Degasperis–Procesi equation can be regarded as a model for the motion of shallow water waves, and its asymptotic accuracy is the same as for the Camassa–Holm shallow water equation [23, 29]. An inverse scattering approach for computing n -peakon solutions to (1.2) was presented in [37]. The traveling wave solutions to (1.2) were investigated in [44, 32]. The multisoliton solutions to (1.2) and their peakon limits were studied in [38].

The Cauchy problem of the Degasperis–Procesi equation is locally well-posed [48] for initial data $u_0 \in H^s(\mathbb{R})$ with $s > \frac{3}{2}$. Analogous to the Camassa–Holm equation, it not only has global strong solutions [35, 51] but also blow-up solutions in finite time [35, 48]. Meanwhile, the Degasperis–Procesi equation also has global weak solutions with initial data $u_0 \in H^1(\mathbb{R})$ (cf. [50, 51]) and global entropy weak solutions in the class of $L^1(\mathbb{R}) \cap BV(\mathbb{R})$ and the class of $L^2(\mathbb{R}) \cap L^4(\mathbb{R})$; cf. [5].

Despite the similarities to the Camassa–Holm equation, we would like to point out that these two equations are truly different. One of the important features of the Degasperis–Procesi equation is that it not only has peakon solitons [19] and periodic peakons [49] but also shock peakons [6, 36] and periodic shock waves [25]. On the other hand, the isospectral problem for the Degasperis–Procesi equation is the third-order equation in the Lax pair [19], while the isospectral problem for the Camassa–Holm equation is the second-order equation [3]. Another indication of the fact that there is no simple transformation of the Degasperis–Procesi equation into the Camassa–Holm equation is the entirely different form of conservation laws for these two equations [3, 19]. Furthermore, the Camassa–Holm equation is a reexpression of geodesic flow on the diffeomorphism group [12] or on the Bott–Virasoro group [40]. Up to now, no geometric derivation of the Degasperis–Procesi equation has been available.

Recently, several new global existence and blow-up results for strong solutions to the Degasperis–Procesi equation were presented in [35]. It is proved that the first blow-up must occur as wave breaking and shock waves possibly appear afterward [35].

Global weak solution and blow-up structure for this equation were investigated in [24]. Initial boundary value problems for the Degasperis–Procesi equation were also discussed in [26].

In general, it is quite difficult to avoid energy dissipation mechanisms in a real world. Ott and Sudan [41] ever investigated how the KdV equation was modified by the presence of dissipation and the effect of such dissipation on the solitary solution of the KdV equation. Ghidaglia [28] studied the long time behavior of solutions to the weakly dissipative KdV equation as a finite-dimensional dynamical system. Recently, Wu and Yin [46] discussed the blow-up, blow-up rate, and decay of the solution of the weakly dissipative periodic Camassa–Holm equation.

In this paper, we would like to consider the following dissipative Degasperis–Procesi equation:

$$u_t - u_{txx} + 4uu_x + L(u) = 3u_x u_{xx} + uu_{xxx}, \quad t > 0, \quad x \in \mathbb{R},$$

where $L(u)$ is a dissipative term, and L can be a differential operator or a quasi-differential operator according to different physical situations. We are interested in the effect of the weakly dissipative term on the Degasperis–Procesi equation. In particular, we study the following weakly dissipative Degasperis–Procesi equation:

$$(1.3) \quad \begin{cases} y_t + uy_x + 3u_x y + \lambda y = 0, & t > 0, \quad x \in \mathbb{R}, \\ y = u - u_{xx}, & t > 0, \quad x \in \mathbb{R}, \\ u(0, x) = u_0(x), & x \in \mathbb{R}, \end{cases}$$

where $\lambda y = \lambda(1 - \partial_x^2)u$ is the weakly dissipative term, and $\lambda > 0$ is a constant.

We find that the behaviors of (1.3) are similar to the Degasperis–Procesi equation in a finite interval of time, such as the local well-posedness and the blow-up phenomena. But there are considerable differences between (1.3) and the Degasperis–Procesi equation in their long time behaviors. Global solution of (1.3) decays to zero as time goes to infinity provided the potential $y_0 = (1 - \partial_x^2)u_0$ is of one sign. This long time behavior is an important feature that the Degasperis–Procesi equation does not possess. It is known that the Degasperis–Procesi equation has peaked traveling wave solutions [19]. Theorem 4.1 in what follows shows that any global solution decays in the H^1 -norm. This means that there are no traveling wave solutions of (1.3). This is also another considerable difference between (1.3) and the Degasperis–Procesi equation (1.2) in their long time behaviors.

It is very interesting that (1.3) has the same blow-up rate as the Degasperis–Procesi equation does when the blow-up occurs. This fact shows that the blow-up rate of the Degasperis–Procesi equation is not affected by the weakly dissipative term. But the occurrence of blow-up of (1.2) is affected by the dissipative parameter.

It should be noticed that the weakly dissipative term breaks the conservation laws of the following Degasperis–Procesi equation [19]:

$$E_1(u) = \int_{\mathbb{R}} y dx, \quad E_2(u) = \int_{\mathbb{R}} y v dx, \quad E_3(u) = \int_{\mathbb{R}} u^3 dx,$$

where $y = (1 - \partial_x^2)u$ and $v = (4 - \partial_x^2)^{-1}u$, which play an important role in the study of the Degasperis–Procesi equation [5, 24, 25, 35, 48, 49, 50, 51].

Notation. As above and henceforth, we denote by $*$ the convolution. We write \hat{f} as the Fourier transform of f . We also use (\cdot, \cdot) to represent the standard inner product in $L^2(\mathbb{R})$. For $1 \leq p \leq \infty$, the norm in the Lebesgue space L^p will be written $\|\cdot\|_{L^p}$, while $\|\cdot\|_r, r \geq 0$, will stand for the norm in the classical Sobolev spaces $H^r(\mathbb{R})$.

2. Local well-posedness and blow-up. In this section, we establish the local well-posedness of (1.3), derive a precise blow-up scenario, and present two blow-up criteria for strong solutions to (1.3).

Note that if $p(x) = \frac{1}{2}e^{-|x|}$, $x \in \mathbb{R}$, then $(1 - \partial_x^2)^{-1}f = p * f$ for all $f \in L^2(\mathbb{R})$ and $p * y = u$. Using this identity, we can rewrite (1.3) as follows:

$$(2.1) \quad \begin{cases} u_t + uu_x + \partial_x p * (\frac{3}{2}u^2) + \lambda u = 0, & t > 0, \quad x \in \mathbb{R}, \\ u(0, x) = u_0(x), & x \in \mathbb{R}. \end{cases}$$

Using Kato’s semigroup approach along lines very similar to [48], one obtains the local well-posedness of (1.3) or (2.1).

THEOREM 2.1. *Given $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$, there exist a maximal $T = T(\lambda, \|u_0\|_r) > 0$, and a unique solution u to (1.3) (or (2.1)), such that*

$$u = u(\cdot, u_0) \in C([0, T]; H^r(\mathbb{R})) \cap C^1([0, T]; H^{r-1}(\mathbb{R})),$$

and the solution depends continuously on the initial data; i.e., the mapping $u \rightarrow u(\cdot, u_0) : H^r(\mathbb{R}) \rightarrow C([0, T]; H^r(\mathbb{R})) \cap C^1([0, T]; H^{r-1}(\mathbb{R}))$ is continuous. Moreover, T may be chosen independent of r in the following sense. If

$$u = u(\cdot, u_0) \in C([0, T]; H^r(\mathbb{R})) \cap C^1([0, T]; H^{r-1}(\mathbb{R}))$$

to (1.3) (or (2.1)), and if $u_0 \in H^{r'}(\mathbb{R})$ for some $r' \neq r, r' > \frac{3}{2}$, then

$$u \in C([0, T]; H^{r'}(\mathbb{R})) \cap C^1([0, T]; H^{r'-1}(\mathbb{R}))$$

with the same T .

The following results are proved only with regard to $r = 3$, since we can obtain the same conclusion for the general case $r > \frac{3}{2}$ by using Theorem 2.1 and a simple density argument.

We now present a precise blow-up scenario for strong solutions to (1.3).

THEOREM 2.2. *Given $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$, the solution of (1.3) (or (2.1)) blows up in a finite time $T > 0$ if and only if*

$$\liminf_{t \rightarrow T} \{ \inf_{x \in \mathbb{R}} u_x(t, x) \} = -\infty.$$

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) (or (2.1)) with initial data $u_0 \in H^3(\mathbb{R})$. By (1.3), we have

$$(2.2) \quad \begin{aligned} \frac{d}{dt} \int_{\mathbb{R}} y^2 dx &= 2 \int_{\mathbb{R}} yy_t dx \\ &= -2 \int_{\mathbb{R}} uyy_x dx - 6 \int_{\mathbb{R}} u_x y^2 dx - 2\lambda \int_{\mathbb{R}} y^2 dx \\ &= -5 \int_{\mathbb{R}} u_x y^2 dx - 2\lambda \int_{\mathbb{R}} y^2 dx. \end{aligned}$$

If $u_0 \in H^4(\mathbb{R})$, then we can obtain by (1.3)

$$\begin{aligned}
 (2.3) \quad \frac{d}{dt} \int_{\mathbb{R}} y_x^2 dx &= 2 \int_{\mathbb{R}} y_x y_{xt} dx \\
 &= -8 \int_{\mathbb{R}} u_x y_x^2 dx - 2 \int_{\mathbb{R}} u y_x y_{xx} dx - 6 \int_{\mathbb{R}} u_{xx} y y_x dx - 2\lambda \int_{\mathbb{R}} y_x^2 dx \\
 &= -7 \int_{\mathbb{R}} u_x y_x^2 dx + 3 \int_{\mathbb{R}} u_x y^2 dx - 2\lambda \int_{\mathbb{R}} y_x^2 dx.
 \end{aligned}$$

As for $u_0 \in H^3(\mathbb{R})$, we will show that (2.3) still holds. In fact, we can approximate u_0 in $H^3(\mathbb{R})$ by function $u_0^n \in H^4(\mathbb{R})$. Moreover, we write $u^n = u^n(\cdot, u_0^n)$ for the solution of (1.3) with initial data u_0^n .

By Theorem 2.1, we know that

$$u^n \in C([0, T_n]; H^4(\mathbb{R})) \cap C^1([0, T_n]; H^3(\mathbb{R})), \quad n \geq 1,$$

$$y^n = u^n - u_{xx}^n \in C([0, T_n]; H^2(\mathbb{R})) \cap C^1([0, T_n]; H^1(\mathbb{R})), \quad n \geq 1,$$

$u^n \rightarrow u$ in $H^3(\mathbb{R})$, and $T_n \rightarrow T$ as $n \rightarrow \infty$.

Due to $u_0^n \in H^4(\mathbb{R})$, we have by (2.3)

$$\frac{d}{dt} \int_{\mathbb{R}} (y_x^n)^2 dx = -7 \int_{\mathbb{R}} u_x^n (y_x^n)^2 dx + 3 \int_{\mathbb{R}} u_x^n (y^n)^2 dx - 2\lambda \int_{\mathbb{R}} (y_x^n)^2 dx.$$

Since $u^n \rightarrow u$ in $H^3(\mathbb{R})$ as $n \rightarrow \infty$, it follows that $u_x^n \rightarrow u_x$ in $L^\infty(\mathbb{R})$ as $n \rightarrow \infty$. Also note that $y^n \rightarrow y$ in $H^1(\mathbb{R})$ and $y_x^n \rightarrow y_x$ in $L^2(\mathbb{R})$ as $n \rightarrow \infty$. Letting n go to infinity in the above equation, we can easily deduce that (2.3) holds for $u_0 \in H^3(\mathbb{R})$.

Adding (2.2) and (2.3), we get

$$\begin{aligned}
 (2.4) \quad \frac{d}{dt} \left(\int_{\mathbb{R}} y^2 dx + \int_{\mathbb{R}} y_x^2 dx \right) &= -7 \int_{\mathbb{R}} u_x y_x^2 dx - 2 \int_{\mathbb{R}} u_x y^2 dx \\
 &\quad - 2\lambda \left(\int_{\mathbb{R}} y^2 dx + \int_{\mathbb{R}} y_x^2 dx \right).
 \end{aligned}$$

If u_x is bounded from below on $[0, T)$, then there exists a positive constant k such that $u_x \geq -k$ on $[0, T)$. Thus, we get by (2.4) and Gronwall's inequality

$$\|y\|_1^2 \leq \exp\{(7k - 2\lambda)t\} \|y(0)\|_1^2.$$

The above inequality implies that the H^3 -norm of the solution u of (1.3) does not blow up in finite time. \square

Next, we present two blow-up criteria for (1.3) guaranteeing the occurrence of this phenomenon. Let us first prove several useful lemmas.

Consider the following differential equation:

$$(2.5) \quad \begin{cases} q_t = u(t, q), & t \in [0, T), \\ q(0, x) = x, & x \in \mathbb{R}. \end{cases}$$

The system (2.5) is essential in deriving invariance properties for solutions of the Degasperis–Procesi equation [25, 35, 51] and the Camassa–Holm equation [7, 9, 13]. It is thus natural to expect it would also be useful in the present context.

Applying classical results in the theory of ordinary differential equations, one can obtain the following two results on q which are crucial in the proof of global existence and blow-up solutions.

LEMMA 2.1 (see [51]). *Let $u_0 \in H^r(\mathbb{R})$, $r \geq 3$, and let $T > 0$ be the maximal existence time of the corresponding solution u to (1.3). Then (2.5) has a unique solution $q \in C^1([0, T) \times \mathbb{R}, \mathbb{R})$. Moreover, the map $q(t, \cdot)$ is an increasing diffeomorphism of \mathbb{R} with*

$$q_x(t, x) = \exp\left(\int_0^t u_x(s, q(s, x)) ds\right) > 0 \quad \forall (t, x) \in [0, T) \times \mathbb{R}.$$

LEMMA 2.2. *Let $u_0 \in H^r(\mathbb{R})$, $r \geq 3$, and let $T > 0$ be the maximal existence time of the corresponding solution u to (1.3). Setting $y := u - u_{xx}$, we have*

$$y(t, q(t, x))q_x^3(t, x) = y_0(x) \exp\{-\lambda t\} \quad \forall (t, x) \in [0, T) \times \mathbb{R}.$$

Proof. Differentiation of the system (2.5) with respect to x yields

$$\begin{cases} \frac{d}{dt}q_x = u_x(t, q)q_x, & t \in (0, T), \\ q_x(0, x) = 1, & x \in \mathbb{R}. \end{cases}$$

Let $g(t, x) = y(t, q(t, x))q_x^3(t, x)$. By Lemma 2.1, we infer from (1.3) and (2.5) that

$$\frac{d}{dt}g(t, x) = -\lambda g(t, x).$$

Integrating the above relation with respect to $t < T$ on $[0, t]$ yields the desired result. \square

LEMMA 2.3. *If $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$, then as long as the solution $u(t, x)$ given by Theorem 2.1 exists, we have*

$$\int_{\mathbb{R}} y(t, x)v(t, x)dx = \exp\{-2\lambda t\} \int_{\mathbb{R}} y_0(x)v_0(x)dx,$$

where $y(t, x) = u(t, x) - u_{xx}(t, x)$ and $v(t, x) = (4 - \partial_x^2)^{-1}u$. Moreover, we have

$$\|u(t)\|_{L^2}^2 \leq 4 \exp\{-2\lambda t\} \|u_0\|_{L^2}^2.$$

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) (or (2.1)) with initial data $u_0 \in H^3(\mathbb{R})$. By (1.3), we have

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\mathbb{R}} yv dx &= \frac{1}{2} \int_{\mathbb{R}} y_t v dx + \frac{1}{2} \int_{\mathbb{R}} yv_t dx = \int_{\mathbb{R}} y_t v dx \\ &= - \int_{\mathbb{R}} vy_x u dx - 3 \int_{\mathbb{R}} vy u_x dx - \lambda \int_{\mathbb{R}} yv dx \\ &= - \int_{\mathbb{R}} v(yu)_x dx - 2 \int_{\mathbb{R}} vy u_x dx - \lambda \int_{\mathbb{R}} yv dx. \end{aligned}$$

Using the relations $y = u - u_{xx}$ and $4v - v_{xx} = u$, it yields

$$\begin{aligned}
\int_{\mathbb{R}} v(yu)_x dx &= - \int_{\mathbb{R}} v_x y u dx = - \int_{\mathbb{R}} v_x u^2 dx + \int_{\mathbb{R}} v_x u u_{xx} dx \\
&= - \int_{\mathbb{R}} v_x u^2 dx - \int_{\mathbb{R}} (v_x u)_x u_x dx \\
&= - \int_{\mathbb{R}} v_x u^2 dx - \int_{\mathbb{R}} v_{xx} u u_x dx - \int_{\mathbb{R}} v_x u_x^2 dx \\
&= - \int_{\mathbb{R}} v_x u^2 dx + \frac{1}{2} \int_{\mathbb{R}} v_{xxx} u^2 dx - \int_{\mathbb{R}} v_x u_x^2 dx \\
&= - \int_{\mathbb{R}} v_x u^2 dx + \frac{1}{2} \int_{\mathbb{R}} (4v_x - u_x) u^2 dx - \int_{\mathbb{R}} v_x u_x^2 dx \\
&= \int_{\mathbb{R}} v_x u^2 dx - \int_{\mathbb{R}} v_x u_x^2 dx.
\end{aligned}$$

On the other hand,

$$\begin{aligned}
2 \int_{\mathbb{R}} v y u_x dx &= 2 \int_{\mathbb{R}} v u u_x dx - 2 \int_{\mathbb{R}} v u_{xx} u_x dx \\
&= - \int_{\mathbb{R}} v_x u^2 dx + \int_{\mathbb{R}} v_x u_x^2 dx.
\end{aligned}$$

Combining the above three relations, we deduce that

$$\frac{1}{2} \frac{d}{dt} \int_{\mathbb{R}} y v dx = -\lambda \int_{\mathbb{R}} y v dx.$$

Consequently, this implies the first desired result. In view of the proved equality, it then follows that

$$\begin{aligned}
\|u(t)\|_{L^2}^2 &= \|\hat{u}(t)\|_{L^2}^2 \leq 4 \int_{\mathbb{R}} \frac{1 + \xi^2}{4 + \xi^2} |\hat{u}(t, \xi)|^2 d\xi \\
&= 4(\hat{y}(t), \hat{v}(t)) = 4(y(t), v(t)) \\
&= 4 \exp\{-2\lambda t\} (y_0, v_0) = 4 \exp\{-2\lambda t\} (\hat{y}_0, \hat{v}_0) \\
&\leq 4 \exp\{-2\lambda t\} \int_{\mathbb{R}} \frac{1 + \xi^2}{4 + \xi^2} |\hat{u}_0(\xi)|^2 d\xi \\
&\leq 4 \exp\{-2\lambda t\} \|\hat{u}_0\|_{L^2}^2 = 4 \exp\{-2\lambda t\} \|u_0\|_{L^2}^2.
\end{aligned}$$

This completes the proof of the lemma. \square

LEMMA 2.4. Assume $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$. Let T be the maximal existence time of the solution u to (1.3) guaranteed by Theorem 2.1. Then we have

$$\|u(t)\|_{L^\infty} \leq \exp\{-\lambda t\} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right) \quad \forall t \in [0, T].$$

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) with the initial data $u_0 \in H^3(\mathbb{R})$. By (2.5), we get

$$\begin{aligned} \frac{du(t, q(t, x))}{dt} &= u_t(t, q(t, x)) + u_x(t, q(t, x)) \frac{dq(t, x)}{dt} \\ &= (u_t + uu_x)(t, q(t, x)). \end{aligned}$$

By (2.1), we have

$$u_t + uu_x = -3p * (uu_x) - \lambda u.$$

Note that

$$\begin{aligned} -3p * (uu_x) &= -\frac{3}{2} \int_{-\infty}^{+\infty} e^{-|x-\eta|} uu_\eta d\eta \\ &= -\frac{3}{2} \int_{-\infty}^x e^{-x+\eta} uu_\eta d\eta - \frac{3}{2} \int_x^{+\infty} e^{x-\eta} uu_\eta d\eta \\ &= \frac{3}{4} \int_{-\infty}^x e^{-|x-\eta|} u^2 d\eta - \frac{3}{4} \int_x^{+\infty} e^{-|x-\eta|} u^2 d\eta. \end{aligned}$$

It then follows that

$$-\frac{3}{4} \int_{q(t,x)}^{+\infty} e^{-|q(t,x)-\eta|} u^2 d\eta \leq \frac{du(t, q(t, x))}{dt} + \lambda u(t, q(t, x)) \leq \frac{3}{4} \int_{-\infty}^{q(t,x)} e^{-|q(t,x)-\eta|} u^2 d\eta.$$

It thus transpires that

$$\left| \frac{du(t, q(t, x))}{dt} + \lambda u(t, q(t, x)) \right| \leq \frac{3}{4} \int_{-\infty}^{+\infty} e^{-|q(t,x)-\eta|} u^2 d\eta \leq \frac{3}{4} \|u(t)\|_{L^2}^2.$$

In view of Lemma 2.3, we have

$$\left| \frac{du(t, q(t, x))}{dt} + \lambda u(t, q(t, x)) \right| \leq 3 \exp\{-2\lambda t\} \|u_0\|_{L^2}^2.$$

Integrating the above inequality with respect to $t < T$ on $[0, t]$ yields

$$|\exp\{\lambda t\} u(t, q(t, x)) - u_0(x)| \leq \frac{3}{\lambda} \|u_0\|_{L^2}^2.$$

Thus,

$$\begin{aligned} (2.6) \quad |u(t, q(t, x))| &\leq \|u(t, q(t, x))\|_{L^\infty} \\ &\leq \exp\{-\lambda t\} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right). \end{aligned}$$

Using the Sobolev embedding to ensure the uniform boundedness of $u_x(s, \eta)$ for $(s, \eta) \in [0, t] \times \mathbb{R}$ with $t \in [0, T]$, in view of Lemma 2.1, we get for every $t \in [0, T]$ a constant $C(t) > 0$ such that

$$e^{-C(t)} \leq q_x(t, x) \leq e^{C(t)}, \quad x \in \mathbb{R}.$$

We deduce from the above equation that the function $q(t, \cdot)$ is strictly increasing on \mathbb{R} with $\lim_{x \rightarrow \pm\infty} q(t, x) = \pm\infty$ as long as $t \in [0, T)$. Thus, by (2.6) we can obtain

$$\|u(t, x)\|_{L^\infty} = \|u(t, q(t, x))\|_{L^\infty} \leq \exp\{-\lambda t\} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right).$$

This completes the proof of Lemma 2.4. \square

We are now in the position to present the first blow-up result.

THEOREM 2.3. *Let $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$, and assume that there exists $x_0 \in \mathbb{R}$ such that*

$$u_0'(x_0) < -\frac{1}{2}\lambda - \frac{1}{2}\sqrt{\lambda^2 + 6 \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right)^2}.$$

Then the corresponding solution of (1.3) blows up in finite time.

Proof. Let $T > 0$ be the existence time of the solution $u(t, \cdot)$ of (1.3) (or (2.1)) with the initial data $u_0 \in H^3(\mathbb{R})$. Differentiating (2.1) with respect to x , in view of $\partial_x^2 p * f = p * f - f$, we get

$$(2.7) \quad u_{tx} = -u_x^2 - uu_{xx} + \frac{3}{2}u^2 - p * \left(\frac{3}{2}u^2 \right) - \lambda u_x.$$

Note that

$$\begin{aligned} \frac{du_x(t, q(t, x))}{dt} &= u_{xt}(t, q(t, x)) + u_{xx}(t, q(t, x)) \frac{dq(t, x)}{dt} \\ &= u_{tx}(t, q(t, x)) + u(t, q(t, x))u_{xx}(t, q(t, x)). \end{aligned}$$

Thus, we have

$$(2.8) \quad \begin{aligned} \frac{du_x(t, q(t, x))}{dt} &= -u_x^2(t, q(t, x)) + \frac{3}{2}u^2(t, q(t, x)) \\ &\quad - p * \left(\frac{3}{2}u^2(t, q(t, x)) \right) - \lambda u_x(t, q(t, x)). \end{aligned}$$

In view of $p * \left(\frac{3}{2}u^2 \right) (t, q(t, x)) \geq 0$, we infer from (2.8) and Lemma 2.4 that

$$\begin{aligned} \frac{du_x(t, q(t, x))}{dt} &\leq -u_x^2(t, q(t, x)) - \lambda u_x(t, q(t, x)) + \frac{3}{2}u^2(t, q(t, x)) \\ &\leq -u_x^2(t, q(t, x)) - \lambda u_x(t, q(t, x)) + \frac{3}{2} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right)^2. \end{aligned}$$

Set $m(t) = u_x(t, q(t, x_0))$ and

$$\alpha^2 = \frac{3}{2} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right)^2.$$

Then we obtain

$$\begin{aligned} \frac{dm(t)}{dt} &\leq -m^2(t) - \lambda m(t) + \alpha^2 \\ &= -\frac{1}{4} \left(2m(t) + \lambda - \sqrt{\lambda^2 + 4\alpha^2} \right) \left(2m(t) + \lambda + \sqrt{\lambda^2 + 4\alpha^2} \right). \end{aligned}$$

From the hypothesis, we have $m(0) < -\frac{1}{2}\lambda - \frac{1}{2}\sqrt{\lambda^2 + 4\alpha^2}$, and thus $\frac{dm}{dt}|_{t=0} < 0$. By continuity with respect to t of $m(t)$, we have $\frac{dm}{dt} < 0$ for all $t \in [0, T)$. Therefore, $m(t) < -\frac{1}{2}\lambda - \frac{1}{2}\sqrt{\lambda^2 + 4\alpha^2}$ for all $t \in [0, T)$. Thus, we can solve the above inequality to obtain

$$\frac{2m(0) + \lambda + \sqrt{\lambda^2 + 4\alpha^2}}{2m(0) + \lambda - \sqrt{\lambda^2 + 4\alpha^2}} \exp\{\sqrt{\lambda^2 + 4\alpha^2} t\} - 1 \leq \frac{2\sqrt{\lambda^2 + 4\alpha^2}}{2m(t) + \lambda - \sqrt{\lambda^2 + 4\alpha^2}} \leq 0.$$

Since

$$0 < \frac{2m(0) + \lambda + \sqrt{\lambda^2 + 4\alpha^2}}{2m(0) + \lambda - \sqrt{\lambda^2 + 4\alpha^2}} < 1,$$

there exists T satisfying

$$T \leq \frac{1}{\sqrt{\lambda^2 + 4\alpha^2}} \ln \left(\frac{2m(0) + \lambda - \sqrt{\lambda^2 + 4\alpha^2}}{2m(0) + \lambda + \sqrt{\lambda^2 + 4\alpha^2}} \right)$$

such that $\lim_{t \uparrow T} m(t) = -\infty$. Hence, the above theorem is proved according to Theorem 2.2. \square

We give the second criterion that guarantees the blow-up of the solutions of (1.3).

THEOREM 2.4. *Assume that $u_0 \in H^r(\mathbb{R})$ ($r > \frac{3}{2}$) is odd, and $u'_0(0) < -\lambda$. Then the corresponding solution of (1.3) blows up in finite time.*

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) (or (2.1)) with initial data $u_0 \in H^3(\mathbb{R})$. As one can check, the function

$$v(t, x) := -u(t, -x), \quad t \in [0, T), \quad x \in \mathbb{R},$$

is also a solution of (1.3) in $C([0, T); H^3(\mathbb{R})) \cap C^1([0, T); H^2(\mathbb{R}))$ with initial data u_0 . By uniqueness we conclude that $v \equiv u$, and therefore $u(t, \cdot)$ is odd for any $t \in [0, T)$. In particular, by continuity with respect to the spatial variable of u and u_{xx} , we get

$$u(t, 0) = u_{xx}(t, 0) = 0, \quad t \in [0, T).$$

Define $h(t) := u_x(t, 0)$ for $t \in [0, T)$. Note that $h \in C^1([0, T), \mathbb{R})$. From (2.7), we get

$$\begin{aligned} \frac{dh}{dt}(t) &= -h^2(t) - \lambda h(t) - p * \left(\frac{3}{2} u^2 \right) \\ &\leq -(h(t) + \lambda) h(t), \quad t \in [0, T). \end{aligned}$$

From the hypothesis, we have $h(0) < -\lambda$. Therefore, $h(t) < -\lambda$ for all $t \in [0, T)$. Solving the above inequality, we get

$$1 - \frac{h(0)}{h(0) + \lambda} \exp\{-\lambda t\} \leq \frac{\lambda}{h(t) + \lambda} \leq 0.$$

Since

$$\frac{h(0)}{h(0) + \lambda} > 1,$$

we conclude that there exists T and

$$T \leq \frac{1}{\lambda} \ln \frac{h(0)}{h(0) + \lambda}$$

such that $\lim_{t \uparrow T} h(t) = -\infty$. We complete the proof of the theorem by Theorem 2.2. \square

3. Blow-up rate. In this section, we give more insight into the blow-up mechanism for the wave-breaking solutions to (1.3).

LEMMA 3.1 (see [10]). *Let $T > 0$ and $v \in C^1([0, T]; H^2(\mathbb{R}))$; then for every $t \in [0, T)$ there exists at least one point $\xi(t) \in \mathbb{R}$ with*

$$m(t) := \inf_{x \in \mathbb{R}} v_x(t, x) = v_x(t, \xi(t)).$$

The function $m(t)$ is a.e. differentiable on $(0, T)$ with

$$\frac{dm}{dt} = v_{tx}(t, \xi(t)) \text{ a.e. on } (0, T).$$

THEOREM 3.1. *Let $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$, and let $T > 0$ be the maximal existence time of the corresponding solution to (1.3). If T is finite, we have*

$$\lim_{t \rightarrow T} (T - t) \inf_{x \in \mathbb{R}} u_x(t, x) = -1.$$

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) (or (2.1)) with initial data $u_0 \in H^3(\mathbb{R})$. We already know by Theorem 2.2 that

$$(3.1) \quad \lim_{t \rightarrow T} \inf \{ \inf_{x \in \mathbb{R}} u_x(t, x) \} = -\infty.$$

We know that $u(x, t) \in C^1([0, T]; H^2(\mathbb{R}))$ by Theorem 2.1; then we infer from Lemma 3.1 that, for every $t \in [0, T)$, there exists at least one point $\xi(t) \in \mathbb{R}$ with $u_x(t, \xi(t)) = \inf_{x \in \mathbb{R}} u_x(t, x)$. Let

$$m(t) = u_x(t, \xi(t)) = \inf_{x \in \mathbb{R}} u_x(t, x);$$

then $u_{xx}(t, \xi(t)) = 0$ for all $t \in [0, T)$. From (2.7) we have

$$(3.2) \quad \frac{dm}{dt} + m^2(t) + \lambda m(t) = \frac{3}{2}u^2(t, \xi(t)) - p * \left(\frac{3}{2}u^2(t, \xi(t)) \right) \text{ a.e. on } (0, T).$$

By Lemmas 2.3–2.4, we have for $t \in [0, T)$ that

$$\begin{aligned} & \left| \frac{3}{2}u^2(t, \xi(t)) - p * \left(\frac{3}{2}u^2(t, \xi(t)) \right) \right| \\ & \leq \left| \frac{3}{2}u^2(t, \xi(t)) \right| + \left| \frac{3}{4} \int_{\mathbb{R}} e^{-|\xi(t)-\eta|} u^2(t, \eta) d\eta \right| \\ & \leq \frac{3}{2} \|u\|_{L^\infty}^2 + \frac{3}{4} \|u\|_{L^2}^2 \\ & \leq \frac{3}{2} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right)^2 + 3 \|u_0\|_{L^2}^2. \end{aligned}$$

Set

$$\beta = \frac{3}{2} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right)^2 + 3 \|u_0\|_{L^2}^2.$$

We infer from (3.2) that

$$\left| \frac{dm}{dt} + m^2(t) + \lambda m(t) \right| \leq \beta \text{ a.e. on } (0, T).$$

Hence,

$$(3.3) \quad -\beta - \frac{1}{4}\lambda^2 \leq \frac{dm}{dt} + \left(m(t) + \frac{1}{2}\lambda \right)^2 \leq \beta + \frac{1}{4}\lambda^2 \text{ a.e. on } (0, T).$$

Let $\varepsilon \in (0, 1)$. Since $\lim_{t \rightarrow T} \inf(m(t) + \frac{1}{2}\lambda) = -\infty$ (by (3.1)), there is some $t_0 \in (0, T)$ with $m(t_0) + \frac{1}{2}\lambda < 0$ and

$$\left(m(t_0) + \frac{1}{2}\lambda \right)^2 > \frac{1}{\varepsilon} \left(\beta + \frac{1}{4}\lambda^2 \right).$$

We claim that

$$(3.4) \quad \left(m(t) + \frac{1}{2}\lambda \right)^2 > \frac{1}{\varepsilon} \left(\beta + \frac{1}{4}\lambda^2 \right), \quad t \in [t_0, T).$$

In fact, since $m(t)$ is locally Lipschitz (it belongs to $W_{loc}^{1,\infty}(\mathbb{R})$ by Theorem 2.1), there is some $\delta > 0$ such that

$$\left(m(t) + \frac{1}{2}\lambda \right)^2 > \frac{1}{\varepsilon} \left(\beta + \frac{1}{4}\lambda^2 \right), \quad t \in (t_0, t_0 + \delta).$$

From (3.3), we have

$$\frac{dm}{dt} < (\varepsilon - 1) \left(m(t) + \frac{1}{2}\lambda \right)^2 < 0, \quad t \in (t_0, t_0 + \delta) \text{ a.e. on } (0, T).$$

Being locally Lipschitz, the function $m(t)$ is absolutely continuous. Therefore, by integrating the above relation on $[t_0, t_0 + \delta]$, we obtain that $m(t_0 + \delta) \leq m(t_0)$. Thus,

$$m(t_0 + \delta) + \frac{1}{2}\lambda \leq m(t_0) + \frac{1}{2}\lambda < 0.$$

By the above inequality, we have

$$\left(m(t_0 + \delta) + \frac{1}{2}\lambda \right)^2 \geq \left(m(t_0) + \frac{1}{2}\lambda \right)^2 > \frac{1}{\varepsilon} \left(\beta + \frac{1}{4}\lambda^2 \right).$$

The relation (3.4) is proved by a continuous extension.

A combination of (3.3) with (3.4) enables us to infer that

$$(3.5) \quad -1 - \varepsilon < \frac{\frac{dm}{dt}}{\left(m(t) + \frac{1}{2}\lambda \right)^2} < -1 + \varepsilon \text{ a.e. on } (t_0, T).$$

For $t \in (t_0, T)$, integrating (3.5) on (t, T) , we obtain

$$-1 - \varepsilon < \frac{1}{\left(m(t) + \frac{1}{2}\lambda \right) (T - t)} < -1 + \varepsilon, \quad t \in (t_0, T).$$

Letting ε go to zero, we obtain

$$\lim_{t \rightarrow T} \left[\left(m(t) + \frac{1}{2} \lambda \right) (T - t) \right] = -1,$$

that is,

$$\lim_{t \rightarrow T} (T - t)m(t) = -1.$$

This completes the proof of the theorem. \square

Remark 3.1. Although the occurrence of blow-up of strong solutions to (1.3) is affected by the dissipative parameter (see Theorems 2.3–2.4), Theorem 3.1 shows that the blow-up rate of strong solutions to the Degasperis–Procesi equation [24] is not affected by the weakly dissipative term.

4. Global solution and its decay. In this section we will show that there exist global strong solutions to (1.3) and these global solutions decay to zero as time goes to infinity provided the initial data u_0 satisfy certain sign conditions.

THEOREM 4.1. *Assume $u_0 \in H^r(\mathbb{R})$, $r > \frac{3}{2}$. If $y_0 = u_0 - u_{0,xx}$ does not change sign on \mathbb{R} , then (1.3) (or (2.1)) has a global strong solution*

$$u = u(\cdot, u_0) \in C([0, \infty); H^r(\mathbb{R})) \cap C^1([0, \infty); H^{r-1}(\mathbb{R})).$$

Moreover, the global solution decays to 0 in the H^1 -norm and the H^3 -norm as time goes to infinity.

Proof. Let $T > 0$ be the maximal time of existence of the solution u to (1.3) (or (2.1)) with initial data $u_0 \in H^3(\mathbb{R})$. We first consider the case where $y_0 \geq 0$ on \mathbb{R} . If $y_0 \geq 0$, then Lemmas 2.1–2.2 ensure that $y \geq 0$ for all $t \in [0, T)$. Using $u = p * y$ and the positivity of p , we infer that $u(t, \cdot) \geq 0$ for all $t \geq 0$. Note that

$$u(t, x) = \frac{e^{-x}}{2} \int_{-\infty}^x e^\eta y(t, \eta) d\eta + \frac{e^x}{2} \int_x^\infty e^{-\eta} y(t, \eta) d\eta$$

and

$$u_x(t, x) = -\frac{e^{-x}}{2} \int_{-\infty}^x e^\eta y(\eta) d\eta + \frac{e^x}{2} \int_x^\infty e^{-\eta} y(\eta) d\eta.$$

From the above two equations, we deduce that

$$\begin{aligned} (4.1) \quad u(t, x) + u_x(t, x) &= e^x \int_x^\infty e^{-\eta} y(t, \eta) d\eta, \\ u(t, x) - u_x(t, x) &= e^{-x} \int_{-\infty}^x e^\eta y(t, \eta) d\eta. \end{aligned}$$

By (4.1) and $y \geq 0$ for all $t \in [0, T)$, we obtain, for $t \in [0, T)$,

$$(4.2) \quad |u_x(t, x)| \leq u(t, x) \quad \forall (t, x) \in [0, T) \times \mathbb{R}.$$

By Lemma 2.4, we have

$$(4.3) \quad \begin{aligned} |u_x(t, x)| &\leq u(t, x) \\ &\leq \exp\{-\lambda t\} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right) \quad \forall (t, x) \in [0, T) \times \mathbb{R}. \end{aligned}$$

The above inequality and Theorem 2.2 imply that $T = \infty$. This proves that the solution u exists globally in time.

Multiplying (1.3) by u and integrating by parts, in view of (4.2) and Lemmas 2.3–2.4, we get

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\mathbb{R}} (u^2 + u_x^2) dx + \lambda \int_{\mathbb{R}} (u^2 + u_x^2) dx \\ &= -4 \int_{\mathbb{R}} u^2 u_x dx + 3 \int_{\mathbb{R}} uu_x u_{xx} dx + \int_{\mathbb{R}} u^2 u_{xxx} dx \\ &= -\frac{1}{2} \int_{\mathbb{R}} u_x^3 dx \leq \frac{1}{2} \int_{\mathbb{R}} u^3 dx \leq \frac{1}{2} \|u\|_{L^\infty} \int_{\mathbb{R}} u^2 dx \\ &\leq 2 \exp\{-3\lambda t\} \left(\frac{3}{\lambda} \|u_0\|_{L^2}^2 + \|u_0\|_{L^\infty} \right) \|u_0\|_{L^2}^2. \end{aligned}$$

Integrating the above inequality with respect to t , we have

$$\|u(t)\|_1^2 \leq \exp\{-2\lambda t\} \left(\frac{12}{\lambda^2} \|u_0\|_{L^2}^4 + \frac{4}{\lambda} \|u_0\|_{L^2}^2 \|u_0\|_{L^\infty} + \|u_0\|_1^2 \right).$$

This shows that the corresponding global solution with $y_0 \geq 0$ decays to 0 in the H^1 -norm.

By (4.3), we obtain that $-u_x \leq \frac{\lambda}{7}$ for sufficiently large t . This yields, in combination with (2.4) and Gronwall's inequality,

$$\|y\|_1^2 \leq e^{-\lambda t} \|y_0\|_1^2$$

for large t . The above inequality implies that the corresponding global solution with $y_0 \geq 0$ decays to 0 in the H^3 -norm. This completes the proof of the theorem with the assumption $y_0 \geq 0$ on \mathbb{R} . In the case when $y_0(x) \leq 0$ on \mathbb{R} , one can repeat the above proof to get the desired result. \square

Remark 4.1. Note that the global solution to the Degasperis–Procesi equation does not generally decay to zero as time goes to infinity [35]. Theorem 4.1 shows that there is a considerable difference between (1.3) and the Degasperis–Procesi equation in their long time behaviors. More precisely, the energy dissipation will affect the long time behavior of global solutions to the Degasperis–Procesi equation.

Remark 4.2. It is well known that the Degasperis–Procesi equation has peaked traveling wave solutions. Theorem 4.1 shows that global H^3 -solutions with y_0 of one sign decay in the H^1 -norm and the H^3 -norm. Lemma 2.4 shows that any global solution decays in the L^∞ -norm. This means that there are no traveling wave solutions of the dissipative equation (1.3). This is also another considerable difference between (1.3) and the Degasperis–Procesi equation in their long time behaviors.

Acknowledgment. The authors thank the referees for their valuable comments and suggestions.

REFERENCES

- [1] R. BEALS, D. SATTINGER, AND J. SZMIGIELSKI, *Acoustic scattering and the extended Korteweg–de Vries hierarchy*, Adv. Math., 140 (1998), pp. 190–206.

- [2] A. BRESSAN AND A. CONSTANTIN, *Global conservative solutions of the Camassa-Holm equation*, Arch. Ration. Mech. Anal., 183 (2007), pp. 215–239.
- [3] R. CAMASSA AND D. HOLM, *An integrable shallow water equation with peaked solitons*, Phys. Rev. Lett., 71 (1993), pp. 1661–1664.
- [4] R. CAMASSA, D. HOLM, AND J. HYMAN, *A new integrable shallow water equation*, Adv. in Appl. Mech., 31 (1994), pp. 1–33.
- [5] G. M. COCLITE AND K. H. KARLSEN, *On the well-posedness of the Degasperis-Procesi equation*, J. Funct. Anal., 233 (2006), pp. 60–91.
- [6] G. M. COCLITE, K. H. KARLSEN, AND N. H. RISEBRO, *Numerical schemes for computing discontinuous solutions of the Degasperis-Procesi equation*, IMA J. Numer. Anal., 28 (2008), pp. 80–105.
- [7] A. CONSTANTIN, *Global existence of solutions and breaking waves for a shallow water equation: A geometric approach*, Ann. Inst. Fourier (Grenoble), 50 (2000), pp. 321–362.
- [8] A. CONSTANTIN, *On the scattering problem for the Camassa-Holm equation*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 457 (2001), pp. 953–970.
- [9] A. CONSTANTIN AND J. ESCHER, *Global existence and blow-up for a shallow water equation*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 26 (1998), pp. 303–328.
- [10] A. CONSTANTIN AND J. ESCHER, *Wave breaking for nonlinear nonlocal shallow water equations*, Acta Math., 181 (1998), pp. 229–243.
- [11] A. CONSTANTIN AND J. ESCHER, *Global weak solutions for a shallow water equation*, Indiana Univ. Math. J., 47 (1998), pp. 1527–1545.
- [12] A. CONSTANTIN AND B. KOLEV, *Geodesic flow on the diffeomorphism group of the circle*, Comment. Math. Helv., 78 (2003), pp. 787–804.
- [13] A. CONSTANTIN AND H. P. MCKEAN, *A shallow water equation on the circle*, Comm. Pure Appl. Math., 52 (1999), pp. 949–982.
- [14] A. CONSTANTIN AND L. MOLINET, *Global weak solutions for a shallow water equation*, Comm. Math. Phys., 211 (2000), pp. 45–61.
- [15] A. CONSTANTIN AND W. STRAUSS, *Stability of a class of solitary waves in compressible elastic rods*, Phys. Lett. A, 270 (2000), pp. 140–148.
- [16] A. CONSTANTIN AND W. A. STRAUSS, *Stability of peakons*, Comm. Pure Appl. Math., 53 (2000), pp. 603–610.
- [17] A. CONSTANTIN AND W. A. STRAUSS, *Stability of the Camassa-Holm solitons*, J. Nonlinear Sci., 12 (2002), pp. 415–422.
- [18] H. H. DAI, *Model equations for nonlinear dispersive waves in a compressible Mooney-Rivlin rod*, Acta Mech., 127 (1998), pp. 193–207.
- [19] A. DEGASPERIS, D. D. HOLM, AND A. N. W. HONE, *A new integral equation with peakon solutions*, Theoret. and Math. Phys., 133 (2002), pp. 1463–1474.
- [20] A. DEGASPERIS AND M. PROCESI, *Asymptotic integrability*, in Symmetry and Perturbation Theory, A. Degasperis and G. Gaeta, eds., World Scientific, River Edge, NJ, 1999, pp. 23–37.
- [21] P. G. DRAZIN AND R. S. JOHNSON, *Solitons: An Introduction*, Cambridge University Press, Cambridge, UK, New York, 1989.
- [22] H. R. DULLIN, G. A. GOTTFWALD, AND D. D. HOLM, *An integrable shallow water equation with linear and nonlinear dispersion*, Phys. Rev. Lett., 87 (2001), pp. 4501–4504.
- [23] H. R. DULLIN, G. A. GOTTFWALD, AND D. D. HOLM, *Camassa-Holm, Korteweg-de Vries-5 and other asymptotically equivalent equations for shallow water waves*, Fluid Dynam. Res., 33 (2003), pp. 73–95.
- [24] J. ESCHER, Y. LIU, AND Z. YIN, *Global weak solutions and blow-up structure for the Degasperis-Procesi equation*, J. Funct. Anal., 241 (2006), pp. 457–485.
- [25] J. ESCHER, Y. LIU, AND Z. YIN, *Shock waves and blow-up phenomena for the periodic Degasperis-Procesi equation*, Indiana Univ. Math. J., 56 (2007), pp. 87–117.
- [26] J. ESCHER AND Z. YIN, *On the initial boundary value problems for the Degasperis-Procesi equation*, Phys. Lett. A, 368 (2007), pp. 69–76.
- [27] B. FUCHSSTEINER AND A. FOKAS, *Symplectic structures, their Bäcklund transformation and hereditary symmetries*, Phys. D, 4 (1981), pp. 47–66.
- [28] J. M. GHIDAGLIA, *Weakly damped forced Korteweg-de Vries equations behave as a finite dimensional dynamical system in the long time*, J. Differential Equations, 74 (1988), pp. 369–390.
- [29] R. I. IVANOV, *Water waves and integrability*, Philos. Trans. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci., 365 (2007), pp. 2267–2280.
- [30] R. S. JOHNSON, *Camassa-Holm, Korteweg-de Vries and related models for water waves*, J. Fluid Mech., 455 (2002), pp. 63–82.
- [31] C. KENIG, G. PONCE, AND L. VEGA, *Well-posedness and scattering results for the generalized*

- Korteweg-de Vries equation via the contraction principle*, Comm. Pure Appl. Math., 46 (1993), pp. 527–620.
- [32] J. LENELLS, *Traveling wave solutions of the Degasperis-Procesi equation*, J. Math. Anal. Appl., 306 (2005), pp. 72–82.
- [33] J. LENELLS, *Conservation laws of the Camassa-Holm equation*, J. Phys. A, 38 (2005), pp. 869–880.
- [34] Y. LI AND P. OLVER, *Well-posedness and blow-up solutions for an integrable nonlinearly dispersive model wave equation*, J. Differential Equations, 162 (2000), pp. 27–63.
- [35] Y. LIU AND Z. YIN, *Global existence and blow-up phenomena for the Degasperis-Procesi equation*, Comm. Math. Phys., 267 (2006), pp. 801–820.
- [36] H. LUNDMARK, *Formation and dynamics of shock waves in the Degasperis-Procesi equation*, J. Nonlinear Sci., 17 (2007), pp. 169–198.
- [37] H. LUNDMARK AND J. SZMIGIELSKI, *Multi-peakon solutions of the Degasperis-Procesi equation*, Inverse Problems, 19 (2003), pp. 1241–1245.
- [38] Y. MATSUNO, *Multisoliton solutions of the Degasperis-Procesi equation and their peakon limit*, Inverse Problems, 21 (2005), pp. 1553–1570.
- [39] H. P. MCKEAN, *Integrable systems and algebraic curves*, in Global Analysis, Lecture Notes in Math. 755, Springer, Berlin, 1979, pp. 83–200.
- [40] G. MISIOLEK, *A shallow water equation as a geodesic flow on the Bott-Virasoro group*, J. Geom. Phys., 24 (1998), pp. 203–208.
- [41] E. OTT AND R. N. SUDAN, *Damping of solitary waves*, Phys. Fluids, 13 (1970), pp. 1432–1434.
- [42] G. RODRIGUEZ-BLANCO, *On the Cauchy problem for the Camassa-Holm equation*, Nonlinear Anal., 46 (2001), pp. 309–327.
- [43] T. TAO, *Low-regularity global solutions to nonlinear dispersive equations*, in Surveys in Analysis and Operator Theory (Canberra, 2001), Proc. Centre Math. Appl. Austral. Nat. Univ. 40, Australian National University, Canberra, Australia, 2002, pp. 19–48.
- [44] V. O. VAKHNENKO AND E. J. PARKES, *Periodic and solitary-wave solutions of the Degasperis-Procesi equation*, Chaos Solitons Fractals, 20 (2004), pp. 1059–1073.
- [45] G. B. WHITHAM, *Linear and Nonlinear Waves*, John Wiley and Sons, New York, 1980.
- [46] S. WU AND Z. YIN, *Blow-up, blow-up rate and decay of the solution of the weakly dissipative Camassa-Holm equation*, J. Math. Phys., 47 (2006), 013504.
- [47] Z. XIN AND P. ZHANG, *On the weak solutions to a shallow water equation*, Comm. Pure Appl. Math., 53 (2000), pp. 1411–1433.
- [48] Z. YIN, *On the Cauchy problem for an integrable equation with peakon solutions*, Illinois J. Math., 47 (2003), pp. 649–666.
- [49] Z. YIN, *Global existence for a new periodic integrable equation*, J. Math. Anal. Appl., 283 (2003), pp. 129–139.
- [50] Z. YIN, *Global weak solutions to a new periodic integrable equation with peakon solutions*, J. Funct. Anal., 212 (2004), pp. 182–194.
- [51] Z. YIN, *Global solutions to a new integrable equation with peakons*, Indiana Univ. Math. J., 53 (2004), pp. 1189–1210.

NONLINEAR STABILITY OF STATIONARY SOLUTIONS FOR SURFACE DIFFUSION WITH BOUNDARY CONDITIONS*

HARALD GARCKE[†], KAZUO ITO[‡], AND YOSHIHITO KOHSAKA[§]

Abstract. The volume-preserving fourth order surface diffusion flow has constant mean curvature hypersurfaces as stationary solutions. We show nonlinear stability of certain stationary curves in the plane which meet an exterior boundary with a prescribed contact angle. Methods include semigroup theory, energy arguments, geometric analysis, and variational calculus.

Key words. surface diffusion, nonlinear stability, energy method, variational calculus

AMS subject classifications. 35B35, 35G30, 35K55, 35R35, 53C44

DOI. 10.1137/070694752

1. Introduction.

The surface diffusion flow

$$(1.1) \quad V = -\Delta_S \kappa$$

is a geometrical evolution law which describes the surface dynamics for phase interfaces, when mass diffusion occurs only within the interface. Here V is the normal velocity of the evolving surface, Δ_S is the surface Laplacian, and κ is the mean curvature of the surface. The flow (1.1) was first proposed by Mullins [20] in works concerned with thermal grooving. A derivation of (1.1) within rational thermodynamics was given by Davi and Gurtin [7]. In [22], Cahn and Taylor showed that (1.1) is the H^{-1} -gradient flow of the area functional, and in [5], Cahn, Elliott, and Novick-Cohen used formal asymptotics to derive (1.1) as the sharp interface limit of the Cahn–Hilliard equation with degenerate mobility. Further, the motion given by (1.1) has the significant geometrical properties that for closed embedded hypersurfaces the enclosed volume is preserved and surface area decreases in time (see, e.g., [10, 12]). The evolution law (1.1) leads to a fourth order parabolic equation which is in contrast to the second order mean curvature flow $V = \kappa$. We remark that the mean curvature flow is also area decreasing but changes the enclosed volume.

In this paper we study the motion by surface diffusion for curves in cases where the interface intersects an external boundary. More precisely, we consider the following problem. Let Ω be an open bounded domain in \mathbb{R}^2 . We look for evolving curves $\Gamma = \{\Gamma_t\}_{t \geq 0}$ lying in Ω with $\partial\Gamma \subset \partial\Omega$ and satisfying

$$(1.2) \quad V = -\kappa_{ss}$$

for all points on Γ_t with the boundary conditions

$$(1.3) \quad \begin{cases} \Gamma_t \perp \partial\Omega & (90^\circ\text{-angle condition}), \\ \kappa_s = 0 & (\text{no-flux condition}) \end{cases}$$

*Received by the editors June 18, 2007; accepted for publication (in revised form) February 4, 2008; published electronically May 28, 2008. This work was supported by the Regensburger Universitätsstiftung Hans Vielberth and the Research Fellowship of the Japan Society for the Promotion of Young Scientists.

<http://www.siam.org/journals/sima/40-2/69475.html>

[†]NWFI-Mathematik, Universitaet Regensburg, 93040 Regensburg, Germany (harald.garcke@mathematik.uni-regensburg.de).

[‡]Advanced Algorithm & Systems, Ebisu IS Building, Tokyo 150-0013, Japan (k.ito@aas-ri.co.jp).

[§]Muroran Institute of Technology, Muroran 050-8585, Japan (kohsaka@mmm.muroran-it.ac.jp).

at $\Gamma_t \cap \partial\Omega$, where a subscript s denotes differentiation with respect to the arc-length parameter of the evolving curve Γ_t . The boundary conditions (1.3) are the natural boundary conditions when viewing the flow as the H^{-1} -gradient flow of the length functional. It is not difficult to show that under the surface diffusion flow (1.2) with the boundary conditions (1.3) the areas enclosed by Γ_t and $\partial\Omega$ are preserved and the length of Γ_t decreases in time. We also find that an arc of a circle or a line segment is stationary under (1.2) and (1.3). Our goal in this paper is to show a nonlinear stability result for stationary solutions to (1.2) and (1.3). A proof of such a result is difficult due to the area-preserving property and due to the fact that highly nonlinear boundary conditions appear. We remark that for nonlinear boundary conditions satisfactory stability results are not available within the context of semigroup theory. We also remark that it is not possible to use methods based on maximum or comparison principles which have been used for mean curvature flow; see [8, 9].

For closed curves evolving by surface diffusion, Elliott and Garcke [10] showed a global existence result in the case that the initial curve is close to a circle. In addition, they proved nonlinear stability of circles under surface diffusion. Escher, Mayer, and Simonett [12] generalized the result in [10] to the higher-dimensional case. For evolving curves which come into contact with the outer boundary, Garcke, Ito, and Kohsaka [13] studied the linearized stability of stationary curves for (1.2) and (1.3). They derived a linearized stability criterion by extending the work for mean curvature flow of [8, 9, 15] to motion by surface diffusion. For three evolving curves with a triple junction in the case that the outer boundary $\partial\Omega$ is a rectangle [11, 16] or a triangle [17], global existence results when the initial curve is a small perturbation of a certain stationary curve have been shown. Also, nonlinear stability of this stationary curve can be shown.

Since the proof of nonlinear stability will depend heavily on the linear stability criterion derived in [13], we will now state it in detail. Let Γ^* be a stationary curve parameterized by X^* such that

$$\Gamma^* = \{X^*(\sigma) \mid \sigma \in [l_-, l_+]\},$$

where σ is the arc-length parameter along Γ^* and $X^*(l_\pm) \in \partial\Omega$. Further, we denote by κ^* the curvature of Γ^* and by h_\pm^* the curvature of $\partial\Omega$ at $X^*(l_\pm)$, where we assume the sign convention that h_\pm^* is negative if Ω is convex. Then the linearized stability criterion requires that

$$(1.4) \quad I^*[w, w] = \int_{l_-}^{l_+} \{w_\sigma^2 - (\kappa^*)^2 w^2\} d\sigma + h_+^*(w^2|_{\sigma=l_+}) + h_-^*(w^2|_{\sigma=l_-})$$

is positive for all $w \in H^1(\Gamma^*)$ with mean value zero. In [13] this criterion was derived by studying the stability of the linearized problem. The same bilinear form also appears if one computes the second variation of the length functional taking boundary contacts into account; see, e.g., Vogel [23]. We refer the reader to section 7 of [13] for several examples in which the linearized stability criterion has been applied. In the papers [2, 3] numerical results on the stability of stationary solutions for surface diffusion are presented.

Our methods to obtain a nonlinear stability result are the following. First we introduce new curvilinear coordinates in order to derive an appropriate parameterization for which we can formulate (1.2) and (1.3) in a PDE setting. We then prove a local existence result, where the local existence time depends only on the $C^{2+\alpha}$ -norm

($0 < \alpha < 1$) of the initial curve. This is very helpful for a global existence result because we need a priori estimates only up to two spatial derivatives. In fact, by applying an energy method as in [6, 10, 16, 17] to a resulting evolution equation for the curvature, we can derive an a priori estimate of the L^2 -norm of κ_s , which implies the boundedness of the $C^{2+\alpha}$ -norm ($0 < \alpha < 1/2$) of the solution for $t > 0$. In the derivation of this a priori estimate, the linearized stability criterion developed in [13] is used. In addition, we need to understand the set of stationary solutions. We can use a result by Vogel [23] which guarantees that linearly stable stationary solutions are strict local minimizers of the length functional under an area constraint. We also show that in the neighborhood of the linearly stable stationary solution other stationary solutions can be represented as a one parameter family, where the parameter is the enclosed area. This implies that the linearly stable stationary solution is isolated, a fact which will be important in order to study the long time behavior of solutions.

This paper proceeds as follows. In section 2, a parameterization established in [13] is employed for the geometric evolution equation (1.2) with boundary conditions (1.3). As a consequence, we obtain a nonlinear fourth order parabolic PDE with nonlinear boundary conditions. We show a local existence result for this nonlinear parabolic problem. For the reader's convenience we show an essential part of the proof of the local existence result in Appendix A. In section 3, an evolution equation for the curvature is derived together with some geometric identities. The evolution equation for the curvature allows it to apply an energy method as in [6, 10, 16, 17]. In section 4, we first derive a priori estimates for the length of Γ_t and the L^2 -norm of κ_s when Γ_t is close to a linearly stable stationary curve. These estimates imply the boundedness of the $C^{2+\alpha}$ -norm ($0 < \alpha < 1/2$) of the solution for $t > 0$, so that the global existence result is proven when the initial curve is close to a linearly stable stationary curve. Finally, in section 5, we show nonlinear stability of linearly stable stationary curves.

2. Local existence and uniqueness. In order to derive local existence and uniqueness for the geometric evolution equation (1.2) with the boundary conditions (1.3), we employ a parameterization which was established in [13]. We remark that our parameterization describes the curves close to a stationary curve by a modified distance function over a fixed interval and leads to a single PDE in contrast to parameterizations used, e.g., in Bronsard and Reitich [4] which involve vector-valued functions and lead to a system of PDEs. We also want to avoid arc-length parameterizations, as they lead to time-dependent parameter intervals, which is less convenient for semigroup theory. For the reader's convenience, we give a detailed derivation of our parameterization in the following.

Let $\Omega \subset \mathbb{R}^2$ be a domain such that

$$\Omega = \{x \in \mathbb{R}^2 \mid \psi(x) < 0\}, \quad \partial\Omega = \{x \in \mathbb{R}^2 \mid \psi(x) = 0\}$$

with a smooth function $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$ fulfilling $\nabla\psi(x) \neq 0$ for x with $\psi(x) = 0$.

Also, let Γ^* be a stationary curve under the flow (1.2) and (1.3); i.e., Γ^* has constant curvature κ^* . We now introduce an arc-length parameterization of Γ^* in the form

$$\Gamma^* = \{\Phi^*(\sigma) \mid \sigma \in [l_-, l_+]\},$$

where Φ^* is a mapping from $[l_-, l_+]$ to \mathbb{R}^2 and $l_+ - l_-$ is the total length of Γ^* . Note that we can extend Γ^* naturally either to the full circle when Γ^* is an arc of a circle

or to a straight line when Γ^* is a line segment. We set

$$\bar{l} := \begin{cases} \pi/|\kappa^*| & \text{if } \kappa^* \neq 0, \\ +\infty & \text{if } \kappa^* = 0; \end{cases}$$

i.e., \bar{l} is the length of the extension of Γ^* to a half circle if $\kappa_* \neq 0$. Without loss of generality, we can assume $[l_-, l_+] \subset (-\bar{l}, \bar{l})$. Define

$$\begin{cases} \xi_+(q) := \max\{\sigma \in (-\bar{l}, \bar{l}) \mid \Phi^*(\sigma) + qN^*(\sigma) \in \Omega\}, \\ \xi_-(q) := \min\{\sigma \in (-\bar{l}, \bar{l}) \mid \Phi^*(\sigma) + qN^*(\sigma) \in \Omega\}, \end{cases}$$

where $N^*(\sigma)$ is a unit normal vector of Γ^* at σ and is obtained by rotating the unit tangent vector $T^*(\sigma)$ of Γ^* by $\pi/2$. Above, q is a parameter with $q \in (-\bar{d}, \bar{d})$ for a small and given $\bar{d} > 0$. It holds that $\psi(\Phi^*(\xi_\pm(q)) + qN^*(\xi_\pm(q))) = 0$ and $\xi_\pm(0) = l_\pm$. Using the implicit function theorem, we see that $\xi_+(q)$ and $\xi_-(q)$ are smooth. Let

$$\Psi(\sigma, q) := \Phi^*(\xi(\sigma, q)) + qN^*(\xi(\sigma, q))$$

with

$$\xi(\sigma, q) := \xi_-(q) + \frac{\sigma - l_-}{l_+ - l_-}(\xi_+(q) - \xi_-(q)).$$

It is not difficult to check that $\xi(l_\pm, q) = \xi_\pm(q)$ and $\xi(\sigma, 0) = \sigma$.

In addition, one derives that $\Psi : (l_-, l_+) \times (-\bar{d}, \bar{d}) \rightarrow \Omega$ parameterizes the intersection W of a tubular neighborhood around the extended Γ^* with Ω . We now consider functions $\rho : [l_-, l_+] \rightarrow (-\bar{d}, \bar{d})$ and obtain $\Psi(\sigma, \rho(\sigma)) \in W$ for $\sigma \in (l_-, l_+)$. Then we define $\Phi(\sigma) := \Psi(\sigma, \rho(\sigma))$ for $\sigma \in [l_-, l_+]$, which is a parameterization of a curve Γ . An evolving curve is now given by

$$(2.1) \quad \Gamma_t := \{\Phi(\sigma, t) \mid \sigma \in [l_-, l_+]\}$$

with $\Phi(\sigma, t) := \Psi(\sigma, \rho(\sigma, t))$ for a function $\rho = \rho(\sigma, t)$. We note that $|\rho(\sigma, t)| < \bar{d}$ guarantees that $\Phi(\sigma, t) = \Psi(\sigma, \rho(\sigma, t)) \in W$ for $\sigma \in (l_-, l_+)$ and $t > 0$. We remark that $\rho \equiv 0$ corresponds to the stationary curve Γ^* .

Let us now express (1.2) and (1.3) with the help of parameterizations which have the form (2.1). For the arc-length parameter s of Γ_t , we have

$$(2.2) \quad \frac{ds}{d\sigma} = |\Phi_\sigma| = \sqrt{|\Psi_\sigma|^2 + 2(\Psi_\sigma, \Psi_q)_{\mathbb{R}^2} \rho_\sigma + |\Psi_q|^2 \rho_\sigma^2} =: J(\rho).$$

By $|\cdot|$ and $(\cdot, \cdot)_{\mathbb{R}^2}$ we denote the norm and the inner product in \mathbb{R}^2 , respectively. Then we find

$$T = \frac{1}{J(\rho)} \Phi_\sigma, \quad N = \frac{1}{J(\rho)} R\Phi_\sigma,$$

where T and N are the unit tangent and unit normal to Γ_t , respectively, and R is the rotation by the angle $\pi/2$. The normal velocity V of Γ_t is given by

$$V = (\Phi_t, N)_{\mathbb{R}^2} = \frac{1}{J(\rho)} (\Phi_t, R\Phi_\sigma)_{\mathbb{R}^2} = \frac{1}{J(\rho)} (\Psi_q, R\Psi_\sigma)_{\mathbb{R}^2} \rho_t.$$

Further, the Laplace–Beltrami operator $\Delta(\rho)$ on Γ_t is given via (2.2) as

$$(2.3) \quad \Delta(\rho) = \partial_s^2 = \frac{1}{J(\rho)} \partial_\sigma \left(\frac{1}{J(\rho)} \partial_\sigma \right) = \frac{1}{(J(\rho))^2} \partial_\sigma^2 + \frac{1}{J(\rho)} \left(\partial_\sigma \frac{1}{J(\rho)} \right) \partial_\sigma.$$

Then the curvature κ of Γ_t can be derived by using $\Delta(\rho)$ as

$$(2.4) \quad \begin{aligned} \kappa(\rho) &= (\Delta(\rho)\Phi, N)_{\mathbb{R}^2} = \frac{1}{(J(\rho))^3} (\Phi_{\sigma\sigma}, R\Phi_\sigma)_{\mathbb{R}^2} \\ &= \frac{1}{(J(\rho))^3} \left[(\Psi_q, R\Psi_\sigma)_{\mathbb{R}^2} \rho_{\sigma\sigma} + \{2(\Psi_{\sigma q}, R\Psi_\sigma)_{\mathbb{R}^2} + (\Psi_{\sigma\sigma}, R\Psi_q)_{\mathbb{R}^2}\} \rho_\sigma \right. \\ &\quad \left. + \{(\Psi_{qq}, R\Psi_\sigma)_{\mathbb{R}^2} + 2(\Psi_{\sigma q}, R\Psi_q)_{\mathbb{R}^2} + (\Psi_{qq}, R\Psi_q)_{\mathbb{R}^2} \rho_\sigma\} \rho_\sigma^2 \right. \\ &\quad \left. + (\Psi_{\sigma\sigma}, R\Psi_\sigma)_{\mathbb{R}^2} \right]. \end{aligned}$$

Furthermore, we note that the Neumann boundary condition $(\Phi_\sigma, T_{\partial\Omega})_{\mathbb{R}^2} = 0$ on $\partial\Omega$ is equivalent to the condition $(R\Phi_\sigma, \nabla\psi(\Phi))_{\mathbb{R}^2} = 0$ on $\partial\Omega$. Then we compute that the parameterization of the Neumann boundary condition is

$$(R\Psi_\sigma + R\Psi_q \rho_\sigma, \nabla\psi(\Psi))_{\mathbb{R}^2} = 0 \text{ at } \sigma = l_\pm.$$

As a consequence, we conclude that the problem (1.2) and (1.3) is represented by

$$(2.5) \quad \begin{cases} \rho_t = -L(\rho)\Delta(\rho)\kappa(\rho) & \text{for } \sigma \in (l_-, l_+), t > 0, \\ (R\Psi_\sigma + R\Psi_q \rho_\sigma, \nabla\psi(\Psi))_{\mathbb{R}^2} = 0 & \text{at } \sigma = l_\pm, \\ \partial_\sigma \kappa(\rho) = 0 & \text{at } \sigma = l_\pm. \end{cases}$$

Here $L(\rho) := J(\rho)/(\Psi_q, R\Psi_\sigma)_{\mathbb{R}^2}$; $\Delta(\rho)$ and $\kappa(\rho)$ are given by (2.3) and (2.4), respectively.

Let $\mathcal{I} = [l_-, l_+]$ and $\mathcal{Q}_{t_0, t_1} = \mathcal{I} \times (t_0, t_1]$ for $0 \leq t_0 < t_1 < \infty$. For $0 < \alpha < 1$, we define the function space

$$\mathcal{Y}(\overline{\mathcal{Q}_{t_0, t_1}}) = \{\rho \in C^{2+\alpha, 0}(\overline{\mathcal{Q}_{t_0, t_1}}) \cap C^{4+\alpha, 1}(\mathcal{Q}_{t_0, t_1}) \mid \|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{t_0, t_1}})} < \infty\}$$

with the norm

$$\begin{aligned} \|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{t_0, t_1}})} &= \sup_{t_0 \leq t \leq t_1} \|\rho(\cdot, t)\|_{C^{2+\alpha}(\mathcal{I})} + \sup_{t_0 < t \leq t_1} (t - t_0)^{1/2} \|\partial_\sigma^4 \rho(\cdot, t)\|_{C^\alpha(\mathcal{I})} \\ &\quad + \sup_{t_0 < t \leq t_1} (t - t_0)^{1/2} \|\rho_t(\cdot, t)\|_{C^\alpha(\mathcal{I})}, \end{aligned}$$

where $\overline{\mathcal{Q}_{t_0, t_1}}$ is the closure of \mathcal{Q}_{t_0, t_1} .

Now we are ready to state a local existence theorem.

THEOREM 2.1 (local existence). *Let $\alpha \in (0, 1)$ and let us assume that $\rho_0 \in C^{2+\alpha}(\mathcal{I})$ with $\|\rho_0\|_{C^0(\mathcal{I})} < \bar{d}$ fulfills*

$$(R\Psi_\sigma + R\Psi_q \rho_\sigma, \nabla\psi(\Psi))_{\mathbb{R}^2} = 0 \text{ at } \sigma = l_\pm.$$

Then there exists a $T_0 = T_0(1/\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}) > 0$ such that the problem (2.5) with $\rho(\cdot, 0) = \rho_0$ has a unique solution in $\mathcal{Y}(\overline{\mathcal{Q}_{0, T_0}})$.

This theorem is proved by applying similar arguments as in [16]. Since we have to take care of the boundary conditions in a different way, we will sketch the proof in Appendix A.

Remark 2.2. Applying the local existence results used in [4] directly, we obtain solutions for $C^{4+\alpha}$ -initial curves. But this makes it difficult to derive a global existence result because we need a priori estimates for higher order derivatives. Thus our local existence result is an improvement of the one obtained in [4].

Remark 2.3. By using a bootstrapping argument as in [16, Theorem 3.6, Remark 3.7], it can be shown that the solution ρ established in Theorem 2.1 is smooth for $t \in (0, T_0]$.

3. An evolution equation for curvature. In order to show nonlinear stability of solutions for which the linearized stability criterion of [13] is fulfilled, we apply an energy method similar to the one used in [6, 10, 16, 17]. For this approach it is important to derive an evolution equation for the curvature. Such an equation will be useful for the derivation of a priori estimates with the help of the linearized stability criterion.

For the above-mentioned purpose, we employ a parameterization of the evolving curve Γ_t by arc length contrary to the one stated in section 2. Let X be a smooth mapping so that $X(\cdot, t)$ is an arc-length parameterization of Γ_t , i.e.,

$$\Gamma_t := \{X(s, t) \mid s \in [r_-(t), r_+(t)]\}$$

for any $t > 0$, where r_+ and r_- are smooth in t . In particular, $X(r_\pm(t), t) \in \partial\Omega$ and $r_+(t) - r_-(t) = L[\Gamma_t]$, where $L[\Gamma_t]$ denotes the total length of Γ_t . Let $N (= N(s, t))$ be the unit normal vector of Γ_t , which is represented as

$$N(s, t) = \begin{pmatrix} \cos \theta(s, t) \\ \sin \theta(s, t) \end{pmatrix}.$$

Also, let $T (= T(s, t))$ and $\kappa (= \kappa(s, t))$ be the unit tangent vector of Γ_t and the curvature of Γ_t , respectively. Note that the unit tangent vector T is obtained by rotating the unit normal vector N by $-\pi/2$. Then, using $\theta_s = \kappa$, we have

$$(3.1) \quad \begin{cases} N_s = -\kappa T, & T_s = \kappa N, \\ N_t = -\theta_t T, & T_t = \theta_t N. \end{cases}$$

In addition, set

$$V := (X_t, N)_{\mathbb{R}^2}, \quad v := (X_t, T)_{\mathbb{R}^2}.$$

Note that V and v are the normal velocity and the tangent velocity of X , respectively. Then it follows that

$$(3.2) \quad X_t = VN + vT.$$

Differentiating (3.2) with respect to s and using (3.1), we have

$$\begin{aligned} X_{ts} &= V_s N + VN_s + v_s T + vT_s \\ &= (V_s + \kappa v)N + (-\kappa V + v_s)T. \end{aligned}$$

This implies the following lemma.

LEMMA 3.1. *Let X be a smooth arc-length parameterization as above. Then*

$$\theta_t = V_s + \kappa v, \quad v_s = \kappa V.$$

Proof. Since $X_{ts} = X_{st}$ and $X_s = T$, it follows from (3.1) that

$$\theta_t N = (V_s + \kappa v)N + (-\kappa V + v_s)T.$$

Thus we obtain the desired results. \square

By Lemma 3.1, we have the following formula for the time derivative of curvature.

LEMMA 3.2. *Let X be a smooth arc-length parameterization as above. Then*

$$\kappa_t = V_{ss} + \kappa^2 V + \kappa_s v.$$

Proof. By $\theta_s = \kappa$ and Lemma 3.1, we derive

$$\kappa_t = \theta_{st} = \theta_{ts} = (V_s + \kappa v)_s = V_{ss} + \kappa v_s + \kappa_s v = V_{ss} + \kappa^2 V + \kappa_s v.$$

This completes the proof. \square

By the assumption that Γ_t touches $\partial\Omega$ with the angle $\pi/2$, we have

$$\psi(X(r_{\pm}(t), t)) = 0, \quad (\nabla\psi(X), N)_{\mathbb{R}^2} = 0 \quad \text{at } s = r_{\pm}(t).$$

Then we derive the following lemma.

LEMMA 3.3. *Let X be a smooth arc-length parameterization as above. Then*

$$v(r_{\pm}(t), t) + r'_{\pm}(t) = 0.$$

Proof. Differentiating $\psi(X(r_{\pm}(t), t)) = 0$ with respect to t and using $(\nabla\psi(X), N)_{\mathbb{R}^2} = 0$ at $s = r_{\pm}(t)$, we have at $s = r_{\pm}(t)$

$$\begin{aligned} 0 &= (\nabla\psi(X), X_s r'_{\pm} + X_t)_{\mathbb{R}^2} = (\nabla\psi(X), X_s r'_{\pm} + VN + vT)_{\mathbb{R}^2} \\ &= (v + r'_{\pm})(\nabla\psi(X), T)_{\mathbb{R}^2} = \pm (v + r'_{\pm})|\nabla\psi(X)|. \end{aligned}$$

The last identity is derived with the help of $T = \pm \nabla\psi(X)/|\nabla\psi(X)|$ at $s = r_{\pm}(t)$. Since $|\nabla\psi(X)| \neq 0$, we obtain the desired result. \square

Now we can present an evolution equation for the curvature.

PROPOSITION 3.4 (evolution equation for the curvature). *Let evolving curves $\Gamma = \{\Gamma_t\}_{t \geq 0}$ be lying in Ω with $\partial\Gamma \subset \partial\Omega$. Then a smooth solution of*

$$(3.3) \quad V = -\kappa_{ss} \quad \text{on } \Gamma_t$$

with the boundary conditions

$$(3.4) \quad \begin{cases} \angle(\Gamma_t, \partial\Omega) = \pi/2 & \text{at } \Gamma_t \cap \partial\Omega, \\ \kappa_s = 0 & \text{at } \Gamma_t \cap \partial\Omega \end{cases}$$

fulfills for $t > 0$

$$(3.5) \quad \kappa_t = -\kappa_{ssss} - \kappa^2 \kappa_{ss} + \kappa_s v \quad \text{on } \Gamma_t$$

and

$$(3.6) \quad \begin{cases} \kappa_s = 0 & \text{at } \Gamma_t \cap \partial\Omega, \\ (\partial_s \pm h_{\pm})\kappa_{ss} = 0 & \text{at } \Gamma_t \cap \partial\Omega. \end{cases}$$

Here h_{\pm} is the curvature of $\partial\Omega$ at the points $X(r_{\pm}(t), t) \in \Gamma_t \cap \partial\Omega$ with the sign convention that $h_{\pm} \leq 0$ if Ω is convex.

Proof. We immediately obtain (3.5) from (3.3) and Lemma 3.2. Next we show (3.6). Differentiating $(\nabla\psi(X), N)_{\mathbb{R}^2} = 0$ at $s = r_{\pm}(t)$ with respect to t and using (3.1), (3.2), Lemma 3.1, and Lemma 3.3, we have at $s = r_{\pm}(t)$

$$\begin{aligned} 0 &= ([D^2\psi(X)](X_s r'_{\pm} + X_t), N)_{\mathbb{R}^2} + (\nabla\psi(X), N_s r'_{\pm} + N_t)_{\mathbb{R}^2} \\ &= (v + r'_{\pm})([D^2\psi(X)]T, N)_{\mathbb{R}^2} + V([D^2\psi(X)]N, N)_{\mathbb{R}^2} \\ &\quad - \kappa r'_{\pm}(\nabla\psi(X), T)_{\mathbb{R}^2} - \theta_t(\nabla\psi(X), T)_{\mathbb{R}^2} \\ &= V([D^2\psi(X)]T_{\partial\Omega}(X), T_{\partial\Omega}(X))_{\mathbb{R}^2} - V_s(\nabla\psi(X), T)_{\mathbb{R}^2} \\ &\quad - \kappa(v + r'_{\pm})(\nabla\psi(X), T)_{\mathbb{R}^2} \\ &= V([D^2\psi(X)]T_{\partial\Omega}(X), T_{\partial\Omega}(X))_{\mathbb{R}^2} \mp V_s|\nabla\psi(X)|. \end{aligned}$$

Here $D^2\psi$ is the Hessian matrix of ψ . Then we observe that

$$\kappa_{\partial\Omega}(X) = -\frac{1}{|\nabla\psi(X)|}([D^2\psi(X)]T_{\partial\Omega}(X), T_{\partial\Omega}(X))_{\mathbb{R}^2},$$

so that

$$V_s \pm h_{\pm}V = 0 \quad \text{at } s = r_{\pm}(t),$$

where h_{\pm} are given by $h_{\pm} := \kappa_{\partial\Omega}(X(r_{\pm}(t), t))$. This completes the proof. \square

4. A priori estimates and global existence. We now derive basic evolution formulas for length and $\int_{\Gamma_t} \kappa_s^2 ds$.

LEMMA 4.1. *A smooth solution of (3.3)–(3.4) fulfills*

$$\begin{aligned} \text{(i)} \quad &\frac{d}{dt}L[\Gamma_t] = -\int_{\Gamma_t} \kappa_s^2 ds, \\ \text{(ii)} \quad &\frac{d}{dt}\int_{\Gamma_t} \kappa_s^2 ds = -2\left\{\int_{\Gamma_t} V_s^2 ds - \int_{\Gamma_t} \kappa^2 V^2 ds + h_+(V^2|_{s=r_+(t)}) \right. \\ &\quad \left. + h_-(V^2|_{s=r_-(t)})\right\} + \int_{\Gamma_t} \kappa_s^2 \kappa V ds, \end{aligned}$$

where h_{\pm} is evaluated at $X(r_{\pm}(t), t)$.

Proof. Recalling $L[\Gamma_t] = r_+(t) - r_-(t)$ and using Lemmas 3.1 and 3.3, we have

$$\begin{aligned} \frac{d}{dt}L[\Gamma_t] &= r'_+(t) - r'_-(t) = -v(r_+(t), t) + v(r_-(t), t) = -\int_{\Gamma_t} v_s ds \\ &= -\int_{\Gamma_t} \kappa V ds = \int_{\Gamma_t} \kappa \kappa_{ss} ds = -\int_{\Gamma_t} \kappa_s^2 ds. \end{aligned}$$

The last term is derived using integration by parts and $\kappa_s = 0$ at $\Gamma_t \cap \partial\Omega$.

In order to prove (ii), we compute

$$(4.1) \quad \int_{\Gamma_t} \kappa_s(\kappa_t)_s ds = \int_{\Gamma_t} \kappa_s(-\kappa_{ssss} - \kappa^2 \kappa_{ss} + \kappa_s v)_s ds.$$

Since $\kappa_{ts} = \kappa_{st}$ and $\kappa_s = 0$ at $\Gamma_t \cap \partial\Omega$, we have

$$\text{(left-hand side of (4.1))} = \int_{\Gamma_t} \kappa_s \kappa_{st} ds = \frac{1}{2} \int_{\Gamma_t} (\kappa_s^2)_t ds = \frac{1}{2} \frac{d}{dt} \int_{\Gamma_t} \kappa_s^2 ds.$$

On the other hand, by means of integration by parts and using (3.6), we derive

$$\begin{aligned} \text{(right-hand side of (4.1))} &= - \int_{\Gamma_t} \kappa_{ss} (-\kappa_{ssss} - \kappa^2 \kappa_{ss} + \kappa_s v) ds \\ &= \int_{\Gamma_t} \kappa_{ss} \kappa_{ssss} ds + \int_{\Gamma_t} \kappa^2 \kappa_{ss}^2 ds - \int_{\Gamma_t} \kappa_{ss} \kappa_s v ds \\ &= -h_+(\kappa_{ss}^2|_{s=r_+(t)}) - h_-(\kappa_{ss}^2|_{s=r_-(t)}) - \int_{\Gamma_t} \kappa_{ss}^2 ds \\ &\quad + \int_{\Gamma_t} \kappa^2 \kappa_{ss}^2 ds + \frac{1}{2} \int_{\Gamma_t} \kappa_s^2 v_s ds. \end{aligned}$$

Thus it follows from $V = -\kappa_{ss}$ and $v_s = \kappa V$ that

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\Gamma_t} \kappa_s^2 ds &= - \left\{ \int_{\Gamma_t} V_s^2 ds - \int_{\Gamma_t} \kappa^2 V^2 ds + h_+(V^2|_{s=r_+(t)}) + h_-(V^2|_{s=r_-(t)}) \right\} \\ &\quad + \frac{1}{2} \int_{\Gamma_t} \kappa_s^2 \kappa V ds. \end{aligned}$$

This completes the proof. \square

Let us define the bilinear form I as

$$(4.2) \quad I[w, w] = \int_{r_-}^{r_+} (w_s^2 - \kappa_{av}^2 w^2) ds + h_+(w^2|_{s=r_+}) + h_-(w^2|_{s=r_-})$$

for $w \in H^1(\Gamma_t)$ with

$$\int_{r_-}^{r_+} w ds = 0.$$

Here s is the arc-length parameter along Γ_t , which belongs to the interval $[r_-, r_+]$ with $L[\Gamma_t] = r_+ - r_-$; h_{\pm} is the curvature of $\partial\Omega$ at $\Gamma_t \cap \partial\Omega$; and κ_{av} is the averaged curvature of Γ_t defined as

$$\kappa_{av} = \frac{1}{L[\Gamma_t]} \int_{r_-}^{r_+} \kappa ds.$$

Since $V = -\kappa_{ss}$ and $\kappa_s = 0$ at $\Gamma_t \cap \partial\Omega$, it holds that

$$(4.3) \quad \int_{\Gamma_t} V ds = 0.$$

Then we can rewrite Lemma 4.1(ii) as

$$(4.4) \quad \frac{d}{dt} \int_{\Gamma_t} \kappa_s^2 ds + 2I[V, V] = -2 \int_{\Gamma_t} (\kappa_{av}^2 - \kappa^2) V^2 ds + \int_{\Gamma_t} \kappa_s^2 \kappa V ds.$$

Remark 4.2. Although the 90° -angle condition in (1.3) is the natural boundary condition when considering the gradient flow of the length functional, we can also consider the case where the prescribed angle is not 90° . In this case, the angle condition is represented as

$$(\nabla\psi(X), N)_{\mathbb{R}^2} = |\nabla\psi(X)| \cos \Theta_\pm \quad \text{at } s = r_\pm(t)$$

for $\Theta_\pm \in (0, \pi)$. Then the identity in Lemma 3.3 is replaced by

$$v + r'_\pm = \mp V \cot \Theta_\pm \quad \text{at } s = r_\pm(t),$$

and the second boundary condition of (3.6) is replaced by

$$\{\partial_s + (\pm h_\pm \csc \Theta_\pm \mp \kappa \cot \Theta_\pm)\} \kappa_{ss} = 0 \quad \text{at } s = r_\pm(t).$$

Further, setting

$$\begin{aligned} I_\Theta[w, w] &= \int_{r_-}^{r_+} (w_s^2 - \kappa_{av}^2 w^2) ds + \{h_+ \csc \Theta_+ - (\kappa|_{s=r_+}) \cot \Theta_+\} (w^2|_{s=r_+}) \\ &\quad + \{h_- \csc \Theta_- - (\kappa|_{s=r_-}) \cot \Theta_-\} (w^2|_{s=r_-}) \end{aligned}$$

instead of I in (4.2), we have

$$\frac{d}{dt} \int_{\Gamma_t} \kappa_s^2 ds + 2I_\Theta[V, V] = -2 \int_{\Gamma_t} (\kappa_{av}^2 - \kappa^2) V^2 ds + \int_{\Gamma_t} \kappa_s^2 \kappa V ds.$$

We remark that for the non- 90° -angle condition the stationary solutions with the property that I_Θ is positive are strict local minimizers of an energy also involving a wetting energy at the boundary if one takes an area constraint into account (see Vogel [23]). On the other hand,

$$\frac{d}{dt} L[\Gamma_t] = - \int_{\Gamma_t} \kappa_s^2 ds - (V|_{s=r_+(t)}) \cot \Theta_+ - (V|_{s=r_-(t)}) \cot \Theta_-,$$

so that the monotonicity of $L[\Gamma_t]$ with respect to t is no longer fulfilled.

The following lemmas will be crucial in order to derive an a priori estimate.

LEMMA 4.3. *A smooth solution of (3.3)–(3.4) fulfills*

(i) $\left| \int_{\Gamma_t} \kappa_s^2 \kappa \kappa_{ss} ds \right| \leq \frac{1}{3} L[\Gamma_t] \|\kappa_s\|_{L^2(\Gamma_t)}^2 \|\kappa_{ss}\|_{L^2(\Gamma_t)}^2,$

(ii) $\|\kappa - \kappa_{av}\|_{C^0(\Gamma_t)} \leq L[\Gamma_t]^{1/2} \|\kappa_s\|_{L^2(\Gamma_t)}.$

Proof. We first prove (i). Since $\kappa_s = 0$ at $\Gamma_t \cap \partial\Omega$, we get

$$\int_{\Gamma_t} \kappa_s^2 \kappa \kappa_{ss} ds = -\frac{1}{3} \int_{\Gamma_t} \kappa_s^4 ds.$$

Then it follows that

$$\begin{aligned} \left| \int_{\Gamma_t} \kappa_s^4 ds \right| &\leq \|\kappa_s\|_{L^2(\Gamma_t)}^2 \|\kappa_s\|_{L^\infty(\Gamma_t)}^2 \\ &\leq L[\Gamma_t] \|\kappa_s\|_{L^2(\Gamma_t)}^2 \|\kappa_{ss}\|_{L^2(\Gamma_t)}^2. \end{aligned}$$

The last term is derived by using a Poincaré inequality since $\kappa_s = 0$ at $\Gamma_t \cap \partial\Omega$.

Next we prove (ii). Since

$$\int_{\Gamma_t} (\kappa - \kappa_{av}) ds = 0,$$

for each $t > 0$, there is a $r_0 (= r_0(t)) \in (r_-(t), r_+(t))$ such that $\kappa(r_0, t) - \kappa_{av}(t) = 0$. This implies that

$$|\kappa(s, \cdot) - \kappa_{av}| = \left| \int_{r_0}^s (\kappa - \kappa_{av})_s ds \right| = \left| \int_{r_0}^s \kappa_s ds \right| \leq \int_{\Gamma_t} |\kappa_s| ds \leq L[\Gamma_t]^{1/2} \|\kappa_s\|_{L^2(\Gamma_t)}.$$

Thus we have the desired result. \square

We remind the reader that for functions w_1, w_2 with mean values zero we can define the H^{-1} -inner product via

$$(w_1, w_2)_{-1} = \int_{l_-}^{l_+} u_{1,\sigma} u_{2,\sigma} d\sigma,$$

where u_i is the solution of $-u_{i,\sigma\sigma} = w_i$ in (l_-, l_+) and $u_{i,\sigma} = 0$ at $\sigma = l_{\pm}$. According to [13], the bilinear form I^* as stated in the introduction (see (1.4)) is positive, provided that the maximal eigenvalue λ for the linearized problem to (1.2) and (1.3) is negative. In [13] it was shown that $I^*[w, w] \geq (-\lambda)(w, w)_{-1}$ for all $w \in H^1(\Gamma^*)$ with mean value zero. We now want to derive a perturbation of this result. Let us denote $L = L[\Gamma]$ and $L^* = L[\Gamma^*] (= l_+ - l_-)$. Then we have the following lemma, which implies a lower bound for I when the parameters κ_{av} , h_{\pm} , and L are close to κ^* , h_{\pm}^* , and L^* , respectively.

LEMMA 4.4. (i) *Let λ be the maximal eigenvalue of the linearized problem. For $\varepsilon > 0$ there exists $\delta > 0$ such that*

$$I[w, w] > (-\lambda - \varepsilon)(w, w)_{-1}$$

for $w \in H^1(\Gamma)$ with mean value zero, provided that

$$|\kappa_{av} - \kappa^*| < \delta, \quad |h_{\pm} - h_{\pm}^*| < \delta, \quad |L - L^*| < \delta.$$

(ii) *There exists $\mu > 0$ such that*

$$\mu \|w_s\|_{L^2(\Gamma)}^2 \leq I[w, w] + (w, w)_{-1}$$

for $w \in H^1(\Gamma)$ with mean value zero.

Proof. The largest eigenvalue λ corresponding to the bilinear form I depends continuously on L , κ_{av} , and h_{\pm} . In the case that $L = L^*$, $\kappa_{av} = \kappa^*$, and $h_{\pm} = h_{\pm}^*$ we obtain (i) with $\varepsilon = 0$, and hence (i) follows from a straightforward perturbation argument; compare [13] for similar arguments. Arguing as in the proof of Lemma 5.3 in [13], we obtain (ii). \square

It is significant to obtain a positive lower bound of $L[\Gamma]$ in terms of ρ . The following lemma implies that L^* is a local minimum of $L[\Gamma]$, provided that I^* is positive.

LEMMA 4.5. *Let Γ^* be a stationary curve such that the bilinear form I^* is positive and let $\rho \in C^1(\mathcal{I})$ be a function describing a curve Γ close to Γ^* as in section 2. Assume that a curve Γ encloses the same area as Γ^* . Then there exist constants $\bar{c}, \gamma^* > 0$ such that*

$$L[\Gamma] \geq L^* + \bar{c} \|\rho\|_{H^1(\mathcal{I})}^2$$

if $\|\rho\|_{C^1(\mathcal{I})} < \gamma^*$.

Proof. This follows as in the proof of Theorem 2.1 of Vogel [23] (see (2.14) and the inequality after (2.19) in [23]). \square

By virtue of Lemma 4.5, we have an a priori estimate of $L[\Gamma_t]$ and can derive useful estimates concerning κ_{av} and h_{\pm} .

LEMMA 4.6. *Let the assumptions of Lemma 4.5 hold for a stationary curve Γ^* and all curves Γ_t , $t \in [0, T]$, described by $\rho(t) \in C^1(\mathcal{I})$ for the parameterization in section 2. Assume in particular that $\|\rho(t)\|_{C^1(\mathcal{I})} < \gamma^*$ for $t \in [0, T]$, where γ^* is as in Lemma 4.5. We then obtain the following:*

- (i) $L[\Gamma_0] \geq L[\Gamma_t] \geq L^*$ for all $t \in [0, T]$.
- (ii) There exist $K_1, K_2 > 0$ such that for $t \in [0, T]$

$$|\kappa_{av}(t) - \kappa^*| \leq K_1 |L[\Gamma_t] - L^*|, \quad |h_{\pm}(t) - h_{\pm}^*| \leq K_2 |L[\Gamma_t] - L^*|.$$

Proof. (i) follows from Lemma 4.1(i) and Lemma 4.5. To prove (ii), we compute

$$\kappa_{av} = \frac{1}{L[\Gamma_t]} \int_{\Gamma_t} \kappa \, ds = \frac{1}{L[\Gamma_t]} \int_{\Gamma_t} \theta_s \, ds = \frac{1}{L[\Gamma_t]} (\theta_+ - \theta_-).$$

A similar computation gives

$$\kappa^* = \frac{1}{L^*} (\theta_+^* - \theta_-^*).$$

Then we have

$$\begin{aligned} |\kappa_{av} - \kappa^*| &= \left| \frac{1}{L[\Gamma_t]} (\theta_+ - \theta_-) - \frac{1}{L^*} (\theta_+^* - \theta_-^*) \right| \\ &= \frac{1}{L[\Gamma_t] L^*} |L^* (\theta_+ - \theta_-) - L[\Gamma_t] (\theta_+^* - \theta_-^*)| \\ &\leq \left(\frac{1}{L^*} \right)^2 \left\{ |L^* (\theta_+ - \theta_- - (\theta_+^* - \theta_-^*))| + |L^* - L[\Gamma_t]| |\theta_+^* - \theta_-^*| \right\}. \end{aligned}$$

By means of the mean value theorem, the smoothness of $\partial\Omega$, and the $\pi/2$ angle condition, we see that the quantity $|\theta_+ - \theta_+^*| + |\theta_-^* - \theta_-|$ is estimated by $\|\rho\|_{C^0(\mathcal{I})}$. Using Lemma 4.5 and an embedding result, we obtain the first inequality in (ii).

Recall that $\kappa_{\partial\Omega}(X)$ is represented by

$$\kappa_{\partial\Omega}(X) = -\frac{1}{|\nabla\psi(X)|} ([D^2\psi(X)]T_{\partial\Omega}(X), T_{\partial\Omega}(X))_{\mathbb{R}^2}.$$

Since this expression does not depend on derivatives of ρ , the mean value theorem implies that the quantity $|h_{\pm} - h_{\pm}^*|$ is estimated by $\|\rho\|_{C^0(\mathcal{I})}$. Using Lemma 4.5 and an embedding result, we derive the second inequality in (ii). \square

Using Lemma 4.4, we obtain the existence of constants $\delta^* > 0$ and $\mu^* > 0$ such that

$$(4.5) \quad I[w, w] > -\frac{\lambda}{2} (w, w)_{-1} + \mu^* \|w_s\|_{L^2(\Gamma_t)}^2$$

for $w \in H^1(\Gamma_t)$ with mean value zero, provided that

$$(4.6) \quad |\kappa_{av}(t) - \kappa^*| < \delta^*, \quad |h_{\pm}(t) - h_{\pm}^*| < \delta^*, \quad |L[\Gamma_t] - L^*| < \delta^*.$$

We are now in a position to derive a priori estimates for solutions of (2.5) if the solution is close to Γ^* .

PROPOSITION 4.7. *Let the assumptions of Lemma 4.5 hold for a stationary curve Γ^* and a curve Γ_t described by $\rho(t) \in C^1(\mathcal{I})$ for the parameterization in section 2. Assume that for $t \in (0, T]$*

$$(4.7) \quad \|\rho(t)\|_{C^1(\mathcal{I})} < \gamma^* \quad \text{and} \quad |L[\Gamma_t] - L^*| \leq \frac{\delta^*}{1 + K_1 + K_2} \quad (=:\delta_1^*),$$

where γ^* is as in Lemma 4.5, K_1 and K_2 are as in Lemma 4.6, and δ^* is as in (4.6). Then there is a constant $\delta_1 > 0$ such that, if $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ for $t \in (0, T]$, it holds that

$$\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 + \mu^* \int_{t_0}^t \|V_s(\tau)\|_{L^2(\Gamma_t)}^2 d\tau \leq \|\kappa_s(t_0)\|_{L^2(\Gamma_t)}^2$$

for $t \in [t_0, T]$ with $t_0 > 0$, where μ^* is as in (4.5).

Proof. By (4.4), we have

$$\begin{aligned} & \frac{d}{dt} \|\kappa_s\|_{L^2(\Gamma_t)}^2 + 2I[V, V] \\ &= -2 \int_{\Gamma_t} (\kappa_{av}^2 - \kappa^2) V^2 ds + \int_{\Gamma_t} \kappa_s^2 \kappa V ds \\ &= 2 \int_{\Gamma_t} (\kappa - \kappa_{av})^2 V^2 ds + 4\kappa_{av} \int_{\Gamma_t} (\kappa - \kappa_{av}) V^2 ds + \int_{\Gamma_t} \kappa_s^2 \kappa V ds. \end{aligned}$$

By virtue of (4.7) and Lemma 4.6(ii), we also see that $\kappa_{av}(t)$, $h_{\pm}(t)$, and $L[\Gamma_t]$ satisfy (4.6). Then it follows from Lemma 4.3, Lemma 4.6(i), and (4.5) that there exist $C_1, C_2 > 0$ such that

$$\begin{aligned} & \frac{d}{dt} \|\kappa_s\|_{L^2(\Gamma_t)}^2 + (-\lambda) (V, V)_{-1} + 2\mu^* \|V_s\|_{L^2(\Gamma_t)}^2 \\ & \leq C_1 \|V\|_{L^2(\Gamma_t)}^2 \|\kappa_s\|_{L^2(\Gamma_t)}^2 + C_2 (\delta^* + |\kappa^*|) \|V\|_{L^2(\Gamma_t)}^2 \|\kappa_s\|_{L^2(\Gamma_t)}. \end{aligned}$$

Since $\|V\|_{L^\infty(\Gamma_t)} \leq C \|V_s\|_{L^2(\Gamma_t)}$ by virtue of (4.3), we derive $\|V\|_{L^2(\Gamma_t)} \leq \tilde{C} \|V_s\|_{L^2(\Gamma_t)}$. By means of this fact and $(-\lambda) (V, V)_{-1} \geq 0$, we are led to

$$(4.8) \quad \frac{d}{dt} \|\kappa_s\|_{L^2(\Gamma_t)}^2 + \{2\mu^* - \tilde{C}_1 \|\kappa_s\|_{L^2(\Gamma_t)}^2 - \tilde{C}_2 (\delta^* + |\kappa^*|) \|\kappa_s\|_{L^2(\Gamma_t)}\} \|V_s\|_{L^2(\Gamma_t)}^2 \leq 0.$$

Then we choose δ_1 such that

$$0 < \delta_1 < \min \left\{ \frac{\mu^*}{2\tilde{C}_1}, \left(\frac{\mu^*}{2\tilde{C}_2(\delta^* + |\kappa^*|)} \right)^2 \right\}.$$

Assuming $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ for $t \in (0, T]$, it follows that

$$(4.9) \quad \frac{d}{dt} \|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 + \mu^* \|V_s(t)\|_{L^2(\Gamma_t)}^2 \leq 0.$$

Integrating (4.9) with respect to t in the interval $[t_0, t]$, we derive the desired result. \square

Now we arrive at the main result in this section.

THEOREM 4.8 (global existence). *Let Γ^* be a stationary curve such that the bilinear form I^* is positive. Also, let $\rho_0 \in C^{2+\alpha}(\mathcal{I})$ be a function describing a curve Γ_0 , which is close to Γ^* as in section 2 and satisfies $\Gamma_0 \perp \partial\Omega$. Assume that a curve Γ_0 includes the same area as Γ^* . Then there exist constants $\gamma_0 > 0$ and $\delta_0 > 0$ such that if $\|\rho_0\|_{C^1(\mathcal{I})} < \gamma_0$ and $L[\Gamma_0] - L^* < \delta_0$, the problem (2.5) admits a unique global-in-time solution ρ with*

$$\|\rho(t)\|_{C^1(\mathcal{I})} < \gamma_0 \quad \text{and} \quad L[\Gamma_t] - L^* < \delta_0 \quad \text{for } t \geq 0,$$

where Γ_t is the curve parameterized by $\Psi(\sigma, \rho(\sigma, t))$ in section 2.

Proof. Choose γ_0 and δ_0 satisfying

$$(4.10) \quad 0 < \gamma_0 < \frac{\gamma^*}{2}, \quad 0 < \delta_0 < \frac{\delta_1^*}{2},$$

where γ^* is as in Lemma 4.5 and δ_1^* is as in (4.7). Assume that the initial curve Γ_0 satisfies $\|\rho_0\|_{C^1(\mathcal{I})} < \gamma_0$ and $L[\Gamma_0] - L^* < \delta_0$. Then Lemma 4.5 and an embedding result imply that

$$(4.11) \quad \|\rho_0\|_{C^0(\mathcal{I})} \leq C(L[\Gamma_0] - L^*) < C\delta_0.$$

Further, Lemma 4.6(i) implies that for $t > 0$

$$(4.12) \quad L[\Gamma_t] - L^* \leq L[\Gamma_0] - L^* < \delta_0.$$

We now prove that $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ for each time t in the existence interval of the solution, where δ_1 is as in Proposition 4.7. Let $0 < \beta < \alpha < 1/2$. By Theorem 2.1, we can construct a unique local-in-time solution for $\rho_0 \in C^{2+\beta}(\mathcal{I})$ and obtain the estimate

$$(4.13) \quad \|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}}_0, T_0)} \leq K_0,$$

where K_0 is a constant, which depends on $\|\rho_0\|_{C^{2+\beta}(\mathcal{I})}$ increasingly, and T_0 is the local existence time, which depends on $1/\|\rho_0\|_{C^{2+\beta}(\mathcal{I})}$ increasingly (for details, see Appendix B). According to the interpolation inequality for Hölder spaces and (4.11), we have

$$(4.14) \quad \|\rho_0\|_{C^{2+\beta}(\mathcal{I})} \leq C(\|\rho_0\|_{C^0(\mathcal{I})})^{\frac{\alpha-\beta}{2+\alpha}} (\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})})^{\frac{2+\beta}{2+\alpha}} \leq \tilde{C}\delta_0^{\frac{\alpha-\beta}{2+\alpha}}.$$

Set $t_0 := \delta_0^{\frac{\alpha-\beta}{2+\alpha}} > 0$. Then it follows from (4.13), (4.14), and the definition of $\mathcal{Y}(\overline{\mathcal{Q}}_0, T_0)$ that there exist $C > 0$ and $\nu > 0$ such that

$$\|\kappa_s(t_0)\|_{L^2(\Gamma_{t_0})}^2 \leq C\delta_0^\nu.$$

Since $\|\rho(t)\|_{C^1(\mathcal{I})}$ is continuous for $t \in [0, T_0]$, we see that $\|\rho(t)\|_{C^1(\mathcal{I})} < \gamma^*$ for $t \in [0, T]$ with a $T \in (0, T_0]$. Further, by (4.10) and (4.12), we have $L[\Gamma_t] - L^* < \delta_1^*$ for $t > 0$. Choose δ_0 such that $t_0 < T$ and $C\delta_0^\nu < \delta_1$. Then, by applying a similar argument to [10, Proof of Theorem 6.1] together with Proposition 4.7, we obtain that $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ for $t \in [t_0, T]$.

Next we prove that $\|\rho(t)\|_{C^1(\mathcal{I})} < \gamma_0$ for $t \in [t_0, T]$. By Lemma 4.5 and (4.12), it holds that for $t \in [0, T]$

$$(4.15) \quad \bar{c}\|\rho(t)\|_{H^1(\mathcal{I})} \leq L[\Gamma_t] - L^* < \delta_0.$$

Then, by the embedding inequality and (4.15), we see that $\|\rho(t)\|_{C^0(\mathcal{I})} \leq C\delta_0$ for $t \in [0, T]$. On the other hand, it follows from Lemma 4.3(ii) and Lemma 4.6(ii) that there exists $C > 0$ such that for $t \in [t_0, T]$

$$(4.16) \quad \|\kappa(t)\|_{C^0(\Gamma_t)} \leq \|\kappa(t) - \kappa_{av}(t)\|_{C^0(\Gamma_t)} + |\kappa_{av}(t) - \kappa^*| + |\kappa^*| \leq C(\delta_1 + \delta_0) + |\kappa^*|.$$

Thus, by virtue of (4.15), (4.16), and $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ for $t \in [t_0, T]$, we derive the boundedness of $\|\rho(t)\|_{H^3(\mathcal{I})}$ for $t \in [t_0, T]$, which implies the boundedness of $\|\rho(t)\|_{C^{2+\alpha}(\mathcal{I})}$ for $\alpha \in (0, 1/2)$. Then, by the interpolation inequality for Hölder spaces, we have

$$\|\rho(t)\|_{C^1(\mathcal{I})} \leq C(\|\rho(t)\|_{C^0(\mathcal{I})})^{\frac{1+\alpha}{2+\alpha}} (\|\rho(t)\|_{C^{2+\alpha}(\mathcal{I})})^{\frac{1}{2+\alpha}} \leq \tilde{C}\delta_0^{\frac{1+\alpha}{2+\alpha}}$$

for $t \in [t_0, T]$. Choosing δ_0 such that $\tilde{C}\delta_0^{\frac{1+\alpha}{2+\alpha}} < \gamma_0$, we obtain $\|\rho(t)\|_{C^1(\mathcal{I})} < \gamma_0$ for $t \in [t_0, T]$.

Finally, let us derive the existence of a unique global-in-time solution. Repeating the above argument until the local existence time T_0 , we see that Γ_t satisfies

$$(4.17) \quad \|\rho(t)\|_{C^1(\mathcal{I})} < \gamma_0, \quad L[\Gamma_t] - L^* < \delta_0, \quad \|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$$

for $t \in [t_0, T_0]$. This implies that Γ_{T_0} satisfies the same conditions as those fulfilled by Γ_0 and the boundedness of $\|\rho(T_0)\|_{C^{2+\alpha}(\mathcal{I})}$ for $\alpha \in (0, 1/2)$ is guaranteed. Thus, due to Theorem 2.1, the solution of (2.5) can be extended over $t = T_0$ by a fixed amount of time. Further, by applying the same argument as we did in the first half of this proof, we have the estimates (4.17) for each time t in the extended existence interval of the solution. This procedure can be iterated as many times as we want, so that a unique global-in-time solution of (2.5) with $\rho(\cdot, 0) = \rho_0$ can be obtained. \square

5. Stability of stationary curves. The following theorem shows nonlinear stability of the stationary curve Γ^* when the bilinear form I^* is positive.

THEOREM 5.1 (nonlinear stability). *Let the assumption of Theorem 4.8 hold. Then*

$$\|\rho(t)\|_{H^3(\mathcal{I})} \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Proof. We apply a method similar to the one used in [10, Proof of Theorem 6.4]. By Lemma 4.1(i), we see that

$$\int_0^\infty \|\kappa_s(\tau)\|_{L^2(\Gamma_\tau)}^2 d\tau \leq L[\Gamma_0].$$

This implies that for any $\varepsilon \in (0, \delta_1)$ there exists a sufficiently large $t_\varepsilon > 0$ such that

$$\|\kappa_s(t_\varepsilon)\|_{L^2(\Gamma_{t_\varepsilon})}^2 < \varepsilon.$$

According to the proof of Theorem 4.8, it holds that $\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 < \delta_1$ as long as the solution exists. Thus, applying Proposition 4.7 for $t \in [t_\varepsilon, \infty)$, we have

$$\|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 + \mu^* \int_{t_\varepsilon}^t \|V_s(\tau)\|_{L^2(\Gamma_\tau)}^2 d\tau \leq \|\kappa_s(t_\varepsilon)\|_{L^2(\Gamma_{t_\varepsilon})}^2 < \varepsilon.$$

This means that

$$(5.1) \quad \|\kappa_s(t)\|_{L^2(\Gamma_t)}^2 \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

By (5.1) and Lemma 4.3(ii), we also see that

$$(5.2) \quad \|\kappa(\cdot, t) - \kappa_{av}(t)\|_{C^0(\Gamma_t)} \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

On the other hand, by virtue of Lemma 4.5 and Lemma 4.6(i), we obtain the boundedness of $\|\rho(t)\|_{H^1(\mathcal{I})}$. Using Lemma 4.6 and (5.2), we also have the boundedness of $\|\kappa(t)\|_{L^2(\Gamma_t)}$. Then the boundedness of $\|\rho(t)\|_{H^1(\mathcal{I})}$ and $\|\kappa(t)\|_{L^2(\Gamma_t)}$ imply the boundedness of $\|\rho(t)\|_{H^2(\mathcal{I})}$. Since it follows from the boundedness of $\|\rho(t)\|_{H^2(\mathcal{I})}$ and (5.1) that $\|\rho(t)\|_{H^3(\mathcal{I})}$ is bounded, there exists a sequence $\{t_n\}_{n \in \mathbb{N}}$ and $\tilde{\rho}$ such that

$$\rho(t_n) \rightarrow \tilde{\rho} \quad \text{in } C^{2+\alpha}(\mathcal{I}) \quad \text{as } n \rightarrow \infty.$$

By virtue of (5.2), $\tilde{\rho}$ satisfies $\tilde{\kappa} - \tilde{\kappa}_{av} = 0$. The solution of the problem

$$\kappa = \kappa_{av}, \quad \angle(\Gamma, \partial\Omega) = \pi/2, \quad \text{Area}[\Gamma] = \text{Area}[\Gamma^*]$$

is unique in the C^0 -neighborhood of Γ^* and given by $\rho \equiv 0$ (see Theorem 5.2). Since $\tilde{\rho}$ is a solution of this problem, we obtain $\tilde{\rho} \equiv 0$. In particular, we get

$$L[\Gamma_{t_n}] \rightarrow L[\Gamma^*] = L^* \quad \text{as } n \rightarrow \infty.$$

We remark that Γ_{t_n} and Γ^* are the curves described by $\rho = \rho(t_n)$ and $\rho \equiv 0$ for the parameterization in section 2, respectively. Then, by the fact that $L[\Gamma_t]$ decreases in time, we obtain that

$$L[\Gamma_t] \rightarrow L^* \quad \text{as } t \rightarrow \infty.$$

Applying Lemma 4.5, we have

$$\|\rho(t)\|_{H^1(\mathcal{I})}^2 \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Hence, using this fact together with both (5.1) and (5.2), we obtain the desired result. \square

It remains to prove the following result. We refer the reader to Grosse-Brauckmann [14] for a similar proof in the case of a different boundary condition.

THEOREM 5.2. *Let Γ^* be a stationary curve such that the bilinear form I^* is positive and let Γ be a curve described by ρ for the parameterization in section 2. Then there exists a C^2 -neighborhood of Γ^* such that $\rho \equiv 0$ is the unique solution of the problem*

$$(5.3) \quad \kappa = \kappa_{av}, \quad \angle(\Gamma, \partial\Omega) = \pi/2, \quad \text{Area}[\Gamma] = \text{Area}[\Gamma^*].$$

Proof. We use the following implicit function theorem (see Zeidler [24, Theorem 4.B]).

Suppose that the following hold:

- (i) *The mapping $F : U(x_0, y_0) \subset X \times Y \rightarrow Z$ is defined on an open neighborhood $U(x_0, y_0)$ of (x_0, y_0) , and $F(x_0, y_0) = 0$, where X, Y , and Z are Banach spaces over \mathbb{R} .*

(ii) F_y exists as the partial Fréchet derivative on $U(x_0, y_0)$ and

$$F_y(x_0, y_0) : Y \rightarrow Z$$

is bijective.

(iii) F and F_y are continuous at (x_0, y_0) .

Then the following holds true: There exist positive numbers r_0 and r such that, for every $x \in X$ satisfying $\|x - x_0\| < r_0$, there is exactly one $y(x) \in Y$ for which $\|y(x) - y_0\| \leq r$ and $F(x, y(x)) = 0$.

We use this theorem for

$$\begin{aligned} X &:= \{\rho \in C^2(\mathcal{I}) \mid \rho = \text{const.}\}, \\ Y &:= \left\{ \rho \in C^2(\mathcal{I}) \mid \int_{l_-}^{l_+} \rho \, d\sigma = 0 \right\}, \\ Z &:= \left\{ \rho \in C^0(\mathcal{I}) \mid \int_{l_-}^{l_+} \rho \, d\sigma = 0 \right\} \times \mathbb{R}^2 \end{aligned}$$

and

$$F(m, u) := \left(\kappa - \kappa_{av}, \angle(\partial\Omega, \Gamma_t)_+ - \frac{\pi}{2}, \angle(\partial\Omega, \Gamma_t)_- - \frac{\pi}{2} \right),$$

where κ is computed for the curve that we get by taking $\rho = u + m$ in section 2. The expression $\angle(\partial\Omega, \Gamma_t)_\pm$ denotes the angles with the outer boundary at the two boundary points. The derivative $F_u(0, 0)$ is (by a similar computation as in [13]) given by

$$\begin{aligned} &F_u(0, 0)(v) \\ &= \left((\partial_\sigma^2 + \kappa^2)v - \frac{1}{l_+ - l_-} \int_{l_-}^{l_+} (\partial_\sigma^2 + \kappa^2)v \, d\sigma, (\partial_\sigma + h_+)v(l_+), (\partial_\sigma - h_-)v(l_-) \right). \end{aligned}$$

The fact that I^* is positive implies that $F_u(0, 0)$ is invertible (using regularity theory for ODEs). Straightforward computations show that F and F_u are continuous at $(0, 0)$.

Hence, for $m \in X$ small, we find exactly one $u(m)$ such that

$$F(m, u(m)) = 0.$$

Let us define

$$\rho_m = u(m) + m$$

and let Γ_m be a curve described by ρ_m for the parameterization in section 2. Then we have

$$\begin{aligned} \text{Area}[\Gamma_m] &= \text{Area}[\Gamma^*] + \int_{l_-}^{l_+} (u(m) + m) \, d\sigma + \mathcal{O}(\|u(m) + m\|_{C^2(\mathcal{I})}^2) \\ &= \text{Area}[\Gamma^*] + (l_+ - l_-)m + \mathcal{O}(\|u(m) + m\|_{C^2(\mathcal{I})}^2). \end{aligned}$$

This implies that for $m \neq 0$

$$(5.4) \quad |\text{Area}[\Gamma_m] - \text{Area}[\Gamma^*]| \neq 0$$

if $\|(m, u(m))\|_{C^2(\mathcal{I})}$ is small enough. We now represent a solution ρ of (5.3) with $\|\rho\|_{C^2(\mathcal{I})}$ small as $\rho = u + m$, where $u = \rho - \rho_{av}$ and $m = \rho_{av}$ with

$$\rho_{av} = \frac{1}{l_+ - l_-} \int_{l_-}^{l_+} \rho \, d\sigma.$$

Then we see that $F(m, u) = 0$. Due to the area-preserving property and (5.4), we obtain $m = 0$ and $u \equiv 0$, which implies that $\rho \equiv 0$. This proves the theorem. \square

Appendix A. Proof of Theorem 2.1. The problem (2.5) is an initial boundary value problem for a quasi-linear parabolic PDE which has the form

$$(A.1) \quad \begin{cases} \rho_t = -\frac{1}{(J(\rho))^4} \partial_\sigma^4 \rho + a(\rho, \partial_\sigma \rho, \partial_\sigma^2 \rho) \partial_\sigma^3 \rho + f(\rho, \partial_\sigma \rho, \partial_\sigma^2 \rho) & \text{in } \mathcal{Q}_{0,T}, \\ b_1(\rho) \partial_\sigma \rho + g_1(\rho) = 0 & \text{at } \sigma = l_\pm, \\ b_2(\rho, \partial_\sigma \rho) \partial_\sigma^3 \rho + g_2(\rho, \partial_\sigma \rho, \partial_\sigma^2 \rho) = 0 & \text{at } \sigma = l_\pm, \\ \rho|_{t=0} = \rho_0 & \text{in } \mathcal{I}, \end{cases}$$

where a, f, b_i , and g_i ($i = 1, 2$) are smooth functions with respect to $\rho, \partial_\sigma \rho$, and $\partial_\sigma^2 \rho$; and g_i ($i = 1, 2$) satisfy $\|g_1(t)\|_{C^0(\mathcal{I})} = \mathcal{O}(\|\rho(t)\|_{C^0(\mathcal{I})})$ and $\|g_2(t)\|_{C^0(\mathcal{I})} = \mathcal{O}(\|\rho(t)\|_{C^{2+\alpha}(\mathcal{I})})$ when $\|\rho\|_{C^{2+\alpha}(\mathcal{I})} \rightarrow 0$. In order to prove Theorem 2.1, we apply a fixed point argument. Let

$$\mathcal{D} := \{\rho \in \mathcal{Y}(\overline{\mathcal{Q}_{0,T}}) \mid \rho(\cdot, 0) = \rho_0, \|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \leq K\}$$

for positive constants K and T , and define a mapping \mathcal{P} as

$$\mathcal{P} : \mathcal{D} \ni \bar{\rho} \mapsto \rho \in \mathcal{Y}(\overline{\mathcal{Q}_{0,T}}),$$

where ρ is the unique solution of the linearized problem

$$(A.2) \quad \begin{cases} \rho_t = \mathcal{A}\rho + F(\sigma, t) & \text{for } (\sigma, t) \in \mathcal{Q}_{0,T}, \\ \mathcal{B}_1 \rho = G_1(\sigma, t) & \text{at } \sigma = l_\pm, t \in (0, T], \\ \mathcal{B}_2 \rho = G_2(\sigma, t) & \text{at } \sigma = l_\pm, t \in (0, T], \\ \rho(\sigma, 0) = \rho_0 & \text{for } \sigma \in \mathcal{I}. \end{cases}$$

Here the linearized operators $\mathcal{A}, \mathcal{B}_1$, and \mathcal{B}_2 around the initial data $\rho_0 \in C^{2+\alpha}(\mathcal{I})$ are given by

$$\begin{aligned} \mathcal{A} &= -\frac{1}{(J(\rho_0))^4} \partial_\sigma^4 + a(\rho_0, \partial_\sigma \rho_0, \partial_\sigma^2 \rho_0) \partial_\sigma^3, \\ \mathcal{B}_1 &= b_1(\rho_0) \partial_\sigma, \quad \mathcal{B}_2 = b_2(\rho_0, \partial_\sigma \rho_0) \partial_\sigma^3, \end{aligned}$$

and for given $\bar{\rho} \in \mathcal{D}$

$$\begin{aligned} F(\sigma, t) &= -\left\{ \frac{1}{(J(\bar{\rho}))^4} - \frac{1}{(J(\rho_0))^4} \right\} \partial_\sigma^4 \bar{\rho} \\ &\quad + \{a(\bar{\rho}, \partial_\sigma \bar{\rho}, \partial_\sigma^2 \bar{\rho}) - a(\rho_0, \partial_\sigma \rho_0, \partial_\sigma^2 \rho_0)\} \partial_\sigma^3 \bar{\rho} \\ &\quad + f(\bar{\rho}, \partial_\sigma \bar{\rho}, \partial_\sigma^2 \bar{\rho}), \\ G_1(\sigma, t) &= -\{b_1(\bar{\rho}) - b_1(\rho_0)\} \partial_\sigma \bar{\rho} - g_1(\bar{\rho}), \\ G_2(\sigma, t) &= -\{b_2(\bar{\rho}, \partial_\sigma \bar{\rho}) - b_2(\rho_0, \partial_\sigma \rho_0)\} \partial_\sigma^3 \bar{\rho} - g_2(\bar{\rho}, \partial_\sigma \bar{\rho}, \partial_\sigma^2 \bar{\rho}). \end{aligned}$$

The existence of a unique solution for the linearized problem (A.2) in $\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})$ is proved by applying the optimal regularity theory for analytic semigroups to the linearized problem (A.2) (see [18]). If the mapping \mathcal{P} is a contraction on \mathcal{D} for suitable constants K and T depending on $\|\rho_0\|_{C^{2+\alpha}(\mathcal{T})}$, \mathcal{P} has a unique fixed point in \mathcal{D} which is a unique solution of the nonlinear problem (A.1). Thus we show that the mapping \mathcal{P} is a contraction on \mathcal{D} . In order to prove this fact, the following lemma is crucial.

LEMMA A.1. (i) *Assume that $\bar{\rho} \in \mathcal{D}$ and that ρ is a solution of the linearized problem (A.2). Then there exist positive constants M_0 and N such that*

$$\|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \leq M_0 + NT^{\frac{\alpha}{4}}.$$

In particular, M_0 depends on $\|\rho_0\|_{C^{2+\alpha}(\mathcal{T})}$ increasingly, and N depends on K increasingly.

(ii) *Assume that $\bar{\rho}_1, \bar{\rho}_2 \in \mathcal{D}$ and that ρ_1, ρ_2 are solutions of the linearized problem (A.2). Then there exists a positive constant N such that*

$$\|\rho_1 - \rho_2\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \leq NT^{\frac{\alpha}{4}} \|\bar{\rho}_1 - \bar{\rho}_2\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})}.$$

In particular, N depends on K increasingly.

A method to prove this lemma is to use the optimal regularity theory of analytic semigroups as in [18]. We prove this lemma in the next section.

Lemma A.1 implies that if we take

$$K = 2M_0, \quad T_0 = \min\left\{ \left(\frac{K}{2N}\right)^{4/\alpha}, \left(\frac{1}{2N}\right)^{4/\alpha} \right\},$$

it follows that for $T \leq T_0$

$$\|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \leq K, \quad \|\rho_1 - \rho_2\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \leq \frac{1}{2} \|\bar{\rho}_1 - \bar{\rho}_2\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})}.$$

This means that \mathcal{P} maps \mathcal{D} into itself and is a contraction on \mathcal{D} for $T \leq T_0$. Thus the proof of Theorem 2.1 is completed.

Appendix B. Proof of Lemma A.1. We prove only Lemma A.1(i). Applying a similar argument, we can also derive Lemma A.1(ii). It is convenient to introduce the following estimate without proof.

LEMMA B.1 (see [18, section 2]). *For $k \in \mathbb{N}$, $\beta_1, \beta_2 \in (0, 1)$, and a sectorial operator A , there exists a constant $C = C(k, \beta_1, \beta_2, A)$ such that*

$$(B.1) \quad \|t^{k-\beta_1+\beta_2} A^k e^{tA}\|_{L(D_A(\beta_1, \infty), D_A(\beta_2, \infty))} \leq C \quad \text{for } 0 < t \leq 1.$$

The statement also holds for $k = 0$, provided that $\beta_1 \leq \beta_2$.

Define $X := C(\mathcal{I})$ and

$$D(A) := \{u \in C^4(\mathcal{I}) \mid \mathcal{B}_1 u(l_\pm) = \mathcal{B}_2 u(l_\pm) = 0\}.$$

Then $A : X \supset D(A) \ni u \mapsto \mathcal{A}u \in X$ is the realization of \mathcal{A} in X . It is known that A is a sectorial operator in X (see [21]).

Let ρ be a unique solution of the linearized problem (A.2). In order to reduce the inhomogeneous problem to a homogeneous problem at the boundaries, we introduce an auxiliary function ζ defined as

$$\begin{aligned} \zeta(\sigma, t) := & \left\{ \frac{(\sigma - l_-)G_1(l_-, t)}{b_1(\rho_0)|_{\sigma=l_-}} + \frac{(\sigma - l_-)^3 G_2(l_-, t)}{3! b_2(\rho_0, \partial_\sigma \rho_0)|_{\sigma=l_-}} \right\} \eta(\sigma) \\ & + \left\{ \frac{(\sigma - l_+)G_1(l_+, t)}{b_1(\rho_0)|_{\sigma=l_+}} + \frac{(\sigma - l_+)^3 G_2(l_+, t)}{3! b_2(\rho_0, \partial_\sigma \rho_0)|_{\sigma=l_+}} \right\} \hat{\eta}(\sigma), \end{aligned}$$

where $\eta, \hat{\eta} \in C^\infty(\mathcal{I})$ are cut-off functions satisfying

$$\begin{cases} \eta'(\sigma) < 0, & \hat{\eta}'(\sigma) > 0 & \text{for } \sigma \in (l_- + L^*/4, l_+ - L^*/4), \\ \eta(\sigma) \equiv 1, & \hat{\eta}(\sigma) \equiv 0 & \text{for } \sigma \in [l_-, l_- + L^*/4], \\ \eta(\sigma) \equiv 0, & \hat{\eta}(\sigma) \equiv 1 & \text{for } \sigma \in [l_+ - L^*/4, l_+]. \end{cases}$$

Then it follows that $\rho - \zeta$ fulfills homogeneous boundary conditions. Since A is sectorial, we represent $\rho - \zeta$ with the help of a variant of the variation of constants formula and the analytic semigroup e^{tA} . By a simple computation, we obtain for $0 \leq t \leq T$

$$\rho(\cdot, t) = \rho_1(\cdot, t) + \rho_2(\cdot, t) + \rho_3(\cdot, t),$$

where

$$\begin{aligned} \rho_1(\cdot, t) &= e^{tA} \{\rho_0 - \zeta(\cdot, 0)\}, \\ \rho_2(\cdot, t) &= \int_0^t e^{(t-r)A} \{F(\cdot, r) + \mathcal{A}\zeta(\cdot, r)\} dr, \\ \rho_3(\cdot, t) &= -A \int_0^t e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, 0)\} dr + \zeta(\cdot, 0). \end{aligned}$$

Applying the theory of analytic semigroups as in [18], we have (see below)

$$(B.2) \quad \begin{cases} \|\rho_1\|_{\mathcal{Y}(\overline{\mathcal{Q}}_{0,T})} \leq C_0 \|\rho_0 - \zeta(\cdot, 0)\|_{D_A(\frac{2+\alpha}{4}, \infty)}, \\ \|\rho_2\|_{\mathcal{Y}(\overline{\mathcal{Q}}_{0,T})} \leq C_0 \sup_{0 < \delta < T} \delta^{\frac{1}{2}} \sup_{t \in [\delta, T]} \|F(\cdot, t) + \mathcal{A}\zeta(\cdot, t)\|_{D_A(\frac{\alpha}{4}, \infty)}, \\ \|\rho_3\|_{\mathcal{Y}(\overline{\mathcal{Q}}_{0,T})} \leq C_0 + C_{0,K} T^{\frac{1}{4}}. \end{cases}$$

In particular, it is verified that a constant C_0 increases with $\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}$ and that a constant $C_{0,K}$ increases with $\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}$ and K . Once (B.2) is proven, it follows from characterization of interpolation spaces $D_A(\beta, \infty)$ (see, e.g., [1, 18, 19]) and

the definition of F that

$$\begin{aligned} \|\rho\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} &\leq \|\rho_1\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} + \|\rho_2\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} + \|\rho_3\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} \\ &\leq \tilde{C}_0 \|\rho_0 - \zeta(\cdot, 0)\|_{C^{2+\alpha}(\mathcal{I})} \\ &\quad + \tilde{C}_0 \sup_{0 < \delta < T} \delta^{\frac{1}{2}} \sup_{t \in [\delta, T]} \|F(\cdot, t) + \mathcal{A}\zeta(\cdot, t)\|_{C^\alpha(\mathcal{I})} \\ &\quad + \tilde{C}_0 + \tilde{C}_{0,K} T^{\frac{1}{4}} \\ &\leq M_0 + N_{0,K} T^{\frac{\alpha}{4}} + N_{0,K} T^{\frac{1}{4}}, \end{aligned}$$

where \tilde{C}_0 and M_0 depend on $\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}$ increasingly, and $\tilde{C}_{0,K}$ and $N_{0,K}$ depend on $\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}$ and K increasingly. This completes the proof of Lemma A.1(i). Thus we give the proof of (B.2) in detail.

First let us explain about the estimates for ρ_1 and ρ_2 . Using (B.1) with $k = 0$ and $\beta_1 = \beta_2 = (2 + \alpha)/4$ to ρ_1 , and with $k = 1$, $\beta_1 = (2 + \alpha)/4$, and $\beta_2 = \alpha/4$ to $\partial\rho_1/\partial t = A\rho_1$, we are led to the estimate of ρ_1 easily. Since $F + \mathcal{A}\zeta \in L^\infty((0, T]; D_A(\frac{\alpha}{4}, \infty))$, applying the same argument as [18, section 4.3.2] to ρ_2 in $[\varepsilon, T]$ ($\varepsilon \in (0, T)$), we have an estimate for ρ_2 . Let us consider the estimate for ρ_3 . Since ζ is less regular, we cannot derive the desired estimate for ρ_3 if we use only (B.1) to ρ_3 directly. Set

$$(B.3) \quad z(t) = \int_0^t e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, 0)\} dr.$$

Then z satisfies

$$\begin{aligned} \rho_3(\cdot, t) &= -Az(t) + \zeta(\cdot, 0) = -\frac{d}{dt}z(t) + \zeta(\cdot, t), \\ \frac{d}{dt}\rho_3(\cdot, t) &= -A\frac{d}{dt}z(t) = A\{\rho_3(\cdot, t) - \zeta(\cdot, t)\}. \end{aligned}$$

This means that if we obtain the estimates for dz/dt , we have the desired estimates for ρ_3 . In fact, the estimate for $\|\rho_3\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})}$ is given by

$$\begin{aligned} \|\rho_3\|_{\mathcal{Y}(\overline{\mathcal{Q}_{0,T}})} &\leq \|\zeta(\cdot, 0)\|_{C^{2+\alpha}(\overline{\mathcal{Q}_{0,T}})} + \|\zeta(\cdot, t) - \zeta(\cdot, 0)\|_{C^{2+\alpha}(\overline{\mathcal{Q}_{0,T}})} \\ &\quad + \sum_{i=1}^3 \sup_{0 < t < T} t^{\frac{1}{2}} \|\mathcal{A}\zeta(\cdot, t)\|_{C^\alpha(\overline{\mathcal{Q}_{0,T}})} \\ &\quad + \tilde{C} \left(\|\dot{z}(t)\|_{D_A(\frac{2+\alpha}{4}, \infty)} + \sup_{0 < \delta < T} \delta^{\frac{1}{2}} \sup_{t \in [\delta, T]} \|A\dot{z}(t)\|_{D_A(\frac{\alpha}{4}, \infty)} \right). \end{aligned}$$

Here and hereafter we use \dot{z} instead of dz/dt to simplify the notation. For the function z , we have the following estimates.

LEMMA B.2. *Let z be a function represented by (B.3). Then there exists a constant N , which depends on $\|\rho_0\|_{C^{2+\alpha}(\mathcal{I})}$, α , and K , such that*

$$(B.4) \quad \begin{cases} \|\dot{z}(t)\|_{D_A(\frac{2+\alpha}{4}, \infty)} \leq NT^{\frac{1}{4}}, \\ \sup_{0 < \delta < T} \delta^{\frac{1}{2}} \sup_{t \in [\delta, T]} \|A\dot{z}(t)\|_{D_A(\frac{\alpha}{4}, \infty)} \leq NT^{\frac{1}{4}}. \end{cases}$$

Proof. The proof of the first estimate of (B.4) is similar to arguments in [16, Appendix]. We prove only the second estimate of (B.4). For $t \geq \varepsilon$ with $\varepsilon \in (0, T)$, we have

$$\begin{aligned} \dot{z}(t) &= e^{(t-\varepsilon/2)A} \dot{z}(\varepsilon/2) + \int_{\varepsilon/2}^t A e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, t)\} dr \\ &\quad + e^{(t-\varepsilon/2)A} \{\zeta(\cdot, t) - \zeta(\cdot, \varepsilon/2)\}. \end{aligned}$$

This implies that

$$\begin{aligned} \|A\dot{z}(t)\|_{D_A(\frac{\alpha}{4}, \infty)} &\leq \|Ae^{(t-\varepsilon/2)A} \dot{z}(\varepsilon/2)\|_{D_A(\frac{\alpha}{4}, \infty)} \\ &\quad + \left\| \int_{\varepsilon/2}^t A^2 e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, t)\} dr \right\|_{D_A(\frac{\alpha}{4}, \infty)} \\ &\quad + \|Ae^{(t-\varepsilon/2)A} \{\zeta(\cdot, t) - \zeta(\cdot, \varepsilon/2)\}\|_{D_A(\frac{\alpha}{4}, \infty)} \\ &=: I_1(t) + I_2(t) + I_3(t). \end{aligned}$$

Let us first derive the estimate of $I_1(t)$. It follows that for $t \geq \varepsilon$

$$(B.5) \quad I_1(t) \leq C_0(t - \varepsilon/2)^{-\frac{\alpha}{4}} \|A\dot{z}(\varepsilon/2)\| \leq C_0(\varepsilon/2)^{-\frac{\alpha}{4}} \|A\dot{z}(\varepsilon/2)\|.$$

Thus it is necessary to obtain an estimate of $\|A\dot{z}(t)\|$. Since $\dot{z}(0) = 0$, we see that

$$\|A\dot{z}(t)\| \leq \int_0^t \|A^2 e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, t)\}\| dr + \|Ae^{tA} \{\zeta(\cdot, t) - \zeta(\cdot, 0)\}\|.$$

We now recall the definition of ζ . Then we have to estimate each term. We show the estimate only for the term including the function

$$\hat{\zeta}(\sigma, t) := (\sigma - l_-)^3 G_2(l_-, t) \eta(\sigma).$$

The idea for the estimation of the other terms is similar. Set

$$J_1(t) := \int_0^t \|A^2 e^{(t-r)A} \{\hat{\zeta}(\cdot, \sigma) - \hat{\zeta}(\cdot, t)\}\| dr,$$

$$J_2(t) := \|Ae^{tA} \{\hat{\zeta}(\cdot, t) - \hat{\zeta}(\cdot, 0)\}\|.$$

Let us derive the estimate of $J_1(t)$. For $t > r$ we have

$$\begin{aligned} &|G_2(\cdot, t) - G_2(\cdot, r)| \\ &\leq |b_2(\bar{\rho}(\cdot, t), \partial_\sigma \bar{\rho}(\cdot, t)) - b_2(\rho_0, \partial_\sigma \rho_0)| |\partial_\sigma^3 \bar{\rho}(\cdot, t) - \partial_\sigma^3 \bar{\rho}(\cdot, r)| \\ &\quad + |b_2(\bar{\rho}(\cdot, t), \partial_\sigma \bar{\rho}(\cdot, t)) - b_2(\bar{\rho}(\cdot, r), \partial_\sigma \bar{\rho}(\cdot, r))| |\partial_\sigma^3 \bar{\rho}(\cdot, r)| \\ &\quad + |g_2(\bar{\rho}(\cdot, t), \partial_\sigma \bar{\rho}(\cdot, t), \partial_\sigma^2 \bar{\rho}(\cdot, t)) - g_2(\bar{\rho}(\cdot, r), \partial_\sigma \bar{\rho}(\cdot, r), \partial_\sigma^2 \bar{\rho}(\cdot, r))| \\ &\leq C_K \left\{ t^{\frac{1+\alpha}{4}} \cdot r^{-\frac{1}{2}} (t-r)^{\frac{1+\alpha}{4}} + r^{-\frac{1}{2}} (t-r)^{\frac{3+\alpha}{4}} \cdot r^{-\frac{1}{4}} + r^{-\frac{1}{2}} (t-r)^{\frac{2+\alpha}{4}} \right\}. \end{aligned}$$

This fact and characterization of interpolation spaces $D_A(\beta, \infty)$ imply that

$$\begin{aligned}
 J_1(t) &\leq C_0 \int_0^t (t-r)^{\frac{3}{4}-2} \|(\sigma - l_-)^3 \eta\|_{D_A(\frac{3}{4}, \infty)} |G_2(l_-, t) - G_2(l_-, r)| dr \\
 &\leq C_{0,K} \int_0^t (t-r)^{\frac{3}{4}-2} \left\{ t^{\frac{1+\alpha}{4}} \cdot r^{-\frac{1}{2}} (t-r)^{\frac{1+\alpha}{4}} \right. \\
 &\quad \left. + r^{-\frac{1}{2}} (t-r)^{\frac{3+\alpha}{4}} \cdot r^{-\frac{1}{4}} + r^{-\frac{1}{2}} (t-r)^{\frac{2+\alpha}{4}} \right\} dr \\
 &\leq C_{0,K,\alpha} (t^{\frac{1+\alpha}{4}} + t^{\frac{1}{4}} + t^{\frac{1}{4}}) t^{\frac{\alpha}{4}-\frac{1}{2}} \\
 &\leq \tilde{C}_{0,K,\alpha} (t^{\frac{1+\alpha}{4}} + t^{\frac{1}{4}}) t^{\frac{\alpha}{4}-\frac{1}{2}}.
 \end{aligned}$$

Applying the similar argument to $J_2(t)$, we are led to

$$\begin{aligned}
 J_2(t) &\leq C_0 t^{\frac{3}{4}-1} \|(\sigma - l_-)^3 \eta\|_{D_A(\frac{3}{4}, \infty)} |G_2(l_-, t) - G_2(l_-, 0)| \\
 &\leq C_{0,K} t^{\frac{3}{4}-1} (t^{\frac{1+\alpha}{4}} \cdot K t^{-\frac{1}{4}} + t^{\frac{\alpha}{4}}) \\
 &\leq \tilde{C}_{0,K} t^{\frac{1}{4}} \cdot t^{\frac{\alpha}{4}-\frac{1}{2}}.
 \end{aligned}$$

Since the estimates for the other terms are also obtained similarly, we have

$$\|A\dot{z}(t)\| \leq C_{0,K,\alpha} T^{\frac{1}{4}} \cdot t^{\frac{\alpha}{4}-\frac{1}{2}}.$$

It follows from (B.5) that

$$I_1(t) \leq C_{0,K,\alpha} T^{\frac{1}{4}} \cdot (\varepsilon/2)^{-\frac{1}{2}}.$$

Let us derive the estimate for $I_2(t)$. Set

$$w(t) := \int_{\varepsilon/2}^t A^2 e^{(t-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, t)\} dr.$$

In order to obtain the estimate of $\|w\|_{D_A(\frac{\alpha}{4}, \infty)}$, we recall the definition of $\|\cdot\|_{D_A(\frac{\alpha}{4}, \infty)}$. Since the estimate of $\|w\|$ is similar to that of $J_1(t)$, we consider only the estimate of the seminorm. According to the definition, we see that

$$\begin{aligned}
 [w]_{D_A(\frac{\alpha}{4}, \infty)} &= \sup_{0 < \tau < 1} \|\tau^{1-\frac{\alpha}{4}} A e^{\tau A} w\| \\
 &\leq \sup_{0 < \tau < 1} \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t \|A^3 e^{(t+\tau-r)A} \{\zeta(\cdot, r) - \zeta(\cdot, t)\}\| dr.
 \end{aligned}$$

We show the estimate only for the term including $\hat{\zeta}(\sigma, t)$. In fact we obtain

$$\begin{aligned}
& \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t \|A^3 e^{(t+\tau-r)A} \{\hat{\zeta}(\cdot, r) - \hat{\zeta}(\cdot, t)\}\| dr \\
& \leq C_0 \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t (t+\tau-r)^{\frac{3}{4}-3} \|(\sigma-l_-)^3 \eta\|_{D_A(\frac{3}{4}, \infty)} |G_2(l_-, t) - G_2(l_-, r)| dr \\
& \leq C_{0,K} \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t (t+\tau-r)^{\frac{3}{4}-3} \{t^{\frac{1+\alpha}{4}} \cdot (\varepsilon/2)^{-\frac{1}{2}} (t-r)^{\frac{1+\alpha}{4}} \\
& \quad + (\varepsilon/2)^{-\frac{1}{2}} (t-r)^{\frac{3+\alpha}{4}} \cdot r^{-\frac{1}{4}} + (\varepsilon/2)^{-\frac{1}{2}} (t-r)^{\frac{2+\alpha}{4}}\} dr \\
& \leq C_{0,K} \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t (t+\tau-r)^{\frac{\alpha}{4}-2} dr \cdot (t^{\frac{1+\alpha}{4}} + t^{\frac{1}{4}}) \cdot (\varepsilon/2)^{-\frac{1}{2}} \\
& \quad + C_{0,K} \tau^{1-\frac{\alpha}{4}} \int_{\varepsilon/2}^t (t+\tau-r)^{\frac{\alpha}{4}-2} (r-\varepsilon/2)^{-\frac{1}{4}} dr \cdot (t-\varepsilon/2)^{\frac{1}{2}} \cdot (\varepsilon/2)^{-\frac{1}{2}} \\
& \leq C_{0,K,\alpha} \tau^{1-\frac{\alpha}{4}} \cdot \tau^{\frac{\alpha}{4}-1} \{T^{\frac{1}{4}} + (t-\varepsilon/2)^{\frac{1}{4}}\} \cdot (\varepsilon/2)^{-\frac{1}{2}} \\
& \leq C_{0,K,\alpha} T^{\frac{1}{4}} \cdot (\varepsilon/2)^{-\frac{1}{2}}.
\end{aligned}$$

As a consequence, we are led to

$$I_2(t) \leq C_{0,K,\alpha} T^{\frac{1}{4}} \cdot (\varepsilon/2)^{-\frac{1}{2}}.$$

The estimate of $I_3(t)$ is omitted, since we can readily obtain it by using (B.1) together with the estimate of $|G_2(\cdot, t) - G_2(\cdot, r)|$.

Consequently, we have

$$\|A\dot{z}(t)\|_{D_A(\frac{\alpha}{4}, \infty)} \leq C_{0,K,\alpha} T^{\frac{1}{4}} \cdot \varepsilon^{-\frac{1}{2}} \quad \text{for } \varepsilon \leq t \leq T.$$

This completes the proof of the second estimate of (B.4). \square

REFERENCES

- [1] P. ACQUISTAPACE AND B. TERRENI, *Hölder classes with boundary conditions as interpolation spaces*, Math. Z., 195 (1987), pp. 451–471.
- [2] J. W. BARRETT, H. GARCKE, AND R. NÜRNBERG, *A parametric finite element method for fourth order geometric evolution equations*, J. Comput. Phys., 222 (2007), pp. 441–467.
- [3] J. W. BARRETT, H. GARCKE, AND R. NÜRNBERG, *On the variational approximation of combined second and fourth order geometric evolution equations*, SIAM J. Sci. Comput., 29 (2007), pp. 1006–1041.
- [4] L. BRONSARD AND F. REITICH, *On three-phase boundary motion and the singular limit of a vector-valued Ginzburg-Landau equation*, Arch. Rational Mech. Anal., 124 (1993), pp. 355–379.
- [5] J. W. CAHN, C. M. ELLIOTT, AND A. NOVICK-COHEN, *The Cahn-Hilliard equation with a concentration dependent mobility: Motion by minus the Laplacian of the mean curvature*, European J. Appl. Math., 7 (1996), pp. 287–301.
- [6] X. CHEN, *The Hele-Shaw problem and area-preserving curve-shortening motions*, Arch. Rational Mech. Anal., 123 (1993), pp. 117–151.
- [7] F. DAVI AND M. GURTIN, *On the motion of a phase interface by surface diffusion*, Z. Angew. Math. Phys., 41 (1990), pp. 782–811.

- [8] S.-I. EI, M.-H. SATO, AND E. YANAGIDA, *Stability of stationary interfaces with contact angle in a generalized mean curvature flow*, Amer. J. Math., 118 (1996), pp. 653–687.
- [9] S.-I. EI AND E. YANAGIDA, *Stability of stationary interfaces in a generalized mean curvature flow*, J. Fac. Sci. Univ. Tokyo Sect. IA Math., 40 (1993), pp. 651–661.
- [10] C. M. ELLIOTT AND H. GARCKE, *Existence results for diffusive surface motion laws*, Adv. Math. Sci. Appl., 7 (1997), pp. 465–488.
- [11] J. ESCHER, H. GARCKE, AND K. ITO, *Exponential stability for a mirror-symmetric three phase boundary motion by surface diffusion*, Math. Nachr., 257 (2003), pp. 3–15.
- [12] J. ESCHER, U. F. MAYER, AND G. SIMONETT, *The surface diffusion flow for immersed hypersurfaces*, SIAM J. Math. Anal., 29 (1998), pp. 1419–1433.
- [13] H. GARCKE, K. ITO, AND Y. KOHSAKA, *Linearized stability analysis of stationary solutions for surface diffusion with boundary conditions*, SIAM J. Math. Anal., 36 (2005), pp. 1031–1056.
- [14] K. GROSSE-BRAUCKMANN, *Stable constant mean curvature surfaces minimize area*, Pacific J. Math., 175 (1996), pp. 527–534.
- [15] R. IKOTA AND E. YANAGIDA, *A stability criterion for stationary curves to the curvature-driven motion with a triple junction*, Differential Integral Equations, 16 (2003), pp. 707–726.
- [16] K. ITO AND Y. KOHSAKA, *Three phase boundary motion by surface diffusion: Stability of a mirror symmetric stationary solution*, Interfaces Free Bound., 3 (2001), pp. 45–80.
- [17] K. ITO AND Y. KOHSAKA, *Three phase boundary motion by surface diffusion in triangular domain*, Adv. Math. Sci. Appl., 11 (2001), pp. 753–779.
- [18] A. LUNARDI, *Analytic Semigroups and Optimal Regularity in Parabolic Problems*, Birkhäuser, Basel, 1995.
- [19] A. LUNARDI, E. SINISTRARI, AND W. VON WAHL, *A semigroup approach to the time dependent parabolic initial-boundary value problem*, Differential Integral Equations, 5 (1992), pp. 1275–1306.
- [20] W. W. MULLINS, *Theory of thermal grooving*, J. Appl. Phys., 28 (1957), pp. 333–339.
- [21] H. B. STEWART, *Generation of analytic semigroups by strongly elliptic operators under general boundary conditions*, Trans. Amer. Math. Soc., 259 (1980), pp. 299–310.
- [22] J. E. TAYLOR AND J. W. CAHN, *Linking anisotropic sharp and diffuse surface motion laws via gradient flows*, J. Statist. Phys., 77 (1994), pp. 183–197.
- [23] T. VOGEL, *Sufficient conditions for capillary surfaces to be energy minima*, Pacific J. Math., 194 (2000), pp. 469–489.
- [24] E. ZEIDLER, *Nonlinear Functional Analysis and its Applications I*, Springer-Verlag, New York, 1986.

DECAY ASYMPTOTICS OF THE VISCOUS CAMASSA–HOLM EQUATIONS IN THE PLANE*

CLAYTON BJORLAND†

Abstract. We consider the vorticity formulation of the two-dimensional viscous Camassa–Holm equations in the whole space. We establish global existence for solutions corresponding to initial data in L^1 and describe the large time behavior of solutions with sufficiently small and localized initial data. We calculate the rate at which such solutions approach an “unfiltered” Oseen vortex by computing the rate at which the solution of a scaled vorticity problem approaches the solution to a corresponding linearized equation.

Key words. Camassa–Holm, Navier–Stokes- α , vorticity, decay

AMS subject classifications. 35Q35, 76B03

DOI. 10.1137/070684070

1. Introduction. The viscous Camassa–Holm equations (VCHE) are commonly written

$$(1.1) \quad \begin{aligned} v_t + u \cdot \nabla v + v \cdot (\nabla u)^T + \nabla p &= \Delta v, \\ H_\alpha^{-1}(u) &= u - \alpha^2 \Delta u = v, \\ \nabla \cdot u &= 0, \\ v(0) &= v_0. \end{aligned}$$

Here H_α is the Helmholtz operator with constant α defined by solving the PDE $u - \alpha^2 \Delta u = v$. This will be referred to as the *filter* associated with the VCHE. For a derivation of these equations from variational principles see [11] or [15]; for a derivation based on modifying the Navier–Stokes system see [7]. The significance of these equations is in the combination of a close relation to the famous Navier–Stokes equations and easily computable bounds on solutions. In particular, in dimensions 2–4, the VCHE admit smooth global solutions which satisfy a modified Kelvin circulation theorem where circulation is conserved around “loops” moving with the filtered flow. The filter is responsible for the smoothing effect on solutions, and the nonlinear term $v \cdot (\nabla u)^T$ brings the solution into compliance with the Kelvin circulation theorem. These solutions are well suited for numerical and analytic calculations and retain many properties displayed by solutions of the Navier–Stokes equations. For proofs of global existence and uniqueness in dimensions 2–4 see [2], [7], [8], [13], and [15]. Decay of energy and higher norms is considered in [2]. Numerical literature is outside the scope of this paper, and the reader is referred to [12] for a survey of the VCHE role in computational turbulence models and a more complete bibliography.

The relation between the VCHE and the Navier–Stokes equations is particularly visible in the vorticity form of the equations found by taking the curl ($\nabla \times v = \tilde{v}$) of

*Received by the editors March 1, 2007; accepted for publication (in revised form) February 7, 2008; published electronically May 28, 2008. This research was partially supported by the NSF under grant OISE-0630623.

<http://www.siam.org/journals/sima/40-2/68407.html>

†Department of Mathematics, UC Santa Cruz, Santa Cruz, CA 95064 (cbjorland@math.ucsc.edu).

(1.1). In two dimensions the vorticity form is

$$(1.2) \quad \begin{aligned} \tilde{v}_t + u \cdot \nabla \tilde{v} &= \Delta \tilde{v}, \\ B(H_\alpha(\tilde{v})) &= u, \\ \tilde{v}(0) &= \tilde{v}_0. \end{aligned}$$

In the above equation, B represents convolution with the well-known Biot–Savart kernel which reconstructs the velocity from the vorticity; see (2.1). The aim of this paper is to further explore the relationship between solutions of the Navier–Stokes equation and the VCHE by describing the way a solution of (1.2) approaches the fixed point *zero*, i.e., computing the first and second order decay asymptotics for solutions with small initial data.

Similar study of the asymptotic behavior of the two-dimensional (2-D) vorticity equation for the Navier–Stokes equation can be found in [5], [9], and [10]. In [9] the asymptotics are calculated by applying invariant manifold techniques to the semiflow governing the vorticity problem. Their approach is to scale the vorticity problem into coordinates which are particularly well suited for studying the large time behavior of the Navier–Stokes equation and then apply the invariant manifold theorem in [6] to construct an invariant manifold in the phase space of the scaled problem and foliate the phase space locally, near the fixed point *zero*. The manifolds constructed give insight into the behavior of solutions near the fixed point, and, among other results, the authors calculate the asymptotics through the interaction and properties of these manifolds.

The close relation of the Navier–Stokes equations and the VCHE gives hope that a similar program may be carried out for the 2-D VCHE, especially when comparing the vorticity equations. In fact, such attempts are met with resistance from the filter in the VCHE. In a functional setting the filter eases problems by smoothing the solution, but in a dynamical setting such as this the filter adds complication to the problem. In particular, the filter does not scale well with the other parts of the equation, and the resulting nonlinear term has dependence on the scaled time variable not present in the case of the Navier–Stokes equations. This time dependence carries through into the semigroups generated by the scaled equation and complicates the invariant manifold construction. Specifically, the theorem in [6] cannot be applied. It is possible to add another equation to the system to express it in an autonomous form (does not depend explicitly on the scaled time variable), but the semigroup generated is still dependent on time and does not commute ($\Theta(s+t) \neq \Theta(s)\Theta(t)$). It is not immediately clear how to construct invariant manifolds with such a time dependence or even if such a structure exists.

The invariant manifold theorem in [6] is based on solving Lyapunov–Perron-type equations which are generated through recursive application of the semiflow. To work around the time dependence of this system we construct an infinite family of systems by stepping the original system forward in time a fixed length and work with the corresponding semigroups. These semigroups do not commute but can be composed to reconstruct the flow of the solution. Following in spirit [6] the semiflow generated by these systems is decomposed into a linear term, a nonlinear term for which we can find uniform Lipschitz bounds, and a forcing function which decays sufficiently fast. Applying this decomposition recursively we find a discrete Lyapunov–Perron-type system which is solvable in a rapidly decaying space. The existence of a solution to this discrete system implies decay properties of the difference system which in turn allow

us to compute the asymptotics we desire. Although the notion of invariant manifolds is lost we still retain enough structure to complete the asymptotic calculations. In essence we rework part of [6] with weaker hypotheses reflecting our situation. To calculate the decay of solutions we linearize the scaled VCHE around a fixed point sufficiently close to zero and determined by the initial data considered, subtract the linearized equation from the scaled one to get a system which measures distance from the linear solution, and apply the newly constructed decay theorem. This procedure, in theory, can be continued to arbitrary orders of asymptotics by considering sufficiently localized data. We consider only the cases where the governing ODEs are linear.

The work is based in large part on the work in [6] and [9]; this is reflected in the notation we use and the statements of many theorems. Following this introduction we introduce the majority of our notation and a few useful lemmas relating to the VCHE. Section 3 is dedicated to proving the existence and uniqueness of solutions to the vorticity problem in two dimensions. This section is somewhat based on the work in [1], where similar results were proved for the Navier–Stokes equations, but for us the work is simplified thanks to the smoothing properties of the filter in functional settings. In section 4 we introduce the scaled variables, prove bounds for solutions with these variables in weighted spaces, and discuss other properties of the scaled system. In this section we also discuss the linearized system and the action of the linear operator \mathcal{L} (the scaled form of Δ) on the weighted spaces. Section 5 contains the decay theorems based on the invariant manifold structure. Sections 6 and 7 hold the theorems and computations involving the first and second order asymptotic, respectively. The theorems in section 7 are very similar to section 6, and many are stated with little or no proof when they have nearly identical analogues in section 6. To end this introduction we state the main conclusions of this paper; proofs are contained in sections 3, 6, and 7, respectively. The spaces $L^2(m)$ are weighted spaces; see (2.2) below.

THEOREM 1.1. *Given initial conditions $\tilde{v}_0 \in L^1(\mathbb{R}^2)$, there exists a unique global solution $\tilde{v} \in L^\infty([0, \infty); L^1(\mathbb{R}^2))$ to the PDE (1.2). This solution satisfies the decay bound*

$$|\tilde{v}(t)|_q \leq Ct^{-(1-\frac{1}{q})}|\tilde{v}(0)|_1.$$

THEOREM 1.2. *For any $\mu \in (0, \frac{1}{2})$, there exists a $r_0 > 0$ so that for any initial data $\tilde{v}_0 \in L^2(2)$ with $\|\tilde{v}_0\|_2 \leq r_0$ the solution of (1.2) satisfies*

$$|\tilde{v}(\cdot, t) - a(\Omega(\cdot, t))|_p \leq C(1+t)^{-1-\mu+\frac{1}{p}},$$

where $a = \int_{\mathbb{R}^2} \tilde{v}_0 dx$ and

$$\Omega(x, t) = \frac{1}{4\pi(1+t)} e^{\frac{-|x|^2}{4(1+t)}}.$$

THEOREM 1.3. *For any $\mu \in (\frac{1}{2}, 1)$, there exists a $r_0 > 0$ so that for any initial data $\tilde{v}_0 \in L^2(3)$ with $\|\tilde{v}_0\|_3 \leq r_0$ the solution of (1.2) satisfies*

$$|\tilde{v}(\cdot, t) - a(\Omega(\cdot, t) - \sum_{i=1,2} b_i(\partial_i \Omega(x, t)))|_p \leq C(1+t)^{-1-\mu+\frac{1}{p}},$$

where $b_i = \int x_i \tilde{v}_0 dx$ and a, Ω are as in the previous theorem.

2. Notation and preliminaries. Throughout this paper we use \mathbb{N} to refer to the natural numbers including 0. Standard Lebesgue spaces will be denoted $L^p(\mathbb{R}^2)$ or $(L^p$ for short) with the norm $|\cdot|_p = (\int |\cdot|^p dx)^{1/p}$. Other Banach spaces defined within will use the norm $\|\cdot\|_X$.

To denote the curl of a vector field we use the tilde. For example, the curl of a vector field v is given by $\nabla \times v = \tilde{v}$. For a divergence-free vector field the curl can be undone through convolution with the well-known Biot–Savart kernel $x^\perp / (2\pi|x|^2)$, $x^\perp = (-x_2, x_1)^T$. We denote this convolution as

$$(2.1) \quad B(\tilde{w}) = \frac{1}{2\pi} \int_{\mathbb{R}^2} \frac{(x-y)^\perp}{|x-y|^2} \tilde{w}(y) dy$$

so that $B(\tilde{v}) = v$ for divergence-free vector fields.

We define $\langle \xi \rangle = (1 + |\xi|^2)^{1/2}$ and make frequent use of the weighted Hilbert spaces $L^2(m)$ defined by

$$(2.2) \quad L^2(m) = \{f \in L^2(\mathbb{R}^2) \mid \|f\|_m < \infty \text{ and } \nabla \cdot f = 0\},$$

$$\|f\|_m^2 = \int_{\mathbb{R}^2} \langle \xi \rangle^{2m}(\xi) |f(\xi)|^2 d\xi.$$

The remaining part of this section contains bounds that will be useful later in the paper. To begin we recall a lemma from [9] concerning the Biot–Savart operator B defined by (2.1) and the curl of divergence-free vector fields.

LEMMA 2.1. *Let*

$$v = B(\tilde{v}) = \frac{1}{2\pi} \int_{\mathbb{R}^2} \frac{(x-y)^\perp}{|x-y|^2} \tilde{v}(y) dy.$$

(i) *If $1 < p < 2 < q < \infty$, $\frac{1}{q} = \frac{1}{p} - \frac{1}{2}$, and $\tilde{v} \in L^p(\mathbb{R}^2)$, then there exists a $C > 0$ such that $|v|_q \leq C|\tilde{v}|_p$.*

(ii) *If $1 \leq p < 2 < q \leq \infty$, $\frac{1}{2} = \frac{\alpha}{p} + \frac{(1-\alpha)}{q}$, where $\alpha \in (0, 1)$, and $\tilde{v} \in L^p \cap L^q(\mathbb{R}^2)$, then there exists a $C > 0$ such that $|v|_\infty \leq C|\tilde{v}|_p^\alpha |\tilde{v}|_q^{1-\alpha}$.*

(iii) *If $1 < p < \infty$ and $\tilde{v} \in L^p(\mathbb{R}^2)$, then there exists a $C > 0$ such that $|\nabla v|_p \leq C|\tilde{v}|_p$.*

Proof. This is Lemma 2.1 in [9]. □

The following bounds for the heat kernel will also be useful in proving the existence of solutions in section 3.

LEMMA 2.2. *Let Φ be the fundamental solution to the heat equation. Then*

$$|\Phi(t)|_p \leq \frac{C(n, p)}{t^{(1-\frac{1}{p})\frac{n}{2}}},$$

$$|\nabla \Phi(t)|_p \leq \frac{C(n, p)}{t^{\frac{1}{2} + (1-\frac{1}{p})\frac{n}{2}}}.$$

Proof. By direct calculation

$$|\Phi(t)|_p^p = \frac{1}{(4\pi t)^{\frac{pn}{2}}} \int_{\mathbb{R}^n} e^{-\frac{p|x|^2}{4t}} dx = \frac{1}{p^{\frac{n}{2}} (4\pi t)^{\frac{(p-1)n}{2}}}.$$

This proves the first bound. For the second, start by differentiating the heat kernel and then take the L^p norm

$$|\nabla\Phi(t)|_p^p \leq \frac{C(n,p)}{t^{p(1+\frac{n}{2})}} \int_{\mathbb{R}^n} |x|^p e^{-\frac{p|x|^2}{4t}} dx.$$

A change of variables and integration now proves the second bound. □

We also recall some facts about the Helmholtz operator H_α .

LEMMA 2.3. *Let $v \in L^r(\mathbb{R}^n)$, $r \in [1, \infty)$. Then there exists a solution $u \in W^{1,r}(\mathbb{R}^n)$ to the Helmholtz equation*

$$H_\alpha^{-1}(u) = u - \alpha^2 \Delta u = v$$

which satisfies the bounds

$$\begin{aligned} |u|_p &\leq |v|_p \text{ for } p \in [1, \infty], \\ |u|_p &\leq \frac{C(n,p,r)}{\alpha^{1+\gamma}} |v|_r \text{ for } \gamma = \frac{n}{2} \left(\frac{1}{r} - \frac{1}{p} \right) < 1, \\ |\nabla u|_p &\leq \frac{C(n,p,r)}{\alpha^{\frac{3}{2}+\gamma}} |v|_r \text{ for } \gamma = \frac{n}{2} \left(\frac{1}{r} - \frac{1}{p} \right) < \frac{1}{2}. \end{aligned}$$

Proof. This follows by standard elliptic theory. □

The last lemma in this section concerns the inclusion $L^2(m) \hookrightarrow L^1$ and will be useful in proving many estimates later in this paper.

LEMMA 2.4. *If $f \in L^2(m)$ for some $m > 1$, then $f \in L^q$ for $1 \leq q \leq 2$.*

Proof. The bound $\|f\|_2 \leq \|f\|_m$ is immediate from the definition of $L^2(m)$. Since $\langle \xi \rangle^{-m}$ is integrable in \mathbb{R}^2 if $m > 1$ we have $\|f\|_1 \leq C\|f\|_m$. Interpolation finishes the proof. □

3. Vorticity problem for the VCHE. This section contains proofs for the existence and uniqueness of solutions to the vorticity form of the VCHE (1.2) with data in $\tilde{v} \in L^1(\mathbb{R}^2)$. Decay of these solutions is provided by the optimal smoothing results in [4].

Instead of working directly with (1.2) we prove existence and uniqueness for the mild form of the PDE by solving the integral equation

$$(3.1) \quad \tilde{v}(t) = e^{\Delta t} \tilde{v}_0 - \int_0^t \nabla\Phi(t-s) * [B(H_\alpha(\tilde{v})) \otimes \tilde{v}](s) ds.$$

The corresponding result for the Navier–Stokes equation was proved in [1] using a fixed point argument in a subspace of L^1 and then extending the solution operator to L^1 ; see also [3] and [14]. This approach was necessary because it is difficult to write the Navier–Stokes equations such that a fixed point argument can be applied in L^1 . In the case of the VCHE the filter provides enough leverage to apply a fixed point argument directly. We start by proving a bound on the bilinear term and then follow with the existence of a global solution.

LEMMA 3.1. *The bilinear form*

$$b(\tilde{v}, \tilde{w}) : L^\infty([0, T], L^1(\mathbb{R}^2)) \times L^\infty([0, T], L^1(\mathbb{R}^2)) \rightarrow L^\infty([0, T], L^1(\mathbb{R}^2))$$

defined by

$$b(\tilde{v}, \tilde{w}) = \int_0^t \nabla\Phi(t-s) * [B(H_\alpha(\tilde{v})) \otimes \tilde{w}](s) ds$$

satisfies the bound

$$\sup_{t \in [0, T]} |b(\tilde{v}, \tilde{w})(t)|_1 \leq C(T) \left(\sup_{t \in [0, T]} |\tilde{v}(t)|_1 \right) \left(\sup_{t \in [0, T]} |\tilde{w}(t)|_1 \right),$$

where $C(T) \rightarrow 0$ as $T \rightarrow 0$.

Proof. First apply Young’s inequality and then Hölder’s inequality to the bilinear form

$$|b(\tilde{v}, \tilde{w})(t)|_1 \leq \int_0^t |\nabla\Phi(t-s)|_1 |B(H_\alpha(\tilde{v}))(s)|_\infty |\tilde{w}(s)|_1 ds.$$

Lemmas 2.1 and 2.3 give the bound $|B(H_\alpha(\tilde{v}))(s)|_\infty \leq C|\tilde{v}(s)|_1$, so

$$|b(\tilde{v}, \tilde{w})(t)|_1 \leq C \int_0^t |\nabla\Phi(t-s)|_1 ds \left(\sup_{s \in [0, t]} |\tilde{v}(s)|_1 \right) \left(\sup_{s \in [0, t]} |\tilde{w}(s)|_1 \right).$$

By Lemma 2.2, $|\nabla\Phi(t-s)|_1 \leq C/\sqrt{t-s}$ and the right-hand side is integrable over finite intervals. Taking the supremum over $t \in [0, T]$ yields

$$\sup_{t \in [0, T]} |b(\tilde{v}, \tilde{w})(t)|_1 \leq CT^{\frac{1}{2}} \left(\sup_{t \in [0, T]} |\tilde{v}(t)|_1 \right) \left(\sup_{t \in [0, T]} |\tilde{w}(t)|_1 \right).$$

This estimate concludes the proof. \square

THEOREM 3.2. *Given initial data $\tilde{v}_0 \in L^1(\mathbb{R}^2)$ there exists a unique global solution $\tilde{v} \in L^\infty([0, \infty); L^1(\mathbb{R}^2))$ to the integral equation (3.1). For any $p \in [1, \infty]$, this solution satisfies the decay bound*

$$(3.2) \quad |\tilde{v}(t)|_p \leq Ct^{-(1-1/p)} |\tilde{v}(0)|_1.$$

Proof. Lemma 3.1 with a standard fixed point argument gives the existence of a mild solution in some possibly small time interval. The length of the time interval depends on the L^1 norm of the initial data. After applying the optimal smoothing results in [4] (see also [16]) the L^1 norm of the solution does not increase beyond the L^1 norm of the data which implies the existence of a global solution. Indeed, the vorticity equation for the VCHE is a viscously damped conservation law, so after applying Theorem 1 in [4] we establish (3.2). In particular, $|v(t)|_1 \leq C|v(0)|_1$, which establishes global existence.

It remains to establish uniqueness. Let \tilde{v} and \tilde{w} be two solutions of (3.1) corresponding to the same initial data $\tilde{v}_0 \in L^1(\mathbb{R}^2)$. After adding and subtracting cross terms we see that

$$\begin{aligned} (\tilde{v} - \tilde{w})(t) &= - \int_0^t \nabla\Phi(t-s) * [B(H(\tilde{v}, \alpha)) \otimes (\tilde{v} - \tilde{w})](s) ds \\ &\quad + \int_0^t \nabla\Phi(t-s) * [B(H((\tilde{v} - \tilde{w}), \alpha)) \otimes \tilde{w}](s) ds. \end{aligned}$$

Similar to the proof of Lemma 3.1, apply Lemmas 2.1 and 2.3 with Young’s inequality:

$$|\tilde{v}(t) - \tilde{w}(t)|_1 \leq C \int_0^t |\nabla\Phi(t-s)|_1 (|\tilde{v}(s)|_1 + |\tilde{w}(s)|_1) |\tilde{v}(s) - \tilde{w}(s)|_1 ds.$$

Apply (3.2) to obtain

$$|\tilde{v}(t) - \tilde{w}(t)|_1 \leq C |\tilde{v}_0|_1 \int_0^t |\nabla\Phi(t-s)|_1 |\tilde{v}(s) - \tilde{w}(s)|_1 ds.$$

From here the Gronwall inequality with Lemma 2.2 is used to establish uniqueness and conclude the proof. \square

This establishes Theorem 1.1.

4. The scaled equations. In this section we introduce scaled variables and rewrite the vorticity equation for the VCHE in these variables, preparing it for use with the theorems in section 5. An existence theorem for the scaled VCHE and related filter equations in the weighted spaces $L^2(m)$ is provided, and we discuss the action of the linear operator \mathcal{L} (scaled Laplacian) on the weighted spaces $L^2(m)$.

The scaled variables are defined as

$$\begin{aligned} \xi &= \frac{x}{\sqrt{1+t}}, & \tau &= \ln(1+t), \\ v(x,t) &= \frac{1}{\sqrt{1+t}} w(\xi,\tau), & u(x,t) &= \frac{1}{\sqrt{1+t}} \omega(\xi,\tau), \\ \tilde{v}(x,t) &= \frac{1}{1+t} \tilde{w}(\xi,\tau), & \tilde{u}(x,t) &= \frac{1}{1+t} \tilde{\omega}(\xi,\tau). \end{aligned}$$

It has been shown in [5], [9], and [10] that these variables are very useful when studying the large time behavior of the Navier–Stokes equation in vorticity form. Under these variables the vorticity form of the VCHE (1.2) becomes

$$(4.1) \quad \tilde{w}_\tau = \mathcal{L}\tilde{w} - \omega \cdot \nabla_\xi \tilde{w}, \quad \tilde{w}(0) = \tilde{w}_0,$$

$$(4.2) \quad \omega = B(\mathcal{H}_{\alpha,\tau}(\tilde{w})),$$

where B is again convolution with the Biot–Savart kernel as in (2.1), \mathcal{L} is the linear operator $\mathcal{L} = \Delta_\xi + \frac{1}{2}\xi \cdot \nabla_\xi + I$, and $\mathcal{H}_{\alpha,\tau}$ is the operator defined by solving the scaled Helmholtz equation $\tilde{\omega} - \alpha^2 e^{-\tau} \Delta_\xi \tilde{\omega} = \tilde{w}$.

The first goal of the section is to show how the filter $\mathcal{H}_{\alpha,\tau}$ acts on the weighted spaces $L^2(m)$; in particular we will show that the Helmholtz equation has a unique solution in these spaces.

LEMMA 4.1. *If $w \in L^2(m)$ and $\omega \in L^2(m)$ are related by the scaled Helmholtz equation $w = \omega - \alpha^2 e^{-\tau} \Delta_\xi \omega$, then $\|\omega(\tau)\|_m^2 \leq C \|w(\tau)\|_m^2$. In the case $4m^2 \alpha^2 < 1$, $C = 1$; otherwise $C = 2(1 + 2^{m-1}(4m^2 \alpha^2)^m)$. This lemma shows how the operator $\mathcal{H}_{\alpha,\tau} : L^2(m) \rightarrow L^2(m)$ is bounded.*

Proof. The case of $m = 0$ follows from the well-known linear elliptic theory; the bound is

$$\|w\|_0^2 = \|\omega\|_0^2 + 2\alpha^2 e^{-\tau} \|\nabla\omega\|_0^2 + \alpha^4 e^{-2\tau} \|\Delta\omega\|_0^2.$$

For $m > 0$ we proceed formally, noting that the following calculations can be applied to a dense set of smooth functions. First square the PDE and then multiply by $\langle \xi \rangle^{2m}$; after integration we have

$$\|w\|_m^2 = \|\omega\|_m^2 + \alpha^4 e^{-2\tau} \|\Delta\omega\|_m^2 - 2\alpha^2 e^{-\tau} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \omega \Delta\omega \, d\xi.$$

Using integration by parts

$$\begin{aligned} \int_{\mathbb{R}^2} (1 + |\xi|^2)^m \omega \Delta\omega \, d\xi &= 2m \int_{\mathbb{R}^2} (\langle \xi \rangle^{2m-2} + (m-1)\langle \xi \rangle^{2m-4} |\xi|^2) \omega^2 \, d\xi \\ &\quad - \|\nabla\omega\|_m^2 \end{aligned}$$

leaves

$$\begin{aligned} \|w\|_m^2 &= \|\omega\|_m^2 + \alpha^4 e^{-2\tau} \|\Delta\omega\|_m^2 + 2\alpha^2 e^{-\tau} \|\nabla\omega\|_m^2 \\ &\quad - 4m\alpha^2 e^{-\tau} \|\omega\|_{m-1}^2 - 4m(m-1)\alpha^2 e^{-\tau} \|\xi|\omega\|_{m-2}^2. \end{aligned}$$

The bound

$$-e^{-\tau} \langle \xi \rangle^{2m-4} (\langle \xi \rangle^2 + (m-1)|\xi|^2) \geq -2m \langle \xi \rangle^{2m}$$

shows that

$$\|w\|_m^2 \geq (1 - 4m^2\alpha^2) \|\omega\|_m^2 + \alpha^4 e^{-2\tau} \|\Delta\omega\|_m^2 + 2\alpha^2 e^{-\tau} \|\nabla\omega\|_m^2$$

and proves the result if $4m^2\alpha^2 < 1$.

If $4m^2\alpha^2 \geq 1$, set $\beta^2 = 8m^2\alpha^2 - 1$ so that for $|\xi| \geq \beta$ we have $\langle \xi \rangle^{2m} \geq 8m^2\alpha^2 \langle \xi \rangle^{2m-2}$. If $B(\beta)$ is the ball with radius β ,

$$\begin{aligned} -4m^2\alpha^2 \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} \omega^2 \, d\xi &\geq -4m^2\alpha^2 (8m^2\alpha^2)^{m-1} \int_{B(\beta)} \omega^2 \, d\xi \\ &\quad - \frac{1}{2} \int_{B^c(\beta)} \langle \xi \rangle^{2m} \omega^2 \, d\xi. \end{aligned}$$

Applying the case $m = 0$ and the bound $\|\tilde{w}\|_2^2 \leq \|\tilde{w}\|_m^2$ allows

$$- \int_{B(\beta)} \omega^2 \, d\xi \geq - \int_{\mathbb{R}^2} \omega^2 \, d\xi \geq - \int_{\mathbb{R}^2} w^2 \, d\xi \geq -\|w\|_m^2.$$

Also,

$$- \frac{1}{2} \int_{B^c(\beta)} \langle \xi \rangle^{2m} \omega^2 \, d\xi \geq -\frac{1}{2} \|\omega\|_m^2.$$

Considering all of this leaves, in the case $4m^2\alpha^2 \geq 1$,

$$(4.3) \quad C\|w\|_m^2 \geq \frac{1}{2} \|\omega\|_m^2 + \alpha^4 e^{-2\tau} \|\Delta\omega\|_m^2 + 2\alpha^2 e^{-\tau} \|\nabla\omega\|_m^2,$$

where $C = 1 + 2^{m-1}(4m^2\alpha^2)^m$. \square

The above proof is some indication that the filter is not well suited for the weighted spaces. There is still a smoothing effect, but as can be seen from (4.3) this effect decreases as τ becomes large.

THEOREM 4.2. *Given $w \in L^2(m)$, there exists a unique solution $\omega \in L^2(m)$ to the scaled Helmholtz equations $w = \omega - \alpha^2 e^{-\tau} \Delta_\xi \omega$. This proves that the operator $\mathcal{H}_{\alpha,\tau} : L^2(m) \rightarrow L^2(m)$ is well defined.*

Proof. The rough estimate $\langle \xi \rangle^{2m} > 1$ implies $L^2(m) \subset L^2(0)$, so if $w \in L^2(m)$, it is well known that there is a unique $\omega \in L^2(0)$ solving the equation. Lemma 4.1 shows that $\omega \in L^2(m)$. \square

We now turn our attention to the scaled VCHE. Using the strongly continuous semigroup $e^{\tau \mathcal{L}}$ generated by \mathcal{L} in $L^2(m)$ (see the appendix of [9]) we write the mild form of the scaled vorticity problem:

$$(4.4) \quad \tilde{w}(\tau) = e^{\tau \mathcal{L}} \tilde{w}_0 - \int_0^\tau e^{-\frac{1}{2}(\tau-s)} \nabla \cdot e^{(\tau-s)\mathcal{L}}(\omega(s)\tilde{w}(s)) ds,$$

$$\tilde{\omega} = \mathcal{H}_{\alpha,\tau}(\tilde{w}).$$

The following lemma, based on Lemma 3.1 in [9], provides an estimate on the bilinear form which will be used to prove the existence of a local solution to (4.4).

LEMMA 4.3. *Given $\tilde{w}_1, \tilde{w}_2 \in C^0([0, T]; L^2(m))$, define*

$$R(\tilde{w}_1, \tilde{w}_2)(\tau) = \int_0^\tau e^{-\frac{1}{2}(\tau-s)} \nabla \cdot e^{(\tau-s)\mathcal{L}}(\omega_1(s)\tilde{w}_2(s)) ds,$$

where $\omega_1 = w_1 - \alpha^2 e^{-\tau} \Delta w_1$ and w_1 is obtained from \tilde{w}_1 via the Biot–Savart law. Then $R \in C^0([0, T], L^2(m))$, and there exists $C_0 = C_0(m, T) > 0$ such that

$$\sup_{0 \leq \tau \leq T} \|R(\tilde{w}_1, \tilde{w}_2)(\tau)\|_m \leq C_0 \left(\sup_{0 \leq \tau \leq T} \|\tilde{w}_1(\tau)\|_m \right) \left(\sup_{0 \leq \tau \leq T} \|\tilde{w}_2(\tau)\|_m \right).$$

Moreover, the constant C_0 becomes arbitrarily small as T tends to zero.

Proof. We rely on an estimate of the semigroup $e^{\tau \mathcal{L}}$ proved in the appendix of [9]; if $r \in (\frac{2}{m+1}, 2)$, then

$$|\langle \xi \rangle^m \nabla \cdot e^{\tau \mathcal{L}} u|_2 \leq \frac{C}{a(\tau)^{(\frac{1}{r}-\frac{1}{2})+\frac{1}{2}}} |\langle \xi \rangle^m u|_r,$$

where $a(\tau) = 1 - e^{-\tau}$. This allows

$$|\langle \xi \rangle^m R(\tilde{w}_1, \tilde{w}_2)(\tau)|_2 \leq C \int_0^\tau \frac{1}{a(\tau-s)^{\frac{1}{r}}} |\langle \xi \rangle^m \omega_1(s)\tilde{w}_2(s)|_r ds.$$

Hölder’s inequality, Lemma 2.1, and the inclusion $L^2(m) \hookrightarrow L^r(\mathbb{R}^2)$ provide the bound $|\langle \xi \rangle^m \omega_1(s)\tilde{w}_2(s)|_r \leq C \|\tilde{w}_2\|_m \|\tilde{\omega}_1\|_m$. Apply Lemma 4.1 and note that $a(\tau-s)^{-1/r}$ is integrable from 0 to τ to finish the proof. \square

In addition to providing existence and uniqueness the following theorem shows how we can control the $L^2(m)$ norm of a solution to the integral equation (4.4) by controlling the $L^2(m)$ norm of the initial data. Our proof is modeled after the proof of Theorem 3.2 in [9].

THEOREM 4.4. *Given $\tilde{w}_0 \in L^2(m)$ for some $m > 1$, there exists a global solution $\tilde{w} \in C^0([0, \infty), L^2(m))$ to the integral equation (4.4) with $\tilde{w}(0) = \tilde{w}_0$. Moreover, there exists a constant $C_1 = C_1(\|\tilde{w}_0\|_m)$ such that*

$$(4.5) \quad \|\tilde{w}(\tau)\|_m \leq C_1$$

and $C_1 \rightarrow 0$ as $\|\tilde{w}_0\|_m \rightarrow 0$.

Proof. The previous lemma and a fixed point argument give local in time existence of a unique solution. Moreover, there exists a $T > 0$ such that

$$(4.6) \quad \sup_{0 \leq \tau \leq T} \|\tilde{w}(\tau)\|_m \leq 2\|\tilde{w}_0\|_m.$$

By scaling (3.2) and using the fact that $L^2(m) \hookrightarrow L^1$ for $m > 1$ we see that this solution satisfies, for all $p \in [1, \infty]$,

$$(4.7) \quad |\tilde{w}(\tau)|_p \leq \frac{C_p \|\tilde{w}(0)\|_m}{a(\tau)^{1-\frac{1}{p}}}.$$

We will now establish (4.5), which will imply global in time existence. Multiplying (4.1) by $\langle \xi \rangle^{2m} \tilde{w}$ and integrating we find that

$$\begin{aligned} \frac{1}{2} \frac{d}{d\tau} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi &= \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} NL(\tilde{w}) d\xi, \\ NL(\tilde{w}) &= \tilde{w} \Delta \tilde{w} + \frac{\tilde{w}}{2} (\xi \cdot \nabla) \tilde{w} + \tilde{w}^2 - \tilde{w}(\omega \cdot \nabla) \tilde{w}. \end{aligned}$$

Integration by parts and the bound $|\xi| \leq \langle \xi \rangle$ give the following estimates:

$$\begin{aligned} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w} \Delta \tilde{w} d\xi &\leq - \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \nabla \tilde{w}^2 d\xi \\ &\quad + 2m^2 \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} \tilde{w}^2 d\xi, \\ \frac{1}{2} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w} (\xi \cdot \nabla) \tilde{w} &= -\frac{1}{2} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi \\ &\quad - \frac{1}{2} m \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} |\xi|^2 \tilde{w}^2 d\xi \\ - \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w} (\omega \cdot \nabla) \tilde{w} d\xi &= m \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} (\xi \cdot \omega) \tilde{w}^2 d\xi. \end{aligned}$$

Furthermore, given $\epsilon > 0$, there exists a $C_\epsilon > 0$ such that

$$\langle \xi \rangle^{2m-2} \leq \epsilon \langle \xi \rangle^{2m} + C_\epsilon$$

and we can bound

$$2m^2 \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} \tilde{w}^2 d\xi \leq \epsilon \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi + C_\epsilon \int_{\mathbb{R}^2} \tilde{w}^2 d\xi.$$

Similarly,

$$m \int_{\mathbb{R}^2} \langle \xi \rangle^{2m-2} (\xi \cdot \omega) \tilde{w}^2 d\xi \leq \epsilon \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi + C_\epsilon |\omega|_\infty^{2m} \int_{\mathbb{R}^2} \tilde{w}^2 d\xi.$$

Putting these bounds together yields

$$\frac{1}{2} \frac{d}{d\tau} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi \leq - \left(\frac{1}{2} - 2\epsilon \right) \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi + C_\epsilon (1 + |\omega|_\infty^{2m}) \int_{\mathbb{R}^2} \tilde{w}^2 d\xi.$$

Write $\delta = \frac{1}{2} - 4\epsilon$; then

$$\frac{d}{d\tau} \left(e^{\delta\tau} \int_{\mathbb{R}^2} \langle \xi \rangle^{2m} \tilde{w}^2 d\xi \right) \leq C_\epsilon e^{\delta\tau} (1 + |\omega|_\infty^{2m}) \int_{\mathbb{R}^2} \tilde{w}^2 d\xi.$$

This inequality implies (4.5), provided

$$(4.8) \quad \sup_{0 \leq \tau \leq \infty} \left((1 + |\omega|_\infty^{2m}) \int_{\mathbb{R}^2} \tilde{w}^2 d\xi \right) \leq C(\|\tilde{w}_0\|_m),$$

where $C(\|\tilde{w}_0\|_m) \rightarrow 0$ as $\|\tilde{w}_0\|_m \rightarrow 0$. We now establish this estimate.

When T is such that (4.6) holds, estimate (4.7) with $p = \infty$ and then $p = 2$ implies (4.8) when the supremum is taken over $\tau \in [T, \infty)$. When $\tau \leq T$, Lemmas 2.1 and 2.3 provide the bound $|\omega(\tau)|_\infty \leq C(T)|\tilde{w}(\tau)|_2$; with (4.6) this implies $(1 + |\omega(\tau)|_\infty) \leq C(T)\|\tilde{w}_0\|_m$ for all $\tau \leq T$. In addition (4.6) implies $|\tilde{w}(\tau)|_2 \leq C(\|\tilde{w}_0\|_m)$ when $\tau \leq T$. This concludes the proof. \square

We spend the remaining portion of this section recalling facts about the operator on the spaces $L^2(m)$ useful to our discussion. The operator \mathcal{L} is studied closely in the appendix of [9], and the reader is referred there for proofs of the statements which are not immediate.

The spectrum of \mathcal{L} in the space $L^2(m)$ is

$$\sigma(\mathcal{L}) = \left\{ \lambda \in \mathbb{C} \mid \operatorname{Re}(\lambda) \leq \frac{1-m}{2} \right\} \cup \left\{ -\frac{k}{2} \mid k \in \mathbb{N} \right\}.$$

The operator \mathcal{L} generates a strongly continuous semigroup $e^{\tau\mathcal{L}}$ on the space $L^2(m)$. The spectrum of \mathcal{L} acting on $L^2(m)$ has $m - 1$ isolated eigenvalues $\lambda_j = \frac{-j}{2}$ for $0 \leq j \leq m - 2$. This allows a spectral decomposition that we will use throughout the remaining parts.

DEFINITION 4.5. Let $X_1^m \subset L^2(m)$ be the finite subspace spanned by the eigenvectors associated with the eigenvalues λ_j , $0 \leq j \leq m - 2$, and $X_2^m = L^2(m) - X_1^m$. Define P_1^m as the spectral projection onto X_1^m and P_2^m the projection onto X_2^m .

Note that P_1^m and P_2^m are guaranteed to exist because X_1^m and X_2^m are closed subspaces of the Hilbert space $L^2(m)$.

LEMMA 4.6. The projections P_1^m and P_2^m defined above satisfy the following bounds for $\tilde{w} \in L^2(m)$:

$$\|(e^{\mathcal{L}} P_1^m)^{-j} \tilde{w}\|_m \leq C_1 e^{\frac{j(m-2)}{2}} \|\tilde{w}\|_m,$$

$$\|(H_{\alpha,\tau} e^{\mathcal{L}} P_2^m)^j \tilde{w}\|_m \leq C_2 e^{\frac{-j(m-1)}{2}} \|\tilde{w}\|_m.$$

Proof. This follows from the definitions of the projections. \square

Of particular interest to us are the first two eigenvalues and the associated eigenvectors.

DEFINITION 4.7. *The first two eigenvalues of \mathcal{L} acting on $L^2(m)$ are $\lambda_0 = 1$ and $\lambda_1 = \frac{-1}{2}$. For the single eigenvector associated with λ_0 we write $G(\xi) = \frac{1}{4\pi} e^{\frac{-|\xi|}{4}}$. For the two eigenvalues associated with λ_2 we write $F_i = \partial_i G(\xi) = -\frac{\xi}{2} G(\xi)$, $i = 1, 2$.*

The eigenvalue G is the scaled ‘‘Oseen vortex’’ and plays an important role in studying the Navier–Stokes system in two dimensions. It is a stationary solution of the scaled PDE, a fact that is easily checked by finding the associated velocity field (v^G) with the Biot–Savart kernel and computing $v^G \cdot \nabla G = 0$. For reference,

$$v^G(\xi) = \frac{1}{2\pi} \frac{1 - e^{\frac{-|\xi|^2}{4}}}{|\xi|^2} |\xi|^{-2} \begin{pmatrix} -\xi_2 \\ \xi_1 \end{pmatrix}.$$

In the scaled VCHE a similar statement is true when we consider G as the filtered and scaled vorticity; this suggests that in describing the behavior of solutions we will also need to consider the following ‘‘unfiltered’’ eigenvectors.

DEFINITION 4.8.

$$\Gamma(\xi, \tau) = G(\xi) - \alpha^2 e^{-\tau} \Delta G(\xi),$$

$$\Lambda_i(\xi, \tau) = F_i(\xi) - \alpha^2 e^{-\tau} \Delta F_i(\xi).$$

Through straightforward calculations one can check that

$$v^G \cdot \nabla \Gamma = 0, \quad \partial_\tau \Gamma = \mathcal{L} \Gamma = \alpha^2 e^{-\tau} \Delta \Gamma.$$

It is then clear that Γ is a solution of the VCHE. This solution is stationary from the perspective of the filtered flow ω and plays a similar role as G in the Navier–Stokes equations.

By linearizing the PDE (4.1) about the ‘‘fixed’’ point $a\Gamma$ ($a \in \mathbb{R}$) we obtain

$$(4.9) \quad \tilde{\psi}_\tau = \mathcal{L} \tilde{\psi} - a\eta \cdot \nabla \Gamma - av^G \cdot \nabla \tilde{\psi}, \quad \tilde{\psi}(0) = \tilde{\psi}_0,$$

$$\eta = B(\mathcal{H}_{\alpha,\tau}(\tilde{\psi})).$$

This linear PDE has strong global solutions, a fact which can be established following the steps in Theorem 4.4; the proof will be omitted here. Note that $a\Gamma$ is a solution to (4.9).

Letting v^{F_i} denote the velocity field associated with F_i , two other important relations are

$$v^{F_i} \cdot \nabla \Gamma + v^G \cdot \nabla \Lambda_i = 0, \quad \partial_t \Lambda_i - \mathcal{L} \Lambda_i + \frac{1}{2} \Lambda_i = 0.$$

The first is quickly checked by differentiating the relation $v^G \cdot \nabla \Gamma = 0$, and the second can be checked directly. Together, these show that $a\Gamma + b_1 e^{\frac{-\tau}{2}} \Lambda_1 + b_2 e^{\frac{-\tau}{2}} \Lambda_2$ is a solution to the linearized equation (4.9).

Subtracting (4.9) from (4.1) we find that

$$(4.10) \quad (\tilde{w} - \tilde{\psi})_\tau = \mathcal{L}(\tilde{w} - \tilde{\psi}) - \omega \cdot \nabla(\tilde{w} - \tilde{\psi}) - (\omega - av^G) \cdot \nabla \tilde{\psi} + a\eta \cdot \nabla \Gamma.$$

This system will be studied in the final sections as the foundation for our asymptotic results.

5. Semigroup theorems. This section contains theorems on collections of semigroups which will be used in the following section to prove our decay results. These theorems are based on the Lyapunov–Perron approach of constructing invariant manifolds in [6] but uses relaxed hypotheses. In particular we allow for the treatment of semigroups which may be time dependent and therefore noncommutative. As our hypotheses are relaxed we are not able to fully construct invariant manifolds. This section instead culminates with a decay theorem.

Unless otherwise specified, in this section let X denote a Hilbert space with norm $\|\cdot\|_X$ and $\{\Theta_n(f, \tau)\}$ a collection of semigroups which map X into itself. That is, each $\Theta_n : X \times \mathbb{R}^+ \rightarrow X$ is a semigroup. The idea we will keep in the back of our mind is a noncommutative semigroup $\Theta(\cdot, t)$, where $\Theta_n(\cdot, \tau)$ is the action of the semigroup when $t = n + \tau$. The semigroup Θ can then be reconstructed by composing the semigroups Θ_n . Indeed, if $t = n + \tau$, then

$$\Theta(f, t) = \Theta_n(\Theta_{n-1}(\cdots \Theta_2(\Theta_1(f, 1), 1), 1), \tau).$$

Conversely we make a definition for a natural flow through the collection.

DEFINITION 5.1. *If $\{\Theta_n(f, \tau)\}$ is a collection of semigroups which, for each $\tau \in [0, 1]$, map X into itself, we call the function $\Theta : X \times \mathbb{R}^+ \rightarrow X$ defined by*

$$\Theta(f, t) = \Theta_n(\Theta_{n-1}(\cdots \Theta_2(\Theta_1(f, 1), 1), 1), \tau),$$

where $t = n + \tau$ and $\tau \in [0, 1]$, the natural flow of the collection through f .

That $\Theta(f, t)$ is well defined for each f follows from the well-defined properties of each Θ_n .

Throughout this section we assume the following uniform conditions on the collection of semigroups; these conditions are relaxations of (H.1)–(H.4) in [6].

(H.1) Define $Lip(\Theta(\cdot, t)) := \sup_n Lip(\Theta_n(\cdot, t))$. Each $\Theta_n(f, t)$ is continuous in $(f, t) \in X \times [0, 1]$ and

$$\sup_{0 \leq t \leq 1} Lip(\Theta(\cdot, t)) = D < \infty.$$

(H.2) Each $\Theta_n(f, 1)$ can be decomposed as $\Theta_n(\cdot, 1) = L + R_n + S_n$, where $L : X \rightarrow X$ is a bounded linear operator, each $R_n : X \rightarrow X$ is a global Lipschitz map satisfying $R_n(0) = 0$, and $S_n \in X$.

(H.3) There are subspaces X_i , $i = 1, 2$, of X and continuous projections P_i , $i = 1, 2$, such that $P_1 + P_2 = I$, $X_1 \oplus X_2 = X$. L leaves X_i invariant and commutes with P_i , $i = 1, 2$. Denoting by $L_i : X_i \rightarrow X_i$ the restriction of L on X_i , L_1 has bounded inverse and there exist constants $\alpha_1 > \alpha_2 \geq 0$, C_1 , and C_2 such that

$$|L_1^{-k} P_1| \leq C_1 \alpha_1^{-k},$$

$$|L_2^k P_1| \leq C_1 \alpha_1^k.$$

We will also write L_i for LP_i when there is no confusion.

(H.4) Define $Lip(R) := \sup_n Lip(R_n)$. Let α_1 and α_2 be as in (H.3); there exists $\alpha_1 > \gamma_1 > \gamma_2 > \alpha_2$ such that the R_n satisfy

$$\frac{C_1}{\alpha_1 - \gamma} + \frac{C_2}{\gamma - \alpha_2} < \frac{1}{Lip(R)}$$

for all $\gamma \in (\gamma_2, \gamma_1)$.

(H.5) With α_1 as in (H.3), S_n satisfies

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \|S_n\|_X < \ln \alpha_1.$$

DEFINITION 5.2. We call a sequence $\{f_n\}_{n \in \mathbb{N}} \subset X$ a discrete positive semiorbit of $\{\Theta_n\}$ through f_0 if it satisfies the relation $f_n = Lf_{n-1} + R_{n-1}(f_{n-1}) + S_{n-1}$ for all $n \geq 1$.

For a given initial position $f_0 \in X$ we find the discrete positive semiorbit $\{f_n\}_{n \in \mathbb{N}}$ it generates using this definition recursively:

$$(5.1) \quad f_n = L^n f_0 + \sum_{j=0}^{n-1} L^{n-j-1} (R_j(f_j) + S_j).$$

We now prove a lemma based on Lemma 3.3 from [6] which will provide a link between the semiflows $\{\Theta_n\}$ and a discrete Lyapunov–Perron system.

LEMMA 5.3. We assume conditions (H.1)–(H.5) are satisfied. Let $\{f_n\}_{n \in \mathbb{N}} \subset X$ satisfy

$$(5.2) \quad \limsup_{n \rightarrow \infty} \frac{1}{n} \ln \|f_n\|_X < \ln \alpha_1.$$

Then the sequence $\{f_n\}_{n \in \mathbb{N}}$ is a positive semiorbit of $\{\Theta_n\}$ if and only if it satisfies, for all $n \in \mathbb{N}$,

$$(5.3) \quad f_n = L_2^n f_0 - \sum_{n \leq j} L_1^{(n-j-1)} (R_j(f_j) + S_j) + \sum_{0 \leq j < n} L_2^{(n-j-1)} (R_j(f_j) + S_j).$$

Proof. Using the iterative relation (5.1) with the projection P_2 shows that

$$P_2 f_n = L_2^n f_0 + \sum_{0 \leq j < n} L_2^{(n-j-1)} (R_j(f_j) + S_j).$$

Likewise, if $m > n$, we can write (5.1) in the form

$$f_m = L^{n-m} f_n + \sum_{n \leq j < m} L_1^{(m-j-1)} (R_j(f_j) + S_j)$$

so that

$$P_1 f_n = L_1^{n-m} f_m - \sum_{n \leq j < m} L_1^{(n-j-1)} (R_j(f_j) + S_j).$$

We will show that the right-hand side of the above equation converges as $m \rightarrow \infty$; then it remains only to add these projections together to finish the proof. Condition (H.3) and the Lipschitz property of R_n allow

$$\left\| \sum_{n \leq j < m} L_1^{(n-j-1)} (R_j(f_j) + S_j) \right\|_X \leq C_1 \sum_{n \leq j < m} \alpha_1^{n-j-1} (\text{Lip}(R) \|f_j\|_X + \|S_j\|_X).$$

The bound (5.2) and assumption (H.5) give convergence of this sum as $m \rightarrow \infty$ as well as the following limit:

$$\begin{aligned} \limsup_{m \rightarrow \infty} \|L_1^{n-m} f_m\|_X &\leq C_1 \alpha_1^n \limsup_{m \rightarrow \infty} (\alpha_1^{-m} \|f_m\|_X) \\ &= 0. \quad \square \end{aligned}$$

The next step is to show that for any initial data f_0 there exists a unique solution to the system (5.3) in a weighted space where all elements satisfy (5.2); then through the previous lemma we can deduce that this solution is a discrete positive semiorbit. Systems such as (5.3) are called Lyapunov–Perron equations and used in the Lyapunov–Perron approach to invariant manifolds. The existence theorem we prove uses a fixed point argument and is similar to Theorem 2.1 in [6].

DEFINITION 5.4. Fix $\gamma \in (\gamma_2, \gamma_1)$, where γ_1 and γ_2 are as in (H.4), and let E_n^γ be the Banach space equal to X as a vector space but equipped with the norm $\|\cdot\|_{E_n^\gamma} = \gamma^{-n} \|\cdot\|_X$. E^γ is the sequence space $f = \{f_n\}_{n \in \mathbb{N}}$, $f_n \in E_n^\gamma$ equipped with the norm $\|f\|_{E^\gamma} = \sup_{n \in \mathbb{N}} \|f_n\|_{E_n^\gamma}$.

As $\gamma_1 < \alpha_1$ by assumption this definition is compatible with (5.2).

THEOREM 5.5. We assume conditions (H.1)–(H.5) are satisfied. Pick $\gamma \in (\gamma_2, \gamma_1)$. Given initial data $x \in X_2$, there exists a unique solution $\{f_n\} \in E^\mu$ to (5.3).

Proof. Define $J : E^\gamma \rightarrow E^\gamma$ componentwise as

$$J_n(\{f_n\}) = L_2^n x + \sum_{0 \leq j < n} L_2^{(n-1-j)}(R_j(f_j) + S_j) - \sum_{n \leq j} L_1^{(n-1-j)}(R_j(f_j) + S_j).$$

To apply a fixed point theorem we need to check that J is well defined and is a contraction map. In that direction we start with the bound

$$\begin{aligned} \|J_n(\{f_n\})\|_{E_n^\gamma} &\leq \|L_2^n x\|_{E_n^\gamma} + \sum_{0 \leq j < n} \|L_2^{(n-1-j)}(R_j(f_j) + S_j)\|_{E_n^\gamma} \\ &\quad - \sum_{n \leq j} \|L_1^{(n-1-j)}(R_j(f_j) + S_j)\|_{E_n^\gamma}. \end{aligned}$$

The first term on the right-hand side is bounded with the help of (H.3), the definition of the space E_n^γ , and the condition $\gamma > \alpha_2$:

$$\begin{aligned} \|L_2^n x\|_{E_n^\gamma} &\leq C_2 \alpha_2^n \gamma^{-n} \|x\|_X \\ &\leq C_2 \|x\|_X. \end{aligned}$$

Assumption (H.5) and $\gamma < \alpha_1$ imply $\{S_n\} \in E^\gamma$; then the second term is bounded using the same ideas but adding the uniform Lipschitz constant of R_j :

$$\begin{aligned} \sum_{0 \leq j < n} \|L_2^{(n-1-j)}(R_j(f_j) + S_j)\|_{E_n^\gamma} &\leq C_2 \alpha_2^{-1} \sum_{0 \leq j < n} (\alpha_2/\gamma)^{n-j} \|(R_j(f_j) + S_j)\|_{E_j^\gamma} \\ &\leq \frac{C_2}{\gamma - \alpha_2} (\text{Lip}(R)\|\{f_n\}\|_{E^\gamma} + \|\{S_n\}\|_{E^\gamma}). \end{aligned}$$

Similarly,

$$\sum_{n \leq j} \|L_1^{(n-1-j)}(R_j(f_j) + S_j)\|_{E_n^\gamma} \leq \frac{C_1}{\alpha_1 - \gamma} (\text{Lip}(R)\|\{f_n\}\|_{E^\gamma} + \|\{S_n\}\|_{E^\gamma}).$$

These bounds show that $\|J_n(\{f_n\})\|_{E^\gamma}$ is bounded independent of n and therefore J maps E^γ into itself. Given another sequence $g_n \in E^\gamma$, following nearly the same steps one finds the bound

$$\|J_n(\{f_n\}) - J_n(\{g_n\})\|_{E^\gamma} \leq \left(\frac{C_1}{\alpha_1 - \gamma} + \frac{C_2}{\gamma - \alpha_2} \right) Lip(R) \|\{f_n - g_n\}\|_{E^\gamma}.$$

Assumption (H.4) guarantees J is a contraction map, and a standard fixed point argument finishes the theorem. \square

We would like to remark here that the system (5.3) does not “see” the component of the initial data in X_1 but instead picks out the correct component to form a solution which satisfies (5.2). Thus, by solving the system one creates a map $h : X_2 \rightarrow X_1$ defined $h(x) = P_1 f_0$ such that the natural flow through $x + h(x)$ decays rapidly. This map is important when constructing invariant manifolds and foliations through the Lyapunov–Perron approach (see, for example, [6]). This next corollary makes this remark precise.

COROLLARY 5.6. *Let $\{\Theta_n\}$ be a collection of semigroups satisfying (H.1)–(H.5). For each $x \in X_2$ there exists $h(x) \in X_1$ such that the discrete positive orbit through $x + h(x)$ satisfies (5.2). This defines a map $h : X_2 \rightarrow X_1$.*

Proof. Define h by using x as initial conditions in the previous theorem to find a sequence $\{f_n\}$; then let $h(x) = P_1 f_0$. As the solution is unique so h is well defined, the remaining properties follow quickly from Lemma 5.3. \square

We are now in a position to prove the theorem, which is the goal of this section and will be the basis for decay estimates in the following sections.

THEOREM 5.7. *Let $\{\Theta_n\}$ be a collection of semigroups satisfying (H.1)–(H.5) and h as in the previous corollary. For each $\gamma \in (\gamma_2, \gamma_1)$ and initial condition $x \in X_2$ denote by $\Theta(x + h(x), t)$ the natural flow of $\{\Theta_n\}$ through $x + h(x)$. This satisfies*

$$(5.4) \quad \limsup_{t \rightarrow \infty} \frac{1}{t} \ln \|\Theta(x + h(x), t)\|_X < \gamma.$$

Proof. $\{\Theta(x + h(x), n)\}$ is a discrete positive orbit of $\{\Theta_n\}$ through $x + h(x)$. Using Theorem 5.5 and its corollary we obtain

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \|\Theta(x + h(x), n)\|_X < \gamma.$$

We now apply the Lipschitz property of the semiflows Θ_n given by (H.1). If $n \in \mathbb{N}$ and $0 \leq \sigma < 1$ are such that $t = n + \sigma$,

$$\begin{aligned} \frac{1}{t} \ln \|\Theta(x + h(x), t)\|_X &\leq \frac{1}{t} \ln(D \|\Theta(x + h(x), n)\|_X) \\ &\leq \frac{1}{n} \ln \|\Theta(x + h(x), n)\|_X + \frac{1}{t} \ln D. \end{aligned}$$

Taking the limit superior first as $n \rightarrow \infty$ and then as $t \rightarrow \infty$ in the above expression finishes the proof. \square

6. First order asymptotic behavior. For the first order calculations we work in the space $L^2(2)$ where the operator \mathcal{L} has a single isolated eigenvalue 0 and corresponding eigenvector G (see Definition 4.7). This eigenvector spans the subspace X_1 as in Definition 4.5. The projection onto this subspace is defined as follows: let

$a = \int \tilde{w}_0 d\xi$; then $P_1(\tilde{w}_0) = aG$. Before we proceed it is important to remark that $\int \tilde{w} d\xi$ is a conserved quantity under the flow of (4.9), so $P_1(\tilde{w}(\tau)) = aG$. Indeed, since $\mathcal{L}\tilde{w} = \nabla \cdot (\nabla\tilde{w} + \frac{\xi}{2})$,

$$(6.1) \quad \frac{1}{2} \frac{d}{dt} \int \tilde{w} d\xi = \int \nabla \cdot \left(\nabla\tilde{w} + \frac{\xi}{2} \cdot \tilde{w} - \omega\tilde{w} \right) d\xi = 0.$$

Fix \tilde{w}_0 (we keep this fixed throughout the calculations) and consider the system (4.10) with initial data $\tilde{w}_0 - a\Gamma(\cdot, 0)$, where $a = \int \tilde{w}_0 d\xi$; later we will require $\|\tilde{w}_0\|_m$ to be sufficiently small. With these initial conditions $\tilde{\psi} = a\Gamma$ is a “stationary” solution of the linear system (4.9). Writing $\tilde{f} = \tilde{w} - a\Gamma$, $\phi = B\mathcal{H}_{\alpha,\tau}(\tilde{f})$, (4.10) becomes

$$(6.2) \quad \tilde{f}_\tau = \mathcal{L}\tilde{f} - \omega \cdot \nabla\tilde{f} - a\phi \cdot \nabla\Gamma.$$

The associated integral equation is

$$(6.3) \quad \tilde{f} = e^{\tau\mathcal{L}}P_2(\tilde{w}_0) - \int_0^\tau e^{\frac{-1}{2}(\tau-\sigma)}\nabla \cdot e^{(\tau-\sigma)\mathcal{L}}(\omega\tilde{f} + a\phi\Gamma)(\sigma) d\sigma.$$

Although it is not obvious at a quick glance, the above system changes with τ . The dependence on τ is buried in the ϕ term (in the filter relation) and destroys the commutative property of the generated semiflow. Commutativity is a very useful property when dealing with semiflows, but we are able to proceed in a fashion consistent with the Lyapunov–Perron approach to constructing invariant manifolds using the results of section 5. To make this work we need a collection of semiflows with uniform bounds which can be put together to reconstruct the original flow. The following system takes us in that direction, and we now change our focus to finding properties of \tilde{f}_n which solve the following system, defined for all $n \in \mathbb{N}$:

$$(6.4) \quad \tilde{f}_n = e^{\tau\mathcal{L}}\tilde{f}_{n,0} - \int_0^\tau e^{\frac{-1}{2}(\tau-\sigma)}\nabla \cdot e^{(\tau-\sigma)\mathcal{L}}(\omega(n+\sigma)\tilde{f}_n(\sigma) + a\phi_m(\sigma)\Gamma(n+\sigma)) d\sigma,$$

$$\phi_n(\sigma) = B\mathcal{H}_{\alpha,\sigma+n}(\tilde{f}_n)(\sigma).$$

Note that this system has a global solution; existence can be proved analogously to Theorem 4.4. It should also be noted here that the above system, (6.2), (6.3), and therefore the flow Θ_n depend on \tilde{w}_0 , which we consider fixed when we write down the system (a depends on \tilde{w}_0); this relation is suppressed in the notation. Most solutions of this equation will have no meaning, but by choosing $\tilde{f}_{0,0} = P_2(\tilde{w}_0)$ and then $\tilde{f}_{n,0} = \tilde{f}_{n-1}(1)$ we are able to recover information about the semiflow corresponding to \tilde{w}_0 .

DEFINITION 6.1. *Throughout the remainder of this section let $\Theta_n(\tilde{f}_{n,0}, \tau)$ denote the global solution to the system (6.4) with initial data $\tilde{f}_{n,0}$.*

The parameter n allows us to keep track of our progress in time. For example, if $0 \leq \sigma < 1$ is such that $\tau = n + \sigma$ and $\Theta(\tilde{f}_0, \tau)$ is the semiflow of (6.3), then

$$(6.5) \quad \Theta(\tilde{f}_0, \tau) = \Theta_n(\Theta_{n-1}(\Theta_{n-2}(\cdots\Theta_0(\tilde{f}_0, 1), 1), 1), \sigma).$$

We will now prove uniform properties of these semiflows.

LEMMA 6.2. *The semiflows $\Theta_n(\tilde{f}_0, \tau)$, $n \in \mathbb{N}$, are all C^1 in $L^2(2) \times \mathbb{R}^+$. There exist a constant $r_0 > 0$ (possibly small) and $D > 0$ such that for all $\|\tilde{w}_0\|_2 < r_0$ and $n \in \mathbb{N}$ the flow Θ_n satisfies the following Lipschitz property:*

$$(6.6) \quad \sup_{n \in \mathbb{N}} \sup_{0 \leq \tau < 1} Lip(\Theta_n(\cdot, \tau)) = D < \infty.$$

This bound holds as $r_0 \rightarrow 0$.

Proof. Since the bound obtained in Lemma 4.1 is independent of τ we can prove this lemma for all $n \in \mathbb{N}$ at once. That the semiflow is C^1 is a classical result. Consider the semiflows $\tilde{f}(\tau) = \Theta_n(\tilde{f}_0, \tau)$ and $\tilde{g}(\tau) = \Theta_n(\tilde{g}_0, \tau)$ found from initial data \tilde{f}_0 and \tilde{g}_0 , respectively, and the corresponding filtered flows ϕ and γ . Subtracting we have

$$\|\Theta_n(\tilde{f}, \tau) - \Theta_n(\tilde{g}, \tau)\|_2 \leq \|e^{\tau\mathcal{L}}(\tilde{f}_0 - \tilde{g}_0)\|_2 + I(\tau) + J(\tau),$$

where

$$I(\tau) = \left\| \int_0^\tau e^{-\frac{1}{2}(\tau-\sigma)} \nabla \cdot e^{(\tau-\sigma)\mathcal{L}}(\omega(n+\sigma)(\tilde{f}(\sigma) - \tilde{g}(\sigma))) d\sigma \right\|_2,$$

$$J(\tau) = \left\| \int_0^\tau e^{-\frac{1}{2}(\tau-\sigma)} \nabla \cdot e^{(\tau-\sigma)\mathcal{L}}((\phi(\sigma) - \gamma(s))a\Gamma(n+\sigma)) ds \right\|_2.$$

Similar to the steps in Lemma 4.3 we obtain

$$I(\tau) + J(\tau) \leq C(\tau) \left(\sup_{0 \leq \sigma < \tau} \|\tilde{f}(\sigma) - \tilde{g}(\sigma)\|_2 \right) \cdot \left(\sup_{0 \leq \sigma < \tau} (\|\tilde{w}(n+\sigma)\|_m + a\|\Gamma(n+\sigma)\|_2) \right)$$

with $C(\tau)$ a continuous function and hence bounded on $\tau \in [0, 1]$ by a constant C . The bound (4.5) combined with the definition of a ($a \leq \|\tilde{w}_0\|_2$) allows us to pick $r_0 > 0$ small enough so that, for all $\sigma > 0$,

$$\|\tilde{w}(n+\sigma)\|_2 + a\|\Gamma(n+\sigma)\|_2 < \frac{1}{2C}.$$

After taking the supremum over $\tau \in [0, 1]$ we have

$$\begin{aligned} \sup_{0 \leq \tau < 1} \|\Theta_n(\tilde{f}, \tau) - \Theta_n(\tilde{g}, \tau)\|_2 &\leq \sup_{0 \leq \tau < 1} \|e^{\tau\mathcal{L}}(\tilde{f}_0 - \tilde{g}_0)\|_2 \\ &\leq D\|\tilde{f}_0 - \tilde{g}_0\|_2. \end{aligned}$$

This is the bound (6.6). \square

LEMMA 6.3. *There exists a constant $r_0 > 0$ (possibly small) such that for all $\|\tilde{w}_0\|_2 < r_0$ and $n \in \mathbb{N}$ the flow Θ_n can be decomposed as $\Theta_n(\tilde{f}_0, 1) = e^{\mathcal{L}}\tilde{f}_0 + R_n(\tilde{f}_0)$, where $R_n(\cdot)$ is Lipschitz as a function from $L^2(2)$ to itself. The uniform Lipschitz constant $Lip(R) := \sup_{n \in \mathbb{N}} Lip(R_n(\cdot))$ satisfies the following conditions:*

(i) *For any $\mu \in (0, \frac{1}{2})$, r_0 may be chosen so that, for all $i \in \mathbb{N}$,*

$$(6.7) \quad \frac{C_1}{1 - e^{-\mu}} + \frac{C_2}{e^{-\mu} - e^{-1/2}} < \frac{1}{Lip(R)}.$$

(ii) *This bound holds as $r_0 \rightarrow 0$.*

Proof. As in the previous proof, consider the semiflows $\tilde{f}(\tau) = \Theta_n(\tilde{f}_0, \tau)$ and $\tilde{g}(\tau) = \Theta_n(\tilde{g}_0, \tau)$. Define

$$\begin{aligned} R_n(\tilde{f}_0) &= \Theta_n(\tilde{f}_0, 1) - e^{\mathcal{L}}\tilde{f}_0 \\ &= - \int_0^1 e^{-1/2(1-\sigma)} \nabla \cdot e^{(1-\sigma)\mathcal{L}}(\omega(n+\sigma)\tilde{f}(\sigma) + a\phi(\sigma)\Gamma(n+\sigma)) d\sigma. \end{aligned}$$

As in Lemmas 4.3 or 6.2, after adding and subtracting cross terms,

$$\begin{aligned} \|R_n(\tilde{f}_0) - R_n(\tilde{g}_0)\|_2 &\leq C \left(\sup_{0 \leq \sigma < \tau} \|\tilde{f}(\sigma) - \tilde{g}(\sigma)\|_2 \right) \\ &\quad \cdot \left(\sup_{0 \leq \sigma < \tau} (\|\tilde{w}(n + \sigma)\|_2 + a\|\Gamma(n + \sigma)\|_2) \right). \end{aligned}$$

Appealing to (6.6),

$$\begin{aligned} \|R_n(\tilde{f}_0) - R_n(\tilde{g}_0)\|_2 &\leq CD \left(\sup_{0 \leq \sigma < \tau} \|\tilde{f}_0 - \tilde{g}_0\|_2 \right) \\ &\quad \cdot \left(\sup_{0 \leq \sigma < \tau} (\|\tilde{w}(n + \sigma)\|_2 + a\|\Gamma(n + \sigma)\|_m) \right). \end{aligned}$$

This shows that

$$Lip(R_i(\cdot)) \leq CD \left(\sup_{0 \leq \sigma < \infty} (\|\tilde{w}(\sigma)\|_2 + a\|\Gamma(\sigma)\|_2) \right)$$

and, with the help of (4.5), satisfies the Lipschitz condition. The bound (6.7) is established by taking r_0 sufficiently small. \square

For the remainder of this section we assume that r_0 is small enough so the conclusions of Lemma 6.2 and 6.3 hold.

THEOREM 6.4. *Pick $\mu \in (0, \frac{1}{2})$ and choose $r_0 > 0$ to satisfy the conclusions of Lemmas 6.2 and 6.3. Let P_i^2 and X_i^2 , $i = 1, 2$, be as in Definition 4.5. Given \tilde{w}_0 such that $\|\tilde{w}_0\|_2 \leq r_0$ and initial data $\tilde{f}_0 \in X_2$, there exists a unique global solution $\Theta(\tilde{f}_0, \tau) \in C^0([0, \infty), L^2(2))$ of (6.3). This solution satisfies $P_1\Theta(\tilde{f}_0, \tau) = 0$ and*

$$(6.8) \quad \limsup_{\tau \rightarrow \infty} \frac{1}{t} \ln \|\Theta(\tilde{f}_0, \tau)\|_2 < -\mu.$$

Proof. The existence of a unique global solution can be argued as in Theorem 4.4; again let $\Theta_n(\tilde{f}_{n,0}, \tau)$ denote the global solution with initial data $\tilde{f}_{n,0}$. We now check the assumptions for Theorem 5.7. Lemma 6.2 shows that the collection $\{\Theta_n\}$ satisfies assumption (H.1) from section 6. Assumption (H.3) is satisfied by Lemma 4.6. As the left-hand side of (6.7) is continuous and the inequality is strict we can find a small neighborhood of μ , say $(\gamma_2, \gamma_1) \subset (0, \frac{1}{2})$, such that the inequality is satisfied for all numbers in this neighborhood. Taking $\alpha_1 = \frac{1}{2}$ and $\alpha_2 = 0$ and applying Lemma 6.3 shows that the collection $\{\Theta_n\}$ satisfies assumptions (H.2), (H.4), and (H.5) with $S_n = 0$.

Applying Theorem 5.7 gives the existence of an element $h(\tilde{f}_0) \in X_1$ such that

$$\limsup_{n \rightarrow \infty} \frac{1}{\tau} \ln \|\Theta(y, \tau)\|_2 < -\mu,$$

where $y = \tilde{f}_0 + h(\tilde{f}_0)$ and Θ is constructed as in (6.5); in other words it is the solution of (6.3). To finish the proof we need only argue that $h(\tilde{f}_0) = 0$, which can be inferred from the ‘‘conservation of mass’’ property of (6.2) and the decay implied by the above inequality. Indeed, using (6.2) we see that $\int \tilde{f} d\xi$ is a conserved property:

$$(6.9) \quad \frac{1}{2} \frac{d}{dt} \int \tilde{f} d\xi = \int \nabla \cdot (\nabla \tilde{f} + \xi \tilde{f} - \omega \tilde{f} - a\phi\Gamma) d\xi = 0.$$

The orthogonal relation $X_1^2 \perp X_2^2$ allows (again writing $y = \tilde{f}_0 + h(\tilde{f}_0)$)

$$\|P_1^2\Theta(y, \tau)\|_2 \leq \|\Theta(y, \tau)\|_2$$

and, as $\|P_1\Theta(y, \tau)\|_2$, is constant for all $\tau > 0$; (6.9) implies that $P_1\Theta(y, \tau) = h(\tilde{f}_0) = 0$. \square

We now prove the theorem, which is the goal of this section.

THEOREM 6.5. *Pick $\mu \in (0, \frac{1}{2})$ and choose $r_0 > 0$ to satisfy the conclusions of Lemmas 6.2 and 6.3. Given initial data \tilde{w}_0 such that $\|\tilde{w}_0\|_2 \leq r_0$, the solution $\tilde{w}(\tau)$ of the scaled VCHE given by Theorem 4.4 is subject to the following decay estimate:*

$$\|\tilde{w}(\tau) - a\Gamma(\tau)\|_2 \leq Ce^{-\mu\tau},$$

where $a = \int \tilde{w}_0 d\xi$.

Proof. In the previous theorem take $\tilde{f}_0 = P_2\tilde{w}_0$; then $\tilde{f}(\tau) = \tilde{w}(\tau) - a\Gamma(\tau)$. The decay follows from (6.8). \square

Finally, as a corollary to this theorem we prove the result listed as Theorem 1.2 in the introduction.

COROLLARY 6.6. *For any $\mu \in (0, \frac{1}{2})$, there exists a $r_0 > 0$ so that, for any initial data $\tilde{v}_0 \in L^2(2)$ such that $\|\tilde{v}_0\|_2 \leq r_0$, the solution of (1.2) given by Theorem 3.2 satisfies*

$$|\tilde{v}(\cdot, t) - a(\Omega(\cdot, t) - \alpha^2\Delta\Omega(\cdot, t))|_p \leq C(1+t)^{-1-\mu+\frac{1}{p}},$$

where $a = \int_{\mathbb{R}^2} \tilde{v}_0 dx$ and

$$\Omega(x, t) = \frac{1}{4\pi(1+t)} e^{\frac{-|x|^2}{4(1+t)}}.$$

Proof. This is the result of the previous theorem in unscaled coordinates. Let

$$(6.10) \quad \Omega(x, t) = \frac{1}{(1+t)} G\left(\frac{x}{\sqrt{1+t}}\right) = \frac{1}{4\pi(1+t)} e^{\frac{-|x|^2}{4(1+t)}};$$

then

$$\frac{1}{(1+t)} \Gamma\left(\frac{x}{\sqrt{1+t}}, \ln(1+t)\right) = \Omega(x, t) - \alpha^2\Delta\Omega(x, t).$$

Thanks to the above theorem and the inclusion $L^2(2) \hookrightarrow L^p$ for when $1 \leq p \leq 2$,

$$\begin{aligned} &|\tilde{v}(\cdot, t) - a(\Omega(\cdot, t) - \alpha^2\Delta\Omega(\cdot, t))|_p \\ &\leq (1+t)^{-1+\frac{1}{p}} |\tilde{w}(\cdot, \ln(1+t)) - a\Gamma(\cdot, \ln(1+t))|_p \\ &\leq C(1+t)^{-1+\frac{1}{p}} \|\tilde{w}(\cdot, \ln(1+t)) - a\Gamma(\cdot, \ln(1+t))\|_2 \\ &\leq C(1+t)^{-1-\mu+\frac{1}{p}}. \quad \square \end{aligned}$$

As $\Delta\Omega$ decays faster than $(1+t)^{-1-\mu+1/p}$ we can include it on the right-hand side; this is how Theorem 1.2 is stated.

7. Second order asymptotic. This section is similar in spirit to the previous section, and many of the proofs are omitted because they are nearly the same as the proofs given in the previous section. We work in the space $L^2(3)$ where the operator \mathcal{L} has two isolated eigenvalues, 0 and $\frac{-1}{2}$, and three corresponding eigenvectors, G and F_i , $i = 1, 2$ (see Definition 4.7). Together these eigenvalues span the subspace X_1^3 given by Definition 4.5. Any $\tilde{w}_0 \in L^2(3)$ can be written as $\tilde{w}_0 = aG + b_1F_1 + b_2F_2 + \tilde{g}$, where $a = \int \tilde{w}_0 d\xi$, $b_i = \int \xi_i \tilde{w}_0 d\xi$, and $\tilde{g} \in X_2$.

In addition to the ‘‘conservation of mass’’ property $\frac{d}{dt} \int \tilde{w} d\xi = 0$ used in the previous section (proof of Theorem 5.7), solutions to the scaled VCHE (4.1) also satisfy the scaled form of conservation of the first moments, $\frac{d}{dt} \int \xi_i \tilde{w} d\xi = -\frac{1}{2} \int \xi_i \tilde{w} d\xi$. Indeed, let $i = 1, 2$ and $j \neq i$; then

$$\xi_i \mathcal{L} \tilde{w} + \frac{1}{2} \xi_i \tilde{w} = \partial_j \left(\xi_i \partial_i \tilde{w} + \frac{1}{2} \xi_i^2 \tilde{w} - \tilde{w} \right) + \partial_j \left(\xi_i \partial_j \tilde{w} + \frac{1}{2} \xi_i \xi_j \tilde{w} \right),$$

and it is clear that $\int_{\mathbb{R}^2} \xi_i \mathcal{L} \tilde{w} d\xi = -\frac{1}{2} \int_{\mathbb{R}^2} \xi_i \tilde{w} d\xi$. For the nonlinear term, note that if $\nabla \cdot \omega = 0$ and $\tilde{\omega} = \partial_1 \omega_2 - \partial_2 \omega_1$, then

$$\xi_i \omega \cdot \nabla \tilde{\omega} = \partial_i (\xi_i \omega_i \tilde{\omega} - \omega_i \omega_j) + \partial_j \left(\xi_i \omega_j \tilde{\omega} + \frac{1}{2} (\omega_i^2 - \omega_j^2) \right).$$

Similarly,

$$\begin{aligned} \xi_i \omega \cdot \nabla \partial_i^2 \tilde{\omega} &= \partial_i (\xi_i \omega \cdot \nabla \partial_j \tilde{\omega}) - \omega \cdot \nabla \partial_i \tilde{\omega} - \xi_i \partial_i \omega \cdot \nabla \partial_i \tilde{\omega} \\ &= \partial_i (\xi_i \omega \cdot \nabla \partial_j \tilde{\omega}) - \nabla (\omega \partial_i \tilde{\omega}) - \partial_i (\xi_i \partial_i \omega_i \partial_i \tilde{\omega} - \partial_i \omega_i \partial_i \omega_j) \\ &\quad - \partial_j \left(\xi_i \partial_i \omega_j \partial_i \tilde{\omega} + \frac{1}{2} ((\partial_i \omega_i)^2 - (\partial_i \omega_j)^2) \right) \end{aligned}$$

and

$$\begin{aligned} \xi_i \omega \cdot \nabla \partial_j^2 \tilde{\omega} &= \partial_j (\xi_i \omega \cdot \nabla \partial_j \tilde{\omega}) - \partial_i (\xi_i \partial_j \omega_i \partial_j \tilde{\omega} - \partial_j \omega_i \partial_j \omega_j) \\ &\quad - \partial_j \left(\xi_i \partial_j \omega_j \partial_j \tilde{\omega} + \frac{1}{2} ((\partial_j \omega_i)^2 - (\partial_j \omega_j)^2) \right) \end{aligned}$$

so that $\int_{\mathbb{R}^2} \xi_i \omega \cdot \nabla \tilde{w} d\xi = 0$.

Fix $\tilde{w}_0 \in L^2(m)$ and consider the system (4.10) with initial data $P_2^m \tilde{w}_0$. Let $a = \int \tilde{w}_0 d\xi$, $b_i = \int \xi_i \tilde{w}_0 d\xi$; then $\tilde{\psi}(\tau) = a\Gamma + e^{\frac{-\tau}{2}} (b_1 \Lambda_1 + b_2 \Lambda_2)$ is the solution to the linear equation (4.9). Write $\tilde{f} = \tilde{w} - a\Gamma - e^{\frac{-\tau}{2}} (b_1 \Lambda_1 + b_2 \Lambda_2)$, $\phi = B\mathcal{H}_{\alpha, \tau}(\tilde{f})$, and after simplification,

$$(7.1) \quad \tilde{f}_\tau = \mathcal{L} \tilde{f} - \omega \cdot \nabla \tilde{f} - \phi \cdot \nabla \psi - e^{-\tau} (b_1 v^{F_1} \cdot \nabla \Lambda_1 - b_2 v^{F_2} \cdot \nabla \Lambda_2).$$

In mild form \tilde{f} satisfies the integral equation

$$\begin{aligned} \tilde{f} &= e^{\tau \mathcal{L}} \tilde{f}_0 - \int_0^\tau e^{-1/2(\tau-\sigma)} \nabla \\ &\quad \cdot e^{(\tau-\sigma) \mathcal{L}} (\omega \tilde{f} + \phi \cdot \nabla \psi + e^{-\tau} (b_1 v^{F_1} \cdot \nabla \Lambda_1 - b_2 v^{F_2} \cdot \nabla \Lambda_2))(\sigma) d\sigma. \end{aligned}$$

The key difference between this system and (6.3) of the previous section is the ‘‘forcing term’’ $e^{-\tau} (b_1 v^{F_1} \cdot \nabla \Lambda_1 - b_2 v^{F_2} \cdot \nabla \Lambda_2)$.

As in the previous section fix \tilde{w}_0 and for any $n \in \mathbb{N}$ consider the system

$$(7.2) \quad \begin{aligned} \tilde{f}_n &= e^{\tau \mathcal{L}} \tilde{f}_{0,n} - \int_0^\tau e^{-1/2(\tau-\sigma)} \nabla \cdot e^{(\tau-\sigma)\mathcal{L}} (\omega(n+\sigma) \tilde{f}(\sigma) + a\phi(\sigma)\psi(n+\sigma)) d\sigma \\ &\quad - \int_0^\tau e^{-1/2(\tau-\sigma)} \nabla \cdot e^{(\tau-\sigma)\mathcal{L}} (e^{-(n+\sigma)} (b_1 v^{F_1} \cdot \nabla \Lambda_1 - b_2 v^{F_2} \cdot \nabla \Lambda_2)(n+\sigma)) d\sigma, \\ \phi(\sigma) &= B\mathcal{H}_{\alpha,\sigma+\tau}(\tilde{f})(\sigma), \quad \tilde{f}(0) = \tilde{f}_0 \in X_2. \end{aligned}$$

We redefine Θ_n for this section.

DEFINITION 7.1. Throughout the remainder of this section let $\Theta_n(\tilde{f}_{n,0}, \tau)$ denote the global solution to the system (7.2) with initial data $\tilde{f}_{n,0}$.

LEMMA 7.2. For each $n \in \mathbb{N}$, the semiflows $\Theta_n(\tilde{f}_0, \tau)$ are C^1 in $L^2(2) \times \mathbb{R}^+$. There exist a constant $r_0 > 0$ (possibly small) and $D > 0$ such that for all $\|\tilde{w}_0\|_3 < r_0$ and $n \in \mathbb{N}$ the flow Θ_n satisfies the following Lipschitz property:

$$(7.3) \quad \sup_{n \in \mathbb{N}} \sup_{0 \leq \tau < 1} \text{Lip}(\Theta_n(\cdot)(\tau)) = D < \infty.$$

This bound holds as $r_0 \rightarrow 0$.

Proof. The proof is nearly identical to the proof of Lemma 6.2. \square

LEMMA 7.3. There exists a constant $r_0 > 0$ (possibly small) such that for all $\|\tilde{w}_0\|_3 < r_0$ and $n \in \mathbb{N}$ the flow Θ_n can be decomposed as $\Theta_n(\tilde{f}_0, 1) = e^{\mathcal{L}} \tilde{f}_0 + R_n(\tilde{f}_0) + S_n$, where $R_n(\cdot)$ is Lipschitz as a function from $L^2(3)$ to itself and $S_n \in L^2(3)$ satisfies

$$(7.4) \quad \limsup_{n \rightarrow \infty} \frac{1}{n} \ln \|S_n\|_3 < -1.$$

The Lipschitz constant $\text{Lip}(R) := \sup_{n \in \mathbb{N}} \text{Lip}(R_n(\cdot))$ can be made arbitrarily small and satisfies the following conditions:

(i) For any $\mu \in (\frac{1}{2}, 1)$, r_0 may be chosen so that, for all $n \in \mathbb{N}$,

$$(7.5) \quad \frac{C_1}{e^{-1/2} - e^{-\mu}} + \frac{C_2}{e^{-\mu} - e^{-1}} < \frac{1}{\text{Lip}(R)}.$$

(ii) This bound holds as $r_0 \rightarrow 0$.

Proof. Define

$$S_n = - \int_0^\tau e^{-1/2(\tau-\sigma)} \nabla \cdot e^{(\tau-\sigma)\mathcal{L}} (e^{-(n+\sigma)} (b_1^2 v^{F_1} \cdot \nabla \Lambda_1 - b_2^2 v^{F_2} \cdot \nabla \Lambda_2)(n+\sigma)) d\sigma$$

and $R_n(\tilde{f}_0) = \Theta_n(\tilde{f}_0)(1) - e^{\mathcal{L}} \tilde{f}_0 - S_n$. Proving (7.4) is similar to the proof of Lemma 4.3:

$$(7.6) \quad \|S_n\|_3 \leq e^{-n} (b_1 \|v^{F_1}\|_3 \|\Lambda_1\|_3 + b_2 \|v^{F_2}\|_3 \|\Lambda_2\|_3).$$

The remaining statements in this proof follow as in the proof of Lemma 6.3. \square

THEOREM 7.4. Pick $\mu \in (\frac{1}{2}, 1)$ and choose $r_0 > 0$ to satisfy the conclusions of Lemmas 6.2 and 6.3. Let P_i^3 and X_i^3 , $i = 1, 2$, be as in Definition 4.5. Given \tilde{w}_0 such that $\|\tilde{w}_0\|_3 \leq r_0$ and initial data $\tilde{f}_0 \in X_2$, there exists a unique global solution $\Theta(\tilde{f}_0, \tau) \in C^0([0, \infty), L^2(2))$ of (6.3). This solution satisfies $P_1^3 \Theta(\tilde{f}_0, \tau) = 0$ and

$$(7.7) \quad \limsup_{\tau \rightarrow \infty} \frac{1}{t} \ln \|\Theta(\tilde{f}_0, \tau)\|_2 < -\mu.$$

Proof. The proof progresses similarly to the proof of Theorem 7.4. The main difference is in the argument for the statement $P_1^3\Theta(\tilde{f}_0, \tau) = 0$. This is accomplished by decomposing $P_1\Theta(\tilde{f}_0, \tau)$ onto the first three eigenvectors of \mathcal{L} using conservation of mass (discussed at the beginning of section 6) and conservation of first moments (discussed at the start of section 7):

$$P_1^3\Theta(\tilde{f}_0, \tau) = aG + e^{-\frac{1}{2}\tau}(b_1F_1 + b_2F_2).$$

That a and b_i must be zero follows by comparing this expression to the decay implied by (7.7). \square

THEOREM 7.5. *Pick $\mu \in (\frac{1}{2}, 1)$ and choose $r_0 > 0$ to satisfy the conclusions of Lemmas 7.2 and 7.3. Given initial data \tilde{w}_0 such that $\|\tilde{w}_0\|_3 \leq r_0$, the solution $\tilde{w}(\tau)$ of the scaled VCHE given by Theorem 4.4 is subject to the following decay estimate:*

$$\|\tilde{w}(\tau) - a\Gamma(\tau) - e^{-\frac{\tau}{2}}(b_1\Lambda_1 + b_2\Lambda_2)\|_2 \leq Ce^{-\mu\tau},$$

where $a = \int \tilde{w}_0 d\xi$ and $b_i = \int \xi_i \tilde{w}_0 d\xi$.

Proof. In the previous theorem take $\tilde{f}_0 = P_2\tilde{w}_0$; then $\tilde{f}(\tau) = \tilde{w}(\tau) - a\Gamma(\tau) - e^{-\frac{\tau}{2}}(b_1\Lambda_1 + b_2\Lambda_2)$. The decay then follows from (6.8). \square

The following corollary proves Theorem 1.3, which was stated after removing the fast decaying term.

COROLLARY 7.6. *For any $\mu \in (\frac{1}{2}, 1)$, there exists a $r_0 > 0$ so that for any initial data $\tilde{v}_0 \in L^2(3)$ with $\|\tilde{v}_0\|_3 \leq r_0$ the solution of (1.2) given by Theorem 3.2 satisfies*

$$|\tilde{v}(\cdot, t) - a(\Omega(\cdot, t) - \alpha^2\Delta\Omega(\cdot, t)) - \sum_{i=1,2} b_i(\partial_i\Omega(x, t) - \alpha^2\Delta\partial_i\Omega(x, t))|_p \leq C(1+t)^{-1-\mu+\frac{1}{p}},$$

where $a = \int \tilde{v}_0 dx$ and $b_i = \int x_i \tilde{v}_0 dx$ and Ω is defined by (6.10).

Proof. This is the result of the previous theorem in unscaled coordinates. Let $\Omega(x, t) = \frac{1}{(1+t)}G(\frac{x}{\sqrt{1+t}})$; then

$$\frac{1}{(1+t)}\Gamma\left(\frac{x}{\sqrt{1+t}}, \ln(1+t)\right) = \Omega(x, t) - \alpha^2\Delta\Omega(x, t),$$

$$e^{-\frac{\tau}{2}}\frac{1}{(1+t)}\Lambda_i\left(\frac{x}{\sqrt{1+t}}, \ln(1+t)\right) = \partial_i\Omega(x, t) - \alpha^2\Delta\partial_i\Omega(x, t).$$

The rest follows as in Corollary 6.6. \square

Acknowledgments. The author would like to thank M. E. Schonbek for suggesting the problem and the anonymous referees for insightful suggestions and comments.

REFERENCES

[1] M. BEN-ARTZI, *Global solutions of two-dimensional Navier-Stokes and Euler equations*, Arch. Rational Mech. Anal., 128 (1994), pp. 329–358.
 [2] C. BJORLAND AND M. E. SCHONBEK, *On questions of decay and existence for the viscous Camassa-Holm equations*, Ann. Inst. H. Poincaré Anal. Non Linéaire, to appear.
 [3] H. BREZIS, *Remarks on the preceding paper by M. Ben-Artzi: “Global solutions of two-dimensional Navier-Stokes and Euler equations”* [Arch. Rational Mech. Anal., 128 (1994), pp. 329–358], Arch. Rational Mech. Anal., 128 (1994), pp. 359–360.
 [4] E. A. CARLEN AND M. LOSS, *Optimal smoothing and decay estimates for viscously damped conservation laws, with applications to the 2-D Navier-Stokes equation*, Duke Math. J., 81 (1995), pp. 135–157.

- [5] A. CARPIO, *Asymptotic behavior for the vorticity equations in dimensions two and three*, Comm. Partial Differential Equations, 19 (1994), pp. 827–872.
- [6] X. CHEN, J. K. HALE, AND B. TAN, *Invariant foliations for C^1 semigroups in Banach spaces*, J. Differential Equations, 139 (1997), pp. 283–318.
- [7] C. FOIAS, D. D. HOLM, AND E. S. TITI, *The Navier-Stokes-alpha model of fluid turbulence*, Phys. D, 152/153 (2001), pp. 505–519.
- [8] C. FOIAS, D. D. HOLM, AND E. S. TITI, *The three dimensional viscous Camassa-Holm equations, and their relation to the Navier-Stokes equations and turbulence theory*, J. Dynam. Differential Equations, 14 (2002), pp. 1–35.
- [9] T. GALLAY AND C. E. WAYNE, *Invariant manifolds and the long-time asymptotics of the Navier-Stokes and vorticity equations on \mathbb{R}^2* , Arch. Ration. Mech. Anal., 163 (2002), pp. 209–258.
- [10] Y. GIGA AND T. KAMBE, *Large time behavior of the vorticity of two-dimensional viscous flow and its application to vortex formation*, Comm. Math. Phys., 117 (1988), pp. 549–568.
- [11] D. D. HOLM, J. E. MARSDEN, AND T. S. RATIU, *The Euler-Poincaré equations and semidirect products with applications to continuum theories*, Adv. Math., 137 (1998), pp. 1–81.
- [12] D. D. HOLM AND E. S. TITI, *Computational models of turbulence: The LANS- α model and the role of global analysis*, SIAM News, 38 (2005).
- [13] A. A. ILYIN AND E. S. TITI, *Attractors for the two-dimensional Navier-Stokes- α model: An α -dependence study*, J. Dynam. Differential Equations, 15 (2003), pp. 751–778.
- [14] T. KATO, *The Navier-Stokes equation for an incompressible fluid in \mathbf{R}^2 with a measure as the initial vorticity*, Differential Integral Equations, 7 (1994), pp. 949–966.
- [15] J. E. MARSDEN AND S. SHKOLLER, *Global well-posedness for the Lagrangian averaged Navier-Stokes (LANS- α) equations on bounded domains*, R. Soc. Lond. Philos. Trans. Ser. A Math. Phys. Eng. Sci., 359 (2001), pp. 1449–1468.
- [16] H. OSADA, *Diffusion processes with generators of generalized divergence form*, J. Math. Kyoto Univ., 27 (1987), pp. 597–619.

RIGOROUS DERIVATION OF INCOMPRESSIBLE e-MHD EQUATIONS FROM COMPRESSIBLE EULER–MAXWELL EQUATIONS*

YUE-JUN PENG[†] AND SHU WANG[‡]

Abstract. We derive incompressible e-MHD equations from compressible Euler–Maxwell equations via the quasi-neutral regime. Under the assumption that the initial data are well prepared for the electric density, electric velocity, and magnetic field (but not necessarily for the electric field), the convergence of the solutions of the compressible Euler–Maxwell equations in a torus to the solutions of the incompressible e-MHD equations is justified rigorously by studies on a weighted energy. One of the main ingredients for establishing uniform a priori estimates is to use the curl-div decomposition of the gradient and the wave-type equations of the Maxwell equations.

Key words. Euler–Maxwell equations, incompressible electron magnetohydrodynamics equations, quasi-neutral limit, weighted energy

AMS subject classifications. 35B40, 35C20, 35L60, 35Q35

DOI. 10.1137/070686056

1. Introduction. Let n and u be the density and velocity vector of the electric particles in a plasma and E and B be, respectively, the electric field and magnetic field. They are vector functions of a three-dimensional position vector $x \in \mathbb{T}$ and of the time $t > 0$, where $\mathbb{T} = (\mathbb{R}/2\pi\mathbb{Z})^3$ is the torus. The fields E and B are coupled to the electron density through the Maxwell equations and act on the electrons via the Lorentz force. We assume that in the plasma the ions are nonmoving and become a uniform background with a fixed unit density. This implies that the density of ions is equal to 1 and the velocity of ions vanishes. Under these assumptions, the dynamics of the compressible electrons obey the (scaled) one-fluid Euler–Maxwell system:

$$(1.1) \quad \partial_t n + \operatorname{div}(nu) = 0,$$

$$(1.2) \quad \partial_t(nu) + \operatorname{div}(nu \otimes u) + \nabla p(n) = -n(E + \gamma u \times B),$$

$$(1.3) \quad \gamma \epsilon \partial_t E - \nabla \times B = \gamma nu, \quad \gamma \partial_t B + \nabla \times E = 0,$$

$$(1.4) \quad \epsilon \operatorname{div} E = 1 - n, \quad \operatorname{div} B = 0$$

for $x \in \mathbb{T}$ and $t > 0$ subject to initial conditions

$$(1.5) \quad (n, u, E, B)(t = 0) = (n_0^\epsilon, u_0^\epsilon, E_0^\epsilon, B_0^\epsilon)$$

for $x \in \mathbb{T}$. In the above equations, $p = p(n)$ is the pressure, assumed to be smooth and strictly increasing for $n > 0$, $j = nu$ is the current density, and $E + \gamma u \times B$ represents

*Received by the editors March 22, 2007; accepted for publication (in revised form) February 4, 2008; published electronically May 28, 2008.

<http://www.siam.org/journals/sima/40-2/68605.html>

[†]Laboratoire de Mathématiques, CNRS UMR 6620, Université Blaise Pascal (Clermont-Ferrand 2), 63177 Aubière cedex, France (peng@math.univ-bpclermont.fr).

[‡]College of Applied Sciences, Beijing University of Technology, PingLeYuan100, Chaoyang District, Beijing 100022, People's Republic of China (wangshu@bjut.edu.cn). This author's research was partially supported by the NSFC (grant 10771009) and BSFC (grant 1082001) of China, the Educational Ministry of China (grant NCET-04-0203), and the Personal Ministry of China.

the Lorentz force. Equations (1.1)–(1.2) are the mass and momentum balance laws for the electrons, respectively, while (1.3)–(1.4) are the Maxwell equations. It is easy to see that equations (1.4) are redundant with equations (1.3) as soon as they are satisfied by the initial data. However, we keep them in the system because this redundancy may be lost in the asymptotic limit.

Equations (1.1)–(1.4) can be viewed as a fluid version of the Vlasov–Maxwell system for a plasma describing the evolution of the electron phase density [3]. Their form is similar to the two-fluid Euler–Maxwell system (see [1, 5, 19, 21]). However, the natures of the one-fluid and the two-fluid Euler–Maxwell systems are different. In the two-fluid model, the vanishing velocity of ions implies the vanishing electric field and the quasi neutrality of the plasma; i.e., the density of ions is equal to the density of electrons. Thus, the one-fluid equations (1.1)–(1.4) cannot be derived from the two-fluid equations.

The dimensionless parameters ϵ and γ can be chosen independently of each other, according to the desired scaling. Physically, ϵ and γ can be chosen to be proportional to the Debye length and $\frac{1}{c}$, where $c = (\epsilon_0 \nu_0)^{-\frac{1}{2}}$ is the speed of light, with ϵ_0 and ν_0 being the vacuum permittivity and permeability (see [5, pp. 349–351]). The limit $\epsilon \rightarrow 0$ leads to $n = 1$, which is called the quasi neutrality of the plasma (here the constant 1 denotes the unit density of nonmoving ions). Thus, the limit $\epsilon \rightarrow 0$ is called the quasi-neutral limit. Also, the limit $\gamma \rightarrow 0$ is physically called the nonrelativistic limit. For the other physical meaning of the dimensionless parameters ϵ and γ , see [3, 10].

In the present paper, we concentrate on the so-called quasi-neutral regime. Hence here we consider only the following quasi-neutral scaling: $\gamma = O(1)$ and $\epsilon \ll 1$.

Now, setting formally $\epsilon = 0$ in the system (1.1)–(1.4), one can arrive at the so-called electron magnetohydrodynamics (e-MHD) equations as follows [3]:

$$(1.6) \quad \partial_t u + u \cdot \nabla u + E = -\gamma u \times B,$$

$$(1.7) \quad -\nabla \times B = \gamma u, \quad \operatorname{div} B = 0,$$

$$(1.8) \quad \gamma \partial_t B + \nabla \times E = 0, \quad n = 1.$$

The e-MHD system (1.6)–(1.8) is incompressible, i.e., $\operatorname{div} u = 0$, which is precisely the limit equation obtained from (1.1) by using the fact that $n = 1$ in (1.8). However, this incompressible condition need not be written separately since it can be obtained by the first equation in (1.7).

The main purpose of this paper is to prove rigorously the above formal limit for smooth solutions of the Euler–Maxwell system (1.1)–(1.4) on time intervals, on which a smooth solution of the incompressible e-MHD equations exists. The precise statement is given in section 2.

From the view of the singular perturbation theory, the quasi-neutral limit $\epsilon \rightarrow 0$ in the Euler–Maxwell system is a problem of singular perturbation for hyperbolic systems (see [8, 11, 12, 22]). However, it is very different from the theory of singular low Mach number limit for symmetrizable hyperbolic systems by Klainerman and Majda in [11, 12]. In the latter case, the essential singularity can be cancelled using a symmetrizer of hyperbolic systems. For the quasi-neutral limit in the Euler–Maxwell system, besides the singularities in the Maxwell equations, there exists an extra singularity caused by the coupling electromagnetic field (source term) in the Euler equations, which cannot be overcome by using the symmetric technique of hyperbolic systems. Hence the singular limit theory for symmetrizable hyperbolic systems developed by Klainerman

and Majda [11, 12] or extended further by Schochet [22] cannot be applied here to obtain the uniform a priori estimates of the solution with respect to ϵ .

In this paper, we control these singularities and then derive rigorously the e-MHD system from the Euler–Maxwell system by an elaborate energy method based on studies on an ϵ -weighted energy. Let $(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)$ be the solution to the Euler–Maxwell problem and (u^0, E^0, B^0) be the solution to the incompressible e-MHD equations. Our basic idea to establish the uniform a priori estimates is the need to obtain an estimate not just for $n^\epsilon - 1 - \epsilon \operatorname{div} E^0$ but also for $\frac{n^\epsilon - 1 - \epsilon \operatorname{div} E^0}{\sqrt{\epsilon}}$ and then for $\sqrt{\epsilon}(E^\epsilon - E^0)$. In particular, the order of derivatives that we need to estimate for the latter is one less than for the other main quantities $(n^\epsilon - 1 - \epsilon \operatorname{div} E^0, u^\epsilon - u^0, B^\epsilon - B^0)$. These estimates cannot be obtained by straightforward Sobolev energy estimates for which the same order of derivatives for both is needed. They are achieved due to the dissipation structure of the equation for $\frac{n^\epsilon - 1 - \epsilon \operatorname{div} E^0}{\sqrt{\epsilon}}$ (see (4.21)). Finally, the desired estimates are obtained through a priori estimates of vorticity and divergence (see Lemmas 4.2–4.3).

There have been a lot of studies on the Euler–Poisson equations and their asymptotic analysis contrarily to the study on the Euler–Maxwell equations. See [2, 4, 7, 8, 20, 23, 24, 25] and the references therein. The first mathematical study of the Euler–Maxwell equations with an extra relaxation term is due to Chen, Jerome, and Wang [6], where a global existence result to weak solutions in the one-dimensional case is established by the fractional step Godunov scheme together with a compensated compactness argument. Paper [6] also exhibits some applications of the model (1.1)–(1.4) in the semiconductor theory. Since then little progress has been made on the Euler–Maxwell equations. Recently the convergence of the Euler–Maxwell system to the compressible Euler–Poisson system has been proved in [17] via the nonrelativistic limit, in which general initial data are allowed by performing an initial layer analysis. This limit corresponds to $\gamma \rightarrow 0$ and $\epsilon > 0$ being fixed. The convergence of the compressible Euler–Maxwell equations to the incompressible Euler equations is justified as $\gamma = \epsilon \rightarrow 0$ (see [18]). Finally, we mention that a related asymptotic limit problem on the Vlasov–Maxwell system is discussed in [2, 3].

We stress that, in this paper, the convergence of the electron density, the current velocity vector, and the magnetic field vector of the Euler–Maxwell systems is strong, whereas the convergence of the electric field vector is only weak. That is why we do not need the assumption that the initial electric field of Euler–Maxwell systems tends to the initial value of the electric field of the limit system. The latter is not arbitrarily given but is determined by the initial data of the limit system (2.6)–(2.9) (see Proposition 2.1). This is different from the situation in [18], where the initial data are prepared. Also, by checking the proof in section 4, we see that it is possible to extend the results of this paper ($\gamma > 0$ being fixed and $\epsilon \rightarrow 0$) to the case $\gamma \rightarrow 0$ and $\epsilon \rightarrow 0$ without any relation between γ and ϵ . This limit is still governed by the incompressible Euler equations. To see this more clearly, the parameter γ is kept in the estimates.

Notation and preliminary results.

(1) Throughout this paper, $\nabla = \nabla_x$ is the gradient, $\nabla \cdot$ is the divergence operator, and $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ and β , etc., are multi-indices. We denote by $H^s(\mathbb{T})$ the standard Sobolev space in torus \mathbb{T} , which is defined by Fourier transform, namely, $f \in H^s(\mathbb{T})$ if and only if

$$\|f\|_s^2 = (2\pi)^3 \sum_{k \in \mathbb{Z}^3} (1 + |k|^2)^s |(\mathcal{F}f)(k)|^2 < +\infty,$$

where $(\mathcal{F}f)(k) = \int_{\mathbb{T}} f(x)e^{-ikx}dx$ is the Fourier transform of $f \in H^s(\mathbb{T})$.

(2) Recall the following basic Moser-type calculus inequalities [11, 12]: for $f, g, v \in H^s$ and any nonnegative multi index $\alpha, |\alpha| \leq s$,

(i) $\|D_x^\alpha(fg)\|_{L^2} \leq C_s(\|f\|_{L^\infty}\|D_x^s g\|_{L^2} + \|g\|_{L^\infty}\|D_x^s f\|_{L^2}), \quad s \geq 0,$

(ii) $\|D_x^\alpha(fg) - fD_x^\alpha g\|_{L^2} \leq C_s(\|D_x f\|_{L^\infty}\|D_x^{s-1}g\|_{L^2} + \|g\|_{L^\infty}\|D_x^s f\|_{L^2}), \quad s \geq 1.$

(3) The following vector analysis formulas will be repeatedly used (see [5]):

(1.9) $\operatorname{div}(f \times g) = \nabla \times f \cdot g - \nabla \times g \cdot f,$

(1.10) $f \cdot \nabla g = (\nabla \times g) \times f + \nabla(f \cdot g) - \nabla f \cdot g,$

(1.11) $f \cdot \nabla f = (\nabla \times f) \times f + \nabla\left(\frac{|f|^2}{2}\right),$

(1.12) $\nabla \times (f \times g) = f \operatorname{div}g - g \operatorname{div}f + (g \cdot \nabla)f - (f \cdot \nabla)g.$

2. Well-posedness of the e-MHD and main results. For smooth solutions of the Euler–Maxwell system (1.1)–(1.5) with $n > 0$, (1.2) is equivalent to

$$\partial_t u + (u \cdot \nabla)u + \nabla h(n) = -(E + \gamma u \times B),$$

where the enthalpy $h(n)$ is defined by

$$h(n) = \int_1^n \frac{p'(s)}{s} ds.$$

Thus, regarding ϵ as a singular perturbation parameter, we can rewrite the problem (1.1)–(1.5) as

(2.1) $\partial_t n^\epsilon + \operatorname{div}(n^\epsilon u^\epsilon) = 0,$

(2.2) $\partial_t u^\epsilon + (u^\epsilon \cdot \nabla)u^\epsilon + \nabla h(n^\epsilon) = -(E^\epsilon + \gamma u^\epsilon \times B^\epsilon),$

(2.3) $\gamma \epsilon \partial_t E^\epsilon - \nabla \times B^\epsilon = \gamma n^\epsilon u^\epsilon, \quad \gamma \partial_t B^\epsilon + \nabla \times E^\epsilon = 0,$

(2.4) $\epsilon \operatorname{div}E^\epsilon = 1 - n^\epsilon, \quad \operatorname{div}B^\epsilon = 0,$

(2.5) $(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)(t = 0) = (n_0^\epsilon, u_0^\epsilon, E_0^\epsilon, B_0^\epsilon),$

where γ is a positive constant of order one. This means that the magnetic field does not vanish in the limiting process.

We also rewrite the limit system (1.6)–(1.8) as

(2.6) $\partial_t u^0 + u^0 \cdot \nabla u^0 + E^0 = -\gamma u^0 \times B^0,$

(2.7) $-\nabla \times B^0 = \gamma u^0, \quad \operatorname{div}B^0 = 0,$

(2.8) $\gamma \partial_t B^0 + \nabla \times E^0 = 0, \quad n^0 = 1.$

Following the idea of [3], introduce the general vorticity

$$\omega^0 = \nabla \times (u^0 - \gamma A^0),$$

where A^0 is the magnetic potential such that

$$\nabla \times A^0 = B^0 \quad \text{and} \quad \operatorname{div} A^0 = 0.$$

Noting $\nabla \cdot u^0 = 0$ and using the identity (1.11), we can write the e-MHD equations (2.6)–(2.8) in a different way:

$$\partial_t \omega^0 + u^0 \cdot \nabla \omega^0 - \omega^0 \cdot \nabla u^0 = 0, \quad -\Delta u^0 + \gamma^2 u^0 = \nabla \times \omega^0.$$

Therefore the existence results are the same as for the incompressible Euler equations [14]. In particular, we have local smooth solutions of the e-MHD equations for the smooth initial data given by

$$(2.9) \quad u^0(t=0) = u_0^0, \quad B^0(t=0) = B_0^0.$$

PROPOSITION 2.1 (see [3]). *Assume that $u_0^0, B_0^0 \in C^\infty(\mathbb{T})$ satisfy*

$$(2.10) \quad -\nabla \times B_0^0 = \gamma u_0^0, \quad \operatorname{div} B_0^0 = 0.$$

Then there exist $0 < T_ \leq \infty$ (if $d = 2$, $T_* = \infty$), the maximal existence time, and a unique smooth solution $(u^0, B^0, E^0) \in C^\infty(\mathbb{T} \times [0, T_*))$ of the incompressible e-MHD equations (2.6)–(2.9) defined on $[0, T_*)$.*

For the convergence of the compressible Euler–Maxwell system (2.1)–(2.5), our main result is stated as follows.

THEOREM 2.1. *Let $s_0 > \frac{3}{2} + 2$ and $\gamma > 0$ be fixed. Let $u_0^0, B_0^0 \in C^\infty(\mathbb{T})$ satisfy (2.10) and $n_0^\epsilon, E_0^\epsilon, B_0^\epsilon \in C^\infty(\mathbb{T})$ satisfy*

$$\epsilon \operatorname{div} E_0^\epsilon = 1 - n_0^\epsilon, \quad \operatorname{div} B_0^\epsilon = 0.$$

Assume that

$$(2.11) \quad \|(n_0^\epsilon - 1, u_0^\epsilon - u_0^0, B_0^\epsilon - B_0^0)\|_{H^{s_0}(\mathbb{T})} + \|\sqrt{\epsilon} E_0^\epsilon\|_{H^{s_0}(\mathbb{T})} \leq C\sqrt{\epsilon}$$

for some positive constant C independent of ϵ . Let T_ , $0 < T_* \leq \infty$ ($d = 2, T_* = \infty$), be the maximal existence time of the smooth solution $(u^0, B^0, E^0) \in C^\infty(\mathbb{T} \times [0, T_*))$ of the incompressible e-MHD equations (2.6)–(2.9). Then, for any $T_0 < T_*$, there exist constants $\epsilon_0(T_0) > 0$ and $\tilde{M}(T_0) > 0$, depending only upon T_0 and the initial data, such that the Euler–Maxwell system (2.1)–(2.5) has a classical smooth solution $(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)$, defined on $[0, T_0]$, satisfying*

$$(2.12) \quad \|(n^\epsilon - 1, u^\epsilon - u^0, B^\epsilon - B^0)(\cdot, t)\|_{H^{s_0}(\mathbb{T})} + \|\sqrt{\epsilon} E^\epsilon(\cdot, t)\|_{H^{s_0-1}(\mathbb{T})} \leq \tilde{M}(T_0)\sqrt{\epsilon}$$

for all $0 < \epsilon \leq \epsilon_0$ and $0 \leq t \leq T_0$. As a result, the sequence $(n^\epsilon, u^\epsilon, B^\epsilon)_{\epsilon > 0}$ converges strongly to $(1, u^0, B^0)$ in $L^\infty(0, T_0; H^{s_0}(\mathbb{T}))$. Furthermore, the sequence $(E^\epsilon)_{\epsilon > 0}$ converges to E^0 in $W^{-1, \infty}(0, T_0; H^{s_0-1}(\mathbb{T}))$.

Remark 2.1. Condition (2.11) means that the initial data are well prepared for $(n^\epsilon, u^\epsilon, B^\epsilon)$ but not for E^ϵ , which are only bounded. This is sufficient to conclude the convergence of E^ϵ to E^0 in some weak sense. Since all the terms involving E^ϵ in the Euler–Maxwell system (2.1)–(2.5) are linear, we can pass to the limit in the system (in the sense of distributions, for instance). The fact that the electric field $E^\epsilon(\cdot, t)$ is bounded uniformly in ϵ for all time $t \in [0, T_0]$ is not surprising because, generally speaking, this property should be maintained over time when the initial electric field is bounded uniformly in ϵ and the density n^ϵ has a better convergence rate. As in

[8, 13, 22], some new techniques are required to study the strong convergence of the electric field when its initial value is also prepared. This is related to the initial layer analysis and will be discussed in the future.

Remark 2.2. The convergence rate $O(\sqrt{\epsilon})$ in Theorem 2.1 is probably not optimal. It should be possible to improve it by constructing better approximate solutions under the same assumption. Moreover, if the initial error for $(n^\epsilon, u^\epsilon, B^\epsilon)$ is better than $O(\sqrt{\epsilon})$, then so should the convergence rate. On the other hand, the convergence rate should be better in a weaker norm than $\|\cdot\|_{s_0}$. For example, the estimate (2.12) on E^ϵ and the first equation in (2.4) already give

$$\|n^\epsilon(\cdot, t) - 1\|_{H^{s_0-2}(\mathbb{T})} = O(\epsilon).$$

Remark 2.3. The results of this paper hold in the whole space \mathbb{R}^3 . Indeed, the key point of the proof of Theorem 2.1 is to establish the uniform a priori estimates (2.12) based on the study of a weighted energy combined with singular perturbation methods. This is dealt with in the same way as that of the papers [11, 12] by Klainerman and Majda. Here neither the compactness of \mathbb{T} nor a Poincaré-type inequality is used.

3. Derivation of error equations and local existence. Let $(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)$ be the unknown solution to the problem (2.1)–(2.5) and (u^0, E^0, B^0) be the solution to the incompressible e-MHD equations defined on $[0, T_*)$ given by Proposition 2.1. Denote this by

$$(3.1) \quad (N^\epsilon, U^\epsilon, F^\epsilon, G^\epsilon) = (n^\epsilon - 1 + \epsilon \operatorname{div} E^0, u^\epsilon - u^0, E^\epsilon - E^0, B^\epsilon - B^0),$$

which satisfies the following problem:

$$(3.2) \quad \left\{ \begin{array}{l} \partial_t N^\epsilon + \operatorname{div}((N^\epsilon + 1 - \epsilon \operatorname{div} E^0)U^\epsilon + N^\epsilon u^0) = \epsilon(\partial_t \operatorname{div} E^0 + \operatorname{div}(u^0 \operatorname{div} E^0)), \\ \partial_t U^\epsilon + [(U^\epsilon + u^0) \cdot \nabla]U^\epsilon + (U^\epsilon \cdot \nabla)u^0 + F^\epsilon + \nabla(h(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \\ \quad - h(1 - \epsilon \operatorname{div} E^0)) \\ \quad = -\gamma((U^\epsilon + u^0) \times G^\epsilon + U^\epsilon \times B^0) + \epsilon h'(1 - \epsilon \operatorname{div} E^0)\nabla(\operatorname{div} E^0), \\ \epsilon \gamma \partial_t F^\epsilon - \nabla \times G^\epsilon = \gamma((N^\epsilon + 1 - \epsilon \operatorname{div} E^0)U^\epsilon + N^\epsilon u^0) - \epsilon \gamma(\partial_t E^0 + u^0 \operatorname{div} E^0), \\ \gamma \partial_t G^\epsilon + \nabla \times F^\epsilon = 0, \\ \epsilon \operatorname{div} F^\epsilon = -N^\epsilon, \quad \operatorname{div} G^\epsilon = 0, \\ (N^\epsilon, U^\epsilon, F^\epsilon, G^\epsilon)|_{t=0} \\ \quad = (n_0^\epsilon - 1 + \epsilon \operatorname{div} E^0(t=0), u_0^\epsilon - u_0^0, E_0^\epsilon - E^0(t=0), B_0^\epsilon - B_0^0). \end{array} \right.$$

Note that we use $1 - \epsilon \operatorname{div} E^0$ as the approximate solution of the density n^ϵ instead of the formal approximate solution 1. This implies that the equation $\epsilon \operatorname{div} F^\epsilon = -N^\epsilon$ is homogeneous. This is a key point in the following convergence analysis. Otherwise, an extra term $\operatorname{div} E^0$ appears in the divergence equation (4.21), which is just bounded but does not converge to 0 as $\epsilon \rightarrow 0$. If so, by the techniques used in this paper, the desired convergence rate $O(\sqrt{\epsilon})$ cannot be obtained.

Set

$$W_I^\epsilon = \begin{pmatrix} N^\epsilon \\ U^\epsilon \end{pmatrix}, \quad W_{II}^\epsilon = \begin{pmatrix} F^\epsilon \\ G^\epsilon \end{pmatrix}, \quad W^\epsilon = \begin{pmatrix} W_I^\epsilon \\ W_{II}^\epsilon \end{pmatrix} = \begin{pmatrix} N^\epsilon \\ U^\epsilon \\ F^\epsilon \\ G^\epsilon \end{pmatrix},$$

$$W_0^\epsilon = \begin{pmatrix} N_0^\epsilon \\ U_0^\epsilon \\ F_0^\epsilon \\ G_0^\epsilon \end{pmatrix} = \begin{pmatrix} n_0^\epsilon - 1 + \epsilon \operatorname{div} E^0(t=0) \\ u_0^\epsilon - u_0^0 \\ E_0^\epsilon - E^0(t=0) \\ B_0^\epsilon - B_0^0 \end{pmatrix}, \quad D_0^\epsilon = \begin{pmatrix} \mathbf{I}_{4 \times 4} & \mathbf{0} \\ \mathbf{0} & \begin{pmatrix} \epsilon \gamma \mathbf{I}_{3 \times 3} & \mathbf{0} \\ \mathbf{0} & \gamma \mathbf{I}_{3 \times 3} \end{pmatrix} \end{pmatrix},$$

$$A_i(W^\epsilon) = \begin{pmatrix} \begin{pmatrix} (U^\epsilon + u^0)_i & (N^\epsilon + 1 - \epsilon \operatorname{div} E^0) e_i^T \\ h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) e_i & (U^\epsilon + u^0)_i \mathbf{I}_{3 \times 3} \end{pmatrix} & \mathbf{0} \\ \mathbf{0} & \begin{pmatrix} \mathbf{0} & B_i \\ B_i^T & \mathbf{0} \end{pmatrix} \end{pmatrix},$$

$$H_1(W_I^\epsilon) = \begin{pmatrix} -\epsilon U^\epsilon \cdot \nabla(\operatorname{div} E^0) \\ (U^\epsilon \cdot \nabla) u^0 - \epsilon(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) - h'(1 - \epsilon \operatorname{div} E^0)) \nabla(\operatorname{div} E^0) \\ 0 \\ 0 \end{pmatrix},$$

$$H_2(F^\epsilon) = \begin{pmatrix} 0 \\ F^\epsilon \\ 0 \\ 0 \end{pmatrix}, \quad H_3(W_I^\epsilon, G^\epsilon) = \begin{pmatrix} 0 \\ (U^\epsilon + u^0) \times G^\epsilon + U^\epsilon \times B^0 \\ -((N^\epsilon + 1 - \epsilon \operatorname{div} E^0) U^\epsilon + N^\epsilon u^0) \\ 0 \end{pmatrix}$$

and

$$R^\epsilon = \begin{pmatrix} R_n^\epsilon \\ R_u^\epsilon \\ R_E^\epsilon \\ R_B^\epsilon \end{pmatrix} = \begin{pmatrix} \partial_t \operatorname{div} E^0 + \operatorname{div}(u^0 \operatorname{div} E^0) \\ h'(1 - \epsilon \operatorname{div} E^0) \nabla(\operatorname{div} E^0) \\ -\gamma(\partial_t E^0 + u^0 \operatorname{div} E^0) \\ 0 \end{pmatrix},$$

where (e_1, e_2, e_3) is the canonical basis of \mathbb{R}^3 , $\mathbf{I}_{d \times d}$ ($d = 3, 4$) is a $d \times d$ unit matrix, y_i denotes the i th component of $y \in \mathbb{R}^3$, and

$$B_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix}, \quad B_2 = \begin{pmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \quad B_3 = \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

From (3.2)_{1,3}, the redundant equations $\epsilon \operatorname{div} F^\epsilon = -N^\epsilon$ and $\operatorname{div} G^\epsilon = 0$ in system (3.2) hold as soon as they are satisfied by the initial data. Thus the problem (3.2) for the

unknown W^ϵ can be rewritten as

$$(3.3) \quad \begin{cases} D_0^\epsilon \partial_t W^\epsilon + \sum_{i=1}^3 A_i(W^\epsilon) \partial_{x_i} W^\epsilon + H_1(W_I^\epsilon) = H_2(F^\epsilon) + \epsilon R^\epsilon - \gamma H_3(W_I^\epsilon, G^\epsilon), \\ W^\epsilon|_{t=0} = W_0^\epsilon, \end{cases}$$

with

$$\epsilon \operatorname{div} F^\epsilon(x, 0) = -N^\epsilon(x, 0), \quad \operatorname{div} G^\epsilon(x, 0) = 0,$$

which can be guaranteed by the assumptions on the initial data.

It is not difficult to see that the equations for W^ϵ in (3.3) are symmetrizable hyperbolic; i.e., if we introduce

$$A_0(W^\epsilon) = \begin{pmatrix} \left(\begin{array}{cc} h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) & 0 \\ 0 & (N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \mathbf{I}_{3 \times 3} \end{array} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{6 \times 6} \end{pmatrix},$$

which is positively definite when $N^\epsilon + 1 - \epsilon \operatorname{div} E^0 \geq M_0 > 0$ for $\epsilon \ll 1$ and $\|N^\epsilon\|_{L^\infty} \leq \frac{1}{2}$, then $A_0 D_0^\epsilon$ and $\tilde{A}_i(W^\epsilon) = A_0(W^\epsilon) A_i(W^\epsilon)$ are symmetric for all $1 \leq i \leq 3$. Note that, for smooth solutions, the Euler–Maxwell system (2.1)–(2.5) is equivalent to that of (3.2) or (3.3). Thus, by the standard existence theory of local smooth solutions for symmetrizable hyperbolic equations (see [9, 15, 16]), we have the following result.

PROPOSITION 3.1. *Let $M > 0$ and $u_0^0, B_0^0 \in C^\infty$ be given and W_0^ϵ satisfy $W_0^\epsilon \in H^s$, $s > \frac{3}{2} + 2$, and $\|N_0^\epsilon\|_{H^s(\mathbb{T})} \leq \delta$ for some given $\delta > 0$ (to be chosen sufficiently small so that $M\delta C_s \leq \frac{1}{2}$, where C_s is the Sobolev embedding constant). Then, for any fixed $\epsilon (\ll 1)$, there exist $0 < T_\epsilon(\delta) \leq \infty$, the maximal existence time, and a unique smooth solution $W^\epsilon \in \bigcap_{l=0}^1 C^l([0, T_\epsilon]; H^{s-l}(\mathbb{T}))$ of the system (3.2) or (3.3) on $[0, T_\epsilon)$ satisfying $\sup_{0 < t < T_\epsilon} \|N^\epsilon(t)\|_{H^s(\mathbb{T})} \leq M\delta$. Moreover, if $T_\epsilon < \infty$, then, for any fixed ϵ (sufficiently small), we have*

$$(3.4) \quad \text{either } \lim_{t \rightarrow T_\epsilon} \|N^\epsilon(t)\|_{H^s(\mathbb{T})} = M\delta \text{ or } \lim_{t \rightarrow T_\epsilon} \|(U^\epsilon, F^\epsilon, G^\epsilon)(\cdot, t)\|_{H^s(\mathbb{T})} = +\infty.$$

Note that the error equations (3.2) or (3.3) from the Euler–Maxwell system are not in the form covered by the well-known Klainerman and Majda theory of singular limits of hyperbolic systems because of the extra singularity caused by the coupling source term $H_2(F^\epsilon)$ in the Euler equations. Hence the singular limit theory for symmetrizable hyperbolic systems developed by Klainerman and Majda [11, 12] or extended further by Schochet [22] cannot be applied here to obtain the uniform a priori estimates of the solution W^ϵ with respect to ϵ .

4. Proof of the main results. In this section, we justify rigorously the convergence of the Euler–Maxwell system to the incompressible e-MHD equations; namely, we prove Theorem 2.1 by using the asymptotic expansion of singular perturbations and careful classical energy methods.

4.1. Convergence rate and uniform a priori estimates. We first establish the convergence rate of the error function $(N^\epsilon, U^\epsilon, B^\epsilon)$ by obtaining the a priori estimates uniformly in ϵ . As a consequence, we obtain the existence of exact solutions

$(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)$ to (2.1)–(2.5) in a time interval independent of ϵ , and the convergence of $(n^\epsilon, u^\epsilon, E^\epsilon, B^\epsilon)$ to $(1, u^0, E^0, B^0)$ as $\epsilon \rightarrow 0$, where (u^0, E^0, B^0) is the solution of the incompressible e-MHD equations (2.6)–(2.9).

In order to justify rigorously the convergence, it suffices to obtain the uniform estimates of the smooth solutions to (3.2) or (3.3) with respect to the parameter ϵ . This is achieved by the elaborate energy method. Since the detailed estimates are very lengthy and involved, we would like to outline some main ingredients here. First, based on the L^2 energy conservation of the Euler–Maxwell system, an ϵ -weighted L^2 energy estimate is derived by noting that there exists a cancellation or some kind of balance between the Euler part and Maxwell part of the Euler–Maxwell system. Second, we introduce a general vorticity of the velocity and the magnetic field and give the corresponding equations of the vorticity and the divergence. Then high order Sobolev energy estimates on the vorticity and the divergence are established by the vector analysis techniques and an elaborate energy method. Next, based on the curl-div decomposition of the gradient, an ϵ -weighted high order Sobolev energy estimate on the density and velocity of the Euler part is derived. Finally, based on the wave-type equation of the Maxwell equations, an ϵ -weighted high order Sobolev energy estimate on the electric field and magnetic field is established. Combining these estimates, an entropy production integral inequality is derived by introducing an ϵ -weighted energy, which yields the desired estimates.

In the following, (\cdot, \cdot) stands for the L^2 inner product of two scalar or vector functions in \mathbb{T} . Also, we denote $\int_{\mathbb{T}}$ by \int and

$$\|\cdot\| = \|\cdot\|_{L^2(\mathbb{T})}, \quad \|\cdot\|_l = \|\cdot\|_{H^l(\mathbb{T})}, \quad l \in \mathbb{N}^*.$$

For convenience, we introduce the ϵ -weighted Sobolev norms

$$\begin{aligned} \|W^\epsilon(t)\|_{l,*} &= \|(N^\epsilon, U^\epsilon, G^\epsilon)(t)\|_l, \\ \| \|W^\epsilon(t)\|_l &= \left(\|W^\epsilon(t)\|_{l,*}^2 + \left\| \left(\frac{N^\epsilon}{\sqrt{\epsilon}}, \sqrt{\epsilon}\gamma F^\epsilon, \epsilon\gamma\partial_t F^\epsilon \right) (t) \right\|_{l-1}^2 \right)^{\frac{1}{2}}, \\ \| \|W^\epsilon\|_{l,T} &= \sup_{0 < t < T} \| \|W^\epsilon(t)\|_l, \quad l \in \mathbb{N}^*. \end{aligned}$$

Note that these norms are different from those used in [18] to prove the convergence of the compressible Euler–Maxwell equations to the incompressible Euler equations.

The key estimate of this paper is contained in the following result.

PROPOSITION 4.1. *Let l be an integer such that $l > \frac{3}{2} + 2$. Assume*

$$(4.1) \quad \| \|W_0^\epsilon\|_l \leq D_1\sqrt{\epsilon}$$

for sufficiently small ϵ and a constant $D_1 > 0$ independent of ϵ . Then, for any $T_0 \in (0, T_)$, there are constants $D_2 > 0$ and $\epsilon_0 > 0$, depending upon T_0 , such that, for all $\epsilon \leq \epsilon_0$, it holds that $T_\epsilon \geq T_0$, and the solution $W^\epsilon(t)$ of (3.2), well defined in $[0, T_\epsilon)$, satisfies*

$$(4.2) \quad \| \|W^\epsilon\|_{l,T_0} \leq D_2\sqrt{\epsilon}.$$

The remainder of section 4 is devoted to the proof of this result.

First, according to the assumption (4.1), using the local existence results in Proposition 3.1 and Sobolev’s embedding lemma, we know that there exists an $\epsilon_0 > 0$ such

that, for all $\epsilon \leq \epsilon_0$, there exists a smooth solution W^ϵ to the system (3.2) or (3.3), defined on $[0, T_\epsilon)$, satisfying

$$(4.3) \quad \sup_{0 < t < T_\epsilon} \|N^\epsilon(t)\|_{L^\infty(\mathbb{T})} \leq C_s \sup_{0 < t < T_\epsilon} \|N^\epsilon(t)\|_{H^s(\mathbb{T})} \leq MC_s D_1 \sqrt{\epsilon} \leq \frac{1}{2}.$$

Here we take $\delta = D_1 \sqrt{\epsilon}$ in Proposition 3.1, and M will be chosen to be a sufficiently large constant independent of ϵ . Hence the rest involves establishing the a priori estimates uniformly with respect to ϵ so as to guarantee $T_\epsilon \geq T_0$ for any given $T_0 < T_*$ and sufficiently small ϵ . Of course, if $T_\epsilon = \infty$, it suffices to obtain the a priori estimates uniformly with respect to ϵ .

In the following, assuming the conditions of Proposition 4.1, we establish a priori estimates by the elaborate energy method in several steps.

For any $T_1 < 1$ independent of ϵ , denote by $T = T^\epsilon = \min\{T_1, T_\epsilon\}$ and by $C > 0$ a constant which depends upon T_0, D_1 but does not depend upon M, T_1, T , and ϵ .

4.2. L^2 -estimates. Based on the L^2 -conservation of solutions to the Euler–Maxwell system, we obtain L^2 -estimates of the error function W^ϵ . Our basic idea is to control the electric field F^ϵ using the special structure between the Euler part and Maxwell part in the Euler–Maxwell system by introducing the extra singular term $\|\frac{N^\epsilon}{\sqrt{\epsilon}}\|$.

LEMMA 4.1. *For all $0 < t < T$ and sufficiently small ϵ , it holds that*

$$(4.4) \quad \int \left\{ |U^\epsilon|^2 + \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds + \epsilon |F^\epsilon|^2 + |G^\epsilon|^2 \right\} (t) dx$$

$$\leq \int \left\{ |U^\epsilon|^2 + \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds + \epsilon |F^\epsilon|^2 + |G^\epsilon|^2 \right\} (t = 0) dx$$

$$+ \int_0^t \left\{ \|\sqrt{\epsilon} F^\epsilon\|^2 + C(\|W^\epsilon\|_{l,*}^2 + \|W^\epsilon\|_{l,*} + 1) \left\| \left(N^\epsilon, U^\epsilon, \nabla U^\epsilon, \gamma G^\epsilon, \frac{N^\epsilon}{\sqrt{\epsilon}} \right) \right\|^2 \right\} (s) ds$$

$$+ C\epsilon^2 + C\epsilon.$$

Proof. Taking the L^2 inner product of the second equation in the error system (3.2) for U^ϵ , by integration by parts, we get

$$\frac{d}{dt}(U^\epsilon, U^\epsilon) + 2(F^\epsilon, U^\epsilon)$$

$$= (\operatorname{div}(U^\epsilon + u^0)U^\epsilon, U^\epsilon) + 2(h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0), \operatorname{div} U^\epsilon)$$

$$(4.5) \quad - 2(U^\epsilon \nabla u^0 + \gamma(U^\epsilon \times B^0 + (U^\epsilon + u^0) \times G^\epsilon) - \epsilon R_u^\epsilon, U^\epsilon).$$

Now we estimate each term on the right-hand side of (4.5).

For the first and third terms, using the property of the approximate solution (u^0, E^0, B^0) , Cauchy–Schwarz’s inequality, and Sobolev’s lemma, we get

$$(4.6) \quad (\operatorname{div}(U^\epsilon + u^0)U^\epsilon, U^\epsilon) \leq C(\|W^\epsilon(t)\|_{l,*} + 1)\|U^\epsilon\|^2$$

and

$$(4.7) \quad \begin{aligned} & -2(U^\epsilon \nabla u^0 + \gamma(U^\epsilon \times B^0 + (U^\epsilon + u^0) \times G^\epsilon) - \epsilon R_u^\epsilon, U^\epsilon) \\ & \leq C(\|W^\epsilon(t)\|_{l,*} + 1)\|(U^\epsilon, \gamma G^\epsilon)\|^2 + C\epsilon^2. \end{aligned}$$

For the second term, noting that the first equation in (3.2) can be rewritten as

$$\operatorname{div} U^\epsilon = -(\partial_t N^\epsilon + \operatorname{div}(N^\epsilon(U^\epsilon + u^0))) + \epsilon(\operatorname{div}(\operatorname{div} E^0 U^\epsilon) + R_n^\epsilon),$$

from (4.3) we have, for sufficiently small ϵ ,

$$(4.8) \quad \begin{aligned} & (h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0), \operatorname{div} U^\epsilon) \\ & = -\int (h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0))(\partial_t N^\epsilon + \operatorname{div}(N^\epsilon(U^\epsilon + u^0))) dx \\ & \quad + \epsilon \int (h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0))(\operatorname{div}(\operatorname{div} E^0 U^\epsilon) + R_n^\epsilon) dx \\ & = -\frac{d}{dt} \int \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds dx \\ & \quad + \int \int_0^{N^\epsilon} (h'(1 - \epsilon \operatorname{div} E^0 + s) - h'(1 - \epsilon \operatorname{div} E^0)) \partial_t (1 - \epsilon \operatorname{div} E^0) ds dx \\ & \quad - \int (h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0)) \operatorname{div}(N^\epsilon(U^\epsilon + u^0)) dx \\ & \quad + \epsilon \int (h(1 - \epsilon \operatorname{div} E^0 + N^\epsilon) - h(1 - \epsilon \operatorname{div} E^0))(\operatorname{div}(\operatorname{div} E^0 U^\epsilon) + R_n^\epsilon) dx \\ & \leq -\frac{d}{dt} \int \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds dx \\ & \quad + C(\|W^\epsilon(t)\|_{l,*} + 1)\|(N^\epsilon, U^\epsilon, \nabla U^\epsilon)\|^2 + C\epsilon^2. \end{aligned}$$

Combining (4.5) with (4.6)–(4.8), we have

$$(4.9) \quad \begin{aligned} & \frac{d}{dt} \left[(U^\epsilon, U^\epsilon) + \int \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds dx \right] + 2(F^\epsilon, U^\epsilon) \\ & \leq C(\|W^\epsilon(t)\|_{l,*} + 1)\|(N^\epsilon, U^\epsilon, \nabla U^\epsilon, \gamma G^\epsilon)\|^2 + C\epsilon^2. \end{aligned}$$

Multiplying the first equation of the Maxwell system in the error system (3.2) by $\frac{1}{\gamma} F^\epsilon$ and the second one by $\frac{1}{\gamma} G^\epsilon$, by integration by parts, we get

$$(4.10) \quad \begin{aligned} & \frac{d}{dt} (\epsilon \|F^\epsilon\|^2 + \|G^\epsilon\|^2) + \frac{2}{\gamma} \int (\nabla \times F^\epsilon \cdot G^\epsilon - \nabla \times G^\epsilon \cdot F^\epsilon) dx - 2(U^\epsilon, F^\epsilon) \\ & = 2(N^\epsilon(U^\epsilon + u^0), F^\epsilon) - 2(\epsilon \operatorname{div} E^0 U^\epsilon, F^\epsilon) + \frac{2}{\gamma} (\epsilon R_E^\epsilon, F^\epsilon). \end{aligned}$$

On one hand, using the vector analysis formulas (1.9), the term $O(\frac{1}{\gamma})$ appearing in Sobolev’s energy estimates vanishes, i.e.,

$$(4.11) \quad \int (\nabla \times F^\epsilon \cdot G^\epsilon - \nabla \times G^\epsilon \cdot F^\epsilon) dx = \int \operatorname{div}(F^\epsilon \times G^\epsilon) dx = 0.$$

On the other hand, using Young’s inequality, we have

$$(4.12) \quad \begin{aligned} (N^\epsilon(U^\epsilon + u^0), F^\epsilon) &\leq \epsilon \|F^\epsilon\|^2 + \frac{1}{4\epsilon} \|N^\epsilon(U^\epsilon + u^0)\|^2 \\ &\leq \epsilon \|F^\epsilon\|^2 + C(\|U^\epsilon(t)\|_l^2 + 1) \left\| \frac{N^\epsilon}{\sqrt{\epsilon}} \right\|^2, \end{aligned}$$

where a singular term $\| \frac{N^\epsilon}{\sqrt{\epsilon}} \|$ appears, and

$$(4.13) \quad -2(\epsilon \operatorname{div} E^0 U^\epsilon, F^\epsilon) + \frac{2}{\gamma} (\epsilon R_E^\epsilon, F^\epsilon) \leq \epsilon \|F^\epsilon\|^2 + C\epsilon \|U^\epsilon\|^2 + C\epsilon.$$

Thus, combining (4.10) with (4.11)–(4.13), we get

$$(4.14) \quad \begin{aligned} &\frac{d}{dt} (\epsilon \|F^\epsilon\|^2 + \|G^\epsilon\|^2) - 2(U^\epsilon, F^\epsilon) \\ &\leq \epsilon \|F^\epsilon\|^2 + C\epsilon \|U^\epsilon\|^2 + C(\|U^\epsilon(t)\|_l^2 + 1) \left\| \frac{N^\epsilon}{\sqrt{\epsilon}} \right\|^2 + C\epsilon. \end{aligned}$$

It follows from (4.9) and (4.14) that

$$\begin{aligned} &\frac{d}{dt} \int \left\{ |U^\epsilon|^2 + \int_0^{N^\epsilon} (h(1 - \epsilon \operatorname{div} E^0 + s) - h(1 - \epsilon \operatorname{div} E^0)) ds + \epsilon |F^\epsilon|^2 + |G^\epsilon|^2 \right\} (t) dx \\ &\leq \|\sqrt{\epsilon} F^\epsilon\|^2 + C(\|W^\epsilon(t)\|_{l,*}^2 + \|W^\epsilon(t)\|_{l,*} + 1) \left\| \left(N^\epsilon, U^\epsilon, \nabla U^\epsilon, \gamma G^\epsilon, \frac{N^\epsilon}{\sqrt{\epsilon}} \right) \right\|^2 \\ &\quad + C\epsilon^2 + C\epsilon. \end{aligned}$$

This completes the proof of Lemma 4.1. \square

Note that it follows from the estimate (4.4) that our problem is the need to obtain an estimate not just for N^ϵ but also for $\frac{N^\epsilon}{\sqrt{\epsilon}}$. In particular, the order of derivatives that we need to estimate for the latter is one less than for the other main quantities $(N^\epsilon, U^\epsilon, G^\epsilon)$. However, straightforward Sobolev energy estimates of higher order would require the same number for both. Hence the above method cannot be generalized directly to the high order Sobolev energy estimates. Here the idea is to establish the uniform estimates for $\| \frac{N^\epsilon}{\sqrt{\epsilon}}(t) \|_{l-1}$ through a priori estimates of vorticity and divergence. To this end, we require the equations of vorticity and divergence.

4.3. Derivation of the vorticity and divergence equations. Taking *curl* on the momentum equation and using $\gamma \partial_t G^\epsilon + \nabla \times F^\epsilon = 0$ of the magnetic field in the error system (3.2), we have

$$(4.15) \quad \begin{aligned} &\partial_t (\nabla \times U^\epsilon) + \nabla \times ((U^\epsilon + u^0) \cdot \nabla) U^\epsilon + \nabla \times ((U^\epsilon \cdot \nabla) u^0) - \gamma \partial_t G^\epsilon \\ &= -\gamma \nabla \times ((U^\epsilon + u^0) \times G^\epsilon + U^\epsilon \times B^0). \end{aligned}$$

Since $\operatorname{div}G^\epsilon = 0$, there exists a vector function \mathcal{G}^ϵ such that

$$(4.16) \quad G^\epsilon = \nabla \times \mathcal{G}^\epsilon.$$

Using the vector analysis formulas (1.10) and (1.11), we get

$$(4.17) \quad \nabla \times ((U^\epsilon + u^0) \cdot \nabla)U^\epsilon = \nabla \times ((\nabla \times U^\epsilon) \times (U^\epsilon + u^0)) - \nabla \times (\nabla u^0 \cdot U^\epsilon).$$

Then, putting (4.16) and (4.17) into the equations (4.15), we obtain

$$(4.18) \quad \partial_t(\nabla \times (U^\epsilon - \gamma \mathcal{G}^\epsilon)) + \nabla \times ((\nabla \times (U^\epsilon - \gamma \mathcal{G}^\epsilon)) \times (U^\epsilon + u^0)) = \mathcal{J}_1^\epsilon,$$

where

$$\mathcal{J}_1^\epsilon = \nabla \times (\nabla u^0 \cdot U^\epsilon) - \nabla \times ((U^\epsilon \cdot \nabla)u^0) - \gamma \nabla \times (U^\epsilon \times B^0)$$

satisfies

$$(4.19) \quad \|\mathcal{J}_1^\epsilon\|_{l-1} \leq C\|U^\epsilon\|_l.$$

Here we require $u^0 \in H^{l+1}$.

Next, introduce the general vorticity

$$\omega^\epsilon = \nabla \times (U^\epsilon - \gamma \mathcal{G}^\epsilon);$$

then it follows from (4.18) that ω^ϵ satisfies the following vorticity equation:

$$(4.20) \quad \partial_t \omega^\epsilon + (U^\epsilon + u^0) \cdot \nabla \omega^\epsilon - \omega^\epsilon \cdot \nabla (U^\epsilon + u^0) + \omega^\epsilon \operatorname{div}(U^\epsilon + u^0) = \mathcal{J}_1^\epsilon.$$

Here we have again used the vector formulas (1.12).

Taking *div* on the momentum equations in the error system (3.2) and using the equation $\epsilon \operatorname{div}F^\epsilon = -N^\epsilon$, we can obtain the following divergence equation:

$$(4.21) \quad \begin{aligned} & \partial_t \operatorname{div}U^\epsilon + \operatorname{div}([(U^\epsilon + u^0) \cdot \nabla]U^\epsilon) - \frac{N^\epsilon}{\epsilon} + \Delta(h(N^\epsilon + 1 - \epsilon \operatorname{div}E^0) - h(1 - \epsilon \operatorname{div}E^0)) \\ & = -\gamma \operatorname{div}((U^\epsilon + u^0) \times G^\epsilon) + \mathcal{J}_2^\epsilon, \end{aligned}$$

where

$$\mathcal{J}_2^\epsilon = -\epsilon \operatorname{div}(R_u^\epsilon) - \operatorname{div}((U^\epsilon \cdot \nabla)u^0) - \gamma \operatorname{div}(U^\epsilon \times B^0)$$

satisfies

$$(4.22) \quad \|\mathcal{J}_2^\epsilon\|_{l-1} \leq C\|U^\epsilon\|_l + C\epsilon.$$

4.4. Estimates of the vorticity and divergence. Now we control $\|\nabla \times U^\epsilon\|_{l-1}$ using the vorticity equations (4.20).

LEMMA 4.2. *For any $0 < t < T$ and $l > \frac{3}{2} + 2$, it holds that*

$$(4.23) \quad \begin{aligned} & \|\nabla \times U^\epsilon(t)\|_{l-1}^2 \\ & \leq C(\|\nabla \times U^\epsilon(t=0)\|_{l-1}^2 + \gamma^2 \|G^\epsilon(t=0)\|_{l-1}^2) + C\gamma^2 \|G^\epsilon(t)\|_{l-1}^2 \\ & + C \int_0^t (\|W^\epsilon(s)\|_{l,*} + 1)(\|\nabla \times U^\epsilon(s)\|_{l-1}^2 + \gamma^2 \|G^\epsilon(s)\|_{l-1}^2) ds. \end{aligned}$$

Proof. Let $\alpha \in \mathbb{N}^3$ with $|\alpha| \leq l - 1$. Taking ∂_x^α on (4.20) and multiplying the resulting equation by $\partial_x^\alpha \omega^\epsilon$, by integration by parts, we have the basic Friedrich energy equation

$$(4.24) \quad \begin{aligned} \frac{d}{dt} \|\partial_x^\alpha \omega^\epsilon\|^2 &= (\operatorname{div}(U^\epsilon + u^0) \partial_x^\alpha \omega^\epsilon, \partial_x^\alpha \omega) + 2(H_\alpha^{(1)}, \partial_x^\alpha \omega) \\ &+ 2\left(\partial_x^\alpha (\omega^\epsilon \cdot \nabla(U^\epsilon + u^0)) - \omega^\epsilon \operatorname{div}(U^\epsilon + u^0) + \mathcal{J}_1^\epsilon\right), \partial_x^\alpha \omega^\epsilon \end{aligned}$$

where the commutator $H_\alpha^{(1)}$ is defined by

$$H_\alpha^{(1)} = -[\partial_x^\alpha ((U^\epsilon + u^0) \cdot \nabla \omega^\epsilon) - (U^\epsilon + u^0) \cdot \nabla \partial_x^\alpha \omega^\epsilon],$$

which can be estimated as follows:

$$(4.25) \quad \begin{aligned} \|H_\alpha^{(1)}\| &\leq C(\|\nabla(U^\epsilon + u^0)\|_{L^\infty} \|\partial_x^{l-2} \nabla \omega^\epsilon\| + \|\nabla \omega^\epsilon\|_{L^\infty} \|\partial_x^{l-1}(U^\epsilon + u^0)\|) \\ &\leq C(\|\nabla(U^\epsilon + u^0)\|_{l-1} \|\partial_x^{l-2} \nabla \omega^\epsilon\| + \|\nabla \omega^\epsilon\|_{l-2} \|\partial_x^{s-1}(U^\epsilon + u^0)\|) \\ &\leq C(\|W^\epsilon(t)\|_{l,*} + 1) \|\omega^\epsilon\|_{l-1}. \end{aligned}$$

Here we have used Sobolev's lemma and $l > \frac{3}{2} + 2$.

Using the estimates (4.25) for $H_\alpha^{(1)}$ and (4.19) for \mathcal{J}_1^ϵ , we have, with the aid of Cauchy-Schwarz's inequality and Sobolev's lemma,

$$(4.26) \quad (\operatorname{div}(U^\epsilon + u^0) \partial_x^\alpha \omega^\epsilon, \partial_x^\alpha \omega) \leq C(\|W^\epsilon(t)\|_{l,*} + 1) \|\omega^\epsilon\|_{l-1}^2,$$

$$(4.27) \quad 2(H_\alpha^{(1)}, \partial_x^\alpha \omega) \leq C(\|W^\epsilon(t)\|_{l,*} + 1) \|\omega^\epsilon\|_{l-1}^2$$

and

$$(4.28) \quad \begin{aligned} &2\left(\partial_x^\alpha (\omega^\epsilon \cdot \nabla(U^\epsilon + u^0)) - \omega^\epsilon \operatorname{div}(U^\epsilon + u^0) + \mathcal{J}_1^\epsilon\right), \partial_x^\alpha \omega^\epsilon \\ &\leq C(\|W^\epsilon(t)\|_{l,*} + 1) \|\omega^\epsilon\|_{l-1}^2. \end{aligned}$$

Combining (4.24) with (4.26), (4.27), and (4.28), we get

$$\frac{d}{dt} \|\omega^\epsilon\|_{l-1}^2 \leq C(\|W^\epsilon(t)\|_{l,*} + 1) \|\omega^\epsilon\|_{l-1}^2,$$

which yields, for any $0 < t < T$,

$$(4.29) \quad \|\omega^\epsilon(t)\|_{l-1}^2 \leq C\|\omega^\epsilon(t=0)\|_{l-1}^2 + C \int_0^t (\|W^\epsilon(s)\|_{l,*} + 1) \|\omega^\epsilon(s)\|_{l-1}^2 ds.$$

Using the definition of ω^ϵ , we get

$$(4.30) \quad \|\nabla \times U^\epsilon\|_{l-1}^2 \leq 2\|\omega^\epsilon\|_{l-1}^2 + 2\gamma^2 \|G^\epsilon\|_{l-1}^2$$

and

$$(4.31) \quad \|\omega^\epsilon(t)\|_{l-1}^2 \leq 2\|\nabla \times U^\epsilon(t)\|_{l-1}^2 + 2\gamma^2 \|G^\epsilon\|_{l-1}^2.$$

Then (4.29)–(4.31) give the estimate (4.23).

The proof of Lemma 4.2 is complete. \square

Next, we estimate $\|\operatorname{div}U^\epsilon\|_{l-1}$ using the divergence equation (4.21). The estimate is contained in the following lemma, whose proof is long and is postponed to the appendix.

LEMMA 4.3. *Let $\alpha \in \mathbb{N}^3$ with $|\alpha| \leq l-1$ and $l > \frac{3}{2} + 2$. Then, for any $0 < t < T$, we have*

$$\begin{aligned}
& \left[\|\partial_x^\alpha \operatorname{div}U^\epsilon\|^2 + \frac{1}{\epsilon} \left(\frac{1}{1 - \epsilon \operatorname{div}E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) \right. \\
& \quad \left. + \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div}E^0)}{N^\epsilon + 1 - \epsilon \operatorname{div}E^0} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \right] (t) \\
\leq & C \left[\|\partial_x^\alpha \operatorname{div}U^\epsilon\|^2 + \frac{1}{\epsilon} \left(\frac{1}{1 - \epsilon \operatorname{div}E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) \right. \\
& \quad \left. + \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div}E^0)}{N^\epsilon + 1 - \epsilon \operatorname{div}E^0} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \right] (t=0) \\
& + C \int_0^t (\|W^\epsilon(s)\|_{l,*} + 1) \|W^\epsilon(s)\|_{l,*}^2 ds \\
(4.32) \quad & + C\gamma^2 \int_0^t (\|U^\epsilon(s)\|_l^2 + 1) \|G^\epsilon(s)\|_l^2 ds + C\epsilon^2 + C\epsilon.
\end{aligned}$$

4.5. High order energy estimates on the electromagnetic field. Finally, we derive the high order energy estimates on the electromagnetic field. For the electric field, we establish its estimates by using the wave formulas of the Maxwell equations as follows.

LEMMA 4.4. *For any $0 < t \leq T$ and $l > \frac{3}{2} + 2$, it holds that*

$$\begin{aligned}
& (\epsilon^2 \gamma^2 \|\partial_t F^\epsilon\|_{l-1}^2 + \epsilon \gamma^2 \|F^\epsilon\|_{l-1}^2 + \epsilon \|\nabla \times F^\epsilon\|_{l-1}^2)(t) \\
& \leq (\epsilon^2 \gamma^2 \|\partial_t F^\epsilon\|_{l-1}^2 + \epsilon \gamma^2 \|F^\epsilon\|_{l-1}^2 + \epsilon \|\nabla \times F^\epsilon\|_{l-1}^2)(t=0) \\
(4.33) \quad & + C(\gamma^2 + 1) \int_0^t (\|W^\epsilon(t)\|_l^4 + \|W^\epsilon(t)\|_l^2 + 1) \|W^\epsilon(t)\|_l^2 ds + C\gamma^2 \epsilon^2.
\end{aligned}$$

Proof. It follows from the Maxwell equation in the error equation (3.2) that F^ϵ satisfies the following wave-type equation:

$$(4.34) \quad \epsilon \gamma \partial_{tt} F^\epsilon - \frac{1}{\gamma} \Delta F^\epsilon + \frac{1}{\gamma} \nabla \operatorname{div} F^\epsilon + \gamma F^\epsilon = -\gamma N^\epsilon F^\epsilon + \gamma \epsilon \operatorname{div} E^0 F^\epsilon + \mathcal{J}_3^\epsilon,$$

where

$$\begin{aligned}
\mathcal{J}_3 &= \gamma(N^\epsilon + 1 - \epsilon \operatorname{div}E^0) \left(-[(U^\epsilon + u^0) \cdot \nabla] U^\epsilon - (U^\epsilon \cdot \nabla) u^0 \right. \\
& \quad \left. - \nabla(h(N^\epsilon + 1 - \epsilon \operatorname{div}E^0) - h(1 - \epsilon \operatorname{div}E^0)) - \gamma((U^\epsilon + u^0) \times G^\epsilon + U^\epsilon \times B^0) \right. \\
& \quad \left. + \epsilon h'(1 - \epsilon \operatorname{div}E^0) \nabla(\operatorname{div}E^0) \right) + \gamma \partial_t(N^\epsilon + 1 - \epsilon \operatorname{div}E^0) U^\epsilon \\
(4.35) \quad & + \gamma u^0 \partial_t N^\epsilon + \gamma \partial_t u^0 N^\epsilon + \epsilon \partial_t(\partial_t E^0 + u^0 \operatorname{div}E^0)
\end{aligned}$$

satisfies

$$(4.36) \quad \|\mathcal{J}_3\|_{l-1} \leq C\gamma(\|W^\epsilon(t)\|_l^2 + \|W^\epsilon(t)\|_l + 1)\|W^\epsilon(t)\|_l + C\gamma\epsilon.$$

Let $\alpha \in \mathbb{N}^3$ with $|\alpha| \leq l - 1$. Taking ∂_x^α on (4.34), multiplying the resulting equation by $\partial_t \partial_x^\alpha F^\epsilon$, by integration by parts, we get

$$(4.37) \quad \begin{aligned} & \frac{1}{2} \frac{d}{dt} \int \left(\epsilon\gamma |\partial_t \partial_x^\alpha F^\epsilon|^2 + \gamma |\partial_x^\alpha F^\epsilon|^2 + \frac{1}{\gamma} |\partial_x^\alpha (\nabla \times F^\epsilon)|^2 \right) dx \\ &= -\gamma \int \partial_x^\alpha (N^\epsilon F^\epsilon) \partial_t \partial_x^\alpha F^\epsilon dx + \gamma\epsilon \int \partial_x^\alpha (\operatorname{div} E^0 F^\epsilon) \partial_t \partial_x^\alpha F^\epsilon dx \\ &+ \int \partial_x^\alpha \mathcal{J}_3 \partial_t \partial_x^\alpha F^\epsilon dx. \end{aligned}$$

In the following we control the right-hand side of (4.37) by $\gamma \|F^\epsilon\|_{l-1}^2$ and $\gamma\epsilon \|\partial_t F^\epsilon\|_{l-1}^2$. First, it holds that

$$(4.38) \quad \begin{aligned} & -\gamma \int \partial_x^\alpha (N^\epsilon F^\epsilon) \partial_t \partial_x^\alpha F^\epsilon dx \\ &= -\gamma \int N^\epsilon \partial_x^\alpha F^\epsilon \partial_t \partial_x^\alpha F^\epsilon dx - \gamma \int \mathcal{H}_\alpha^{(10)} \partial_t \partial_x^\alpha F^\epsilon dx \\ &\leq \frac{\gamma}{6} \int |\partial_x^\alpha F^\epsilon|^2 dx + 6\gamma \|N^\epsilon\|_{L^\infty}^2 \int |\partial_t \partial_x^\alpha F^\epsilon|^2 dx + \gamma \|\mathcal{H}_\alpha^{(10)}\| \|\partial_t \partial_x^\alpha F^\epsilon\| \\ &\leq \frac{\gamma}{6} \int |\partial_x^\alpha F^\epsilon|^2 dx + C\gamma \|N^\epsilon\|_{l-1}^2 \int |\partial_t \partial_x^\alpha F^\epsilon|^2 dx + C\gamma \|N^\epsilon\|_{l-1} \|F^\epsilon\|_{l-1} \|\partial_t \partial_x^\alpha F^\epsilon\| \\ &\leq \frac{\gamma}{3} \|F^\epsilon\|_{l-1}^2 + C\gamma \|N^\epsilon\|_{l-1}^2 \|\partial_t F^\epsilon\|_{l-1}^2 \\ &= \frac{\gamma}{3} \|F^\epsilon\|_{l-1}^2 + C \left\| \frac{N^\epsilon}{\sqrt{\epsilon}} \right\|_{l-1}^2 \epsilon\gamma \|\partial_t F^\epsilon\|_{l-1}^2, \end{aligned}$$

where the commutator

$$\mathcal{H}_\alpha^{(10)} = \partial_x^\alpha (N^\epsilon F^\epsilon) - N^\epsilon \partial_x^\alpha F^\epsilon$$

can be estimated by

$$\begin{aligned} \|\mathcal{H}_\alpha^{(10)}\| &\leq C \|\nabla N^\epsilon\|_{L^\infty} \|\partial_x^{l-2} F^\epsilon\| + C \|F^\epsilon\|_{L^\infty} \|\partial_x^{l-1} N^\epsilon\| \\ &\leq C \|\nabla N^\epsilon\|_{l-2} \|\partial_x^{l-2} F^\epsilon\| + C \|F^\epsilon\|_{l-1} \|\partial_x^{l-1} N^\epsilon\| \\ &\leq C \|N^\epsilon\|_{l-1} \|F^\epsilon\|_{l-1}. \end{aligned}$$

Then, as in the establishment of (4.38), we can get

$$(4.39) \quad \gamma\epsilon \int \partial_x^\alpha (\operatorname{div} E^0 F^\epsilon) \partial_t \partial_x^\alpha F^\epsilon dx \leq C\epsilon\gamma (\|F^\epsilon\|_{l-1}^2 + \|\partial_t F^\epsilon\|_{l-1}^2).$$

Finally, by Cauchy–Schwarz’s inequality, we have

$$(4.40) \quad \int \partial_x^\alpha \mathcal{J}_3 \partial_t \partial_x^\alpha F^\epsilon dx \leq \frac{C}{\epsilon\gamma} \|\partial_x^\alpha \mathcal{J}_3\|^2 + C\epsilon\gamma \|\partial_x^\alpha \partial_t F^\epsilon\|^2.$$

Combining (4.37) with (4.38), (4.39), and (4.40), we get

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} \int \left(\epsilon \gamma |\partial_t \partial_x^\alpha F^\epsilon|^2 + \gamma |\partial_x^\alpha F^\epsilon|^2 + \frac{1}{\gamma} |\partial_x^\alpha (\nabla \times F^\epsilon)|^2 \right) dx \\
 & \leq \frac{\gamma}{3} \|F^\epsilon\|_{l-1}^2 + C \left\| \frac{N^\epsilon}{\sqrt{\epsilon}} \right\|_{l-1}^2 \epsilon \gamma \|\partial_t F^\epsilon\|_{l-1}^2 + C \epsilon \gamma (\|F^\epsilon\|_{l-1}^2 + \|\partial_t F^\epsilon\|_{l-1}^2) \\
 (4.41) \quad & + \frac{C}{\epsilon \gamma} \|\partial_x^\alpha \mathcal{J}_3\|^2 + C \epsilon \gamma \|\partial_x^\alpha \partial_t F^\epsilon\|^2,
 \end{aligned}$$

from which we can easily obtain (4.33).

The proof of Lemma 4.4 is complete. \square

We estimate the magnetic field by using the estimates on the electric field and the curl-div decomposition technique of the gradient for the magnetic field.

LEMMA 4.5. *For any $0 < t \leq T$ and $l > \frac{3}{2} + 2$, it holds that*

$$(4.42) \quad \|\nabla G^\epsilon\|_{l-1} \leq \|\epsilon \gamma \partial_t F^\epsilon\|_{l-1} + C \gamma (\|N^\epsilon\|_{l-1} + 1) \|U^\epsilon\|_{l-1} + C \gamma \|N^\epsilon\|_{l-1} + C \gamma \epsilon.$$

Proof. Because $\operatorname{div} G^\epsilon = 0$, by the curl-div decomposition formulas of the gradient for the magnetic field

$$(4.43) \quad \|\nabla G^\epsilon\|_{l-1} \leq \|\nabla \times G^\epsilon\|_{l-1} + \|\operatorname{div} G^\epsilon\|_{l-1},$$

it suffices to control $\|\nabla \times G^\epsilon\|_{l-1}$.

Using the third equation in the error system (3.2), we get

$$\begin{aligned}
 \|\nabla \times G^\epsilon\|_{l-1} & \leq \|\epsilon \gamma \partial_t F^\epsilon\|_{l-1} + \gamma \|(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)U^\epsilon + N^\epsilon u^0\|_{l-1} + C \gamma \epsilon \\
 (4.44) \quad & \leq \|\epsilon \gamma \partial_t F^\epsilon\|_{l-1} + C \gamma (\|N^\epsilon\|_{l-1} + 1) \|U^\epsilon\|_{l-1} + C \gamma \|N^\epsilon\|_{l-1} + C \gamma \epsilon.
 \end{aligned}$$

Combining (4.43) and (4.44), we obtain (4.42).

The proof of Lemma 4.5 is complete. \square

4.6. The end of proof of Proposition 4.1. Introduce an ϵ -weighted Sobolev-type energy function

$$\Gamma^\epsilon(t) = \|W^\epsilon(t)\|_l^2.$$

Then it follows from (4.4), (4.23), (4.32), (4.33), and (4.42) that there exists an $\epsilon_0 > 0$, depending only upon T_0 , such that, for any $0 < \epsilon \leq \epsilon_0$ and any $0 < t < T$,

$$(4.45) \quad \Gamma^\epsilon(t) \leq C \Gamma^\epsilon(t=0) + C \int_0^t ((\Gamma^\epsilon + \sqrt{\Gamma^\epsilon} + 1)\Gamma^\epsilon)(s) ds + C \epsilon.$$

Then, applying Gronwall's lemma and using Proposition 3.1 and (3.4), it follows from (4.45) and $\Gamma^\epsilon(t=0) \leq C \epsilon$ that there exist a $0 < T_1 < 1$ and an $\epsilon_0 > 0$ such that $T_\epsilon \geq T_1$ for all $0 < \epsilon \leq \epsilon_0$ and that there exists an ϵ_0 sufficiently small such that, for any $\epsilon \leq \epsilon_0$ and $0 < t < T$,

$$\Gamma^\epsilon(t) \leq C \epsilon,$$

which gives the desired a priori estimate (4.2). Moreover, by the standard continuous induction method, we can extend $T_\epsilon \geq T_0$ for any $T_0 < T_*$.

The proof of Proposition 4.1 is complete. \square

4.7. Proof of Theorem 2.1. According to the definitions of the error functions $N^\epsilon, U^\epsilon, F^\epsilon, G^\epsilon$, the regularities of u^0, E^0, B^0 , and the error system (3.2), it follows from the assumption (2.11) in Theorem 2.1 that

$$\begin{aligned} \|\sqrt{\epsilon}F^\epsilon(t=0)\|_{s_0-1} &\leq C\sqrt{\epsilon}, \\ \left\| \frac{N^\epsilon}{\sqrt{\epsilon}}(t=0) \right\|_{s_0-1} &= \sqrt{\epsilon}\|\operatorname{div}F^\epsilon(t=0)\|_{s_0-1} \leq \sqrt{\epsilon}\|F^\epsilon(t=0)\|_{s_0} \leq C\sqrt{\epsilon}. \end{aligned}$$

Hence, the assumption (4.1) in Proposition 4.1 holds. Thus, the results of Proposition 4.1 imply (2.12). As a consequence, $(n^\epsilon, u^\epsilon, B^\epsilon)$ converges strongly to $(1, u^0, B^0)$ in $L^\infty(0, T_0; H^{s_0}(\mathbb{T}))$. It follows from (2.2), (2.6), and the uniqueness of solutions to the limit problem (2.6)–(2.9) that E^ϵ converges to E^0 in $W^{-1,\infty}(0, T_0; H^{s_0-1}(\mathbb{T}))$.

The proof of Theorem 2.1 is complete. \square

Appendix A. Proof of Lemma 4.3. Taking ∂_x^α on (4.21) and taking the L^2 inner product of the resulting equation with $\partial_x^\alpha \operatorname{div}U^\epsilon$, by integration by parts, we have the following energy equation:

$$\begin{aligned} \text{(A.1)} \quad &\frac{d}{dt} \|\partial_x^\alpha \operatorname{div}U^\epsilon\|^2 \\ &= (\operatorname{div}(U^\epsilon + u^0)\partial_x^\alpha \operatorname{div}U^\epsilon, \partial_x^\alpha \operatorname{div}U^\epsilon) + 2(\mathcal{H}_\alpha^{(2)}, \partial_x^\alpha \operatorname{div}U^\epsilon) \\ &\quad + \frac{2}{\epsilon}(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \operatorname{div}U^\epsilon) - 2(\partial_x^\alpha \Delta(h(N^\epsilon + 1 - \epsilon \operatorname{div}E^0) - h(1 - \epsilon \operatorname{div}E^0)), \partial_x^\alpha \operatorname{div}U^\epsilon) \\ &\quad - 2\gamma(\partial_x^\alpha \operatorname{div}((U^\epsilon + u^0) \times G^\epsilon), \partial_x^\alpha \operatorname{div}U^\epsilon) + 2(\partial_x^\alpha \mathcal{J}_2^\epsilon, \partial_x^\alpha \operatorname{div}U^\epsilon), \end{aligned}$$

where the commutator is defined by

$$\mathcal{H}_\alpha^{(2)} = -\{\partial_x^\alpha \operatorname{div}([(U^\epsilon + u^0) \cdot \nabla]U^\epsilon) - [(U^\epsilon + u^0) \cdot \nabla]\partial_x^\alpha \operatorname{div}U^\epsilon\},$$

which can be estimated as follows:

$$\begin{aligned} \|\mathcal{H}_\alpha^{(2)}\| &\leq C\|\nabla(U^\epsilon + u^0)\|_{L^\infty}\|\partial_x^{l-1}\nabla U^\epsilon\| + C\|\nabla U^\epsilon\|_{L^\infty}\|\partial_x^l(U^\epsilon + u^0)\| \\ &\leq C\|\nabla(U^\epsilon + u^0)\|_{l-1}\|\partial_x^{l-1}\nabla U^\epsilon\| + C\|\nabla U^\epsilon\|_{l-1}\|\partial_x^l(U^\epsilon + u^0)\| \\ &\leq C(\|U^\epsilon\|_l + 1)\|U^\epsilon\|_l. \end{aligned}$$

Hence, by Cauchy–Schwarz’s inequality, Sobolev’s lemma, and using the estimate (4.22) for \mathcal{J}_2^ϵ , we obtain

$$\begin{aligned} &(\operatorname{div}(U^\epsilon + u^0)\partial_x^\alpha \operatorname{div}U^\epsilon, \partial_x^\alpha \operatorname{div}U^\epsilon) + 2(\mathcal{H}_\alpha^{(2)}, \partial_x^\alpha \operatorname{div}U^\epsilon) \\ &\leq C(\|U^\epsilon\|_l + 1)\|U^\epsilon\|_l^2 + C\|\mathcal{H}_\alpha^{(2)}\|\|\partial_x^\alpha \operatorname{div}U^\epsilon\| \\ \text{(A.2)} \quad &\leq C(\|W^\epsilon(t)\|_{l,*} + 1)\|U^\epsilon\|_l^2 \end{aligned}$$

and

$$\begin{aligned} &-2\gamma(\partial_x^\alpha \operatorname{div}((U^\epsilon + u^0) \times G^\epsilon), \partial_x^\alpha \operatorname{div}U^\epsilon) + 2(\partial_x^\alpha \mathcal{J}_2^\epsilon, \partial_x^\alpha \operatorname{div}U^\epsilon) \\ \text{(A.3)} \quad &\leq C\gamma^2(\|U^\epsilon\|_l^2 + 1)\|G^\epsilon\|_l^2 + \|W^\epsilon(t)\|_{l,*}^2 + C\epsilon^2. \end{aligned}$$

The rest involves dealing with the other two terms on the right-hand side of (A.1), which are more difficult to control. To do this, we rewrite the density equation in (3.2) into the following two formulations:

$$(A.4) \quad \operatorname{div} U^\epsilon = -\frac{\partial_t N^\epsilon + \operatorname{div}(N^\epsilon(U^\epsilon + u^0)) - \epsilon U^\epsilon \cdot \nabla(\operatorname{div} E^0) - \epsilon R_n^\epsilon}{1 - \epsilon \operatorname{div} E^0}$$

and

$$(A.5) \quad \operatorname{div} U^\epsilon = -\frac{\partial_t N^\epsilon + (U^\epsilon + u^0) \cdot \nabla N^\epsilon + N^\epsilon \operatorname{div} u^0 - \epsilon U^\epsilon \cdot \nabla(\operatorname{div} E^0) - \epsilon R_n^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0}.$$

In the following, we use the first formulation (A.4) to estimate the electric field term and the second one (A.5) to estimate the nonlinear pressure term integral in order to avoid the presence of the $(l+1)$ th order derivative of the velocity because we are now in H^l energy estimates.

First, we control the singular term $O(\frac{1}{\epsilon})$ in (A.1). We have

$$\begin{aligned} I_1 &= \frac{2}{\epsilon} (\partial_x^\alpha N^\epsilon, \partial_x^\alpha \operatorname{div} U^\epsilon) \\ &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{\partial_t N^\epsilon + \operatorname{div}(N^\epsilon(U^\epsilon + u^0)) - \epsilon U^\epsilon \cdot \nabla(\operatorname{div} E^0) - \epsilon R_n^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right) \\ &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right) - \frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{\operatorname{div}(N^\epsilon(U^\epsilon + u^0))}{1 - \epsilon \operatorname{div} E^0} \right) \right) \\ &\quad + 2 \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{U^\epsilon \cdot \nabla(\operatorname{div} E^0)}{1 - \epsilon \operatorname{div} E^0} \right) \right) + 2 \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{R_n^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right) \\ (A.6) \quad &= \sum_{i=1}^4 I_1^{(i)}, \end{aligned}$$

in which each term can be estimated as follows.

For $I_1^{(1)}$, by Cauchy–Schwarz’s inequality and using the fact that

$$\left\| \partial_t \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \leq C\epsilon,$$

we have

$$\begin{aligned} I_1^{(1)} &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right) \\ &= -\frac{2}{\epsilon} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha \partial_t N^\epsilon \right) - \frac{2}{\epsilon} (\mathcal{H}_\alpha^{(2)}, \partial_x^\alpha N^\epsilon) \\ &= -\frac{1}{\epsilon} \frac{d}{dt} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) + \frac{1}{\epsilon} \left(\partial_t \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) \\ &\quad - \frac{2}{\epsilon} (\mathcal{H}_\alpha^{(2)}, \partial_x^\alpha N^\epsilon) \\ (A.7) \quad &\leq -\frac{1}{\epsilon} \frac{d}{dt} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) + C \|\partial_x^\alpha N^\epsilon\|^2 + \frac{C}{\epsilon} \|\mathcal{H}_\alpha^{(3)}\| \|\partial_x^\alpha N^\epsilon\|, \end{aligned}$$

where the commutator is defined by

$$\mathcal{H}_\alpha^{(3)} = \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) - \frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha \partial_t N^\epsilon,$$

which can be controlled as follows:

$$\begin{aligned} \|\mathcal{H}_\alpha^{(3)}\| &\leq C \left\| \nabla \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \|\partial_x^{l-2} \partial_t N^\epsilon\| + C \|\partial_t N^\epsilon\|_{L^\infty} \left\| \partial_x^{l-1} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\| \\ &\leq C\epsilon \|\partial_t N^\epsilon\|_{l-2} \end{aligned}$$

$$(A.8) \quad \leq C\epsilon \|(N^\epsilon, U^\epsilon)\|_{l-1} + C\epsilon (\|U^\epsilon\|_{l-1} + 1) \|N^\epsilon\|_{l-1} + C\epsilon^2.$$

Here we have used $\|\nabla(\frac{1}{1-\epsilon \operatorname{div} E^0})\|_{L^\infty} \leq C\epsilon$ and $l > \frac{3}{2} + 2$. Combining (A.7) and (A.8), we obtain

$$(A.9) \quad I_1^{(1)} \leq -\frac{1}{\epsilon} \frac{d}{dt} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) + C(\|U^\epsilon\|_{l-1} + 1) \|W^\epsilon(t)\|_{l,*}^2 + C\epsilon^2.$$

For $I_1^{(2)}$, a direct calculation yields

$$\begin{aligned} (A.10) \quad I_1^{(2)} &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{\operatorname{div}(N^\epsilon(U^\epsilon + u^0))}{1 - \epsilon \operatorname{div} E^0} \right) \right) \\ &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha (\operatorname{div}(N^\epsilon(U^\epsilon + u^0))) \right) - \frac{2}{\epsilon} (\partial_x^\alpha N^\epsilon, \mathcal{H}_\alpha^{(4)}) \\ &= -\frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha ((U^\epsilon + u^0) \cdot \nabla N^\epsilon + N^\epsilon \operatorname{div}(U^\epsilon + u^0)) \right) \\ &\quad - \frac{2}{\epsilon} (\partial_x^\alpha N^\epsilon, \mathcal{H}_\alpha^{(4)}) \\ &= \frac{1}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \operatorname{div} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} (U^\epsilon + u^0) \right) \partial_x^\alpha N^\epsilon \right) \\ &\quad - \frac{2}{\epsilon} \left(\partial_x^\alpha N^\epsilon, \frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon \operatorname{div}(U^\epsilon + u^0) \right) \\ &\quad - \frac{2}{\epsilon} (\partial_x^\alpha N^\epsilon, \mathcal{H}_\alpha^{(4)}) - \frac{2}{\epsilon} \sum_{i=5}^{i=6} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \mathcal{H}_\alpha^{(i)} \right), \end{aligned}$$

where the commutators $\mathcal{H}_\alpha^{(i)}$, $i = 4, 5, 6$, are defined by

$$\mathcal{H}_\alpha^{(4)} = \partial_x^\alpha \left(\frac{\operatorname{div}(N^\epsilon(U^\epsilon + u^0))}{1 - \epsilon \operatorname{div} E^0} \right) - \frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha (\operatorname{div}(N^\epsilon(U^\epsilon + u^0))),$$

$$\mathcal{H}_\alpha^{(5)} = \partial_x^\alpha ((U^\epsilon + u^0) \cdot \nabla N^\epsilon) - (U^\epsilon + u^0) \cdot \partial_x^\alpha \nabla N^\epsilon,$$

$$\mathcal{H}_\alpha^{(6)} = \partial_x^\alpha (N^\epsilon \operatorname{div}(U^\epsilon + u^0)) - \partial_x^\alpha N^\epsilon \operatorname{div}(U^\epsilon + u^0),$$

which, with the aid of Sobolev's lemma, can be estimated, respectively, as follows:

$$\begin{aligned}
\|\mathcal{H}_\alpha^{(4)}\| &\leq C \left\| \nabla \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \|\partial_x^{l-2}(\operatorname{div}(N^\epsilon(U^\epsilon + u^0)))\| \\
&\quad + C \|\operatorname{div}(N^\epsilon(U^\epsilon + u^0))\|_{L^\infty} \left\| \partial_x^{l-1} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\| \\
&\leq C \left\| \nabla \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \|\partial_x^{l-2}(\operatorname{div}(N^\epsilon(U^\epsilon + u^0)))\| \\
&\quad + C \|\operatorname{div}(N^\epsilon(U^\epsilon + u^0))\|_{l-2} \left\| \partial_x^{l-1} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \right) \right\| \\
\text{(A.11)} \quad &\leq C\epsilon(\|U^\epsilon\|_l + 1)\|N^\epsilon\|_{l-1},
\end{aligned}$$

$$\begin{aligned}
\|\mathcal{H}_\alpha^{(5)}\| &\leq C \|\nabla(U^\epsilon + u^0)\|_{L^\infty} \|\partial_x^{l-2} \nabla N^\epsilon\| + C \|\nabla N^\epsilon\|_{L^\infty} \|\partial_x^{l-1}(U^\epsilon + u^0)\| \\
&\leq C \|\nabla(U^\epsilon + u^0)\|_{l-1} \|\partial_x^{l-2} \nabla N^\epsilon\| + C \|\nabla N^\epsilon\|_{l-2} \|\partial_x^{l-1}(U^\epsilon + u^0)\| \\
\text{(A.12)} \quad &\leq C(\|U^\epsilon\|_l + 1)\|N^\epsilon\|_{l-1},
\end{aligned}$$

$$\begin{aligned}
\|\mathcal{H}_\alpha^{(6)}\| &\leq C \|\nabla \operatorname{div}(U^\epsilon + u^0)\|_{L^\infty} \|\partial_x^{l-2} N^\epsilon\| + C \|N^\epsilon\|_{L^\infty} \|\partial_x^{l-1} \operatorname{div}(U^\epsilon + u^0)\| \\
&\leq C \|\nabla \operatorname{div}(U^\epsilon + u^0)\|_{l-2} \|\partial_x^{l-2} N^\epsilon\| + C \|N^\epsilon\|_{l-1} \|\partial_x^{l-1} \operatorname{div}(U^\epsilon + u^0)\| \\
\text{(A.13)} \quad &\leq C(\|U^\epsilon\|_l + 1)\|N^\epsilon\|_{l-1}.
\end{aligned}$$

Here we have again used $\|\nabla(\frac{1}{1 - \epsilon \operatorname{div} E^0})\|_{L^\infty} \leq C\epsilon$ and $l > \frac{3}{2} + 2$.

Thus, combining (A.10) together with (A.11)–(A.13), we get, with the aid of Cauchy–Schwarz's inequality, that

$$\text{(A.14)} \quad I_1^{(2)} \leq \frac{C}{\epsilon} (\|U^\epsilon\|_l + 1) \|N^\epsilon\|_{l-1}^2.$$

Also, by Cauchy–Schwarz's inequality, $I_1^{(3)}$ can be estimated as follows:

$$\text{(A.15)} \quad I_1^{(3)} = 2 \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{U^\epsilon \cdot \nabla(\operatorname{div} E^0)}{1 - \epsilon \operatorname{div} E^0} \right) \right) \leq C \|W^\epsilon(t)\|_{l,*}^2.$$

Finally, we estimate $I_1^{(4)}$. Since R_n^ϵ is not small with respect to ϵ (only uniformly bounded), we must use $\|\frac{N^\epsilon}{\sqrt{\epsilon}}\|_{l-1}$ to control $I_1^{(4)}$. Thus, we rewrite $I_1^{(4)}$ as

$$I_1^{(4)} = 2 \left(\partial_x^\alpha N^\epsilon, \partial_x^\alpha \left(\frac{R_n^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right) = 2 \left(\frac{\partial_x^\alpha N^\epsilon}{\sqrt{\epsilon}}, \sqrt{\epsilon} \partial_x^\alpha \left(\frac{R_n^\epsilon}{1 - \epsilon \operatorname{div} E^0} \right) \right).$$

By Cauchy–Schwarz's inequality and the uniform bound on R_n^ϵ , we have

$$\text{(A.16)} \quad I_1^{(4)} \leq \frac{C}{\epsilon} \|N^\epsilon\|_{l-1}^2 + C\epsilon.$$

Combining (A.6) with (A.9), (A.14), (A.15), and (A.16), we obtain

$$\text{(A.17)} \quad I_1 \leq -\frac{1}{\epsilon} \frac{d}{dt} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) + C(\|U^\epsilon\|_l + 1) \|W^\epsilon(t)\|_l^2 + C\epsilon.$$

Next, we control the pressure term I_2 . It holds that

$$\begin{aligned}
I_2 &= -2(\partial_x^\alpha \Delta(h(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) - h(1 - \epsilon \operatorname{div} E^0)), \partial_x^\alpha \operatorname{div} U^\epsilon) \\
&= -2\left(\partial_x^\alpha \operatorname{div}(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \nabla N^\epsilon \right. \\
&\quad \left. - \epsilon(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) - h'(1 - \epsilon \operatorname{div} E^0)) \nabla(\operatorname{div} E^0)), \partial_x^\alpha \operatorname{div} U^\epsilon\right) \\
&= -2\left(\operatorname{div}(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon), \partial_x^\alpha \operatorname{div} U^\epsilon\right) - 2(\mathcal{H}_\alpha^{(7)}, \partial_x^\alpha \operatorname{div} U^\epsilon) \\
&\quad + 2\epsilon\left(\partial_x^\alpha \operatorname{div}((h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) - h'(1 - \epsilon \operatorname{div} E^0)) \nabla(\operatorname{div} E^0)), \partial_x^\alpha \operatorname{div} U^\epsilon\right) \\
\text{(A.18)} &= I_2^{(1)} + I_2^{(2)} + I_2^{(3)},
\end{aligned}$$

where

$$\begin{aligned}
I_2^{(1)} &= -2\left(\operatorname{div}(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \nabla \partial_x^\alpha N^\epsilon), \partial_x^\alpha \operatorname{div} U^\epsilon\right), \\
I_2^{(2)} &= -2(\mathcal{H}_\alpha^{(7)}, \partial_x^\alpha \operatorname{div} U^\epsilon), \\
I_2^{(3)} &= 2\epsilon\left(\partial_x^\alpha \operatorname{div}((h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) - h'(1 - \epsilon \operatorname{div} E^0)) \nabla(\operatorname{div} E^0)), \partial_x^\alpha \operatorname{div} U^\epsilon\right)
\end{aligned}$$

and the commutator $\mathcal{H}_\alpha^{(7)}$ is defined by

$$\begin{aligned}
\mathcal{H}_\alpha^{(7)} &= \partial_x^\alpha \operatorname{div}(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \nabla N^\epsilon) - \operatorname{div}(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon) \\
&= \partial_x^\alpha (h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \operatorname{div} \nabla N^\epsilon) - h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \operatorname{div} \nabla N^\epsilon \\
&\quad + \partial_x^\alpha (\nabla h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \nabla N^\epsilon) - \nabla h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon.
\end{aligned}$$

Noting that $I_2^{(3)}$ does contain only the l th order derivatives of the error function W_l^ϵ , it can be easily estimated by

$$\text{(A.19)} \quad I_2^{(3)} \leq \|U^\epsilon\|_l^2 + C\|N^\epsilon\|_l^2.$$

Here we have used the regularity of the limit system, i.e., $\|E^0\|_{l+2} \leq C$. The estimate techniques of the commutator yield

$$\begin{aligned}
\|\mathcal{H}_\alpha^{(7)}\| &\leq C\|\nabla h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)\|_{L^\infty} \|\partial_x^{l-2} \operatorname{div} \nabla N^\epsilon\| \\
&\quad + C\|\operatorname{div} \nabla N^\epsilon\|_{L^\infty} \|\partial_x^{l-1} h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)\| \\
&\quad + C\|\nabla^2 h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)\|_{L^\infty} \|\partial_x^{l-2} \nabla N^\epsilon\| \\
&\quad + C\|\nabla N^\epsilon\|_{L^\infty} \|\partial_x^{l-1} \nabla h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)\| \\
&\leq C(\|N^\epsilon\|_l + 1)\|N^\epsilon\|_l,
\end{aligned}$$

which implies by Cauchy–Schwarz’s inequality that

$$\text{(A.20)} \quad I_2^{(2)} \leq C(\|N^\epsilon\|_l + 1)\|W^\epsilon(t)\|_{l,*}^2.$$

In the following, we estimate $I_2^{(1)}$. By the relation (A.5) between the density and the divergence, we have

$$\begin{aligned}
I_2^{(1)} &= 2 \left(\operatorname{div} \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon \right), \right. \\
&\quad \left. \partial_x^\alpha \left(\frac{\partial_t N^\epsilon + (U^\epsilon + u^0) \cdot \nabla N^\epsilon - \epsilon U^\epsilon \cdot \nabla (\operatorname{div} E^0) + N^\epsilon \operatorname{div} u^0 - \epsilon R_n^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
&= -2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \nabla \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
&\quad - 2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \nabla \partial_x^\alpha \left(\frac{(U^\epsilon + u^0) \cdot \nabla N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
&\quad - 2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \right. \\
&\quad \left. \nabla \partial_x^\alpha \left(\frac{-\epsilon U^\epsilon \cdot \nabla (\operatorname{div} E^0) + N^\epsilon \operatorname{div} u^0 - \epsilon R_n^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right)
\end{aligned}$$

$$(A.21) \quad = \sum_{i=1}^3 I_{21}^{(i)}.$$

Now we estimate each term of $I_2^{(1)}$.

For $I_{21}^{(1)}$ and $I_{21}^{(2)}$, we have

$$\begin{aligned}
I_{21}^{(1)} &= -2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \nabla \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
&= -\frac{d}{dt} \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \\
&\quad + \left(\partial_t \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)} \right) \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \\
&\quad - 2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \mathcal{H}_\alpha^{(8)} \right) \\
&\leq -\frac{d}{dt} \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{2(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right)
\end{aligned}$$

$$(A.22) \quad + C(\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*}^2$$

and

$$\begin{aligned}
 I_{21}^{(2)} &= -2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \nabla \partial_x^\alpha \left(\frac{(U^\epsilon + u^0) \cdot \nabla N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
 &= \left(\operatorname{div} \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) (U^\epsilon + u^0)}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \\
 &\quad - 2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \mathcal{H}_\alpha^{(9)} \right) \\
 \text{(A.23)} \quad &\leq C (\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*}^2
 \end{aligned}$$

because the commutators

$$\mathcal{H}_\alpha^{(8)} = \nabla \partial_x^\alpha \left(\frac{\partial_t N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) - \frac{\partial_t \nabla \partial_x^\alpha N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0}$$

and

$$\mathcal{H}_\alpha^{(9)} = \nabla \partial_x^\alpha \left(\frac{(U^\epsilon + u^0) \cdot \nabla N^\epsilon}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) - \frac{(U^\epsilon + u^0) \cdot \nabla (\partial_x^\alpha \nabla N^\epsilon)}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0}$$

can be estimated, respectively, by

$$\begin{aligned}
 \|\mathcal{H}_\alpha^{(8)}\| &\leq C \left\| \nabla \left(\frac{1}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \|\partial_x^{l-1} \partial_t N^\epsilon\| \\
 &\quad + C \|\partial_t N^\epsilon\|_{L^\infty} \left\| \partial_x^l \left(\frac{1}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right\| \\
 &\leq C (\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*} + C \epsilon^2
 \end{aligned}$$

and

$$\begin{aligned}
 \|\mathcal{H}_\alpha^{(9)}\| &\leq C \left\| \nabla \left(\frac{U^\epsilon + u^0}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right\|_{L^\infty} \|\partial_x^{l-1} \nabla N^\epsilon\| \\
 &\quad + C \|\nabla N^\epsilon\|_{L^\infty} \left\| \partial_x^l \left(\frac{U^\epsilon + u^0}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right\| \\
 &\leq C (\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*} + C \epsilon^2.
 \end{aligned}$$

For $I_{21}^{(3)}$, by Cauchy–Schwarz’s inequality and the regularities of u^0, E^0 , we have

$$\begin{aligned}
 I_{21}^{(3)} &= -2 \left(h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0) \partial_x^\alpha \nabla N^\epsilon, \right. \\
 &\quad \left. \nabla \partial_x^\alpha \left(\frac{-\epsilon U^\epsilon \cdot \nabla (\operatorname{div} E^0) + N^\epsilon \operatorname{div} u^0 - \epsilon R_n^c}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \right) \right) \\
 \text{(A.24)} \quad &\leq C \|W^\epsilon(t)\|_{l,*}^3 + C \|W^\epsilon(t)\|_{l,*}^2 + C \epsilon^2.
 \end{aligned}$$

Combining (A.21) with (A.22)–(A.24), we get

$$(A.25) \quad \begin{aligned} I_2^{(1)} \leq & -\frac{d}{dt} \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \\ & + C(\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*}^2 + C\epsilon^2. \end{aligned}$$

Similarly, combining (A.18) with (A.19), (A.20), and (A.25), we get

$$(A.26) \quad \begin{aligned} I_2 \leq & -\frac{d}{dt} \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \\ & + C(\|W^\epsilon(t)\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*}^2 + C\epsilon^2, \end{aligned}$$

and combining (A.1) with (A.2), (A.3), (A.17), and (A.26), we get

$$\begin{aligned} & \frac{d}{dt} \left[\|\partial_x^\alpha \operatorname{div} U^\epsilon\|^2 + \frac{1}{\epsilon} \left(\frac{1}{1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha N^\epsilon, \partial_x^\alpha N^\epsilon \right) \right. \\ & \left. + \left(\frac{h'(N^\epsilon + 1 - \epsilon \operatorname{div} E^0)}{N^\epsilon + 1 - \epsilon \operatorname{div} E^0} \partial_x^\alpha \nabla N^\epsilon, \partial_x^\alpha \nabla N^\epsilon \right) \right] \\ & \leq C(\|W^\epsilon\|_{l,*} + 1) \|W^\epsilon(t)\|_{l,*}^2 + C\gamma^2(\|U^\epsilon\|_l^2 + 1) \|G^\epsilon\|_l^2 + C\epsilon^2 + C\epsilon, \end{aligned}$$

which yields (4.32).

This ends the proof of Lemma 4.3. \square

Acknowledgments. The authors are very grateful to both referees for their constructive comments and helpful suggestions, which considerably improved the presentation of the paper. The second author would like to express his gratitude for the hospitality of the Laboratory of Mathematics at Blaise Pascal University, France, when he visited there in May through July 2006.

REFERENCES

- [1] C. BESSE, P. DEGOND, F. DELUZET, J. CLAUDEL, G. GALLICE, AND C. TESSIERAS, *A model hierarchy for ionospheric plasma modeling*, Math. Models Methods Appl. Sci., 14 (2004), pp. 393–415.
- [2] Y. BRENIER, *Convergence of the Vlasov-Poisson system to the incompressible Euler equations*, Comm. Partial Differential Equations, 25 (2000), pp. 737–754.
- [3] Y. BRENIER, N. J. MAUSER, AND M. PUEL, *Incompressible Euler and e-MHD as scaling limits of the Vlasov-Maxwell system*, Commun. Math. Sci., 1 (2003), pp. 437–447.
- [4] H. BRÉZIS, F. GOLSE, AND R. SENTIS, *Analyse asymptotique de l'équation de Poisson couplée à la relation de Boltzmann. Quasi-neutralité des plasmas*, C. R. Acad. Sci. Paris Sér. I Math., 321 (1995), pp. 953–959.
- [5] F. CHEN, *Introduction to Plasma Physics and Controlled Fusion*, Vol. 1, Plenum Press, New York, 1984.
- [6] G. Q. CHEN, J. W. JEROME, AND D. H. WANG, *Compressible Euler-Maxwell equations*, Transport Theory Statist. Phys., 29 (2000), pp. 311–331.
- [7] S. CORDIER AND E. GRENIER, *Quasineutral limit of an Euler-Poisson system arising from plasma physics*, Comm. Partial Differential Equations, 25 (2000), pp. 1099–1113.
- [8] E. GRENIER, *Pseudo-differential energy estimates of singular perturbations*, Comm. Pure Appl. Math., 50 (1997), pp. 821–865.
- [9] T. KATO, *Nonstationary flow of viscous and ideal fluids in \mathbb{R}^3* , J. Funct. Anal., 9 (1972), pp. 296–305.

- [10] A. KINGSEP, K. CHUKBAR, AND V. YANKOV, *Electron magnetohydrodynamics*, in Review of Plasma Physics 16, B. B. Kadomtsev, ed., Consultants Bureau, New York, 1990, pp. 243–291.
- [11] S. KLAINERMAN AND A. MAJDA, *Singular limits of quasilinear hyperbolic systems with large parameters and the incompressible limit of compressible fluids*, Comm. Pure Appl. Math., 34 (1981), pp. 481–524.
- [12] S. KLAINERMAN AND A. MAJDA, *Compressible and incompressible fluids*, Comm. Pure Appl. Math., 35 (1982), pp. 629–651.
- [13] H. O. KREISS, J. LORENZ, AND M. J. NAUGHTON, *Convergence of the solutions of the compressible to the solutions of the incompressible Navier-Stokes equations*, Adv. in Appl. Math., 12 (1991), pp. 187–214.
- [14] P. L. LIONS, *Mathematical Topics in Fluid Mechanics, Vol. 1: Incompressible Models*, Oxford Lecture Ser. Math. Appl. 3, Oxford University Press, New York, 1996.
- [15] A. MAJDA, *Compressible Fluid Flow and Systems of Conservation Laws in Several Space Variables*, Springer-Verlag, New York, 1984.
- [16] F. J. MCGRATH, *Nonstationary plane flow of viscous and ideal fluids*, Arch. Rational Mech. Anal., 27 (1968), pp. 229–348.
- [17] Y. J. PENG AND S. WANG, *Convergence of compressible Euler-Maxwell equations to compressible Euler-Poisson equations*, Chinese Ann. Math. Ser. B, 28 (2007), pp. 583–602.
- [18] Y. J. PENG AND S. WANG, *Convergence of compressible Euler-Maxwell equations to incompressible Euler equations*, Comm. Partial Differential Equations, 33 (2008), pp. 349–376.
- [19] Y. J. PENG AND S. WANG, *Asymptotic expansions in two-fluid compressible Euler-Maxwell equations with small parameters*, Discrete Contin. Dyn. Syst., to appear.
- [20] Y. J. PENG AND Y. G. WANG, *Convergence of compressible Euler-Poisson equations to incompressible type Euler equations*, Asymptot. Anal., 41 (2005), pp. 141–160.
- [21] H. RISHBETH AND O. K. GARRIOTT, *Introduction to Ionospheric Physics*, Academic Press, New York, 1969.
- [22] S. SCHOCHET, *Fast singular limits of hyperbolic PDEs*, J. Differential Equations, 114 (1994), pp. 476–512.
- [23] M. SLEMROD AND N. STERNBERG, *Quasi-neutral limit for the Euler-Poisson system*, J. Nonlinear Sci., 11 (2001), pp. 193–209.
- [24] S. WANG, *Quasineutral limit of Euler-Poisson system with and without viscosity*, Comm. Partial Differential Equations, 29 (2004), pp. 419–456.
- [25] S. WANG AND S. JIANG, *The convergence of the Navier-Stokes-Poisson system to the incompressible Euler equations*, Comm. Partial Differential Equations, 31 (2006), pp. 571–591.

SYMMETRY-BREAKING BIFURCATION IN NONLINEAR SCHRÖDINGER/GROSS–PITAEVSKII EQUATIONS*

E. W. KIRR[†], P. G. KEVREKIDIS[‡], E. SHLIZERMAN[§], AND M. I. WEINSTEIN[¶]

Abstract. We consider a class of nonlinear Schrödinger/Gross–Pitaevskii (NLS-GP) equations, i.e., NLS with a linear potential. NLS-GP plays an important role in the mathematical modeling of nonlinear optical as well as macroscopic quantum phenomena (BEC). We obtain conditions for a symmetry-breaking bifurcation in a symmetric family of states as \mathcal{N} , the squared L^2 norm (particle number, optical power), is increased. The bifurcating asymmetric state is a “mixed mode” which, near the bifurcation point, is approximately a superposition of symmetric and antisymmetric modes. In the special case where the linear potential is a double well with well-separation L , we estimate $\mathcal{N}_{cr}(L)$, the symmetry breaking threshold. Along the “lowest energy” symmetric branch, there is an exchange of stability from the symmetric to the asymmetric branch as \mathcal{N} is increased beyond \mathcal{N}_{cr} .

Key words. nonlinear Schrödinger, Gross–Pitaevskii, soliton, bound state

AMS subject classifications. 35Q55, 37K45, 37K50

DOI. 10.1137/060678427

1. Introduction. Symmetry breaking is a ubiquitous and important phenomenon which arises in a wide range of physical systems. In this paper, we consider a class of partial differential equations (PDEs), which are invariant under a symmetry group. For sufficiently small values of a parameter, \mathcal{N} , the preferred (dynamically stable) stationary state of the system is invariant under this symmetry group. However, above a critical parameter, \mathcal{N}_{cr} , although the group-invariant state persists, the new preferred state of the system is a state which (i) exists only for $\mathcal{N} > \mathcal{N}_{cr}$ and (ii) is no longer invariant. That is, symmetry is broken and there is an exchange of stability.

Physical examples of symmetry breaking include liquid crystals [31], quantum dots [34], semiconductor lasers [13], and pattern dynamics [28]. This article focuses on spontaneous symmetry breaking as a phenomenon in nonlinear optics [4, 20, 18], as well as in the macroscopic quantum setting of Bose–Einstein condensation (BEC) [1]. Here, the governing equations are PDEs of nonlinear Schrödinger/Gross–Pitaevskii (NLS-GP) type. Symmetry breaking has been observed experimentally in optics for two-component spatial optical vector solitons (i.e., for self-guided laser beams in Kerr media and focusing cubic nonlinearities) in [4], as well as for the electric field distribution between multiple wells of a photorefractive crystal in [20, 18]. In BECs,

*Received by the editors December 22, 2006; accepted for publication (in revised form) July 25, 2007; published electronically June 6, 2008.

<http://www.siam.org/journals/sima/40-2/67842.html>

[†]Department of Mathematics, University of Illinois, Urbana-Champaign, Urbana, IL 61801 (ekirr@math.uiuc.edu). The work of this author was partially supported by grants DMS-0405921 and DMS-060372.

[‡]Department of Mathematics and Statistics, University of Massachusetts, Amherst, MA 01003 (kevrekid@math.umass.edu). The work of this author was partially supported by DMS-0204585, NSF-CAREER, and DMS-0505663.

[§]Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot, Israel (eli.shlizerman@weizmann.ac.il). Part of this research was done while this author was a visiting graduate student in the Department of Applied Physics and Applied Mathematics at Columbia University.

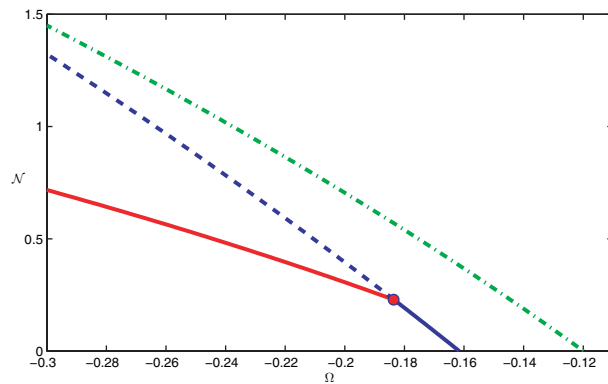
[¶]Department of Applied Physics and Applied Mathematics, Columbia University, New York, NY 10027 (miw2103@columbia.edu). The work of this author was partially supported by DMS-0412305 and DMS-0530853.

an effective double well formed by a combined (parabolic) magnetic trapping and a (periodic) optical trapping of the atoms may have similar effects [1] and lead to “macroscopic quantum self-trapping.”

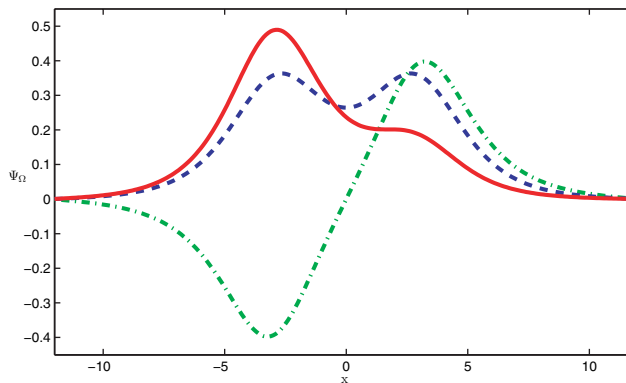
Symmetry breaking in ground states of the three-dimensional NLS-GP equation, with an attractive nonlinearity of Hartree type and a symmetric double-well linear potential, was considered in Aschbacher et. al. [3]; see also Remark 2.2. An example of a double-well potential is the function $V_L(x)$, plotted for the one-dimensional case, in Figure 2b of Example 2.1. The spectral properties of $H_L = -\Delta + V_L(x)$ are discussed in the appendix in section 8. For L sufficiently large, H_L has (at least) two eigenpairs, (Ω_0, ψ_0) and (Ω_1, ψ_1) . The normalized ground state eigenfunction ψ_0 is of even parity in x (“symmetric”) and the normalized excited state ψ_1 is of odd parity in x (“antisymmetric”). Moreover, $\psi_0(x) \sim 2^{-\frac{1}{2}}(\psi_\omega(x+L) + \psi_\omega(x-L))$ and $\psi_1(x) \sim 2^{-\frac{1}{2}}(\psi_\omega(x+L) - \psi_\omega(x-L))$.

Ground states of NLS-GP are positive nonlinear bound states, arising as *minimizers* of \mathcal{H} , the NLS-GP Hamiltonian energy subject to fixed \mathcal{N} , the squared L^2 norm. For the class of equations considered in [3], ground states exist for any $\mathcal{N} > 0$. It is proved that for sufficiently large \mathcal{N} , any ground state is concentrated in only one of the wells, i.e., symmetry is broken. The analysis in [3] is an asymptotic study of a variational problem (see Remark 2.2) showing that if \mathcal{N} is sufficiently large, then it is energetically preferable for the ground state to localize in a single well. In contrast, for small L^2 norm the ground state is even and bimodal having the symmetries of the ground state of the *linear* Schrödinger operator with symmetric double-well potential. For macroscopic quantum systems, the squared L^2 norm, denoted by \mathcal{N} , is the particle number, while in optics it is the optical power. An attractive nonlinearity corresponds to the case of negative scattering length in BEC and focusing Kerr nonlinearity in optics.

An alternative approach to symmetry breaking in NLS-GP is via bifurcation theory. It follows from [26, 25] that a family of “nonlinear ground states” bifurcates from the zero solution ($\mathcal{N} = 0$) at the ground state energy of the Schrödinger operator with a linear double-well potential. This nonlinear ground state branch consists of states having the same bimodal symmetry of the linear ground state. In this article we prove, under suitable conditions, that there is a secondary bifurcation to an asymmetric state at a critical $\mathcal{N} = \mathcal{N}_{cr} > 0$. The bifurcating asymmetric state is a “mixed mode” (see, for example, [8]) which, near the bifurcation point, is approximately a superposition of symmetric and antisymmetric modes. As \mathcal{N} is increased beyond \mathcal{N}_{cr} , the asymmetric state tends to concentrate in only one of the two wells. Since the double-well potential is symmetric, the bifurcating state is doubly degenerate; there are two families of asymmetric states associated with energy concentration in each of the wells. Figure 1 shows a typical bifurcation diagram demonstrating symmetry breaking for the NLS-GP system with a double-well potential. At the bifurcation point \mathcal{N}_{cr} (marked by a circle in Figure 1(a)), the symmetric (even with respect to coordinate x_1) ground state becomes unstable and a stable asymmetric (neither even nor odd parity) “mixed (degenerate) state” emanates from it. Further, we show that there is a transfer or exchange of stability which takes place at \mathcal{N}_{cr} ; for $\mathcal{N} < \mathcal{N}_{cr}$ the symmetric state is stable, while for $\mathcal{N} > \mathcal{N}_{cr}$ the asymmetric state is stable. Since our method is based on local bifurcation analysis we do not require that the states we consider satisfy a minimization principle, as in [3]. Thus, quite generally, symmetry breaking occurs as a consequence of the (finite-dimensional) *normal form* arising in systems with certain symmetry properties. Although we can treat a large class of problems for which there is no minimization principle, our analysis, at present, is



(a)



(b)

FIG. 1. (a) Bifurcation diagram for bound state solutions $\psi(x, t) = e^{-i\Omega t} \Psi_\Omega(x)$ of NLS-GP equation (2.1), with double-well potential (6.1) and cubic nonlinearity. Double-well parameters are $s = 1, L = 6$. Ω denotes the (nonlinear) frequency and $\mathcal{N} = \mathcal{N}[\Psi_\Omega]$, the squared L^2 norm of Ψ_Ω (optical power or particle number). The first bifurcation is from the zero state at the ground state energy of the double well. This state is an even function of x (symmetric). A secondary bifurcation, to an asymmetric state, at $\mathcal{N} = \mathcal{N}_{cr}$, is marked by a (red) circle. For $\mathcal{N} < \mathcal{N}_{cr}$ the symmetric state ((blue) solid line) is nonlinearly dynamically stable. For $\mathcal{N} > \mathcal{N}_{cr}$ the symmetric state is unstable ((blue) dashed line). The stable asymmetric state, appearing for $\mathcal{N} > \mathcal{N}_{cr}$, is marked by a (red) solid line. The (unstable) odd in x (antisymmetric) state is marked by a (green) dashed-dotted line. (b) Bound state solutions $\Psi_\Omega(x)$ plotted for a level set of $\mathcal{N}[\Psi_\Omega] = 0.5$.

restricted to small norm. As we shall see, a bifurcation occurring at small norm can be ensured, for example, by taking the distance between wells in the double well to be sufficiently large.

In [14] the precise transition point to symmetry breaking, \mathcal{N}_{cr} , of the ground state and the transfer of its stability to an asymmetric ground state was considered (by geometric dynamical systems methods) in the exactly solvable NLS-GP, with a double-well potential consisting of two Dirac delta functions separated by a distance L . Additionally, the behavior of the function $\mathcal{N}_{cr}(L)$ was considered. Another solvable model was examined by numerical means in [23]. A study of dynamics for nonlinear double wells appeared in [27].

We study $\mathcal{N}_{cr}(L)$ in general. The value at which symmetry breaking occurs, $\mathcal{N}_{cr}(L)$, is closely related to the spectral properties of the linearization of NLS-GP

about the symmetric branch. Indeed, so long as the linearization of NLS-GP at the symmetric state has no nonsymmetric null space, the symmetric state is locally unique by the implicit function theorem [24]. The mechanism for symmetry breaking is the first appearance of an antisymmetric element in the null space of the linearization for some $\mathcal{N} = \mathcal{N}_{cr}$. This is demonstrated for a finite-dimensional Galerkin approximation of NLS-GP in [20, 17]. The present work extends and generalizes this analysis to the full infinite-dimensional problem using the Lyapunov–Schmidt method [24]. Control of the corrections to the finite-dimensional approximation requires small norm of the states considered. Since, as anticipated by the Galerkin approximation, \mathcal{N}_{cr} is proportional to the distance between the lowest eigenvalues of the double well, which is exponentially small in L , our results apply to double wells with separation L , for L sufficiently large.

Our bifurcation results can be viewed as a very detailed study of a class of $O(2) \times \mathbb{Z}_2$ symmetric dynamical systems with four degrees of freedom in which there is a bifurcation of mixed states. We know of no classification of such systems. We show, under the nondegeneracy condition (4.3), that they are reducible to systems with two degrees of freedom and $\mathbb{Z}_2 \times \mathbb{Z}_2$ symmetry which have been classified in, for example, [8, 11]. In addition to being a self-contained treatment, our study of the nature of the solution set on the manifold of constant \mathcal{N} , as \mathcal{N} is varied, shows the key role of \mathcal{N} , a physical NLS-GP time-invariant, in determining the bifurcation diagram, which encodes the dynamic stability properties.

The article is organized as follows. In section 2 we introduce the NLS-GP model and give a technical formulation of the bifurcation problem. In section 3 we study a finite-dimensional truncation of the bifurcation problem, identifying a relevant bifurcation point. In section 4, we prove the persistence of this symmetry-breaking bifurcation in the full NLS-GP problem for $\mathcal{N} \geq \mathcal{N}_{cr}$. Moreover, in section 5 we show that the lowest energy symmetric state becomes dynamically unstable at \mathcal{N}_{cr} and the bifurcating asymmetric state is the dynamically stable ground state for $\mathcal{N} > \mathcal{N}_{cr}$.

The main results are stated in Theorem 4.1, Corollary 4.1, and Theorem 5.1. In particular, we obtain an asymptotic formula for the critical particle number (optical power) for symmetry breaking in NLS-GP,

$$(1.1) \quad \mathcal{N}_{cr} = \frac{\Omega_1 - \Omega_0}{\Xi[\psi_0, \psi_1]} + \mathcal{O}\left(\frac{(\Omega_1 - \Omega_0)^2}{\Xi[\psi_0, \psi_1]^3}\right).$$

Here, (Ω_0, ψ_0) and (Ω_1, ψ_1) are eigenvalue-eigenfunction pairs of the *linear* Schrödinger operator $H = -\Delta + V$, where Ω_0 and Ω_1 are separated from other spectrum and Ξ is a positive constant, given by (4.2), depending on ψ_0 and ψ_1 . (Section 8 is an appendix where we discuss the basic spectral properties of $-\Delta + V$, where V is a double-well potential.) *The most important case is where $\Omega_0 < \Omega_1$ are the lowest two energies (linear ground and first excited states).* For double wells with separation L , we have $\mathcal{N}_{cr} = \mathcal{N}_{cr}(L)$, depending on the eigenvalue spacing $\Omega_0(L) - \Omega_1(L)$, which is exponentially small if L is large and Ξ is of order one. Thus, for large L , the bifurcation occurs at small L^2 norm. This is the weakly nonlinear regime in which the corrections to the finite-dimensional model can be controlled perturbatively. A local bifurcation diagram of this type will occur for *any* simple even-odd symmetric pair of simple eigenvalues of H in the weakly nonlinear regime, so long as the eigenfrequencies are separated from the rest of the spectrum of H ; see Proposition 4.1 and the gap condition (4.13). Therefore, a similar phenomenon occurs for higher order, nearly degenerate pairs of eigenstates of the double wells, arising from isolated single

wells with multiple eigenstates. Section 6 contains numerical results validating our theoretical analysis.

2. Technical formulation. Consider the initial value problem for the time-dependent NLS-GP equation

$$(2.1) \quad i\partial_t \psi = H\psi + g(x)K[\psi\bar{\psi}] \psi, \quad \psi(x, 0) \text{ given,}$$

$$(2.2) \quad H = -\Delta + V(x).$$

We assume the following.

(H1) The initial value problem for NLS-GP is well-posed in the space $C^0([0, \infty); H^1(\mathbb{R}^n))$.

(H2) The potential $V(x)$ is assumed to be real valued, smooth, and rapidly decaying as $|x| \rightarrow \infty$. We also assume symmetry with respect to a hyperplane which, without loss of generality, can be taken to be $\{x_1 = 0\}$:

$$(2.3) \quad V(x_1, x_2, \dots, x_n) = V(-x_1, x_2, \dots, x_n).$$

We assume the nonlinear term, $K[\psi\bar{\psi}]$, to be attractive, cubic (local or non-local), and symmetric with respect to the same hyperplane. Specifically, we assume the following hypotheses on the nonlinear term:

(H3) (a) $g(x_1, x_2, \dots, x_n) = g(-x_1, x_2, \dots, x_n)$ (symmetry).

(b) $g(x) < 0$ (attractive/focusing).

(c) $K[h] = \int K(x-y)h(y)dy$,
 $K(x_1, x_2, \dots, x_n) = K(-x_1, x_2, \dots, x_n)$, $K > 0$.

(d) Consider the map $N : H^2 \times H^2 \times H^2 \mapsto L^2$ defined by

$$(2.4) \quad N(\phi_1, \phi_2, \phi_3) = gK[\phi_1\bar{\phi}_2]\phi_3.$$

We also write $N(u) = N(u, u, u)$ and note that $\partial_u N(u) = N(\cdot, u, u) + N(u, \cdot, u) + N(u, u, \cdot)$. We assume that there exists a constant $k > 0$ such that

$$(2.5) \quad \|N(\phi_0, \phi_1, \phi_2)\|_{L^2} \leq k \|\phi_1\|_{H^2} \|\phi_2\|_{H^2} \|\phi_3\|_{H^2}.$$

Remark 2.1. The symmetry restrictions in (H3) insure that the nonlinear term is equivariant with respect to the action of $O(2) \times \mathbb{Z}_2$, namely,

$$\begin{aligned} N(e^{i\theta} u) &= e^{i\theta} N(u), \quad 0 \leq \theta < 2\pi, \\ N(\bar{u}) &= \overline{N(u)}, \\ N(Ru) &= RN(u), \end{aligned}$$

where R is the reflection with respect to the $x_1 = 0$ hyperplane. Although, the nonlinearity taken in (H3) is cubic, our analysis can be easily extended to more general cases satisfying the $O(2) \times \mathbb{Z}_2$ equivariance.

We now give several illustrative and important examples of NLS-GP.

Example 1. GP equation for BECs with negative scattering length $g(x) \equiv -1$, $K(x) = \delta(x)$.

Example 2. NLS equation for optical media with a nonlocal kernel $g(x) \equiv \pm 1$, $K(x) = A \exp(-x^2/\sigma^2)$ [22]. See also [6] for similar considerations in BECs.

Example 3. Photorefractive (saturable) nonlinearities and optically induced potentials [7]. The relevant symmetry breaking phenomenology is experimentally observable, as shown in [20]. Also contained therein is the finite-dimensional Galerkin approach of section 3.

Nonlinear bound states. Nonlinear bound states are solutions of NLS-GP of the form

$$(2.6) \quad \psi(x, t) = e^{-i\Omega t} \Psi_\Omega(x),$$

where $\Psi_\Omega \in H^1(\mathbb{R}^n)$ solves

$$(2.7) \quad H \Psi_\Omega + g(x) K[|\Psi_\Omega|^2] \Psi_\Omega - \Omega \Psi_\Omega = 0, \quad \Psi_\Omega \in H^1.$$

If the potential $V(x)$ is such that the operator $H = -\Delta + V(x)$ has a discrete eigenvalue, Ω_* , and a corresponding eigenstate ψ_* , then for energies Ω near Ω_* , one expects small amplitude *nonlinear* bound states, which are to leading order small multiples of ψ_* . This is the standard setting of bifurcation from a simple eigenvalue [24], which follows from the implicit function theorem.

THEOREM 2.1 (see [21, 25, 26]). *Let $(\Psi, \Omega) = (\psi_*, \Omega_*)$ be a simple eigenpair of the eigenvalue problem $H\Psi = \Omega\Psi$, i.e., $\dim\{\rho : (H - \Omega_*)\rho = 0\} = 1$. Then there exists a unique smooth curve of nontrivial solutions $\alpha \mapsto (\Psi(\cdot; \alpha), \Omega(\alpha))$, defined in a neighborhood of $\alpha = 0$, such that*

$$(2.8) \quad \Psi_\Omega = \alpha (\psi_* + \mathcal{O}(|\alpha|^2)), \quad \Omega = \Omega_* + \mathcal{O}(|\alpha|^2), \quad \alpha \rightarrow 0.$$

Remark 2.2. For a large class of problems, a nonlinear ground state can be characterized variationally as a constrained minimum of the NLS-GP energy subject to fixed squared L^2 norm. Define the NLS-GP Hamiltonian energy functional

$$(2.9) \quad \mathcal{H}_{NLS-GP}[\Phi] \equiv \int |\nabla\Phi|^2 + V|\Phi|^2 dy + \frac{1}{2} \int g(y)|\Phi(y)|^2 K[|\Phi|^2] dy$$

and the particle number (optical power)

$$(2.10) \quad \mathcal{N}[\Phi] = \int |\Phi|^2 dy.$$

In particular, the following can be proved.

THEOREM 2.2. *Let $I_\lambda = \inf_{\mathcal{N}[f]=\lambda} \mathcal{H}[f]$. If $-\infty < I_\lambda < 0$, then the minimum is attained at a positive solution of (2.7). Here, $\Omega = \Omega(\lambda)$ is a Lagrange multiplier for the constrained variational problem.*

In [3] the nonlinear Hartree equation is studied; $K[h] = |y|^{-1} \star h$, $g \equiv -1$. It is proved that if $V(x)$ is a double-well potential, then for λ sufficiently large, the minimizer does not have the same symmetry as the linear ground state. By uniqueness, ensured by the implicit function theorem, for small \mathcal{N} , the minimizer has the same symmetry as that as the linear ground state and has the expansion (2.8); see [3] and section 4.

We make the following assumptions.

Spectral assumptions on H .

- (H4) H has a pair of simple eigenvalues Ω_0 and Ω_1 . ψ_0 and ψ_1 , the corresponding (real-valued) eigenfunctions are, respectively, even (symmetric) and odd (anti-symmetric) in x_1 .

Example 2.1 (the basic example: *double-well potentials*). Our basic example of $V(x)$ is a double-well potential consisting of two identical potential wells and separated by a distance L . Thus, assume symmetry with respect to the hyperplane which, without loss of generality, can be taken to be $\{x_1 = 0\}$:

$$(2.11) \quad V(x_1, x_2, \dots, x_n) = V(-x_1, x_2, \dots, x_n).$$

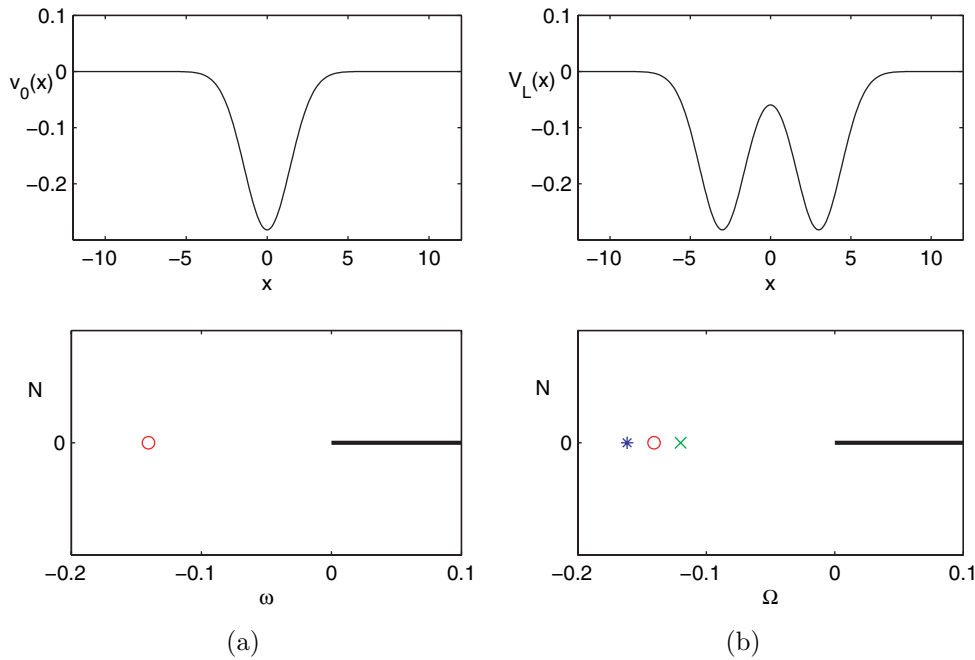


FIG. 2. This figure demonstrates a single- and a double-well potential and the spectrum of H and H_L , respectively. (a) shows a single-well potential and under it the spectrum of H , with an eigenvalue marked by a (red) mark “o” at ω and continuous spectrum marked by a (black) line for energies $\omega \geq 0$. (b) shows the double well centered at $\pm L$ and the spectrum of H_L underneath. The eigenvalues Ω_0 and Ω_1 are each marked by a (blue) mark “*” and a (green) mark “x,” respectively, on either side of the location ω —(red) mark “o.” The continuous spectrum is marked by a (black) line for energies $\Omega \geq 0$.

In the appendix, we discuss how spectral hypothesis (H4) is shown for a double wells with well-separation parameter L . The simplest example, in one space dimension, is obtained as follows. Start with a single potential well (rapidly decaying as $|x| \rightarrow \infty$), $v_0(x)$, having exactly one eigenvalue, ω , $H_0\psi_\omega = (-\Delta + v_0(x))\psi_\omega = \omega\psi_\omega$; see Figure 2(a). Center this well at $x = -L$ and place an identical well centered at $x = L$. Denote by $V_L(x)$ the resulting *double-well potential* and by H_L the Schrödinger operator:

$$(2.12) \quad H_L = -\Delta + V_L(x).$$

There exists $L > L_0$ such that for $L > L_0$, H_L has a pair of eigenvalues $\Omega_0 = \Omega_0(L)$ and $\Omega_1 = \Omega_1(L)$, $\Omega_0 < \Omega_1$, and corresponding eigenfunctions ψ_0 and ψ_1 ; see Figure 2(b). ψ_0 is symmetric with respect to $x = 0$ and ψ_1 is antisymmetric with respect to $x = 0$. Moreover, for L sufficiently large, $|\Omega_0 - \Omega_1| = \mathcal{O}(e^{-\kappa L})$, $\kappa > 0$; see [12] and see also the appendix.

The construction can be generalized. If $-\Delta + v_0(x)$ has m bound states, then forming a double well V_L with L sufficiently large, $H_L = -\Delta + V_L$ will have m pairs of eigenvalues: $(\Omega_{2j}, \Omega_{2j+1})$, $j = 0, \dots, m - 1$, eigenfunctions ψ_{2j} (symmetric), and ψ_{2j+1} antisymmetric.

By Theorem 2.1, for small \mathcal{N} , there exists a unique nontrivial nonlinear bound state bifurcating from the zero solution at the ground state energy, Ω_0 , of H . By uniqueness, ensured by the implicit function theorem, these small amplitude nonlinear bound states have the same symmetries as the double well; they are bimodal. We also

know from [3] that for sufficiently large \mathcal{N} the ground state has broken symmetry. We now seek to elucidate the transition from the regime of \mathcal{N} small to \mathcal{N} large.

We work in the general setting of hypotheses (H1)–(H4). Define spectral projections onto the bound and continuous spectral parts of H :

$$(2.13) \quad P_0 = (\psi_0, \cdot)\psi_0, \quad P_1 = (\psi_1, \cdot)\psi_1, \quad \tilde{P} = I - P_0 - P_1.$$

Here,

$$(2.14) \quad (f, g) = \int \bar{f}g \, dx.$$

We decompose the solutions of (2.7) according to

$$(2.15) \quad \Psi_\Omega = c_0\psi_0 + c_1\psi_1 + \eta, \quad \eta = \tilde{P}\eta.$$

We next substitute the expression (2.15) into (2.7) and then act with projections P_0 , P_1 , and \tilde{P} on the resulting equation. Using the symmetry and antisymmetry properties of the eigenstates, we obtain three equations which are equivalent to the PDE (2.7):

$$(2.16) \quad (\Omega_0 - \Omega) c_0 + a_{0000}|c_0|^2 c_0 + (a_{0110} + a_{0011}) |c_1|^2 c_0, \\ + a_{0011} c_1^2 \bar{c}_0 + (\psi_0 g, \mathcal{R}(c_0, c_1, \eta)) = 0,$$

$$(2.17) \quad (\Omega_1 - \Omega) c_1 + a_{1111}|c_1|^2 c_1 + (a_{1010} + a_{1001}) |c_0|^2 c_1 \\ + a_{1010} c_0^2 \bar{c}_1 + (\psi_1 g, \mathcal{R}(c_0, c_1, \eta)) = 0,$$

$$(2.18) \quad (H - \Omega) \eta = -\tilde{P} g [F(\cdot; c_0, c_1) + \mathcal{R}(c_0, c_1, \eta)],$$

where $F(\cdot, c_0, c_1)$ is independent of η and $\mathcal{R}(c_0, c_1, \eta)$ contains linear, quadratic, and cubic terms in η . The coefficients a_{klmn} are defined by

$$(2.19) \quad a_{klmn} = (\psi_k, gK[\psi_l \bar{\psi}_m] \psi_n).$$

We shall study the character of the set of solutions of the system (2.16)–(2.18) restricted to the level set

$$(2.20) \quad \int |\Psi_\Omega|^2 dx = \mathcal{N} \iff |c_0|^2 + |c_1|^2 + \int |\eta|^2 dx = \mathcal{N}$$

as \mathcal{N} varies.

Let Ω_0 and Ω_1 denote the two lowest eigenvalues of H_L . We prove (Theorem 4.1, Corollary 4.1, Theorem 5.1) the following.

- There exist two solution branches, parametrized by \mathcal{N} , which bifurcate from the zero solution at the eigenvalues Ω_0 and Ω_1 .
- Along the branch (Ω, Ψ_Ω) , emanating from the solution $(\Omega = \Omega_0, \Psi = 0)$, there is a symmetry-breaking bifurcation at $\mathcal{N} = \mathcal{N}_{crit} > 0$. In particular, let u_{crit} denote the solution of (2.7) corresponding to the value $\mathcal{N} = \mathcal{N}_{crit}$. Then, in a neighborhood u_{crit} , for $\mathcal{N} < \mathcal{N}_{crit}$ there is only one solution of (2.7), the symmetric ground state, while for $\mathcal{N} > \mathcal{N}_{crit}$ there are two solutions one symmetric and a second asymmetric.
- *Exchange of stability* at the bifurcation point: For $\mathcal{N} < \mathcal{N}_{crit}$ the symmetric state is dynamically stable, while for $\mathcal{N} > \mathcal{N}_{crit}$ the asymmetric state is stable and the symmetric state is exponentially unstable.

3. Bifurcations in a finite-dimensional approximation. It is illustrative to consider the finite-dimensional approximation to the system (2.16)–(2.18), obtained by neglecting the continuous spectral part, η . Let us first set $\eta = 0$, and therefore $\mathcal{R}(c_0, c_1, 0) = 0$. Under this assumption of no coupling to the continuous spectral part of H , we obtain the finite-dimensional system

$$(3.1) \quad (\Omega_0 - \Omega) c_0 + a_{0000} |c_0|^2 c_0 + (a_{0110} + a_{0011}) |c_1|^2 c_0 + a_{0011} c_1^2 \bar{c}_0 = 0,$$

$$(3.2) \quad (\Omega_1 - \Omega) c_1 + a_{1111} |c_1|^2 c_1 + (a_{1010} + a_{1001}) |c_0|^2 c_1 + a_{1010} c_0^2 \bar{c}_1 = 0.$$

Note that in this approximation the physically meaningful parameter \mathcal{N} is given by

$$(3.3) \quad |c_0|^2 + |c_1|^2 = \mathcal{N}.$$

Our strategy is to first analyze the bifurcation problem for this approximate finite-dimensional system of *algebraic equations*. We then treat the corrections, coming from coupling to the continuous spectral part of H , η , perturbatively. In fact, the analysis in the next section indirectly shows that (3.1)–(3.2) is the universal unfolding of bifurcation problems with four degrees of freedom and $O(2) \times \mathbb{Z}_2$ symmetry; see [8] for related notions.¹

Our point of view, however, is not the *general* theory of (3.1)–(3.2), its bifurcations and its structural stability, but rather what it implies for solutions of the NLS-GP system for given “data”, e.g., spectral properties of the potential $V(x)$, the nonlinearity term defined in terms of $K[\cdot]$ and g . This data, together with structural properties of *NLS-GP*, such as its Hamiltonian structure and the second-order elliptic character of the nonlinear bound state equations, imply natural assumptions as well as constraints on the coefficients a_{jklm} . In this section, where we study the algebraic problem (3.1)–(3.2), we posit and comment briefly on these properties. The first are

$$(3.4) \quad a_{0000} - a_{1001} - 2a_{1010} > 0,$$

$$(3.5) \quad a_{1010} \neq 0.$$

Condition (3.4) will guarantee a symmetry-breaking bifurcation at small amplitude, while (3.5) is a nondegeneracy condition, easily seen to hold for generic assumptions on $V, K[\cdot]$ and g in NLS-GP.

Additional conditions are

$$(3.6) \quad a_{1001} - a_{0000} > 0$$

and

$$(3.7) \quad a_{1111} - a_{0110} - 2a_{0011} > 0,$$

$$(3.8) \quad a_{1111} - a_{0110} > 0,$$

which we shall see are related to secondary bifurcations and the exchange of stability at the bifurcation point.

Note also that

$$(3.9) \quad a_{0000} < 0, \quad a_{1111} < 0$$

since $K[\cdot]$ preserves positivity and $g < 0$; see (2.19) and hypothesis (H3).

¹Since the full problem, with $\eta \neq 0$, adds higher order terms in (3.1)–(3.2), from singularity theory we expect that they will not qualitatively change the bifurcation diagram. We do not make use of results in singularity theory instead we prefer a self-contained treatment yielding a detailed understanding of the effect of the physically important unfolding parameter $\Omega_1 - \Omega_0$.

One can easily construct, using (3.9), the following two *pure mode solutions* of (3.1)-(3.2):

(1) *Approximate nonlinear symmetric branch:*

$$c_1 = 0, \quad c_0 = \sqrt{\frac{\Omega - \Omega_0}{a_{0000}}} e^{i\theta}, \quad 0 \leq \theta < 2\pi, \quad \Omega \leq \Omega_0.$$

In terms of the physical parameter $\mathcal{N} : c_1 = 0, c_0 = \sqrt{\mathcal{N}} e^{i\theta}, \Omega = \Omega_0 + a_{0000}\mathcal{N}$, where $\mathcal{N} \geq 0, 0 \leq \theta < 2\pi$.

(2) *Approximate nonlinear antisymmetric branch:*

$$c_0 = 0, \quad c_1 = \sqrt{\frac{\Omega - \Omega_1}{a_{1111}}} e^{i\theta}, \quad 0 \leq \theta < 2\pi, \quad \Omega \leq \Omega_1.$$

In terms of the physical parameter $\mathcal{N} : c_0 = 0, c_1 = \sqrt{\mathcal{N}} e^{i\theta}, \Omega = \Omega_1 + a_{1111}\mathcal{N}$, where $\mathcal{N} \geq 0, 0 \leq \theta < 2\pi$.

The above two solutions correspond to those of NLS-GP given by Theorem 2.1.

Under assumptions (3.4)-(3.8) the system (3.1)-(3.2) has exactly one more (secondary) branch of solutions which bifurcates from the symmetric branch at

$$\Omega^* = \Omega_0 + \frac{a_{0000}(\Omega_1 - \Omega_0)}{a_{0000} - a_{1001} - 2a_{1010}} < \Omega_0$$

or, equivalently, at

$$(3.10) \quad \mathcal{N}_{cr} = \frac{\Omega_1 - \Omega_0}{a_{0000} - a_{1001} - 2a_{1010}} > 0.$$

In terms of the physical parameter \mathcal{N} , these *mixed mode solutions of the approximate system* are given by

(3) *Approximate nonlinear asymmetric branch:* For $\mathcal{N} \geq \mathcal{N}_{cr}, 0 \leq \theta_0, \theta_1 < 2\pi, \theta_1 - \theta_0 \in \{0, \pi\}$,

$$(3.11) \quad \begin{aligned} c_0 &= \sqrt{\frac{(a_{1111} - a_{0110} - 2a_{0011})\mathcal{N} + (\Omega_1 - \Omega_0)}{a_{1111} - a_{0110} - 2a_{0011} + a_{0000} - a_{1001} - 2a_{1010}}} e^{i\theta_0}, \\ c_1 &= \sqrt{\frac{(a_{0000} - a_{1001} - 2a_{1010})\mathcal{N} - (\Omega_1 - \Omega_0)}{a_{1111} - a_{0110} - 2a_{0011} + a_{0000} - a_{1001} - 2a_{1010}}} e^{i\theta_1}, \\ \Omega &= \Omega_0 + \\ &\frac{(a_{0000}a_{1111} - (a_{1001} + 2a_{1010})(a_{0110} + 2a_{0011}))\mathcal{N} - (a_{0000} - a_{1001} - 2a_{1010})(\Omega_1 - \Omega_0)}{a_{1111} - a_{0110} - 2a_{0011} + a_{0000} - a_{1001} - 2a_{1010}}. \end{aligned}$$

To derive the above solutions, first introduce polar coordinates: $c_j = \rho_j e^{i\theta_j}, j = 1, 2, \Delta\theta = \theta_1 - \theta_0$. Then, (3.1)-(3.2) becomes

$$(3.12) \quad \rho_0 [\Omega_0 - \Omega + a_{0000}\rho_0^2 + (a_{0110} + a_{0011} + a_{0011}e^{i2\Delta\theta})\rho_1^2] = 0,$$

$$(3.13) \quad \rho_1 [\Omega_1 - \Omega + a_{1111}\rho_1^2 + (a_{1001} + a_{1010} + a_{1010}e^{-i2\Delta\theta})\rho_0^2] = 0.$$

Taking the imaginary parts of both equations and using the fact that all coefficients are real, see (2.19), we get $a_{0011} \sin(2\Delta\theta) = 0 = a_{1010} \sin(2\Delta\theta)$. Because of (3.5), the latter are equivalent with $\Delta\theta \in \{0, \pi/2, \pi, 3\pi/2\}$ modulo 2π , note that $a_{0011} = a_{1010}$; see (2.19). The real parts of (3.12)-(3.13) become

$$(3.14) \quad \rho_0 [\Omega_0 - \Omega + a_{0000}\rho_0^2 + (a_{0110} + 2a_{0011})\rho_1^2] = 0,$$

$$(3.15) \quad \rho_1 [\Omega_1 - \Omega + a_{1111}\rho_1^2 + (a_{1001} + 2a_{1010})\rho_0^2] = 0$$

for $\Delta\theta \in \{0, \pi\}$, while for $\Delta\theta \in \{\pi/2, 3\pi/2\}$,

$$(3.16) \quad \rho_0 [\Omega_0 - \Omega + a_{0000}\rho_0^2 + a_{0110}\rho_1^2] = 0,$$

$$(3.17) \quad \rho_1 [\Omega_1 - \Omega + a_{1111}\rho_1^2 + a_{1001}\rho_0^2] = 0.$$

Both systems (3.14)-(3.15) and (3.16)-(3.17), respectively, can be easily analyzed and even reduced to a linear system by the change of variables: $\mathcal{P}_0 = \rho_0^2$, $\mathcal{P}_1 = \rho_1^2$. The first one, under assumption (3.4), gives the mixed mode solution (3) in addition to the pure mode solutions (1) and (2), and no other solutions provided (3.7) holds. The second system gives no other solutions provided (3.6), (3.8) hold. Thus we have that the finite-dimensional Galerkin approximation yields the qualitative structure of the NLS-GP bifurcation diagram of Figure 3. Dynamic stability can also be addressed in the context of the Galerkin-truncated time-dependent Hamiltonian dynamics. We defer our study of stability to our considerations of the full PDE problem in the coming sections.

The study of the systems (3.14)-(3.15) and (3.16)-(3.17) can be put in the context of singularity theory and bifurcations in systems with symmetry [8]. These systems exhibit a $\mathbb{Z}_2 \oplus \mathbb{Z}_2$ symmetry. Actually the argument in the next section shows that under the nondegeneracy condition (3.5), bifurcation problems with four degrees of freedom and $O(2) \times \mathbb{Z}_2$ symmetry are reducible to two distinct systems each with two degrees of freedom and $\mathbb{Z}_2 \oplus \mathbb{Z}_2$ symmetry. The latter have been classified, for example, in [8, Chapter X]. Under assumptions (3.4), (3.7), our first system (3.14)-(3.15) falls in case A subcase (1) [8, Chapter X] while under (3.6), (3.8) our second system falls in their case A subcases (2) or (3). Note that the above-cited book uses $\lambda = \Omega_0 - \Omega$ as bifurcation parameter and $\sigma = \Omega_0 - \Omega_1$ as the unfolding parameter which explains the differences between our bifurcation diagram in Figure 1 and their Figure 4.3 in Chapter X.

Finally, the posited properties of a_{klmn} can be directly verified for general double-well problems with L sufficiently large; see the proof of Corollary 4.1.

4. Bifurcation/symmetry breaking analysis of the PDE. In this section we prove the following theorem.

THEOREM 4.1 (symmetry breaking for NLS-GP). *Consider NLS-GP with hypotheses (H2)–(H4). Let a_{klmn} be given by (2.19) and let*

$$(4.1) \quad \Xi[\psi_0, \psi_1, g] \equiv a_{0000} - a_{1001} - 2a_{1010}$$

$$(4.2) \quad = (\psi_0^2, gK[\psi_0^2]) - (\psi_1^2, gK[\psi_0^2]) - 2(\psi_0\psi_1, gK[\psi_0\psi_1]) > 0,$$

$$(4.3) \quad a_{1010} = (\psi_0\psi_1, gK[\psi_0\psi_1]) \neq 0.$$

Assume

$$(4.4) \quad \frac{\Omega_1 - \Omega_0}{\Xi[\psi_0, \psi_1, g]^2} \text{ is sufficiently small.}$$

Then, there exists $\mathcal{N}_{cr} > 0$ such that

- (i) for any $\mathcal{N} \leq \mathcal{N}_{cr}$, there is (up to the symmetry $u \mapsto u e^{i\gamma}$) a locally unique symmetric state, $u_{\mathcal{N}}^{sym}$.
- (ii) $\mathcal{N} = \mathcal{N}_{cr}$, $u_{\mathcal{N}_{cr}}^{sym}$ is a bifurcation point. For $\mathcal{N} > \mathcal{N}_{cr}$, there are, in a neighborhood of $\mathcal{N} = \mathcal{N}_{cr}$, $u_{\mathcal{N}_{cr}}^{sym}$, two branches of solutions: (a) a continuation of the symmetric branch and (b) a new asymmetric branch.

(iii) The critical \mathcal{N} value for bifurcation is given approximately by

$$\mathcal{N}_{cr} = \frac{\Omega_1 - \Omega_0}{\Xi[\psi_0, \psi_1, g]} \left[1 + \mathcal{O} \left(\frac{\Omega_1 - \Omega_0}{\Xi[\psi_0, \psi_1, g]^2} \right) \right].$$

(iv) The bifurcation from the zero state at $\Omega = \Omega_0$ and the bifurcation described in (i)–(iii) are the only ones along this branch for \mathcal{N} small.

COROLLARY 4.1 (application to double wells, V_L , with separation L). Fix a pair of eigenvalues, (Ω_{2j}, ψ_{2j}) , $(\Omega_{2j+1}, \psi_{2j+1})$, of the linear double-well potential, $V_L(x)$; see Example 2.1. Assume that $K(x) \rightarrow 0$ as $|x| \rightarrow 0$ and $\lim_{L \rightarrow \pm\infty} g(x_1 + L, x_2, \dots) = g_{\pm}(x_2, \dots) < 0$. Then, for the NLS-GP with double-well potential of well-separation L , there exists $\tilde{L} > 0$ such that for all $L \geq \tilde{L}$, there is a symmetry-breaking bifurcation, as described in Theorem 4.1, with $\mathcal{N}_{cr} = \mathcal{N}_{cr}(L; j)$.

Proof of Corollary 4.1. The crux is to verify the spectral properties (H4), the nondegeneracy condition (4.3), and the smallness condition (4.4). The spectral properties of $H_L = -\Delta + V_L$ are handled in Proposition 8.1 of the appendix. In what follows we show that

$$(4.5) \quad \lim_{L \rightarrow \infty} a_{0000} - a_{1001} = 0,$$

$$(4.6) \quad \lim_{L \rightarrow \infty} a_{1010} < 0.$$

These imply both (4.3) and (4.4) since $\Omega_0 - \Omega_1 \rightarrow 0$ as $L \rightarrow \infty$; see Proposition 8.1.

For (4.5), (4.6) we use the large L asymptotic behavior of the eigenfunctions of H_L , given in Proposition 8.1, together with the continuity of the nonlinearity (2.5), to rewrite them as

$$(4.7) \quad \lim_{L \rightarrow \infty} a_{0000} - a_{1001} = \lim_{L \rightarrow \infty} \frac{1}{2} (T_L \psi_{\omega}, gK[(T_L \psi_{\omega} + RT_L \psi_{\omega})^2] RT_L \psi_{\omega}),$$

$$(4.8) \quad \lim_{L \rightarrow \infty} a_{1010} = \lim_{L \rightarrow \infty} \frac{1}{4} (T_L \psi_{\omega}^2 - RT_L \psi_{\omega}^2, gK[T_L \psi_{\omega}^2 - RT_L \psi_{\omega}^2]).$$

To obtain (4.5) approximate $\psi_{\omega} \in H^2$ by smooth functions with compact support ψ_{ω}^0 . For sufficiently large L , $T_L \psi_{\omega}^0$ and $RT_L \psi_{\omega}^0$ have no common support and the scalar product in (4.7) is zero. The continuity of the nonlinearity (2.5) implies (4.5).

Now, the right-hand side in (4.8) has two types of similar terms:

$$\lim_{L \rightarrow \infty} (T_L \psi_{\omega}^2, gK[RT_L \psi_{\omega}^2]) = \lim_{L \rightarrow \infty} (\psi_{\omega}^2, (T_{-L}g)T_{-2L}K[R\psi_{\omega}^2]) = 0,$$

$$\lim_{L \rightarrow \infty} (T_L \psi_{\omega}^2, gK[T_L \psi_{\omega}^2]) = \lim_{L \rightarrow \infty} (\psi_{\omega}^2, (T_{-L}g)K[\psi_{\omega}^2]) = (\psi_{\omega}^2, g_{-}K[\psi_{\omega}^2]) < 0,$$

where we used the symmetries of g and K and their behavior as $x_1 \rightarrow \pm\infty$. This completes the proof of the corollary. \square

To prove Theorem 4.1 we will establish that, under hypotheses (4.2)–(4.4), the character of the solution set (symmetry-breaking bifurcation) of the finite-dimensional

approximation (3.1)-(3.2) persists for the full (infinite-dimensional) problem:

$$(4.9) \quad (\Omega_0 - \Omega) c_0 + a_{0000}|c_0|^2 c_0 + (a_{0110} + a_{0011}) |c_1|^2 c_0 + a_{0011} c_1^2 \bar{c}_0 + (\psi_0 g, \mathcal{R}(c_0, c_1, \eta)) = 0,$$

$$(4.10) \quad (\Omega_1 - \Omega) c_1 + a_{1111}|c_1|^2 c_1 + (a_{1010} + a_{1001}) |c_0|^2 c_1 + a_{1010} c_0^2 \bar{c}_1 + (\psi_1 g, \mathcal{R}(c_0, c_1, \eta)) = 0,$$

$$(4.11) \quad (H - \Omega) \eta = -\tilde{P} g [F(\cdot; c_0, c_1) + \mathcal{R}(c_0, c_1, \eta)], \quad \eta = \tilde{P} \eta,$$

$$(4.12) \quad |c_0|^2 + |c_1|^2 + \int |\eta|^2 = \mathcal{N}.$$

We analyze this system using the Lyapunov–Schmidt-type method. The strategy is to solve (4.11) for η as a functional of c_0, c_1 , and Ω . Then, substituting $\eta = \eta[c_0, c_1, \Omega]$ into (4.9), (4.10), and (4.12), we obtain three closed equations, depending on a parameter \mathcal{N} , for c_0, c_1 , and Ω . This system is a perturbation of the finite-dimensional (truncated) system: (3.1)–(3.3). We then show that under hypotheses (4.2)–(4.4) there is a symmetry-breaking bifurcation. Finally, we show that the terms perturbing the finite-dimensional model have a small and controllable effect on the character of the solution set for a range of \mathcal{N} , which includes the bifurcation point. As seen in the proof, the form of the nonlinearity implies analyticity. However, analyticity is not essential for the arguments, at the root of which is the implicit function theorem, and the methods can be adapted to situations where one has a finite degree of smoothness.

We begin with the following proposition, which characterizes $\eta = \eta[c_0, c_1, \Omega]$.

PROPOSITION 4.1. *Consider (4.11) for η . By (H4) we have the following:*

$$(4.13) \quad \text{Gap condition : } |\Omega_j - \tau| \geq 2d_* \text{ for } j = 0, 1 \text{ and for all } \tau \in \sigma(H) \setminus \{\Omega_0, \Omega_1\}.$$

Then there exist $n_, r_* > 0$, depending on d_* , such that in the open set*

$$(4.14) \quad |c_0| + |c_1| < r_*,$$

$$\|c_0 \psi_0 + c_1 \psi_1 + \eta\|_{H^2} < n_*(d_*),$$

$$(4.15) \quad \text{dist}(\Omega, \sigma(H) \setminus \{\Omega_0, \Omega_1\}) > d_*,$$

the unique solution of (2.18) is given by the real-analytic mapping

$$(4.16) \quad (c_0, c_1, \Omega) \mapsto \eta[c_0, c_1, \Omega],$$

defined on the domain given by (4.14)-(4.15). Moreover, there exists $C_ > 0$ such that*

$$(4.17) \quad \|\eta[c_0, c_1, \Omega]\|_{H^2} \leq C_*(|c_0| + |c_1|)^3.$$

Proof. Consider the map

$$N : H^2 \times H^2 \times H^2 \mapsto L^2,$$

$$N(\phi_0, \phi_1, \phi_2) = gK[\phi_1 \bar{\phi}_2] \phi_3.$$

By assumptions on the nonlinearity (see section 2), there exists a constant $k > 0$ such that

$$(4.18) \quad \|N(\phi_0, \phi_1, \phi_2)\|_{L^2} \leq k \|\phi_1\|_{H^2} \|\phi_2\|_{H^2} \|\phi_3\|_{H^2}.$$

Moreover, since the map is real linear in each component, it is real analytic.²

²The trilinearity follows from the implicit bilinearity of K in formulas (2.16)–(2.18).

Let c_0, c_1 , and Ω be restricted according to (4.14)-(4.15). Equation (2.18) can be rewritten in the form

$$(4.19) \quad \eta + (H - \Omega)^{-1} \tilde{P}N[c_0\psi_0 + c_1\psi_1 + \eta] = 0.$$

Since the spectrum of $H\tilde{P}$ is bounded away from Ω by d_* , the resolvent

$$(H - \Omega)^{-1} \tilde{P} : L^2 \mapsto H^2$$

is a (complex) analytic map and bounded uniformly,

$$(4.20) \quad \|(H - \Omega)^{-1} \tilde{P}\|_{L^2 \mapsto H^2} \leq p(d_*^{-1}),$$

where $p(s) \rightarrow \infty$ as $s \rightarrow \infty$. Consequently the map $F : \mathbb{C}^2 \times \{\Omega \in \mathbb{C} : \text{dist}(\Omega, \sigma(H) \setminus \{\Omega_0, \Omega_1\}) \geq d_*\} \times H^2 \mapsto H^2$ given by

$$(4.21) \quad F(c_0, c_1, \Omega, \eta) = \eta + (H - \Omega)^{-1} \tilde{P}N[c_0\psi_0 + c_1\psi_1 + \eta]$$

is real analytic. Moreover,

$$F(0, 0, \Omega, 0) = 0, \quad D_\eta F(0, 0, \Omega, 0) = I.$$

Applying the implicit function theorem to (4.19), we have that there exists $n_*(\Omega), r_*(\Omega)$ such that whenever $|c_0| + |c_1| < r_*$ and $\|c_0\psi_0 + c_1\psi_1 + \eta\|_{H^2} < n_*$, (4.19) has a unique solution:

$$\eta = \eta(c_0, c_1, \Omega) \in H^2$$

which depends analytically on the parameters c_0, c_1, Ω . By applying the projection operator \tilde{P} to (4.19), which commutes with $(H - \Omega)^{-1}$, we immediately obtain $\tilde{P}\eta = \eta$, i.e., $\eta \in \tilde{P}L^2$.

We now show that n_*, r_* can be chosen independently of Ω , satisfying (4.15). The implicit function theorem can be applied in an open set for which

$$D_\eta F(c_0, c_1, \Omega, \eta) = I + (H - \Omega)^{-1} \tilde{P}D_\eta N[c_0\psi_0 + c_1\psi_1 + \eta]$$

is invertible. For this it suffices to have

$$\|(H - \Omega)^{-1} \tilde{P}D_\eta N[c_0\psi_0 + c_1\psi_1 + \eta]\|_{H^2} \leq Lip < 1.$$

A direct application of (2.5) and (4.20) shows that

$$(4.22) \quad \|(H - \Omega)^{-1} \tilde{P}D_\eta N[c_0\psi_0 + c_1\psi_1 + \eta]\|_{H^2} \leq 3k p(d_*^{-1}) \|c_0\psi_0 + c_1\psi_1 + \eta\|_{H^2}^2.$$

Fix $Lip = 3/4$. Then, a sufficient condition for invertibility is

$$(4.23) \quad 3k p(d_*^{-1}) \|c_0\psi_0 + c_1\psi_1 + \eta\|_{H^2}^2 \leq Lip = 3/4,$$

which allows us to choose $n_* = \frac{1}{2} \sqrt{\frac{1}{kp(d_*^{-1})}}$ independently of Ω .

But, if (4.23) holds, then, from (4.22), the H^2 operator

$$(H - \Omega)^{-1} \tilde{P}N[c_0\psi_0 + c_1\psi_1 + \cdot]$$

is Lipschitz with Lipschitz constant less than or equal to $Lip = 3/4$. The standard contraction principle estimate applied to (4.19) gives

$$(4.24) \quad \begin{aligned} \|\eta\|_{H^2} &\leq \frac{1}{1 - Lip} \|(H - \Omega)^{-1} \tilde{P}N[c_0\psi_0 + c_1\psi_1]\|_{H^2} \\ &\leq 4p(d_*^{-1}) k \|c_0\psi_0 + c_1\psi_1\|_{H^2}^3. \end{aligned}$$

Plugging the above estimate into (4.23) gives

$$\|c_0\psi_0 + c_1\psi_1\|_{H^2} + 4p(d_*^{-1}) k \|c_0\psi_0 + c_1\psi_1\|_{H^2}^3 \leq \frac{1}{2\sqrt{p(d_*^{-1})k}}.$$

Since the left-hand side is continuous in $(c_0, c_1) \in \mathbb{C}^2$ and zero for $c_0 = c_1 = 0$, one can construct $r_* > 0$ depending only on d_* , k such that the above inequality, hence (4.23) and (4.24), all hold whenever $|c_0| + |c_1| \leq r_*$. Finally, (4.17) now follows from (4.24). \square

In particular, for the double-well potential we have the following proposition.

PROPOSITION 4.2. *Let $V = V_L$ denote the double-well potential with well-separation L . There exist $L_* > 0, \varepsilon(L_*) > 0$, and $d_*(L_*) > 0$ such that for $L > L_*$, we have that for (c_0, c_1, Ω) satisfying $\text{dist}(\Omega, \sigma(H) \setminus \{\Omega_0, \Omega_1\}) \geq d_*(L_*)$ and $|c_0| + |c_1| < \varepsilon(L_*)$, $\eta[c_0, c_1, \Omega]$ is defined and analytic and satisfies the bound (4.17) for some $C_* > 0$.*

Proof. Since $\Omega_0, \Omega_1, \psi_0$, and ψ_1 can be controlled uniformly in L large, both d_* and r_* in the previous proposition can be controlled uniformly in L large. \square

Next we study the symmetries of $\eta(c_0, c_1, \Omega)$ and the properties of $\mathcal{R}(c_0, c_1, \eta)$ which we will use in analyzing (2.16)-(2.17). The following result is a direct consequence of the symmetries of (2.18) and Proposition 4.1.

PROPOSITION 4.3. *We have*

$$(4.25) \quad \eta(e^{i\theta}c_0, e^{i\theta}c_1, \Omega) = e^{i\theta}\eta(c_0, c_1, \Omega) \quad \text{for } 0 \leq \theta < 2\pi,$$

$$(4.26) \quad \overline{\eta(c_0, c_1, \Omega)} = \eta(\overline{c_0}, \overline{c_1}, \overline{\Omega}),$$

in particular

$$(4.27) \quad \eta(e^{i\theta}c_0, c_1 = 0, \Omega) = e^{i\theta}\eta(c_0, c_1 = 0, \Omega),$$

$$(4.28) \quad \eta(c_0 = 0, e^{i\theta}c_1, \Omega) = e^{i\theta}\eta(c_0 = 0, c_1, \Omega),$$

$\eta(c_0, 0, \Omega)$ is even in x_1 , $\eta(0, c_1, \Omega)$ is odd in x_1 , and if c_0, c_1 , and Ω are real-valued, then $\eta(c_0, c_1, \Omega)$ is real-valued.

In addition

$$(4.29) \quad \langle \psi_0, \mathcal{R}(c_0, c_1, \eta) \rangle = c_0 f_0(c_0, c_1, \Omega),$$

$$(4.30) \quad \langle \psi_1, \mathcal{R}(c_0, c_1, \eta) \rangle = c_1 f_1(c_0, c_1, \Omega),$$

where, for any $0 \leq \theta < 2\pi$ and $j = 0, 1$,

$$(4.31) \quad f_j(e^{i\theta}c_0, e^{i\theta}c_1, \Omega) = f_j(c_0, c_1, \Omega),$$

$$(4.32) \quad \overline{f_j(c_0, c_1, \Omega)} = f_j(\overline{c_0}, \overline{c_1}, \overline{\Omega}),$$

$$(4.33) \quad |f_j(c_0, c_1, \Omega)| \leq C(|c_0| + |c_1|)^4$$

for some constant $C > 0$. Moreover, both f_0 and f_1 can be written as absolutely convergent power series:

$$(4.34) \quad f_j(c_0, c_1, \Omega) = \sum_{\substack{k+l+m+n \geq 4, \\ k-l+m-n = 0, \\ m+n = \text{even}}} b_{klmn}^j(\Omega) c_0^k \bar{c}_0^l c_1^m \bar{c}_1^n, \quad j = 0, 1,$$

where $b_{klmn}^j(\Omega)$ are real valued when Ω is real valued. In particular, if c_0, c_1 , and Ω are real valued, then $f_j(c_0, c_1, \Omega)$ is real valued and, in polar coordinates, for $c_0, c_1 \neq 0$, we have

$$(4.35) \quad f_j(|c_0|, |c_1|, \Delta\theta, \Omega) = \sum_{k+m \geq 2, p \in \mathbb{Z}} b_{kmp}^j(\Omega) e^{ip2\Delta\theta} |c_0|^{2k} |c_1|^{2m}, \quad j = 0, 1,$$

where $\Delta\theta$ is the phase difference between $c_1 \in \mathbb{C}$ and $c_0 \in \mathbb{C}$.

Proof of Proposition 4.3. We start with (4.25) which clearly implies (4.27)-(4.28). We fix Ω and suppress dependence on it in subsequent notation. From (4.19) we have

$$\begin{aligned} & \eta(e^{i\theta} c_0, e^{i\theta} c_1) \\ &= -(H - \Omega)^{-1} \tilde{P} \\ & \quad N(e^{i\theta} c_0 \psi_0 + e^{i\theta} c_1 \psi_1 + \eta, e^{i\theta} c_0 \psi_0 + e^{i\theta} c_1 \psi_1 + \eta, e^{i\theta} c_0 \psi_0 + e^{i\theta} c_1 \psi_1 + \eta) \\ &= -(H - \Omega)^{-1} \tilde{P} \\ & \quad e^{i\theta} N(c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta, c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta, c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta), \end{aligned}$$

where we used

$$(4.36) \quad N(a\phi_1, b\phi_2, c\phi_3) = \bar{a}bcN(\phi_1, \phi_2, \phi_3).$$

Consequently,

$$e^{-i\theta} \eta(e^{i\theta} c_0, e^{i\theta} c_1) = -(H - \Omega)^{-1} \tilde{P} N[c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta, c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta, c_0 \psi_0 + c_1 \psi_1 + e^{-i\theta} \eta]$$

which shows that both $e^{-i\theta} \eta(e^{i\theta} c_0, e^{i\theta} c_1)$ and $\eta(c_0, c_1)$ satisfy the same equation (4.19). From the uniqueness of the solution proved in Proposition 4.1 we have the relation (4.25).

A similar argument (and use of the complex conjugate) leads to (4.26) and to the parities of $\eta(c_0, 0)$ and $\eta(0, c_1)$.

To prove (4.29) and (4.30), recall that

$$(4.37) \quad \begin{aligned} \mathcal{R}(c_0, c_1, \eta(c_0, c_1, \Omega)) &= N(c_0 \psi_0 + c_1 \psi_1 + \eta, c_0 \psi_0 + c_1 \psi_1 + \eta, c_0 \psi_0 + c_1 \psi_1 + \eta) \\ &\quad - N(c_0 \psi_0 + c_1 \psi_1, c_0 \psi_0 + c_1 \psi_1, c_0 \psi_0 + c_1 \psi_1). \end{aligned}$$

Consider first the case $c_1 = \rho_1 \in \mathbb{R}$. Note that

$$\langle \psi_1 g, \mathcal{R}(c_0, \rho_1 = 0, \eta(c_0, 0)) \rangle = 0.$$

Indeed, for $\rho_1 = 0$, all the functions in the arguments of \mathcal{R} are even functions (in x_1) making \mathcal{R} an even function. Since ψ_1 is odd we get that the above is the integral over the entire space of an odd function, and therefore equal to zero. Since

$\langle \psi_1, \mathcal{R}(c_0, \rho_1, \eta(c_0, \rho_1)) \rangle$ is analytic in $\rho_1 \in \mathbb{R}$ by the composition rule and its Taylor series starts with zero, we get (4.30) for real $c_1 = \rho_1$. To extend the result for complex values c_1 we use the rotational symmetry of \mathcal{R} , namely from (4.25), (4.36), and (4.37) we have

$$\mathcal{R}(e^{i\theta}c_0, e^{i\theta}c_1, \eta(e^{i\theta}c_0, e^{i\theta}c_1, \Omega)) = e^{i\theta}\mathcal{R}(c_0, c_1, \eta(c_0, c_1, \Omega)), \quad 0 \leq \theta < 2\pi,$$

hence (4.30) holds for $c_1 = |c_1|e^{-i\theta}$ by extending f_1 via the equality (4.31).

A similar argument holds for (4.29). Identity (4.33) follows from the definition of \mathcal{R} and (4.17) while identity (4.32) follows from (4.26).

We now turn to a proof of the expansions for f_j : (4.34) and (4.35). Note first that \mathcal{R} is real analytic because in (4.37) N is real linear in each variable and η is real analytic by Proposition 4.1. Hence, both f_0 and f_1 given by (4.29)-(4.30) are real analytic in c_0, c_1 and can be written in power series of the type (4.34). Estimate (4.33) implies that $k + l + m + n \geq 4$, while the rotational invariance (4.31) implies $k - l + m - n = 0$. The following parity argument shows why $m + n$ hence $m - n = l - k$ and $k + l$ are all even. Assume $m + n$ is odd. Note that because of (4.29), b_{klmn}^0 is the scalar product between an even function (in x_1) ψ_0 and the term in the power series of \mathcal{R} in which ψ_1 is repeated $m + n$ times. The latter is an odd function (in x_1) because ψ_1 is an odd function and it is repeated an odd number of times. The scalar product and hence b_{klmn}^0 for $m + n$ odd will be zero. A similar argument holds for b_{klmn}^1 , $m + n$ odd. Finally $b_{klmn}^j(\Omega)$ are real-valued when Ω is real because they are scalar products of real-valued functions.

The form (4.35) of the power series follows directly from (4.34) by expressing c_0 and c_1 in their polar forms: $c_0 = |c_0|e^{i\theta_0}$ and $c_1 = |c_1|e^{i\theta_1}$, $\Delta\theta = \theta_1 - \theta_0$, and using that $m + n, k + l$, and $m - n = -(k - l)$ are all even. The proof of Proposition 4.3 is now complete. \square

4.1. Ground state and excited state branches, prebifurcation. In this section we prove part (i) of Theorem 4.1 as well as a corresponding statement about the excited state. In particular, we show that for sufficiently small amplitude, the only nonlinear bound state families are those arising via bifurcation from the zero state at the eigenvalues Ω_0 and Ω_1 . This is true for general potentials with two bound states. Here, however, we can determine threshold amplitude, \mathcal{N}_{cr} , above which the solution set changes.

A closed system of equations for c_0, c_1 , and Ω , parametrized by \mathcal{N} , is obtained upon substitution of $\eta[c_0, c_1, \Omega]$ (Proposition 4.1) into (4.9)-(4.12). Furthermore, using the structural properties (4.29)-(4.30) of Proposition 4.3, we obtain

$$(4.38) \quad (\Omega_0 - \Omega)c_0 + a_{0000}|c_0|^2c_0 + (a_{0110} + a_{0011})|c_1|^2c_0 + a_{0011}c_1^2\bar{c}_0 + c_0f_0(c_0, c_1, \Omega) = 0,$$

$$(4.39) \quad (\Omega_1 - \Omega)c_1 + a_{1111}|c_1|^2c_1 + (a_{1010} + a_{1001})|c_0|^2c_1 + a_{1010}c_0^2\bar{c}_1 + c_1f_1(c_0, c_1, \Omega) = 0,$$

$$(4.40) \quad |c_0|^2 + |c_1|^2 + \mathcal{O}(|c_0|^2 + |c_1|^2)^3 = \mathcal{N}.$$

This system of equations is valid for $|c_0| + |c_1| < r_$, independent of L (the distance between wells).*

If we choose $c_1 = 0$, then the second equation in the system, (4.39), is satisfied. In this case, a nontrivial solution requires $c_0 \neq 0$. The first equation, (4.38), after

factoring out c_0 becomes

$$(4.41) \quad \Omega_0 - \Omega + a_{0000}|c_0|^2 + f_0(|c_0|, 0, \Omega) = 0,$$

where we used (4.31) to eliminate the phase of the complex quantity c_0 . Since Ω is real, (4.41) becomes one equation with two real parameters $\Omega, |c_0|$. Since the right-hand side of (4.41) vanishes for $\Omega = \Omega_0$ and $|c_0| = 0$ and since the partial derivative of this function with respect to Ω , evaluated at this solution, is nonzero, we have by the implicit function theorem that there is a unique solution

$$(4.42) \quad \Omega = \Omega_g(|c_0|) = \Omega_0 + a_{0000}|c_0|^2 + \mathcal{O}(|c_0|^4).$$

By (4.40), for small amplitudes, the mapping from $|c_0|^2 + |c_1|^2$ to \mathcal{N} is invertible. The family of solutions

$$|c_0| \mapsto (|c_0|e^{i\theta}, |c_1| = 0, \Omega = \Omega_g(|c_0|)), \quad \theta \in [0, 2\pi),$$

defined for $|c_0|$ sufficiently small, corresponds to a family of symmetric nonlinear bound states, $u_{\mathcal{N}}$ with $\|u_{\mathcal{N}}\|_{L^2}^2 = \mathcal{N}$, bifurcating from the zero solution at the linear eigenvalue Ω_0

$$u_{\mathcal{N}} = (|c_0|\psi_0(x) + \eta[|c_0|, 0, \Omega_g(|c_0|)](x)) e^{i\theta_0}, \quad \theta_0 \in [0, 2\pi);$$

see, for example, [21, 25, 26]. Since both ψ_0 and $\eta(|c_0|, 0, \Omega_g)$ are even (in x_1) we infer that $u_{\mathcal{N}}$ is symmetric (even).

Remark 4.1. A similar result holds for the case $c_0 = 0$ leading to the antisymmetric excited state branch.

PROPOSITION 4.4. *For $|c_0| + |c_1|$ sufficiently small, these two branches of solutions are the only nontrivial solutions of (2.7).*

Proof. Indeed, suppose the contrary. By local uniqueness of these branches, ensured by the implicit function theorem, a solution not already lying on one of these branches must have both c_0 and c_1 nonzero. Now, divide the first equation by c_0 , the second equation by c_1 , and subtract the results. Introducing polar coordinates

$$(4.43) \quad c_0 = \rho_0 e^{i\theta_0}, \quad c_1 = \rho_1 e^{i\theta_1}, \quad \Delta\theta = \theta_1 - \theta_0,$$

we obtain from (4.38)

$$(4.44) \quad \begin{aligned} \Omega_1 - \Omega_0 &= a_{0000}\rho_0^2 + (a_{0110} + a_{0011} + a_{0011}e^{i2\Delta\theta})\rho_1^2 + f_0(\rho_0, \rho_1, \Delta\theta, \Omega) \\ &- a_{1111}\rho_1^2 - (a_{1001} + a_{1010} + a_{1010}e^{-i2\Delta\theta})\rho_0^2 - f_1(\rho_0, \rho_1, \Delta\theta, \Omega). \end{aligned}$$

The left-hand side is nonzero while the right-hand side is continuous uniformly for Ω satisfying (4.15) and zero for $\rho_0 = 0 = \rho_1$. Equation (4.44) cannot hold for $|\rho_0| + |\rho_1| < \varepsilon$ where $\varepsilon > 0$ is independent of Ω . This completes the proof of Proposition 4.4. \square

Note, however, that nothing can preclude validity of (4.44) for larger ρ_0 and ρ_1 , possibly leading to a third branch of solutions of (2.7). In what follows, we show that this is indeed the case and the third branch bifurcates from the ground state one at a critical value of $\rho_0 = \rho_0^*$.

4.2. Symmetry-breaking bifurcation along the ground state/symmetric branch. A consequence of the previous section is that there are no bifurcations from the ground state branch for sufficiently small amplitude. We now seek a bifurcating

branch of solutions to (2.16), (4.40), along which $c_0 \cdot c_1 \neq 0$. As argued just above, along such a new branch one must have

$$(4.45) \quad \Omega_0 - \Omega + a_{0000}\rho_0^2 + (a_{0110} + a_{0011} + a_{0011}e^{i2\Delta\theta})\rho_1^2 + f_0(\rho_0, \rho_1, \Delta\theta, \Omega) = 0,$$

$$(4.46) \quad \Omega_1 - \Omega + a_{1111}\rho_1^2 + (a_{1010} + a_{1001} + a_{1010}e^{-i2\Delta\theta})\rho_0^2 + f_1(\rho_0, \rho_1, \Delta\theta, \Omega) = 0.$$

We first derive constraints on $\Delta\theta$. Consider the imaginary parts of the two equations and use the expansions (4.35) and the fact that Ω is real:

$$\begin{aligned} & a_{0011} \sin(2\Delta\theta)\rho_1^2 + \sum_{k,m \geq 1, p \in \mathbb{Z}} b_{kmp}^0(\Omega) \sin(p2\Delta\theta)\rho_0^{2k}\rho_1^{2m} \\ &= \sin(2\Delta\theta)\rho_1^2(a_{0011} + \mathcal{O}(\rho_0^2 + \rho_1^2)) = 0, \\ & -a_{1010} \sin(2\Delta\theta)\rho_0^2 + \sum_{k,m \geq 1, p \in \mathbb{Z}} b_{kmp}^1(\Omega) \sin(p2\Delta\theta)\rho_0^{2k}\rho_1^{2m} \\ &= \sin(2\Delta\theta)\rho_0^2(-a_{1010} + \mathcal{O}(\rho_0^2 + \rho_1^2)) = 0. \end{aligned}$$

Due to the nondegeneracy assumption (4.3) and $a_{0011} = a_{1010}$, see (2.19), both equations hold if and only if $\sin(2\Delta\theta) = 0$ or, equivalently,

$$(4.47) \quad \Delta\theta \in \left\{ 0, \frac{\pi}{2}, \pi, \frac{3\pi}{2} \right\}.$$

Case 1. $\Delta\theta \in \{0, \pi\}$: Here, the system (4.45)-(4.46) is equivalent to the same system of two real equations with three real parameters $\rho_0 \geq 0$, $\rho_1 \geq 0$, and Ω :

$$(4.48) \quad F_0(\rho_0, \rho_1, \Omega) \stackrel{def}{=} \Omega_0 - \Omega + a_{0000}\rho_0^2 + (a_{0110} + 2a_{0011})\rho_1^2 + f_0(\rho_0, \rho_1, \Omega) = 0,$$

$$(4.49) \quad F_1(\rho_0, \rho_1, \Omega) \stackrel{def}{=} \Omega_1 - \Omega + a_{1111}\rho_1^2 + (2a_{1010} + a_{1001})\rho_0^2 + f_1(\rho_0, \rho_1, \Omega) = 0.$$

We begin by seeking the point along the ground state branch $(\rho_0^*, 0, \Omega_g(\rho_0^*))$ from which a new family of solutions of (4.48)-(4.49), parametrized by $\rho_1 \geq 0$, bifurcates; see (4.42).

Recall first that for *any* $\rho_0 \geq 0$ sufficiently small, $F_0(\rho_0, 0, \Omega_g(\rho_0)) = 0$. A candidate for a bifurcation point is $\rho_0^* > 0$ for which, in addition,

$$(4.50) \quad F_1(\rho_0^*, 0, \Omega_g(\rho_0^*)) = 0.$$

Using (4.33) one can check that

$$(4.51) \quad F_1(\rho_0, 0, \Omega_g(\rho_0)) = \Omega_1 - \Omega_0 + (a_{1001} + 2a_{1010} - a_{0000} + \mathcal{O}(\rho_0^2))\rho_0^2 = 0$$

has a solution:

$$(4.52) \quad \rho_0^* = \sqrt{\frac{\Omega_1 - \Omega_0}{|a_{1001} + 2a_{1010} - a_{0000}|}} \left[1 + \mathcal{O}\left(\frac{\Omega_1 - \Omega_0}{|a_{1001} + 2a_{1010} - a_{0000}|^2}\right) \right].$$

We now show that a new family of solutions bifurcates from the symmetric state at $(\rho_0^*, 0, \Omega_g(\rho_0^*))$. This is realized as a unique, one-parameter family of solutions

$$(4.53) \quad \rho_1 \mapsto (\rho_0(\rho_1), \rho_1, \Omega_{asym}(\rho_1))$$

of the equations

$$(4.54) \quad F_0(\rho_0, \rho_1, \Omega) = 0, \quad F_1(\rho_0, \rho_1, \Omega) = 0.$$

To see this, note that by the preceding discussion we have $F_j(\rho_0^*, 0, \Omega_g(\rho_0^*)) = 0, j = 1, 2$. Moreover, the Jacobian

$$\left| \frac{\partial(F_0, F_1)}{\partial(\rho_0, \Omega)}(\rho_0^*, 0, \Omega_g(\rho_0^*)) \right| = 2\rho_0^*(a_{1001} + 2a_{1010} - a_{0000} + \mathcal{O}(\rho_0^{*2}))$$

is nonzero because $\rho_0^* > 0$ and

$$(4.55) \quad a_{1001} + 2a_{1010} - a_{0000} + \mathcal{O}(\rho_0^{*2}) < 0$$

since ρ_0^* solves (4.51) and $\Omega_1 - \Omega_0 > 0$. Therefore, by the implicit function theorem, for small $\rho_1 > 0$, there is a unique solution of the system (4.48)-(4.49):

$$(4.56) \quad \rho_0 = \rho_0(\rho_1) = \rho_0^* + \frac{\rho_1^2}{2\rho_0^*} \left(\frac{a_{0110} + 2a_{0011} - a_{1111}}{a_{1001} + 2a_{1010} - a_{0000}} + \mathcal{O}(\rho_0^{*2}) \right) + \mathcal{O}(\rho_1^4),$$

$$(4.57) \quad \begin{aligned} \Omega &= \Omega_{asym}(\rho_1) = \Omega_g(\rho_0^*) \\ &+ \rho_1^2 \left(a_{1111} + (2a_{1010} + a_{1001}) \frac{a_{0110} + 2a_{0011} - a_{1111}}{a_{1001} + 2a_{1010} - a_{0000}} + \mathcal{O}(\rho_0^{*2}) \right) + \mathcal{O}(\rho_1^4). \end{aligned}$$

Remark 4.2. (1) Due to the equivalence of \mathcal{N} and $\rho_0^2 + \rho_1^2$ as parameters, for small amplitude, we have that symmetry is broken at

$$(4.58) \quad \mathcal{N}_{cr} \sim \frac{\Omega_1 - \Omega_0}{|a_{0000} - a_{1001} - 2a_{1010}|}.$$

(2) Note also that we have the family of solutions

$$(4.59) \quad e^{i\theta} (\rho_0(\rho_1)\psi_0 \pm \rho_1\psi_1 + \eta(\rho_0(\rho_1), \pm\rho_1, \Omega_{asym}(\rho_1))), \quad 0 \leq \theta < 2\pi, \quad \rho_1 > 0.$$

Here the \pm is present because the phase difference $\Delta\theta$ between c_0 and c_1 can be 0 or π ; see (4.47) and immediately below it. Because $\rho_0 \neq 0 \neq \rho_1$ this branch is neither symmetric nor antisymmetric. Thus, symmetry breaking has taken place. In the case of the double well, the \pm sign in (4.59) shows that the bound states on this asymmetric branch tend to localize in one of the two wells but not symmetrically in both; see also [3], [23], [14].

Case 2. $\Delta\theta \in \{\frac{\pi}{2}, \frac{3\pi}{2}\}$: In both cases the system (4.45)-(4.46) is equivalent to the system of two real equations, depending on three real parameters $\rho_0 \geq 0, \rho_1 \geq 0$, and Ω :

$$(4.60) \quad F_0(\rho_0, \rho_1, \Omega) \stackrel{def}{=} \Omega_0 - \Omega + a_{0000}\rho_0^2 + a_{0110}\rho_1^2 + f_0(\rho_0, \rho_1, \Omega) = 0,$$

$$(4.61) \quad F_1(\rho_0, \rho_1, \Omega) \stackrel{def}{=} \Omega_1 - \Omega + a_{1111}\rho_1^2 + a_{1001}\rho_0^2 + f_1(\rho_0, \rho_1, \Omega) = 0.$$

As before, in order to have a second bifurcation of the symmetric branch it is necessary to find a point, $(\rho_0^{**}, 0, \Omega_g(\rho_0^{**}))$, for which

$$(4.62) \quad F_1(\rho_0^{**}, 0, \Omega_g(\rho_0^{**})) = \Omega_1 - \Omega_0 + (a_{1001} - a_{0000})\rho_0^{**2} + \mathcal{O}(\rho_0^{**4}) = 0.$$

To preclude the existence of this bifurcation, one must have that there are no real solutions for ρ_0^{**} small. This is ensured by inequality (3.6): $a_{1001} - a_{0000} > 0$.

In fact (3.6) can be *deduced* from the following general argument and thus we need not hypothesize it. Let L_- and L_+ denote the second-order, self-adjoint Schrödinger operators, defined in section 5, and related to the real and imaginary parts of the linearized NLS-GP equation about a real-valued nonlinear bound state. Consider the full *complex* linearization of the nonlinear bound state equation for NLS-GP. A nontrivial solution of the type in question would correspond to a second element of the null space of L_- . Since zero is the lowest eigenvalue of L_- and the lowest eigenvalue of an elliptic second-order operator is nondegenerate, we have a contradiction. It follows that there is no bifurcation for small amplitude and therefore (3.6) holds.

Remark 4.3. A similar argument holds along the antisymmetric branch. We now denote L_{\pm} as linearized operators associated with the *antisymmetric* (odd parity with respect to x_1) branch. A bifurcation would occur along it if and only if $\dim \ker L_- \geq 2$ (the antisymmetric nonlinear bound state is always in $\ker L_-$ along this branch) or $\dim \ker L_+ \geq 1$. In the regime in which the first two eigenvalues are still separated from the rest of the spectrum, the former means that the lowest eigenvalue of L_- , i.e., the one bifurcating from $\Omega_0 - \Omega_1$, crosses zero and becomes double. This contradicts the nondegeneracy of the ground state of the second-order elliptic operator L_- . Now, $\dim \ker L_+ \geq 1$ is also impossible. This follows from the fact that the nonlinearity is attractive, see (H3), which easily implies that, for the nonlinear antisymmetric bound state Ψ_{Ω} , we have $(\Psi_{\Omega}, L_+ \Psi_{\Omega}) < 0$. Since L_+ is self-adjoint and its lowest eigenvalue has a symmetric eigenvector orthogonal on Ψ_{Ω} , we get that the first two eigenvalues must be strictly negative and none can become zero. A consequence of a lack of bifurcation along the antisymmetric branch are inequalities (3.7), (3.8).

In summary, there are no other bifurcations than those stated in Theorem 4.1 for NLS-GP, where our Lyapunov–Schmidt-type reduction applies, i.e., in the regime covered by Proposition 4.1.

5. Exchange of stability at the bifurcation point. In this section we consider the dynamic stability of the symmetric and asymmetric waves, associated with the branch bifurcating from the zero state at the *ground state frequency*, Ω_0 , of the linear Schrödinger operator $-\Delta + V(x)$; see Figure 1.

The notion of stability with which we work is H^1 -orbital Lyapunov stability.

DEFINITION 5.1. *The family of nonlinear bound states $\{\Psi_{\Omega} e^{-i\Omega t} : \theta \in [0, 2\pi)\}$ is H^1 -orbitally Lyapunov stable if for every $\varepsilon > 0$ there is a $\delta(\varepsilon) > 0$ such that if the initial data u_0 satisfies*

$$\inf_{\theta \in [0, 2\pi)} \|u_0(\cdot) - \Psi_{\Omega}(\cdot)e^{i\theta}\|_{H^1} < \delta,$$

then for all $t \neq 0$, the solution $u(x, t)$ satisfies

$$\inf_{\theta \in [0, 2\pi)} \|u(\cdot, t) - \Psi_{\Omega}(\cdot)e^{i\theta}\|_{H^1} < \varepsilon.$$

In this section we prove the following theorem.

THEOREM 5.1. *Consider the bifurcations elucidated in Theorem 4.1. The symmetric branch is H^1 -orbitally Lyapunov stable for $0 \leq \rho_0 < \rho_0^*$ or, equivalently, $0 < \mathcal{N} < \mathcal{N}_{cr}$. At the bifurcation point $\rho_0 = \rho_0^*$ ($\mathcal{N} = \mathcal{N}_{cr}$), there is an exchange of stability from the symmetric branch to the asymmetric branch. In particular, for $\mathcal{N} > \mathcal{N}_{cr}$ the asymmetric state is stable and the symmetric state is unstable.*

We summarize basic results on stability and instability. Introduce L_+ and L_- , real and imaginary parts, respectively, of the linearized operators about a real-valued state Ψ_Ω :

$$\begin{aligned}
 L_+ &= L_+[\Psi_\Omega] \cdot = (H - \Omega) \cdot + \partial_u N(u, u, u) |_{\Psi_\Omega} \\
 &\equiv (H - \Omega) \cdot + D_u N[\Psi_\Omega](\cdot), \\
 (5.1) \quad L_- &= L_-[\Psi_\Omega] \cdot = (H - \Omega) \cdot + N(\Psi_\Omega, \Psi_\Omega, \cdot).
 \end{aligned}$$

By (2.7) and (2.4), $L_- \Psi_\Omega = 0$.

We state a special case of known results on stability and instability, directly applicable to the symmetric branch which bifurcates from the zero state at the ground state frequency of $-\Delta + V$.

THEOREM 5.2 (see [32, 33, 10, 9, 15]).

- (1) Stability: *Suppose L_+ has exactly one negative eigenvalue and L_- is nonnegative. Assume that*

$$(5.2) \quad \frac{d}{d\Omega} \int |\Psi_\Omega(x)|^2 dx < 0.$$

Then, Ψ_Ω is H^1 -orbitally stable.

- (2) Instability: *Suppose L_- is nonnegative. If $n_-(L_+) \geq 2$, then the linearized dynamics about Ψ_Ω has spatially localized solution which is exponentially growing in time. Moreover, Ψ_Ω is not H^1 -orbitally stable.*

First we claim that along the branch of symmetric solutions bifurcating from the zero solution at frequency Ω_0 , the hypothesis on L_- holds. To see that the operator $L_-[\Psi_\Omega]$ is always nonnegative, consider $L_-[\Psi_{\Omega_0}] = L_-[0] = -\Delta + V - \Omega_0$. Clearly, $L_-[0]$ is a nonnegative operator because Ω_0 is the lowest eigenvalue of $-\Delta + V$. Since clearly we have $L_- \Psi_\Omega = 0$, $0 \in \text{spec}(L_-[\Psi_\Omega])$. Since the lowest eigenvalue is necessarily simple, by continuity there cannot be any negative eigenvalues for Ω sufficiently close to Ω_0 . Finally, if for some Ω , L_- has a negative eigenvalue, then by continuity there would be an Ω_* for which $L_-[\Psi_{\Omega_*}]$ would have a double eigenvalue at zero and no negative spectrum. But this contradicts that the ground state is simple. Therefore, it is the quantity $n_-(L_+)$ which controls whether or not Ψ_Ω is stable.

Next we discuss the slope condition (5.2). It is clear from the construction of the branch $\Omega \mapsto \Psi_\Omega$ that (5.2) holds for Ω near Ω_0 . Suppose now that $\partial_\Omega \int |\Psi_\Omega|^2 = 0$. Then, $\langle \Psi_\Omega, \partial_\Omega \Psi_\Omega \rangle = 0$. As shown below, L_+ has exactly one negative eigenvalue for Ω sufficiently near Ω_0 . It follows that $L_+ \geq 0$ on $\{\Psi_\Omega\}^\perp$ (see [32, 33]). Therefore, we have $(L_+^{\frac{1}{2}} \partial_\Omega \Psi_\Omega, L_+^{\frac{1}{2}} \partial_\Omega \Psi_\Omega) = (L_+ \partial_\Omega \Psi_\Omega, \partial_\Omega \Psi_\Omega) = (\Psi_\Omega, \partial_\Omega \Psi_\Omega) = 0$. Therefore, $L_+^{\frac{1}{2}} \partial_\Omega \Psi_\Omega = 0$, implying $\Psi_\Omega = L_+ \partial_\Omega \Psi_\Omega = 0$, which is a contradiction. It follows that (5.2) holds so long as $L_+ > 0$ on $\{\Psi_\Omega\}^\perp$ and when (5.2) first fails, it does so due to a nontrivial element of the nullspace of L_+ .

Therefore, Ψ_Ω is stable so long as $n_-(L_+)$ does not increase. We shall now show that for $\mathcal{N} < \mathcal{N}_{cr}$, $n_-(L_+[\Psi_\Omega]) = 1$ but that along the symmetric branch for $\mathcal{N} > \mathcal{N}_{cr}$, $n_-(L_+[\Psi_\Omega]) = 2$. Furthermore, we show that along the bifurcating asymmetric branch, the hypotheses of Theorem 5.2 ensuring stability hold.

Remark 5.1. For simplicity we have considered the most important case where there is a transition from dynamical stability to dynamical instability along the symmetric branch, bifurcating from the ground state of H . However, our analysis which actually shows that along any symmetric branch, associated with any of the eigenvalues $\Omega_{2j}, j \geq 0$, of H , there is a critical $\mathcal{N} = \mathcal{N}_{cr}(j)$ such that as \mathcal{N} is increased

through $\mathcal{N}_{cr}(j)$, the number of negative eigenvalues of the linearization about the symmetric state along the j th symmetric branch, $n_-(L_+^{(j)})$, increases by one. By the results in [15, 9, 19], this has implications for the number of unstable modes of higher order ($j \geq 1$) symmetric states.

Consider the spectral problem for $L_+ = L_+[\Psi_\Omega]$:

$$(5.3) \quad L_+[\Psi_\Omega]\phi = \mu\phi.$$

We now formulate a Lyapunov–Schmidt reduction of (5.3) and then relate it to our formulation for nonlinear bound states. We first decompose ϕ relative to the states ψ_0, ψ_1 and their orthogonal complement:

$$\phi = \alpha_0\psi_0 + \alpha_1\psi_1 + \xi, \quad (\psi_j, \xi) = 0, \quad j = 0, 1.$$

Projecting (5.3) onto ψ_0, ψ_1 and onto the range of \tilde{P} we obtain the system

$$(5.4) \quad \langle \psi_0, L_+[\Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1 + \xi) \rangle = \mu\alpha_0,$$

$$(5.5) \quad \langle \psi_1, L_+[\Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1 + \xi) \rangle = \mu\alpha_1,$$

$$(5.6) \quad (H - \Omega)\xi + D_u N[\Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1 + \xi) = \mu\xi.$$

The last equation can be rewritten in the form

$$(5.7) \quad \left[I + (H - \Omega - \mu)^{-1} \tilde{P} D_u N[\Psi_\Omega] \right] \xi = -(H - \Omega - \mu)^{-1} \tilde{P} D_u N[\Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1).$$

The operator on the right-hand side of (5.7) is essentially the Jacobian studied in the proof of Proposition 4.1, evaluated at $\Omega + \mu$. Hence, by the proof of Proposition 4.1, if $\Omega + \mu$ satisfies (4.15) and $\|\Psi_\Omega\|_{H^2} \leq \mathcal{N}_*$, then the operator $I + (H - \Omega - \mu)^{-1} \tilde{P} D_u N[\Psi_\Omega]$ is invertible on H^2 and (5.7) has a unique solution

$$(5.8) \quad \begin{aligned} \xi &\stackrel{def}{=} \xi[\mu, \alpha_0, \alpha_1, \Omega] \\ &\equiv Q[\mu, \Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1) \\ &= -(I + (H - \Omega - \mu)^{-1} \tilde{P} D_u N[\Psi_\Omega])^{-1} (H - \Omega - \mu)^{-1} \tilde{P} D_u N[\Psi_\Omega](\alpha_0\psi_0 + \alpha_1\psi_1) \\ &= \mathcal{O}[(|\rho_0| + |\rho_1|)^2] [\alpha_0\psi_0 + \alpha_1\psi_1]. \end{aligned}$$

The last relation follows from $D_u N[\psi]$ being a quadratic form in $\Psi_\Omega = \rho_0\psi_0 + \rho_1\psi_1 + \mathcal{O}((|\rho_0| + |\rho_1|)^3)$.

Substitution of the expression for ξ as a functional of α_j into (5.4) and (5.5) we get a closed system of two real equations:

$$(5.9) \quad \begin{aligned} (\Omega_0 - \Omega)\alpha_0 + \langle \psi_0, D_u N[\Psi_\Omega] (I + Q[\mu, \Psi_\Omega]) (\alpha_0\psi_0 + \alpha_1\psi_1) \rangle &= \mu \alpha_0, \\ (\Omega_1 - \Omega)\alpha_1 + \langle \psi_1, D_u N[\Psi_\Omega] (I + Q[\mu, \Psi_\Omega]) (\alpha_0\psi_0 + \alpha_1\psi_1) \rangle &= \mu \alpha_1. \end{aligned}$$

The system (5.9) is the Lyapunov–Schmidt reduction of the linear eigenvalue problem for L_+ with eigenvalue parameter μ . *Our next step will be to write it in a form relating it to the linearization of the Lyapunov–Schmidt reduction of the nonlinear problem.*

Remark 5.2. For $\|\Psi_\Omega\|_{H^2} \leq n_*$, the above system is equivalent to the eigenvalue problem for the operator $L_+[\Psi_\Omega]$ with eigenvalue parameter μ as long as (4.15) holds

with Ω replaced by $\Omega + \mu$. This restriction on the spectral parameter, μ , is in fact very mild and has no impact on the analysis. This is because we are primarily interested in μ near zero as we are interested in detecting the crossing of an eigenvalue of L_+ from positive to negative reals as \mathcal{N} is increased beyond some \mathcal{N}_{cr} . The values of μ for which (4.15) does not hold do not play a role in any change of index, $n_-(L_+)$.

First rewrite (5.9) as

$$(5.10) \quad (\Omega_0 - \Omega - \mu)\alpha_0 + \langle \psi_0, D_u N[\Psi_\Omega] (I + Q[0, \Psi_\Omega]) (\alpha_0 \psi_0 + \alpha_1 \psi_1) \rangle + \langle \psi_0, D_u N[\Psi_\Omega] \Delta Q[\mu, \Psi_\Omega] (\alpha_0 \psi_0 + \alpha_1 \psi_1) \rangle = 0,$$

$$(5.11) \quad (\Omega_1 - \Omega - \mu)\alpha_1 + \langle \psi_1, D_u N[\Psi_\Omega] (I + Q[0, \Psi_\Omega]) (\alpha_0 \psi_0 + \alpha_1 \psi_1) \rangle + \langle \psi_1, D_u N[\Psi_\Omega] \Delta Q[\mu, \Psi_\Omega] (\alpha_0 \psi_0 + \alpha_1 \psi_1) \rangle = 0.$$

Here,

$$(5.12) \quad \Delta Q[\mu, \Psi_\Omega] = Q[\mu, \Psi_\Omega] - Q[0, \Psi_\Omega].$$

Note that terms involving ΔQ in (5.10), (5.11) are of size $\mathcal{O}[(\rho_0^2 + \rho_1^2)\mu\alpha_j]$.

PROPOSITION 5.1.

$$(5.13) \quad Q[0, \Psi_\Omega](\alpha_0 \psi_0 + \alpha_1 \psi_1) = \partial_{\rho_0} \eta[\rho_0, \rho_1, \Omega] \alpha_0 + \partial_{\rho_1} \eta[\rho_0, \rho_1, \Omega] \alpha_1.$$

Proof. Recall that η satisfies

$$(5.14) \quad F(\rho_0, \rho_1, \Omega, \eta) \equiv \eta + (H - \Omega)^{-1} \tilde{P}N[\rho_0 \psi_0 + \rho_1 \psi_1 + \eta] = 0.$$

Differentiation with respect to ρ_j , $j = 0, 1$, yields

$$(5.15) \quad \left(I + (H - \Omega)^{-1} \tilde{P}D_u N[\Psi_\Omega] \right) \partial_{\rho_j} \eta = - (H - \Omega)^{-1} \tilde{P}D_u N[\Psi_\Omega] \psi_j,$$

where

$$\Psi_\Omega = \rho_0 \psi_0 + \rho_1 \psi_1 + \eta[\rho_0, \rho_1, \Omega].$$

Thus,

$$(5.16) \quad \partial_{\rho_j} \eta = Q[0, \Psi_\Omega] \psi_j,$$

from which Proposition 5.1 follows. \square

We now use Proposition 5.1 to rewrite the first inner products in (5.10), (5.11). For $k = 0, 1$,

$$(5.17) \quad \begin{aligned} & \langle \psi_k, D_u N[\Psi_\Omega] (I + Q[0, \Psi_\Omega]) (\alpha_0 \psi_0 + \alpha_1 \psi_1) \rangle \\ &= \sum_{j=0}^1 \langle \psi_k, D_u N[\rho_0 \psi_0 + \rho_1 \psi_1 + \eta] (\psi_j + \partial_{\rho_j} \eta) \rangle \alpha_j \\ &= \sum_{j=0}^1 \frac{\partial}{\partial \rho_j} \langle \psi_k, N[\Psi_\Omega] \rangle \alpha_j \\ &= \sum_{j=0}^1 \partial_{\rho_j} \langle \psi_k, N[\rho_0 \psi_0 + \rho_1 \psi_1] \rangle \alpha_j + \partial_{\rho_j} [\rho_k f_k(\rho_0, \rho_1, \Omega)], \end{aligned}$$

where $N[\psi_\Omega] = N[\rho_0\psi_0 + \rho_1\psi_1] + \mathcal{R}$; see (2.16)–(2.18) and (4.29), (4.30). Therefore, the Lyapunov–Schmidt reduction of the eigenvalue problem for L_+ becomes

$$(5.18) \quad (\Omega_0 - \Omega - \mu)\alpha_0 + \sum_{j=0}^1 \partial_{\rho_j} \langle \psi_0, N[\rho_0\psi_0 + \rho_1\psi_1] \rangle \alpha_j + \partial_{\rho_j} [\rho_0 f_0(\rho_0, \rho_1, \Omega)] + \langle \psi_0, D_u N[\Psi_\Omega] \Delta Q[\mu, \Psi_\Omega] (\alpha_0\psi_0 + \alpha_1\psi_1) \rangle = 0,$$

$$(5.19) \quad (\Omega_1 - \Omega - \mu)\alpha_1 + \sum_{j=0}^1 \partial_{\rho_j} \langle \psi_1, N[\rho_0\psi_0 + \rho_1\psi_1] \rangle \alpha_j + \partial_{\rho_j} [\rho_1 f_1(\rho_0, \rho_1, \Omega)] + \langle \psi_1, D_u N[\Psi_\Omega] \Delta Q[\mu, \Psi_\Omega] (\alpha_0\psi_0 + \alpha_1\psi_1) \rangle = 0.$$

This can be written succinctly in a matrix form as

$$(5.20) \quad [M - \mu + \mathcal{C}(\mu)] \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

where

$$(5.21) \quad M = M[\Omega, \rho_0, \rho_1] = \begin{pmatrix} \Omega_0 - \Omega + 3a_{0000}\rho_0^2 + (a_{0110} + 2a_{0011})\rho_1^2 + \partial_{\rho_0}(\rho_0 f_0) & 2(a_{0110} + 2a_{0011})\rho_0\rho_1 + \partial_{\rho_1}(\rho_0 f_0) \\ 2(2a_{1010} + a_{1001})\rho_0\rho_1 + \partial_{\rho_0}(\rho_1 f_1) & (\Omega_1 - \Omega) + 3a_{1111}\rho_1^2 + (2a_{1010} + a_{1001})\rho_0^2 + \partial_{\rho_1}(\rho_1 f_1) \end{pmatrix}$$

and

$$(5.22) \quad \mathcal{C}(\mu)_{lm} = \langle \psi_l, D_u N[\Psi_\Omega] \Delta Q[\mu, \Psi_\Omega] \psi_m \rangle, \quad l, m = 0, 1.$$

Note that

$$(5.23) \quad \mathcal{C}(\mu = 0) = 0.$$

Recall that μ is the spectral parameter for the eigenvalue problem L_+ , (5.3), and we are interested in $n_-(L_+[\Psi_\Omega])$, the number of negative eigenvalues along a family of nonlinear bound states $\Omega \mapsto \Psi_\Omega$. By Theorem 5.2, $n_-(L_+)$ determines the stability or instability of a particular state. This question has now been mapped to the problem of following the roots of

$$(5.24) \quad D(\mu, \rho_0, \rho_1) = \det(\mu I - M - \mathcal{C}(\mu)) = 0,$$

where ρ_0 and ρ_1 are parameters along the different branches of nonlinear bound states. Since $\mathcal{C}(\mu)$, defined in (5.22), is small for small amplitude nonlinear bound states, we expect the roots, μ , to be small perturbations of the eigenvalues of the matrix M . We study these roots along the symmetric ($M = M(\Omega_g(\rho_0), \rho_0, 0)$) and asymmetric branch ($M = M(\Omega_{asym}(\rho_1), \rho_0(\rho_1), \rho_1)$) using the implicit function theorem.

Symmetric branch. Along the symmetric branch:

$$\rho_1 = 0, \quad \rho_0 \geq 0, \quad \Omega = \Omega_g = \Omega_0 + a_{0000}\rho_0^2 + \mathcal{O}(\rho_0^4), \\ \Psi_\Omega = \rho_0\psi_0 + \eta(\rho_0, 0, \Omega) = \text{symmetric}.$$

Thus, $D = D(\mu, \rho_0)$. Moreover, the system (5.20) is diagonal. This is because Q and hence ΔQ preserve parity at a symmetric Ψ_Ω ; see their definitions (5.8) and (5.12). Therefore $\mathcal{C}_{01} = 0 = \mathcal{C}_{10}$, since each is the scalar product of an even and an odd function. Moreover, from (4.35) we get $\frac{\partial f_j}{\partial \rho_1}(\rho_0, 0, \Omega) = 0$, $j = 0, 1$.

Therefore, the matrix $\mu I - M - \mathcal{C}(\mu)$ is diagonal and μ is an eigenvalue of $L_+[\psi_{\Omega_g(\rho_0)}]$ if and only if μ is a root of either

$$(5.25) \quad P_0(\mu, \rho_0) \equiv \mu - M_{00}(\rho_0) - \mathcal{C}_{00}(\mu, \rho_0) = 0$$

or

$$(5.26) \quad P_1(\mu, \rho_0) \equiv \mu - M_{11}(\rho_0) - \mathcal{C}_{11}(\mu, \rho_0) = 0.$$

Both P_0 and P_1 are analytic in μ and ρ_0 and it is easy to check that

$$P_0(0, 0) = 0, \quad \partial_\mu P_0(0, 0) = 1$$

and

$$P_1(\Omega_1 - \Omega_0, 0) = 0, \quad \partial_\mu P_1(\Omega_1 - \Omega_0, 0) = 1.$$

Formally differentiating (5.25) or (5.26) with respect to ρ_0 gives

$$(5.27) \quad \partial_{\rho_0} \mu_j = \frac{\partial_{\rho_0} M_{jj} + \partial_{\rho_0} \mathcal{C}_{jj}}{1 - \partial_\mu \mathcal{C}_{jj}}.$$

By the implicit function theorem (5.25) and (5.26) define, respectively, μ_0 and μ_1 as smooth functions of ρ provided

$$(5.28) \quad |\partial_\mu \mathcal{C}_{jj}| < 1, \quad j = 0, 1.$$

A direct calculation using (5.8) and estimates (2.5), (4.20) shows that in the regime of interest: Ω satisfying (4.15), it suffices to have

$$(5.29) \quad \|\Psi_\Omega\|_{H^2} \leq n_* (9 \max(\|\psi_0\|_{H^2}, \|\psi_1\|_{H^2}))^{-\frac{1}{4}},$$

where n_* is given by Proposition 4.1. The latter can be reduced to an estimate on ρ_0 via the above definition of Ψ_Ω and (4.24) as in the end of the proof of Proposition 4.1.

Therefore, under conditions (4.15) and (5.29), we have a unique solution μ_0 , respectively μ_1 , of (5.25), respectively (5.26). Moreover, the two solutions are analytic in ρ_0 and, for small ρ_0 , we have the following estimates:

$$(5.30) \quad \mu_0 = 2a_{0000}\rho_0^2 + \mathcal{O}(\rho_0^4) < 0,$$

$$(5.31) \quad \mu_1 = \Omega_1 - \Omega_0 + \mathcal{O}(\rho_0^2) > 0,$$

where we used $a_{0000} \equiv g\langle \psi_0^2, K[\psi_0^2] \rangle < 0$ and $\mu_1(\rho_0 = 0) = \Omega_1 - \Omega_0 > 0$.

We claim that μ_1 changes sign for the first time at $\rho_0 = \rho_0^*$. Indeed, by continuity, the sign can only change when $\mu_1 = 0$, i.e., when (5.26) has a solution of the form $(0, \rho_0)$. Since $\mathcal{C}_{11}(0, \rho_0) = 0$ (see (5.23)) (5.26) becomes

$$0 = M_{11}(\rho_0) = \Omega_1 - \Omega_g(\rho_0) + (2a_{1010} + a_{1001})\rho_0^2 + f_1(\rho_0, 0, \Omega_g) = F_1(\rho_0, 0, \Omega_g(\rho_0));$$

see (5.21) and note that $\rho_1 = 0$. But this equation is the same as (4.50), which determines the bifurcation point ρ_0^* . Thus, as expected, $D(\mu, \rho_0) = 0$ has a root

$\rho_1(\rho_0^*) = 0$ or, equivalently, L_+ has a zero eigenvalue at the bifurcation point. Note that the associated null eigenfunction has odd parity in one space dimension and is, more generally, nonsymmetric and changes sign.

To see that $\mu_1(\rho_0)$ changes sign at $\rho_0 = \rho_0^*$ we differentiate (5.26) with respect to ρ_0 at $\rho_0 = \rho_0^*$ and obtain from (5.27) that

$$\partial_{\rho_0}\mu_1 = \frac{\partial_{\rho_0}M_{11} + \partial_{\rho_0}C_{11}}{1 - \partial_{\mu}C_{11}} < 0.$$

This follows because the denominator is positive, by (5.28), while direct calculation gives for the numerator

$$\partial_{\rho_0}M_{11}(\rho_0^*) + \partial_{\rho_0}C_{11}(\rho_0^*) = 2\rho_0^* (a_{1001} + 2a_{1010} - a_{0000} + \mathcal{O}(\rho_0^{*2})) < 0;$$

see (4.55). Therefore μ_1 becomes negative for $\rho_0 > \rho_0^*$ at least when $|\rho_0 - \rho_0^*|$ is small enough.

In conclusion, $L_+[\Omega_g(\rho_0)]$ has exactly one negative eigenvalue for $0 \leq \rho_0 < \rho_0^*$ and two negative eigenvalues for $\rho_0 > \rho_0^*$ and $|\rho_0 - \rho_0^*|$ small. Therefore, following the criteria of [32, 33, 10, 15, 9, 16], the symmetric branch is stable for $0 \leq \rho_0 < \rho_0^*$ and becomes unstable past the bifurcation point.

Asymmetric branch: Stability for $\mathcal{N} > \mathcal{N}_{cr}$. Finally, we study the behavior of the eigenvalue problem (5.20) on the asymmetric branch:

$$(5.32) \quad 0 \leq \rho_1 \ll 1,$$

$$\rho_0 = \rho_0(\rho_1) = \rho_0^* + \frac{\rho_1^2}{2\rho_0^*} \left(\frac{a_{0110} + 2a_{0011} - a_{1111}}{a_{1001} + 2a_{1010} - a_{0000}} + \mathcal{O}(\rho_0^{*2}) \right) + \mathcal{O}(\rho_1^4),$$

$$\begin{aligned} \Omega = \Omega_{asym}(\rho_1) &= \Omega_g(\rho_0^*) \\ &+ \rho_1^2 \left(a_{1111} + (2a_{1010} + a_{1001}) \frac{a_{0110} + 2a_{0011} - a_{1111}}{a_{1001} + 2a_{1010} - a_{0000}} + \mathcal{O}(\rho_0^{*2}) \right) \\ &+ \mathcal{O}(\rho_1^4), \end{aligned} \tag{5.33}$$

$$\Psi_{\Omega} = \rho_0(\rho_1)\psi_0 + \rho_1\psi_1 + \eta(\rho_0(\rho_1), \rho_1, \Omega_{asym}(\rho_1)).$$

The eigenvalues will be given by the zeros of the real-valued function

$$(5.34) \quad D(\mu, \rho_1) = \det(\mu I - M(\rho_1) - \mathcal{C}(\mu, \rho_1)),$$

which is analytic in μ and ρ_1 for $\Omega + \mu$ satisfying (4.15) and ρ_1 small. Note that at $\rho_1 = 0$ we are still on the symmetric branch at the bifurcation point $\rho_0 = \rho_0^*$. Hence, the matrix is diagonal and

$$(5.35) \quad D(\mu, 0) = P_0(\mu, \rho_0^*)P_1(\mu, \rho_0^*),$$

where P_j , $j = 0, 1$, are defined in (5.25)-(5.26). In the previous subsection we showed that each $P_j(\cdot, \rho_0^*)$ has exactly one zero, μ_j , on the interval $-\infty < \mu < d_* - \Omega_g(\rho_0^*) > 0$. The zeros were simple by our implicit function theorem application in which

$$(5.36) \quad \partial_{\mu}P_j(\mu_j, \rho_0^*) = 1 - \partial_{\mu}C_{jj} > 0;$$

see (5.28). In addition one can easily deduce that $\lim_{\mu \rightarrow -\infty} P_j(\mu, \rho_0^*) = -\infty$ by using the definitions (5.22), (5.12) and the fact that $\|(H - \Omega - \mu)^{-1}\|_{L^2 \rightarrow H^2} \xrightarrow{\mu \rightarrow -\infty} 0$ which implies $\|Q[\mu, \Psi_{\Omega}]\|_{H^2 \rightarrow H^2} \xrightarrow{\mu \rightarrow -\infty} 0$.

Consequently, $D(\cdot, 0)$ has exactly two simple zeros $\mu_0 < 0$ and $\mu_1 = 0$ on the interval $-\infty < \mu \leq (-d_* - \Omega_g(\rho_0^*)) / 2 > 0$, which are both simple and $\lim_{\mu \rightarrow -\infty} D(\mu, 0) = \infty$. It is well known, and a consequence of continuity arguments and of the implicit function theorem, that the previous statement continues to hold for small perturbations. More precisely, there exists $\varepsilon > 0$ such that whenever $|\rho_1| < \varepsilon$, $D(\cdot, \rho_1)$ has exactly two zeros $\mu_0(\rho_1) < 0$ and $\mu_1(\rho_1)$ on the interval $-\infty < \mu \leq (-d_* - \Omega_g(\rho_0^*)) / 2 > 0$, which are both simple and analytic in ρ_1 .

Since we are interested in $n_-(L_+)$, the number of negative eigenvalues of L_+ , we still need to determine the sign of $\mu_1(\rho_1)$. In what follows we will show that its derivatives satisfy

$$(5.37) \quad \partial_{\rho_1} \mu_1(0) = 0, \quad \partial_{\rho_1}^2 \mu_1(0) > 0.$$

We can then conclude that for $0 < \rho_1 \ll 1$, $\mu_1(\rho_1) > 0$ and L_+ has exactly one (simple) negative eigenvalue, $\mu_0(\rho_1)$. Therefore, the asymmetric branch is stable.

We now prove (5.37). By differentiating

$$(5.38) \quad D(\mu_1(\rho_1), \rho_1) = 0$$

once with respect to ρ_1 at $\rho_1 = 0$ we get

$$\partial_\mu D(0, 0) \partial_{\rho_1} \mu_1(0) + \partial_{\rho_1} D(0, 0) = 0.$$

Using (5.35) we obtain

$$(5.39) \quad \partial_\mu D(0, 0) = P_0(0, \rho_0^*) \partial_\mu P_1(\mu_1 = 0, \rho_0^*) > 0,$$

where we used (5.36) and $P_0(0, \rho_0^*) = -M_{00}(\rho_0^*) > 0$. Using (5.34) and (5.23) we obtain

$$(5.40) \quad \partial_{\rho_1} D(0, 0) = \frac{\partial \det(M)}{\partial \rho_1}(\rho_1 = 0) = \det 10 + \det 01,$$

where $\det ij$ is the determinant evaluated at $\rho_1 = 0$ of the matrix obtained from M by differentiating the first row i times, respectively the second row j times. $\det ij$ can be evaluated using (4.34), (4.50), and (5.33).

Note that the second row of $\det 10$ is zero and therefore $\det 10 = 0$. Furthermore, $\det 01$ is zero because its second column is zero. Therefore, by (5.40) we have $\partial_{\rho_1} \mu_1(0) = 0$.

We now calculate $\partial_{\rho_1}^2 \mu_1(\rho_1 = 0)$. Differentiate (5.38) twice with respect to ρ_1 at $\rho_1 = 0$ and use $\partial_{\rho_1} \mu_1(0) = 0$ to obtain

$$\partial_\mu D(0, 0) \partial_{\rho_1}^2 \mu_1(0) + \partial_{\rho_1}^2 D(0, 0) = 0$$

which implies, by (5.39),

$$\text{sign}(\partial_{\rho_1}^2 \mu_1(0)) = -\text{sign}(\partial_{\rho_1}^2 D(0, 0)).$$

But, as before, (5.34) and (5.23) imply

$$\partial_{\rho_1}^2 D(0, 0) = \frac{\partial^2 \det(M)}{\partial \rho_1^2}(0) = \det 20 + 2 \det 11 + \det 02 < 0.$$

The last inequality is a consequence of the following argument. First, $\det 20 = 0$, since its second row is zero. A direct calculation using the definition of M and relations (5.32) show

$$\begin{aligned} \det 11 &= -4(a_{0110} + 2a_{0011})(2a_{1010} + a_{1001})\rho_0^{*2} + \mathcal{O}(\rho_0^{*4}), \\ \det 02 &= 8a_{0000}a_{1111}\rho_0^{*2} + \mathcal{O}(\rho_0^{*4}). \end{aligned}$$

By hypothesis (4.2) and by (3.7), shown in Remark 4.3, we have

$$a_{0000} > a_{1001} + 2a_{1010}, \quad a_{1111} > a_{0110} + 2a_{0011}.$$

This implies, for ρ_0^* sufficiently small, that

$$2 \det 11 + \det 02 < 0.$$

Therefore, $\partial_{\rho_1}^2 \mu_1(0) > 0$ and the proof of Theorem 5.1 is now complete.

Remark 5.3. We remark that, in the special case of $g < 0$ a constant and $V = V_L$ a double-well potential with well-separation parameter L , for L large all coefficients $a_{klmn} = a_{klmn}(L)$ converge to the same value $g\alpha^2 < 0$. This implies

$$2 \det 11 + \det 02 = (-64g^2\alpha^4 + \mathcal{O}(e^{-\tau L}))\rho_0^{*2} + \mathcal{O}(\rho_0^{*4}) < 0.$$

6. Numerical study of symmetry breaking.

6.1. Symmetry-breaking bifurcation for fixed well-separation, L . In this section we numerically compute the bifurcation diagram for the lowest energy nonlinear bound state branch for NLS-GP equation (2.1) and compare these results to the predictions of the finite-dimensional approximation equations (3.12), (3.13). Specifically, we numerically compute the bifurcation structure of (2.1) for a double-well potential, $V_L(x)$, of the form

$$(6.1) \quad V(x) = V_0 \left[\frac{1}{\sqrt{4\pi s^2}} \exp\left(-\frac{(x - L/2)^2}{4s^2}\right) + \frac{1}{\sqrt{4\pi s^2}} \exp\left(-\frac{(x + L/2)^2}{4s^2}\right) \right].$$

The potential for $V_0 = -1$, $s = 1$, and $L = 6$ has two discrete eigenvalues $\Omega_0 = -0.1616$ and $\Omega_1 = -0.12$ and a continuous spectral part for $\Omega > 0$. The linear eigenstates can also be obtained and used to numerically compute the coefficients of the finite-dimensional decomposition of (3.12), (3.13) as $a_{0000} = -0.09397$, $a_{1111} = -0.10375$, $a_{0011} = a_{1010} = a_{1001} = a_{0110} = -0.08836$ (for $g = -1$). Then, using (3.10), we can compute the approximate threshold in \mathcal{N} for bifurcation of an asymmetric branch (and the destabilization of the symmetric one):

$$\mathcal{N}_{cr} \sim \mathcal{N}_{cr}^{(0)} = 0.24331, \quad \Omega_{cr} \sim \Omega_{cr}^{(0)} \equiv \Omega_0 + a_{0000} \mathcal{N}_{cr}^{(0)} = -0.18447.$$

We expect good agreement because the values of s and L suggest the regime of large L , where our rigorous theory holds.

Using numerical fixed-point iterations (in particular Newton’s method), we obtain the branches of the nonlinear eigenvalue problem (2.7). To study the stability of a solution, u_0 , of (2.7), consider the evolution of a small perturbation of it:

$$(6.2) \quad u = e^{-i\Omega t} \left[u_0(x) + \left(p(x)e^{\lambda t} + q(x)e^{\bar{\lambda}t} \right) \right].$$

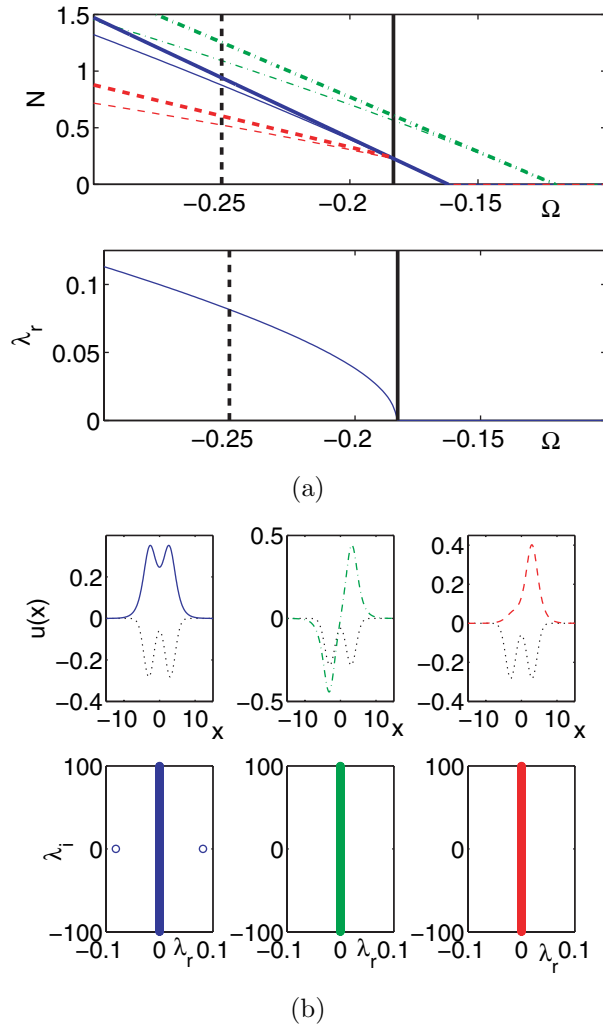


FIG. 3. The figure shows the typical numerical bifurcation results for the cubic case and their comparison with the finite-dimensional analysis of section 3. Panel (a) shows the bifurcation diagram in the top subplot and the relevant real eigenvalues in the bottom subplot. In the top, the solid (blue) line represents the symmetric branch, the dashed-dotted (green) line represents the antisymmetric branch, while the dashed (red) line represents the bifurcating asymmetric branch. The thin lines indicate the numerical findings, while the thick ones show the corresponding finite-dimensional, weakly nonlinear predictions. The solid vertical (black) line indicates the critical point (of $\Omega \approx -0.1835$) obtained numerically. The dashed vertical (black) line is a guide to the eye for the case with $\Omega = -0.25$, whose detailed results are shown in panel (b). The bottom subplot of panel (a) shows the real eigenvalue (as a function of Ω) of the symmetric branch that becomes unstable for $\Omega < -0.1835$. Panel (b) shows, using the same symbolism as panel (a), the symmetric (left), antisymmetric (middle), and asymmetric (right) branches and their linearization eigenvalues (bottom subplots) for $\Omega = -0.25$. The potential is shown by a dotted black line.

Keeping only linear terms in p, q , we obtain a linear evolution equation, whose normal modes satisfy a linear eigenvalue problem with spectral parameter, which we denote by λ and eigenvector $(p(x), \bar{q}(x))^T$.

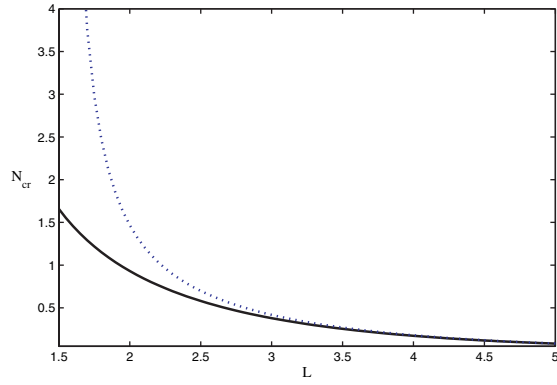
Our computations for the simplest case of the cubic nonlinearity with $K[\psi\bar{\psi}] = \psi\bar{\psi}$ are shown in Figure 3 (for $g(x) = -1$). In particular, the top subplot of panel (a) shows

the full numerical results by thin lines (solid for the symmetric solution, dashed for the bifurcating asymmetric, and dashed-dotted for the antisymmetric one) and compares them with the predictions based on the finite-dimensional truncation of (3.12)-(3.13) shown by the corresponding thick lines. The approximate threshold values \mathcal{N}_{cr} and Ω_{cr} are found numerically to be $\Omega_{cr}^{(0)} \approx -0.1835$ and $\mathcal{N}_{cr}^{(0)} \approx 0.229$. This suggests a relative error in its evaluation by the finite-dimensional reduction of less than 1%. This critical point is indicated by a solid vertical line in panel (a). For $\Omega > \Omega_{cr}^{(0)}$, there exist two branches in the problem, namely the one that bifurcates from the symmetric linear state (this branch exists for $\Omega < \Omega_0$) and the one that bifurcates from the antisymmetric linear state (and, hence, exists for $\Omega < \Omega_1$). For $\Omega < \Omega_{cr}^{(0)}$, the symmetric branch becomes unstable due to a real eigenvalue (see bottom subplot of panel (a)), signalling the emergence of a new branch, namely the asymmetric one. All three branches are shown for $\Omega = -0.25$ (indicated by dashed vertical line in panel (a)) in panel (b) and their corresponding linearization spectrum (λ_r, λ_i) is shown for the eigenvalues $\lambda = \lambda_r + i\lambda_i$.

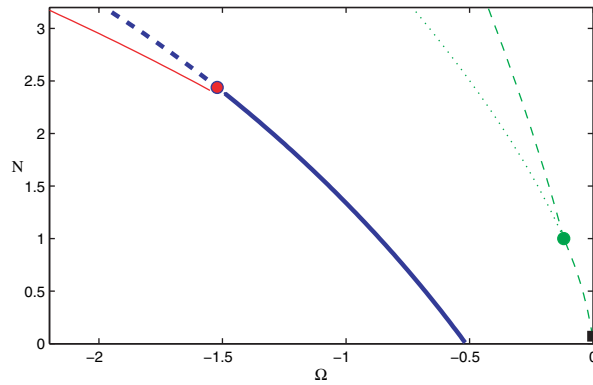
6.2. Symmetry-breaking threshold $\mathcal{N}_{cr}(L)$ as L varies. We now investigate the limits of validity of $\mathcal{N}_{cr}^{(0)}(L)$ as an approximation to $\mathcal{N}_{cr}(L)$ by varying the distance L between the potential wells (6.1). For L large, $\mathcal{N}_{cr}^{(0)}$, given by (3.10), is close to the actual $\mathcal{N}_{cr}(L)$, the exact threshold. In this case the eigenvalues of $-\partial_x^2 + V_L(x)$, $\Omega_0(L)$, and $\Omega_1(L)$ are close to each other; see Example 2.1. Therefore, the bifurcation occurs for small \mathcal{N} and one is in the regime of validity of Theorem 4.1. In Figure 4 we display a comparison between the estimate for \mathcal{N}_{cr} based on the finite-dimensional truncation, $\mathcal{N}_{cr}^{(0)}$, and the actual \mathcal{N}_{cr} . For large L the two values are close to each other. As L is decreased the wells approach one another and eventually, at $L = 0$, merge to form a single-well potential. As L is decreased, the eigenvalues of the linear bound states $\Omega_0(L)$ and $\Omega_1(L)$ move farther apart. For some value of L , L_d , the eigenvalue of the excited state, $\Omega_1(L)$, merges at $\Omega = 0$ into the continuous spectrum (and becomes a complex *scattering resonance*). For $L < L_d$ the estimate $\mathcal{N}_{cr}^{(0)}$ is not correct. In fact, $\mathcal{N}_{cr}^{(0)}(L) \rightarrow \infty$, while the actual value of $\mathcal{N}_{cr}(L)$ appears to remain finite. In Figure 4(a) we observe that for $L < 2$, $\mathcal{N}^{(0)}$ and \mathcal{N}_{cr} diverge from one another and eventually the approximation $\mathcal{N}_{cr}^{(0)}(L)$ tends to infinity, while the actual $\mathcal{N}_{cr}(L)$ remains finite. In Figure 4(b) we show a bifurcation diagram for small L , in which a second discrete (excited state) eigenvalue of $-\partial_x^2 + V_L$, Ω_1 , does not exist. Yet there exists a symmetry-breaking threshold, \mathcal{N}_{cr} , along the symmetric branch.

There are interesting observations to make with regard to antisymmetric solutions. Although there is no linear antisymmetric state, from which to bifurcate in the linear (zero amplitude) limit, we do observe a bifurcation of antisymmetric solutions. The black square in Figure 4(b) marks a *strictly positive* threshold value of the squared L^2 norm, $\mathcal{N}_{cr}^{excited} > 0$, at which an antisymmetric nonlinear bound state emerges from zero frequency. The linearization, L_+ , about this *excited state branch* has two negative eigenvalues: with corresponding symmetric and antisymmetric eigenstates. The bifurcation along this antisymmetric branch at a larger value of \mathcal{N} to an “asymmetric” state (see Figure 4(b)) is the result of a third eigenvalue of L_+ (with corresponding *even* parity eigenfunction) emerging from the continuous spectrum, as \mathcal{N} is increased, and hitting zero at some critical value $\mathcal{N}_{cr}^{excited, asym} > \mathcal{N}_{cr}^{excited}$.

Finally, note that we expect, at least in the regime of weak nonlinearity, that branches of nonlinear bound states originating from nonground state eigenvalues to be unstable. There can be various mechanisms: linear [10, 15, 9, 16] as well as nonlinear.



(a)



(b)

FIG. 4. The figure demonstrates the validity of $N_{cr}^{(0)}(L)$ as an approximation to $N_{cr}(L)$. Panel (a) compares the linear finite-dimensional estimation for the bifurcation point $N_{cr}^{(0)}(L)$ and the actual numerical bifurcation point N_{cr} . The computations are for the double-well potential (6.1), $V_0 = -1$, and $s = 1/4$ and cubic nonlinearity. The curve $N_{cr}(L)$ is marked by a solid (black) line and the curve $N_{cr}^{(0)}(L)$ is marked by a dotted (blue) line. Panel (b) shows a numerical bifurcation diagram for the double-well potential (6.1), $V_0 = -1$, $s = 1/4$, and $L = 1.3$. The bifurcation point N_{cr} is marked by a (red) circle. For $N < N_{cr}$ the ground state, marked by a thick (blue) solid line, is stable. For $N > N_{cr}$ the ground state is unstable and marked by a thick (blue) dashed line. The stable asymmetric state which appears for $N > N_{cr}$ is marked by a thin (red) solid line. The antisymmetric state Ω_1^N is marked by a thin (light green) dashed line. The point N for which the antisymmetric state appears in the discrete spectrum is marked by a (black) square. Notice that in this bifurcation diagram there is also a bifurcation from the antisymmetric branch. The state which bifurcates from the antisymmetric state is marked by a (dark green) thin dotted line. See the text for remarks on the instability of these latter branches.

Concerning the latter, if $2\Omega_1 - \Omega_0 > 0$, then it is known that the excited state is unstable due to resonant coupling to the radiation modes; see, for example, [29, 30]. Linearly this is manifested by an exponential instability, computed via perturbation theory of an embedded eigenvalue of the linearization about the excited state at the bifurcation point (zero amplitude) [5]. If $2\Omega_1 - \Omega_0 < 0$, which is the case in the example of Figure 1, the excited state is linearly stable for small amplitude, but it is nonlinearly unstable due to higher nonlinear order coupling to radiation modes [35].

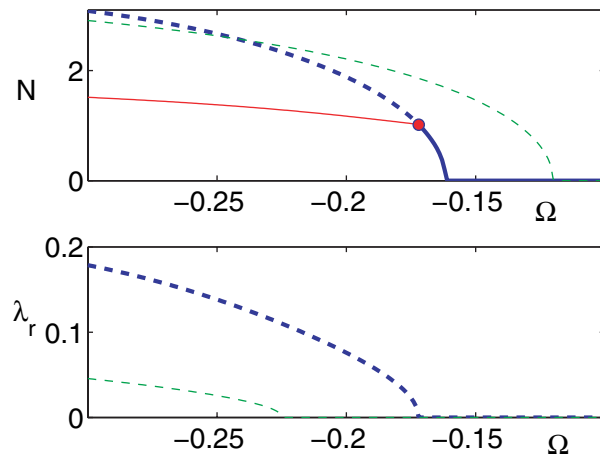


FIG. 5. Same as panel (a) of Figure 3 but for the quintic nonlinearity. This serves to illustrate the analogies between the bifurcation pictures but also their differences (shifted critical point and also partial instability of the antisymmetric branch).

6.3. More general nonlinearities. To simplify the analysis in this paper, we assumed a cubic nonlinearity in NLS-GP; see (H3). Our approach is quite general and more general nonlinearities can be treated. That is, a more general finite-dimensional Galerkin approximation can be derived and its normal form/symmetry/bifurcation theory can be developed.

In this subsection we present numerical computations for general power law nonlinearities such as $K[\psi\bar{\psi}] = (\psi\bar{\psi})^p$ and discuss phenomena analogous to the cubic case, $p = 1$.

Our numerical results for the case of $p = 2$, $K[\psi\bar{\psi}]\psi = |\psi|^4\psi$, are presented in Figure 5. The curves are analogous to those of panel (a) of Figure 3. The bifurcation diagram for this higher order nonlinearity is similar to that of the cubic case. However, the critical point for the emergence of the asymmetric branch is now shifted to $\Omega_{cr} \approx -0.1725$, closer to the linear limit. We have also examined the case of $p = 3$, $K[\psi\bar{\psi}]\psi = |\psi|^6\psi$, finding that the relevant critical point is further shifted toward the linear limit, $\Omega_{cr} = -0.168$. Using our methods one could identify, in the double-well case with large separation, $\Omega_{cr}(p; L)$ such that for all $\Omega < \Omega_{cr}(p; L)$, the symmetric branch is unstable. We also note in passing that bifurcation diagrams for higher values of p may also bear additional (to the shift in Ω_{cr}) differences from the cubic case; one such example in Figure 5 is given by the presence of a linear instability (due to a complex eigenvalue quartet emerging for $\Omega < -0.224$) for the antisymmetric branch. The latter was found to be linearly stable in the cubic case of Figure 3.

6.4. Nonlocal nonlinearities. Finally, we consider the case of nonlocal nonlinearities depending on a parameter ϵ , the range of the nonlocal interaction. In particular, consider the case of a nonlocal nonlinearity of the form

$$(6.3) \quad K[\psi\bar{\psi}] = \int_{-\infty}^{\infty} \mathcal{K}(x-y)\psi(y)\bar{\psi}(y)dy,$$

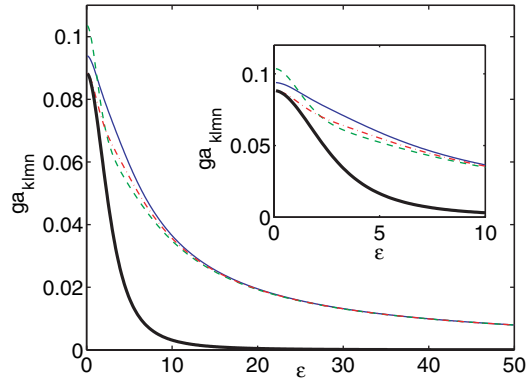
where

$$(6.4) \quad \mathcal{K}(x-y) = \frac{1}{2\pi\epsilon^2} e^{-\frac{(x-y)^2}{2\epsilon^2}}.$$

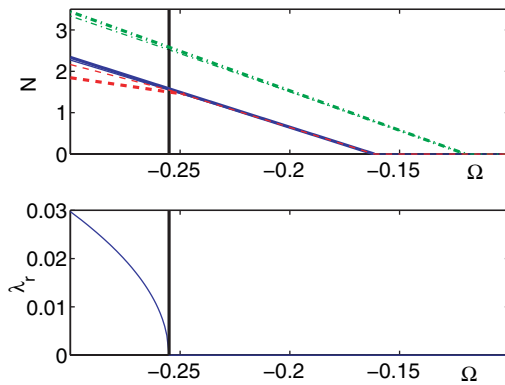
Here, $\epsilon > 0$ is a parameter controlling the range of the nonlocal interaction. As ϵ tends to 0, $\mathcal{K}(x - y) \rightarrow \delta(x - y)$ and we recover the “local” cubic limit. The form of the finite-dimensional reduction, of (3.12), (3.13), does not change; the only modification is that the coefficients a_{klmn} are now functions of the range of the interaction ϵ . The dependence of the coefficients a_{klmn} on ϵ is displayed in panel (a) of Figure 6. The solid (blue) line shows $|a_{0000}|$, the dashed (green) one corresponds to $|a_{1111}|$, the dashed-dotted (red) one corresponds to $|a_{1001}| = |a_{0110}|$ (g is constant, K even), while the thick solid (black) one corresponds to $|a_{0101}| = |a_{0011}| = |a_{1010}|$. Notice in the inset how the coefficients asymptote smoothly to their “local” limit. Additionally, note the expected asymptotic relation $a_{1001} = a_{0011}$. Also note the significant (decaying) dependence of the relevant coefficients on the range of the interaction. The nature of this dependence indicates that while the character of the bifurcation may be the same as in the case of local nonlinearities, its details (such as the location of the critical points) depend sensitively on the range of the nonlocal interaction. This is illustrated in panel (b) for the specific case of $\epsilon = 5$. In this panel (which is analogous to panel (a) of Figure 3, but for the nonlocal case) the critical point for emergence of the asymmetric branch/instability of the symmetric branch is shifted to $\Omega_{cr} = -0.2466$ (and the corresponding $\mathcal{N}_{cr} = 1.4353$) in comparison to the numerically obtained value of $\Omega_{cr} \approx -0.256$; the relative error in the identification of the critical point (by the finite-dimensional reduction) is in this case of the order of 3.7%. This can be attributed to the more nonlinear character (i.e., occurring for higher value of \mathcal{N}_{cr}) nature of the bifurcation. However, as the finite-dimensional approximation still yields a reliable estimate for the location of the critical point, in panel (c) we use it to obtain an approximation to the location of the critical point $(\Omega_{cr}, \mathcal{N}_{cr})$ as a function of the nonlocality parameter ϵ .

7. Concluding remarks. We have obtained rigorous results on the spontaneous symmetry-breaking bifurcation for a large class of NLS-GP equations and studied in detail the case of double-well potentials. Our analysis of the symmetry-breaking bifurcation and the exchange of stability is based on an expansion which, to leading order in amplitude, is a superposition of a symmetric-antisymmetric pair of eigenstates of the linear Hamiltonian, H , whose energies are separated (gap condition (4.13)) from all other spectra of H . This gap condition holds for double wells with sufficiently large L , but breaks down as L decreases. Nevertheless, numerical studies show the existence of a finite threshold for symmetry breaking; see the discussion above on the variation of $\mathcal{N}_{cr}(L)$ with L . A theory encompassing this phenomenon is of interest and is currently under investigation.

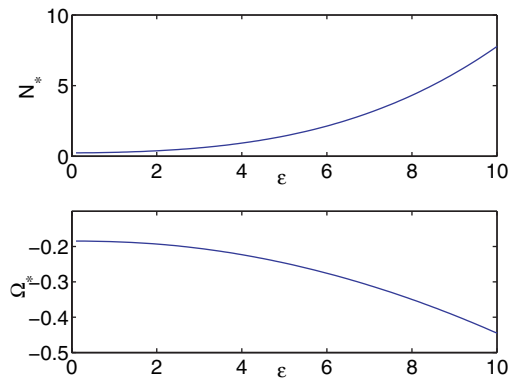
Finally, we remark that our analysis can be naturally extended to treat cases of general multi-wells, identical or not, since the methods involve a strategy for analysis of the weakly nonlinear regime, given spectral assumptions on the linear limit. An example is the case of a symmetric triple well, studied in [18], where the finite-dimensional Galerkin analysis has been implemented, revealing a rich bifurcation picture, but with no symmetry-breaking bifurcation in the symmetric branch. In such multi-well cases, naturally the dimension of the Galerkin approximation needs to be increased accordingly (e.g., 3-dimensional for 3-wells, 4-dimensional for 4-wells, etc.) introducing greater complexity into the global bifurcation structure. Formally, the derivation may be systematically extended to the case of infinitely many wells, constituting the so-called tight-binding approximation [2], although a rigorous derivation of such lattice equations may pose a considerable challenge.



(a)



(b)



(c)

FIG. 6. This figure shows the nonlocal analogue of Figure 3. Panel (a) shows the dependence of the (absolute value of the) coefficients of the finite-dimensional approximation on the nonlocality parameter ϵ ($\epsilon = 0$ denotes the “local” nonlinearity limit). The solid (blue) line denotes a_{0000} , the dashed (green) denotes a_{1111} , the dashed-dotted (red) denotes a_{0110} , while the thick solid (black) one denotes a_{0101} . Panel (b) is analogous to panel (a) of Figure 3, but now shown for the nonlocal case, with the nonlocality parameter $\epsilon = 5$. Finally, panel (c) shows the dependence of the critical point of the finite-dimensional bifurcation (N_*, Ω_*) on the nonlocality parameter ϵ .

8. Appendix. Double wells. In this discussion, we are going to follow the analysis of [12]. Consider a (single-well) real-valued potential $v_0(x)$ on \mathbb{R}^n such that $v_0(x) \in L^r + L^\infty_\epsilon$ for all $1 \leq r \leq q$ where $q \geq \max(n/2, 2)$ for $n \neq 4$, $q > 2$ for $n = 4$. Then, multiplication by v_0 defines a compact operator from H^2 to L^2 and

$$H_0 = -\Delta + v_0(x)$$

is a self-adjoint operator on L^2 with domain H^2 .

Consider now the double-well potential

$$V_L = T_L v_0 T_{-L} + R T_L v_0 T_{-L} R,$$

where T_L and R are the unitary operators

$$\begin{aligned} T_L g(x_1, x_2, \dots, x_n) &= g(x_1 + L, x_2, \dots, x_n), \\ R g(x_1, x_2, \dots, x_n) &= g(-x_1, x_2, \dots, x_n), \end{aligned}$$

and the self-adjoint operator

$$H_L = -\Delta + V_L(x).$$

PROPOSITION 8.1. *Assume that $\omega < 0$ is a nondegenerate eigenvalue of H_0 separated from the rest of the spectrum of H_0 by a distance greater than $2d_*$. Denote by ψ_ω its corresponding eigenvector, $\|\psi_\omega\|_{L^2} = 1$. Then there exists $L_0 > 0$ such that for $L \geq L_0$ the following are true.*

- (i) H_L has exactly two eigenvalues $\Omega_0(L)$ and $\Omega_1(L)$ nearer to ω than $2d_*$. Moreover, $\lim_{L \rightarrow \infty} \Omega_j(L) = \omega$, $j = 0, 1$.
- (ii) One can choose the normalized eigenvectors $\psi_j(L)$, $\|\psi_j(L)\|_{L^2} = 1$, corresponding to the eigenvalues $\Omega_j(L)$, $j = 0, 1$, such that they satisfy

$$\lim_{L \rightarrow \infty} \|\psi_j(L) - (T_L \psi_\omega + (-1)^j R T_L \psi_\omega) / \sqrt{2}\|_{H^2} = 0, \quad j = 0, 1.$$

- (iii) If P_j^L are the orthogonal projections in L^2 onto $\psi_j(L)$, $j = 0, 1$, and $\tilde{P}_L = Id - P_0^L - P_1^L$, then there exists $d > 0$ independent of L such that

$$\|(H_L - \Omega)^{-1} \tilde{P}_L\|_{L^2 \mapsto H^2} \geq d \quad \text{for all } L \geq L_0 \text{ and } |\Omega - \omega| \leq d_*.$$

Proof. For (i) we refer the reader to [12]. The L^2 convergence in (ii) has also been proved there. The H^2 convergence follows from the following compactness argument. Let

$$\psi_j^L = n_L \psi_j(L), \quad j = 0, 1,$$

where n_L is such that $\|\psi_j^L\|_{H^2} = 1$, $j = 0, 1$. From the eigenvector equations: $(H_L - \Omega(L))\psi^L = 0$, where we dropped the index $j = 0, 1$ and the convergence $\Omega(L) \rightarrow \omega$, see part (i), we get

$$(8.1) \quad \lim_{L \rightarrow \infty} \|(-\Delta - \omega + V_L)\psi^L\|_{L^2} = 0.$$

Denote

$$(8.2) \quad g_L = (-\Delta - \omega)\psi^L \in L^2.$$

Since $-\Delta - \omega : H^2 \mapsto L^2$ is bounded, there exists a constant $C > 0$ independent of L such that

$$\|g_L\|_{L^2} \leq C.$$

Since $\omega < 0$, $-\Delta - \omega : H^2 \mapsto L^2$ has a continuous inverse, then (8.1) is equivalent to

$$g_L + V_L(-\Delta - \omega)^{-1}g_L \rightarrow 0 \text{ in } L^2.$$

By expanding V_L we get

$$(8.3) \quad g_L + T_L v_0(-\Delta - \omega)^{-1}T_{-L}g_L + RT_L v_0(-\Delta - \omega)^{-1}T_{-L}Rg_L \rightarrow 0.$$

But $v_0(-\Delta - \omega)^{-1} : L^2 \mapsto L^2$ is compact while the translation and reflection operators are unitary. These and the uniform boundedness of g_L lead to the existence of $\psi \in L^2$ and $\tilde{\psi} \in L^2$ and a subsequence of g_L , which we will redenote by g_L , such that

$$(8.4) \quad \lim_{L \rightarrow \infty} \|v_0(-\Delta - \omega)^{-1}T_{-L}g_L - \psi\|_{L^2} = 0 \text{ and } \lim_{L \rightarrow \infty} \|v_0(-\Delta - \omega)^{-1}T_{-L}Rg_L - \tilde{\psi}\|_{L^2} = 0.$$

By plugging in (8.3) and multiplying to the left by T_{-L} we get

$$\lim_{L \rightarrow \infty} \|T_{-L}g_L + \psi + RT_{2L}\tilde{\psi}\|_{L^2} = 0.$$

But $RT_{2L}\tilde{\psi}$ converges weakly to zero, hence $T_{-L}g_L$ converges weakly to $-\psi$. By plugging now in (8.4) and using compactness we get

$$\psi + v_0(-\Delta - \omega)^{-1}\psi = 0.$$

The latter shows that $(-\Delta - \omega)^{-1}\psi$ is an eigenvector of $-\Delta + v_0$ corresponding to the eigenvalue ω . By nondegeneracy of ω we get

$$(8.5) \quad \psi = -n(-\Delta - \omega)\psi_\omega,$$

where n is a constant. A similar argument shows that

$$(8.6) \quad \tilde{\psi} = -\tilde{n}(-\Delta - \omega)\psi_\omega,$$

where \tilde{n} is a constant.

Combining (8.1)–(8.6) we get

$$(8.7) \quad \lim_{L \rightarrow \infty} \|(-\Delta - \omega)(\psi^L - nT_L\psi_\omega - \tilde{n}RT_L\psi_\omega)\|_{L^2} = 0$$

which by the continuity of $(-\Delta - \omega)^{-1} : L^2 \mapsto H^2$ implies

$$\lim_{L \rightarrow \infty} \|\psi^L - nT_L\psi_\omega - \tilde{n}RT_L\psi_\omega\|_{H^2} = 0.$$

Using now that $\|\psi^L\|_{H^2} = 1$ and that the rescaled ψ_j^L , such that it has norm 1 in L^2 , converges to $(T_L\psi_\omega + (-1)^j RT_L\psi_\omega)/\sqrt{2}$, we get the conclusion of part (ii) for a subsequence first, then, by uniqueness of the limit, for all $L \rightarrow \infty$.

For part (iii), it suffices to show that there are no sequences $(\Omega_L, \psi^L) \in [\omega - d_*, \omega + d_*] \times H^2$ with $\|\psi^L\|_{H^2} = 1$ and $\psi^L \perp \psi_j(L)$, $j = 0, 1$, in L^2 such that

$$(8.8) \quad \lim_{L \rightarrow \infty} \|(H_L - \Omega_L)\psi^L\|_{L^2} = 0.$$

The spectral estimate

$$\|(H_L - \Omega_L)\psi^L\|_{L^2} \geq \text{dist}(\Omega_L, \sigma(H_L) \setminus \{\Omega_0(L), \Omega_1(L)\}) \|\psi^L\|_{L^2} \geq d_* \|\psi^L\|_{L^2},$$

combined with (8.8), implies

$$(8.9) \quad \lim_{L \rightarrow \infty} \|\psi^L\|_{L^2} = 0.$$

In principle, we can now employ the compactness argument in part (ii) to get

$$(8.10) \quad \lim_{L \rightarrow \infty} \|\psi^L\|_{H^2} = 0$$

which will contradict $\|\psi^L\|_{H^2} = 1$. More precisely, (8.8), (8.9) imply

$$\lim_{L \rightarrow \infty} \|(-\Delta - \omega - d_* + V_L)\psi^L\|_{L^2} = 0$$

which, by repeating the argument after (8.1) with ω replaced by $\omega + d_*$, gives

$$\lim_{L \rightarrow \infty} \|\psi^L + T_L \psi_{\omega+d_*} + RT_L \tilde{\psi}_{\omega+d_*}\|_{H^2} = 0,$$

where $\psi_{\omega+d_*}$ and $\tilde{\psi}_{\omega+d_*}$ are eigenvectors of $-\Delta + v_0$ corresponding to eigenvalue $\omega + d_*$. But the latter is not actually an eigenvalue, hence $\psi_{\omega+d_*} = 0$ and $\tilde{\psi}_{\omega+d_*} = 0$. These show (8.10) and finishes the proof of part (iii).

The proposition is now completely proven. \square

Acknowledgments. The authors acknowledge the support of the US National Science Foundation, Division of Mathematical Sciences (DMS). EK acknowledges valuable discussions with A. Malkin, B. Sandstede, and E. Lerman. PGK acknowledges valuable discussions with T. Kapitula and Z. Chen. Finally, the authors thank the referees for their careful reading and comments on the submitted manuscript.

REFERENCES

- [1] M. ALBIEZ, R. GATI, J. FÖLLING, S. HUNSMANN, M. CRISTIANI, AND M. K. OBERTHALER, *Direct observation of tunneling and nonlinear self-trapping in a single bosonic Josephson junction*, Phys. Rev. Lett., 95 (2005), p. 010402.
- [2] G. L. ALFIMOV, P. G. KEVREKIDIS, V. V. KONOTOP, AND M. SALERNO, *Wannier functions analysis of the nonlinear Schrödinger equation with a periodic potential*, Phys. Rev. E, 66 (2002), p. 046608.
- [3] W. H. ASCHBACHER, J. FRÖHLICH, G. M. GRAF, K. SCHNEE, AND M. TROYER, *Symmetry breaking regime in the nonlinear Hartree equation*, J. Math. Phys., 43 (2002), pp. 3879–3891.
- [4] C. CAMBOURNAC, T. SYLVESTRE, H. MAILLOTTE, B. VANDERLINDEN, P. KOCKAERT, PH. EMPLIT, AND M. HAELTERMAN, *Symmetry-breaking instability of multimode vector solitons*, Phys. Rev. Lett., 89 (2002), p. 083901.
- [5] S. CUCCAGNA, D. PELINOVSKY, AND V. VOUGALTER, *Spectra of positive and negative energies in the linearized NLS problem*, Comm. Pure Appl. Math., 58 (2005), pp. 1–29.
- [6] B. DECONINCK AND J. NATHAN KUTZ, *Singular instability of exact stationary solutions of the nonlocal Gross-Pitaevskii equation*, Phys. Lett. A, 319 (2003), pp. 97–103.
- [7] J. FLEISCHER, G. BARTAL, O. COHEN, T. SCHWARTZ, O. MANELA, B. FREEDMAN, M. SEGEV, H. BULJAN, AND N. EFREMIDIS, *Spatial photonics in nonlinear waveguide arrays*, Optics Express, 13 (2005), pp. 1780–1796.
- [8] M. GOLUBITSKY AND D. G. SCHAEFFER, *Singularities and Groups in Bifurcation Theory*, Vol. 1–2, Appl. Math. Sci. 51, Springer, New York, 1985.
- [9] M. G. GRILLAKIS, *Linearized instability for nonlinear Schrödinger and Klein-Gordon equations*, Comm. Pure Appl. Math., 41 (1988), pp. 747–774.

- [10] M. G. GRILLAKIS, J. SHATAH, AND W. A. STRAUSS, *Stability theory of solitary waves in the presence of symmetry I*, J. Funct. Anal., 74 (1987), pp. 160–197.
- [11] J. GUCKENHEIMER AND P. HOLMES, *Nonlinear Oscillations, Dynamical Systems, and Bifurcation of Vector Fields*, Appl. Math. Sci. 42, Springer, New York, 1983.
- [12] E. M. HARRELL, *Double wells*, Comm. Math. Phys., 75 (1980), pp. 239–261.
- [13] T. HEIL, I. FISCHER, AND W. ELSÄSSER, *Chaos synchronization and spontaneous symmetry-breaking in symmetrically delay-coupled semiconductor lasers*, Phys. Rev. Lett., 86 (2001), pp. 795–798.
- [14] R. K. JACKSON AND M. I. WEINSTEIN, *Geometric analysis of bifurcation and symmetry breaking in a Gross-Pitaevskii equation*, J. Statist. Phys., 116 (2004), pp. 881–905.
- [15] C. K. R. T. JONES, *An instability mechanism for radially symmetric standing waves of a nonlinear Schrödinger equation*, J. Differential Equations, 71 (1988), pp. 34–62.
- [16] T. KAPITULA, *Stability of waves in perturbed Hamiltonian systems*, Phys. D, 156 (2001), pp. 186–200.
- [17] T. KAPITULA AND P. KEVREKIDIS, *Bose-Einstein condensates in the presence of a magnetic trap and optical lattice: Two-mode approximation*, Nonlinearity, 18 (2005), pp. 2491–2512.
- [18] T. KAPITULA, P. G. KEVREKIDIS, AND Z. CHEN, *Three is a crowd: Solitary waves in photorefractive media with three potential wells*, SIAM J. Appl. Dyn. Syst., 5 (2006), pp. 598–633.
- [19] T. KAPITULA, P. G. KEVREKIDIS, AND B. SANDSTEDE, *Counting eigenvalues via the Krein signature in infinite-dimensional Hamiltonian systems*, Phys. D, 195 (2004), pp. 263–282.
- [20] P. G. KEVREKIDIS, Z. CHEN, B. A. MALOMED, D. J. FRANTZESKAKIS, AND M. I. WEINSTEIN, *Spontaneous symmetry breaking in photonic lattices: Theory and experiment*, Phys. Lett. A, 340 (2005), pp. 275–280.
- [21] E. KIRR AND A. ZARNESCU, *On the asymptotic stability of bound states in 2D cubic Schrödinger equation*, Comm. Math. Phys., 272 (2007), pp. 443–468.
- [22] W. KROLIKOWSKI, O. BANG, N. I. NIKOLOV, D. NESHEV, J. WYLLER, J. J. RASMUSSEN, AND D. EDMUNDSON, *Modulational instability, solitons and beam propagation in spatially non-local nonlinear media*, J. Opt. B Quantum Semiclass. Opt., 6 (2004), pp. S288–S294.
- [23] K. W. MAHMUD, J. N. KUTZ, AND W. P. REINHARDT, *Bose-Einstein condensates in a one-dimensional double square well: Analytical solutions of nonlinear Schrödinger equation*, Phys. Rev. A, 66 (2002), p. 063607.
- [24] L. NIRENBERG, *Lectures on Nonlinear Functional Analysis*, Courant Institute, New York, 1974.
- [25] C.-A. PILLET AND C. E. WAYNE, *Invariant manifolds for a class of dispersive, Hamiltonian partial differential equations*, J. Differential Equations, 141 (1997), pp. 310–326.
- [26] H. A. ROSE AND M. I. WEINSTEIN, *On the bound states of the nonlinear Schrödinger equation with a linear potential*, Phys. D, 30 (1988), pp. 207–218.
- [27] A. SACCHETTI, *Nonlinear time-dependent one-dimensional Schrödinger equation with double-well potential*, SIAM J. Math. Anal., 35 (2003), pp. 1160–1176.
- [28] S. SAWAI, Y. MAEDA, AND Y. SAWADA, *Spontaneous symmetry breaking Turing-type pattern formation in a confined dictyostelium cell mass*, Phys. Rev. Lett., 85 (2000), pp. 2212–2215.
- [29] A. SOFFER AND M. I. WEINSTEIN, *Selection of the ground state for nonlinear Schrödinger equations*, Rev. Math. Phys., 16 (2004), pp. 977–1071.
- [30] A. SOFFER AND M. I. WEINSTEIN, *Theory of nonlinear dispersive waves and selection of the ground state*, Phys. Rev. Lett., 95 (2005), p. 213905.
- [31] A. G. VANAKARAS, D. J. FOTINOS, AND E. T. SAMULSKI, *Tilt, polarity, and spontaneous symmetry breaking in liquid crystals*, Phys. Rev. E, 57 (1998), pp. R4875–R4878.
- [32] M. I. WEINSTEIN, *Modulational stability of ground states of nonlinear Schrödinger equations*, SIAM J. Math. Anal., 16 (1985), pp. 472–491.
- [33] M. I. WEINSTEIN, *Lyapunov stability of ground states of nonlinear dispersive evolution equations*, Comm. Pure Appl. Math., 39 (1986), pp. 51–68.
- [34] C. YANNOULEAS AND U. LANDMAN, *Spontaneous symmetry breaking in single and molecular quantum dots*, Phys. Rev. Lett., 82 (1999), pp. 5325–5328.
- [35] G. ZHOU, *Perturbation expansion and Nth order Fermi golden rule of the nonlinear Schrödinger equations with potential*, J. Math. Phys., 48 (2007), p. 053509.

ON 2×2 CONSERVATION LAWS AT A JUNCTION*

R. M. COLOMBO[†], M. HERTY[‡], AND V. SACHERS[§]

Abstract. This paper deals with 2×2 conservation laws at a junction. For the Cauchy problem, existence, uniqueness, and Lipschitz continuous dependence of the solution from the initial data as well as from the conditions at the junction are proved. The present construction comprehends the case of the p -system used to describe gas flow in networks and hereby unifies different approaches present in the literature. Furthermore, different models for water networks are considered.

Key words. hyperbolic systems of conservation laws, p -system, St. Venant equations

AMS subject classifications. 35L65, 76N10, 34B45

DOI. 10.1137/070690298

1. Introduction. This paper studies the initial value problem consisting of

$$(1.1) \quad \partial_t u_l + \partial_{x_l} f(u_l) = 0, \quad \text{with } l = 1, \dots, n, \quad t \in [0, +\infty[, \quad x \in [0, +\infty[,$$

along n pipes together with

$$(1.2) \quad \Psi(u_1(t, 0+), u_2(t, 0+), \dots, u_n(t, 0+)) = 0$$

at a junction. In other words, we deal with n initial boundary value problems for systems of 2×2 conservation laws, coupled through nonlinear boundary conditions. In this general setting, extending the results in [10], we derive conditions under which the Cauchy problem for (1.1) has a unique solution. Furthermore, the Lipschitz continuous dependence of the solution on initial data and coupling conditions is proved.

We model a junction connecting n rectilinear pipes by n pairwise distinct vectors ν_1, \dots, ν_n parallel to the pipes and such that $\|\nu_l\|$ equals the cross section of the l th pipe. Furthermore, we assign a space coordinate $x_l > 0$ to each pipe. The transport of a specific quantity in the l th pipe is given by (1.1), where u_l is the vector of variables along the l th duct and f is a general nonlinear flux function. Condition (1.2) describes the interaction of the transported quantities at the intersection of the pipes; see [2, 3, 9, 10, 11, 19]. The standard situation of the Cauchy problem on a line is recovered in the case $n = 2$, $\nu_1 + \nu_2 = 0$, and $\Psi(u_1, u_2) = f(u_1) - f(u_2)$; see section 4.1. To simplify the notation, below we denote by x all the coordinates x_l along the various pipes.

Applications of the theoretical results are in the field of fluid flow in networks and in particular in high-pressure gas pipelines in open canals. In recent years, there has been intense research in flow problems on networks; see, e.g., the book on gas networks [23] and the publications of the Pipeline Simulation Interest Group [26]. Most of the proposed models [12, 14, 21, 23, 24, 25, 27, 22, 18] consider each pipe as a

*Received by the editors April 28, 2007; accepted for publication (in revised form) February 14, 2008; published electronically June 6, 2008. This work was supported by DFG Program 1253 and DAAD Vigoni project D/06/19582.

<http://www.siam.org/journals/sima/40-2/69029.html>

[†]Department of Mathematics, Brescia University, 25133 Brescia, Italy (rinaldo@ing.unibs.it).

[‡]Mathematik, RWTH Aachen, 52056 Aachen, Germany (herty@mathc.rwth-aachen.de).

[§]Department of Mathematics, TU Kaiserslautern, 67653 Kaiserslautern, Germany (sachers@mathematik.uni-kl.de).

one-dimensional domain and use balance laws to describe the dynamics. The validity of one-dimensional models is the subject of intense discussions; we refer the reader to [26] for more details.

In this context, the most challenging and interesting point is the coupling condition at pipe-to-pipe intersections. In an engineering context, this coupling is typically modeled through tables prescribing suitable relations depending on, e.g., the geometry of the pipe, the material, and flow conditions [12, 21]. Among the first mathematical treatments of this situation are [2, 3, 9, 10, 19]. The current presentation considers the *subsonic* case as in [2, 3, 9, 10, 11, 19] and typical physical conditions [14]. Extending the presentations in [4, 16, 20], we consider solutions possibly containing shocks.

Our purpose is to present a general framework for coupling conditions and prove well posedness. Indeed, we unify and extend the approaches [2, 3, 9, 10, 16, 19]. First, the present framework includes the currently used one-dimensional isothermal model for gas flow as well as the shallow-water equations for flows in open channels. In particular, we extend the results in [10] covering not only the isentropic Euler equations but a general 2×2 system. More importantly, the present result allows us to consider any coupling conditions specified through any (possibly nonlinear) function Ψ , the sole constraint being condition (2.2) below. Hence, the present results also extend previous works on open channel flow with gate control or pumps (see [16]), as well as the model for a kink in a pipe introduced in [19]. Finally, within this general setting, we also prove the Lipschitz continuous dependence of the solutions from the condition Ψ at the junction; see (3.2) in Theorem 3.2. As a consequence, in all the cases mentioned above, it is possible to prove the existence of an optimal control.

Numerically, we show that different one-dimensional coupling conditions lead to qualitatively different solutions. The comparison is in the context of the present theory, i.e., for one-dimensional models. Two-dimensional situations are considered, for instance, in [17, 18]. For a comparison with results of the engineering community we refer the reader to the pressure loss tables [12, 21].

The paper is organized as follows. Section 2 is devoted to the Riemann problem and extends [9]. In section 3, the well posedness of the Cauchy problem is stated. Applications of this result to gas flow in pipes as well as flow in open canals are collected in section 4. Section 5 contains the detailed constructions and proofs.

2. The Riemann problem at a junction. Throughout, we refer the reader to [6] for the general theory of hyperbolic systems of conservation laws. Let $\Omega \subseteq \mathbb{R}^2$ be a nonempty open set. Fix a flow $f \in \mathbf{C}^4(\Omega; \mathbb{R}^2)$ satisfying the following assumption:

- (F) There exists a $\bar{u} \in \Omega$ such that $Df(\bar{u})$ admits a strictly negative eigenvalue $\lambda_1(\bar{u})$ and a strictly positive one $\lambda_2(\bar{u})$, the corresponding eigenvectors are linearly independent, and each characteristic field is either genuinely nonlinear or linearly degenerate.

Under this condition, (1.1) generates a standard Riemann semigroup; see [6, Chapter 8]. By *Riemann problem at the junction* we mean the problem

$$(2.1) \quad \begin{cases} \partial_t u_l + \partial_x f(u_l) = 0, & t \in \mathbb{R}^+, \quad l \in \{1, \dots, n\}, \\ u_l(0, x) = \bar{u}_l, & x \in \mathbb{R}^+, \quad u_l \in \Omega, \end{cases}$$

where $\bar{u}_1, \dots, \bar{u}_n$ are constant states in Ω . For $l = 1, \dots, n$, u_l has two components; i.e., $u_l = (u_{l,1}, u_{l,2})$ denotes the densities of the conserved quantities in the l th tube. Here and in what follows, $\mathbb{R}^+ = [0, +\infty[$.

DEFINITION 2.1. Fix a map $\Psi \in \mathbf{C}^1(\Omega^n; \mathbb{R}^n)$. A Ψ -solution to the Riemann problem (2.1) is a function $u: \mathbb{R}^+ \times \mathbb{R}^+ \mapsto \Omega^n$ such that the following hold:

(L) For $l = 1, \dots, n$, the function $(t, x) \mapsto u_l(t, x)$ is self-similar and coincides with the restriction to $x \in \mathbb{R}^+$ of the Lax solution to the standard Riemann problem

$$\begin{cases} \partial_t u_l + \partial_x f(u_l) = 0, \\ u_l(0, x) = \begin{cases} \bar{u}_l & \text{if } x > 0, \\ u_l(1, 0+) & \text{if } x < 0. \end{cases} \end{cases}$$

(Ψ) The trace $u(t, 0+)$ of u at the junction satisfies (1.2) for a.e. $t > 0$. Given an entropy-entropy flux pair (E, F) , the Ψ-solution is entropic at the junction if the following holds:

(E) At the junction, entropy may not decrease; i.e., for a.e. $t > 0$

$$\sum_{l=1}^n \|\nu_l\| F(u_l(t, 0+)) \leq 0.$$

For later use, with a slight abuse of notation, we denote

$$\begin{aligned} \|u\| &= \sum_{l=1}^n \|u_l\| && \text{for } u \in \Omega^n, \\ \|u\|_{\mathbf{L}^1} &= \int_{\mathbb{R}^+} \|u(x)\| dx && \text{for } u \in \mathbf{L}^1(\mathbb{R}^+; \Omega^n), \\ \text{TV}(u) &= \sum_{l=1}^n \text{TV}(u_l) && \text{for } u \in \mathbf{BV}(\mathbb{R}^+; \Omega^n). \end{aligned}$$

The following proposition yields the well posedness of the Riemann problem and the continuous dependence of the solution to the Riemann problem from the initial state and from the function Ψ. These results are used in section 3 to prove well posedness of the Cauchy problem by the wave tracking algorithm. The proofs are deferred to section 5.

PROPOSITION 2.2. Let $n \in \mathbb{N}$ with $n \geq 2$. Fix the pairwise distinct vectors ν_1, \dots, ν_n in $\mathbb{R}^3 \setminus \{0\}$ and an n -tuple of constant states $\hat{u} \in \Omega^n$ giving a stationary solution to the Riemann problem (2.1) in the sense of Definition 2.1. Assume that for $l = 1, \dots, n$, (F) holds in \hat{u}_l . If $\Psi \in \mathbf{C}^1(\Omega^n; \mathbb{R}^n)$ satisfies

$$(2.2) \quad \det \begin{bmatrix} D_1 \Psi(\hat{u}) r_2(\hat{u}_1) & D_2 \Psi(\hat{u}) r_2(\hat{u}_2) & \dots & D_n \Psi(\hat{u}) r_2(\hat{u}_n) \end{bmatrix} \neq 0,$$

where $D_l \Psi = D_{u_l} \Psi$, then there exist positive δ, K such that the following hold:

1. For all $\bar{u} \in \Omega^n$ satisfying $\|\bar{u} - \hat{u}\| < \delta$, the Riemann problem (2.1) admits a unique self-similar solution $(t, x) \mapsto (\mathcal{R}^\Psi(\bar{u}))(t, x)$ in the sense of Definition 2.1.
2. If (1.1) admits an entropy-entropy flux pair (E, F) , requiring that

$$(2.3) \quad \sum_{l=1}^n \|\nu_l\| F(\hat{u}_l) < 0$$

ensures that the solution $(t, x) \mapsto (\mathcal{R}^\Psi(\bar{u}))(t, x)$ is also entropic.

3. If $\bar{u}, \bar{w} \in \Omega^n$ both satisfy $\|\bar{u} - \hat{u}\| < \delta$ and $\|\bar{w} - \hat{u}\| < \delta$, then the traces at the junction of the corresponding solutions to (2.1) satisfy

$$(2.4) \quad \left\| \left((\mathcal{R}^\Psi(\bar{u})) (t, 0+) \right) - \left((\mathcal{R}^\Psi(\bar{w})) (t, 0+) \right) \right\| \leq K \cdot \|\bar{u} - \bar{w}\|.$$

4. For any $\tilde{\Psi} \in \mathbf{C}^1(\Omega^n; \mathbb{R}^n)$ with $\|\tilde{\Psi} - \Psi\|_{\mathbf{C}^1} < \delta$, $\tilde{\Psi}$ also satisfies (2.2) and for all $\bar{u} \in \Omega^n$ satisfying $\|\bar{u} - \hat{u}\| < \delta$,

$$\left\| \left(\mathcal{R}^{\tilde{\Psi}}(\bar{u}) \right) (t) - \left(\mathcal{R}^\Psi(\bar{u}) \right) (t) \right\|_{\mathbf{L}^1} \leq K \cdot \left\| \tilde{\Psi} - \Psi \right\|_{\mathbf{C}^1} \cdot t.$$

In the previous proposition, $r_2(\hat{u})$ is the right eigenvector of $Df(\hat{u})$ corresponding to the second characteristic field.

We remark that, for subsonic initial data, we obtain here a unique Ψ -solution, without any additional condition, such as **(E)**.

3. The Cauchy problem at an intersection. Next we consider the Cauchy problem at a junction. First, we give a definition of a solution which naturally extends Definition 2.1 to the Cauchy problem. Then we prove the existence of solutions for initial data with small total variation and the continuous dependence on the coupling condition Ψ .

DEFINITION 3.1. Fix $\hat{u} \in \Omega^n$ and $T \in]0, +\infty]$. A weak Ψ -solution to

$$(3.1) \quad \begin{cases} \partial_t u_l + \partial_x f(u_l) = 0, & t \in \mathbb{R}^+, & l \in \{1, \dots, n\}, \\ u(0, x) = u_o(x), & x \in \mathbb{R}^+, & u_o \in \hat{u} + \mathbf{L}^1(\mathbb{R}^+; \Omega^n) \end{cases}$$

on $[0, T]$ is a map $u \in \mathbf{C}^0([0, T]; \hat{u} + \mathbf{L}^1(\mathbb{R}^+; \Omega^n))$ such that the following hold:

(W) For all $\varphi \in \mathbf{C}_c^\infty(]-\infty, T[\times \mathbb{R}^+; \mathbb{R})$ and for $l = 1, \dots, n$

$$\int_0^T \int_{\mathbb{R}^+} (u_l \partial_t \varphi + f(u_l) \partial_x \varphi) \, dx \, dt + \int_{\mathbb{R}^+} u_{o,l}(x) \varphi(0, x) \, dx = 0.$$

(Ψ) The condition at the junction is met: for a.e. $t \in \mathbb{R}^+$, $\Psi(u(t, 0+)) = 0$.

If (1.1) admits an entropy-entropy flux pair (E, F) , then the weak Ψ -solution u is entropic if for all $\varphi \in \mathbf{C}_c^\infty(]-\infty, T[\times \mathbb{R}^+; \mathbb{R}^+)$

$$\sum_{l=1}^n \left(\int_0^T \int_{\mathbb{R}^+} (E(u_l) \partial_t \varphi + F(u_l) \partial_x \varphi) \, dx \, dt + \int_{\mathbb{R}^+} E(u_o) \varphi(0, x) \, dx \right) \|\nu_l\| \geq 0.$$

We are now ready to state the main result of this paper, namely the well posedness of the Cauchy problem for (3.1) at the junction.

THEOREM 3.2. Let $n \in \mathbb{N}$, $n \geq 2$. Fix the pairwise distinct vectors ν_1, \dots, ν_n in $\mathbb{R}^3 \setminus \{0\}$. Fix an n -tuple of states $\bar{u} \in \Omega^n$ such that f satisfies **(F)** at \bar{u} and the Riemann problem (2.1) with initial datum \bar{u} admits the stationary solution in the sense of Definition 2.1. Let $\Psi \in \mathbf{C}^1(\Omega^n; \mathbb{R}^n)$ satisfy (2.2). Then there exist positive δ, L and a map $S: [0, +\infty[\times \mathcal{D} \rightarrow \mathcal{D}$ such that

1. $\mathcal{D} \supseteq \{u \in \bar{u} + \mathbf{L}^1(\mathbb{R}^+; \Omega^n): \text{TV}(u) \leq \delta\}$;
2. for $u \in \mathcal{D}$, $S_0 u = u$ and for $s, t \geq 0$, $S_s S_t u = S_{s+t} u$;
3. for $u, w \in \mathcal{D}$ and $s, t \geq 0$, $\|S_t u - S_s w\|_{\mathbf{L}^1} \leq L \cdot (\|u - w\|_{\mathbf{L}^1} + |t - s|)$;

4. if $u \in \mathcal{D}$ is piecewise constant, then for $t > 0$ sufficiently small, $S_t u$ coincides with the juxtaposition of the solutions to Riemann problems centered at the points of jumps or at the junction.

Moreover, for every $u \in \mathcal{D}$, the map $t \mapsto S_t u$ is a Ψ -solution to the Cauchy problem (3.1) according to Definition 3.1.

For any $\tilde{\Psi} \in \mathbf{C}^1(\Omega^n; \mathbb{R}^n)$ with $\|\tilde{\Psi} - \Psi\|_{\mathbf{C}^1} < \delta$, $\tilde{\Psi}$ generates a semigroup of solutions on \mathcal{D} and for $u \in \mathcal{D}$,

$$(3.2) \quad \left\| S_t^{\tilde{\Psi}} \bar{u} - S_t^{\Psi} \bar{u} \right\|_{\mathbf{L}^1} \leq L \cdot \left\| \tilde{\Psi} - \Psi \right\|_{\mathbf{C}^1} \cdot t.$$

If (1.1) admits an entropy-entropy flux pair (E, F) and \bar{u} is strictly entropic in the sense of (2.3), then the Ψ -solution $t \mapsto S_t u$ is entropic at the junction.

The proof is deferred to section 5.

4. Gas networks and open channels. A widely used model for gas flow in pipe networks is the system of isothermal Euler equations; see [23] and the references therein. In this section, we discuss the coupling conditions [2, 10] in the context of the presented theory in the case of a general p -system. By p -system we mean

$$(4.1) \quad \begin{cases} \partial_t \rho_l + \partial_x q_l = 0, & t \in \mathbb{R}^+, \\ \partial_t q_l + \partial_x \left(\frac{q_l^2}{\rho_l} + p(\rho_l) \right) = 0, & x \in \mathbb{R}^+, \\ & l \in \{1, \dots, n\}, \\ & (\rho_l, q_l) \in \mathring{\mathbb{R}}^+ \times \mathbb{R}, \end{cases}$$

where $\rho > 0$ is the mass density of a given fluid, q its linear momentum density, and $p = p(\rho)$ the pressure law, which we assume to satisfy the following:

- (P) $p \in \mathbf{C}^2(\mathbb{R}^+; \mathbb{R}^+)$, $p(0) = 0$ and for all $\rho \in \mathbb{R}^+$, $p'(\rho) > 0$, $p''(\rho) \geq 0$.

In the context of gas pipelines, the pressure law typically chosen is $p(\rho) = a^2 \rho$, where the sound speed a depends on the gas type and temperature [23].

As is well known, (4.1) is equipped with the (mathematical) entropy-entropy flux pair

$$\begin{aligned} E(\rho, q) &= \frac{q^2}{2\rho} + \rho \int_{\rho_*}^{\rho} \frac{p(r)}{r^2} dr \quad (\text{total energy}), \\ F(\rho, q) &= \frac{q}{\rho} \cdot (E(\rho, q) + p(\rho)) \quad (\text{flow of the total energy}) \end{aligned}$$

for a $\rho_* > 0$. Choosing an initial datum $\bar{u} = (\bar{\rho}, \bar{q})$ in the *subsonic* region

$$\Omega = \left\{ (\rho, q) \in \mathring{\mathbb{R}}^+ \times \mathbb{R} : \lambda_1(\rho, q) < 0 < \lambda_2(\rho, q) \right\}$$

ensures that (F) holds at \bar{u} . Recall the standard relations

$$(4.2) \quad \begin{aligned} \lambda_1(\rho, q) &= (q/\rho) - \sqrt{p'(\rho)}, & \lambda_2(\rho, q) &= (q/\rho) + \sqrt{p'(\rho)}, \\ r_1(\rho, q) &= \begin{bmatrix} -1 \\ -\lambda_1(\rho, q) \end{bmatrix}, & r_2(\rho, q) &= \begin{bmatrix} 1 \\ \lambda_2(\rho, q) \end{bmatrix}. \end{aligned}$$

Below we consider different coupling conditions that appeared in the literature. Remark that the geometry of the junction implicitly enters into all the relations below

through the choice of the x_l coordinates. Explicitly, only the area $\|\nu_l\|$ of the cross section appears.

Below we consider the conditions presented in [3, 4, 10, 14, 16, 19, 23]. Any of them prescribes the conservation of mass, so that the first component in (1.2) reads

$$(4.3) \quad \sum_{l=1}^n \|\nu_l\| q_l = 0.$$

With a slight abuse of notation, we use Ψ to denote the remaining $n - 1$ conditions.

- (a) We prescribe the equal momentum flow for all connected pipes; see [9] and in the case of open channels [16]:

$$P(\rho_i(t, 0+), q_i(t, 0+)) = P(\rho_j(t, 0+), q_j(t, 0+)) \quad \forall i \neq j.$$

- (b) We prescribe a single pressure at the pipe-to-pipe intersection; see [2, 3] and [14] in the engineering literature. For open canals, see [15]:

$$p(\rho_i(t, 0+)) = p(\rho_j(t, 0+)) \quad \forall i \neq j.$$

- (c) In case of only two connected pipes, a further coupling condition is proposed in [19] (see (4.7) for the definition of $k(\theta)$):

$$P(\rho_1(t, 0+), q_1(t, 0+)) + k(\theta) = P(\rho_2(t, 0+), q_2(t, 0+)).$$

Some remarks are in order. Most of the engineering literature uses (b) modified by so-called *minor loss* factors. These factors are listed in tables and depend on additional information; see [12, 21].

On the other hand, (a) is not so commonly used in the engineering literature but yields the \mathbf{L}^1 continuity across stationary transonic shocks (see [10, Example 2.3]), which does not hold when (b) or (c) is adopted. A numerical study of a two-dimensional situation for the p -system can be found in [18] and for the Euler system in [17].

4.1. Equal linear momentum flow for the p -system. We consider the setting in [10], i.e., a junction among n pipes, each modeled through the p -system (4.1) with a general pressure law and with the corresponding function Ψ in (1.2) given by (4.3) together with

$$(4.4) \quad \Psi(\rho, q) = \begin{bmatrix} P(\rho_1, q_1) & - & P(\rho_2, q_2) \\ \vdots & \vdots & \vdots \\ P(\rho_{n-1}, q_{n-1}) & - & P(\rho_n, q_n) \end{bmatrix},$$

where $P(\rho, q) = q^2/\rho + p(\rho)$ is the flow of the linear momentum. The well posedness of the Cauchy problem for (4.1) at a junction was proved in [9, Theorem 3.3].

Remark that this choice of Ψ allows us to consider the problem at the junction as an extension of the standard Cauchy problem. Indeed, setting $n = 2$ and $\nu_1 + \nu_2 = 0$, then (4.1)–(4.4) reduces to the usual situation.

In the general case, the total linear momentum Q varies by

$$Q(t_2) - Q(t_1) = \int_{t_1}^{t_2} \sum_{l=1}^n P(\rho_l(t, 0+), q_l(t, 0+)) \nu_l dt = \left(\int_{t_1}^{t_2} P_*(t) dt \right) \sum_{l=1}^n \nu_l$$

(see [9, section 1]), where $P_*(t)$ is the trace of $P(\rho_l(t, 0+), q_l(t, 0+))$, which is independent from l by (1.2)–(4.4). Note that the right-hand side in the equation above depends explicitly on the geometry of the junction.

4.2. Equal pressure for the p -system. In [3], the condition at the junction amounts to mass conservation and to pressure equality. Here we extend that approach to the case of n ducts with a general pressure law, so that (1.2) consists of (4.3) with

$$(4.5) \quad \Psi(\rho, q) = \begin{bmatrix} p(\rho_1) & - & p(\rho_2) \\ \vdots & \vdots & \vdots \\ p(\rho_{n-1}) & - & p(\rho_n) \end{bmatrix}.$$

Using [9, Lemma 4.4], the determinant in condition (2.2) evaluates to

$$(4.6) \quad (-1)^{n+1} \prod_{i=1}^n \lambda_2(\hat{\rho}_i, \hat{q}_i) \cdot \sum_{i=1}^n \prod_{j \neq i} \frac{p'(\hat{\rho}_j)}{\lambda_2(\hat{\rho}_j, \hat{q}_j)}.$$

Since $(\hat{\rho}, \hat{q}) \in \Omega^n$, we have $\lambda_2(\hat{\rho}_i, \hat{q}_i) > 0$, while the assumption **(P)** on the pressure law implies $p'(\hat{\rho}_i) > 0$. Thus, Theorem 3.2 applies, yielding well posedness in the case of n pipes with a general pressure law.

4.3. Two pipes with friction for the p -system. The case studied in [19] corresponds to (4.1) with $p(\rho) = \rho$, $n = 2$, $\nu_1 = [-1 \ 0]$, and $\nu_2 = [\cos \theta \ \sin \theta]$. Here the angular dependence is modeled explicitly and taken into account in the coupling conditions. This situation mimics a kink forming an angle θ in a pipe. Due to the kink, the linear momentum is assumed to vary by a factor k such that

$$(4.7) \quad k = \sqrt{2(1 - \cos \theta)}$$

for $\theta \in [0, \pi/2]$. In the case of pipes possibly with different cross sections $\|\nu_l\|$ and a general pressure law satisfying **(P)**, at the junction we obtain condition (4.3) with

$$(4.8) \quad \Psi(\rho, q) = \left(\frac{q_1^2}{\rho_1} + p(\rho_1) + f k q_1 \right) - \left(\frac{q_2^2}{\rho_2} - p(\rho_2) \right),$$

which reduces to the case considered in [19] when $p(\rho) = \rho$ and equal cross sections. The parameter f denotes a nonnegative empirical friction coefficient. The condition (2.2) is

$$\lambda_2(\hat{\rho}_2, \hat{q}_2) \lambda_2(\hat{\rho}_1, \hat{q}_1) (\lambda_2(\hat{\rho}_2, \hat{q}_2) + \lambda_2(\hat{\rho}_1, \hat{q}_1) + f k) \neq 0.$$

Hence, Theorem 3.2 applies, yielding well posedness for general pressure laws.

4.4. Equal momentum flow for open canals. The model presented in [20, formulae (2.3)–(2.7)] for a node among n open canals reads

$$(4.9) \quad \begin{cases} \partial_t A_l + \partial_x(A_l V_l) = 0, \\ \partial_t V_l + \partial_x \left(\frac{1}{2} V_l^2 + g h(A_l) \right) = 0, \end{cases} \quad \begin{aligned} t &\in \mathbb{R}^+, \\ x &\in \mathbb{R}^+, \\ l &\in \{1, \dots, n\}, \\ (A_l, V_l) &\in \mathbb{R}^+ \times \mathbb{R}, \end{aligned}$$

where V_l is the water speed in the l th canal, A_l is the vertical cross section occupied by the water, g is gravity, and h is the water level. Here, differently from [20], we assume

that the canals' beds are all at the same height above sea level, which is acceptable in a neighborhood of the node among the pipes.

Other descriptions for the dynamics of open canals are found in the literature. In particular, [15, section 6.1] presents a different model, based on [13], that fits into our framework of section 4.2.

The coupling condition at the node is given by [20, formulae (2.9)–(2.16)]:

$$\Psi(A, V) = \begin{bmatrix} \sum_{l=1}^n A_l V_l \\ \frac{1}{2}V_1^2 + gh(A_1) - \frac{1}{2}V_2^2 - gh(A_2) \\ \vdots \quad \quad \quad \vdots \\ \frac{1}{2}V_{n-1}^2 + gh(A_{n-1}) - \frac{1}{2}V_n^2 - gh(A_n) \end{bmatrix}.$$

The first component in Ψ ensures the conservation of water. For a motivation for the other components, we refer the reader to [20, formulae (2.9)–(2.16)].

This condition is analogous to the equal linear momentum flow condition (4.4) discussed in section 4.1.

Choosing an initial datum $\bar{u} = (\bar{A}, \bar{V})$ such that $\bar{V} < \sqrt{\bar{A}gh'(\bar{A})}$ ensures that **(F)** is fulfilled at \bar{u} . In the present case we have

$$\begin{aligned} \lambda_1(A, V) &= V - \sqrt{Agh'(A)}, & \lambda_2(\rho, q) &= V + \sqrt{Agh'(A)}, \\ r_1(A, V) &= \begin{bmatrix} \sqrt{A} \\ -\sqrt{gh'(A)} \end{bmatrix}, & r_2(A, V) &= \begin{bmatrix} \sqrt{A} \\ \sqrt{gh'(A)} \end{bmatrix}. \end{aligned}$$

The determinant in condition (2.2) evaluates therefore to

$$(-1)^{n+1} \left(\prod_{i=1}^n \lambda_2(\bar{A}_i, \bar{V}_i) \right) \cdot \sum_{i=1}^n \sqrt{\bar{A}_i} \prod_{j \neq i} \sqrt{gh'(\bar{A}_j)}.$$

Here $A > 0$ by assumption, as well as $\sqrt{gh'(A)} > 0$ by the monotonicity of $h(A)$, which ensures that $\lambda_2(A, V) > 0$. Thus, Theorem 3.2 can be applied, yielding the well posedness for a junction of n open canals.

4.5. Numerical examples for the p -system. This section is devoted to comparisons among the different coupling conditions at the junction in the case of the p -system (4.1). Throughout, we use the γ -law: $p(\rho) = p_* \cdot (\rho/\rho_*)^\gamma$ with $\gamma = 1.4$, $\rho_* = 1$, and $p_* = 1$, which clearly satisfies **(P)**. In the case of the coupling conditions (4.8), we set $f \equiv 1$.

In general, stationary solutions for (4.4) fail to be stationary solutions for (4.5) or (4.8). Therefore, we perturb below *static* solutions, which are stationary for all coupling conditions and allow the comparisons.

Below the initial data $u_l(x) = (\rho_l(x), q_l(x))$ along the l th duct attains at most two values, say $u_l^\infty = \lim_{x \rightarrow +\infty} u_l(x)$ and $u_l^0 = \lim_{x \rightarrow 0+} u_l(x)$. When waves hit the junction, we solve numerically condition (1.2) using Newton's method and obtain the traces u_l^+ of the solution at the junction. The solution to (2.1) is then computed

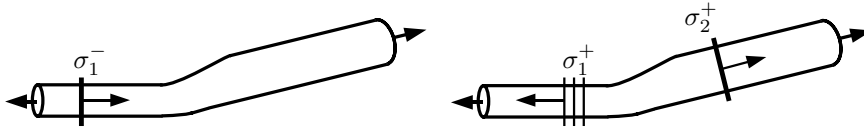


FIG. 1. *Kink with a small angle. Left: a shock approaches a kink with a small angle. Right: a rarefaction is reflected and a shock is refracted.*

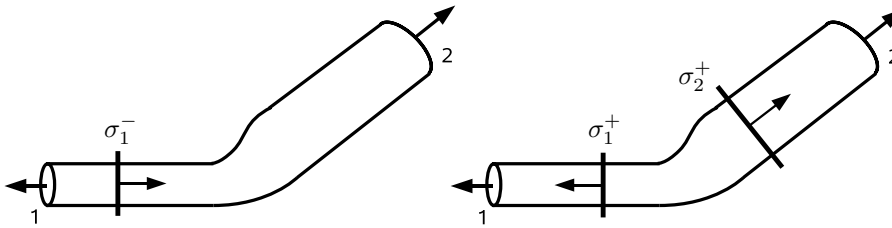


FIG. 2. *Kink with a large angle. Left: a shock approaches a kink. Right: according to condition (4.8), a shock is reflected and a shock is refracted.*

solving a classical Riemann problem between the states u_l^+ (on the left) and u_l^∞ (on the right). Due to the chosen directions of the space variables, waves approaching the junction belong to the first family and those exiting it to the second.

We selected three different examples. In the case of the coupling conditions (4.4) and (4.5), only the cross section of the connected pipes appears explicitly. In the case of (4.8), the angular dependence is taken into account explicitly.

Remark that the numerical values provided in the tables below are expressed in the coordinates x_l adapted to the junction. In particular, the column *Wave speed* is the modulus of the propagation speed of the wave; its x and y components depend on the direction of the pipe. In case of rarefactions, the column *Wave speed* displays the minimal and maximal moduli of the propagation speeds.

4.5.1. Two pipes, different cross sections, possibly different angles.

Consider a junction between $n = 2$ horizontal pipes having cross sections $\|\nu_1\| = 1$ and $\|\nu_2\| = 2$. The situation is depicted in Figures 1 and 2. We choose the cases $\theta = 0, \theta = \pi/4, \theta = \pi/16,$ and $\theta = \pi/32$. In each of these cases, a shock with right state $u_1^0 = [1.1000, -0.1253]$ propagating along pipe 1 hits the junction. At first, we compare conditions (4.4) and (4.5). In Table 1, the first column refers to the coupling condition; in the case of (4.8), Table 2 displays the angle θ . The type of the resulting wave is in the fifth column, where S and R refer to 2-(Lax-)shocks and 2-rarefaction waves, respectively.

These numerical integrations suggest that, according to (4.8), there exists an angle θ_* at which no wave is reflected. Clearly, θ_* depends on the initial states, on the ducts' sections, and on the empirical factor f .

4.5.2. Three pipes, different cross sections.

Consider $n = 3$ pipes identified by $\nu_1 = [-1 - \sqrt{2} \ 0]^T, \nu_2 = [1 \ 1]^T,$ and $\nu_3 = [1 \ 0]^T$ and having cross sections $\|\nu_l\|$. The situation is depicted in Figure 3. Again, we assume that the flow is initially at rest, i.e., $\hat{q}_l = 0$ for $l = 1, 2, 3$. We consider the case of a shock approaching the junction along pipe 1. Due to the choice of the initial data and to the geometry of the junction, the numerical solutions to both coupling conditions (4.4) and (4.5) in fact yield the same results. The final states and the corresponding waves are in Table 3.

TABLE 1

Comparison of results obtained by condition (4.4), (4.8), $\theta = 0$, and (4.5) for two connected pipes with different cross sections.

Ψ	l	u_l^∞	u_l^+	Wave	Wave speed
(4.4), (4.8), $\theta \equiv 0$	1	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.0553 \\ -0.1728 \end{bmatrix}$	R	$\begin{bmatrix} +1.0322 \\ +1.0921 \end{bmatrix}$
	2	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0701 \\ +0.0864 \end{bmatrix}$	S	+1.2324
(4.5)	1	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.0658 \\ -0.1618 \end{bmatrix}$	R	$\begin{bmatrix} +1.0466 \\ +1.0921 \end{bmatrix}$
	2	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0658 \\ +0.0809 \end{bmatrix}$	S	+1.2294

TABLE 2

Comparison for the situation of a kink of angle θ . The case $\theta = 0$ can be found in Table 1.

Ψ	l	u_l^∞	u_l^+	Wave	Wave speed
(4.8) $\theta = \frac{\pi}{32}$	1	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.0645 \\ -0.1633 \end{bmatrix}$	R	$\begin{bmatrix} +1.0447 \\ +1.0921 \end{bmatrix}$
	2	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0664 \\ +0.0816 \end{bmatrix}$	S	+1.1298
(4.8) $\theta = \frac{\pi}{16}$	1	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.0725 \\ -0.1548 \end{bmatrix}$	R	$\begin{bmatrix} +1.0556 \\ +1.0921 \end{bmatrix}$
	2	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0631 \\ +0.0774 \end{bmatrix}$	S	+1.2275
(4.8) $\theta = \frac{\pi}{4}$	1	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.1056 \\ -0.1192 \end{bmatrix}$	S	+1.0958
	2	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0489 \\ +0.0596 \end{bmatrix}$	S	+1.2177

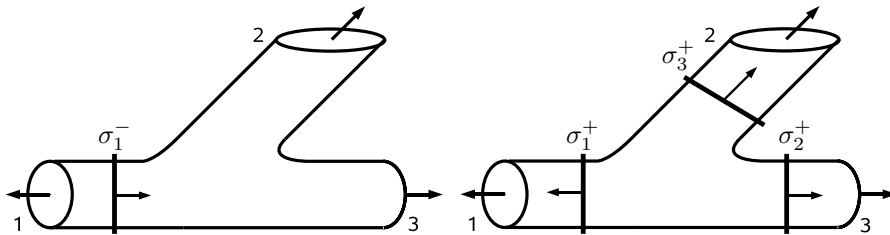


FIG. 3. A shock hitting a T-junction.

4.5.3. Four pipes, different cross sections. Finally, we consider a junction with $n = 4$ pipes defined by $\nu_1 = [-1 \ 0]$, $\nu_2 = [0 \ 1]$, $\nu_3 = [0 \ -1]$, and $\nu_4 = [1 \ 0]$. The situation is depicted in Figure 4. Initially the gas flow is at rest, i.e., $\hat{q}_l = 0$, and $\hat{\rho}_l = 1$ for $l = 1, \dots, 4$. We let three 1-Lax-shocks of different strengths collide simultaneously at the junction along pipes 1, 3, and 4.

It is remarkable that the two coupling conditions yield qualitatively different results. The wave reflected in tube 1 is a rarefaction according to (4.4) and a shock according to (4.5). However, the propagation speeds in the two cases are close to each other, coherently with the \mathbf{L}^1 continuous dependence.

Table 4 displays, for each duct, the type of wave arising after the interaction at the junction, and its speed is reported below for the two different coupling conditions (4.4) and (4.5).

TABLE 3
Numerical results for section 4.5.2.

(4.4) and (4.5)	$l = 1$	$l = 2$	$l = 3$
u_l^∞	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$
u_l^0	$\begin{bmatrix} +1.2000 \\ -0.2642 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$
u_l^+	$\begin{bmatrix} +1.2002 \\ -0.2640 \end{bmatrix}$	$\begin{bmatrix} +1.2002 \\ +0.2640 \end{bmatrix}$	$\begin{bmatrix} +1.2002 \\ +0.2640 \end{bmatrix}$
Wave type	S	S	S
Wave speed	+1.0071	+1.3206	+1.3206

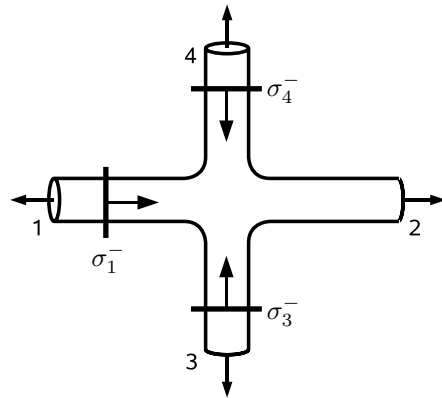


FIG. 4. Three Lax shocks hit the junction simultaneously.

TABLE 4
Numerical results for section 4.5.3.

	$l = 1$	$l = 2$	$l = 3$	$l = 4$
u_l^∞	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$
u_l^0	$\begin{bmatrix} +1.1500 \\ -0.1931 \end{bmatrix}$	$\begin{bmatrix} +1.0000 \\ +0.0000 \end{bmatrix}$	$\begin{bmatrix} +1.1000 \\ -0.1253 \end{bmatrix}$	$\begin{bmatrix} +1.0500 \\ -0.0609 \end{bmatrix}$
Coupling condition (4.4)				
u_l^+	$\begin{bmatrix} +1.1396 \\ -0.2039 \end{bmatrix}$	$\begin{bmatrix} +1.1441 \\ +0.1848 \end{bmatrix}$	$\begin{bmatrix} +1.1625 \\ -0.0545 \end{bmatrix}$	$\begin{bmatrix} +1.1611 \\ +0.0736 \end{bmatrix}$
Wave type	R	S	S	S
Wave speed	$\begin{bmatrix} +1.0357 \\ +1.0489 \end{bmatrix}$	+1.2831	+1.1328	+1.2113
Coupling condition (4.5)				
u_l^+	$\begin{bmatrix} +1.1520 \\ -0.1909 \end{bmatrix}$	$\begin{bmatrix} +1.1520 \\ +0.1957 \end{bmatrix}$	$\begin{bmatrix} +1.1520 \\ -0.0667 \end{bmatrix}$	$\begin{bmatrix} +1.1520 \\ +0.0620 \end{bmatrix}$
Wave type	S	S	S	S
Wave speed	+1.0501	+1.2884	+1.1260	+1.2053

5. Technical details. As a general reference on the theory of hyperbolic systems of conservation laws, we refer the reader to [6]. Denote by $\sigma \mapsto \mathcal{L}_i(u_o, \sigma)$ the i th Lax curve through u_o for $i = 1, 2$. As usual, $\mathcal{O}(1)$ denotes a sufficiently large constant dependent only on f restricted to a neighborhood of the initial states.

Proof of Proposition 2.2. It is sufficient to show that for all \bar{u} sufficiently near \hat{u} the nonlinear system $\Psi(\mathcal{L}_2(\bar{u}_1, \sigma_1), \dots, \mathcal{L}_2(\bar{u}_n, \sigma_n)) = 0$ admits a unique n -tuple of solution $\sigma_1, \dots, \sigma_n$. Indeed, condition (2.2) allows us to use the implicit function theorem.

The Lipschitz estimate (2.4) follows immediately from the regularity of the implicit function.

To prove the latter statement, it is sufficient to consider a single pipe, say the first one. Let u , respectively, \tilde{u} , be the solutions to (2.1) with condition Ψ , respectively, $\tilde{\Psi}$, at the junction. If u and \tilde{u} contain a single shock, then

$$\begin{aligned} & \|u - \tilde{u}\| \\ &= \int_0^{\min\{\Lambda^1, \Lambda^2\}t} \|u(x) - \tilde{u}(x)\| dx \\ & \quad + \int_{\min\{\Lambda^1, \Lambda^2\}t}^{\max\{\Lambda^1, \Lambda^2\}t} \|u(x) - \tilde{u}(x)\| dx \\ &= \mathcal{O}(1) \cdot \|u(t, 0+) - \tilde{u}(t, 0+)\| \cdot t + \mathcal{O}(1) \cdot |\Lambda^2 - \Lambda^1| \cdot t \\ &= \mathcal{O}(1) \cdot \left\| \tilde{\Psi} - \Psi \right\|_{\mathbf{C}^1} \cdot t + \mathcal{O}(1) \cdot \|u(t, 0+) - \tilde{u}(t, 0+)\| \cdot t \\ &= \mathcal{O}(1) \cdot \left\| \tilde{\Psi} - \Psi \right\|_{\mathbf{C}^1} \cdot t. \end{aligned}$$

The case of one or both solutions containing a rarefaction is similar; see also the proof of [5, Corollary 2.5]. \square

We now pass to the proof of Theorem 3.2. To do this, we first use wave front tracking to construct approximate solutions to the Cauchy problem (3.1) adapting the wave front tracking technique; see [6, Chapter 7].

Let $\hat{\delta} > 0$ be such that the closed sphere $\overline{B(\hat{u}_l, \hat{\delta})} \subset \Omega$ for $l = 1, \dots, n$, and introduce the compact set $\mathcal{B} = \prod_{l=1}^n \overline{B(\hat{u}_l, \hat{\delta})}$. We omit the proof of the following simple estimate.

LEMMA 5.1. *For all $u \in \overline{B(\hat{u}, \hat{\delta})}$ there exists $C > 0$ such that if $\|u - \hat{u}\| < \hat{\delta}$ and $\|\mathcal{L}_i(u; \sigma) - u\| \leq \hat{\delta}$, then*

$$\frac{1}{C} \cdot |\sigma| \leq \|\mathcal{L}_i(u; \sigma) - u\| \leq C \cdot |\sigma|.$$

Fix $\varepsilon > 0$. Approximate the initial datum u_o with a sequence $u_{o,\varepsilon}$ of piecewise constant initial data each having a finite number of discontinuities so that $\lim_{\varepsilon \rightarrow 0} \|u_{o,\varepsilon} - u_o\|_{\mathbf{L}^1} = 0$. Then, at the junction and at each point of jump in the approximate initial datum along the pipes, we solve the corresponding Riemann problem according to Definition 2.1. If the total variation of the initial datum is sufficiently small, then Proposition 2.2 ensures the existence and uniqueness of solutions to the Riemann problem. We approximate each rarefaction wave with a rarefaction fan,

i.e., by means of (nonentropic) shock waves traveling at the characteristic speed of the state to the right of the shock and with size at most ε .

This construction can be extended up to the first time \bar{t}_1 at which two waves interact in a pipe or a wave hits the junction. Clearly, at time \bar{t}_1 the functions so constructed are piecewise constant with a finite number of discontinuities. Hence, at any subsequent interaction or collision with the junction, we repeat the previous construction with the following provisions:

1. no more than two waves interact at the same point or at the junction;
2. a rarefaction fan of the i th family produced by the interaction between an i th rarefaction and any other wave is *not* split any further;
3. when the product of the strengths of two interacting waves falls below a threshold $\tilde{\varepsilon}$, then we let the waves cross each other, their size being unaltered, and introduce a *nonphysical* wave with speed $\hat{\lambda}$, with $\hat{\lambda} > \sup_u \lambda_2(u)$; see [6, Chapter 7] and the refinement [1].

In the present case, we have to complete the above algorithm stating how the Riemann problem at the junction is to be solved. At time $t = 0$ and whenever a physical wave with size greater than $\tilde{\varepsilon}$ hits the junction, the accurate solver is used; i.e., the exact solution as in Definition 2.1 is approximated by replacing rarefaction waves with rarefaction fans. When a wave with strength smaller than $\tilde{\varepsilon}$ hits the junction, then we let it be reflected into a nonphysical wave with speed $\hat{\lambda}$, and no wave in any other pipe is produced.

Repeating recursively this procedure, we construct a wave front tracking sequence of approximate solutions u_ε in the sense of [6, Definition 7.1].

At interactions of waves in a pipe, we have the following classical result.

LEMMA 5.2. *There exists a constant K with the following property.*

1. *If there is an interaction in a pipe between two waves σ_1^- and σ_2^- , respectively, of the first and second family, producing the waves σ_1^+ and σ_2^+ (see Figure 5, left), then*

$$(5.1) \quad \left| \sigma_1^+ - \sigma_1^- \right| + \left| \sigma_2^+ - \sigma_2^- \right| \leq K \cdot \left| \sigma_1^- \sigma_2^- \right|.$$

2. *If there is an interaction in a pipe between two waves σ'_i and σ''_i of the same i th family producing waves of total size σ_1^+ and σ_2^+ (see Figure 5, right, for the case $i = 2$), then*

$$\left| \sigma_1^+ - (\sigma''_1 + \sigma'_1) \right| + \left| \sigma_2^+ \right| \leq K \cdot \left| \sigma'_1 \sigma''_1 \right| \quad \text{if } i = 1,$$

$$\left| \sigma_1^+ \right| + \left| \sigma_2^+ - (\sigma''_2 + \sigma'_2) \right| \leq K \cdot \left| \sigma'_2 \sigma''_2 \right| \quad \text{if } i = 2.$$

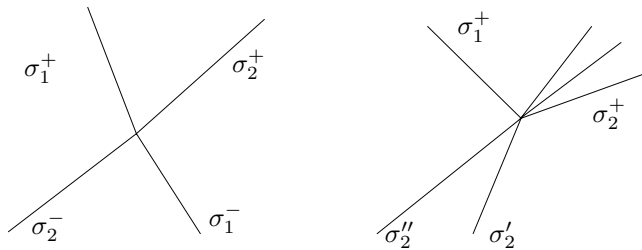


FIG. 5. Notation for the standard interaction estimates in Lemma 5.2.

3. If there is an interaction in a pipe between two physical waves σ_1^- and σ_2^- producing a nonphysical wave σ_3^+ (see Figure 6, left), then

$$|\sigma_3^+| \leq K \cdot |\sigma_1^- \sigma_2^-|.$$

4. If there is an interaction in a pipe between a physical wave σ and a nonphysical wave σ_3^- producing a physical wave σ and a nonphysical wave σ_3^+ (see Figure 6, right), then

$$|\sigma_3^+| - |\sigma_3^-| \leq K \cdot |\sigma \sigma_3^-|.$$

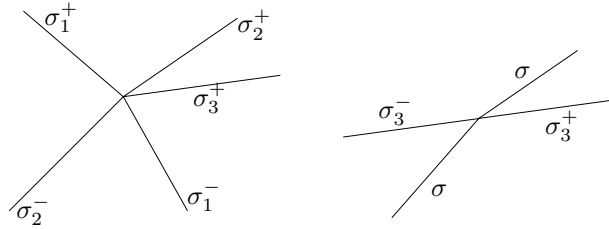


FIG. 6. Left: a nonphysical wave arises. Right: a nonphysical wave hits a physical one.

For a proof of this result see [6, Chapter 7]. By construction, nonphysical waves cannot interact with the junction or with other nonphysical waves. In the case of the junction, we have the following result, with notation as depicted in Figure 7.

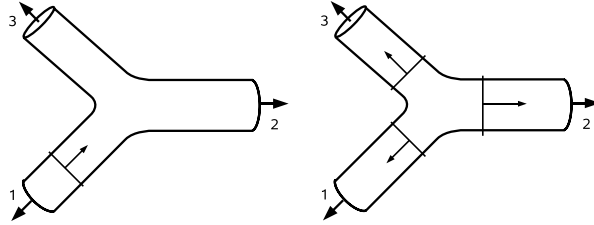


FIG. 7. Notation for Proposition 5.3. Left: before the interaction. Right: after the interaction.

PROPOSITION 5.3. *There exist $\delta_J > 0$ and $K_J \geq 1$ with the following property. For any $\bar{u} \in \mathcal{B}$ that yields a stationary solution to the Riemann problem (2.1), and for any 1-waves $\sigma_l^- \in]-\delta_J, \delta_J[$ hitting the junction and producing the 2-waves σ_l^+ ,*

$$(5.2) \quad \sum_{l=1}^n |\sigma_l^+| \leq K_J \cdot |\sigma_l^-|.$$

The proof follows immediately from (2.4) and Lemma 5.1.

Define $\tilde{K} = 2KK_J + 1$. Fix a wave front tracking approximate solution u_ε . For $t > 0$ and $l \in \{1, \dots, n\}$, we denote with $\{x_{l,\alpha} : \alpha \in \mathcal{J}_l(u)\}$ the set of the positions of the discontinuities of the approximate solution u in the l th pipe and with $\sigma_{l,1,\alpha}$, $\sigma_{l,2,\alpha}$, $\sigma_{l,3,\alpha}$ the strengths of the waves, respectively, of the first family, of the second family, and of the nonphysical waves at $x_{l,\alpha}$. Introduce the Glimm-type functionals

$$(5.3) \quad V(t) = \sum_{l=1}^n \sum_{\alpha \in \mathcal{J}_l} \left(2K_J \cdot |\sigma_{l,1,\alpha}| + |\sigma_{l,2,\alpha}| + |\sigma_{l,3,\alpha}| \right),$$

$$Q(t) = \sum_{l=1}^n \sum \left\{ |\sigma_{l,i,\alpha} \sigma_{l,j,\beta}| : (\sigma_{l,i,\alpha}, \sigma_{l,j,\beta}) \in \mathcal{A}_l \right\},$$

$$\Upsilon(t) = V(t) + \check{K} \cdot Q(t),$$

where \mathcal{A}_l denotes the set of approaching waves in the l th pipe; see [6, section 7.3].

The functionals above are well defined for every $t > 0$ at which no interaction takes place. Now suppose that at a time $\tau > 0$ there is an interaction between the wave $\sigma_{\bar{l},1}$ of the first family and the junction. In general, this interaction produces n waves $\sigma_{l,2}^+$ of the second family. Thus

$$\Delta V(\tau) \leq \left(\sum_{l=1}^n |\sigma_{l,2}^+| \right) - 2 K_J |\sigma_{\bar{l},1}| \leq -K_J \cdot |\sigma_{\bar{l},1}|,$$

$$\Delta Q(\tau) \leq \sum_{l=1}^n \left(|\sigma_{l,2}^+| \sum_{i,\alpha} |\sigma_{l,i,\alpha}| \right) \leq K_J \cdot V(\tau-) \cdot |\sigma_{\bar{l},1}|,$$

$$\Delta \Upsilon(\tau) \leq K_J \cdot (\check{K} \cdot V(\tau-) - 1) \cdot |\sigma_{\bar{l},1}|.$$

Now suppose an interaction between two waves $\sigma_{l,i,\alpha}, \sigma'_{l,j,\beta}$ happens in a pipe at time τ . By Lemma 5.2 we deduce that

$$\Delta V(\tau) \leq 2 \cdot K \cdot K_J \cdot |\sigma_{l,i,\alpha} \sigma'_{l,j,\beta}|,$$

$$\Delta Q(\tau) \leq K \cdot |\sigma_{l,i,\alpha} \sigma'_{l,j,\beta}| \cdot V(\tau-) - |\sigma_{l,i,\alpha} \sigma'_{l,j,\beta}|,$$

$$\Delta \Upsilon(\tau) \leq |\sigma_{l,i,\alpha} \sigma'_{l,j,\beta}| \left(K \cdot (2 \cdot K_J + \check{K} \cdot V(\tau-)) - \check{K} \right).$$

We have thus proved the following basic result.

PROPOSITION 5.4. *Let $\delta = \min\{1/(2\check{K} + 1), 1/(2K\check{K} + 1), \hat{\delta}, \delta_J\}$. At any interaction time $\tau > 0$, if $V(\tau-) < \delta$, then $\Delta \Upsilon(\tau) < 0$ with Υ defined in (5.4).*

Proof of Theorem 3.2. Let δ be as in Proposition 5.4, and define

$$\tilde{\mathcal{D}} = \left\{ u \in \hat{u} + \mathbf{L}^1 \left(\mathbb{R}^+; (\mathring{\mathbb{R}}^+ \times \mathbb{R})^n \right) : u \in \mathbf{PC} \text{ and } \Upsilon(u) \leq \delta \right\};$$

here \mathbf{PC} denotes the set of piecewise constant functions with finitely many jumps. It is immediate to prove that there exists a suitable $C_1 > 0$ such that $\frac{1}{C_1} \text{TV}(u)(t, \cdot) \leq V(t) \leq C_1 \text{TV}(u)(t, \cdot)$ for all $u \in \tilde{\mathcal{D}}$. Any initial data in $\tilde{\mathcal{D}}$ yields an approximate solution to (1.1) attaining values in $\tilde{\mathcal{D}}$ by Proposition 5.4.

We now pass to the \mathbf{L}^1 -Lipschitz continuous dependence of the approximate solutions from the initial datum. Consider two wave front tracking approximate solutions u_1 and u_2 . Define the functional

$$(5.4) \quad \Phi(u_1, u_2) = \sum_{l=1}^n \sum_{i=1}^2 \int_0^{+\infty} |s_{l,i}(x)| W_{l,i}(x) dx,$$

where $s_{l,i}(x)$ measures the strengths of the i th shock wave in the l th pipe at point x (see [6, Chapter 8]) and the weights $W_{l,i}$ are defined by

$$W_{l,1}(x) = \hat{K} \cdot (1 + \kappa_1 A_{l,i}(x) + \kappa_1 \kappa_2 (\Upsilon(u_1) + \Upsilon(u_2))),$$

$$W_{l,2}(x) = 1 + \kappa_1 A_{l,i}(x) + \kappa_1 \kappa_2 (\Upsilon(u_1) + \Upsilon(u_2))$$

for suitable positive constants κ_1, κ_2 chosen as in [6, formula (8.7)] and $\hat{K} = 1 + (\max_l W_{l,2}) K_J \hat{\lambda} / |\inf \lambda_1|$. Here Υ is the functional defined in (5.4), while the $A_{l,i}$ are defined by

$$A_{l,i}(x) = \sum \left\{ |\sigma_{l,k_\alpha,\alpha}| : \begin{array}{l} x_\alpha < x, i < k_\alpha \leq 2 \\ x_\alpha > x, 1 \leq k_\alpha < i \end{array} \right\}$$

$$+ \begin{cases} \left\{ \sum \left\{ |\sigma_{l,i,\alpha}| : \begin{array}{l} x_\alpha < x, \alpha \in \mathcal{J}_l(u_1) \\ x_\alpha > x, \alpha \in \mathcal{J}_l(u_2) \end{array} \right\} \right\} & \text{if } s_{l,i}(x) < 0, \\ \left\{ \sum \left\{ |\sigma_{l,i,\alpha}| : \begin{array}{l} x_\alpha < x, \alpha \in \mathcal{J}_l(u_2) \\ x_\alpha > x, \alpha \in \mathcal{J}_l(u_1) \end{array} \right\} \right\} & \text{if } s_{l,i}(x) \geq 0; \end{cases}$$

see [6, Chapter 8]. We first fix κ_1, κ_2 , so that δ in the definition of $\tilde{\mathcal{D}}$ can be chosen to make $W_{l,i}(x) \geq 1$ and uniformly bounded for every $l \in \{1, \dots, n\}$, $i \in \{1, 2\}$, and $x \geq 0$. Hence the functional Φ is equivalent to \mathbf{L}^1 distance:

$$\frac{1}{C_2} \cdot \|u_1 - u_2\|_{\mathbf{L}^1} \leq \Phi(u_1, u_2) \leq C_2 \cdot \|u_1 - u_2\|_{\mathbf{L}^1}$$

for a positive constant C_2 . The same calculations as in [6, Chapter 8] show that, at any time $t > 0$ when an interaction happens neither in u_1 nor in u_2 ,

$$\begin{aligned} & \frac{d}{dt} \Phi(u_1(t), u_2(t)) \\ & \leq C_3 \varepsilon + \sum_l \sum_i |s_{l,i}(0+)| W_{l,i} \lambda_{l,i}(0+) \\ & \leq C_3 \varepsilon - \left(\sum_l |s_{l,1}(0+)| \right) \hat{K} \left| \inf_l \lambda_{l,1} \right| + K_J \left(\sum_l |s_{l,1}(0+)| \right) \max_l W_{l,2} \hat{\lambda} \\ & \leq C_3 \varepsilon + \sum_l |s_{l,1}(0+)| \left(K_J \max_l W_{l,2} \hat{\lambda} - \hat{K} \left| \inf_l \lambda_{l,1} \right| \right) \\ & \leq C_3 \varepsilon, \end{aligned}$$

where C_3 is a suitable positive constant depending only on a bound on the total variation of the initial data. Above we used the analogue of Proposition 5.3 for shock curves, i.e., if $\Psi(S_2(S_1(u, q_{l,1}), q_{l,2})) = 0$, then $\sum_l |q_{l,2}| \leq K_J \sum_l |q_{l,1}|$, for a suitable K_J .

If $t > 0$ is an interaction time for u_1 or u_2 , then, by Proposition 5.4, $\Delta[\Upsilon(u_1(t)) + \Upsilon(u_2(t))] < 0$ and, choosing κ_2 large enough, we obtain

$$\Delta \Phi(u_1(t), u_2(t)) < 0.$$

Thus, $\Phi(u_1(t), u_2(t)) - \Phi(u_1(s), u_2(s)) \leq C_2 \varepsilon(t - s)$ for every $0 \leq s \leq t$. The proof is now completed using the standard arguments in [6, Chapter 8].

The proof that the semigroup trajectory does indeed yield a solution to (3.1) and, in particular, that (Ψ) is satisfied on the traces is exactly as that of [7, Proposition 5.3].

We now pass to the stability estimate (3.2). Its proof is similar to those of [5, Theorem 2.1] or [8, Theorem 3.1] and is based on [6, Theorem 2.9], which we recall for convenience: for every Lipschitz map $w: [0, T] \rightarrow \mathcal{D}$ and every Lipschitz semigroup $S: [0, +\infty[\times \mathcal{D} \rightarrow \mathcal{D}$, the following estimate holds:

$$\|w(T) - S_T w(0)\|_{\mathbf{L}^1} \leq L \cdot \int_0^T \left(\liminf_{h \rightarrow 0^+} \frac{1}{h} \|w(t+h) - S_h w(t)\|_{\mathbf{L}^1} \right) dt,$$

L being the Lipschitz constant of S . In the present case, we are led to

$$\left\| S_t^{\tilde{\Psi}} u - S_t^{\Psi} u \right\|_{\mathbf{L}^1} \leq L^{\Psi} \cdot \int_0^t \left(\liminf_{h \rightarrow 0^+} \frac{1}{h} \left\| S_h^{\tilde{\Psi}} S_{\tau}^{\Psi} u - S_h^{\Psi} S_{\tau}^{\Psi} u \right\| \right) d\tau.$$

It remains to estimate $\|S_h^{\tilde{\Psi}} v - S_h^{\Psi} v\|_{\mathbf{L}^1}$ for $v = S_{\tau}^{\Psi} u$. We use the wave front tracking approximations $v^{\Psi, \varepsilon}(h, \cdot) = S_h^{\Psi, \varepsilon} u$ and $v^{\tilde{\Psi}, \varepsilon}(h, \cdot) = S_h^{\tilde{\Psi}, \varepsilon} u$. For $h > 0$ sufficiently small, we can assume that there is at most a single interaction of the waves of $u \in \mathcal{D}$ with the intersection. Then $S_h^{\tilde{\Psi}, \varepsilon}$ coincides with the Riemann solver of Proposition 2.2, and estimate (2.4) can be applied:

$$\left\| S_h^{\tilde{\Psi}} u - S_h^{\Psi} u \right\| \leq L^{\Psi} \|\Psi_1 - \Psi_2\|_{\mathbf{C}^1} h.$$

Since the right-hand side is independent of ε and since $S_h^{\tilde{\Psi}, \varepsilon} u$ converges in \mathbf{L}^1 to $S_h^{\tilde{\Psi}} u$, we obtain

$$\left\| S_h^{\tilde{\Psi}} u - S_h^{\Psi} u \right\|_{\mathbf{L}^1} \leq \int_0^t \left(\liminf_{h \rightarrow 0^+} \frac{1}{h} \cdot \left\| \Psi - \tilde{\Psi} \right\|_{\mathbf{C}^1} \cdot h \right) d\tau \leq L^{\Psi} \cdot \left\| \Psi - \tilde{\Psi} \right\|_{\mathbf{C}^1} \cdot t,$$

completing the proof. \square

Acknowledgments. We thank Graziano Guerra for helpful discussions. This work was started while the first author was visiting TU Kaiserslautern.

REFERENCES

[1] P. BAITI AND H. K. JENSSSEN, *On the front-tracking algorithm*, J. Math. Anal. Appl., 217 (1998), pp. 395–404.
 [2] M. K. BANDA, M. HERTY, AND A. KLAR, *Coupling conditions for gas networks governed by the isothermal Euler equations*, Netw. Heterog. Media, 1 (2006), pp. 275–294.
 [3] M. K. BANDA, M. HERTY, AND A. KLAR, *Gas flow in pipeline networks*, Netw. Heterog. Media, 1 (2006), pp. 41–56.
 [4] G. BASTIN, B. HAUT, J.-M. CORON, AND B. D’ANDRÉA-NOVEL, *Lyapunov stability analysis of networks of scalar conservation laws*, Netw. Heterog. Media, 2 (2007), pp. 751–759.
 [5] S. BIANCHINI AND R. M. COLOMBO, *On the stability of the standard Riemann semigroup*, Proc. Amer. Math. Soc., 130 (2002), pp. 1961–1973.
 [6] A. BRESSAN, *Hyperbolic Systems of Conservation Laws. The One-Dimensional Cauchy Problem*, Oxford Lecture Ser. Math. Appl. 20, Oxford University Press, Oxford, UK, 2000.
 [7] R. M. COLOMBO AND A. CORLI, *Sonic hyperbolic phase transitions and Chapman-Jouguet detonations*, J. Differential Equations, 184 (2002), pp. 321–347.
 [8] R. M. COLOMBO AND A. CORLI, *Stability of the Riemann semigroup with respect to the kinetic condition*, Quart. Appl. Math., 62 (2004), pp. 541–551.

- [9] R. M. COLOMBO AND M. GARAVELLO, *A well posed Riemann problem for the p -system at a junction*, *Netw. Heterog. Media*, 1 (2006), pp. 495–511.
- [10] R. M. COLOMBO AND M. GARAVELLO, *On the Cauchy problem for the p -system at a junction*, *SIAM J. Math. Anal.*, 39 (2008), pp. 1456–1471.
- [11] R. M. COLOMBO AND C. MAURI, *Euler system for compressible fluids at a junction*, *J. Hyperbolic Differ. Equ.*, to appear.
- [12] CRANE VALVE GROUP, *Flow of Fluids through Valves, Fittings and Pipes*, Technical report 410, Crane Technical paper, 1998.
- [13] A. J. C. B. DE SAINT-VENANT, *Theorie du mouvement non-permanent des eaux avec application aux crues des rivières et à l'introduction des marées dans leur lit*, *Comptes Rendus Academie des Sciences*, 73 (1871), pp. 148–154, 237–240.
- [14] K. EHRHARDT AND M. STEINBACH, *Nonlinear gas optimization in gas networks*, in *Modeling, Simulation and Optimization of Complex Processes*, H. G. Bock, E. Kostina, H. X. Pu, and R. Ranacher, eds., Springer, Berlin, 2005, pp. 139–148.
- [15] M. GUGAT, *Nodal control of conservation laws on networks*, in *Control and Boundary Analysis*, *Lect. Notes Pure Appl. Math.* 240, Chapman and Hall/CRC, Boca Raton, FL, 2005, pp. 201–215.
- [16] M. GUGAT AND G. LEUGERING, *Global boundary controllability of the de St. Venant equations between steady states*, *Ann. Inst. H. Poincaré Anal. Non Linéaire*, 20 (2003), pp. 1–11.
- [17] M. HERTY, *Coupling conditions for networked systems of Euler equations*, *SIAM J. Sci. Comput.*, 30 (2008), pp. 1596–1612.
- [18] M. HERTY AND M. SEAID, *Simulation of transient gas flow at pipe-to-pipe intersections*, *Internat. J. Numer. Methods Fluids*, 56 (2008), pp. 485–506.
- [19] H. HOLDEN AND N. H. RISEBRO, *Riemann problems with a kink*, *SIAM J. Math. Anal.*, 30 (1999), pp. 497–515.
- [20] G. LEUGERING AND E. J. P. G. SCHMIDT, *On the modelling and stabilization of flows in networks of open canals*, *SIAM J. Control Optim.*, 41 (2002), pp. 164–180.
- [21] MAPRESS. *Mapress Pressfitting System*, Technical report, Mapress GmbH & Co. KG, Langenfeld, Germany, <http://www.mapress.de> (2002).
- [22] A. MARTIN, M. MÖLLER, AND S. MORITZ, *Mixed integer models for the stationary case of gas network optimization*, *Math. Program.*, 105 (2006), pp. 563–582.
- [23] A. J. OSIADACZ, *Simulation and Analysis of Gas Networks*, Gulf Publishing, Houston, TX, 1989.
- [24] A. J. OSIADACZ, *Different transient models—limitations, advantages and disadvantages*, in *Proceedings of the 28th Annual Meeting of PSIG (Pipeline Simulation Interest Group)*, 1996.
- [25] A. OSIADACZ, *Simulation of transient gas flows in networks*, *Internat. J. Numer. Methods Fluids*, 4 (1984), pp. 13–24.
- [26] PIPELINE SIMULATION INTEREST GROUP, <http://www.psig.org/>.
- [27] J. ZHOU AND M. A. ADEWUMI, *Simulation of transients in natural gas pipelines using hybrid TVD schemes*, *Internat. J. Numer. Methods Fluids*, 32 (2000), pp. 407–437.

A VARIATIONAL PRINCIPLE FOR HARDENING ELASTOPLASTICITY*

ULISSE STEFANELLI†

Abstract. We present a variational principle governing the quasi-static evolution of a linearized elastoplastic material. In the case of linear hardening, the novel characterization allows us to recover and partly extend some known results and proves itself to be especially well suited for discussing general approximation and convergence issues. In particular, the variational principle is exploited in order to prove in a novel setting the convergence of time and space-time discretizations as well as to provide some possible a posteriori error control.

Key words. variational principle, elastoplasticity, approximation

AMS subject classifications. 35K55, 49S05, 74C05

DOI. 10.1137/070692571

1. Introduction. The primal initial-boundary value problem of elastoplasticity consists in determining the generalized deformation state of a material subject to external mechanical actions. In particular, starting from some initial state and for a given load and traction, one shall determine the displacement u of the body from the reference configuration, the inelastic (plastic) part p of its strain, and, possibly, a vector of internal hardening variables ξ . In the small deformation regime and within the frame of associative elastoplasticity, the problem is classically formulated in a variational form as that of finding the absolutely continuous trajectory $t \in [0, T] \mapsto y(t) \in Y$ (Y is a Banach space) such that

$$(1.1) \quad \partial\psi(\dot{y}) + Ay \ni \ell \quad \text{a.e. in } (0, T), \quad y(0) = y_0,$$

where $y = (u, p, \xi)$ stands for the vector of unknown fields, $A : Y \rightarrow Y^*$ (dual) is linear, continuous, and symmetric, and $\psi : Y \rightarrow [0, \infty]$ is the positively 1-homogeneous and convex dissipation potential (∂ is the classical subdifferential in the sense of convex analysis; see below). Moreover, $\ell : [0, T] \rightarrow Y^*$ is a given and suitably smooth generalized load (possibly including surface tractions), and y_0 represents the initial state. The reader is referred to section 2 for some brief mechanical motivation as well as to the classical monographs by Duvaut and Lions [7], Han and Reddy [13], Lemaitre and Chaboche [18], and Simo and Hughes [42] for a comprehensive collection of results.

The aim of this paper is that of investigating a global-in-time variational formulation of problem (1.1). In particular, we shall introduce the functional $\mathcal{F} : W^{1,1}(0, T; Y) \rightarrow [0, \infty]$ on trajectories as

$$\mathcal{F}(y) = \int_0^T (\psi(\dot{y}) + \psi^*(\ell - Ay) - \langle \ell - Ay, \dot{y} \rangle),$$

*Received by the editors May 21, 2007; accepted for publication (in revised form) January 17, 2008; published electronically July 3, 2008. This research was performed during a visit to the Seminar for Applied Mathematics at ETH Zürich and the Institute of Mathematics of the University of Zürich under the sponsorship of the STM CNR 2006 program and the Swiss National Science Foundation. The support and hospitality of both institutions is gratefully acknowledged.

<http://www.siam.org/journals/sima/40-2/69257.html>

†IMATI - CNR, v. Ferrata 1, I-27100, Pavia, Italy (ulisse.stefanelli@imati.cnr.it).

where ψ^* stands for the conjugate $\psi^*(w) = \sup_{v \in Y} (\langle w, v \rangle - \psi(v))$ of ψ and $\langle \cdot, \cdot \rangle$ denotes the duality pairing between Y^* and Y . The starting point of this analysis relies on the fact that solutions of (1.1) and minimizers of \mathcal{F} fulfilling the initial condition and making \mathcal{F} zero coincide, namely (see Theorem 3.1),

$$(1.2) \quad y \text{ solves (1.1) iff } \mathcal{F}(y) = \min \mathcal{F} = 0 \text{ and } y(0) = y_0.$$

This variational characterization has a clear mechanical interpretation. Indeed, since ψ is positively 1-homogeneous, its conjugate ψ^* turns out to be the indicator function of the convex set $\partial\psi(0)$. Hence, $\min \mathcal{F}$ is actually a constrained minimization problem, and we have

$$(1.3) \quad \mathcal{F}(y) = 0 \quad \text{iff} \quad \begin{cases} \ell - Ay \in \partial\psi(0) \quad \text{a.e. in } (0, T), \\ \varphi(T, y(T)) + \int_0^T \psi(\dot{y}) = \varphi(0, y(0)) - \int_0^T \langle \dot{\ell}, y \rangle, \end{cases}$$

where we have used the notation $(t, y) \mapsto \varphi(t, y) = \frac{1}{2} \langle Ay, y \rangle - \langle \ell(t), y \rangle$. The first relation above expresses the so-called *local stability* [26] of the trajectory, whereas the second is nothing but the energy balance at time T . More precisely, $\varphi(t, y)$ denotes the complementary energy at time t for the state y , $\int_0^T \psi(\dot{y})$ represents the dissipation of the system on $[0, T]$, and $-\int_0^T \langle \dot{\ell}, y \rangle$ is related to external actions on $[0, T]$. Hence, minimizing \mathcal{F} consists in selecting the (only) stable trajectory which conserves the energy. In this regards, the reader is referred to the pioneering papers by Moreau [33, 34, 35].

The interest of variational characterization (1.2) of the differential problem (1.1) relies on the possibility of exploiting the general tools from the calculus of variations. Some care is, however, required. Indeed, although \mathcal{F} is convex and lower semicontinuous with respect to the weak topology of $W^{1,1}(0, T; Y)$, the functional generally fails to be coercive. Moreover, one is not just asked to minimize \mathcal{F} but also to prove that the minimum is 0. These considerations suggest that the direct method is hardly applicable in order to get solutions to (1.1) via the characterization in (1.2).

The first issue of this paper is that of exploiting the variational principle in (1.2) in order to address general approximation procedures. Since solutions and minimizers coincide, a natural tool in order to frame an abstract approach to limiting procedures within (1.1) is that of considering the corresponding minimum problems via Γ -convergence [11]. As the value of the functional is directly quantified to be 0 on the minimizers, what is actually needed here for passing to limits are so-called Γ -lim inf inequalities only, and the latter are generally easily available. We shall specifically focus on the case of linear hardening elastoplasticity and apply the above-mentioned perspective in order to recover in a unified and more transparent frame and partly generalize some convergence results for conformal finite elements (Theorem 5.3), time discretizations (Theorem 6.5), and fully discrete space-time approximations (Theorem 7.1). In particular, for time discretization we develop a discrete version of the variational principle (1.2) in the same spirit of the theory of *variational integrators* [25] (see subsection 6.1). This connection entails also some generalized view at the classical discrete time schemes (see subsection 6.5).

A second novel point of the present variational approach consists in the possibility of exploiting \mathcal{F} in order to estimate a posteriori some approximation error. By letting $\mathcal{F}(y) = 0$, we will check that (Corollary 4.5)

$$\max_{[0, T]} \frac{1}{2} \langle A(y - v), y - v \rangle \leq \mathcal{F}(v) \quad \forall v \in W^{1,1}(0, T; Y), \quad v(0) = y_0.$$

If A shows some coercivity (which is precisely the case of linearized hardening; see subsection 2.4), and v is the outcome of some approximation procedure, the estimate above may serve as the basis for some a posteriori estimation procedure, possibly headed to adaptivity (see subsection 6.7). Let us stress that the latter and (1.3) entail that the distance of a (stable) trajectory from the solution to (1.1) can be uniformly estimated by means of its energy production along the path.

The variational characterization in (1.2) is strictly linked to the celebrated principle by Brezis, Ekeland, and Nayroles [3, 4, 38, 39] for the gradient flow of a convex functional $j : Y \rightarrow (-\infty, \infty]$ in a Hilbert space Y , namely,

$$(1.4) \quad \dot{y} + \partial j(y) \ni \ell \quad \text{a.e. in } (0, T).$$

In particular, by using the convexity of j (see details in section 3.1), this inclusion may be equivalently rewritten as the scalar relation

$$(1.5) \quad j(y) + j^*(\ell - \dot{y}) - \langle \ell - \dot{y}, y \rangle = 0 \quad \text{a.e. in } (0, T),$$

where j^* is the conjugate of j . As the above left-hand side is a.e. nonnegative for all trajectories y and it is 0 iff y solves (1.4), the latter is equivalent to minimize the global functional

$$\mathcal{I}(y) = \int_0^T \left(j(y) + j^*(\ell - \dot{y}) - \langle \ell - \dot{y}, y \rangle \right)$$

and check that $\mathcal{I}(y) = 0$. Since its introduction, the latter principle has continuously attracted attention. In particular, it has been exploited in the direction of proving the existence [40, 2, 41, 10, 9] (note that the above-mentioned obstructions to the application of the direct methods again appear) and the description of long-time dynamics [17]. Moreover, the Brezis–Ekeland–Nayroles approach has been adapted to the case of second order [21, 22] and doubly nonlinear equations [45] and to the study of variable time step discretizations of (1.4) [43] as well.

The variational characterization (1.2) stems basically from the same idea as in (1.5). Namely, the differential inclusion in (1.1) is equivalently rewritten by convexity of ψ as

$$\psi(\dot{y}) + \psi^*(\ell - Ay) - \langle \ell - Ay, \dot{y} \rangle = 0 \quad \text{a.e. in } (0, T),$$

and (1.2) follows by noting that the above left-hand side is a.e. nonnegative for all trajectories y and it is 0 iff y is fulfilling the inclusion in (1.1). On the other hand, apart from this conceptual analogy, one has indeed to mention that the characterization (1.2) has little in common with the original Brezis–Ekeland–Nayroles principle. In particular, the two principles turn out to be different even in the case of a quadratic functional ψ on a Hilbert space.

One has to mention that, of course, (1.2) is not the only possible global-in-time variational characterization of (1.1). Besides minimizing the L^2 space-time norm of the residual (which might be a little interesting since the order of the problem is doubled), one has at least to mention Visintin [48], where generalized solutions are obtained as minimal elements of a certain partial-order relation on the trajectories, and the recent contribution by Mielke and Ortiz [27], where the functional

$$(1.6) \quad y \mapsto e^{-T/\varepsilon} \varphi(T, y(T)) + \int_0^T e^{-t/\varepsilon} \left(\psi(\dot{y}) + \frac{1}{\varepsilon} \varphi(t, y) \right)$$

is minimized among trajectories with $y(0) = y_0$. Under extra-smoothness conditions on ψ (not fulfilled in the current frame), the Euler–Lagrange equations of the latter functional are

$$\begin{aligned} -\varepsilon D^2\psi(\dot{y})\ddot{y} + D\psi(\dot{y}) + Ay &= \ell, \\ y(0) = y_0, \quad D\psi(\dot{y}(T)) + Ay(T) &= \ell(T). \end{aligned}$$

In particular, minimizing the functional in (1.6) consists in performing a suitable elliptic (in time) regularization of the problem. In the specific case of ψ positively 1-homogeneous, the limit $\varepsilon \rightarrow 0$ can be carried out, and the minimizers of the functional in (1.6) are proved to converge to the solution of (1.1). The latter approach is quite different from that of (1.2). On the one hand, it is much more general as it naturally applies to the nonsmooth case as well (no derivatives of ψ and ϕ are involved). The results of [27] have been recently extended in the direction of discretizations and relaxation in [31].

2. Mechanical model. Let us provide the reader with a brief introduction to the mechanical setting under consideration. Our aim is just that of recalling some essential features of the models as well as their variational formulation. In particular, we restrain from reporting here an extensive discussion on associative elastoplasticity as the latter can be easily recovered from the many contributions on the subject. The reader is particularly referred to the mentioned monographs for some comprehensive presentation.

2.1. Preliminaries. We will denote by $\mathbb{R}_{\text{sym}}^{3 \times 3}$ the space of symmetric 3×3 tensors endowed with the natural scalar product $a : b := \text{tr}(ab) = a_{ij}b_{ij}$ (summation convention). The space $\mathbb{R}_{\text{sym}}^{3 \times 3}$ is orthogonally decomposed as $\mathbb{R}_{\text{sym}}^{3 \times 3} = \mathbb{R}_{\text{dev}}^{3 \times 3} \oplus \mathbb{R}1_2$, where $\mathbb{R}1_2$ is the subspace spanned by the identity 2-tensor 1_2 and $\mathbb{R}_{\text{dev}}^{3 \times 3}$ is the subspace of deviatoric symmetric 3×3 tensors. In particular, for all $a \in \mathbb{R}_{\text{sym}}^{3 \times 3}$, we have that $a = a_{\text{dev}} + \text{tr}(a)1_2/3$.

We shall assume the reference configuration Ω to be a nonempty, bounded, and connected open set in \mathbb{R}^3 with a Lipschitz continuous boundary. The space dimension 3 plays essentially no role throughout the analysis, and we would be in the position of reformulating our results in \mathbb{R}^d with no particular intricacy. Our unknown variables are the displacement of the body $u \in \mathbb{R}^3$, the plastic strain $p \in \mathbb{R}_{\text{dev}}^{3 \times 3}$, and a vector of internal variables $\xi \in \mathbb{R}^m$ ($m \in \mathbb{N}$) which will describe the hardening of the material. We will denote by $\varepsilon(u)$ the standard symmetric gradient.

2.2. Constitutive relation. Moving within the small-strain regime, we additively decompose the linearized deformation $\varepsilon(u)$ into the elastic strain e and the inelastic (or plastic) strain p as

$$\varepsilon(u) = e + p.$$

Let \mathbb{C} be the elasticity tensor. By regarding the latter as a symmetric positive definite linear map $\mathbb{C} : \mathbb{R}_{\text{sym}}^{3 \times 3} \rightarrow \mathbb{R}_{\text{sym}}^{3 \times 3}$, we shall assume that the orthogonal subspaces $\mathbb{R}_{\text{dev}}^{3 \times 3}$ and $\mathbb{R}1_2$ are invariant under \mathbb{C} . This amounts to saying that indeed

$$\mathbb{C}a = \mathbb{C}_{\text{dev}}a_{\text{dev}} + \kappa \text{tr}(a)1_2$$

for a given $\mathbb{C}_{\text{dev}} : \mathbb{R}_{\text{dev}}^{3 \times 3} \rightarrow \mathbb{R}_{\text{dev}}^{3 \times 3}$, a constant κ , and all $a \in \mathbb{R}_{\text{sym}}^{3 \times 3}$. The case of isotropic materials is given by $\mathbb{C}_{\text{dev}} = 2G(1_4 - 1_2 \otimes 1_2/3)$, and G and κ are, respectively, the

shear and the bulk moduli. The latter decomposition is not exploited in our analysis, but it is clearly suggested by the mechanical application. Moreover, we shall introduce two linear symmetric positive semidefinite hardening moduli $\mathbb{H}_p : \mathbb{R}_{\text{dev}}^{3 \times 3} \rightarrow \mathbb{R}_{\text{dev}}^{3 \times 3}$ and $\mathbb{H}_\xi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ (to be identified with a fourth order tensor and a matrix, respectively) and define the Helmholtz free energy $W : \mathbb{R}_{\text{sym}}^{3 \times 3} \times \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m \rightarrow [0, \infty)$ of the material as

$$W(\varepsilon(u), p, \xi) := \frac{1}{2}(\varepsilon(u) - p) : \mathbb{C}(\varepsilon(u) - p) + \frac{1}{2}p : \mathbb{H}_p p + \frac{1}{2}\xi^T \cdot \mathbb{H}_\xi \xi.$$

The generalized stresses (σ, η) are conjugate to the above-defined generalized strains (e, ξ) via the energy W . In particular, the material is classically assumed to show elastic response,

$$(2.1) \quad \sigma = \frac{\partial W}{\partial e} = \mathbb{C}e = \mathbb{C}(\varepsilon(u) - p),$$

and the thermodynamic force η driving the evolution of the internal variables ξ is defined as

$$(2.2) \quad \eta = -\frac{\partial W}{\partial \xi} = -\mathbb{H}_\xi \xi.$$

Moreover, by moving within the frame of associative elastoplasticity, we assume the existence of a function $R : \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m \rightarrow [0, \infty]$ convex, positively 1-homogeneous, and lower semicontinuous such that

$$(2.3) \quad \partial R(\dot{p}, \dot{\xi}) \ni \begin{pmatrix} \sigma - \mathbb{H}_p p \\ \eta \end{pmatrix}.$$

In particular, R is asked to be the support function of a convex set $C^* \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m$, i.e., $R(p) = \sup_{q \in C^*} q : p$. We will indicate with R^* its conjugate, namely, the indicator function of C^* given by $R^*(q) = 0$ if $q \in C^*$ and $R^*(q) = \infty$ otherwise. Moreover, we let C be the domain of R , namely, $C = D(R) = \{(p, \xi) \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m : R(p, \xi) < \infty\}$.

Finally, the above material relations (2.1)–(2.3) can be condensed as the following constitutive material law:

$$(2.4) \quad \partial R(\dot{p}, \dot{\xi}) + \begin{pmatrix} (\mathbb{C} + \mathbb{H}_p)p \\ \mathbb{H}_\xi \xi \end{pmatrix} \ni \begin{pmatrix} \mathbb{C}\varepsilon(u) \\ 0 \end{pmatrix},$$

which in turn can be rephrased in the form of (1.1) by letting

$$(2.5) \quad y = (p, \xi), \quad Y = \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m, \quad \psi = R, \\ A(p, \xi) = ((\mathbb{C} + \mathbb{H}_p)p, \mathbb{H}_\xi \xi), \quad \ell = (\mathbb{C}\varepsilon(u), 0).$$

Let us close this subsection by explicitly mentioning three classical linear hardening models [13, Ex. 4.8, p. 88]:

Linear kinematic hardening. Choose $\mathbb{H}_p = h_p 1_4$, where $h_p > 0$, and $\mathbb{H}_\xi = 0$. In this case the internal variable ξ is not evolving and shall be removed from the set of unknowns.

Linear isotropic hardening. Choose $\mathbb{H}_p = 0$, $m = 1$, and $\mathbb{H}_\xi = h_\xi > 0$. Moreover, let $D(R) = \{(p, \xi) \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R} : |p| \leq \xi\}$.

Linear combined kinematic-isotropic hardening. Let $\mathbb{H}_p = h_p 1_4$, $m = 1$, and $\mathbb{H}_\xi = h_\xi$, where $h_p, h_\xi > 0$. Moreover, let $D(R) = \{(p, \xi) \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R} : |p| \leq \xi\}$.

It is beyond the purpose of this introduction to discuss and justify the above-mentioned material models. The reader should check the cited references for comments on their relevance within applications and some mechanical motivation.

2.3. Variational formulation of the quasi-static evolution. Let us now move to the consideration of the full equilibrium problem. To this aim, we assume that the boundary $\partial\Omega$ is partitioned in two disjoint open sets Γ_{tr} and Γ_{Dir} , with $\partial\Gamma_{\text{tr}} = \partial\Gamma_{\text{Dir}}$ (in $\partial\Omega$). We ask Γ_{Dir} to be such that there exists a positive constant c_{Korn} depending on Γ_{Dir} and Ω such that the Korn inequality

$$(2.6) \quad c_{\text{Korn}} \|u\|_{H^1(\Omega; \mathbb{R}^3)}^2 \leq \|u\|_{L^2(\Gamma_{\text{Dir}}; \mathbb{R}^3)}^2 + \|\varepsilon(u)\|_{L^2(\Omega; \mathbb{R}^{3 \times 3}_{\text{sym}})}^2$$

holds true for all $u \in H^1(\Omega; \mathbb{R}^3)$. It would indeed suffice to impose Γ_{Dir} to have a positive surface measure (see, e.g., [7, Thm. 3.1, p. 110]).

For the sake of simplicity, we will prescribe homogeneous Dirichlet boundary conditions on Γ_{Dir} (our analysis extends with little notational intricacy to the case of non-homogeneous Dirichlet boundary conditions as well). On Γ_{tr} some time-dependent traction will be prescribed instead.

As for the full quasi-static evolution of the material, we shall couple the constitutive relation (2.4) with the equilibrium equation

$$(2.7) \quad \operatorname{div} \sigma + f = 0 \quad \text{in } \Omega.$$

Here we assume to be given the body force $f : [0, T] \rightarrow L^2(\Omega; \mathbb{R}^3)$ and a surface traction $g : [0, T] \rightarrow L^2(\Gamma_{\text{tr}}; \mathbb{R}^3)$.

Then one can rephrase the problem into the form of (1.1) by choosing

$$(2.8) \quad y = (u, p, \xi),$$

$$Y = \left\{ (u, p, \xi) \in H^1(\Omega; \mathbb{R}^3) \times L^2(\Omega; \mathbb{R}^{3 \times 3}_{\text{dev}}) \times L^2(\Omega; \mathbb{R}^m) \right.$$

$$(2.9) \quad \left. \text{such that } u = 0 \text{ on } \Gamma_{\text{Dir}} \right\},$$

$$\langle A(u, p, \xi), (v, q, z) \rangle = \int_{\Omega} \left((\varepsilon(u) - p) : \mathbb{C}(\varepsilon(v) - q) + p : \mathbb{H}_p q + \xi^T \cdot \mathbb{H}_{\xi} z \right)$$

$$(2.10) \quad \forall (v, q, z) \in Y,$$

$$(2.11) \quad \psi(u, p, \xi) = \int_{\Omega} R(p, \xi)$$

and defining the total load $\ell : [0, T] \rightarrow Y^*$ as

$$\langle \ell(t), (u, p, \xi) \rangle = \int_{\Omega} f \cdot u + \int_{\Gamma_{\text{tr}}} g \cdot u \, d\mathcal{H}^2 \quad \forall u \in H^1(\Omega; \mathbb{R}^3), \quad t \in [0, T],$$

where \mathcal{H}^2 is the 2-dimensional Hausdorff measure.

2.4. The coercivity of A . Let us close this introductory discussion by explicitly commenting on the coercivity of the bilinear form induced by A . We shall recall some sufficient conditions on \mathbb{H}_p , \mathbb{H}_{ξ} , and R in such a way that there exists a constant $\alpha > 0$ such that

$$(2.12) \quad \langle Ay, y \rangle \geq \alpha |y|^2 \quad \forall y \in D(\psi),$$

where $|\cdot|$ is the norm in Y . This issue is fairly classical [13, sec. 7.3, p. 167], and we discuss it here for the sake of completeness only.

Of course, (2.12) holds (and even for all $y \in Y$) whenever \mathbb{H}_p and \mathbb{H}_{ξ} are positive definite (this is the case of the above-mentioned linear combined kinematic-isotropic hardening).

As we have already observed, in the case $\mathbb{H}_\xi = 0$, the problem naturally reduces to the pair (u, p) only. Up to this reduction, (2.12) holds (again for all $y \in Y$) when \mathbb{H}_p is positive definite. This is exactly the case of linear kinematic hardening.

On the other hand, in the case $\mathbb{H}_p = 0$, the plastic strain will still evolve, and one has (2.12) if $D(R)$ is bounded in the p -direction for all ξ , namely, if [13, eq. (7.51)]

$$(2.13) \quad D(R) \subset \{(p, \xi) \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m : \beta |p|^2 \leq \xi^T \cdot \mathbb{H}_\xi \xi \text{ for some constant } \beta > 0\},$$

which is clearly the case for linear isotropic hardening.

Some generalization of the latter condition could in principle be considered for the case when \mathbb{H}_p and \mathbb{H}_ξ are only semidefinite. In particular, (2.12) holds if one assumes (2.13) and

$$\xi \neq 0 \text{ and } \xi^T \cdot \mathbb{H}_\xi \xi = 0 \Rightarrow R(p, \xi) = \infty \quad \forall p \in \mathbb{R}_{\text{dev}}^{3 \times 3}.$$

Let us mention that the most critical case in the class of (2.4) is $\mathbb{H}_p = 0, \mathbb{H}_\xi = 0$ where actually no hardening takes place. This is the situation of *perfect plasticity* for which the Sobolev space framework above is not appropriate, and one would consider the space $BD(\Omega)$ of functions of bounded deformations instead [6]. We shall make clear that, even if our variational characterization covers the case of perfect plasticity, the subsequent approximation results apply to the linear hardening situation only.

3. Characterization.

3.1. General assumptions. Let us start by recalling notation and enlisting the basic assumptions for the following. First of all, we will ask that

$$(3.1) \quad Y \text{ is a separable and reflexive Banach space.}$$

We will use the symbols $|\cdot|$ for the norm of Y and $\langle \cdot, \cdot \rangle$ for the duality pairing between Y^* (dual) and Y . The norm in Y^* will be denoted by $|\cdot|_*$ instead.

We introduce the functional

$$(3.2) \quad \begin{aligned} \psi : Y \rightarrow [0, \infty] & \text{ proper, convex, lower semicontinuous,} \\ & \text{and positively 1-homogeneous.} \end{aligned}$$

Equivalently, ψ is required to be the support function of a convex and closed set $C^* \subset Y^*$ containing 0, namely,

$$(3.3) \quad \psi(y) = \sup\{\langle y^*, y \rangle : y^* \in C^*\}.$$

We shall define $C = D(\psi)$. Hence, the conjugate $\psi^* : Y^* \rightarrow [0, \infty]$, which is classically defined as $\psi^*(y^*) = \sup_{y \in Y} (\langle y^*, y \rangle - \psi(y))$, is the indicator function of the convex set C^* , namely, $\psi^*(y^*) = 0$ if $y^* \in C^*$ and $\psi^*(y^*) = \infty$ otherwise. Let us remark that ψ fulfills the triangle inequality $\psi(a) \leq \psi(b) + \psi(c)$ whenever $a = b + c$.

We shall use the symbol ∂ in order to denote the usual subdifferential in the sense of convex analysis, namely,

$$y^* \in \partial\psi(y) \text{ iff } y \in D(\psi) \text{ and } \langle y^*, w - y \rangle \leq \psi(w) - \psi(y) \quad \forall w \in Y.$$

Similarly, we define

$$\begin{aligned} y \in \partial\psi^*(y^*) & \text{ iff } y^* \in D(\psi^*) \text{ and } \langle w^* - y^*, y \rangle \leq \psi^*(w^*) - \psi^*(y^*) \quad \forall w^* \in Y^* \\ & \text{ iff } y^* \in C^* \text{ and } \langle w^* - y^*, y \rangle \leq 0 \quad \forall w^* \in C^*. \end{aligned}$$

Finally, we recall Fenchel's inequality

$$\psi(y) + \psi^*(y^*) \geq \langle y^*, y \rangle \quad \forall y \in Y, y^* \in Y^*,$$

and remark that equality holds iff $y^* \in \partial\psi(y)$ (or, equivalently, $y \in \partial^*\psi^*(y^*)$).

As for the operator A we require that

$$(3.4) \quad A : Y \rightarrow Y^* \quad \text{linear, continuous, and symmetric}$$

and define the function

$$y \rightarrow \phi(y) = \frac{1}{2} \langle Ay, y \rangle,$$

so that $A = D\phi$. Moreover, we will ask ϕ to be coercive on $C = D(\psi)$; namely, we assume that there exists a positive constant α such that

$$(3.5) \quad \phi(y) \geq \frac{\alpha}{2} |y|^2 \quad \forall y \in C.$$

As we have already commented in subsection 2.4, the latter coercivity is fulfilled in the situation of elastoplastic evolution with linear kinematic, isotropic, or combined kinematic-isotropic hardening and will turn out to be sufficient for both of the forthcoming characterization results.

On the other hand, the following uniqueness-type results will be checked under some stronger coercivity frame, and we will ask for

$$(3.6) \quad \phi(y) \geq \frac{\alpha}{2} |y|^2 \quad \forall y \in C - C.$$

Clearly, condition (3.6) is fulfilled when ϕ happens to be coercive on the whole space Y . The latter applies in particular to the case of linear kinematic and combined kinematic-isotropic hardening elastoplasticity. In this case, ϕ defines an equivalent (squared) norm in Y .

We shall make use of the following notation:

$$\chi(y) = \phi(y) + |y|^2 \quad \forall y \in Y.$$

Indeed the latter choice is just motivated by simplicity and could be replaced as well by any other $\chi : Y \rightarrow [0, \infty)$ such that $\chi(y) = 0$ iff $y = 0$ and that $y \mapsto \chi(y) - \phi(y)$ is lower semicontinuous.

Finally, we shall fix data such that

$$(3.7) \quad \ell \in L^\infty(0, T; Y^*), \quad y_0 \in C.$$

The restriction on the choice of the initial datum in C is motivated by the coercivity assumption on ϕ in (3.5). On the other hand, we shall explicitly mention that the usual choice for y_0 in elastoplasticity is $y_0 = 0$.

In what follows, the above assumptions (3.1)–(3.5) and (3.7) will be tacitly assumed (unless explicitly stated). It should be clear, however, that the above choice is motivated by the sake of simplicity. Indeed, most of the following results still hold under suitably weaker assumptions, as we shall comment.

3.2. The functional. Let the *Lagrangian* $L : (0, T) \times Y \times Y \rightarrow [0, \infty]$ be defined as

$$(3.8) \quad \begin{aligned} L(t, y, p) &= \psi(p) + \psi^*(\ell(t) - Ay) - \langle \ell(t) - Ay, p \rangle \\ &\text{for a.e. } t \in (0, T), \forall y, p \in Y, \end{aligned}$$

and the functional $F : W^{1,1}(0, T; Y) \rightarrow [0, \infty]$ as

$$(3.9) \quad F(y) = \int_0^T L(t, y(t), \dot{y}(t)) \, dt + \chi(y(0) - y_0).$$

Now, by simply using the chain rule, we obtain

$$F(y) = \int_0^T \left(\psi(\dot{y}) + \psi^*(\ell - Ay) - \langle \ell, \dot{y} \rangle \right) + \phi(y(T)) - \phi(y(0)) + \chi(y(0) - y_0).$$

A first remark is that, by exploiting the particular form of χ ,

$$(3.10) \quad \begin{aligned} F(y) &= \int_0^T \left(\psi(\dot{y}) + \psi^*(\ell - Ay) - \langle \ell, \dot{y} \rangle \right) \\ &\quad + \phi(y(T)) + \phi(y_0) - \langle Ay(0), y_0 \rangle + |y(0) - y_0|^2. \end{aligned}$$

In particular, F is clearly convex.

3.3. The characterization. Let us state here our variational principle.

THEOREM 3.1 (variational principle). $y \in W^{1,1}(0, T; Y)$ solves (1.1) iff

$$F(y) = 0 = \min F.$$

Proof. Owing to Fenchel's inequality we have

$$L(t, y, p) = 0 \quad \text{iff } \ell(t) - Ay \in \partial\psi(p)$$

and, clearly, $\chi(y(0) - y_0) = 0$ iff $y(0) = y_0$. Hence, all solutions y of (1.1) are such that $F(y) = 0$ and vice versa. \square

Let us remark that the latter variational characterization result holds in much greater generality. The proof made no use of the separability and reflexivity of Y nor of the linearity of A (besides its being single-valued and such that $t \mapsto Ay(t)$ is measurable). Moreover, the positive 1-homogeneity of ψ is unessential [44]. In particular, the variational approach of Theorem 3.1 can be directly extended to a variety of different dissipative systems possibly including viscous evolution as well. We shall address this perspective in a forthcoming contribution.

We have already observed that F is convex. Moreover, F is lower semicontinuous with respect to the weak topology of $W^{1,1}(0, T; Y)$, since all weakly convergent sequences in $W^{1,1}(0, T; Y)$ are pointwise weakly convergent as well. Hence, one could be tempted to use the direct method in order to get the existence of minimizers, i.e., solutions to (1.1). As we commented in the introduction, this seems to be no trivial task.

First of all, the functional F need not be coercive with respect to the weak topology of $W^{1,1}(0, T; Y)$. Indeed, the functional ψ may degenerate and hence not control the norm of its argument. Moreover, even in the case when ψ is nondegenerate, the homogeneity assumption just entails that the sublevels of F are bounded in $W^{1,1}(0, T; Y)$ and no weak compactness follows.

Second, even assuming coercivity in the weak topology of $W^{1,1}(0, T; Y)$, one would still need to prove that the minimum 0 is attained.

3.4. The variational principle for hardening elastoplasticity. By referring to the notations of section 2, let us now present the actual form of the functional F for the case of the constitutive relation for linearized elastoplastic materials with linear hardening (see (2.5)). In this case the functional reads

$$\begin{aligned} F(p, \xi) &= \int_0^T \left(R(\dot{p}, \dot{\xi}) + R^*(\mathbb{C}(\varepsilon(u) - p) - \mathbb{H}_p p, -\mathbb{H}_\xi \xi) \right) \\ &\quad - \int_0^T \left((\mathbb{C}(\varepsilon(u) - p) - \mathbb{H}_p p) : \dot{p} - \xi^T \cdot \mathbb{H}_\xi \dot{\xi} \right) \\ &\quad + \frac{1}{2} (p(0) - p_0) : (\mathbb{C} + \mathbb{H}_p)(p(0) - p_0) + \frac{1}{2} (\xi(0) - \xi_0)^T \cdot \mathbb{H}_\xi (\xi(0) - \xi_0) \\ &\quad + |(p(0), \xi(0)) - (p_0, \xi_0)|^2 \end{aligned}$$

for some given initial datum $(p_0, \xi_0) \in \mathbb{R}_{\text{dev}}^{3 \times 3} \times \mathbb{R}^m$ and $\varepsilon(u) \in L^\infty(0, T; \mathbb{R}_{\text{sym}}^{3 \times 3})$.

In the situation of the quasi-static evolution, for some given initial datum $(u_0, p_0, \xi_0) \in Y$, a load $f \in L^\infty(0, T; L^2(\Omega; \mathbb{R}^3))$, and a traction $g \in L^\infty(0, T; L^2(\Gamma_{\text{tr}}; \mathbb{R}^3))$, the functional reads (see (2.8)–(2.11))

$$\begin{aligned} F(u, p, \xi) &= \int_0^T \int_\Omega \left(R(\dot{p}, \dot{\xi}) + R^*(\mathbb{C}(\varepsilon(u) - p) - \mathbb{H}_p p, -\mathbb{H}_\xi \xi) \right) \\ &\quad - \int_0^T \int_\Omega f \cdot \dot{u} - \int_0^T \int_{\Gamma_{\text{tr}}} g \cdot \dot{u} \, d\mathcal{H}^2 \\ &\quad + \int_0^T \int_\Omega \left((\varepsilon(u) - p) : \mathbb{C}(\varepsilon(\dot{u}) - \dot{p}) + p : \mathbb{H}_p \dot{p} + \xi^T \cdot \mathbb{H}_\xi \dot{\xi} \right) \\ &\quad + \frac{1}{2} \int_\Omega (\varepsilon(u(0) - u_0) - (p(0) - p_0)) : \mathbb{C}(\varepsilon(u(0) - u_0) - (p(0) - p_0)) \\ &\quad + \frac{1}{2} \int_\Omega \left((p(0) - p_0) : \mathbb{H}_p (p(0) - p_0) + (\xi(0) - \xi_0)^T \cdot \mathbb{H}_\xi (\xi(0) - \xi_0) \right) \\ &\quad + \frac{1}{2} \int_\Omega |(u(0), p(0), \xi(0)) - (u_0, p_0, \xi_0)|^2 \end{aligned}$$

for all points $(u, p, \xi) \in Y$ such that

$$\begin{aligned} \int_\Omega (\varepsilon(u) - p) : \mathbb{C} \varepsilon(v) &= \int_\Omega f \cdot v - \int_{\Gamma_{\text{tr}}} g \cdot v \, d\mathcal{H}^2 \\ \forall v \in H^1(\Omega; \mathbb{R}^3), \quad \text{with } v &= 0 \text{ on } \Gamma_{\text{Dir}}, \quad \text{a.e. in } (0, T), \end{aligned}$$

and $F(u, p, \xi) = \infty$ otherwise.

4. Properties of the minimizers. For the sake of illustrating the variational principle of Theorem 3.1, we shall collect here some properties of the trajectories belonging to the domain of the functional F and, in particular, of the minimizers.

4.1. Trajectories are in C .

LEMMA 4.1. *Let $F(y) < \infty$. Then $y(t) \in C$ for all $t \in [0, T]$.*

Proof. Since $F(y) < \infty$ we have $\dot{y} \in C$ a.e. in $(0, T)$. Hence, for all $t \in [0, T]$, we have $\int_0^t \dot{y} \in C$ by Jensen's inequality. On the other hand, $y_0 \in C$ and $y(t) = y_0 + \int_0^t \dot{y}$. The assertion follows by recalling that C is a cone. \square

4.2. Stability at regular points of ℓ . Assume that $\ell : [0, T] \rightarrow Y^*$ is given, and let the set of *stable* states $S(t) \subset Y$ for $t \in [0, T]$ be defined as

$$S(t) = \{y \in Y : \phi(y) - \langle \ell(t), y \rangle \leq \phi(w) - \langle \ell(t), w \rangle + \psi(w - y) \quad \forall w \in Y\}.$$

In particular, one has $y \in S(t)$ iff it minimizes the convex functional $G : w \mapsto \phi(w) - \langle \ell(t), w \rangle + \psi(w - y)$, namely,

$$y \in S(t) \quad \text{iff} \quad 0 \in \partial G(y) = Ay - \ell(t) + \partial\psi(0).$$

Let us now check that indeed $\partial\psi(0) = C^*$. To this end, owing to (3.3), given $y^* \in C^*$ one has $\langle y^*, w \rangle \leq \psi(w)$ for all $w \in Y$. Namely, $C^* \subset \partial\psi(0)$. On the other hand, assume that there exists $\xi \in \partial\psi(0) \setminus C^*$. Hence, by the Hahn–Banach theorem one finds $y \in Y$, $\alpha \in \mathbb{R}$, and $\varepsilon > 0$ such that

$$\psi(y) \stackrel{\xi \in \partial\psi(0)}{\geq} \langle \xi, y \rangle \geq \alpha + \varepsilon > \alpha - \varepsilon \geq \langle y^*, y \rangle \quad \forall y^* \in C^*,$$

and, by passing to the supremum as $y^* \in C^*$, we obtain a contradiction. We may summarize this discussion as follows:

$$y \in S(t) \quad \text{iff} \quad \ell(t) - Ay \in C^*.$$

LEMMA 4.2 (stability of the minimizers). *Let ℓ be either left- or right-weakly continuous at some point $t \in [0, T]$ and $F(y) < \infty$. Then $y(t) \in S(t)$.*

Proof. Since $F(y) < \infty$, we have $\ell(t) - Ay(t) \in C^*$ for all $t \in (0, T) \setminus N$, where $|N| = 0$. Choose a sequence $t_k \in (0, T) \setminus N$ such that $t_k \rightarrow t$ (from the left or from the right) and $\ell(t_k) \rightarrow \ell(t)$ weakly in Y^* . Hence, $\ell(t_k) - Ay(t_k) \rightarrow \ell(t) - Ay(t)$ weakly in Y^* and $\ell(t) - Ay(t) \in C^*$. \square

In particular, if ℓ happens to be right-continuous at 0, the functional F will not attain the minimum value 0 unless the initial datum y_0 is stable, namely, $y_0 \in S(0)$.

4.3. Equivalent formulations. Letting now $\ell \in W^{1,1}(0, T; Y^*)$, problem (1.1) admits some alternative equivalent formulations [26, sec. 2.1]. We explicitly mention that $y \in W^{1,1}(0, T; Y)$ is said to be an *energetic solution* if it solves the *energetic formulation* [32] of (1.1), namely,

$$(4.1) \quad y(t) \in S(t) \quad \forall t \in [0, T],$$

$$(4.2) \quad \phi(y(t)) - \langle \ell(t), y(t) \rangle + \int_0^t \psi(\dot{y}) = \phi(y(0)) - \langle \ell(0), y(0) \rangle - \int_0^t \langle \dot{\ell}, y \rangle$$

$$(4.3) \quad y(0) = y_0.$$

Mielke and Theil [32] proved that the latter is equivalent to (1.1) and hence, owing to the characterization of Theorem 3.1, to $F(y) = 0 = \min F$ (note that the analysis in [32] is much more general and is in particular allowing discontinuous in time evolutions by introducing the above notion of energetic solutions in the frame of functions of a bounded variation). For the aim of pointing out some features of our variational approach, we shall present here a direct proof of this fact.

LEMMA 4.3 (equivalence with the energetic formulation). *Let $\ell \in W^{1,1}(0, T; Y^*)$. Then $F(y) = 0 = \min F$ iff y fulfills (4.1)–(4.3).*

Proof. Owing to Lemma 4.2, we readily have that the stability condition (4.1) holds iff $\psi^*(\ell - Ay) = 0$ a.e.

Let y be such that $F(y) = 0$. Then (4.1) and (4.3) hold and $L(t, y(t), \dot{y}(t)) = 0$ for a.e. $t \in (0, T)$. In particular, for all $t \in [0, T]$,

$$\begin{aligned} 0 &= \int_0^t L(s, y(s), \dot{y}(s)) \, ds = \int_0^t \left(\psi(\dot{y}) - \langle \ell, \dot{y} \rangle \right) + \phi(y) \Big|_0^t \\ &= (\phi(y) - \langle \ell, y \rangle) \Big|_0^t + \int_0^t \psi(\dot{y}) + \int_0^t \langle \dot{\ell}, y \rangle, \end{aligned}$$

so that the energy equality (4.2) holds for all $t \in [0, T]$.

On the contrary, let $y \in W^{1,1}(0, T; Y)$ fulfill (4.1)–(4.3). Then $\chi(y(0) - y_0) = 0$ and $\psi^*(\ell - Ay) = 0$ a.e. (see above). Hence $F(y) = 0$ follows from the energy equality (4.2) at time T and an integration by parts. \square

Let us mention that the last lemma proves in particular that the energy equality (4.2) could be equivalently enforced *at the final time* T only. Moreover, it proves that, as already commented in the introduction, all *stable trajectories* $t \mapsto y(t)$ (i.e., trajectories such that $y(t) \in S(t)$ for all $t \in [0, T]$) are such that the following energy inequality holds:

$$\phi(y(t)) - \langle \ell(t), y(t) \rangle + \int_0^t \psi(\dot{y}) \geq \phi(y(0)) - \langle \ell(0), y(0) \rangle - \int_0^t \langle \dot{\ell}, y \rangle \quad \forall t \in [0, T].$$

Hence, we have provided a proof to [26, Prop. 5.7] (note that the referred result is, however, more general as it is concerned with the *BV* situation, the energy is implicitly depending on time, and no linear structure on Y is required).

Before closing this subsection, let us explicitly remark that the above inferred equivalence between formulations has been obtained for the absolutely continuous case $\ell \in W^{1,1}(0, T; Y^*)$ only, whereas the characterization of Theorem 3.1 holds more generally for bounded ℓ .

4.4. The functional controls the uniform distance: Uniqueness. So far, we have simply reformulated known results in a variational fashion. Here we present some novel results instead.

LEMMA 4.4 (uniform distance control via F). *We have*

$$(4.4) \quad \begin{aligned} \eta(1 - \eta) \max_{t \in [0, T]} \phi(u(t) - v(t)) &\leq \eta F(u) + (1 - \eta) F(v) \\ \forall u, v \in W^{1,1}(0, T; Y), \eta \in [0, 1]. \end{aligned}$$

Proof. The statement follows from the quadratic character of ϕ . Fix $t \in [0, T]$, and define $F^t : W^{1,1}(0, t; Y) \rightarrow [0, \infty]$ as

$$\begin{aligned} F^t(y) &= \int_0^t L(s, y(s), \dot{y}(s)) \, ds + \chi(y(0) - y_0) \\ &= \int_0^t \left(\psi(\dot{y}) + \psi^*(\ell - Ay) - \langle \ell, \dot{y} \rangle \right) \\ &\quad + \phi(y(t)) + \phi(y_0) - \langle Ay(0), y_0 \rangle + |y(0) - y_0|^2. \end{aligned}$$

Then, clearly, $y \mapsto G^t(y) = F^t(y) - \phi(y(t))$ is convex. Hence, by letting $w = \eta u + (1 - \eta)v$ we have

$$\begin{aligned} 0 &\leq F^t(w) \\ &\leq \eta(G^t(u) + \phi(u(t))) + (1 - \eta)(G^t(v) + \phi(v(t))) - \eta(1 - \eta)\phi(u(t) - v(t)), \end{aligned}$$

whence the assertion follows. \square

The latter lemma exploits the quadratic character of ϕ only. In particular, no coercivity for ϕ is assumed. It should, however, be clear that its application (as the title of the lemma indeed suggests) will always be referred to the situation where the stronger (3.6) is required, namely, when the left-hand side of (4.4) controls

$$\eta(1 - \eta) \frac{\alpha}{2} \max_{[0, T]} |u - v|^2.$$

We now present two immediate corollaries of Lemma 4.4.

COROLLARY 4.5 (uniform distance from the minimizer). *Let $F(y) = 0$. Then*

$$\max_{t \in [0, T]} \phi(y(t) - v(t)) \leq F(v).$$

This corollary encodes an interesting novel feature of our variational approach, for it provides a possible a posteriori error estimator to be used within approximation procedures. It is interesting to remark that the uniform distance of any stable trajectory from the minimizer is controlled by means of its energy production along the path only. Again, although Corollary 4.5 holds under no coercivity assumptions of ϕ , let us mention that its application will be restricted to the frame of (3.6). Finally, we have uniqueness of the minimizers of F attaining the value 0.

COROLLARY 4.6 (uniqueness). *Assume (3.6). Then there exists at most one trajectory y such that $F(y) = 0$.*

4.5. Lipschitz bound. Throughout the remainder of the paper we shall tacitly assume that

$$(4.5) \quad \ell \in W^{1, \infty}(0, T; Y^*), \quad y_0 \in S(0).$$

As already commented after Lemma 4.2, the above restriction on the initial datum is mandatory whenever ℓ admits a weak right limit in 0.

As for ℓ , the extra Lipschitz continuity assumption is motivated by the rate independence of the problem (every absolutely continuous datum can be time-rescaled to a Lipschitz continuous datum) and the following well known result.

LEMMA 4.7 (Lipschitz bound). *Assume (3.6), and let $\ell \in W^{1, \infty}(0, T; Y^*)$ and $F(y) = 0$. Then*

$$(4.6) \quad \|\dot{y}\|_{L^\infty(0, T; Y)} \leq \frac{1}{\alpha} \|\dot{\ell}\|_{L^\infty(0, T; Y^*)} \quad \text{a.e. in } (0, T).$$

The proof of the lemma is exactly the classical one [32, Thm. 7.5] but formulated by means of our variational arguments. We provide it for the sake of completeness.

Proof. Let $0 \leq s < t \leq T$ be fixed. Since $L(y, \dot{y}) = 0$ a.e. we have

$$\int_s^t \left(\psi(\dot{y}) + \langle \dot{\ell}, y \rangle \right) + (\phi(y) - \langle \ell, y \rangle) \Big|_s^t = 0.$$

On the other hand, owing to the strong monotonicity of A and the fact that $y(s) \in S(s)$ (see Lemma 4.2), one obtains

$$\begin{aligned} \phi(y(t) - y(s)) &\leq \phi(y(t)) - \langle \ell(s), y(t) \rangle + \psi(y(t) - y(s)) - \phi(y(s)) - \langle \ell(s), y(s) \rangle \\ &= (\phi(y) - \langle \ell, y \rangle) \Big|_s^t + \int_s^t \langle \dot{\ell}(r), y(t) \rangle \, dr + \psi(y(t) - y(s)). \end{aligned}$$

By taking the sum of these two relations and recalling that, by Jensen,

$$\psi(y(t) - y(s)) = \psi \left(\int_s^t \dot{y} \right) \leq \int_s^t \psi(\dot{y}),$$

we get

$$\phi(y(t) - y(s)) \leq \int_s^t \langle \dot{\ell}(r), y(t) - y(r) \rangle dr.$$

Finally, an application of some extended Gronwall lemma (see [26, Thm. 3.4]) entails that

$$\frac{\alpha}{2} |y(t) - y(s)|^2 \leq \frac{1}{2} \|\dot{\ell}\|_{L^\infty(0,T;Y^*)} |y(t) - y(s)| (t - s),$$

and the assertion follows. \square

5. Space approximation and stability under data perturbation. We now apply the characterization results of Theorem 3.1 to the approximation of solutions of (1.1). As already commented in the introduction, we shall proceed via Γ -convergence [11]. The reader is referred to the monographs by Attouch [1] and Dal Maso [5] for some comprehensive discussion on this topic. Indeed, since Theorem 3.1 directly quantifies the value of the minimum to be 0, what is actually needed for passing to limits are Γ -liminf inequalities only. We shall illustrate this fact by discussing the simple case of stability under data perturbation first.

LEMMA 5.1 (stability under data perturbation). *Assume (3.6), and let $\ell_h \rightarrow \ell$ strongly in $L^1(0, T; Y^*)$ being uniformly Lipschitz continuous and $y_{0,h} \rightarrow y_0$. Moreover, let $F_h : W^{1,1}(0, T; Y) \rightarrow [0, \infty]$ be defined as*

$$F_h(y) = \int_0^T \left(\psi(\dot{y}) + \psi^*(\ell_h - Ay) - \langle \ell_h - Ay, \dot{y} \rangle \right) + \chi(y(0) - y_{0,h}),$$

and let $F_h(y_h) = 0$. Then $y_h \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ and $F(y) = 0$.

Proof. Owing to Lemma 4.7, we find a (not relabeled) subsequence y_h such that $y_h \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$. Hence, we have by lower semicontinuity

$$\begin{aligned} 0 \leq F(y) &\leq \liminf_{h \rightarrow 0} \left(\int_0^T \left(\psi(\dot{y}_h) + \psi^*(\ell_h - Ay_h) - \langle \ell_h, \dot{y}_h \rangle \right) \right. \\ &\quad \left. + \phi(y_h(T)) + \phi(y_{0,h}) - \langle Ay_h(0), y_{0,h} \rangle + |y_h(0) - y_{0,h}|^2 \right) \\ &= \liminf_{h \rightarrow 0} F_h(y_h) = 0. \end{aligned}$$

Hence, $F(y) = 0$, y is unique, and the assertion follows from the fact that the whole sequence converges. \square

5.1. Preliminaries on functional convergence. In order to move to more general approximation situations, we are forced to discuss a suitable functional convergence notion. We limit ourselves in introducing the relevant definitions, referring to the mentioned monographs for all of the necessary details.

Recall that Y is a real reflexive Banach space. Letting $f_n, f : Y \rightarrow (-\infty, \infty]$ be convex, proper, and lower semicontinuous, we say that $f_n \rightarrow f$ in the *Mosco sense* in Y [1, 36] iff, for all $y \in Y$,

$$f(y) \leq \liminf_{n \rightarrow \infty} f_n(y_n) \quad \forall y_n \rightarrow y \text{ weakly in } Y,$$

$$\exists y_n \rightarrow y \text{ strongly in } Y \text{ such that } f(y) = \limsup_{n \rightarrow \infty} f_n(y_n).$$

In particular, $f_n \rightarrow f$ in the Mosco sense iff $f_n \rightarrow f$ in the sense of Γ -convergence with respect to both the weak and the strong topology in Y .

We will consider the situation of approximating functionals ψ_h . By [1, Thm. 3.18, p. 295], we have $\psi_h \rightarrow \psi$ in the Mosco sense in Y iff $\psi_h^* \rightarrow \psi^*$ in the Mosco sense in Y^* . By assuming the functionals ψ_h to be positively 1-homogeneous, it turns out that the Mosco convergence $\psi_h \rightarrow \psi$ in Y is equivalent to the Mosco convergence of sets $C_h^* \rightarrow C^*$ in Y^* , which reads

$$C_n^* \ni y_n^* \rightarrow y^* \text{ weakly in } Y^* \Rightarrow y^* \in C^*,$$

$$\forall y^* \in C^*, \exists y_n^* \in C_n^* : y_n^* \rightarrow y^* \text{ strongly in } Y^*.$$

Finally, we repeatedly use a lemma from [45] which we report it here, for the sake of completeness.

LEMMA 5.2 ([45, Cor. 4.4]). *Let $p \in [1, \infty]$ and $f_h, f : Y \rightarrow (-\infty, \infty]$ be convex, proper, and lower semicontinuous such that*

$$f(y) \leq \inf \left\{ \liminf_{h \rightarrow 0} f_h(y_h) : y_h \rightarrow y \text{ weakly in } Y \right\} \quad \forall y \in Y.$$

Moreover, let $y_h \rightarrow y$ weakly in $W^{1,p}(0, T; Y)$ (weakly star if $p = \infty$). Then we have

$$\int_0^T f(y(t)) \, dt \leq \liminf_{h \rightarrow 0} \int_0^T f_h(y_h(t)) \, dt.$$

5.2. Space approximations. We now move to the analysis of some space approximation situation, indeed specifically tailored for the case of conformal finite elements. Let us list here our assumptions for the sake of later referencing.

We assume to be given

- (5.1) $Y_h \subset Y$ closed subspaces such that $h \mapsto Y_h$ increases and $\cup_{h>0} Y_h$ is dense in Y ,
- (5.2) $\phi_h(y) = \phi(y)$ if $y \in Y_h$ and $\phi_h(y) = \infty$ otherwise.
- (5.3) $\psi_h : Y \rightarrow (-\infty, \infty]$ convex, proper, and lower semicontinuous,
- (5.4) ψ_h positively 1-homogeneous,
- (5.5) $\psi_h \rightarrow \psi$ in the Mosco sense in Y ,
- (5.6) $\phi(y) \geq \frac{\alpha}{2}|y|^2 \quad \forall y \in C_h - C_h$ where $C_h = D(\psi_h)$,
- (5.7) $\ell_h \rightarrow \ell$ pointwise strongly in Y^* ,
- (5.8) ℓ_h uniformly Lipschitz continuous,
- (5.9) $y_{0,h} \in Y_h, y_{0,h} \rightarrow y_0$.

We shall mention that within the frame of conformal finite elements methods the subspaces Y_h are obviously taken to be finite-dimensional and that the approximating functionals ϕ_h and ψ_h are usually the restrictions of the functionals ϕ and ψ on the subspace Y_h . This is exactly our choice here for ϕ_h . In particular, one shall observe that $\phi_h \rightarrow \phi$ in the Mosco sense in Y , $D(\partial\phi_h) = Y_h$, and that

$$(5.10) \quad A_h y = \partial\phi_h(y) = \partial\phi(y) = D\phi(y) = Ay \quad \forall y \in Y_h.$$

As for ψ_h we are allowing some extra freedom (let us remark, however, that (5.6) follows from (3.6) as soon as ψ_h is the restriction of ψ to Y_h since, in this case, $C_h = C \cap Y_h$). On the other hand, we are asking ψ_h to be positively 1-homogeneous; namely, we are considering the case of some rate-independent approximation of (1.1) only. The reader is referred instead to Efendiev and Mielke [8], Mielke, Rossi, and Savaré [29, 28], Toader and Zanini [47], and Zanini [49] for some results in the direction of rate-dependent approximation of rate-independent processes.

Finally, we shall (re)define the approximating functionals as $F_h : W^{1,1}(0, T; Y) \rightarrow [0, \infty]$ as

$$F_h(y) = \int_0^T \left(\psi_h(\dot{y}) + \psi_h^*(\ell_h - A_h y) - \langle \ell_h - A_h y, \dot{y} \rangle \right) + \chi_h(y(0) - y_{0,h}),$$

where $A_h = \partial\phi_h$ and $\chi_h(\cdot) = \phi_h(\cdot) + |\cdot|^2$. We have the following.

THEOREM 5.3 (convergence of space approximations). *Assume (5.1)–(5.9), and let $F_h(y_h) = 0$. Then $y_h \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ and $F(y) = 0$.*

Proof. By Lemma 4.7, we find a (not relabeled) subsequence $y_n \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ and weakly pointwise. Since $F_h(y_h) = 0$ we readily check that $y(t) \in Y_h$ for all $t \in [0, T]$. In particular, $A_h y_h = Ay_h$ for all $t \in [0, T]$ owing to (5.10). Hence, by lower semicontinuity,

$$\begin{aligned} 0 \leq F(y) &\leq \liminf_{h \rightarrow 0} \left(\int_0^T \left(\psi_h(\dot{y}_h) + \psi_h^*(\ell_h - Ay_h) - \langle \ell_h, \dot{y}_h \rangle \right) \right. \\ &\quad \left. + \phi(y_h(T)) + \phi(y_{0,h}) - \langle Ay_h(0), y_{0,h} \rangle + |y_h(0) - y_{0,h}|^2 \right) \\ &= \liminf_{h \rightarrow 0} \left(\int_0^T \left(\psi_h(\dot{y}_h) + \psi_h^*(\ell_h - Ay_h) - \langle \ell_h, \dot{y}_h \rangle \right) \right. \\ &\quad \left. + \phi_h(y_h(T)) - \phi_h(y_h(0)) \right) \\ &= \liminf_{h \rightarrow 0} F_h(y_h) = 0. \end{aligned}$$

Note that the integral terms containing ψ and ψ^* pass to the \liminf by means of Lemma 5.2. \square

By inspecting the proof of Theorem 5.3 (which of course generalizes Lemma 5.1), one realizes that, whenever the weak-star precompactness in $W^{1,\infty}(0, T; Y)$ of the sequence y_h is assumed, the convergence statement holds more generally in the case $F_h(y_h) \rightarrow 0$. Namely, by directly asking for the above-mentioned compactness, one could consider the convergence of some approximated solutions y_h such that, possibly, $F_h(y_h) > 0$. We rephrase this fact in the following statement.

LEMMA 5.4 (Γ -lim inf inequality for F_h). *Assume (5.1)–(5.3), (5.5)–(5.7), and (5.9). Then*

$$F(u) \leq \inf \left\{ \liminf_{h \rightarrow 0} F_h(y_h) : y_h \rightarrow y \text{ weakly star in } W^{1,\infty}(0, T; Y) \right\}.$$

Note that the homogeneity of ψ_h , the uniform convexity of ϕ_h , and the Lipschitz continuity of ℓ_h play no role here.

Finally, again by looking carefully to the proof of Theorem 5.3, one could wonder if the requirement on the Mosco convergence of ψ_h could be weakened. Indeed, what we are actually using is only that

$$(5.11) \quad \psi \leq \Gamma\text{-}\liminf_{h \rightarrow 0} \psi_h \quad \text{and} \quad \psi^* \leq \Gamma\text{-}\liminf_{h \rightarrow 0} \psi_h^*$$

with respect to the weak topologies of Y and Y^* , respectively. On the other hand, in our specific situation, [45, Lem. 4.1] entails that (5.11) and the fact that $\psi_h \rightarrow \psi$ in the Mosco sense in Y are equivalent.

This observation motivates once again the belief that Mosco convergence is the right frame in order to pass to limits within rate-independent problems. For the sake of completeness, let us recall that a first result in the direction of the approximation of the play operator (Y Hilbert and A coercive on Y) under the Hausdorff convergence of the characteristic sets $C_h^* = D(\psi_h^*)$ is contained in [16, Thm. 3.12, p. 34], whereas the extension of this result to the more general situation of Mosco converging sets as well as some application to parabolic PDEs with hysteresis is discussed in [46]. More recently, Mielke, Roubíček, and Stefanelli [30] addressed in full generality the issue of Γ -convergence and relaxation for the energetic solutions of rate-independent processes. An alternative convergence proof in the specific case of convex energies is obtained by means of the Brezis–Ekeland–Nayroles approach in [45].

6. Time discretization. Assume now that we are given the partitions $P_n = \{0 = t_n^0 < t_n^1 < \dots < t_n^{N_n} = T\}$, and denote by $\tau_n^i = t_n^i - t_n^{i-1}$ the i th time step and by $\tau_n = \max_{1 \leq i \leq N_n} \tau_n^i$ the diameter of the n th partition. No constraints are imposed on the possible choice of the time steps throughout this analysis besides $\tau_n \rightarrow 0$ as $n \rightarrow \infty$. Moreover, let the parameter $\theta \in [1/2, 1]$ be given.

In the following we will make extensive use of the following notation: Letting $v = (v^0, \dots, v^{N_n})$ be a vector, we will denote by \widehat{v}_n and \bar{v}_n two functions of the time interval $[0, T]$ which interpolate the values of the vector v piecewise linearly and backward constantly on the partition P_n , respectively. Namely,

$$\begin{aligned} \widehat{v}_n(0) &= v^0, \quad \widehat{v}_n(t) = \gamma_n^i(t)v^i + (1 - \gamma_n^i(t))v^{i-1}, \\ \bar{v}_n(0) &= v^0, \quad \bar{v}_n(t) = v^i, \quad \text{for } t \in (t_n^{i-1}, t_n^i], \quad i = 1, \dots, N_n \end{aligned}$$

where

$$\gamma_n^i(t) = (t - t_n^{i-1})/\tau_n^i \quad \text{for } t \in (t_n^{i-1}, t_n^i], \quad i = 1, \dots, N_n.$$

Moreover, we let $\delta v^i = (v^i - v^{i-1})/\tau_n^i$ for $i = 1, \dots, N_n$ (so that $\dot{\widehat{v}}_n = \overline{\delta v_n}$) and denote by v_θ the vector with components $v_\theta^i = \theta v^i + (1 - \theta)v^{i-1}$.

Recall that $\ell \in W^{1,\infty}(0, T; Y^*)$ and $y_0 \in S(0)$. We shall be concerned with the so-called θ -scheme for problem (1.1):

$$(6.1) \quad \partial\psi\left(\frac{y_n^i - y_n^{i-1}}{\tau_n^i}\right) + A(\theta y_n^i + (1-\theta)y_n^{i-1}) \ni \ell(\theta t_n^i + (1-\theta)t_n^{i-1})$$

for $i = 1, \dots, N_n$,

$$(6.2) \quad y_n^0 = y_0.$$

One usually refers to the latter as the backward or implicit Euler scheme for the choice $\theta = 1$ and as the Crank–Nicolson scheme for $\theta = 1/2$.

Owing to the above-introduced notation, the latter scheme can be equivalently rewritten as

$$(6.3) \quad \partial\psi(y_n^i - y_n^{i-1}) + Ay_{n,\theta}^i \ni \ell(t_{n,\theta}^i) \quad \text{for } i = 1, \dots, N_n, \quad y_n^0 = y_0.$$

Clearly, the θ -scheme (6.3) is rate-independent. Namely, no time step appears in (6.3), and the choice of the partition affects the solution via the values of the load ℓ only. In this concern, our focus on variable time-step partition could be simplified by considering proper rescaled loads ℓ instead. We shall, however, keep up with it, especially in order to underline the possibility of adapting the partition according to some a posteriori analysis (see subsection 6.7).

Before moving on, let us comment that, for all n , the latter scheme is a unique solution. Indeed, given $y_n^{i-1} \in C$, it suffices to (uniquely) solve iteratively the incremental problem

$$(6.4) \quad y_n^i \in \text{Arg min}_{y \in Y} \left(\theta\phi(y) - \langle \ell(t_{n,\theta}^i) - (1-\theta)Ay_n^{i-1}, y \rangle + \psi(y - y_n^{i-1}) \right).$$

Note that, since $y_n^{i-1} \in C$, the functional under minimization turns out to be uniformly convex. Hence, by (3.5), the minimum problem has a unique solution. In particular, exactly as in Lemma 4.1 we have the following.

LEMMA 6.1. $y_n^i \in C$ for all $i = 0, 1, \dots, N_n$.

A crucial observation is that, as in the continuous case, the discrete trajectories show some sort of stability as well.

LEMMA 6.2 (stability of the discrete trajectories). *We have*

$$(6.5) \quad y_n^i \in \text{Arg min}_{y \in Y} \left(\theta\phi(y) - \langle \ell(t_{n,\theta}^i) - (1-\theta)Ay_n^{i-1}, y \rangle + \psi(y - y_n^i) \right)$$

for $i = 1, \dots, N_n$.

In particular, if $\theta = 1$, then $y_n^i \in S(t_n^i)$.

Proof. From the incremental formulation (6.4) and the triangle inequality for ψ , we get, for all $y \in Y$,

$$\begin{aligned} & \theta\phi(y_n^i) - \langle \ell(t_{n,\theta}^i) - (1-\theta)Ay_n^{i-1}, y_n^i \rangle + \psi(y_n^i - y_n^{i-1}) \\ & \leq \theta\phi(y) - \langle \ell(t_{n,\theta}^i) - (1-\theta)Ay_n^{i-1}, y \rangle + \psi(y - y_n^{i-1}) \\ & \leq \theta\phi(y) - \langle \ell(t_{n,\theta}^i) - (1-\theta)Ay_n^{i-1}, y \rangle + \psi(y - y_n^i) + \psi(y_n^i - y_n^{i-1}), \end{aligned}$$

whence the assertion follows. \square

Again as in the continuous case, we readily check that

$$(6.6) \quad (6.5) \text{ holds iff } \ell(t_{n,\theta}^i) - Ay_{n,\theta}^i \in C^*.$$

6.1. The discrete variational principle. We shall now present a discrete version of the variational principle of Theorem 3.1.

We define $L_n^{\theta,i}(y, z) : Y \times Y \rightarrow [0, \infty]$ as

$$L_n^{\theta,i}(y, z) = \psi\left(\frac{y-z}{\tau_n^i}\right) + \psi^*\left(\ell(t_{n,\theta}^i) - A(\theta y + (1-\theta)z)\right) - \left\langle \ell(t_{n,\theta}^i) - A(\theta y + (1-\theta)z), \frac{y-z}{\tau_n^i} \right\rangle$$

and the functionals $F_n^\theta : Y^{N_n+1} \rightarrow [0, \infty]$ as

$$F_n^\theta(y_n^0, \dots, y_n^{N_n}) = \sum_{i=1}^{N_n} \tau_n^i L_n^{\theta,i}(y_n^i, y_n^{i-1}) + \chi(y_n^0 - y_0).$$

LEMMA 6.3 (discrete variational principle). $(y_n^0, \dots, y_n^{N_n})$ solves (6.3) iff $F_n^\theta(y_n^0, \dots, y_n^{N_n}) = 0 = \min F_n^\theta$.

Proof. Analogously to the continuous case, we have, for all $i = 1, \dots, N_n$,

$$\partial\psi(\delta y_n^i) + Ay_{n,\theta}^i \ni \ell(t_{n,\theta}^i) \quad \text{iff} \quad L_n^{\theta,i}(y_n^i, y_n^{i-1}) = 0,$$

and $y_n^0 = y_0$ iff $\chi(y_n^0 - y_0) = 0$. \square

Let us observe that the functional F_n^θ is convex and lower semicontinuous. Moreover, by the homogeneity of ψ (see (3.2)), F_n^θ is actually independent of the time steps. In fact, we have

$$F_n^\theta(y_n^0, \dots, y_n^{N_n}) = \sum_{i=1}^{N_n} \left(\psi(y_n^i - y_n^{i-1}) + \psi^*\left(\ell(t_{n,\theta}^i) - Ay_{n,\theta}^i\right) - \left\langle \ell(t_{n,\theta}^i) - Ay_{n,\theta}^i, y_n^i - y_n^{i-1} \right\rangle \right) + \chi(y_n^0 - y_0).$$

The idea of dealing with time discretizations via a discrete variational principle closely relates our analysis to the theory of so-called *variational integrators*. The latter are numerical schemes stemming from the approximation of the action functional in Lagrangian mechanics. By referring the reader to the monograph [12] and the survey [25], we shall restrain here from giving a detailed presentation of the subject and limit ourselves to some (necessarily sketchy) considerations. By letting $(t, y, p) \in [0, T] \times \mathbb{R}^m \times \mathbb{R}^m \mapsto \mathcal{L}(t, y, p)$ denote the Lagrangian of a (finite-dimensional, for simplicity) system, the Hamilton principle asserts that the actual trajectory $t \mapsto y(t)$ of the system minimizes the action functional

$$y \mapsto \int_0^T \mathcal{L}(t, y(t), \dot{y}(t)) \, dt$$

among all curves with prescribed end points, thus solving the Lagrange equations

$$(6.7) \quad \partial_{y_i} \mathcal{L} - \frac{d}{dt} \partial_{p_i} \mathcal{L} = 0 \quad \text{for } i = 1, \dots, m.$$

Hence, a natural idea is that of deriving numerical schemes for Lagrangian mechanics by applying some quadrature procedure to the action functional, i.e., discretizing

Hamilton’s principle. The resulting discrete schemes show comparable performance with respect to other methods but generally enjoy some interesting extra (and often crucial) properties such as the conservation of suitable quantities [20]. Variational integrators have been intensively applied in finite-dimensional contexts and, more recently, to the situation of nonlinear wave equations [24] and nonequilibrium elasticity [19].

The present analysis may bear some resemblance to the above-mentioned theory. Indeed, the formulation of the θ -scheme in the case $\theta = 1/2$ stems exactly from the midpoint quadrature of the functional F as

$$\begin{aligned} & \int_{t_n^{i-1}}^{t_n^i} L(t, \widehat{y}(t), \dot{\widehat{y}}(t)) \, dt \\ &= \tau_n^i L(t_{n,1/2}^i, \widehat{y}(t_{n,1/2}^i), \dot{\widehat{y}}(t_{n,1/2}^i)) \\ &= \psi(y^i - y^{i-1}) + \psi^* \left(\ell(t_{n,1/2}^i) - A \left(\frac{y^i + y^{i-1}}{2} \right) \right) \\ &\quad - \left\langle \ell(t_{n,1/2}^i) - A \left(\frac{y^i + y^{i-1}}{2} \right), y_n^i - y_n^{i-1} \right\rangle, \end{aligned}$$

where \widehat{y} is taken to be piecewise affine on the partition P_n .

On the other hand, our focus here is quite different. First of all, we are not dealing with the Hamilton principle (end points are not fixed) as we are not aimed at solving the Euler–Lagrange equations for F (i.e., solving (6.7)). Second, we are specifically interested at infinite-dimensional situations, namely, PDEs. Finally, the only choice of θ which is directly related with a quadrature of F is $\theta = 1/2$, and we are not considering higher-order schemes.

Before closing this discussion, let us mention that some Γ -convergence techniques have been recently exploited in the (finite-dimensional) frame of variational integrators by Müller and Ortiz [37] (see also [23]).

6.2. Stability of the θ -scheme. It has been known since Han and Reddy [15, 13] that the choice $\theta < 1/2$ in (6.3) leads to an unconditionally unstable scheme and that, on the contrary, for $\theta \in [1/2, 1]$ the θ -scheme is stable in $H^1(0, T; Y)$ when Y is a Hilbert space and the partitions are chosen to be uniform.

Here we shall provide an alternative stability proof by taking into account the Banach-space frame.

LEMMA 6.4 (stability). *Assume (3.6), and let $\theta \in [1/2, 1]$. Then the solution to the θ -scheme (6.3) fulfills*

$$(6.8) \quad \|\widehat{y}_{n,\theta}\|_{L^\infty(0,T;Y)} \leq \frac{1}{\alpha} \|\dot{\ell}\|_{L^\infty(0,T;Y^*)} \quad \text{if } \theta = 1 \text{ or } \theta = \frac{1}{2}.$$

Moreover, for constant time steps,

$$(6.9) \quad \|\dot{\widehat{y}}_{n,\theta}\|_{L^\infty(0,T;Y)} \leq \frac{1}{\alpha(2\theta - 1)} \|\dot{\ell}\|_{L^\infty(0,T;Y^*)} \quad \text{if } \frac{1}{2} < \theta < 1.$$

Our argument coincides with that of [32, Thm. 4.4] in the case of the implicit Euler scheme, i.e., $\theta = 1$, and it is an extension of the latter for the case $1/2 < \theta < 1$. Here we do not play with the variational inequality by choosing suitable tests but use the scalar relations $L_n^{\theta,i}(y_n^i, y_n^{i-1}) = 0$ instead (this, however, makes no substantial

difference since the latter scalar relations are exactly the outcome of the test on the variational inequality in [32, Thm. 4.4]).

The stability proof for the Crank–Nicolson scheme $\theta = 1/2$ is quite different from former arguments and stems as a direct outcome of our variational approach. In both cases $\theta = 1$ and $\theta = 1/2$, the stability constant $1/\alpha$ is sharp (see Lemma 4.7).

We complement this analysis by providing the stability for the θ -scheme for $1/2 < \theta < 1$ in the case of constant time steps (likely with a nonoptimal, although explicit, stability constant).

Note that the stability estimates (6.8)–(6.9) do not hold in the classical parabolic situation (i.e., ψ quadratic on a Hilbert space) unless the very restrictive compatibility assumption $Ay_0 = \ell(0)$ is made (and in this case, the proof below just goes through). Evidence of this failure is available even in the simplest scalar situation ($Y = \mathbb{R}$) of problem $\dot{y}(t) + y(t) = 0$ for $t > 0$ and $y(0) = y_0$. Indeed, the corresponding θ -scheme fulfills (6.8)–(6.9) iff $y_0 = 0$ (and hence $y \equiv 0$). The current rate-independent situation turns out to be better behaved since the restriction $Ay_0 = \ell(0)$ of the parabolic case is replaced by the weaker $y_0 \in S(0)$ (see (4.5)).

Proof. Let us prove the stability of the Crank–Nicolson scheme $\theta = 1/2$ first. For this aim, it suffices to recall that

$$\begin{aligned} 0 &= F_n^{1/2}(y_n^0, \dots, y_n^{N_n}) \\ &= \int_0^T \left(\psi(\dot{\hat{y}}_n) + \psi^*(\hat{\ell}_n - A\hat{y}_n) - \langle \hat{\ell}_n - A\hat{y}_n, \dot{\hat{y}}_n \rangle \right) + \chi(\hat{y}_n(0) - y_0). \end{aligned}$$

Hence \hat{y}_n minimizes the functional F where ℓ is replaced by $\hat{\ell}_n$. The stability estimate follows from Lemma 4.7.

Let us now move to the case $1/2 < \theta \leq 1$. Relation (6.5) applied at level $i - 1$ for some $i = 2, \dots, N_n$ along with the choice $y = y_n^i$ entails that

$$\begin{aligned} &\theta\phi(y_n^i - y_n^{i-1}) + \theta\phi(y_n^{i-1}) - \langle \ell(t_{n,\theta}^{i-1}) - (1 - \theta)Ay_n^{i-2}, y_n^{i-1} \rangle \\ &\leq \theta\phi(y_n^i) - \langle \ell(t_{n,\theta}^{i-1}) - (1 - \theta)Ay_n^{i-2}, y_n^i \rangle + \psi(y_n^i - y_n^{i-1}), \end{aligned}$$

where the extra term $\theta\phi(y_n^i - y_n^{i-1})$ is obtained from the fact ϕ is quadratic. Hence, we have

$$\begin{aligned} &\theta\phi(y_n^i - y_n^{i-1}) + \theta\phi(y_n^{i-1}) - \theta\phi(y_n^i) \\ &\leq \langle \ell(t_{n,\theta}^{i-1}), y_n^{i-1} - y_n^i \rangle + (1 - \theta)\langle A(y_n^{i-2} - y_n^{i-1}), y_n^i - y_n^{i-1} \rangle + \psi(y_n^i - y_n^{i-1}) \\ &\quad + (1 - \theta)\langle Ay_n^{i-1}, y_n^i - y_n^{i-1} \rangle \\ &= \langle \ell(t_{n,\theta}^{i-1}), y_n^{i-1} - y_n^i \rangle + (1 - \theta)\langle A(y_n^{i-2} - y_n^{i-1}), y_n^i - y_n^{i-1} \rangle + \psi(y_n^i - y_n^{i-1}) \\ &\quad - (1 - \theta)\left(\phi(y_n^{i-1}) + \phi(y_n^i - y_n^{i-1}) - \phi(y_n^i) \right), \end{aligned}$$

so that

$$(6.10) \quad \phi(e_n^i) + \phi(y_n^{i-1}) - \phi(y_n^i) \leq -\langle \ell(t_{n,\theta}^{i-1}), e_n^i \rangle + (\theta - 1)\langle Ae_n^{i-1}, e_n^i \rangle + \psi(e_n^i),$$

where we have used $e_n^i = y_n^i - y_n^{i-1}$ in order to shorten notations.

Next, from $L_n^{\theta,i}(y_n^i, y_n^{i-1}) = 0$ for $i = 1, \dots, N_n$, we obtain

$$\begin{aligned} 0 &= \psi(e_n^i) - \langle \ell(t_{n,\theta}^i) - Ay_{n,\theta}^i, e_n^i \rangle \\ &= \psi(e_n^i) - \langle \ell(t_{n,\theta}^i), e_n^i \rangle + \theta\left(\phi(y_n^i) + \phi(e_n^i) - \phi(y_n^{i-1}) \right) \\ &\quad - (1 - \theta)\left(\phi(y_n^{i-1}) + \phi(e_n^i) - \phi(y_n^i) \right). \end{aligned}$$

In particular, we have checked that

$$(6.11) \quad \psi(e_n^i) + \phi(y_n^i) - \phi(y_n^{i-1}) + (2\theta - 1)\phi(e_n^i) = \langle \ell(t_{n,\theta}^i), e_n^i \rangle.$$

We take the sum between the latter and (6.10) and get

$$2\theta\phi(e_n^i) \leq \langle \ell(t_{n,\theta}^i) - \ell(t_{n,\theta}^{i-1}), e_n^i \rangle + (\theta - 1)\langle Ae_n^{i-1}, e_n^i \rangle$$

or, equivalently,

$$(6.12) \quad \langle Ae_n^{i,\theta}, e_n^i \rangle \leq \langle \ell(t_{n,\theta}^i) - \ell(t_{n,\theta}^{i-1}), e_n^i \rangle.$$

Now, if $\theta = 1$, we conclude that

$$|e_n^i| \leq \frac{1}{\alpha} |\ell(t_{n,\theta}^i) - \ell(t_{n,\theta}^{i-1})|_*,$$

and the assertion follows.

In case $1/2 \leq \theta < 1$ and for a constant time-step partition, one proceeds from (6.12) by computing

$$\begin{aligned} & \tau_n \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)} \sqrt{\frac{2}{\alpha}} \sqrt{\phi(e_n^i)} \\ & \geq \langle \ell(t_{n,\theta}^i) - \ell(t_{n,\theta}^{i-1}), e_n^i \rangle \\ & \geq \theta \langle Ae_n^i, e_n^i \rangle + (1 - \theta) \langle Ae_n^{i-1}, e_n^i \rangle \\ & = (2\theta - 1) \langle Ae_n^i, e_n^i \rangle + (1 - \theta) \langle A(e_n^i + e_n^{i-1}), e_n^i \rangle \\ & \geq 2(2\theta - 1)\phi(e_n^i) + (1 - \theta) (\phi(e_n^i) - \phi(e_n^{i-1})). \end{aligned}$$

Note that the coefficient $(2\theta - 1)$ is strictly positive as $\theta > 1/2$. By using the fact that $y_n^0 = y_0 \in S(0)$ (recall (4.5)), we readily check that

$$(6.13) \quad \phi(e_n^1) + \phi(y_n^0) - \langle \ell(0), y_n^0 \rangle \leq \phi(y_n^1) - \langle \ell(0), y_n^1 \rangle + \psi(e_n^1),$$

and, by adding the latter to (6.11) for $i = 1$, we have

$$(6.14) \quad 2\theta\phi(e_n^1) \leq \langle \ell(t_{n,\theta}^1) - \ell(0), e_n^1 \rangle \leq \tau_n \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)} \sqrt{\frac{2}{\alpha}} \sqrt{\phi(e_n^1)}.$$

Let us define

$$\begin{aligned} a_i^2 &= \phi\left(\frac{y_n^i - y_n^{i-1}}{\tau_n}\right) = \phi(e_n^i)/\tau_n^2, \\ C_0 &= \frac{2(2\theta - 1)}{1 - \theta}, \quad C_1 = \frac{1}{1 - \theta} \sqrt{\frac{2}{\alpha}} \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)}, \quad C_2 = \frac{C_1}{C_0}, \end{aligned}$$

so that, owing to (6.13) and (6.14) and by using the fact that $2(2\theta - 1) < 2\theta$,

$$\begin{aligned} (C_0 + 1)a_i^2 - a_{i-1}^2 &\leq C_1 a_i \quad \text{for } i = 2, \dots, N_n, \\ a_1 &\leq \frac{1}{2\theta} \sqrt{\frac{2}{\alpha}} \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)} \leq \frac{1}{2(2\theta - 1)} \sqrt{\frac{2}{\alpha}} \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)} \\ &= \frac{1}{1 - \theta} \sqrt{\frac{2}{\alpha}} \|\hat{\ell}_n\|_{L^\infty(0,T;Y^*)} \frac{1 - \theta}{2(2\theta - 1)} \\ &= \frac{C_1}{C_0} = C_2. \end{aligned}$$

Now, since $(C_0 + 1)C_2^2 - C_1C_2 = C_2^2$, we easily prove by induction that $a_i \leq C_2$, and the assertion follows. \square

6.3. Convergence. We shall prove the weak-star $W^{1,\infty}(0, T; Y)$ convergence for the θ -method. This result has to be compared with that of Han and Reddy [14, Thm. 3.4], where the uniform convergence of the backward constant interpolations is obtained. Our result is weaker than that of [14, Thm. 3.4] since we are not providing strong convergence. On the other hand, we believe our half-page proof to be possibly more transparent than the long argument developed in [14]. Let us moreover mention that, in the Hilbertian case and for A coercive on Y , the strong convergence in $W^{1,p}(0, T; Y)$ for all $p < \infty$ of the Euler method $\theta = 1$ has been proved in [16, Prop. 3.9, p. 33].

THEOREM 6.5 (convergence for the θ -method). *Assume (3.6), and let $F_n^\theta(y_n^0, \dots, y_n^{N_n}) = 0$. Then $\widehat{y}_n \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$, where $F(y) = 0 = \min F$.*

Proof. Owing to Lemma 6.4, we can extract a (not relabeled) subsequence such that $\widehat{y}_n \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ and hence weakly pointwise in Y . Moreover, we clearly have that both \bar{y}_n and $\bar{y}_{n,\theta}$ converge at the same limit weakly star in $L^\infty(0, T; Y)$. Finally, we directly check that $\bar{\ell}_{n,\theta} \rightarrow \ell$ strongly in $L^\infty(0, T; Y^*)$. By observing that, since $\theta \geq 1/2$,

$$\begin{aligned} \tau_n^i \langle A(\theta y_n^i + (1 - \theta)y_n^{i-1}), \delta y_n^i \rangle &= \phi(y_n^i) + (2\theta - 1)\phi(y_n^i - y_n^{i-1}) - \phi(y_n^{i-1}) \\ &\geq \phi(y_n^i) - \phi(y_n^{i-1}), \end{aligned}$$

we compute that

$$\begin{aligned} 0 &= F_n^\theta(y_n^0, \dots, y_n^{N_n}) \\ &\geq \int_0^T \left(\psi(\dot{\widehat{y}}_n) + \psi^*(\bar{\ell}_{n,\theta} - A\bar{y}_{n,\theta}) - \langle \bar{\ell}_{n,\theta}, \dot{\widehat{y}}_n \rangle \right) \\ &\quad + \phi(\widehat{y}_n(T)) - \phi(\widehat{y}_n(0)) + \chi(\widehat{y}_n(0) - y_0) \\ &= \int_0^T \left(\psi(\dot{\widehat{y}}_n) + \psi^*(\bar{\ell}_{n,\theta} - A\bar{y}_{n,\theta}) - \langle \bar{\ell}_{n,\theta}, \dot{\widehat{y}}_n \rangle \right) \\ &\quad + \phi(\widehat{y}_n(T)) + \phi(y_0) - \langle A\widehat{y}_n(0), y_0 \rangle + |\widehat{y}_n(0) - y_0|^2. \end{aligned}$$

Finally, it suffices to pass to the \liminf above as $n \rightarrow \infty$ and exploit lower semicontinuity and the stated convergences in order to obtain $F(y) \leq 0$. Hence, by Theorem 3.1 and Corollary 4.6, y is the only solution to (1.1), and the whole sequence \widehat{y}_n converges. \square

6.4. The functional controls the uniform distance. We shall reproduce at the discrete level the results of subsection 4.6. We begin by showing how to possibly control the uniform distance of two vectors by means of the discrete functional F_n^θ .

LEMMA 6.6 (uniform distance control via F_n^θ). *Let the vectors $u = (u^0, \dots, u^{N_n})$ and $v = (v^0, \dots, v^{N_n}) \in Y^{N_n+1}$ be given. Then*

$$\begin{aligned} &\eta(1 - \eta) \max_{1 \leq i \leq N_n} \phi(u^i - v^i) \\ &\leq \eta F_n^\theta(u^0, \dots, u^{N_n}) + (1 - \eta) F_n^\theta(v^0, \dots, v^{N_n}) \quad \forall \eta \in [0, 1]. \end{aligned}$$

Proof. This proof follows the same lines as that of Corollary 4.4. Let $1 \leq i \leq N_n$ be fixed, and define $F_n^{\theta,i} : Y^{i+1} \rightarrow [0, \infty]$ as

$$\begin{aligned} F_n^{\theta,i}(y^0, \dots, y^i) &= \sum_{j=1}^i \tau_n^j L_n^{\theta,j}(y^j, y^{j-1}) + \chi(y^0 - y_0) \\ &= \sum_{j=1}^i \psi(y^j - y^{j-1}) + \psi^*(\ell(t_{n,\theta}^j) - Ay_\theta^j) - \langle \ell(t_{n,\theta}^j), y^j - y^{j-1} \rangle \\ &\quad + \phi(y^i) + (2\theta - 1) \sum_{j=1}^i \phi(y^j - y^{j-1}) + \phi(y_0) - \langle Ay^0, y_0 \rangle + |y^0 - y_0|^2. \end{aligned}$$

Then, clearly, $y = (y^0, \dots, y^i) \mapsto G_n^{\theta,i}(y) = F_n^{\theta,i}(y) - \phi(y^i)$ is convex. Hence, by letting $w = \eta u + (1 - \eta)v$ for some $\eta \in [0, 1]$, we have

$$\begin{aligned} 0 &\leq F_n^{\theta,i}(w) \\ &\leq \eta(G_n^{\theta,i}(u) + \phi(u^i)) + (1 - \eta)(G_n^{\theta,i}(v) + \phi(v^i)) - \eta(1 - \eta)\phi(u^i - v^i), \end{aligned}$$

whence the assertion follows. \square

Again, note that the latter lemma controls the uniform norm of the distance only if the stronger (3.6) is required. The following corollary of Lemma 6.6 will be the starting point for some possible a posteriori error control procedure (see subsection 6.6).

COROLLARY 6.7 (uniform distance from the minimizer). *Let $F_n^\theta(y^0, \dots, y^{N_n}) = 0$. Then*

$$\max_{1 \leq i \leq N_n} \phi(y^i - v^i) \leq F_n^\theta(v^0, \dots, v^{N_n}) \quad \forall (v^0, \dots, v^{N_n}) \in Y^{N_n+1}.$$

Moreover, we reobtain a proof of the uniqueness of the solution of the θ -method.

COROLLARY 6.8 (uniqueness of the minimizer). *Assume (3.6). Then there exists at most one $y = (y^0, \dots, y^{N_n})$ such that $F_n^\theta(y) = 0$.*

6.5. The generalized θ -method. Although minimizers of F_n^θ and solutions of the θ -scheme (6.3) coincide, minimizing sequences of F_n^θ need not solve (6.3). This extra freedom allows the minimization formulation to capture the convergence of some generalized θ -method, where the relations in (6.3) are not solved exactly but rather are approximated. Namely, we shall look for vectors $u_n = (u_n^0, \dots, u_n^{N_n})$ such that

$$F_n^\theta(u_n) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

instead of $F_n^\theta(u_n) = 0$ for all $n \in \mathbb{N}$.

From the computational viewpoint, note that the θ -scheme consists in solving N_n nonlinear equations in one unknown each, while checking for stationarity for F_n^θ implies the solution of a tridiagonal system of $N_n + 1$ nonlinear equations with (up to) three unknowns each. This entails in particular that minimizing F_n^θ instead of solving (6.3) could be of a scarce interest if one is merely concerned in reproducing the θ -scheme with no error. On the other hand, the issue of solving up to some tolerance turns out to be particularly relevant whenever one is aimed at implementing an optimization procedure for the solution of (6.3). Indeed, one should be prepared to run the algorithm (some descent method, say) until some given tolerance is reached.

Our starting point for a possible convergence analysis of the generalized θ -method is the following classical error control result.

THEOREM 6.9 (Mielke and Theil [32]). *Assume (3.6). Then $\widehat{y}_n \rightarrow y$ uniformly and $F(y) = 0$. In particular,*

$$(6.15) \quad \max_{t \in [0, T]} |(\widehat{y}_n - y)(t)| \leq C_e \sqrt{\tau_n},$$

where C_e depends only on data and is independent of n .

More precisely, in [32] solely the case of the Euler scheme $\theta = 1$ is discussed. However, an easy adaptation of the argument entails the result for $\theta \in [1/2, 1)$ as well.

By explicitly comparing the minimizing sequence $u_n = (u_n^0, \dots, u_n^{N_n})$ with the corresponding solution $(y_n^0, \dots, y_n^{N_n})$ of the θ -method, we have the following.

THEOREM 6.10 (convergence for the generalized θ -method). *Assume (3.6), and let $F_n^\theta(u_n^0, \dots, u_n^{N_n}) \rightarrow 0$. Then $\widehat{u}_n \rightarrow y$ uniformly, where $F(y) = 0$. In particular,*

$$(6.16) \quad \max_{t \in [0, T]} |(\widehat{u}_n - y)(t)| \leq C_e \sqrt{\tau_n} + \left(\frac{2}{\alpha} F_n^\theta(u_n^0, \dots, u_n^{N_n}) \right)^{1/2}.$$

Proof. We have

$$\begin{aligned} \max_{t \in [0, T]} |(y - \widehat{u}_n)(t)| &\leq \max_{t \in [0, T]} |(y - \widehat{y}_n)(t)| + \max_{t \in [0, T]} |(\widehat{y}_n - \widehat{u}_n)(t)| \\ &\leq C_e \sqrt{\tau_n} + \max_{1 \leq i \leq N_n} |u_n^i - y_n^i| \\ &\leq C_e \sqrt{\tau_n} + \left(\frac{2}{\alpha} \max_{1 \leq i \leq N_n} \phi(u_n^i - y_n^i) \right)^{1/2}, \end{aligned}$$

and we conclude by applying Corollary 6.7. \square

6.6. A posteriori error control. Let us now exploit both Corollary 4.5 and Theorem 6.10 in order to provide some possible a posteriori estimates of the approximation error by means of solutions u_n of the generalized θ -method described above.

LEMMA 6.11 (a posteriori error control via F_n^θ). *Assume (3.6), and let $F_n^\theta(u_n^0, \dots, u_n^{N_n}) \sim \tau_n^s$ for some $s > 0$ and $F(y) = 0$. Then*

$$\max_{t \in [0, T]} |(\widehat{u}_n - y)(t)| \sim \tau_n^r, \quad \text{where } 2r = \max\{1, s\}.$$

LEMMA 6.12 (a posteriori error control via F). *Assume (3.6), and let $F(\widehat{u}_n) \sim \tau_n^s$ for some $s > 0$ and $F(y) = 0$. Then $\max_{t \in [0, T]} |(\widehat{u}_n - y)(t)| \sim \tau_n^{s/2}$.*

We are also in the position of proving the weak-star convergence of the time derivatives of solutions u_n of the generalized θ -method by comparing them with the corresponding derivatives of the exact solution of the θ -method.

LEMMA 6.13 (improved convergence for the generalized θ -method). *Assume (3.6), and let $F_n^\theta(u_n^0, \dots, u_n^{N_n}) \sim \tau_n^2$. Then \widehat{u}_n is equibounded in $W^{1, \infty}(0, T; Y)$. In particular, $\widehat{u}_n \rightarrow y$ weakly star in $W^{1, \infty}(0, T; Y)$.*

Proof. Let $(y_n^0, \dots, y_n^{N_n})$ be the solution of the θ -scheme. By exploiting Lemma 6.7, we check that

$$\begin{aligned} |u_n^i - u_n^{i-1}| &\leq |u_n^i - y_n^i| + |y_n^i - y_n^{i-1}| + |y_n^{i-1} - u_n^{i-1}| \\ &\leq \tau_n \|\dot{\widehat{y}}_n\|_{L^\infty(0, T; Y)} + 2 \left(\frac{2}{\alpha} F_n^\theta(u_n^0, \dots, u_n^{N_n}) \right)^{1/2}. \end{aligned}$$

The uniform bound on $\|\widehat{u}_n\|_{W^{1,\infty}(0,T;Y)}$ follows by dividing the latter by τ_n^i , taking the maximum as $1 \leq i \leq N_n$, and recalling Lemma 6.4. \square

6.7. Adaptivity. By assuming (3.6), the above-introduced a posteriori error estimators can be exploited in order to develop an adaptive strategy. In particular, the error control in the uniform norm up to a given tolerance $\text{tol} > 0$

$$\max_{t \in [0,T]} |(y - \widehat{y}_n)(t)| \leq \text{tol}$$

for some piecewise approximation \widehat{y}_n , with $\chi(\widehat{y}_n(0) - y_0) \leq \alpha \text{tol}^2/4$, can be inferred, for instance, by choosing time steps in such a way that

$$\int_{t_{i-1}^i}^{t_n^i} L(t, \widehat{y}_n(t), \dot{\widehat{y}}_n(t)) \leq \frac{\alpha \text{tol}^2}{4N_n},$$

namely, by uniformly distributing the error along the partition.

Alternatively, one could develop an adaptive strategy by considering just computed quantities at the discrete level by asking for

$$\tau_n^i L_n^{\theta,i}(y_n^i, y_n^{i-1}) \leq \frac{\alpha \text{tol}^2}{32N_n} \quad \text{for} \quad \tau_n \leq \frac{\text{tol}^2}{16C_e^2}$$

and exploiting Theorem 6.10.

7. Space-time approximations. Let us combine the results of the previous sections (and use the corresponding notation) in order to state and prove a result on the convergence of full space-time approximations. Our results have to be compared with the former convergence analysis by Han and Reddy [13]. Our approach leads to a convergence proof with respect to a weaker topology. However, it is, on the one hand, slightly more general (some assumptions on the spaces and the functionals—see (H1)–(H2) [13, p. 264]—are not required) and, on the other hand, has a much simpler proof.

THEOREM 7.1 (convergence of space-time approximations). *Assuming (5.1)–(5.9) and that $\theta \in [1/2, 1]$, define $L_{n,h}^{\theta,i}(y, z) : Y \times Y \rightarrow [0, \infty]$ as*

$$\begin{aligned} L_{n,h}^{\theta,i}(y, z) &= \psi_h \left(\frac{y - z}{\tau_n^i} \right) + \psi_h^* \left(\ell_h(t_{n,\theta}^i) - A_h(\theta y + (1 - \theta)z) \right) \\ &\quad - \left\langle \ell_h(t_{n,\theta}^i) - A_h(\theta y + (1 - \theta)z), \frac{y - z}{\tau_n^i} \right\rangle, \end{aligned}$$

where $A_h = \partial\phi_h$, and let the functionals $F_{n,h}^\theta : Y^{N_n+1} \rightarrow [0, \infty]$ be defined as

$$F_{n,h}^\theta(y^0, \dots, y^{N_n}) = \sum_{i=1}^{N_n} \tau_n^i L_{n,h}^{\theta,i}(y^i, y^{i-1}) + \chi_h(y^0 - y_0),$$

where $\chi_h(\cdot) = \phi_h(\cdot) + |\cdot|^2$ (note that $D(F_{n,h}^\theta) \subset Y_h^{N_n+1}$). Finally, let $F_{n,h}(y_h^0, \dots, y_h^{N_n}) = 0$. We have the following:

- (a) $\widehat{y}_{n,h} \rightarrow y_h$ weakly star in $W^{1,\infty}(0, T; Y)$ as $(n, h) \rightarrow (\infty, h)$ and $F_h(y_h) = 0$.
- (b) $\widehat{y}_{n,h} \rightarrow \widehat{y}_n$ weakly star in $W^{1,\infty}(0, T; Y)$ as $(n, h) \rightarrow (n, 0)$ and $F_n^\theta(y_n) = 0$.
- (c) $\widehat{y}_n \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ as $(n, 0) \rightarrow (\infty, 0)$ and $F(y) = 0$.

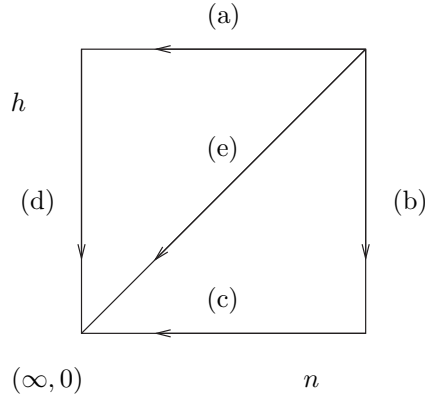


FIG. 1. Convergences for space-time approximations (see Theorem 7.1).

(d) $y_h \rightarrow y$ weakly star $W^{1,\infty}(0, T; Y)$ as $(\infty, h) \rightarrow (\infty, 0)$ and $F(y) = 0$.

(e) $\hat{y}_{n,h} \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$ as $(n, h) \rightarrow (\infty, 0)$ and $F(y) = 0$.

The thesis of the theorem is illustrated in Figure 1. In particular, we aim at showing that the space (or data) and time limit can be taken in any order. Note that limit (c) has been already checked in Theorem 6.5 and that the very same argument yields limit (a) as well (recall that Y_h is closed). Moreover, limit (d) is discussed in Theorem 5.3. So what we are actually left to check are limits (b) and (e) only.

Proof. Limit (b). The assertion follows once we check that, for all $i = 1, \dots, N_n$, if $y_{n,h}^{i-1} \rightarrow y_n^{i-1}$ weakly in Y , one has the weak convergence $y_{n,h}^i \rightarrow y_n^i$ as well. Recall that

$$\begin{aligned} y_{n,h}^i &\in \text{Arg min}_{y \in Y} \left(\theta \phi_h(y) - \langle \ell_h(t_{n,\theta}^i) - (1 - \theta)A_h y_{n,h}^{i-1}, y \rangle + \psi_h(y - y_{n,h}^{i-1}) \right) \\ &= \text{Arg min}_{y \in Y_h} \left(\theta \phi(y) - \langle \ell_h(t_{n,\theta}^i) - (1 - \theta)A y_{n,h}^{i-1}, y \rangle + \psi_h(y - y_{n,h}^{i-1}) \right). \end{aligned}$$

Hence, since we have (5.6), the sequence $y_{n,h}^i$ is weakly precompact, and, up to the extraction of a (not relabeled) subsequence, $y_{n,h}^i \rightarrow \tilde{y}$ weakly in Y . Let us prove that \tilde{y} solves the incremental problem (6.4). Indeed, we have

$$\begin{aligned} 0 &\leq L_n^{\theta,i}(\tilde{y}, y_n^{i-1}) \\ &\leq \liminf_{h \rightarrow 0} \left(\psi_h(y_{n,h}^i - y_{n,h}^{i-1}) + \psi_h^*(\ell_h(t_{n,\theta}^i) - A y_{n,h}^i) \right. \\ &\quad \left. - \langle \ell_h(t_{n,\theta}^i) - A y_{n,h}^i, y_{n,h}^i - y_{n,h}^{i-1} \rangle \right) \\ &= \liminf_{h \rightarrow 0} L_{n,h}^{\theta,i}(y_{n,h}^i, y_{n,h}^{i-1}) = 0, \end{aligned}$$

where we have used the Mosco convergence in (5.5) and the pointwise convergence of ℓ_h (5.7). Since the only solution of (6.4) is y_n^i , we have $\tilde{y} = y_n^i$, and the whole sequence converges.

Let us mention that, if the functionals ψ_h are uniformly linearly bounded (which is quite common in practice), one could prove the latter convergence to be actually strong: Namely, $y_{n,h}^{i-1} \rightarrow y_n^{i-1}$ strongly in Y implies the strong convergence $y_{n,h}^i \rightarrow y_n^i$. Indeed, let w_h and \tilde{w}_h be such that $w_h - y_{n,h}^i \rightarrow 0$ strongly in Y , $\psi_h(w_h - y_{n,h}^i) \rightarrow 0$,

$\tilde{w}_h \in Y_h$ and $\tilde{w}_h - w_h \rightarrow 0$ strongly in Y . Then

$$\begin{aligned} & \theta \phi(y_{n,h}^i) - \langle \ell(t_{n,\theta}^i) - (1 - \theta)A_h y_{n,h}^{i-1}, y_{n,h}^i \rangle + \psi_h(y_{n,h}^i - y_{n,h}^{i-1}) \\ & \leq \theta \phi_h(\tilde{w}_h) - \langle \ell(t_{n,\theta}^i) - (1 - \theta)A y_{n,h}^{i-1}, \tilde{w}_h \rangle + \psi_h(\tilde{w}_h - w_h) + \psi_h(w_h - y_{n,h}^i). \end{aligned}$$

If ψ_h are uniformly linearly bounded above, then $\psi_h(\tilde{w}_h - w_h) \rightarrow 0$, with $h \rightarrow 0$. Then, by passing to the limsup in the latter, we check that $\limsup_{h \rightarrow 0} \phi(y_{n,h}^i) \leq \phi(y_n^i)$, which together with lower semicontinuity gives $\phi(y_{n,h}^i) \rightarrow \phi(y_n^i)$, and the strong convergence follows from the reflexivity of Y .

Limit (e). Lemma 6.4, the uniform Lipschitz continuity of ℓ_h (5.8), and the initial datum convergence (5.9) entail that $\hat{y}_{n,h}$ are uniformly Lipschitz continuous as well. Hence, by extracting a (not relabeled) subsequence, $\hat{y}_{n,h} \rightarrow y$ weakly star in $W^{1,\infty}(0, T; Y)$. In order to check that y solves (1.1), let us remark that, since $\ell_{n,h,\theta}^i = \ell_h(t_{n,\theta}^i)$,

$$\bar{\ell}_{n,h,\theta} \rightarrow \ell \quad \text{strongly in } L^1(0, T; Y^*)$$

and that, by [45, Cor. 4.4],

$$\begin{aligned} \int_0^T \psi(y) & \leq \liminf_{h \rightarrow 0} \int_0^T \psi_h(\hat{y}_{n,h}), \\ \int_0^T \psi^*(\ell - Ay) & \leq \liminf_{h \rightarrow 0} \int_0^T \psi_h^*(\bar{\ell}_{n,h,\theta} - A\bar{y}_{n,h,\theta}) \end{aligned}$$

and compute that

$$\begin{aligned} 0 & \leq F(y) \\ & \leq \liminf_{h \rightarrow 0} \left(\int_0^T (\psi_h(\hat{y}_{n,h}) + \psi_h^*(\bar{\ell}_{n,h,\theta} - A\bar{y}_{n,h,\theta}) - \langle \bar{\ell}_{n,h,\theta}, \hat{y}_{n,h} \rangle) \right. \\ & \quad \left. + \phi(\hat{y}_{n,h}(T)) + \phi(y_{0,h}) - \langle A\hat{y}_{n,h}(0), y_{0,h} \rangle + |\hat{y}_{n,h}(0) - y_{0,h}|^2 \right) \\ & \leq \liminf_{h \rightarrow 0} F_{n,h}^\theta(y_{n,h}^0, \dots, y_{n,h}^{N_n}) = 0, \end{aligned}$$

and we have $F(y) = 0$. \square

We shall conclude by briefly mentioning some further results which can be obtained by suitably adapting to the current fully discretized situation the arguments developed above for time discretizations. First, in the same spirit of Lemma 6.6, one could consider the possibility of estimating the distance of a vector from the minimizer of $F_{n,h}^\theta$ by means of the functional itself. Second, the use of Corollary 4.5 would entail the possibility of an a posteriori error control, and some adaptive strategy along the lines of subsection 6.7 could be considered. Finally, by relying on the known convergence estimates for full space-time discretized problems [13], one could obtain a convergence and an a posteriori error control result for some generalized space-time approximated problem where $F_{n,h}^\theta$ are not exactly minimized and one considers minimizing sequences instead (see subsection 6.5). We shall develop these considerations elsewhere.

REFERENCES

- [1] H. ATTOUCH, *Variational Convergence for Functions and Operators*, Pitman, Boston, 1984.
- [2] G. AUCHMUTY, *Saddle-points and existence-uniqueness for evolution equations*, *Differential Integral Equations*, 6 (1993), pp. 1161–1171.
- [3] H. BREZIS AND I. EKELAND, *Un principe variationnel associé à certaines équations paraboliques. Le cas dépendant du temps*, *C. R. Acad. Sci. Paris Sér. A-B*, 282 (1976), pp. Ai, A1197–A1198.
- [4] H. BREZIS AND I. EKELAND, *Un principe variationnel associé à certaines équations paraboliques. Le cas indépendant du temps*, *C. R. Acad. Sci. Paris Sér. A-B*, 282 (1976), pp. Aii, A971–A974.
- [5] G. DAL MASO, *An Introduction to Γ -Convergence*, *Progr. Nonlinear Differential Equations Appl.* 8, Birkhäuser Boston, Cambridge, MA, 1993.
- [6] G. DAL MASO, A. DESIMONE, AND M. G. MORA, *Quasistatic evolution problems for linearly elastic-perfectly plastic materials*, *Arch. Ration. Mech. Anal.*, 180 (2006), pp. 237–291.
- [7] G. DUVAUT AND J.-L. LIONS, *Inequalities in Mechanics and Physics*, Springer-Verlag, Berlin, 1976.
- [8] M. A. EFENDIEV AND A. MIELKE, *On the rate-independent limit of systems with dry friction and small viscosity*, *J. Convex Anal.*, 13 (2006), pp. 151–167.
- [9] N. GHOUSSOUB, *Selfdual Partial Differential Systems and Their Variational Principles*, Universitext, Springer-Verlag, New York, to appear.
- [10] N. GHOUSSOUB AND L. TZOU, *A variational principle for gradient flows*, *Math. Ann.*, 330 (2004), pp. 519–549.
- [11] E. DE GIORGI AND T. FRANZONI, *Su un tipo di convergenza variazionale*, *Atti Accad. Naz. Lincei Rend. Cl. Sci. Fis. Mat. Natur.*, 58 (1975), pp. 842–850.
- [12] E. HAIRER, CH. LUBICH, AND G. WANNER, *Geometric Numerical Integration*, *Springer Ser. Comput. Math.* 31, 2nd ed., Springer-Verlag, Berlin, 2006.
- [13] W. HAN AND B. D. REDDY, *Plasticity, Mathematical Theory and Numerical Analysis*, Springer-Verlag, New York, 1999.
- [14] W. HAN AND B. D. REDDY, *Convergence of approximations to the primal problem in plasticity under conditions of minimal regularity*, *Numer. Math.*, 87 (2000), pp. 283–315.
- [15] W. HAN AND B. D. REDDY, *Computational plasticity: The variational basis and numerical analysis*, *Comput. Mech. Adv.*, 2 (1995), pp. 283–400.
- [16] P. KREJIĆ, *Hysteresis, Convexity and Dissipation in Hyperbolic Equations*, *GAKUTO Internat. Ser. Math. Sci. Appl.* 8, Gakkotosho, Tokyo, 1996.
- [17] B. LEMAIRE, *An asymptotical variational principle associated with the steepest descent method for a convex function*, *J. Convex Anal.*, 3 (1996), pp. 63–70.
- [18] J. LEMAITRE AND J.-L. CHABOCHE, *Mechanics of Solid Materials*, Cambridge University Press, London, 1990.
- [19] A. LEW, J. E. MARSDEN, M. ORTIZ, AND M. WEST, *Asynchronous variational integrators*, *Arch. Ration. Mech. Anal.*, 167 (2003), pp. 85–146.
- [20] A. LEW, J. E. MARSDEN, M. ORTIZ, AND M. WEST, *Variational time integrators*, *Internat. J. Numer. Methods Engrg.*, 60 (2004), pp. 153–212.
- [21] M. MABROUK, *Un principe variationnel pour une équation non linéaire du second ordre en temps*, *C. R. Acad. Sci. Paris Sér. I Math.*, 332 (2001), pp. 381–386.
- [22] M. MABROUK, *A variational principle for a nonlinear differential equation of second order*, *Adv. Appl. Math.*, 31 (2003), pp. 388–419.
- [23] F. MAGGI AND M. MORINI, *A Γ -convergence result for variational integrators of Lagrangians with quadratic growth*, *ESAIM Control Optim. Calc. Var.*, 10 (2004), pp. 656–665.
- [24] J. E. MARSDEN, G. W. PATRICK, AND S. SHKOLLER, *Multisymplectic geometry, variational integrators, and nonlinear PDEs*, *Comm. Math. Phys.*, 199 (1998), pp. 351–395.
- [25] J. E. MARSDEN AND M. WEST, *Discrete mechanics and variational integrators*, *Acta Numer.*, 10 (2001), pp. 357–514.
- [26] A. MIELKE, *Evolution of rate-independent systems*, in *Handbook of Differential Equations, Evolutionary Equations*, Vol. 2, C. Dafermos and E. Feireisl, eds., Elsevier, New York, 2005, pp. 461–559.
- [27] A. MIELKE AND M. ORTIZ, *A class of minimum principles for characterizing the trajectories and the relaxation of dissipative systems*, *ESAIM Control Optim. Calc. Var.*, 2008, to appear; also available online from <http://www.wias-berlin.de/main/publications/wias-publ/index.cgi.en>.
- [28] A. MIELKE, R. ROSSI, AND G. SAVARÉ, manuscript, 2008.

- [29] A. MIELKE, R. ROSSI, AND G. SAVARÉ, *A metric approach to a class of doubly nonlinear evolution equations and applications*, Ann. Sc. Norm. Super. Pisa Cl. Sci., 5 (2008), pp. 97–169; also available online from <http://www.wias-berlin.de/main/publications/wias-publ/index.cgi.en>.
- [30] A. MIELKE, T. ROUBÍČEK, AND U. STEFANELLI, *Γ -limits and relaxations for rate-independent evolutionary problems*, Calc. Var. Partial Differential Equations, 31 (2008), pp. 387–416.
- [31] A. MIELKE AND U. STEFANELLI, *A Discrete Variational Principle for Rate-Independent Evolution*, Wias preprint 1295, 2008, also available online from <http://www.wias-berlin.de/main/publications/wias-publ/index.cgi.en>.
- [32] A. MIELKE AND F. THEIL, *On rate-independent hysteresis models*, NoDEA Nonlinear Differential Equations Appl., 11 (2004), pp. 151–189.
- [33] J.-J. MOREAU, *La notion de sur-potentielle et les liaisons unilatérales en élastostatique*, C. R. Acad. Sci. Paris Sér. A-B, 267 (1968), pp. A954–A957.
- [34] J.-J. MOREAU, *Sur les lois de frottement, de viscosité et plasticité*, C. R. Acad. Sci. Paris Sér. II Méc. Phys. Chim. Sci. Univers Sci. Terre, 271 (1970), pp. 608–611.
- [35] J.-J. MOREAU, *Sur l'évolution d'un système élasto-visco-plastique*, C. R. Acad. Sci. Paris Sér. A-B, 273 (1971), pp. A118–A121.
- [36] U. MOSCO, *Convergence of convex sets and of solutions of variational inequalities*, Adv. Math., 3 (1969), pp. 510–585.
- [37] S. MÜLLER AND M. ORTIZ, *On the Γ -convergence of discrete dynamics and variational integrators*, J. Nonlinear Sci., 14 (2004), pp. 279–296.
- [38] B. NAYROLES, *Deux théorèmes de minimum pour certains systèmes dissipatifs*, C. R. Acad. Sci. Paris Sér. A-B, 282 (1976), pp. Aiv, A1035–A1038.
- [39] B. NAYROLES, *Un théorème de minimum pour certains systèmes dissipatifs. Variante hilbertienne*, Travaux Sém. Anal. Convexe, 6 (1976), p. 22.
- [40] H. RIOS, *Étude de la question d'existence pour certains problèmes d'évolution par minimisation d'une fonctionnelle convexe*, C. R. Acad. Sci. Paris Sér. A-B, 283 (1976), pp. Ai, A83–A86.
- [41] T. ROUBÍČEK, *Direct method for parabolic problems*, Adv. Math. Sci. Appl., 10 (2000), pp. 57–65.
- [42] J. C. SIMO AND T.J.R. HUGHES, *Computational Inelasticity*, Interdiscip. Appl. Math. 7, Springer-Verlag, New York, 1998.
- [43] U. STEFANELLI, *The Discrete Brezis-Ekeland Principle*, J. Convex. Anal., to appear, also available online from <http://www.imati.cnr.it/ulisse/pubbl.html>.
- [44] U. STEFANELLI, *A variational principle in non-smooth mechanics*, in MFO Workshop: Analysis and Numerics for Rate-Independent Processes, Oberwolfach Rep. 4, European Mathematics Society, Zurich, 2007, pp. 622–624.
- [45] U. STEFANELLI, *The Brezis-Ekeland principle for doubly nonlinear equations*, SIAM J. Control Optim., to appear; also available online from <http://www.imati.cnr.it/ulisse/pubbl.html>.
- [46] U. STEFANELLI, *Some remarks on convergence and approximation for a class of hysteresis problems*, Istit. Lombardo Accad. Sci. Lett. Rend. A, to appear; also available online from <http://www.imati.cnr.it/ulisse/pubbl.html>.
- [47] R. TOADER AND C. ZANINI, *An artificial viscosity approach to quasistatic crack growth*, submitted; also available online from <http://cvgmt.sns.it/>.
- [48] A. VISINTIN, *A new approach to evolution*, C. R. Acad. Sci. Paris Sér. I Math., 332 (2001), pp. 233–238.
- [49] C. ZANINI, *Singular perturbations of finite dimensional gradient flows*, Discrete Contin. Dyn. Syst., 18 (2007), pp. 657–675.

AVERAGES OVER SPHERES FOR KINETIC TRANSPORT EQUATIONS WITH VELOCITY DERIVATIVES IN THE RIGHT-HAND SIDE*

NIKOLAOS BOURNAVEAS[†] AND SUSANA GUTIÉRREZ[‡]

Abstract. We prove estimates in hyperbolic Sobolev spaces $H^{s,\delta}(R^{1+d})$, $d \geq 3$, for velocity averages over spheres of solutions to the kinetic transport equation $\partial_t f + v \cdot \nabla_x f = \Omega_v^{i,j} g$, where $\Omega_v^{i,j} g$ are tangential velocity derivatives of g . Our results extend to all dimensions earlier results of Bournaveas and Perthame in dimension two [*J. Math. Pures Appl.*, 9 (2001), pp. 517–534]. We construct counterexamples to test the optimality of our results.

Key words. velocity-averaging lemmas, kinetic transport equation, hyperbolic Sobolev spaces

AMS subject classifications. 82C70, 35B45, 35F10

DOI. 10.1137/070698415

1. Introduction. Consider the kinetic transport equation

$$\partial_t f + v \cdot \nabla_x f = g,$$

where $f, g : R_t \times R_x^d \times R_v^d \rightarrow R$. It is well known that averaging with respect to the velocity variable v has a smoothing effect which can be measured in classical Sobolev spaces. For example, in all dimensions, the average over the unit ball

$$\rho_B(t, x) = \int_{|v| \leq 1} f(t, x, v) dv$$

is smoother than f by $1/2$ derivatives in L^2 , and we have the estimate

$$\|\rho_B\|_{H^{1/2}(R^{1+d})} \leq C \left(\|f\|_{L^2(R^{1+2d})} + \|g\|_{L^2(R^{1+2d})} \right).$$

Results of this type are known as averaging lemmas. They were first discovered in [24, 23] and developed further by many authors [4, 5, 9, 11, 12, 13, 14, 15, 17, 19, 20, 21, 25, 26, 27, 28, 29, 33, 34, 37]. We refer the reader to [3, 35] for a review and an extensive bibliography.

Averages over spheres

$$\rho_S(t, x) = \int_{S^{d-1}} f(t, x, v) d\sigma(v)$$

were first studied in [7] and later in [8]. These averages appear in the equation of radiative transfer, a phenomenon that describes the scattering of photons in a hot medium [1, 2, 10, 22, 36]. They also appear in certain kinetic models of chemotaxis when the velocity of the cells is normalized to $|v| = 1$ [16]. It turns out that in

*Received by the editors July 26, 2007; accepted for publication (in revised form) February 26, 2008; published electronically July 3, 2008.

<http://www.siam.org/journals/sima/40-2/69841.html>

[†]School of Mathematics, University of Edinburgh, Mayfield Road, Edinburgh EH9 3JZ, UK (n.bournaveas@ed.ac.uk).

[‡]School of Mathematics, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK (s.gutierrez@bham.ac.uk).

dimensions $d \geq 3$ averages over spheres gain $1/2$ derivatives in L^2 [7, Theorem 1] and we have the estimate

$$(1.1) \quad \|\rho_s\|_{H^{1/2}(R^{1+d})} \leq C \left(\|f\|_{L^2(R \times R^d \times S^{d-1})} + \|g\|_{L^2(R \times R^d \times S^{d-1})} \right).$$

Thus, if $d \geq 3$, ρ_s gains the same amount of regularity as averages over balls, although it is an average over a lower-dimensional set. Moreover, the estimate for spheres implies the estimate for balls (see Remark 1 in [7]).

In dimension $d = 2$ averages over spheres gain only $1/4$ derivatives, but the “missing” regularity can be recovered in the so-called hyperbolic Sobolev spaces $H^{s,\delta}(R^{1+d})$ [6, 30, 31, 32, 38]. More precisely [7, Theorem 2],

$$\|\rho_s\|_{H^{1/4,1/4}(R^{1+2})} \leq C \left(\|f\|_{L^2(R \times R^2 \times S^1)} + \|g\|_{L^2(R \times R^2 \times S^1)} \right),$$

where

$$(1.2) \quad \|F\|_{H^{s,\delta}(R^{1+d})} = \left\| w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \widehat{F}(\tau, \xi) \right\|_{L^2(R_\tau \times R_\xi^d)},$$

$w_\pm(\tau, \xi) = 1 + |\tau| \pm |\xi|$, and \widehat{F} is the space-time Fourier transform of F .

In [8] these spaces were used to further improve estimate (1.1). We have [8, Theorem 1], when $d \geq 3$, $d \neq 7$,

$$(1.3a) \quad \|\rho_s\|_{H^{s,\delta}(R^{1+d})} \leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})},$$

provided that

$$(1.3b) \quad s + \delta \leq 1/2, \quad s \leq \min \left\{ \frac{d-1}{4}, 1 \right\}.$$

This allows us to take s larger than the classical $1/2$, provided that we compensate by using a negative δ . Moreover, in regions where $||\tau| - |\xi||$ is small (near the cone $|\tau| = |\xi|$) we have $w_+^s w_-^\delta \simeq w_+^s$, and we gain $s \geq 1/2$ derivatives. When $d = 7$ we have the same estimate but with a logarithmic loss:

$$\|\rho_s\|_{H_{\log}^{s,\delta}(R^{1+d})} \leq C \left(\|f\|_{L^2(R \times R^d \times S^{d-1})} + \|g\|_{L^2(R \times R^d \times S^{d-1})} \right), \quad d = 7,$$

where

$$\|F\|_{H_{\log}^{s,\delta}(R^{1+d})} = \left\| \frac{w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \widehat{F}(\tau, \xi)}{\left(1 + \log \frac{w_+(\tau, \xi)}{w_-(\tau, \xi)}\right)^{1/2}} \right\|_{L^2(R_\tau \times R_\xi^d)}.$$

Now we turn our attention to equations with right-hand sides containing derivatives with respect to the velocity variable:

$$\partial_t f + v \cdot \nabla_x f = \partial_v^m g.$$

This case is of great interest in applications. For example, the Vlasov part of the Vlasov–Maxwell system has this structure with $m = 1$. It is well known (see [24, 23]) that in this case averages of the form $\int_{R^d} f(t, x, v) \phi(v) dv$, where $\phi(v)$ is a smooth cut-off function, gain $\frac{1}{2(m+1)}$ derivatives in L^2 . In particular, if $m = 1$, then the gain is $1/4$ derivatives. The proof of these results uses an integration by parts which removes the velocity derivatives from the function g .

Consider now the case of spheres. In order to be able to integrate by parts on spheres we are forced to restrict ourselves to equations with tangential v -derivatives in the right-hand side:

$$(1.4) \quad \partial_t f + v \cdot \nabla_x f = \Omega_v^{i,j} g,$$

where

$$\Omega_v^{i,j} g = v_i \frac{\partial g}{\partial v_j} - v_j \frac{\partial g}{\partial v_i}, \quad 1 \leq i, j \leq d.$$

In this case, if $d \geq 3$, the optimal gain is $1/4$ derivatives (the same gain as for balls) [7, Theorem 3]. In view of the discussion above, one would expect the gain in $d = 2$ dimensions to be $1/8$ derivatives. However, the tangential derivatives introduce a special structure (see, for example, (2.3) below and the discussion in [8]) reminiscent of the Klainerman null forms structure for the wave equation [30, 31, 32], which allows for better estimates. By using this observation the expected classical gain of $1/8$ derivatives was improved in [7] to $1/7$ derivatives and further improved in [8] to $1/6$ derivatives, and it was also shown that this result is optimal.

If one considers the hyperbolic Sobolev spaces, it is also possible to prove estimates with a total of $s + \delta = 1/4$ derivatives. In this setting, the two-dimensional case has been previously considered in [7, 8]. More precisely, in [7] it was shown that $\rho_s \in H^{1/16, 3/16}(R^{1+2})$. This was improved in [8] to $\rho_s \in H^{1/8, 1/8}(R^{1+2})$, and it was also shown that this is optimal.

In this paper we study the regularity of velocity averages over spheres for solutions of (1.4) in hyperbolic Sobolev spaces in dimensions $d \geq 3$. Recall first the following result from [7].

Let $d \geq 3$, and let $f \in L^2(R_t \times R_x^d \times S_v^{d-1})$ be a solution of (1.4) with $g \in L^2(R_t \times R_x^d \times S_v^{d-1})$. Then the velocity average over the unit sphere $\rho_s(t, x) = \int_{S^{d-1}} f(t, x, v) d\sigma(v)$ satisfies

$$(1.5) \quad \|\rho_s\|_{H^{1/4}(R \times R^d)} \leq C \left[\|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})} \right].$$

We shall improve this result by showing that $\rho_s \in H^{s, \delta}$ for certain choices of (s, δ) depending on the dimension d . More precisely we have the following.

THEOREM 1.1. *Let $d \geq 3$, and let $f \in L^2(R_t \times R_x^d \times S_v^{d-1})$ be a solution of (1.4) with $g \in L^2(R_t \times R_x^d \times S_v^{d-1})$. Define*

$$(1.6) \quad \begin{aligned} (s, \delta) &= \left(\frac{d-1}{6}, -\frac{2d-5}{12} \right) \quad \text{if } d \in \{3, 4, 5\}, \\ (s, \delta) &= \left(\frac{4}{5}, -\frac{11}{20} \right) \quad \text{if } d = 6, \\ (s, \delta) &= \left(1, -\frac{3}{4} \right), \quad \text{if } d \geq 7. \end{aligned}$$

Then the velocity average $\rho_s(t, x) = \int_{S^{d-1}} f(t, x, v) d\sigma(v)$ satisfies

$$(1.7a) \quad \|\rho_s\|_{H^{s, \delta}(R^{1+d})} \leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})}, \quad d \neq 5, 7,$$

$$(1.7b) \quad \|\rho_s\|_{H_{\text{log}}^{s, \delta}(R^{1+d})} \leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})}, \quad d = 5, 7.$$

Notice that in all cases we have $s + \delta = 1/4$, with $s > 1/4$ and $\delta < 0$. By comparing this result to the classical estimate (1.5) we see that it is possible to improve on the

“good derivative” w_+ , provided that we make a sacrifice in the derivative w_- . Observe also that since $\delta = 1/4 - s$ and $s > 1/4$ we have

$$w_+^s w_-^\delta = w_+^{\frac{1}{4}} \left(\frac{w_+}{w_-} \right)^{s-\frac{1}{4}} \geq w_+^{\frac{1}{4}},$$

and therefore estimate (1.7) implies estimate (1.5). Moreover, near the cone $|\tau| = |\xi|$ we have $w_+^s w_-^\delta \simeq w_+^s$, because $w_- \simeq 1$, and therefore our estimate says that we gain $s > 1/4$ derivatives in this region. Notice also that, in contrast to (1.3), we gain even in $d = 3$ dimensions. This is due to the fact that we have not only the usual averaging smoothing effect and the flexibility of the $H^{s,\delta}$ -spaces but also the special structure of the tangential derivatives.

It is possible to extend the validity of estimate (1.7) to a whole region of pairs (s, δ) . We have the following.

THEOREM 1.2. *Let $d \geq 3$, and let $f \in L^2(R_t \times R_x^d \times S_v^{d-1})$ be a solution of (1.4) with $g \in L^2(R_t \times R_x^d \times S_v^{d-1})$. Let $s, \delta \in R$ be such that*

- (1.8a) $s + \delta \leq \frac{1}{4},$
- (1.8b) $s \leq \frac{d-1}{6} \quad \text{if } d \in \{3, 4, 5\},$
- (1.8c) $s \leq \frac{4}{5} \quad \text{if } d = 6,$
- (1.8d) $s \leq 1 \quad \text{if } d \geq 7.$

Then estimates (1.7) are true.

In the region $\{|\tau| > |\xi|\}$ in phase space the symbol of the operator $\partial_t + v \cdot \nabla_x$ satisfies $|\tau + v \cdot \xi| \geq |\tau| - |\xi| > 0$, and therefore it is reasonable to expect better estimates. Indeed, improved estimates in this region are contained in [7, 8]. We shall establish the following improvements to the results of Theorem 1.1. First, estimates (1.7) are valid in $\{|\tau| > |\xi|\}$ with the same s and with $\delta = 0$ (instead of a negative δ).

THEOREM 1.3. *Let $d \geq 3$, and let $f \in L^2(R_t \times R_x^d \times S_v^{d-1})$ be a solution of (1.4) with $g \in L^2(R_t \times R_x^d \times S_v^{d-1})$. Define*

- (1.9a) $s = \frac{d-1}{6} \quad \text{if } d \in \{3, 4, 5\},$
- (1.9b) $s = \frac{4}{5} \quad \text{if } d = 6,$
- (1.9c) $s = 1 \quad \text{if } d \geq 7.$

Then the velocity average $\rho_s(t, x) = \int_{S^{d-1}} f(t, x, v) d\sigma(v)$ satisfies the following estimates in the classical Sobolev spaces $H^s(\{|\tau| > |\xi|\})$.

If $d \neq 5, 7$,

- (1.10a) $\|w_+(\tau, \xi)^s \widehat{\rho}_s(\tau, \xi)\|_{L^2(\{|\tau| > |\xi|\})}$
- (1.10b) $\leq C \left(\|f\|_{L^2(R \times R^d \times S^{d-1})} + \|g\|_{L^2(R \times R^d \times S^{d-1})} \right).$

If $d = 5, 7$,

- (1.10c) $\left\| w_+(\tau, \xi)^s \left[1 + \log \frac{w_+(\tau, \xi)}{w_-(\tau, \xi)} \right]^{-1/2} \widehat{\rho}_s(\tau, \xi) \right\|_{L^2(\{|\tau| > |\xi|\})}$
- (1.10d) $\leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})}.$

Moreover, it is possible to prove $H^{s,\delta}$ -estimates in all dimensions $d \geq 3$ with $s + \delta = 1$ and $s, \delta \geq 0$. Precisely, we have the following.

THEOREM 1.4. *Under the same assumptions and notation as in Theorem 1.3 and with*

$$(1.11a) \quad (s, \delta) = \left(\frac{d-3}{4}, \frac{7-d}{4} \right) \quad \text{if } d \in \{3, 4, 5\},$$

$$(1.11b) \quad (s, \delta) = \left(\frac{3}{4}, \frac{1}{4} \right) \quad \text{if } d = 6,$$

$$(1.11c) \quad (s, \delta) = (1, 0) \quad \text{if } d \geq 7,$$

we have

$$(1.12a) \quad \|w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \widehat{\rho}_s(\tau, \xi)\|_{L^2(\{|\tau| > |\xi|\})}$$

$$(1.12b) \quad \leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})} \quad \text{if } d \neq 5, 7;$$

$$(1.12c) \quad \left\| w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \left[1 + \log \frac{w_+(\tau, \xi)}{w_-(\tau, \xi)} \right]^{-1/2} \widehat{\rho}_s(\tau, \xi) \right\|_{L^2(\{|\tau| > |\xi|\})}$$

$$(1.12d) \quad \leq C \|f\|_{L^2(R \times R^d \times S^{d-1})} + C \|g\|_{L^2(R \times R^d \times S^{d-1})} \quad \text{if } d = 5, 7.$$

Remark. Regarding the optimality of the conditions on s and δ we will show in the last section that all of the upper bounds for s and $s + \delta$ are optimal with the possible exception of the bound $s \leq 4/5$ in dimension $d = 6$. In some cases it is possible to have improved estimates if we allow weights of the form $w_+^{s_1} w_-^{\delta_1} + w_+^{s_2} w_-^{\delta_2}$; see page 664 for details. The logarithmic divergence in some of the estimates is due to the logarithmic divergence in Lemmas 3.2 and 3.3 in [8].

Notation. We define the weights w_\pm by $w_\pm(\tau, \xi) = 1 + \|\tau\| \pm \|\xi\|$. We denote the classical Sobolev spaces by H^s and the hyperbolic Sobolev spaces (see (1.2)) by $H^{s,\delta}$. We use $\widehat{\cdot}$ for the Fourier transform in space-time denoting the dual variables by (τ, ξ) . We will always average the solution f over the unit sphere S^{d-1} in R^d and denote the average by ρ_s .

2. Proofs of Theorems 1.1 and 1.2. In this section we prove Theorems 1.1 and 1.2. We need the following result from [8].

PROPOSITION 2.1. *Let $m > -1$, $l > 1/2$, and define*

$$(2.1) \quad J_l^m(\tau, \xi) = \int_0^\pi \frac{(\sin \theta)^m}{[1 + (\tau + |\xi| \cos \theta)^2]^l} d\theta, \quad \tau \in R, \xi \in R^d.$$

Let $\alpha = \min \{ \frac{m+1-4l}{2}, 0 \}$. Then the integrals $J_l^m(\tau, \xi)$ satisfy the following pointwise estimates:

$$J_l^m(\tau, \xi) \leq C \frac{w_-(\tau, \xi)^{2l-1+\alpha}}{w_+(\tau, \xi)^{2l+\alpha}} \quad \text{if } m + 1 \neq 4l,$$

$$J_l^m(\tau, \xi) \leq C \frac{w_-(\tau, \xi)^{2l-1}}{w_+(\tau, \xi)^{2l}} \left(1 + \log \frac{w_+(\tau, \xi)}{w_-(\tau, \xi)} \right) \quad \text{if } m + 1 = 4l.$$

Proof of Theorem 1.1. By taking the space-time Fourier transform of (1.4) and adding $\lambda \widehat{f}(\tau, \xi, v)$ to both sides, where $\lambda = \lambda(\tau, \xi) > 0$ will be determined later,

we get

$$\widehat{f}(\tau, \xi, v) = \frac{\Omega_v^{i,j} \widehat{g}(\tau, \xi, v)}{\lambda + i(\tau + v \cdot \xi)} + \lambda \frac{\widehat{f}(\tau, \xi, v)}{\lambda + i(\tau + v \cdot \xi)}.$$

By averaging in v and integrating by parts we obtain

$$\begin{aligned} \widehat{\rho}(\tau, \xi) &= \int_{S^{d-1}} \frac{\Omega_v^{i,j} \widehat{g}(\tau, \xi, v)}{\lambda + i(\tau + v \cdot \xi)} d\sigma(v) + \lambda \int_{S^{d-1}} \frac{\widehat{f}(\tau, \xi, v)}{\lambda + i(\tau + v \cdot \xi)} d\sigma(v) \\ &= i \int_{S^{d-1}} \frac{v_i \xi_j - v_j \xi_i}{[\lambda + i(\tau + v \cdot \xi)]^2} \widehat{g}(\tau, \xi, v) d\sigma(v) + \lambda \int_{S^{d-1}} \frac{\widehat{f}(\tau, \xi, v) d\sigma(v)}{\lambda + i(\tau + v \cdot \xi)}. \end{aligned}$$

By using the Cauchy–Schwarz inequality we get

$$(2.2) \quad \begin{aligned} |\widehat{\rho}(\tau, \xi)| &\leq \left(\int_{S^{d-1}} |\widehat{g}(\tau, \xi, v)|^2 d\sigma(v) \right)^{1/2} \cdot K(\tau, \xi) \\ &\quad + \left(\int_{S^{d-1}} |\widehat{f}(\tau, \xi, v)|^2 d\sigma(v) \right)^{1/2} \cdot L(\tau, \xi), \end{aligned}$$

where

$$(2.3) \quad \begin{aligned} K(\tau, \xi) &= \left(\int_{S^{d-1}} \frac{|v_i \xi_j - v_j \xi_i|^2}{|\lambda + i(\tau + v \cdot \xi)|^4} d\sigma(v) \right)^{1/2}, \\ L(\tau, \xi) &= \left(\lambda^2 \int_{S^{d-1}} \frac{d\sigma(v)}{|\lambda + i(\tau + v \cdot \xi)|^2} \right)^{1/2}. \end{aligned}$$

We estimate $K(\tau, \xi)$ and $L(\tau, \xi)$ as follows. Set $\tau' = \frac{\tau}{\lambda}$, $\xi' = \frac{\xi}{\lambda}$, and let θ be the angle between ξ and v . Then

$$(2.4) \quad \begin{aligned} K(\tau, \xi) &= \left(\int_{S^{d-1}} \frac{|v_i \xi_j - v_j \xi_i|^2}{|\lambda + i(\tau + v \cdot \xi)|^4} d\sigma(v) \right)^{1/2} \\ &= \left(\int_{S^{d-1}} \frac{|v_i \xi_j - v_j \xi_i|^2}{|\lambda^2 + (\tau + v \cdot \xi)^2|^2} d\sigma(v) \right)^{1/2} \\ &= \frac{|\xi|}{\lambda^2} \left(\int_{S^{d-1}} \frac{|v_i \frac{\xi_j}{|\xi|} - v_j \frac{\xi_i}{|\xi|}|^2}{|1 + (\tau' + v \cdot \xi')|^2} d\sigma(v) \right)^{1/2} \\ &\simeq \frac{|\xi|}{\lambda^2} \left(\int_0^\pi \frac{(\sin \theta)^2 (\sin \theta)^{d-2}}{|1 + (\tau' + |\xi'| \cos \theta)|^2} d\theta \right)^{1/2} \\ &= \frac{|\xi|}{\lambda^2} J_2^d(\tau', \xi')^{1/2}, \end{aligned}$$

and

$$\begin{aligned}
 L(\tau, \xi) &= \left(\lambda^2 \int_{S^{d-1}} \frac{d\sigma(v)}{\lambda^2 + (\tau + v \cdot \xi)^2} \right)^{1/2} \\
 &= \left(\int_{S^{d-1}} \frac{d\sigma(v)}{1 + (\tau' + v \cdot \xi')^2} \right)^{1/2} \\
 &\simeq \left(\int_0^\pi \frac{(\sin \theta)^{d-2}}{1 + (\tau' + |\xi'| \cos \theta)^2} d\theta \right)^{1/2} \\
 (2.5) \quad &= J_1^{d-2}(\tau', \xi')^{1/2}.
 \end{aligned}$$

We start with the case $d \in \{3, 4, 5\}$. From (2.4) we have $K(\tau, \xi) \leq C \frac{|\xi|}{\lambda^{\frac{3}{2}}} J_2^d(\tau', \xi')^{1/2}$. Proposition 2.1 with $m = d, l = 2$ gives

$$\begin{aligned}
 K(\tau, \xi) &\leq C \frac{|\xi|}{\lambda^2} \frac{(1 + \|\tau' - |\xi'|\|)^{\frac{d-1}{4}}}{(1 + |\tau'| + |\xi'|)^{\frac{d+1}{4}}} \\
 &= C \frac{|\xi|}{\lambda^{\frac{3}{2}}} \frac{(\lambda + \|\tau| - |\xi||)^{\frac{d-1}{4}}}{(\lambda + |\tau| + |\xi|)^{\frac{d+1}{4}}} \\
 &\leq C \frac{(\lambda + \|\tau| - |\xi||)^{\frac{d-1}{4}}}{\lambda^{\frac{3}{2}} (\lambda + |\tau| + |\xi|)^{\frac{d-3}{4}}} \\
 &\leq C \frac{(\lambda + w_-)^{\frac{d-1}{4}}}{\lambda^{\frac{3}{2}} (w_+)^{\frac{d-3}{4}}},
 \end{aligned}$$

where we have set $w_\pm = 1 + \|\tau| \pm |\xi||$ and we have used the fact that $\lambda > 1$ (this will be guaranteed by the choice of λ ; see (2.9)).

Next we deal with $L(\tau, \xi)$. From estimate (2.5) we have $L(\tau, \xi) \leq J_1^{d-2}(\tau', \xi')^{1/2}$, and we now apply Proposition 2.1 with $m = d - 2$ and $l = 1$. Notice that $m + 1 = 4l$ when $d = 5$, and hence there is an extra logarithmic term in this case. For $d \in \{3, 4\}$ we get

$$\begin{aligned}
 L(\tau, \xi) &\leq C \frac{(1 + \|\tau' - |\xi'|\|)^{\frac{d-3}{4}}}{(1 + |\tau'| + |\xi'|)^{\frac{d-1}{4}}} \\
 &= C \frac{\lambda^{1/2} (\lambda + \|\tau| - |\xi||)^{\frac{d-3}{4}}}{(\lambda + |\tau| + |\xi|)^{\frac{d-1}{4}}} \\
 (2.6a) \quad &\leq C \frac{\lambda^{1/2} (\lambda + w_-)^{\frac{d-3}{4}}}{(w_+)^{\frac{d-1}{4}}},
 \end{aligned}$$

and for $d = 5$ we get

$$(2.6b) \quad L(\tau, \xi) \leq C \frac{\lambda^{1/2} (\lambda + w_-)^{\frac{d-3}{4}}}{(w_+)^{\frac{d-1}{4}}} \left(1 + \log \frac{\lambda + |\tau| + |\xi|}{\lambda + \|\tau| - |\xi||} \right)^{1/2}.$$

We choose $\lambda = \lambda(\tau, \xi)$ such that

$$(2.7) \quad \frac{(\lambda + w_-)^{\frac{d-1}{4}}}{\lambda^{\frac{3}{2}} (w_+)^{\frac{d-3}{4}}} = \frac{\lambda^{1/2} (\lambda + w_-)^{\frac{d-3}{4}}}{(w_+)^{\frac{d-1}{4}}}$$

or, equivalently,

$$(2.8) \quad \lambda^4 - w_+ \lambda - w_+ w_- = 0.$$

We need to know that such a λ exists, and we also need an estimate of λ in terms of the weights $1 + \|\tau\| \pm |\xi|$.

Consider the function $f(\lambda) := \lambda^4 - w_+ \lambda - w_+ w_-$. Clearly, $f(1) < 0$. On the other hand,

$$\begin{aligned} f\left(2w_+^{1/3}w_-^{1/6}\right) &= 16w_+^{4/3}w_-^{2/3} - 2w_+^{4/3}w_-^{1/6} - w_+w_- \\ &= \left(8w_+^{4/3}w_-^{2/3} - 2w_+^{4/3}w_-^{1/6}\right) + \left(8w_+^{4/3}w_-^{2/3} - w_+w_-\right) \\ &= 2w_+^{4/3}w_-^{1/6}\left(4w_-^{1/2} - 1\right) + w_+w_-^{2/3}\left(8w_+^{1/3} - w_-^{1/3}\right) \\ &> 0. \end{aligned}$$

It follows that there exists a λ with

$$(2.9) \quad 1 < \lambda < 2w_+^{1/3}w_-^{1/6}$$

such that (2.7) and (2.8) are satisfied. From (2.8) we have $\lambda + w_- = \frac{\lambda^4}{w_+}$; therefore

$$(2.10) \quad \lambda + w_- < \frac{\left(2w_+^{1/3}w_-^{1/6}\right)^4}{w_+} \leq Cw_+^{1/3}w_-^{2/3}.$$

It follows that the right-hand side of (2.7) is bounded by

$$\frac{\left(2w_+^{1/3}w_-^{1/6}\right)^{1/2}\left(w_+^{1/3}w_-^{2/3}\right)^{\frac{d-3}{4}}}{\left(w_+\right)^{\frac{d-1}{4}}} \leq Cw_+^{-\frac{d-1}{6}}w_-^{\frac{2d-5}{12}}.$$

Moreover $1 < \lambda \leq Cw_+$, and therefore in (2.6b)

$$1 + \log \frac{\lambda + \|\tau\| + |\xi|}{\lambda + \|\tau\| - |\xi|} \leq C \left(1 + \log \frac{w_+}{w_-}\right).$$

We conclude that $K(\tau, \xi)$ satisfies

$$K(\tau, \xi) \leq Cw_+^{-\frac{d-1}{6}}w_-^{\frac{2d-5}{12}}, \quad d \in \{3, 4, 5\},$$

while $L(\tau, \xi)$ satisfies

$$L(\tau, \xi) \leq Cw_+^{-\frac{d-1}{6}}w_-^{\frac{2d-5}{12}}, \quad d \in \{3, 4\},$$

and

$$L(\tau, \xi) \leq Cw_+^{-\frac{d-1}{6}}w_-^{\frac{2d-5}{12}} \left(1 + \log \frac{w_+}{w_-}\right)^{1/2}, \quad d = 5.$$

Next we deal with all dimensions $d \geq 7$. From estimate (2.4) we have $K(\tau, \xi) \leq C\frac{|\xi|}{\lambda^2}J_2^d(\tau', \xi')^{1/2}$. Proposition 2.1 with $m = d$ and $l = 2$ now gives

$$\alpha = \min \left\{ \frac{m + 1 - 4l}{2}, 0 \right\} = \min \left\{ \frac{d - 7}{2}, 0 \right\} = 0.$$

Notice, however, that $m + 1 = 4l$ when $d = 7$, and therefore there is a logarithmic loss in this case. We get

$$\begin{aligned} K(\tau, \xi) &\leq C \frac{|\xi|}{\lambda^2} \left(\frac{(1 + \|\tau'\| - |\xi'|)^3}{(1 + |\xi'| + |\tau'|)^4} \right)^{1/2} \\ &= C \frac{|\xi|}{\lambda^2} \cdot \frac{\lambda^2 (\lambda + \|\tau\| - |\xi|)^{3/2}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)^2} \\ &\leq C \frac{(\lambda + \|\tau\| - |\xi|)^{3/2}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)} \quad \text{if } d > 7 \end{aligned}$$

and

$$K(\tau, \xi) \leq C \frac{(\lambda + \|\tau\| - |\xi|)^{3/2}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)} \left(1 + \log \frac{\lambda + |\tau| + |\xi|}{\lambda + \|\tau\| - |\xi|} \right)^{1/2} \quad \text{if } d = 7.$$

From (2.5) we have $L(\tau, \xi) \leq C J_1^{d-2}(\tau', \xi')^{1/2}$. By applying Proposition 2.1, this time with $m = d - 2$, $l = 1$, and $\alpha = \min \left\{ \frac{m+1-4l}{2}, 0 \right\} = 0$ (notice that $m + 1 \neq 4l$, so there is no logarithmic loss for L), we get

$$\begin{aligned} L(\tau, \xi) &\leq C \left(\frac{1 + \|\tau'\| - |\xi'|}{(1 + |\tau'| + |\xi'|)^2} \right)^{1/2} \\ &= C \frac{\lambda^{1/2} (\lambda + \|\tau\| - |\xi|)^{1/2}}{(\lambda + |\tau| + |\xi|)}. \end{aligned}$$

We need to choose λ such that

$$(2.11) \quad \frac{(\lambda + \|\tau\| - |\xi|)^{3/2}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)} = \frac{\lambda^{1/2} (\lambda + \|\tau\| - |\xi|)^{1/2}}{(\lambda + |\tau| + |\xi|)}$$

or, equivalently, $\lambda^2 - \lambda - \|\tau\| - |\xi| = 0$. We choose $\lambda = \frac{1 + \sqrt{1 + 4\|\tau\| - |\xi|}}{2}$. Notice that $\lambda \geq 1$ and that $\lambda \leq C(\sqrt{\|\tau\| - |\xi|})$, and therefore

$$\lambda + \|\tau\| - |\xi| \leq C \left(1 + \sqrt{\|\tau\| - |\xi|} + \|\tau\| - |\xi| \right) \simeq 1 + \|\tau\| - |\xi|.$$

Therefore the right-hand side of (2.11) is bounded by

$$(2.12) \quad \frac{\left(1 + \sqrt{\|\tau\| - |\xi|} \right)^{1/2} (1 + \|\tau\| - |\xi|)^{1/2}}{1 + |\tau| + |\xi|} \leq C \frac{(1 + \|\tau\| - |\xi|)^{3/4}}{1 + |\tau| + |\xi|}.$$

Moreover, since $1 < \lambda \leq C(1 + |\tau| + |\xi|)$, we have

$$(2.13) \quad \log \frac{\lambda + |\tau| + |\xi|}{\lambda + \|\tau\| - |\xi|} \leq C \log \frac{1 + |\tau| + |\xi|}{1 + \|\tau\| - |\xi|}.$$

By putting everything together we conclude that

$$\begin{aligned} K(\tau, \xi) &\leq C \frac{(1 + \|\tau\| - |\xi|)^{3/4}}{1 + |\tau| + |\xi|} && \text{if } d > 7, \\ K(\tau, \xi) &\leq C \frac{(1 + \|\tau\| - |\xi|)^{3/4}}{1 + |\tau| + |\xi|} \left(1 + \log \frac{1 + |\tau| + |\xi|}{1 + \|\tau\| - |\xi|} \right)^{1/2}, && \text{if } d = 7, \end{aligned}$$

and

$$L(\tau, \xi) \leq C \frac{(1 + |\tau| - |\xi|)^{3/4}}{1 + |\tau| + |\xi|}, \quad d \geq 7.$$

This completes the proof in the case $d \geq 7$.

It remains to deal with the intermediate case $d = 6$. In this case we have the extra difficulty that K still behaves as in the lower-dimensional case but L behaves as in the higher-dimensional case. Indeed, we have $K(\tau, \xi) \leq C \frac{|\xi|}{\lambda^2} J_2^d(\tau', \xi')^{1/2}$, and Proposition 2.1 gives $\alpha = \min \left\{ \frac{d-7}{2}, 0 \right\} = -\frac{1}{2}$; therefore

$$\begin{aligned} K(\tau, \xi) &\leq C \frac{|\xi|}{\lambda^2} \cdot \frac{(1 + |\tau'| - |\xi'|)^{5/4}}{(1 + |\tau'| + |\xi'|)^{7/4}} \\ &= \frac{|\xi|}{\lambda^2} \cdot \frac{\lambda^{7/4} (\lambda + |\tau| - |\xi|)^{5/4}}{\lambda^{5/4} (\lambda + |\tau| + |\xi|)^{7/4}} \\ &\leq C \frac{(\lambda + |\tau| - |\xi|)^{5/4}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)^{3/4}}. \end{aligned}$$

For $L(\tau, \xi)$ we have $L(\tau, \xi) \leq C J_1^{d-2}(\tau', \xi')^{1/2}$, and Proposition 2.1 now gives $\alpha = \min \left\{ \frac{d-5}{2}, 0 \right\} = 0$. We get

$$\begin{aligned} L(\tau, \xi) &\leq C \frac{(1 + |\tau'| - |\xi'|)^{1/2}}{1 + |\tau'| + |\xi'|} \\ &= C \frac{\lambda^{1/2} (\lambda + |\tau| - |\xi|)^{1/2}}{\lambda + |\tau| + |\xi|}. \end{aligned}$$

We need to choose λ so that

$$(2.14) \quad \frac{(\lambda + |\tau| - |\xi|)^{5/4}}{\lambda^{3/2} (\lambda + |\tau| + |\xi|)^{3/4}} = \frac{\lambda^{1/2} (\lambda + |\tau| - |\xi|)^{1/2}}{\lambda + |\tau| + |\xi|}$$

or, equivalently,

$$(2.15) \quad \lambda^8 - (\lambda + |\tau| - |\xi|)^3 (\lambda + |\tau| + |\xi|) = 0.$$

To see that such a λ exists and to obtain estimates for it, consider the function

$$f : [1, \infty) \rightarrow \mathbb{R}, \quad f(\lambda) = \lambda^8 - (\lambda + |\tau| - |\xi|)^3 (\lambda + |\tau| + |\xi|).$$

Then, on the one hand, $f(1) \leq 0$ and, on the other hand, with $w_{\pm} = 1 + |\tau| \pm |\xi|$, we have

$$(2.16) \quad f\left(2w_-^{4/5} w_+^{1/5}\right) = 2^8 w_-^{32/5} w_+^{8/5}$$

$$(2.17) \quad - \left(2w_-^{4/5} w_+^{1/5} + |\tau| - |\xi|\right)^3 \left(2w_-^{4/5} w_+^{1/5} + |\tau| + |\xi|\right).$$

We have

$$2w_-^{4/5} w_+^{1/5} + |\xi| - |\tau| \leq 2w_-^{4/5} w_+^{1/5} + w_- \leq 3w_-^{4/5} w_+^{1/5}$$

and

$$2w_-^{4/5} w_+^{1/5} + |\xi| + |\tau| \leq 3w_+,$$

and therefore

$$\begin{aligned} f\left(2w_-^{4/5}w_+^{1/5}\right) &\geq 2^8w_-^{32/5}w_+^{8/5} - \left(3w_-^{4/5}w_+^{1/5}\right)^3 \cdot 3w_+ \\ &= 2^8w_-^{32/5}w_+^{8/5} - 3^4w_-^{12/5}w_+^{8/5} \\ &= w_+^{8/5}\left(2^8w_-^{32/5} - 3^4w_-^{12/5}\right) \\ &> 0. \end{aligned}$$

It follows that there exists a λ with $1 \leq \lambda \leq 2w_-^{4/5}w_+^{1/5}$ such that both (2.14) and (2.15) are satisfied. Notice that for this λ we have

$$(2.18) \quad \lambda + \|\tau\| - \|\xi\| \leq 2w_-^{4/5}w_+^{1/5} + w_- \leq 3w_-^{4/5}w_+^{1/5}.$$

Moreover, from (2.15) we get

$$(2.19) \quad \lambda = (\lambda + \|\tau\| - \|\xi\|)^{3/8} (\lambda + |\tau| + |\xi|)^{1/8}.$$

By using first (2.19) we find that the right-hand side of (2.14) is bounded by

$$\begin{aligned} &\frac{((\lambda + \|\tau\| - \|\xi\|)^{3/8}(\lambda + |\tau| + |\xi|)^{1/8})^{1/2} (\lambda + \|\tau\| - \|\xi\|)^{1/2}}{\lambda + |\tau| + |\xi|} \\ &= \frac{(\lambda + \|\tau\| - \|\xi\|)^{11/16}}{(\lambda + |\tau| + |\xi|)^{15/16}}. \end{aligned}$$

Next we use (2.18) for the numerator and $\lambda \geq 1$ for the denominator to get

$$\dots \leq C \frac{\left(w_-^{4/5}w_+^{1/5}\right)^{11/16}}{w_+^{15/16}} = C \frac{w_-^{11/20}}{w_+^{4/5}}.$$

We conclude that

$$K(\tau, \xi) \leq C \frac{w_-^{11/20}}{w_+^{4/5}}, \quad L(\tau, \xi) \leq C \frac{w_-^{11/20}}{w_+^{4/5}}.$$

This completes the proof of Theorem 1.1. \square

Proof of Theorem 1.2. Let (s_0, δ_0) be defined by

$$(2.20a) \quad (s_0, \delta_0) = \left(\frac{d-1}{6}, -\frac{2d-5}{12}\right) \quad \text{if } d \in \{3, 4, 5\},$$

$$(2.20b) \quad (s_0, \delta_0) = \left(\frac{4}{5}, -\frac{11}{20}\right) \quad \text{if } d = 6,$$

$$(2.20c) \quad (s_0, \delta_0) = \left(1, -\frac{3}{4}\right) \quad \text{if } d \geq 7.$$

Notice that in all cases we have $s_0 + \delta_0 = \frac{1}{4}$. Then the conditions (1.8) in Theorem 1.2 amount to $s + \delta \leq \frac{1}{4}$ and $s \leq s_0$. We have

$$w_+^s w_-^\delta \leq w_+^s w_-^{1/4-s} = w_+^{s_0} w_-^{\delta_0} \frac{w_-^{1/4-\delta_0-s}}{w_+^{s_0-s}} = w_+^{s_0} w_-^{\delta_0} \left(\frac{w_-}{w_+}\right)^{s_0-s} \leq w_+^{s_0} w_-^{\delta_0}.$$

The result then follows from Theorem 1.1. \square

Remark. It is possible to improve the estimates in this section if we allow weights of the form $w_+^{s_1} w_-^{\delta_1} + w_+^{s_2} w_-^{\delta_2}$. Consider, for example, the case $d \in \{3, 4, 5\}$. In the proof of Theorem 1.1 we bounded both K and L by

$$B := \frac{\lambda^{\frac{1}{2}}(\lambda + w_-)^{\frac{d-3}{4}}}{w_+^{\frac{d-1}{4}}},$$

with a logarithmic loss if $d = 5$. From (2.8), $\lambda = (\lambda + w_-)^{1/4} w_+^{1/4}$, and hence

$$B \leq C \frac{(\lambda + w_-)^{\frac{2d-5}{8}}}{w_+^{\frac{2d-3}{8}}}.$$

By using $(\lambda + w_-)^{2d-5} \leq C\lambda^{2d-5} + w_-^{2d-5}$ and $\lambda \leq Cw_-^{1/6} w_+^{1/3}$ (see (2.9)) we get

$$B \leq C \frac{w_-^{\frac{2d-5}{48}}}{w_+^{\frac{d-1}{6}}} + \frac{w_-^{\frac{2d-5}{8}}}{w_+^{\frac{2d-3}{8}}}.$$

Each of these fractions is smaller than

$$\frac{w_-^{\frac{2d-5}{12}}}{w_+^{\frac{d-1}{6}}}.$$

3. Improved estimates in the region $\{|\tau| > |\xi|\}$. In this section we prove Theorems 1.3 and 1.4.

Proof of Theorems 1.3 and 1.4. By arguing, and using the same notation, as in the proof of Theorem 1.1, we have (see (2.2)–(2.5))

$$\begin{aligned} |\widehat{\rho}(\tau, \xi)| &\leq \left(\int_{S^{d-1}} |\widehat{g}(\tau, \xi, v)|^2 d\sigma(v) \right)^{1/2} \cdot K(\tau, \xi) \\ (3.1) \quad &+ \left(\int_{S^{d-1}} |\widehat{f}(\tau, \xi, v)|^2 d\sigma(v) \right)^{1/2} \cdot L(\tau, \xi), \end{aligned}$$

where

$$K(\tau, \xi) \simeq \frac{|\xi|}{\lambda^2} J_2^d(\tau', \xi')^{1/2}, \quad L(\tau, \xi) \simeq J_1^{d-2}(\tau', \xi')^{1/2},$$

with $\tau' = \tau/\lambda$ and $\xi' = \xi/\lambda$.

In what follows we will reduce to proving pointwise estimates for $K(\tau, \xi)$ and $L(\tau, \xi)$. As in the proof of Theorem 1.1, the estimates announced in Theorems 1.3 and 1.4 are a direct consequence of these estimates and (3.1), so we will omit the details here.

Recall from [8] that, when $|\tau| > |\xi|$, the integrals $J_l^m(\tau, \xi)$ defined in (2.1) satisfy

$$\begin{aligned} J_l^m(\tau, \xi) &\simeq \frac{(1 + |\tau| - |\xi|)^\alpha}{(1 + |\tau| + |\xi|)^{2l+\alpha}} \quad \text{if } m + 1 \neq 4l, \\ J_l^m(\tau, \xi) &\simeq \frac{1}{(1 + |\tau| + |\xi|)^{2l}} \left(1 + \log \frac{1 + |\tau| + |\xi|}{1 + |\tau| - |\xi|} \right) \quad \text{if } m + 1 = 4l, \end{aligned}$$

where $\alpha = \min((m + 1 - 4l)/2, 0)$.

We may assume that $|\tau| \geq 1$. First, observe that if $|\tau| \geq 2|\xi|$, then

$$J_l^m(\tau, \xi) \simeq \frac{(1 + |\tau| - |\xi|)^\alpha}{(1 + |\tau| + |\xi|)^{2l+\alpha}} \simeq \frac{1}{|\tau|^{2l}}$$

and $1 + \log \frac{1+|\tau|+|\xi|}{1+|\tau|-|\xi|} \simeq 1$; therefore in all dimensions (use $\lambda = 1$) we have

$$K(\tau, \xi) \simeq |\xi| J_2^d(\tau, \xi)^{1/2} \simeq \frac{|\xi|}{|\tau|^2} \leq C \frac{1}{|\tau|}$$

and

$$L(\tau, \xi) \simeq J_1^{d-2}(\tau, \xi)^{1/2} \leq C \frac{1}{|\tau|}$$

which is stronger than all estimates in (1.10) and (1.12). Thus we may assume that $|\xi| < |\tau| < 2|\xi|$, in which case $1 + |\tau| + |\xi| \simeq |\xi|$.

Suppose that $d \in \{3, 4, 5\}$. Then we have

$$K(\tau, \xi) \simeq \frac{|\xi|}{\lambda^2} J_2^d(\tau', \xi')^{1/2} \leq C \frac{1}{(\lambda + |\tau| - |\xi|)^{\frac{7-d}{4}} |\xi|^{\frac{d-3}{4}}}.$$

Similarly,

$$L(\tau, \xi) \simeq J_1^{d-2}(\tau', \xi')^{1/2} \leq C \frac{\lambda}{(\lambda + |\tau| - |\xi|)^{\frac{5-d}{4}} |\xi|^{\frac{d-1}{4}}} \quad \text{if } d = 3, 4,$$

and the same estimate with a logarithmic loss if $d = 5$. Choose $\lambda = 1$ to get the estimates in Theorem 1.4. On the other hand, there exists a λ such that

$$(3.2) \quad \frac{1}{(\lambda + |\tau| - |\xi|)^{\frac{7-d}{4}} |\xi|^{\frac{d-3}{4}}} = \frac{\lambda}{(\lambda + |\tau| - |\xi|)^{\frac{5-d}{4}} |\xi|^{\frac{d-1}{4}}}.$$

Indeed, this is equivalent to

$$(3.3) \quad \lambda(\lambda + |\tau| - |\xi|)^{\frac{1}{2}} = |\xi|^{\frac{1}{2}},$$

i.e., $\varphi(\lambda) = 0$, where $\varphi(\lambda) = \lambda^3 + (|\tau| - |\xi|)\lambda^2 - |\xi|$. We have $\varphi(0) = -|\xi| < 0$ and $\varphi(|\xi|^{1/3}) = (|\tau| - |\xi|)|\xi|^{2/3} > 0$, and therefore there exists a $\lambda \in (0, |\xi|^{1/3})$ such that (3.2) is satisfied. For this λ , (3.3) gives $\lambda + |\tau| - |\xi| = \frac{|\xi|}{\lambda^2} \geq |\xi|^{1/3}$, and therefore the left-hand side of (3.2) is no greater than

$$\frac{1}{|\xi|^{\frac{1}{3}, \frac{7-d}{4}} |\xi|^{\frac{d-3}{4}}} = \frac{1}{|\xi|^{\frac{d-1}{6}}}.$$

This gives the estimates in Theorem 1.3.

Consider next $d = 6$. By working as above we get

$$K(\tau, \xi) \leq C \frac{1}{(\lambda + |\tau| - |\xi|)^{1/4} |\xi|^{3/4}} \quad \text{and} \quad L(\tau, \xi) \leq C \frac{\lambda}{|\xi|}.$$

Use $\lambda = 1$ to get the estimates of Theorem 1.4. Alternatively, find $\lambda \in (0, |\xi|^{1/5})$ such that

$$\frac{1}{(\lambda + |\tau| - |\xi|)^{1/4} |\xi|^{3/4}} = \frac{\lambda}{|\xi|},$$

and work as above to get the estimates in Theorem 1.3. In the case $d \geq 7$ we are already covered by (1.9c). \square

Remark. We can obtain more estimates in $H^{s,\delta}$ -spaces by interpolating between (1.9) and (1.11). For example, if $d \in \{3, 4, 5\}$, we have that for all $\alpha \in [0, 1]$ estimates (1.10) hold true with $s = \alpha \frac{d-1}{6} + (1-\alpha) \frac{d-3}{4}$ and $\delta = (1-\alpha) \frac{7-d}{4}$, with similar results in all other dimensions.

4. Counterexamples. In this section we discuss the optimality of our results. We use the following notation: We fix $i, j \in \{1, \dots, d\}$, with $i \neq j$. Given a vector $\xi \in R^d$, we denote by $\xi' \in R^{d-2}$ the vector resulting from ξ if the i th and j th coordinates are removed, i.e., $\xi' = (\xi_i, \dots, \xi_{i-1}, \xi_{i+1}, \dots, \xi_{j-1}, \xi_{j+1}, \dots, \xi_n)$. We denote by μ the Lebesgue measure in R^d and by σ the corresponding surface measure on S^{d-1} .

PROPOSITION 4.1. *The condition $s + \delta \leq 1/4$ in Theorem 1.2 is necessary.*

Proof. Fix $N \gg 1$, and define

$$A = \left\{ (\tau, \xi) : 5 \leq \tau \leq 10, -2N^{1/2} \leq \xi_i \leq -N^{1/2}, N/2 \leq \xi_j \leq N, |\xi'| \leq 1 \right\}.$$

Let $\phi : S^{d-1} \rightarrow R$ be defined by $\phi(v) = \phi_i(v_i)\phi_j(v_j)$, where $\phi_i, \phi_j : R \rightarrow [0, 1]$ are smooth cutoff functions such that

$$\begin{aligned} \phi_i(v_i) &= 1 \quad \text{for} \quad -\frac{9}{10} \leq v_i \leq -\frac{3}{5}, \\ \phi_i(v_i) &= 0 \quad \text{for} \quad v_i \geq -\frac{1}{2}, \\ \phi_j(v_j) &= 1 \quad \text{for} \quad \frac{1}{10N^{1/2}} \leq v_j \leq \frac{3}{20N^{1/2}}, \\ \phi_j(v_j) &= 0 \quad \text{for} \quad v_j \leq 0 \quad \text{or} \quad v_j \geq \frac{1}{4N^{1/2}}. \end{aligned}$$

In particular

$$\text{supp } \phi_i \subseteq [-1, -1/2], \quad \text{supp } \phi_j \subseteq [0, 1/4N^{1/2}].$$

For $(\tau, \xi) \in A$ and $u \in \text{supp } \phi$ we have

$$\begin{aligned} \tau + v \cdot \xi &= \tau + v_i \xi_i + v_j \xi_j + v' \cdot \xi' \leq 10 + \frac{18}{10N^{1/2}} + \frac{3}{20N^{1/2}} N + 1 \leq 11 + 2N^{1/2}, \\ \tau + v \cdot \xi &= \tau + v_i \xi_i + v_j \xi_j + v' \cdot \xi' \geq 5 + \frac{3N^{1/2}}{5} + 0 - 1 = 4 + \frac{3N^{1/2}}{5}, \end{aligned}$$

and therefore

$$(4.1) \quad \tau + v \cdot \xi \simeq N^{1/2}.$$

Also,

$$\begin{aligned} v_j \xi_i - v_i \xi_j &\leq 0 + \frac{9N}{10} \leq N, \\ v_j \xi_i - v_i \xi_j &\geq -\frac{6}{20} + \frac{3N}{10} = \frac{3(N-1)}{10}, \end{aligned}$$

and therefore

$$(4.2) \quad v_j \xi_i - v_i \xi_j \simeq N.$$

On the other hand,

$$(4.3) \quad \left| \frac{\Omega_v^{i,j} \phi(v)}{\tau + v \cdot \xi} \right| \leq \frac{|v_i| |\phi_i| \left| \frac{\partial \phi_j}{\partial v_j} \right| + |v_j| |\phi_j| \left| \frac{\partial \phi_i}{\partial v_i} \right|}{|\tau + v \cdot \xi|} \leq C \frac{N^{1/2} + N^{-1/2}}{N^{1/2}} \leq C.$$

We define $f(t, x, v)$ and $g(t, x, v) \in L^2(R \times R^d \times S^{d-1})$ by

$$\widehat{f}(\tau, \xi, v) = \frac{\Omega_v^{i,j} \widehat{g}(\tau, \xi, v)}{i(\tau + v \cdot \xi)} \quad \text{and} \quad \widehat{g}(\tau, \xi, v) = \chi_A(\tau, \xi) \phi(v).$$

Then f and g are well defined (recall (4.1)) and $\partial_t f + v \cdot \nabla_x f = \Omega_v^{i,j} g$. The definition of f and integration by parts yield

$$\begin{aligned} \widehat{\rho}_s(\tau, \xi) &= \int_{S^{d-1}} \widehat{f}(\tau, \xi, v) d\sigma v \\ &= -i \chi_A(\tau, \xi) \int_{S^{d-1}} \frac{\Omega_v^{i,j} \phi(v)}{\tau + v \cdot \xi} d\sigma v \\ &= i \chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_j \xi_i - v_i \xi_j}{(\tau + v \cdot \xi)^2} \phi(v) d\sigma(v). \end{aligned}$$

Since $\phi \geq 0$ and $v_j \xi_i - v_i \xi_j \geq 0$ we have

$$|\widehat{\rho}_s(\tau, \xi)| = \chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_j \xi_i - v_i \xi_j}{(\tau + v \cdot \xi)^2} \phi(v) d\sigma(v)$$

and, by using (4.1) and (4.2),

$$|\widehat{\rho}_s(\tau, \xi)| \simeq \chi_A(\tau, \xi) \int_{S^{d-1}} \phi(v) d\sigma(v) \geq c \chi_A(\tau, \xi) \frac{1}{N^{1/2}}.$$

For $(\tau, \xi) \in A$ we have $1 + \|\tau\| \pm \|\xi\| \simeq N$, and therefore

$$(4.4) \quad \begin{aligned} \|\rho_s\|_{H^{s,\delta}} &= \|(1 + |\tau| + |\xi|)^s (1 + \|\tau\| - \|\xi\|)^\delta \widehat{\rho}_s(\tau, \xi)\|_{L^2} \\ &\geq c N^{s+\delta-1/2} \mu(A)^{1/2}. \end{aligned}$$

On the other hand, by using (4.3),

$$(4.5) \quad \begin{aligned} \|f\|_{L^2(R \times R^d \times S^{d-1})} &= \left(\int_{R^{1+d}} \int_{S^{d-1}} \left| \frac{\chi_A(\tau, \xi) \Omega_v^{i,j} \phi(v)}{i(\tau + v \cdot \xi)} \right|^2 d\sigma(v) d\tau d\xi \right)^{1/2} \\ &\leq C \mu(A)^{1/2} \sigma(\text{supp } \phi)^{1/2} \\ &\leq C \mu(A)^{1/2} \frac{1}{N^{1/4}}. \end{aligned}$$

Finally,

$$\begin{aligned}
 \|g\|_{L^2(R \times R^d \times S^{d-1})} &= \left(\int_{R^{1+d}} \int_{S^{d-1}} \chi_A(\tau, \xi)^2 \phi(v)^2 d\sigma(v) d\tau d\xi \right)^{1/2} \\
 &\leq C\mu(A)^{1/2} \sigma(\text{supp } \phi)^{1/2} \\
 (4.6) \qquad \qquad \qquad &\leq C\mu(A)^{1/2} \frac{1}{N^{1/4}}.
 \end{aligned}$$

If for some (s, δ) we have $\|\rho_s\|_{H^{s,\delta}} \leq C(\|f\|_{L^2} + \|g\|_{L^2})$, then (4.4)–(4.6) imply that $N^{s+\delta-1/2} \leq CN^{-1/4}$, and therefore $s + \delta \leq 1/4$ as required. We get exactly the same result if we have $H_{\log}^{s,\delta}$ instead of $H^{s,\delta}$. \square

PROPOSITION 4.2. *The condition $s \leq 1$ in Theorem 1.2 is necessary.*

Proof. Fix $N \gg 1$. Define

$$A = \{(\tau, \xi) : N \leq \xi_i \leq 2N, -1 \leq \xi_j \leq 0, |\xi'| \leq 1, |\xi| \leq \tau \leq |\xi| + 1\}.$$

Let $\phi : R \rightarrow R$ be a smooth cutoff function, with $0 \leq \phi \leq 1$ and $\phi = 1$ on $[1/20, 1/10]$ and $\phi = 0$ outside $[1/21, 1/9]$. Define

$$\widehat{g}(\tau, \xi, v) = \chi_A(\tau, \xi) \phi(v_i) \phi(v_j), \quad \widehat{f}(\tau, \xi, v) = \frac{\Omega_v^{i,j} \widehat{g}(\tau, \xi, v)}{i(\tau + v \cdot \xi)}.$$

Notice that for $(\tau, \xi) \in A$ and $v \in S^{d-1}$ such that $v_i, v_j \in \text{supp } \phi$ we have

$$\tau + v \cdot \xi \leq |\tau| + |\xi| \leq CN$$

and, since $v_i \xi_i \geq 0$ and $v_j \xi_j \geq -1$,

$$\tau + v \cdot \xi = \tau + v_i \xi_i + v_j \xi_j + v' \cdot \xi' \geq N + 0 - 1 - 1 \geq cN.$$

Therefore $\tau + v \cdot \xi > 0$ and

$$\tau + v \cdot \xi \simeq N.$$

It follows that f is well defined and, moreover,

$$\begin{aligned}
 \left| \widehat{f}(\tau, \xi, v) \right| &\leq \frac{\chi_A(\tau, \xi) [|v_i| \phi(v_i) |\phi'(v_j)| + |v_j| |\phi'(v_i)| \phi(v_j)]}{|\tau + v \cdot \xi|} \\
 &\leq C \frac{\chi_A(\tau, \xi)}{N};
 \end{aligned}$$

therefore

$$\|f\|_{L^2} \leq C \frac{\mu(A)^{1/2}}{N}.$$

Also,

$$\|g\|_{L^2} \leq C\mu(A)^{1/2},$$

and therefore

$$(4.7) \qquad \|f\|_{L^2} + \|g\|_{L^2} \leq C\mu(A)^{1/2}.$$

On the other hand, by using integration by parts we find

$$\widehat{\rho}_s(\tau, \xi) = i\chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_j \xi_i - v_i \xi_j}{(\tau + v \cdot \xi)^2} \phi(v_i) \phi(v_j) d\sigma(v).$$

Since $v_i \geq 0$ and $\xi_j \leq 0$ we have $v_j \xi_i - v_i \xi_j \geq v_j \xi_i \geq \frac{N}{21} \geq cN$, and therefore

$$\begin{aligned} |\widehat{\rho}_s(\tau, \xi)| &= \chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_j \xi_i - v_i \xi_j}{(\tau + v \cdot \xi)^2} \phi(v_i) \phi(v_j) d\sigma(v) \\ &\geq c\chi_A(\tau, \xi) \frac{N}{N^2} \int_{S^{d-1}} \phi(v_i) \phi(v_j) d\sigma(v) \\ &\geq c \frac{\chi_A(\tau, \xi)}{N}. \end{aligned}$$

Moreover, for $(\tau, \xi) \in A$, $1 + |\tau| + |\xi| \simeq N$ and $1 + ||\tau| - |\xi|| \simeq 1$; therefore

$$\begin{aligned} &\| (1 + |\tau| + |\xi|)^s (1 + ||\tau| - |\xi||)^\delta |\widehat{\rho}_s(\tau, \xi)| \|_{L^2} \\ (4.8) \quad &\geq cN^s \frac{\mu(A)^{1/2}}{N} = N^{s-1} \mu(A)^{1/2}. \end{aligned}$$

If $\|\rho_s\|_{H^{s,\delta}} \leq C(\|f\|_{L^2} + \|g\|_{L^2})$ is satisfied, then (4.7) and (4.8) imply that $N^{s-1} \leq C$; therefore $s \leq 1$. We get exactly the same result if we have $H_{\log}^{s,\delta}$ instead of $H^{s,\delta}$. \square

PROPOSITION 4.3. *The condition $s \leq \frac{d-1}{6}$ in Theorem 1.2 is necessary.*

Proof. For simplicity of notation we use $i = 1$ and $j = d$. Then $\xi' = (\xi_2, \dots, \xi_{d-1})$. Fix $N \gg 1$, and define

$$(4.9) \quad A = \{(\tau, \xi) : 1 \leq \tau - |\xi| \leq 2, 0 \leq \xi_1 \leq 1, -2N \leq \xi_d \leq -N, |\xi'| \leq 1\}.$$

Notice that if $(\tau, \xi) \in A$, then

$$\begin{aligned} (4.10) \quad &1 + ||\tau| - |\xi|| \simeq 1, \\ (4.11) \quad &1 + |\tau| + |\xi| \simeq N, \\ &\frac{1}{2} \leq \frac{-\xi_d}{|\xi|} \leq 1. \end{aligned}$$

Let $\phi : R \rightarrow [0, 1]$ be a smooth cutoff function, with $\text{supp } \phi \subseteq [N^{-\frac{1}{3}}, 2N^{-\frac{1}{3}}]$ and $\phi \equiv 1$ on $[\frac{5}{4}N^{-\frac{1}{3}}, \frac{3}{2}N^{-\frac{1}{3}}]$. Let $\psi : R \rightarrow [0, 1]$ be a smooth cutoff function, with $\text{supp } \psi \subseteq [1 - 4N^{-\frac{2}{3}}, 1 - N^{-\frac{2}{3}}]$ and $\psi \equiv 1$ on $[1 - 3N^{-\frac{2}{3}}, 1 - 2N^{-\frac{2}{3}}]$. Define

$$(4.12) \quad B = \{v \in S^{d-1} : v_1 \in \text{supp } \phi, v_d \in \text{supp } \psi\}.$$

For $(\tau, \xi) \in A$ and $v \in B$ we then have $v_d, \xi_1 \geq 0$ and $\xi_d \leq 0$; therefore

$$(4.13) \quad v_d \xi_1 - v_1 \xi_d \geq 0 + v_1 |\xi_d| \geq N^{\frac{2}{3}}.$$

Also

$$(4.14) \quad \tau + v \cdot \xi = \tau - |\xi| + |\xi| \left(1 + v_d \frac{\xi_d}{|\xi|} \right) + \sum_{i=1}^{d-1} v_i \xi_i.$$

We have $\tau - |\xi| \simeq 1$ and $|\sum_{i=1}^{d-1} v_i \xi_i| \leq C1$. Moreover,

$$1 + v_d \frac{\xi_d}{|\xi|} = \left(1 + \frac{\xi_d}{|\xi|}\right) + (1 - v_d) \frac{-\xi_d}{|\xi|} \geq 0 + N^{-\frac{2}{3}} \cdot \frac{1}{2} \geq cN^{-\frac{2}{3}},$$

and therefore

$$|\xi| \left(1 + v_d \frac{\xi_d}{|\xi|}\right) \geq cN^{\frac{1}{3}}.$$

As a consequence,

$$\tau + v \cdot \xi \geq cN^{\frac{1}{3}}.$$

Notice also that

$$|\xi| \left(1 + v_d \frac{\xi_d}{|\xi|}\right) = (|\xi| + \xi_d) + (1 - v_d) \cdot (-\xi_d),$$

with

$$(1 - v_d) \cdot (-\xi_d) \leq 4N^{-\frac{2}{3}} \cdot 2N \leq CN^{\frac{1}{3}}$$

and

$$|\xi| + \xi_d = \frac{|\xi|^2 - \xi_d^2}{|\xi| - \xi_d} = \frac{\sum_{k=1}^{d-1} \xi_k^2}{|\xi| + |\xi_d|} \leq C \frac{1}{N} \leq C.$$

Thus

$$|\xi| \left(1 + v_d \frac{\xi_d}{|\xi|}\right) \leq CN^{\frac{1}{3}}.$$

By using this in (4.14) we get

$$\tau + v \cdot \xi \leq CN^{\frac{1}{3}}.$$

We have shown that, for $(\tau, \xi) \in A$ and $v \in B$,

$$(4.15) \quad \tau + v \cdot \xi \simeq N^{\frac{1}{3}}.$$

Next define f and $g \in L^2(R_t \times R_x^d \times S_v^{d-1})$ by

$$\widehat{f}(\tau, \xi, v) = \frac{\Omega_v^{1,d} \widehat{g}(\tau, \xi, v)}{i(\tau + v \cdot \xi)} \quad \text{and} \quad \widehat{g}(\tau, \xi, v) = \chi_A(\tau, \xi) \phi(v_1) \psi(v_d).$$

We have

$$(4.16) \quad \|g\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} \leq \mu(A)^{1/2} \sigma(B)^{1/2}.$$

Also

$$\begin{aligned} \left| \widehat{f}(\tau, \xi, v) \right| &\leq \frac{\chi_A(\tau, \xi)}{|\tau + v \cdot \xi|} [v_1 \phi(v_1) |\psi'(v_d)| + v_d |\phi'(v_1)| \psi(v_d)] \\ &\leq C \frac{\chi_A(\tau, \xi) \chi_B(v)}{N^{\frac{1}{3}}} \left[N^{-\frac{1}{3}} \cdot 1 \cdot N^{\frac{2}{3}} + 1 \cdot N^{\frac{1}{3}} \cdot 1 \right] \\ &\leq C \chi_A(\tau, \xi) \chi_B(v), \end{aligned}$$

so that

$$(4.17) \quad \|f\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} \leq C\mu(A)^{1/2}\sigma(B)^{1/2}.$$

Now observe that B is contained in the set $\{v \in S^{d-1} : 1 - 4N^{-\frac{2}{3}} \leq v_d \leq 1\}$, which is a “cap” centered at the North pole of the sphere. It is easy to check that its area is $\leq CN^{-\frac{d-1}{3}}$, and therefore $\sigma(B) \leq CN^{-\frac{d-1}{3}}$. Then, from (4.16), (4.17), and previous observation, we get

$$(4.18) \quad \|f\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} + \|g\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} \leq C\mu(A)^{1/2}N^{-(d-1)/6}.$$

For ρ_s we integrate by parts to get

$$\begin{aligned} \widehat{\rho}_s(\tau, \xi) &= \int_{S^{d-1}} \frac{\Omega_v^{1,d}\widehat{g}(\tau, \xi, v)}{i(\tau + v \cdot \xi)} d\sigma(v) \\ &= i\chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_d\xi_1 - v_1\xi_d}{(\tau + v \cdot \xi)^2} \phi(v_1)\psi(v_d) d\sigma(v). \end{aligned}$$

Since $\phi, \psi \geq 0$ and $v_d\xi_1 - v_1\xi_d \geq 0$ we have

$$|\widehat{\rho}_s(\tau, \xi)| = \chi_A(\tau, \xi) \int_{S^{d-1}} \frac{v_d\xi_1 - v_1\xi_d}{(\tau + v \cdot \xi)^2} \phi(v_1)\psi(v_d) d\sigma(v),$$

and therefore, by using (4.13) and (4.15),

$$(4.19) \quad |\widehat{\rho}_s(\tau, \xi)| \geq c\chi_A(\tau, \xi) \int_{S^{d-1}} \phi(v_1)\psi(v_d) d\sigma(v),$$

we claim that

$$(4.20) \quad \int_{S^{d-1}} \phi(v_1)\psi(v_d)d\sigma(v) \geq cN^{-(d-1)/3}.$$

We have (recall that v' denotes the vector $(v_2, \dots, v_{d-1}) \in R^{d-2}$)

$$\begin{aligned} &\int_{S^{d-1}} \phi(v_1)\psi(v_d)d\sigma(v) \\ &= c \int_{v_1^2 + |v'|^2 \leq 1} \phi(v_1)\psi\left(\sqrt{1 - v_1^2 - |v'|^2}\right) dv_1 dv' \\ (4.21) \quad &\geq c \int_X \phi(v_1)\psi\left(\sqrt{1 - v_1^2 - |v'|^2}\right) dv_1 dv', \end{aligned}$$

where $X = \{\frac{5}{4}N^{-1/3} \leq v_1 \leq \frac{3}{2}N^{-1/3}, \frac{\sqrt{39}}{4}N^{-1/3} \leq |v'| \leq \frac{3}{2}N^{-1/3}\}$. Clearly, in the domain of integration in (4.21) we have $\phi(v_1) = 1$. We claim that we also have $\psi(\sqrt{1 - v_1^2 - |v'|^2}) = 1$. Indeed,

$$1 - \sqrt{1 - v_1^2 - |v'|^2} = \frac{v_1^2 + |v'|^2}{1 + \sqrt{1 - v_1^2 - |v'|^2}} \leq \frac{\frac{9}{4}N^{-2/3} + \frac{9}{4}N^{-2/3}}{\frac{3}{2}} = 3N^{-2/3},$$

and therefore

$$1 - 3N^{-2/3} \leq \sqrt{1 - v_1^2 - |v'|^2}.$$

Moreover,

$$1 - \sqrt{1 - v_1^2 - |v'|^2} = \frac{v_1^2 + |v'|^2}{1 + \sqrt{1 - v_1^2 - |v'|^2}} \geq \frac{\frac{25}{16}N^{-2/3} + \frac{39}{16}N^{-2/3}}{2} = 2N^{-2/3},$$

and therefore

$$\sqrt{1 - v_1^2 - |v'|^2} \leq 1 - 2N^{-2/3}.$$

It follows that

$$\int_{S^{d-1}} \phi(v_1)\psi(v_d)d\sigma(v) \geq c \int_X 1 \, dv_1 dv' \geq cN^{-(d-1)/3}.$$

This proves (4.20). By using (4.20) in (4.19) we get

$$|\widehat{\rho_s}(\tau, \xi)| \geq c\chi_A(\tau, \xi)N^{-\frac{d-1}{3}},$$

and therefore, by using (4.10) and (4.11),

$$(4.22) \quad \|w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \widehat{\rho_s}(\tau, \xi)\|_{L^2(R \times R^d)} \geq \mu(A)^{1/2} N^{s - \frac{d-1}{3}}.$$

If the estimate

$$\begin{aligned} & \|w_+(\tau, \xi)^s w_-(\tau, \xi)^\delta \widehat{\rho_s}(\tau, \xi)\|_{L^2(R \times R^d)} \\ & \leq C \left(\|f\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} + \|g\|_{L^2(R_t \times R_x^d \times S_v^{d-1})} \right) \end{aligned}$$

is valid, then it follows from (4.22) and (4.18) that

$$N^{s-(d-1)/3} \leq CN^{-(d-1)/6}$$

for $N \gg 1$ which gives

$$s \leq \frac{d-1}{6}.$$

We have exactly the same bound for s if $H^{s,\delta}$ is replaced by $H_{\log}^{s,\delta}$. \square

Acknowledgments. The authors acknowledge with pleasure several helpful conversations with Benoît Perthame and Luis Vega.

REFERENCES

- [1] C. BARDOS, F. GOLSE, AND B. PERTHAME, *The Rosseland approximation for the radiative transfer equations*, Comm. Pure Appl. Math., 40 (1987), pp. 691–721.
- [2] C. BARDOS, F. GOLSE, B. PERTHAME, AND R. SENTIS, *The nonaccretive radiative transfer equations: Existence of solutions and Rosseland approximation*, J. Funct. Anal., 77 (1988), pp. 434–460.
- [3] F. BOUCHUT, *Introduction to the mathematical theory of kinetic equations*, in Kinetic Equations and Asymptotic Theories, Ser. Appl. Math. 4, Elsevier, New York, 2000.
- [4] F. BOUCHUT, *Hypoelliptic regularity in kinetic equations*, J. Math. Pures Appl., 81 (2002), pp. 1135–1159.

- [5] F. BOUCHUT AND L. DESVILLETES, *Averaging lemmas without time Fourier transform and applications to discretized kinetic equations*, Proc. Roy. Soc. Edinburgh Sect. A 129, 1 (1999), pp. 19–36.
- [6] J. BOURGAIN, *Fourier transform restriction phenomena for certain lattice subsets and applications to nonlinear evolution equations. Part I*, Geom. Funct. Anal., 3 (1993), pp. 107–156; *Part II*, Geom. Funct. Anal., 3 (1993), pp. 209–262.
- [7] N. BOURNAVEAS AND B. PERTHAME, *Averages over spheres for kinetic transport equations; hyperbolic Sobolev spaces and Strichartz inequalities*, J. Math. Pures Appl., 9 (2001), pp. 517–534.
- [8] N. BOURNAVEAS AND S. GUTIERREZ, *On the regularity of averages over spheres for kinetic transport equations in hyperbolic Sobolev spaces*, Rev. Mat. Iberoamericana, 23 (2007), pp. 481–512.
- [9] F. CASTELLA AND B. PERTHAME, *Estimations de Strichartz pour les équations de transport cinétique*, C. R. Acad. Sci. Paris Sér. I, 322 (1996), pp. 535–540.
- [10] S. CHANDRASEKHAR, *Radiative Transfer*, Clarendon Press, Oxford, 1950.
- [11] L. DESVILLETES AND S. MISCHLER, *About the splitting algorithm for Boltzmann and B.G.K. equations*, Math. Models Methods Appl. Sci., 6 (1996), pp. 1079–1101.
- [12] R. DEVORE AND G. PETROVA, *The averaging lemma*, J. Amer. Math. Soc., 14 (2001), pp. 279–296.
- [13] R. J. DIPERNA AND P. L. LIONS, *On the Cauchy problem for the Boltzmann equation: Global existence and weak stability results*, Ann. of Math., 130 (1989), pp. 321–366.
- [14] R. J. DIPERNA AND P. L. LIONS, *Global weak solutions of Vlasov - Maxwell systems*, Comm. Pure Appl. Math., 42 (1989), pp. 729–757.
- [15] R. J. DIPERNA, P. L. LIONS, AND Y. MEYER, *L^p regularity of velocity averages*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 8 (1991), pp. 271–287.
- [16] Y. DOLAK AND C. SCHMEISER, *Kinetic models for chemotaxis: Hydrodynamic limits and spatio-temporal mechanisms*, J. Math. Biol., 51 (2005), pp. 595–615.
- [17] M. ESCOBEDO, P. LAURENCOT, AND S. MISCHLER, *On a kinetic equation for coalescing particles*, Comm. Math. Phys., 246 (2004), pp. 237–267.
- [18] D. FOSCHI AND S. KLAINERMAN, *Bilinear space-time estimates for homogeneous wave equations*, Ann. Sci. École Norm. Sup., 33 (2000), pp. 211–274.
- [19] P. GÉRARD, *Regularity of means of solutions of partial differential equations*, Journées Équations aux Dérivées Partielles, 1987, pp. 1–8.
- [20] P. GÉRARD, *Microlocal defect measures*, Comm. Partial Differential Equations, 16 (1991), pp. 1761–1794.
- [21] P. GÉRARD AND F. GOLSE, *Averaging regularity results for PDEs under transversality assumptions*, Comm. Pure Appl. Math., 45 (1992), pp. 1–26.
- [22] F. GOLSE AND B. PERTHAME, *Generalized solutions of the radiative transfer equations in a singular case*, Comm. Math. Phys., 106 (1986), pp. 211–239.
- [23] F. GOLSE, P. L. LIONS, B. PERTHAME, AND R. SENTIS, *Regularity of the moments of the solution of a transport equation*, J. Funct. Anal., 76 (1988), pp. 110–125.
- [24] F. GOLSE, B. PERTHAME, AND R. SENTIS, *Un résultat de compacité pour les équations du transport et application au calcul de la limite de la valeur propre principale d'un opérateur de transport*, C. R. Acad. Sci. Sér. I, 301 (1985), pp. 341–344.
- [25] F. GOLSE AND L. SAINT-RAYMOND, *The Navier-Stokes limit of the Boltzmann equation for bounded collision kernels*, Invent. Math., 155 (2004), pp. 81–161.
- [26] F. GOLSE AND L. SAINT-RAYMOND, *Velocity averaging in L^1 for the transport equation*, C. R. Math. Acad. Sci. Paris, 334 (2002), pp. 557–562.
- [27] P. E. JABIN AND L. VEGA, *A real space method for averaging lemmas*, J. Math. Pures Appl., 83 (2004), pp. 1309–1351.
- [28] P. E. JABIN AND L. VEGA, *Averaging lemmas and the x-ray transform*, C. R. Math. Acad. Sci. Paris, 337 (2003), pp. 505–510.
- [29] P. E. JABIN AND B. PERTHAME, *Regularity in kinetic formulations via averaging lemmas*, ESAIM Control Optim. Calc. Var., 8 (2002), pp. 761–774.
- [30] S. KLAINERMAN AND M. MACHEDON, *Space-time estimates for null forms and the local existence theorem*, Comm. Pure Appl. Math., 46 (1993), pp. 1221–1268.
- [31] S. KLAINERMAN AND M. MACHEDON, *Smoothing estimates for null forms and applications*, Duke Math. J., 81 (1995), pp. 99–133.
- [32] S. KLAINERMAN AND M. MACHEDON, *Estimates for null forms and the spaces H^s , S_b , δ* , Int. Math. Res. Not., 17 (1996), pp. 853–865.
- [33] P. L. LIONS AND B. PERTHAME, *Lemmes de moments, de moyenne et de dispersion*, C. R. Acad. Sci. Paris Sér. I 314, 11 (1992), pp. 801–806.

- [34] P. L. LIONS, *Régularité optimale des moyennes en vitesses*, C. R. Acad. Sci. Paris Sér. I, 320 (1995), pp. 911–915.
- [35] B. PERTHAME, *Mathematical tools for kinetic equations*, Bull. Amer. Math. Soc. (N.S.), 41 (2004), pp. 205–244.
- [36] B. PERTHAME AND J. L. VÁZQUEZ, *Bounded speed of propagation for solutions to radiative transfer equations*, Comm. Math. Phys., 130 (1990), pp. 457–469.
- [37] B. PERTHAME AND P. E. SOUGANIDIS, *A limiting case for velocity averaging*, Ann. Sci. École Norm. Sup. (4), 31 (1998), pp. 591–598.
- [38] D. TATARU, *The $X^{s,p}$ spaces and unique continuation for solutions to the semilinear wave equation*, Comm. Partial Differential Equations, 21 (1996), pp. 841–887.

STABILITY OF EQUILIBRIA FOR THE STEFAN PROBLEM WITH SURFACE TENSION*

JAN PRÜSS[†] AND GIERI SIMONETT[‡]

Abstract. We characterize the equilibrium states for the two-phase Stefan problem with surface tension and with or without kinetic undercooling, and we analyze their stability in terms of dependence on physical and geometric quantities.

Key words. free boundary problem, phase transitions, surface tension, kinetic undercooling, stability, bifurcation

AMS subject classifications. 35R55, 35B55, 35K55, 80A22

DOI. 10.1137/070700632

1. Introduction. The Stefan problem is a model for phase transitions in solid-liquid systems. In this paper, we consider the two-phase Stefan problem with the modified Gibbs–Thomson law

$$(1.1) \quad u = \sigma H + \delta V \quad \text{on} \quad \Gamma(t), \quad \sigma > 0, \quad \delta \geq 0,$$

and the kinetic condition

$$(1.2) \quad [d\partial_\nu u] = (\ell - [\kappa]u)V \quad \text{on} \quad \Gamma(t).$$

Here $\Gamma(t)$ denotes the unknown moving hypersurface that separates the liquid from the solid phase, u is the temperature, H the mean curvature of $\Gamma(t)$, σ the surface tension coefficient, δ the coefficient of kinetic undercooling, V the normal velocity of $\Gamma(t)$, ℓ the latent heat, $[\kappa]$ the jump of the heat capacities across $\Gamma(t)$, and $[d\partial_\nu u]$ the jump of the heat fluxes across $\Gamma(t)$. Note that in case $\sigma = \delta = 0$, i.e., for the classical Stefan problem, we have $u = 0$ at the interface, and then the kinetic condition becomes the classical Stefan condition.

Under appropriate boundary conditions we will show that spheres (together with constant temperature distributions) are the only equilibrium states for this system, and we will characterize the stability of these equilibria in terms of dependence on physical and geometric quantities.

In order to formulate the Stefan problem we introduce the following notation. Let Ω be a smooth bounded domain in \mathbb{R}^n whose boundary $\partial\Omega$ consists of two disjoint components, an “interior” part J_1 and an “exterior” part J_2 . We think of Ω as a homogeneous medium which is occupied by a liquid and a solid phase, say water and ice, that initially occupy the regions Ω_0^1 and Ω_0^2 , and that are separated by a sharp interface Γ_0 . More precisely, we assume that $\Gamma_0 \subset \Omega$ is a compact closed hypersurface, and that Ω_0^1 and Ω_0^2 are disjoint open sets such that $\bar{\Omega} = \bar{\Omega}_0^1 \cup \bar{\Omega}_0^2$, and such that $\partial\Omega_0^i = J_i \cup \Gamma_0$ for $i = 1, 2$. For the sake of definiteness we consider the open set Ω_0^1 as

*Received by the editors August 20, 2007; accepted for publication (in revised form) March 11, 2008; published electronically July 3, 2008.

<http://www.siam.org/journals/sima/40-2/70063.html>

[†]Institut für Mathematik, Martin-Luther-Universität Halle-Wittenberg, D-60120 Halle, Germany (jan.pruess@mathematik.uni-halle.de).

[‡]Department of Mathematics, Vanderbilt University, Nashville, TN 37240 (simonett@math.vanderbilt.edu). The research of this author was partially supported by NSF grant DMS-0600870.

the region occupied by the liquid phase. Consequently, the component J_1 is in contact with the liquid phase, and J_2 is in contact with the solid phase. The boundaries J_1 and J_2 , corresponding for instance to the walls of a container, are fixed, whereas Γ_0 will change as time evolves, due to solidification or liquidation of the two different phases.

Given $t \geq 0$, let $\Gamma(t)$ be the position of Γ_0 at time t , and let $V(\cdot, t)$ and $H(\cdot, t)$ be the normal velocity and the mean curvature of $\Gamma(t)$. Moreover, let $\Omega_1(t)$ and $\Omega_2(t)$ be the two regions in Ω separated by $\Gamma(t)$. According to our assumption, $\Omega_1(t)$ is the region occupied by the liquid phase, and $\Gamma(t)$ is a sharp interface which separates the liquid from the solid phase. Let $\nu(\cdot, t)$ be the outer unit normal field on $\Gamma(t)$ with respect to $\Omega_1(t)$. We shall use the convention that the normal velocity is positive if $\Omega_1(t)$ is expanding, and that the mean curvature is positive if the intersection of $\Omega_1(t)$ with a small ball centered at $\Gamma(t)$ is convex. Consequently, the normal velocity is positive if the liquid region is growing, ν points into the solid phase, and H is positive for a water ball surrounded by ice, and negative for an ice ball surrounded by water.

Here we concentrate on the case $J_1 = \emptyset$. Let Γ_0 and $u_0^i : \Omega_0^i \rightarrow \mathbb{R}$ be given, where u_0^1 and u_0^2 denote the initial temperatures of the liquid and solid phase, respectively. The strong formulation of the *two-phase Stefan problem with surface tension and kinetic undercooling* consists of finding a family $\Gamma := \{\Gamma(t); t \geq 0\}$ of hypersurfaces and functions $u_i : \cup_{t \geq 0} (\Omega_i(t) \times \{t\}) \rightarrow \mathbb{R}$, satisfying

$$(1.3) \quad \left\{ \begin{array}{ll} \kappa_i \partial_t u_i - d_i \Delta u_i = 0 & \text{in } \Omega_i(t), \\ \partial_\nu u_2 = 0 & \text{on } J_2, \\ u_i = \sigma H_\Gamma + \delta V & \text{on } \Gamma(t), \\ [d\partial_\nu u] = (\ell - [\kappa]u)V & \text{on } \Gamma(t), \\ u_i(0) = u_0^i & \text{in } \Omega_0^i, \\ \Gamma(0) = \Gamma_0, & \end{array} \right.$$

where $\kappa_i \geq 0$ is the heat capacity of phase i , d_i is its thermal conductivity coefficient, $\ell > 0$ is the latent heat per unit mass absorbed or released for melting or solidifying, $\sigma > 0$ is the surface tension, and $\delta \geq 0$ is the speed of kinetic undercooling. Moreover,

$$\begin{aligned} [\kappa] &:= \kappa_2 - \kappa_1, \\ [d\partial_\nu u] &:= d_2 \partial_\nu u_2 - d_1 \partial_\nu u_1 \end{aligned}$$

denote the jump of the heat capacities and the heat fluxes, respectively, across the interface $\Gamma(t)$. Note that $[\kappa] = \kappa_2 - \kappa_1 < 0$ is physically reasonable since in the liquid phase there are more degrees of freedom than in the solid phase; hence, the liquid phase can absorb more energy per unit mass. However, we do not assume $[\kappa] < 0$ in what follows. The condition $u_i = \sigma H_\Gamma$ on the free interface is usually called the Gibbs–Thomson law, and $u_i = \sigma H_\Gamma + \delta V$ the modified Gibbs–Thomson law, or the Gibbs–Thomson law with kinetic undercooling; see [2, 3, 16, 17, 19, 21, 24, 25, 32] for more information.

We refer to [12, 13, 14, 22, 23, 28, 29] for existence and regularity results for the Stefan problem with the Gibbs–Thomson law $u_i = \sigma H_\Gamma$ in case $\kappa_1 = \kappa_2$. The Stefan problem with surface tension and kinetic undercooling in case $\kappa_1 = \kappa_2$ has been studied in [5, 28, 29, 31]; see also [20] for the one-phase case.

It will be shown in [26] that the Stefan problem (1.3) has a unique local solution which is analytic in space and time, provided that the well-posedness condition

$$(1.4) \quad \ell - \sigma[\kappa]H_{\Gamma_0} > 0 \quad \text{in case } \delta = 0$$

is satisfied. On the other hand, if $\delta = 0$ and $\kappa_1 > \kappa_2$, problem (1.3) is not well-posed if H_{Γ_0} is too negative, that is, in case the solid region sharply protrudes into the liquid. Associated to the Stefan problem (1.3) is the energy functional

$$(1.5) \quad E(u(t), \Gamma(t)) := \int_{\Omega} \kappa u \, dx + \ell |\Omega_1(t)| = \int_{\Omega_1(t)} \kappa_1 u_1 \, dx + \int_{\Omega_2(t)} \kappa_1 u_2 \, dx + \ell |\Omega_1(t)|,$$

where $|\Omega_1(t)|$ is the volume of the region $\Omega_1(t)$. If (u, Γ) is a sufficiently smooth solution of (1.3), then we obtain

$$(1.6) \quad \begin{aligned} \frac{d}{dt} E(u(t), \Gamma(t)) &= \int_{\Omega_1(t)} \kappa_1 \partial_t u_1 \, dx + \int_{\Omega_2(t)} \kappa_1 \partial_t u_2 \, dx - [\kappa] \int_{\Gamma(t)} u V \, ds + \ell \int_{\Gamma(t)} V \, ds \\ &= \int_{\Gamma(t)} \left(-[d\partial_\nu u] - [\kappa]uV + \ell V \right) ds = 0, \end{aligned}$$

thus showing that energy is conserved.

If $\kappa_1 = \kappa_2 = 0$ and $\delta = 0$, then the resulting problem is the quasi-stationary Stefan problem with surface tension, which has also been termed the Mullins–Sekerka model (or the Hele–Shaw model with surface tension). Existence, uniqueness, regularity, and global existence of solutions for the quasi-stationary approximation have been investigated in [1, 4, 6, 8, 9, 10, 11, 15]. Existence and global existence of classical solutions for the quasi-stationary approximation with $\sigma > 0$ and $\delta > 0$ have been studied in [33, 20].

A major difficulty in the mathematical treatment of the Stefan problem (1.3) is due to fact that the boundary $\Gamma(t)$, and thus the geometry, is unknown and ever changing. A widely used method to overcome this inherent difficulty is to choose a fixed reference surface Σ and then represent the moving surface $\Gamma(t)$ as the graph of a function (which we will denote by $\rho = \rho(s, t)$) in normal direction of Σ . In this way, one obtains a time-dependent (unknown) diffeomorphism from Σ onto $\Gamma(t)$, and in the next step this diffeomorphism is extended to a diffeomorphism of fixed reference regions D^i onto the unknown domains $\Omega_i(t)$. The treatment of the moving boundary problem (1.3) then proceeds by transforming the equations into a new system of equations defined on the fixed domain $D_1 \cup D_2$ from which both the solution and the parameterizing function ρ have to be determined. In the context of the Stefan problem this approach has first been used by Hanzawa [18].

The same approach has also been used in [10, 11] for the quasi-stationary approximation of the Stefan problem with surface tension, and in [26] for the Stefan problem with surface tension. Once the transformed system has been obtained, one can study the mapping properties of the nonlinearities involved, and in particular, one can determine their linearizations; see [26] for more details.

In this paper, we assume that $\Gamma(t)$ does not touch the fixed boundary $J_2 = \partial\Omega$. Under this assumption, we will characterize all of the equilibrium states (u_1, u_2, Σ) of (1.3). In fact, it is easy to see that the equilibria are precisely given by

$$\Sigma = \bigcup_{j=1}^m S_R(x_j), \quad u_1 = u_2 = \sigma/R,$$

where $S_R(x_j)$ denotes disjoint spheres of the same radius R and centers x_j . This can be seen by the following arguments: the equilibria (u_1, u_2, Σ) of the Stefan problem (1.3)

are given by the system of equations

$$(1.7) \quad \begin{cases} -d_i \Delta u_i = 0 & \text{in } \Omega_i, \\ \partial_\nu u_i = 0 & \text{on } \partial\Omega, \\ u_i = \sigma H_\Sigma & \text{on } \Sigma, \\ [d\partial_\nu u] = 0 & \text{on } \Sigma. \end{cases}$$

Taking the inner product of (1.7)₁ with u_i , the divergence theorem and condition (1.7)₄ yield

$$\int_{\Omega_i} |\nabla u_i|^2 dx = 0;$$

hence u_i is constant on Ω_i . Equation (1.7)₃, in turn, shows that $u_1 = u_2$ and also that $H = u/\sigma$ is constant on Σ . But then, since Ω is bounded, Σ must be a sphere $S_R(x_0)$ centered at some point $x_0 \in \Omega$ with radius $R > 0$, if the phases are connected. Otherwise, again due to the boundedness of Ω , Σ is the union of finitely many spheres of the same radius $R > 0$. Here we concentrate on the case of connected phases. Thus there is an $(n + 1)$ -parameter family of equilibria, the parameters being the n coordinates of the center x_0 and the radius R .

We want to discuss the stability of those equilibria. The linearized problem (associated to the transformed equations) at such an equilibrium state is given by

$$(1.8) \quad \begin{cases} \kappa \partial_t v - d \Delta v = f & \text{in } (\Omega \setminus \Sigma) \times \mathbb{R}_+, \\ \partial_\nu v = 0 & \text{on } \partial\Omega \times \mathbb{R}_+, \\ v = \sigma \mathcal{A}_\Sigma \rho + \delta \partial_t \rho + g & \text{on } \Sigma \times \mathbb{R}_+, \\ l \partial_t \rho - [d\partial_\nu v] = h & \text{on } \Sigma \times \mathbb{R}_+, \\ v(0) = v_0 & \text{in } \Omega \setminus \Sigma, \\ \rho(0) = \rho_0 & \text{on } \Sigma; \end{cases}$$

see [26]. Here, $l = \ell - [\kappa]\sigma/R$, and the operator \mathcal{A}_Σ is given by

$$\mathcal{A}_\Sigma = -\frac{1}{n-1} \left(\frac{n-1}{R^2} + \Delta_\Sigma \right),$$

where Δ_Σ denotes the Laplace–Beltrami operator on Σ . This is the linearization of the mean curvature $H'(0)$ at the sphere Σ ; cf., e.g., Escher and Simonett [11]. Here we use the notation $v = v_1 \chi_{\Omega_1} + v_2 \chi_{\Omega_2}$, where χ_G denotes the characteristic function of a set G , and similarly $\kappa = \kappa_1 \chi_{\Omega_1} + \kappa_2 \chi_{\Omega_2}$ and $d = d_1 \chi_{\Omega_1} + d_2 \chi_{\Omega_2}$. Associated to the linearization (1.8) is the following eigenvalue problem:

$$(1.9) \quad \begin{cases} \lambda \kappa v - d \Delta v = 0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ v = \sigma \mathcal{A}_\Sigma \rho + \lambda \delta \rho & \text{on } \Sigma, \\ \lambda l \rho - [d\partial_\nu v] = 0 & \text{on } \Sigma, \end{cases}$$

where as before $l = \ell - [\kappa]\sigma/R$. We will now state the main results of this paper. We will formulate our results for a domain in \mathbb{R}^n for $n \in \mathbb{N}$, $n > 1$, although the physically relevant dimensions, naturally, are $n = 2, 3$.

THEOREM 1.1. *Suppose that the phases in the Stefan problem are connected. Then the following assertions hold:*

- (a) *The equilibrium states (without boundary contact) for problem (1.3) are given by*

$$(u, \Sigma), \quad \text{where } \Sigma = S_R(x_0) \quad \text{and} \quad u = \sigma/R,$$

with $S_R(x_0) \subset \Omega$ being the sphere with radius R and center x_0 .

- (b) *For $l > 0$, the eigenvalue problem (1.9) has countably many real eigenvalues of finite algebraic multiplicity.*
- (c) *0 is an eigenvalue of (1.9) with geometric multiplicity $(n+1)$. The (geometric) eigenspace is spanned by*

$$(-1, Y_0), (0, Y_1), \dots, (0, Y_n),$$

where $Y_0 = R^2/\sigma$, and where $Y_j, 1 \leq j \leq n$, are the spherical harmonics of degree 1 (normalized by the orthogonality condition $(Y_i|Y_j)_\Sigma = \delta_{ij}$).

- (d) *If $\sigma(\kappa|1)_\Omega \leq l|\Sigma|R^2$, then (1.9) has no positive eigenvalues.*
- (e) *If $\sigma(\kappa|1)_\Omega > l|\Sigma|R^2 > 0$, then (1.9) has exactly one positive, algebraically simple eigenvalue.*
- (f) *If $l < 0$ and $\delta > 0$, then (1.9) has at least one positive eigenvalue.*
- (g) *If $l < 0$ and $\delta = 0$, then the linearized problem (1.8) is not well-posed.*

Proof. The assertion in (a) has been proved above. We refer to Theorem 2.1 for a proof of assertions (b)–(e), and for additional information about the eigenvalue problem (1.9), for the case $l > 0$. The proof of (f) is given at the end of section 5, and (g) follows from [7]. \square

Remark 1.2. (a) If $l < 0$, then all equilibrium states are linearly unstable (and the linearized problem (1.8) is not even well-posed in case $\delta = 0$). Therefore we mainly concentrate on the case $l > 0$. Define then

$$\zeta := \frac{\sigma(\kappa|1)_\Omega}{l|\Sigma|R^2}, \quad \text{where } l = \ell - \frac{\sigma[\kappa]}{R}, \quad (\kappa|1)_\Omega := \int_\Omega \kappa \, dx.$$

According to Theorem 1.1.(d)–(e), we know that all eigenvalues of (1.9) are non-positive if $\zeta \leq 1$, and that there exists exactly one positive simple eigenvalue if $\zeta > 1$. We will refer to the case $\zeta \leq 1$ as a *stability condition*.

Observe that neither the thermal conductivity coefficients d_i nor the kinetic coefficient δ enters this stability condition, as it depends only on the heat capacities κ_i , the latent heat ℓ , the surface tension σ , and on the geometry. In particular, decreasing the size of a ball decreases its stability, as does increasing surface tension; see also Remark 1.5(a). We also mention that the stability condition $\zeta \leq 1$ is always valid in the quasi-stationary case $\kappa_i = 0$, i.e., for the Mullins–Sekerka problem.

(b) It will be shown in the forthcoming paper [27] that solutions for the Stefan problem (1.3) that start out close to an equilibrium (u, Σ) exist globally and converge towards an equilibrium state (u', Σ') as time goes to infinity, provided that $l > 0$ and $\zeta < 1$. This gives justice to the wording *stability condition* for the case $\zeta < 1$. We note again that $\zeta = 0$ if the heat capacities κ_i are zero, that is, in the quasi-stationary case. In this case, global existence and convergence to equilibria were obtained in [11, 20] by using a center-manifold analysis; see also [15] for a different approach in the one-phase case.

(c) If the Gibbs–Thomson condition on the free interface $\Gamma(t)$ is replaced by $u_i = 0$, then (1.3) is called the *(classical) Stefan model*. It should be observed that,

in contrast to the problem with surface tension, the classical Stefan problem does not admit nontrivial equilibrium states.

For $l > 0$, the results in Theorem 1.1 suggest that one eigenvalue, λ_* , crosses the imaginary axis at 0 as the quantity ζ increases and exceeds 1. According to part (c) of Theorem 1.1, 0 is always an eigenvalue with geometric multiplicity $(n + 1)$. This suggests that as the eigenvalue λ_* crosses through 0, the algebraic multiplicity of 0 raises by one, and then drops again as soon as the eigenvalue has crossed. This is exactly what happens, as will be proved in Theorem 2.1.

Another way to view and understand this situation can be gained from considering the following parameter-dependent eigenvalue problem:

$$(1.10) \quad \begin{cases} \lambda_* s \kappa v - d\Delta v = 0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ v = \sigma \mathcal{A}_\Sigma \rho + \lambda_* \delta \rho & \text{on } \Sigma, \\ \lambda_* l \rho - [d\partial_\nu v] = 0 & \text{on } \Sigma. \end{cases}$$

The following result will be proved in section 6.

THEOREM 1.3. *Let $l > 0$ and set $s_0 := l|\Sigma|R^2/\sigma(\kappa|1)_\Omega$. Then the following hold:*

(a) *The eigenvalue problem (1.10) has an analytic curve of solutions*

$$[s \mapsto (\lambda_*(s), v(s), \rho(s))], \quad s \in (s_0 - \varepsilon_0, \infty),$$

such that $\lambda_(s) > 0$ iff $s > s_0$, where ε_0 is an appropriate positive number.*

(b) *$\lambda_*(s)$ crosses the imaginary axis with positive speed at $s = s_0$.*

(c) *$[s \mapsto \lambda_*(s)]$ is strictly increasing.*

(d) *If $\delta > 0$, then $\lambda_*(s)$ is bounded above by $\sigma/\delta R^2$.*

(e) *If $\delta = 0$, then $\lambda_*(s) \rightarrow \infty$ as $s \rightarrow \infty$.*

Clearly, the eigenvalues of the modified problem (1.10) coincide with the eigenvalues of (1.9) if $s = 1$. In case $\zeta > 1$ we have $s_0 < 1$ and see that $\lambda = \lambda_*(1)$ is a (the only) positive eigenvalue of (1.9).

According to (1.5) an equilibrium state $(\sigma/R, S_R(x_0))$ for the Stefan problem (1.3) has energy

$$(1.11) \quad \begin{aligned} \phi(R) &:= E\left(\frac{\sigma}{R}, S_R(x_0)\right) = \frac{\sigma}{R}(\kappa|1)_\Omega + \ell|\Omega_1| \\ &= \frac{\sigma}{R}(\kappa_1|\Omega_1| + \kappa_2|\Omega_2|) + \ell|\Omega_1|, \end{aligned}$$

where $|\Omega_1| = R^n|B|$ and $|\Omega_2| = |\Omega| - R^n|B|$, with $|B|$ the volume of the unit ball. A straightforward computation shows that the function ϕ has a unique minimum. It is attained at the point R_* , where R_* is the unique solution of

$$(1.12) \quad \frac{\sigma(\kappa|1)_\Omega}{R^2} = \left(\ell - \frac{\sigma[\kappa]}{R}\right)|S_R|,$$

with $|S_R| = |S_R(x_0)|$ being the area of the sphere $S_R(x_0)$.

In the following, we denote by R_* the point where ϕ attains its (unique) minimum and by R^* the largest number R such that $\overline{B}_R(x_0) \subset \overline{\Omega}$, and we suppose that $R_* < R^*$. Then we have the following result; see also the stability diagram in Figure 1.

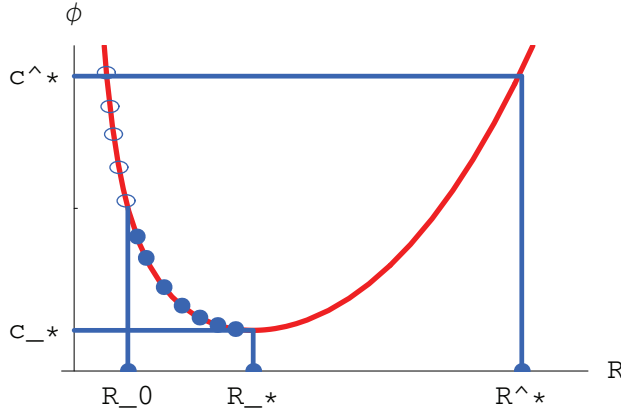


FIG. 1. Stability diagram for $\kappa_1 < \kappa_2$ and $\delta = 0$; circled: ill-posed, dotted: unstable.

COROLLARY 1.4. Let $c_* = \phi(R_*)$ be the minimum value of ϕ , and let $c^* = \phi(R^*)$. Moreover, let $c_0 = E(u_0, \Gamma_0) = \phi(R)$ be a given energy level.

- (a) If $c_0 < c_*$, then problem (1.3) does not admit equilibrium states.
- (b) If $c_* < c_0 < c^*$, then (1.3) admits two branches of equilibrium states. The branch of equilibria $(\sigma/R, S_R(x_0))$ with $0 < R < R_*$ is linearly unstable, whereas the branch with $R_* < R < R^*$ is linearly stable.
- (d) In case $c_0 = c_*$ or $c_0 > c^*$, the Stefan problem (1.3) admits one family of equilibrium states. All equilibria $(\sigma/R, S_R(x_0))$ with $\phi(R) > c^*$ are linearly unstable.
- (e) If $R_0 := [\kappa]\sigma/\ell > 0$ and $\delta = 0$, then the linearized problem is ill-posed for $R < R_0$.

Proof. This follows from Remark 1.2(a), (1.12), and the fact that $\phi'(R)$ is negative for $R < R_*$ and positive for $R_* < R < R^*$. \square

Remark 1.5. (a) We can show that R_* is increasing with respect to σ (and, for that matter, also with respect to $[\kappa]$). In order to see this, let $R_*(\sigma)$ denote the unique solution of (1.12). Then we have $R'_*(\sigma) > 0$. For this we note that (1.12) can be written as

$$(1.13) \quad n\ell|B|R_*^{n+1}(\sigma) - \sigma(\kappa_2|\Omega| + (n-1)[\kappa]|B|R_*^n(\sigma)) = 0.$$

Taking the derivative of this equation with respect to σ yields

$$\begin{aligned} & n(n+1)\ell|B|R_*^{n-1} \left(R_* - \frac{(n-1)\sigma[\kappa]}{(n+1)\ell} \right) R'_*(\sigma) \\ & = \kappa_2|\Omega| + (n-1)[\kappa]|B|R_*^n = (\ell/\sigma)n|B|R_*^{n+1}(\sigma) > 0. \end{aligned}$$

Note that we have used (1.13) for the last equality. It remains to be observed that the parenthetical expression in front of $R'_*(\sigma)$ is always positive. This is clear in case $[\kappa] \leq 0$ and follows from the fact that R_* is always greater than $R_0 = \sigma[\kappa]/\ell$ in case $[\kappa] > 0$. We therefore see that increasing σ increases $R_*(\sigma)$, showing that spheres with a fixed radius R can lose stability as σ increases.

(b) If one considers the case where the domain Ω_1 is occupied by the solid phase, and Ω_2 by the fluid phase, then the third and fourth lines in (1.3) must be replaced

by

$$(1.14) \quad \begin{aligned} u_i &= -\sigma H_\Gamma - \delta V && \text{on } \Gamma(t), \\ -[d\partial_\nu u] &= (\ell + [\kappa]u)V && \text{on } \Gamma(t), \end{aligned}$$

and the energy functional by

$$E(u(t), \Gamma(t)) := \int_\Omega \kappa u \, dx + \ell |\Omega_2(t)| = \int_{\Omega_1(t)} \kappa_1 u_1 \, dx + \int_{\Omega_2(t)} \kappa_1 u_2 \, dx + \ell |\Omega_2(t)|,$$

while all other conventions are left unchanged. Thus, one formally has to switch signs in the normal ν and in ℓ and $[\kappa]$. Then all of the results and assertions stated in this paper remain valid for the equilibrium states $(-\sigma/R, S_R(x_0))$.

The plan of this paper is the following. In section 2 we will state a more general and concise version of Theorem 1.1; its proof will be given in sections 3–5. Finally, in section 6 we will prove Theorem 1.3.

2. Main theorem. In this section we will introduce an appropriate functional analytic setting to study the eigenvalue problem (1.9). We always assume $l > 0$ except when proving (f) of Theorem 1.1.

For the case $\delta = 0$ we define the operator L_0 on $E_0 := L_p(\Omega) \times W_p^{2-2/p}(\Sigma)$ by means of

$$\begin{aligned} D(L_0) &:= \{(v, \rho) \in W_p^2(\Omega \setminus \Sigma) \times W_p^{4-1/p}(\Sigma) : [d\partial_\nu v] \in W_p^{2-2/p}(\Sigma), \\ &\quad \partial_\nu v = 0 \text{ on } \partial\Omega, [v] = 0 \text{ on } \Sigma, v = \sigma \mathcal{A}_\Sigma \rho \text{ on } \Sigma\}, \\ L_0(v, \rho) &:= ((d/\kappa)\Delta v, [(d/l)\partial_\nu v]), \quad (v, \rho) \in D(L_0). \end{aligned}$$

In case $\delta > 0$, we instead set $E_\delta := L_p(\Omega) \times W_p^{1-1/p}(\Sigma)$, and we define L_δ by

$$\begin{aligned} D(L_\delta) &:= \{(v, \rho) \in W_p^2(\Omega \setminus \Sigma) \times W_p^{3-1/p}(\Sigma) : \\ &\quad \partial_\nu v = 0 \text{ on } \partial\Omega, [v] = 0 \text{ on } \Sigma, v - (\delta/l)[d\partial_\nu v] = \sigma \mathcal{A}_\Sigma \rho \text{ on } \Sigma\}, \\ L_\delta(v, \rho) &:= ((d/\kappa)\Delta v, [(d/l)\partial_\nu v]), \quad (v, \rho) \in D(L_\delta). \end{aligned}$$

We remark that L_0 and L_δ differ only by their respective domains of definition. It will be shown in [7] that the operators L_δ generate an analytic semigroup on E_δ . This property, in conjunction with the spectral information contained in the next theorem, will be crucial in proving global existence and convergence of solutions for problem (1.3) that start out close to an equilibrium, which will be provided in [27]; see also Remark 1.2(c).

THEOREM 2.1. *Suppose $1 < p < \infty$, and let $l > 0$. For $\delta \geq 0$ let L_δ be defined as above.*

- (a) *The spectrum of L_δ consists of countably many real eigenvalues of finite algebraic multiplicity and is independent of p .*
- (b) *0 is an eigenvalue of L_δ with geometric multiplicity $(n + 1)$. The null space of L_δ is spanned by*

$$(2.1) \quad (-1, Y_0), (0, Y_1), \dots, (0, Y_n),$$

where $Y_0 = R^2/\sigma$, and where Y_j , $1 \leq j \leq n$, are the spherical harmonics of degree 1 (normalized by the orthogonality condition $(Y_i|Y_j)_\Sigma = \delta_{ij}$).

(c) Suppose that the degeneracy condition

$$(2.2) \quad (\kappa|1)_\Omega := \kappa_1|\Omega_1| + \kappa_2|\Omega_2| = l|\Sigma|R^2/\sigma$$

holds. Then the eigenvalue 0 has algebraic multiplicity $(n + 2)$.

(d) If the degeneracy condition (2.2) does not hold, then 0 is semi-simple; that is, $N(L_\delta^2) = N(L_\delta)$.

(e) If $\sigma(\kappa|1)_\Omega \leq l|\Sigma|R^2$, then L_δ has no positive eigenvalues.

(f) If $\sigma(\kappa|1)_\Omega > l|\Sigma|R^2$, then L_δ has exactly one positive simple eigenvalue.

Proof. (a) By the compact embeddings $D(L_\delta) \hookrightarrow E_\delta$, the spectrum of L_δ consists of eigenvalues of finite algebraic multiplicity. The assertion that all eigenvalues are real will be proved in section 4.

Let $1 < p < \infty$ be fixed, and suppose that λ is an eigenvalue of L_δ with a corresponding eigenfunction (v, ρ) . Then $v \in W_p^2(\Omega \setminus \Sigma)$, and v solves the elliptic transmission problem

$$\begin{cases} \kappa\lambda v - d\Delta v = 0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ [v] = 0 & \text{on } \Sigma, \\ -[d\partial_\nu v] = -\lambda l\rho & \text{on } \Sigma, \end{cases}$$

with $\rho \in W_p^{4-\text{sign}(\delta)-1/p}(\Sigma)$, where $\text{sign}(\delta) = 1$ if $\delta > 0$, and $\text{sign}(\delta) = 0$ if $\delta = 0$. Due to Sobolev's imbedding theorem we have that $\rho \in W_{p_1}^{1-1/p_1}(\Sigma)$, where $p_1 \in (p, \infty)$ is appropriately chosen. Proposition 5.1 then yields $v \in W_{p_1}^2(\Omega \setminus \Sigma)$. Next, we recall that ρ satisfies

$$\sigma\mathcal{A}_\Sigma\rho = v - (\delta/l)[d\partial_\nu v] =: h \quad \text{on } \Sigma.$$

Since $v \in W_{p_1}^2(\Omega \setminus \Sigma)$ we see that $h \in W_{p_1}^{2-\text{sign}(\delta)-1/p_1}(\Sigma)$, and we obtain from the properties of the elliptic differential operator \mathcal{A}_Σ that $\rho \in W_{p_1}^{4-\text{sign}(\delta)-1/p_1}(\Sigma)$. The arguments given above can now be iterated a finite number of times to show that

$$(v, \rho) \in W_q^2(\Omega \setminus \Sigma) \times W_q^{4-\text{sign}(\delta)-1/q}(\Sigma)$$

for any fixed $q > p$. Clearly, this is also true for any $q < p$. We have, thus, shown that the spectrum of L_δ is independent of p . The properties listed in (b)–(d) are proved in section 3, and assertion (e) is shown in section 4 while (f) is established in section 5. \square

PROPOSITION 2.2. *Let $1 < p < \infty$. Suppose that $(\lambda, v, \rho) \in \mathbb{R} \times W_p^2(\Omega \setminus \Sigma) \times W_p^2(\Sigma)$ solves the eigenvalue problem (1.9). Then the functions (v, ρ) are smooth; that is,*

$$v|_{\Omega_i} \in C^\infty(\bar{\Omega}_i), \quad \rho \in C^\infty(\Sigma).$$

Proof. This follows from a similar bootstrapping argument as in the proof of Theorem 2.1(a), based on regularity properties of the elliptic transmission problems (3.4) and (5.2), and regularity properties of the differential operator \mathcal{A}_Σ . \square

Due to Theorem 2.1 and Proposition 2.2 we may restrict our attention to the eigenvalue problem (1.9) in the Hilbert space setting of $L_2(\Omega) \times L_2(\Sigma)$. In the following, we use the notation $(\cdot|\cdot)_\Omega$ and $\|\cdot\|_\Omega$ for the inner product and the norm in $L_2(\Omega)$, respectively, and similarly for $L_2(\Sigma)$.

3. The trivial eigenvalue. Let us first look at the eigenvalue problem (1.9) with $\lambda = 0$. Obviously, here $l \in \mathbb{R}$ can be arbitrary, and also $\delta \in \mathbb{R}$. For this purpose we recall some properties of the operator A_Σ .

PROPOSITION 3.1. *Let $\Sigma = S_R(x_0) \subset \mathbb{R}^n$ be a sphere of radius R and center x_0 , and let*

$$A_\Sigma = -\frac{1}{n-1} \left(\frac{n-1}{R^2} + \Delta_\Sigma \right)$$

be defined on $L_2(\Sigma)$ with domain $W_2^2(\Sigma)$. Then the following assertions hold:

- (a) A_Σ is self-adjoint. Its spectrum consists of countably many eigenvalues $\lambda_k = \frac{1}{(n-1)R^2} (k(k+n-2) - (n-1))$ with $k \geq 0$. The eigenfunctions are given by the spherical harmonics of degree k .
- (b) The kernel of A_Σ is given by $N(A_\Sigma) = \text{span}\{Y_1, \dots, Y_n\}$, where Y_j denotes the spherical harmonics of degree 1 on Σ , normalized by $(Y_i|Y_j)_\Sigma = \delta_{ij}$.
- (c) The range of A_Σ , $R(A_\Sigma)$ is closed, and we have $L_2(\Sigma) = N(A_\Sigma) \oplus R(A_\Sigma)$.
- (d) There is precisely one negative eigenvalue, namely $-1/R^2$, with eigenfunction 1, which is simple.
- (e) A_Σ is positive semi-definite on $L_{2,0}(\Sigma) = \{\rho \in L_2(\Sigma) : (\rho|1)_\Sigma = 0\}$ and positive definite on

$$L_{2,0}(\Sigma) \cap R(A_\Sigma) = \{\rho \in L_2(\Sigma) : (\rho|1)_\Sigma = (\rho|Y_j)_\Sigma = 0, j = 1, \dots, n\}.$$

Proof. We can assume, without loss of generality, that $\Sigma = S_R(0) = R\mathbb{S}^{n-1}$, where \mathbb{S}^{n-1} denotes the standard unit sphere in \mathbb{R}^n . Let $\Phi : \Sigma \rightarrow \mathbb{S}^{n-1}$ be defined by $p \mapsto (1/R)p$. Then Φ is a smooth diffeomorphism of Σ into \mathbb{S}^{n-1} , and one readily verifies that

$$(3.1) \quad (g|h)_{L_2(\Sigma)} = R^{n-1} (\Phi_*g|\Phi_*h)_{L_2(\mathbb{S}^{n-1})}, \quad \Delta_\Sigma = (1/R^2) \Phi^* \Delta_{\mathbb{S}^{n-1}} \Phi_*$$

where Φ^* and Φ_* are the pull-back and push-forward operators, respectively. We then have

$$(3.2) \quad (\lambda - A_\Sigma)\rho = 0 \iff \left(\lambda + \frac{1}{(n-1)R^2} ((n-1) + \Delta_{\mathbb{S}^{n-1}}) \right) \Phi_*\rho = 0,$$

which shows that λ is an eigenvalue of A_Σ iff

$$(3.3) \quad \lambda = \frac{1}{(n-1)R^2} (\mu - (n-1))$$

with μ an eigenvalue of $-\Delta_{\mathbb{S}^{n-1}}$. The assertions in (a)–(b) and (d)–(e) follow now from (3.1)–(3.3) and well-known results for the Laplace–Beltrami operator on \mathbb{S}^{n-1} ; see, for instance, [30, section 31]. Since A_Σ has compact resolvent we conclude that $R(A_\Sigma)$ is closed, and the fact that A_Σ is self-adjoint then implies the remaining assertion in (c). \square

Before we proceed we need the following result on the elliptic transmission problem:

$$(3.4) \quad \begin{cases} -d\Delta v = f & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ [v] = 0 & \text{on } \Sigma, \\ -[d\partial_\nu v] = g & \text{on } \Sigma. \end{cases}$$

PROPOSITION 3.2. *Let $1 < p < \infty$. Then the following hold:*

- (a) *The transmission problem (3.4) has a solution $v \in W_p^2(\Omega \setminus \Sigma)$ if and only if $(f, g) \in L_p(\Omega) \times W_p^{1-1/p}(\Sigma)$ and the compatibility condition*

$$(f|1)_\Omega + (g|1)_\Sigma = 0$$

is satisfied. The solution is unique with the normalization $(\kappa|v)_\Omega = 0$.

- (b) *Let $v = T_0g$ be the unique solution of (3.4) with $f = 0$, $(g|1)_\Sigma = 0$, and $(\kappa|v)_\Omega = 0$. Then T_0 is self-adjoint and positive definite on $L_2(\Sigma)$; that is, there exists a positive constant $c = c(d_i, \Omega_i)$ such that*

$$(T_0g|g)_{L_2(\Sigma)} \geq c \|g\|_{L_2(\Sigma)}^2, \quad g \in W_2^{1/2}(\Sigma).$$

Proof. (a) This proof follows from known results in elliptic theory since the Lopatinskii–Shapiro conditions are satisfied.

- (b) Let $g, h \in W_2^{1/2}(\Sigma)$ be given. Then we have

$$\begin{aligned} (T_0g|h)_\Sigma &= (T_0g|[-d\partial_\nu T_0h])_\Sigma = (d\nabla T_0g|\nabla T_0h)_\Omega \\ &= (-[d\partial_\nu T_0g]|T_0h)_\Sigma = (g|T_0h)_\Sigma, \end{aligned}$$

thus showing that T_0 is symmetric. For $v := T_0g$ the computation above yields

$$(T_0g|g)_\Sigma = (d\nabla v|\nabla v)_\Omega.$$

On the other hand, setting $v_i = v|_{\Omega_i}$ we obtain

$$\begin{aligned} \|g\|_{L_2(\Sigma)} &= \|d_1\partial_\nu v_1 - d_2\partial_\nu v_2\|_{L_2(\Sigma)} \leq c(\|v_1\|_{W_2^2(\Omega_1)} + \|v_2\|_{W_2^2(\Omega_2)}) \\ &\leq c(\|v_1\|_{L_2(\Omega_1)} + \|\Delta v_1\|_{L_2(\Omega_1)} + \|v_2\|_{L_2(\Omega_2)} + \|\Delta v_2\|_{L_2(\Omega_2)}) \\ &= c\|v\|_{L_2(\Omega)} \leq c\|v\|_{W_2^1(\Omega)} \leq c\|\nabla v\|_{L_2(\Omega)} \leq c(T_0g|g)_\Sigma^{1/2}. \end{aligned}$$

Here we used the fact that

$$v = T_0g \in W_2^1(\Omega) \cap W_2^2(\Omega \setminus \Sigma).$$

Moreover, we used that $(\|\cdot\|_{L_2(\Omega_i)} + \|\Delta \cdot\|_{L_2(\Omega_i)})$ defines an equivalent norm on $W_2^2(\Omega_i)$, and also that $\|\nabla u\|_{L_2(\Omega)}$ defines an equivalent norm on $W_2^1(\Omega)$ for all functions $u \in W_2^1(\Omega)$ with $(\kappa|u)_\Omega = 0$. This completes the proof of Proposition 3.2. \square

We are now ready to establish the assertions (b)–(d) of Theorem 2.1.

- (b) Suppose that (v, ρ) is a solution of (1.9) with $\lambda = 0$. Then, taking the inner product of (1.9)₁ with v , the divergence theorem and (1.9)_{2,4} show that v is constant on $\Omega \setminus \Sigma$; hence, v is constant on Ω and $v = \sigma\mathcal{A}_\Sigma\rho$ due to (1.9)₃. A special solution of this problem is $\rho_0 = -R^2v/\sigma$, and the solutions of the corresponding homogeneous equation are the spherical harmonics Y_j on Σ for $j = 1, \dots, n$. Thus we obtain an $(n + 1)$ -dimensional null space spanned by (2.1), which proves Theorem 2.1(b).

This null space is tangent to the $(n+1)$ -dimensional manifold of equilibria, where $(0, Y_j)$ corresponds to the center x_0 , and $(-1, Y_0)$ corresponds to the radius R . Note that the null spaces of L_0 and L_δ , $\delta > 0$, coincide.

- (c) Suppose that (2.2) holds. Then there exists a pair $(v^*, \rho^*) \in N(L_\delta^2) \setminus N(L_\delta)$. Indeed, this can be seen as follows: we first solve (3.4) with $(f, g) = (-\kappa, lR^2/\sigma)$. According to Proposition 3.2, this problem has a unique solution v_0 with $(\kappa|v_0)_\Omega = 0$ since the necessary compatibility condition is precisely (2.2).

Set $v^* = v_0 + l \sum_{j=1}^n \alpha_j T_0 Y_j$. Then v^* satisfies

$$(3.5) \quad \begin{cases} -(d/\kappa)\Delta v^* = -1 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v^* = 0 & \text{on } \partial\Omega, \\ [v^*] = 0 & \text{on } \Sigma, \\ -[(d/l)\partial_\nu v^*] = R^2/\sigma + \sum_{j=1}^n \alpha_j Y_j & \text{on } \Sigma. \end{cases}$$

We now want to solve $v^* = \sigma A_\Sigma \rho + (\delta/l)[d\partial_\nu v^*]$ in terms of ρ ; that is, we consider the problem

$$\sigma A_\Sigma \rho = v^* - (\delta/l)[d\partial_\nu v^*] =: h.$$

According to Proposition 3.1(b)–(c) this problem has a solution ρ^* iff $(h|Y_i)_\Sigma = 0$ for $i = 1, \dots, n$. The conditions $(h|Y_i)_\Sigma = 0$ will then be employed to determine the coefficients α_j . A short computation yields

$$(h|Y_i)_\Sigma = 0 \iff \delta\alpha_i + \sum_{j=1}^n l(T_0 Y_j | Y_i)_\Sigma \alpha_j = -(v_0 | Y_i)_\Sigma, \quad i = 1, \dots, n.$$

Since T_0 is self-adjoint and positive definite on $L_2(\Sigma)$, there exists a unique solution of this system, as we shall see in (6.7). Due to $\sigma A_\Sigma(R^2/\sigma + \sum_{j=1}^n \alpha_j Y_j) = -1$ (see Proposition 3.1), we conclude that $(v^*, \rho^*) \in D(L_\delta^2)$. It is then easy to see that $L_\delta(v^*, \rho^*) \neq (0, 0)$ and $L_\delta^2(v^*, \rho^*) = (0, 0)$. These facts in combination with part (b) show that

$$(3.6) \quad N(L_\delta^2) = N(L_\delta) \oplus \text{span}\{(v^*, \rho^*)\},$$

in the degenerate case where (2.2) holds.

Next we show, still in the degenerate case (2.2), that $N(L_\delta^3) = N(L_\delta^2)$. In fact, if $(v, \rho) \in N(L_\delta^3)$, then $L_\delta(v, \rho) = (v_N, \rho_N) + \beta(v^*, \rho^*)$ for some $(v_N, \rho_N) \in N(L_\delta)$ and some scalar β . For solvability of this equation, the compatibility condition

$$(\kappa(v_N + \beta v^*)|1)_\Omega + l(\rho_N + \beta \rho^*|1)_\Sigma = 0$$

must be valid. Due to the degeneracy condition we have $(\kappa v_N|1)_\Omega + l(\rho_N|1)_\Sigma = 0$, and the compatibility condition is reduced to $\beta\{(\kappa v^*)|1)_\Omega + l(\rho^*|1)_\Sigma\} = 0$. Using the property that $-(d/\kappa)\Delta v^* = -1$ (see (3.5)), we obtain

$$\begin{aligned} -\{(\kappa v^*)|1)_\Omega + l(\rho^*|1)_\Sigma\} &= -(d\Delta v^*|v^*)_\Omega - l(\rho^*|1)_\Sigma \\ &= \|\sqrt{d} \nabla v^*\|_\Omega^2 + ([d\partial_\nu v^*]|v^*)_\Sigma - l(\rho^*|1)_\Sigma \\ &= \|\sqrt{d} \nabla v^*\|_\Omega^2 - l \left(R^2/\sigma + \sum_{j=1}^n \alpha_j Y_j | \sigma A_\Sigma \rho^* \right)_\Sigma + (\delta/l)\|[d\partial_\nu v^*]\|_\Sigma^2 - l(\rho^*|1)_\Sigma \\ &= \|\sqrt{d} \nabla v^*\|_\Omega^2 + (\delta/l)\|[d\partial_\nu v^*]\|_\Sigma^2 \end{aligned}$$

since A_Σ is self-adjoint on $L_2(\Sigma)$ and $\sigma A_\Sigma(R^2/\sigma + \sum_{j=1}^n \alpha_j Y_j) = -1$; see Proposition 3.1. This implies $\beta = 0$, i.e., $(v, \rho) \in N(L_\delta^2)$, thus establishing Theorem 2.1(c).

(d) Let us examine when the eigenvalue $\lambda_0 = 0$ of L_δ is semi-simple. Assume that $(v, \rho) \in D(L_\delta^2)$ is such that $L_\delta^2(v, \rho) = 0$. Then

$$L_\delta(v, \rho) = \alpha_0(-1, Y_0) + \sum_{j=1}^n \alpha_j(0, Y_j).$$

This implies that

$$\left\{ \begin{array}{ll} -(d/\kappa)\Delta v = \alpha_0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ [v] = 0 & \text{on } \Sigma, \\ -[(d/l)\partial_\nu v] = -\sum_{j=0}^n \alpha_j Y_j & \text{on } \Sigma, \\ v = \sigma \mathcal{A}_\Sigma \rho + (\delta/l)[d\partial_\nu v] & \text{on } \Sigma. \end{array} \right.$$

According to Proposition 3.2 we necessarily have

$$\alpha_0(\kappa_1|\Omega_1| + \kappa_2|\Omega_2|) = l \sum_{j=0}^n \alpha_j (Y_j|1)_\Sigma = l\alpha_0|\Sigma|R^2/\sigma,$$

since the mean value of Y_j over Σ is zero for $j \geq 1$. Assuming the nondegeneracy condition

$$(3.7) \quad (\kappa|1)_\Omega := \kappa_1|\Omega_1| + \kappa_2|\Omega_2| \neq l|\Sigma|R^2/\sigma,$$

we conclude that $\alpha_0 = 0$. But then

$$0 = - \int_\Omega d\Delta v v \, dx = \int_\Omega d|\nabla v|^2 \, dx + \int_\Sigma [d\partial_\nu v]v \, ds,$$

which further yields

$$0 = \|\sqrt{d}\nabla v\|_\Omega^2 + l\sigma \sum_{j=1}^n \alpha_j (Y_j|\mathcal{A}_\Sigma \rho)_\Sigma + (\delta/l)\|[d\partial_\nu v]\|_\Sigma^2 = \|\sqrt{d}\nabla v\|_\Omega^2 + (\delta/l)\|[d\partial_\nu v]\|_\Sigma^2$$

since \mathcal{A}_Σ is self-adjoint on $L_2(\Sigma)$ and $A_\Sigma Y_j = 0$ for $j \geq 1$. We conclude that v is constant in Ω and that $0 = [d\partial_\nu v] = l \sum_{j=1}^n \alpha_j Y_j$; hence, $\alpha_j = 0$ for all j . This shows that $\lambda_0 = 0$ is a semi-simple eigenvalue of L_δ , that is, $N(L_\delta^2) = N(L_\delta)$ for $\delta \geq 0$, provided the nondegeneracy condition (3.7) is valid, and this proves the assertion of Theorem 2.1(d).

4. Nontrivial eigenvalues. Now we consider the eigenvalue problem (1.9) for $\lambda \in \mathbb{C}$, $\lambda \neq 0$, in case $l > 0$. Suppose that $\lambda \neq 0$ is an eigenvalue with nontrivial eigenfunction (v, ρ) . Taking the inner product in $L_2(\Omega)$ of the first equation in (1.9) with v and using the divergence theorem, we get

$$\begin{aligned} \lambda \|\sqrt{\kappa}v\|_\Omega^2 &= (d\Delta v|v)_\Omega = -\|\sqrt{d}\nabla v\|_\Omega^2 - ([d\partial_\nu v]|v)_\Sigma \\ &= -\|\sqrt{d}\nabla v\|_\Omega^2 - (\delta/l)\|[d\partial_\nu v]\|_\Sigma^2 - \lambda\sigma(\rho|\mathcal{A}_\Sigma \rho)_\Sigma; \end{aligned}$$

hence, we obtain the identity

$$(4.1) \quad \lambda \left(\|\sqrt{\kappa}v\|_\Omega^2 + l\sigma(\rho|\mathcal{A}_\Sigma \rho)_\Sigma \right) + \|\sqrt{d}\nabla v\|_\Omega^2 + (\delta/l)\|[d\partial_\nu v]\|_\Sigma^2 = 0.$$

If $\text{Im } \lambda \neq 0$, then $\|\sqrt{\kappa}v\|_\Omega^2 + l\sigma(\rho|\mathcal{A}_\Sigma \rho)_\Sigma = 0$; hence, $\|\sqrt{d}\nabla v\|_\Omega^2 + (\delta/l)\|[d\partial_\nu v]\|_\Sigma^2 = 0$. We conclude that v is constant and (1.9)_{1,4} now implies that $(v, \rho) = (0, 0)$ since

$\lambda \neq 0$. Therefore the eigenvalues and eigenfunctions are real, thus establishing Theorem 2.1(a).

Using $\lambda \rho = [d\partial_\nu v]$ and the fact that λ is real, we may rewrite (4.1) as

$$(4.2) \quad \lambda \left(\|\sqrt{\kappa}v\|_\Omega^2 + l\sigma(\rho|_{A_\Sigma\rho})_\Sigma + \lambda l\delta\|\rho\|_\Sigma^2 \right) + \|\sqrt{d}\nabla v\|_\Omega^2 = 0.$$

Integrating the eigenvalue equation (1.9) we obtain

$$\lambda(v|\kappa)_\Omega = (d\Delta v|1)_\Omega = -([d\partial_\nu v]|1)_\Sigma = -\lambda l(\rho|1)_\Sigma;$$

hence, dividing by λ ,

$$(4.3) \quad (v|\kappa)_\Omega + l(\rho|1)_\Sigma = 0.$$

Splitting $\rho = \rho_0 + \bar{\rho}$ and $v = v_0 + \bar{v}$, where $(\rho_0|1)_\Sigma = (v_0|\kappa)_\Omega = 0$, from (4.2) and (4.3) we derive an identity equivalent to (4.2), namely,

$$(4.4) \quad \lambda \left(\|\sqrt{\kappa}v_0\|_\Omega^2 + l\sigma(\rho_0|_{A_\Sigma\rho_0})_\Sigma + \lambda l\delta\|\rho_0\|_\Sigma^2 \right) + \|\sqrt{d}\nabla v_0\|_\Omega^2 + \lambda \left(\lambda\delta + \frac{l|\Sigma|}{(\kappa|1)_\Omega} - \frac{\sigma}{R^2} \right) l|\Sigma|\bar{\rho}^2 = 0.$$

Since A_Σ is positive semi-definite on $L_{2,0}(\Sigma)$, the L_2 -functions with mean zero, we see that in case $\lambda > 0$, (4.4) implies $v = \text{constant}$, and hence $(v, \rho) = (0, 0)$, provided that

$$(4.5) \quad (\kappa|1)_\Omega \leq l|\Sigma|R^2/\sigma.$$

Consequently, (1.9) cannot have positive eigenvalues if the stability condition (4.5) is satisfied, thus proving the assertion of Theorem 2.1(e).

5. The unstable eigenvalue. As far as we know, for $l > 0$ and

$$(5.1) \quad \zeta := \frac{\sigma(\kappa|1)_\Omega}{l|\Sigma|R^2} \leq 1$$

there are no positive eigenvalues; however, the algebraic eigenspace of L_δ rises in dimension by one when ζ becomes 1. This indicates that for $\zeta > 1$ there is exactly one algebraically simple eigenvalue $\lambda_* > 0$. We want to prove that this is indeed the case. In order to do so, we consider the following transmission problem:

$$(5.2) \quad \begin{cases} \lambda\kappa v - d\Delta v = f & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ [v] = 0 & \text{on } \Sigma, \\ -[d\partial_\nu v] = g & \text{on } \Sigma. \end{cases}$$

Then the following result holds.

PROPOSITION 5.1. *Let $1 < p < \infty$ and $\text{Re } \lambda > 0$. Then the following hold:*

- (a) *Problem (5.2) has precisely one solution $v \in W_p^1(\Omega) \cap W_p^2(\Omega \setminus \Sigma)$ iff $(f, g) \in L_p(\Omega) \times W_p^{1-1/p}(\Sigma)$.*
- (b) *Let T_λ be the solution operator for (5.2) with $f = 0$. Given any number $\theta \in (0, \pi)$, there exist positive numbers $\lambda_0 = \lambda_0(\theta, d_i, \kappa_i, \Omega_i)$ and $M_0 = M_0(\theta, d_i, \kappa_i, \Omega_i)$ such that*

$$\|T_\lambda g\|_{L_p(\Sigma)} \leq M_0|\lambda|^{-1/2}\|g\|_{L_p(\Sigma)}$$

for $g \in W_p^{1-1/p}(\Sigma)$, whenever $|\lambda| \geq \lambda_0$ and $|\arg \lambda| \leq \theta$.

(c) For $\lambda > 0$, T_λ is positive definite on $L_2(\Sigma)$; that is, there exists a positive constant $\beta = \beta(d_i, \kappa_i, \Omega_i)$ such that

$$(T_\lambda g|g)_{L_2(\Sigma)} \geq \beta \frac{\sqrt{\lambda}}{1 + \lambda} \|g\|_{L_2(\Sigma)}^2, \quad g \in W_2^{1/2}(\Sigma).$$

Proof. (a) This proof follows from known results in elliptic theory.

(b) Suppose that $\Omega_1 = \mathbb{R}_-^n$ and $\Omega_2 = \mathbb{R}_+^n$, with $\mathbb{R}_\pm^n = \{(x', x_n) \in \mathbb{R}^n : \pm x_n < 0\}$. Then one readily obtains that

$$T_\lambda g|_{\mathbb{R}^{n-1} \times \{0\}} = \mathcal{F}^{-1}(m_\lambda \mathcal{F}g),$$

where \mathcal{F} denotes the Fourier transform in the tangential variables, and where

$$m_\lambda(\xi) = \frac{1}{\sqrt{d_1} \sqrt{\kappa_1 \lambda + d_1 |\xi|^2} + \sqrt{d_2} \sqrt{\kappa_2 \lambda + d_2 |\xi|^2}}.$$

The assertion then follows from Mikhlin's multiplier theorem. The general case can be obtained by the usual procedure of localization.

(c) Let $g, h \in W_2^{1/2}(\Sigma)$ be given. Then we have

$$\begin{aligned} (T_\lambda g|h)_\Sigma &= (T_\lambda g|[-d\partial_\nu T_\lambda h])_\Sigma = \lambda(\kappa T_\lambda g|T_\lambda h)_\Omega + (d \nabla T_\lambda g|\nabla T_\lambda h)_\Omega \\ &= (-[d\partial_\nu T_\lambda g]|T_\lambda h)_\Sigma = (g|T_\lambda h)_\Sigma, \end{aligned}$$

showing that $T_\lambda^* = T_\lambda$, in particular, that T_λ is symmetric for $\lambda > 0$. For $v := T_\lambda g$, $\lambda > 0$, the computation above yields

$$(5.3) \quad (T_\lambda g|g)_\Sigma = \lambda(\kappa v|v)_\Omega + (d \nabla v|\nabla v)_\Omega.$$

Setting $v_i = v|_{\Omega_i}$, we conclude similarly as in the proof of Proposition 3.2 that

$$\begin{aligned} \|g\|_{L_2(\Sigma)} &= \|d_1 \partial_\nu v_1 - d_2 \partial_\nu v_2\|_{L_2(\Sigma)} \leq c(\|d_1 v_1\|_{W_2^2(\Omega_1)} + \|d_2 v_2\|_{W_2^2(\Omega_2)}) \\ &\leq c(\|d_1 v_1\|_{L_2(\Omega_1)} + \|d_1 \Delta v_1\|_{L_2(\Omega_1)} + \|d_2 v_2\|_{L_2(\Omega_2)} + \|d_2 \Delta v_2\|_{L_2(\Omega_2)}) \\ &= c(\|d_1 v_1\|_{L_2(\Omega_1)} + \lambda \|\kappa_1 v_1\|_{L_2(\Omega_1)} + \|d_2 v_2\|_{L_2(\Omega_2)} + \lambda \|\kappa_2 v_2\|_{L_2(\Omega_2)}) \\ &\leq c_\lambda \sqrt{\lambda} \|v\|_{L_2(\Omega)} \leq c_\lambda (\sqrt{\lambda} \|v\|_{L_2(\Omega)} + \|\nabla v\|_{L_2(\Omega)}) \leq c_\lambda (T_\lambda g|g)_\Sigma^{1/2}, \end{aligned}$$

where $c_\lambda = c(d_i, \kappa_i, \Omega_i)(1 + \lambda)/\sqrt{\lambda}$. In the estimates above we have used that $v = T_\lambda g \in W_2^1(\Omega) \cap W_2^2(\Omega \setminus \Sigma)$, and that $(\|\cdot\|_{L_2(\Omega_i)} + \|\Delta \cdot\|_{L_2(\Omega_i)})$ defines an equivalent norm on $W_2^2(\Omega_i)$ and, lastly, we employed (5.3). This completes the proof of Proposition 5.1. \square

We assume now that $\lambda > 0$ is a fixed number. For given $\rho \in W_2^{1/2}(\Sigma)$, let v be the solution of the transmission problem

$$(5.4) \quad \begin{cases} \lambda \kappa v - d \Delta v = 0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ [v] = 0 & \text{on } \Sigma, \\ -[d\partial_\nu v] = -\lambda \rho & \text{on } \Sigma. \end{cases}$$

Then $v = T_\lambda(-\lambda\rho) = -\lambda T_\lambda \rho$ with T_λ being the solution operator introduced above. By inserting this representation of v into the equation $v = \sigma A_\Sigma \rho + \lambda \delta \rho$, we obtain the problem

$$(5.5) \quad \lambda \delta \rho + \lambda T_\lambda \rho + \sigma A_\Sigma \rho = 0,$$

which is equivalent to the eigenvalue problem.

Setting $B_\lambda(s) := \lambda\delta I + \lambda l T_\lambda + s\sigma A_\Sigma$ for $s > 0$ and employing Proposition 5.1(c), we obtain the estimate

$$\begin{aligned} (B_\lambda(s)\rho|\rho)_\Sigma &\geq \lambda(\delta + \gamma l)\|\rho\|_\Sigma^2 + s\sigma(A_\Sigma\rho|\rho)_\Sigma \\ &= \lambda(\delta + \gamma l)\|\rho_0\|_\Sigma^2 + s\sigma(A_\Sigma\rho_0|\rho_0)_\Sigma + \{\lambda(\delta + \gamma l) - s\sigma/R^2\}|\Sigma|\bar{\rho}^2, \end{aligned}$$

where $\gamma = \beta\sqrt{\lambda}/(1 + \lambda)$ and $\rho = \rho_0 + \bar{\rho}$ with $(\rho_0|1)_\Sigma = 0$. Since $(A_\Sigma\rho_0|\rho_0)_\Sigma \geq 0$ we see that all of the terms in the previous line are nonnegative, provided $\lambda(\delta + \gamma l) \geq s\sigma/R^2$, i.e., for small s . Hence, for small $s > 0$, the operator $B_\lambda(s)$ is positive definite, which means that λ cannot be an eigenvalue of (1.9), where in the third line of (1.9) σ is replaced by $s\sigma$. On the other hand, choosing $\rho = 1$ we have

$$(B_\lambda(s)1|1)_\Sigma = \lambda\delta|\Sigma| + \lambda l(T_\lambda 1|1)_\Sigma - s\sigma|\Sigma|/R^2 < 0$$

if s becomes large. Now we set

$$s_* := s_*(\lambda) := \sup\{s > 0 : B_\lambda(s) \text{ is positive definite}\}.$$

Then $B_\lambda(s_*)$ is still semi-definite, but not definite, and hence, by compactness of the resolvent, has a nontrivial kernel. Therefore, for a given $\lambda > 0$ there is an $s_* = s_*(\lambda)$ such that λ is an eigenvalue of (1.9), where σA_Σ is replaced by $s_*\sigma A_\Sigma$ in the third line.

Next, we show that positive eigenvalues are simple. Rewrite (5.5) as

$$\lambda\delta\rho_0 + \lambda l T_\lambda\rho_0 + \sigma A_\Sigma\rho_0 = -\{\lambda\delta + \lambda l T_\lambda 1 - \sigma/R^2\}\bar{\rho}.$$

Since B_λ is positive definite on $L_{2,0}(\Sigma)$, this equation has precisely one solution for given $\bar{\rho}$, which shows that the eigenspace $N(\lambda - L_\delta)$ is at most one-dimensional for any given $\lambda > 0$.

To show that nontrivial eigenvalues are semi-simple, suppose that

$$(\lambda - L_\delta)(v, \rho) = (v_1, \rho_1), \quad (\lambda - L_\delta)(v_1, \rho_1) = 0.$$

Then by Green's formula

$$\begin{aligned} \|\sqrt{\kappa}v_1\|_\Omega^2 &= (\lambda\kappa v - d\Delta v|v_1)_\Omega \\ &= (v|\lambda\kappa v_1 - d\Delta v_1)_\Omega + ([d\partial_\nu v]|v_1)_\Sigma - (v|[d\partial_\nu v_1])_\Sigma \\ &= (\delta/l)([d\partial_\nu v|[d\partial_\nu v_1])_\Sigma + l\sigma(\lambda\rho - \rho_1|A_\Sigma\rho_1)_\Sigma \\ &\quad - (\delta/l)([d\partial_\nu v|[d\partial_\nu v_1])_\Sigma - \lambda l\sigma(A_\Sigma\rho|\rho_1)_\Sigma \\ &= -l\sigma(\rho_1|A_\Sigma\rho_1), \end{aligned}$$

which yields

$$\|\sqrt{\kappa}v_1\|_\Omega^2 + l\sigma(\rho_1|A_\Sigma\rho_1).$$

It follows now from (4.1) that v_1 is constant on $\Omega \setminus \Sigma$. Since $\lambda \neq 0$, we then obtain from (1.9) that $(v_1, \rho_1) = (0, 0)$. Thus any nontrivial eigenvalue is semi-simple, and, in particular, positive eigenvalues are algebraically simple.

We want to show that for $\zeta > 1$ there is precisely one positive eigenvalue $\lambda_* > 0$. For this purpose we fix the parameters d, l, σ, δ as well as R , but replace κ by $s\kappa$

in the first line of (1.9). Fixing $\mu = \lambda s$ and scaling $\rho \mapsto \rho/s$, we obtain the scaling $\sigma \mapsto s\sigma$. The argument given previously then shows that there is $s_* > 0$ such that μ is a simple eigenvalue of the scaled problem; hence, $\lambda_* = \mu/s_* > 0$ is a simple positive eigenvalue for (1.9) with κ replaced by $s_*\kappa$ in the first line. Since $\lambda_* = \lambda_*(s_*)$ is simple, the eigenvalue problem

$$(5.6) \quad \begin{cases} \lambda_* s \kappa v - d\Delta v = 0 & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v = 0 & \text{on } \partial\Omega, \\ v = \sigma \mathcal{A}_\Sigma \rho + \lambda_* \delta \rho & \text{on } \Sigma, \\ \lambda_* l \rho - [d\partial_\nu v] = 0 & \text{on } \Sigma \end{cases}$$

has a smooth (analytic) family $[s \mapsto (\lambda_*(s), v(s), \rho(s))]$ of solutions, which exists as long as $\lambda_*(s)$ remains a simple eigenvalue. As $\zeta(s) := s\sigma(\kappa|1)_\Omega/l|\Sigma|R^2$ approaches the value $\zeta = 1$ from above, we must have $\lambda_*(s) \rightarrow 0$ from the right. This means that at the value

$$(5.7) \quad s = s_0 := \frac{l|\Sigma|R^2}{\sigma(\kappa|1)_\Omega}$$

the eigenvalue $\lambda_*(s)$ passes through the origin, in accordance with the jump of the algebraic multiplicity by 1 of the eigenvalue 0 for L_δ at $\zeta = 1 = \zeta(s_0)$. This shows that there can be only one positive eigenvalue for (5.6), independently of the values of the parameters, and there is precisely one iff $\zeta > 1$.

If $\zeta = \sigma(\kappa|1)_\Omega/l|\Sigma|R^2 > 1$, then we have that $s_0 < 1$. The argument given above shows that the modified eigenvalue problem (5.6) has for each $s > s_0$ exactly one simple eigenvalue. This is, in particular, true for $s = 1$, thus establishing Theorem 2.1(f).

Now we turn our attention to the case $l < 0$. As before, we conclude that the operator L_δ has countably many eigenvalues. We note that the argument given in section 4 also applies to the case $l < 0$ and $\delta = 0$, showing that all eigenvalues of (1.9) are real in this case.

In the following, we assume that $l < 0$ and $\delta > 0$. In order to show Theorem 1.1(f), we consider the operators $B_\lambda := \lambda\delta I + \lambda l T_\lambda + \sigma \mathcal{A}_\Sigma$ for $\lambda > 0$. By Proposition 5.1 we have

$$(B_\lambda \rho | \rho)_\Sigma \geq (\delta\lambda - |l| M_0 \lambda^{1/2} - \sigma/R^2) \|\rho\|_\Sigma^2 \geq \|\rho\|_\Sigma^2$$

provided that $\lambda \geq \mu_0$, for some $\mu_0 \geq \lambda_0$. Hence, B_λ is positive definite for large $\lambda > 0$. On the other hand we have

$$(B_\lambda 1 | 1)_\Sigma = \lambda\delta|\Sigma| - \lambda|l|(T_\lambda 1 | 1)_\Sigma - \sigma|\Sigma|/R^2 \leq \lambda\delta|\Sigma| - \sigma|\Sigma|/R^2.$$

Thus for λ small we see that B_λ is not positive. Let

$$\lambda_* := \inf\{\lambda > 0 : B_\mu \text{ is positive definite for all } \mu \geq \lambda\}.$$

Then B_{λ_*} is still semi-definite, but not definite, and hence, by compactness of the resolvent, has a nontrivial kernel. This shows that λ_* is an eigenvalue of (1.9), proving Theorem 1.1(f).

Remark 5.2. Suppose $l < 0$ and $\delta > 0$.

(a) While it is still true that all nontrivial eigenvalues of (1.9) are semi-simple, we cannot conclude that positive eigenvalues are simple.

(b) We do not know whether all eigenvalues of (1.9) are real if $l < 0$ and $\delta > 0$. We can, however, prove that every sector $[\arg \lambda \leq \theta]$ can only contain finitely many eigenvalues for a fixed $\theta \in (0, \pi)$. This can be shown as follows. Let $\theta \in (\pi/2, \pi)$ be fixed, and suppose that $|\arg \lambda| \leq \theta$. Moreover, let $\alpha \in (0, \pi/2)$ be an arbitrary fixed number. Then one verifies that

$$|\lambda + \mu| \geq \min\{\sin \alpha, \sin(\pi - \theta)\}|\lambda|, \quad \text{whenever } \mu \in \mathbb{R}, \alpha \leq |\arg \lambda| \leq \theta,$$

and this shows that there exists a constant $c > 0$ such that

$$(5.8) \quad |\lambda\delta + \sigma(A_{\Sigma}\rho|_{\rho})_{\Sigma}| \geq c|\lambda|, \quad \alpha \leq |\arg \lambda| \leq \theta.$$

Using that $\sigma(A_{\Sigma}\rho|_{\rho}) \geq -\sigma/R^2\|\rho\|_{\Sigma}^2$, we see that

$$(5.9) \quad |\lambda\delta + \sigma(A_{\Sigma}\rho|_{\rho})_{\Sigma}| \geq (\delta/2)|\lambda|, \quad \|\rho\|_{\Sigma} = 1, \quad \operatorname{Re} \lambda \geq 2\sigma/\delta R^2.$$

Combining (5.8)–(5.9) yields

$$(5.10) \quad |\lambda\delta(\rho|_{\rho})_{\Sigma} + \sigma(A_{\Sigma}\rho|_{\rho})_{\Sigma}| \geq k|\lambda|, \quad \|\rho\|_{\Sigma} = 1, \quad |\arg \lambda| \leq \theta, \quad |\lambda| \geq \eta,$$

where $k = \min\{c, \delta/2\}$ and $\eta = (2\sigma/\delta R^2)(1/\cos \alpha)$. Let λ_0 and M_0 be as in Proposition 5.1. Suppose that $\lambda \in \mathbb{C} \setminus \{0\}$ with $|\arg \lambda| \leq \theta$ is an eigenvalue of (1.9) with eigenfunction (v, ρ) . Then we have

$$(5.11) \quad \lambda\delta\rho + \sigma A_{\Sigma}\rho = \lambda l |T_{\lambda}\rho$$

and we can assume, without loss of generality, that $\|\rho\|_{\Sigma} = 1$. If $|\lambda| \geq \max\{\lambda_0, \eta\}$, then we conclude from (5.10)–(5.11) and Proposition 5.1 that

$$k|\lambda| \leq M_0 l \|\lambda\|^{1/2}$$

and so $|\lambda|$ is bounded by $(M_0 l/k)^2$. Clearly, if $|\lambda| \leq \max\{\lambda_0, \eta\}$, we have a trivial bound. This shows that all possible eigenvalues in the sector $[\arg \lambda \leq \theta]$ are bounded. Since eigenvalues cannot accumulate in a bounded set, we see that (1.9) can only have finitely many eigenvalues in the sector $[\arg \lambda \leq \theta]$.

6. Analysis of the unstable eigenvalue. In this section we analyze the properties of the unstable eigenvalue λ_* of problem (5.6) in case $l > 0$ in more detail. In particular, we study the behavior of $\lambda_*(s)$ and the corresponding eigenfunctions near the critical value 1 of $\zeta(s)$, i.e., for s near s_0 ; see (5.7).

Proof Theorem 1.3. (a) We will first analyze the behavior of $(\lambda_*(s), v(s), \rho(s))$ for s near s_0 . In order to do so we use the following ansatz:

$$(6.1) \quad \begin{aligned} \lambda_*(s) &= (s - s_0)\lambda_1(s), \\ v(s) &= -1 + (s - s_0)\lambda_1(s)v_1(s), \\ \rho(s) &= \rho_0 + (s - s_0)\eta + (s - s_0)\lambda_1(s)(\rho_1(s) + \vec{\beta}(s) \cdot \vec{y}), \\ (v_1(s)|\kappa)_{\Omega} &= 0, \quad (\rho_1(s)|1)_{\Sigma} = (\rho_1(s)|Y_j) = 0, \quad 1 \leq j \leq n, \end{aligned}$$

with

$$(6.2) \quad \rho_0 := R^2/\sigma + \vec{\alpha} \cdot \vec{y}, \quad \eta := (\kappa|1)_{\Omega}/l\Sigma,$$

where $\vec{\alpha}, \vec{\beta}(s) \in \mathbb{R}^n$ and $\vec{y} = (Y_1, \dots, Y_n)$. Setting $r = s - s_0$ and inserting this ansatz into the eigenvalue problem (5.6), we obtain the following system of equations:

$$(6.3) \quad \begin{cases} -d\Delta v_1 = s_0\kappa + r\kappa(1 - s\lambda_1 v_1) & \text{in } \Omega \setminus \Sigma, \\ \partial_\nu v_1 = 0 & \text{on } \partial\Omega, \\ [v_1] = 0 & \text{on } \Sigma, \\ -[d\partial_\nu v_1] = -l\rho_0 - r\eta - r\lambda_1 l(\rho_1 + \vec{\beta} \cdot \vec{y}) & \text{on } \Sigma, \\ \sigma A_\Sigma \rho_1 = \frac{\sigma\eta}{R^2 \lambda_1} - \delta\rho_0 + v_1 - r\delta\eta - r\lambda_1 \delta(\rho_1 + \vec{\beta} \cdot \vec{y}) & \text{on } \Sigma. \end{cases}$$

We first observe that due to (5.7), (6.1)₄, (6.2), and the fact that $(Y_i|1)_\Sigma = 0$, the compatibility condition

$$(6.4) \quad (1|s_0\kappa + r\kappa(1 - s\lambda_1 v_1))_\Omega - l(1|\rho_0 + r\eta + r\lambda_1(\rho_1 + \vec{\beta} \cdot \vec{y}))_\Sigma = 0$$

holds. It is our intention to apply the implicit function theorem to find a smooth (analytic) curve of solutions

$$[s \mapsto (\lambda_1(s), v_1(s), \rho_1(s), \vec{\beta}(s))]$$

of (6.3) for s near s_0 . The idea is to use the $(n + 1)$ orthogonality conditions

$$(A_\Sigma \rho_1|1) = (A_\Sigma \rho_1|Y_j) = 0, \quad 1 \leq j \leq n,$$

to determine the $(n + 1)$ scalar functions λ_1 and β_j . In order to do so, we will first derive an expression for $\vec{\alpha}$ and $\vec{\beta}$. Taking the inner product of (6.3)₅ with Y_j yields

$$(6.5) \quad 0 = -\delta \vec{\alpha} + (v_1|\vec{y})_\Sigma - r\delta\lambda_1 \vec{\beta},$$

where $(v_1|\vec{y})_\Sigma$ denotes the vector in \mathbb{R}^n with components $(v_1|Y_j)_\Sigma$, $1 \leq j \leq n$. Due to Proposition 3.2 we have

$$\begin{aligned} (v_1|Y_j)_\Sigma &= (v_1| - [d\partial_\nu T_0 Y_j])_\Sigma = \int_\Omega d \operatorname{div}(v_1 T_0 Y_j) dx \\ &= (d\nabla v_1|\nabla T_0 Y_j)_\Omega = (-d\Delta v_1|T_0 Y_j)_\Omega + (-[d\partial_\nu v_1]|T_0 Y_j)_\Sigma \\ &= (s_0\kappa + r\kappa(1 - s\lambda_1 v_1)|T_0 Y_j)_\Omega - l(\rho_0 + r\eta + r\lambda_1(\rho_1 + \vec{\beta} \cdot \vec{y})|T_0 Y_j)_\Sigma \\ &= -rs\lambda_1(\kappa v_1|T_0 Y_j)_\Omega - l(\rho_0 + r\eta + r\lambda_1(\rho_1 + \vec{\beta} \cdot \vec{y})|T_0 Y_j)_\Sigma \end{aligned}$$

for $j = 1, \dots, n$. Setting first $r = 0$ yields an equation for $\vec{\alpha}$, namely,

$$0 = -\delta \vec{\alpha} - (lR^2/\sigma)(1|T_0 \vec{y})_\Sigma - l(T_0 \vec{y}|\vec{y})_\Sigma \vec{\alpha},$$

i.e.,

$$(6.6) \quad \vec{\alpha} = -(\delta I + l(T_0 \vec{y}|\vec{y})_\Sigma)^{-1} (lR^2/\sigma)(1|T_0 \vec{y})_\Sigma,$$

where $(T_0 \vec{y}|\vec{y})_\Sigma$ denotes the symmetric matrix with entries $[(T_0 Y_i|Y_j)_\Sigma]_{1 \leq i, j \leq n}$. Here we remind the reader that

$$(6.7) \quad \langle (T_0 \vec{y}|\vec{y})_\Sigma \xi | \xi \rangle = \sum_{i, j=1}^n \xi_i \xi_j (T_0 Y_i|Y_j)_\Sigma = (T_0(\xi \cdot \vec{y}) | (\xi \cdot \vec{y}))_\Sigma \geq c \|\xi\|^2$$

for $\xi \in \mathbb{R}^n$, where $\langle \cdot | \cdot \rangle$ denotes the Euclidean inner product on \mathbb{R}^n . This shows that the matrix $(T_0 \vec{y} | \vec{y})_\Sigma$ is positive definite, and hence $\delta I + l(T_0 \vec{y} | \vec{y})_\Sigma$ is invertible for any $\delta \geq 0$. Next, we obtain for $\vec{\beta}$ that

$$(6.8) \quad \vec{\beta} = -(\delta I + l(T_0 \vec{y} | \vec{y})_\Sigma)^{-1} \{s(\kappa v_1 | T_0 \vec{y})_\Omega + l(\eta/\lambda_1 + \rho_1 | T_0 \vec{y})_\Sigma\}.$$

Thus we have a function $\vec{\beta} = \vec{\beta}(\lambda_1, v_1, \rho_1, s)$. Finally, we obtain an equation for λ_1 by taking the inner product of (6.3)₅ with 1:

$$(6.9) \quad 0 = \sigma\eta|\Sigma|/R^2\lambda_1 - \delta R^2|\Sigma|/\sigma + (v_1|1)_\Sigma - r\delta|\Sigma|\eta.$$

Employing the relation $(\kappa|v_1)_\Omega = 0$ and (6.3)_{1,4} as well as (6.5) yields

$$\begin{aligned} -rs\lambda_1(\kappa v_1|v_1)_\Omega &= (s_0\kappa + r\kappa(1 - s\lambda_1 v_1)|v_1)_\Omega \\ &= (-d\Delta v_1|v_1)_\Omega = \|\sqrt{d}\nabla v_1\|_\Omega^2 + ([d\partial_\nu v_1]|v_1)_\Sigma \\ &= \|\sqrt{d}\nabla v_1\|_\Omega^2 + (l\rho_0 + rl\eta + r\lambda_1 l(\rho_1 + \vec{\beta} \cdot \vec{y})|v_1)_\Sigma \\ &= \|\sqrt{d}\nabla v_1\|_\Omega^2 + \{lR^2/\sigma + rl\eta\}(v_1|1)_\Sigma + l\delta|\vec{\alpha} + r\lambda_1\vec{\beta}|^2 + r\lambda_1 l(\rho_1|v_1)_\Sigma. \end{aligned}$$

This leads to

$$(6.10) \quad 0 = \|\sqrt{d}\nabla v_1\|_\Omega^2 - \{lR^2/\sigma + rl\eta\} \{(\sigma|\Sigma|\eta/R^2\lambda_1) - \delta R^2|\Sigma|/\sigma - r\delta|\Sigma|\eta\} + l\delta|\vec{\alpha} + r\lambda_1\vec{\beta}|^2 + r\lambda_1 \{l(\rho_1|v_1)_\Sigma + s(\kappa v_1|v_1)_\Omega\},$$

where we used (6.9).

Suppose now that v_1 solves the first four equations of (6.3). Then one easily verifies that (6.5) is equivalent to (6.6) and (6.8). Moreover, assuming once again that v_1 satisfies the first four equations of (6.3), and that $\vec{\alpha}$ and $\vec{\beta}$ satisfy (6.6) and (6.8), one verifies that

$$(6.11) \quad (6.9) \iff (6.10).$$

For $r = 0$, that is, for $s = s_0$, we obtain from (6.10)

$$(6.12) \quad \lambda_1(s_0) = l\Sigma\eta/\{\|\sqrt{d}\nabla v_1(s_0)\|_\Omega^2 + l\delta(|\vec{\alpha}|^2 + R^4|\Sigma|/\sigma^2)\},$$

where $v_1(s_0)$ is the unique solution of problem (3.4) with $(f, g) = (s_0\kappa, -l\rho_0)$ and $(\kappa|v_1(s_0))_\Omega = 0$; see Proposition 3.2. This shows that $\lambda_1(s_0)$ is uniquely defined and strictly positive. Moreover, we also know from (6.11) that

$$0 = (\sigma\eta|\Sigma|/R^2\lambda_1(s_0)) - \delta R^2|\Sigma|/\sigma + (v_1(s_0)|1)_\Sigma - r\delta|\Sigma|\eta.$$

We obtain $\rho_1(s_0)$ by solving

$$\sigma A_\Sigma \rho_1 = \sigma\eta/R^2\lambda_1(s_0) - \delta\rho_0 + v_1(s_0)$$

for ρ_1 , which is possible since we chose $\lambda_1(s_0)$ and $\vec{\alpha}$ in such a way that the necessary orthogonality conditions of Proposition 3.1(d) hold. Equation (6.8) shows that the mapping

$$(6.13) \quad I \times \left\{v \in W_2^2(\Omega \setminus \Sigma) : [v] = 0 \text{ on } \Sigma\right\} \times W_2^2(\Sigma) \times \mathbb{R} \rightarrow \mathbb{R}^n, \\ [(\lambda_1, v_1, \rho_1, s) \mapsto \vec{\beta}(\lambda_1, v_1, \rho_1, s)]$$

is analytic where $I \subset \mathbb{R}$ is an open interval that contains $\lambda_1(s_0)$ but does not contain 0.

We are now in a position to apply the implicit function theorem at the point $(v_1(s_0), \rho_1(s_0), s_0)$ to solve the first four equations in (6.3) and (6.9) for (λ_1, v_1) in terms of (ρ_1, s) . We choose the functional analytic setting

$$\begin{aligned} X_1 &:= \{v \in W_2^2(\Omega \setminus \Sigma) : \partial_\nu v = 0 \text{ on } \partial\Omega, [v] = 0 \text{ on } \Sigma, (\kappa|v)_\Omega = 0\}, \\ X_2 &:= \{\rho \in W_2^2(\Sigma) : (\rho|1)_\Sigma = (\rho|Y_j)_\Sigma = 0, 1 \leq j \leq n\}, \\ X &:= \mathbb{R} \times X_1 \times X_2 \times \mathbb{R}, \\ Y &:= \mathbb{R} \times \{(f, g) \in L_2(\Omega) \times W_2^{1/2}(\Sigma) : (f|1)_\Omega + (g|1)_\Sigma = 0\} \end{aligned}$$

and we define $F : V \subset X \rightarrow Y$ by means of

$$F(\lambda_1, v_1, \rho_1, s) := \begin{pmatrix} \sigma\eta|\Sigma|/R^2\lambda_1 - \delta R^2|\Sigma|/\sigma + (v_1|1)_\Sigma - r\delta|\Sigma|\eta \\ -d\Delta v_1 - s_0\kappa - r\kappa(1 - s\lambda_1 v_1) \\ -[d\partial_\nu v_1] + l(\rho_0 + r\eta + r\lambda_1(\rho_1 + \vec{\beta} \cdot \vec{y})) \end{pmatrix},$$

where $V := I \times X_1 \times X_2 \times \mathbb{R}$.

Equation (6.4) implies that F maps V into Y , and (6.13) and the definition of F show that

$$[(\lambda_1, v_1, \rho_1, s) \mapsto F(\lambda_1, v_1, \rho_1, s)] \in C^\omega(V, Y).$$

Clearly, the first four equations of (6.3) together with (6.9) are equivalent to $F(\lambda_1, v_1, \rho_1, s) = (0, 0, 0)$. Since we already know that

$$F(\lambda_1(s_0), v_1(s_0), \rho_1(s_0), s_0) = (0, 0, 0),$$

we are left with verifying that the derivative of F at the point $(v_1(s_0), \rho_1(s_0), s_0)$ w.r.t. (ρ_1, v_1) is an isomorphism, i.e.,

$$\mathbb{D}_1 F(\lambda_1(s_0), v_1(s_0), \rho_1(s_0), s_0) \in \text{Isom}(\mathbb{R} \times X_1, Y).$$

It follows from Proposition 3.2(a) that the problem

$$\mathbb{D}_1 F(\lambda_1(s_0), v_1(s_0), \rho_1(s_0), s_0)(\lambda, w) = (\mu, f, g)$$

has for each $(\mu, f, g) \in Y$ a unique solution $(\lambda, w) \in \mathbb{R} \times X_1$, namely,

$$w = R_0(f, g), \quad \lambda = \frac{R^2 \lambda_1^2(s_0)}{\sigma\eta|\Sigma|} ((R_0(f, g)|1)_\Sigma - \mu),$$

where $w = R_0(f, g)$ is the unique solution of (3.4) with $(\kappa|w)_\Omega = 0$.

The implicit function theorem then yields a neighborhood U of $(\rho_1(s_0), s_0)$ in $X_2 \times \mathbb{R}$ such that

$$(6.14) \quad \begin{aligned} &[(\rho_1, s) \mapsto (\lambda_1(\rho_1, s), v_1(\rho_1, s))] \in C^\omega(U, \mathbb{R} \times X_1) \\ &F(\lambda_1(\rho_1, s), v_1(\rho_1, s), \rho_1, s) = (0, 0, 0), \quad (\rho_1, s) \in U. \end{aligned}$$

Combining all of the above results, we conclude that

$$(6.15) \quad [(\rho_1, s) \mapsto (\lambda_1(\rho_1, s), v_1(\rho_1, s), \vec{\beta}(\rho_1, s))] \in C^\omega(U, \mathbb{R} \times X_1 \times \mathbb{R}^n)$$

and that the functions $(\lambda_1(\rho_1, s), v_1(\rho_1, s), \vec{\alpha}, \vec{\beta}(\rho_1, s))$ satisfy the first four equations of (6.3) as well as (6.5) and (6.9). We now insert these functions into the equation for ρ_1 which gives an equation of the form

$$G(\rho_1, s) := \sigma A_{\Sigma} \rho_1 - \sigma \eta / R^2 \lambda_1(s_0) + \delta \rho_0 - v_1(s_0) - (s - s_0)R(\rho_1, s) = 0,$$

where $[(\rho_1, s) \mapsto R(\rho_1, s)] \in C^\omega(U, L_2(\Sigma))$. By (6.5) and (6.9) we know that

$$G : X_2 \times \mathbb{R} \rightarrow Z := \{g \in L_2(\Sigma) : (g|1)_{\Sigma} = (g|Y_j) = 0, 1 \leq j \leq n\}.$$

Moreover, we also know that $G(\rho_1(s_0), s_0) = 0$. The derivative of G with respect to ρ_1 at $(\rho_1(s_0), s_0)$ is σA_{Σ} , and Proposition 3.1(d) and the implicit function theorem then yield an analytic curve

$$[s \mapsto \rho_1(s)] \in C^\omega((s_0 - \varepsilon_0, s_0 + \varepsilon_0), X_2)$$

such that $G(\rho_1(s), s) = 0$.

Combining all of the results, we obtain an analytic curve of solutions

$$[s \mapsto (\lambda_*(s), v(s), \rho(s))]$$

of (5.6) for $s \in (s_0 - \varepsilon_0, s_0 + \varepsilon_0)$. If $s > s_0$, then the statement in (a) follows from the considerations in section 5.

(b) The proof of part (a) shows that the eigenvalue curve $[s \mapsto \lambda_*(s)]$ is analytic near the critical value $s = s_0$ and crosses the imaginary axis at $s = s_0$ with positive speed $\lambda_1(s_0)$; see (6.12).

(c) We show that $\lambda_*(s)$ is strictly increasing. To see this, we differentiate (5.6) w.r.t. s and form the inner product of the resulting equation in Ω with $v = v(s)$. This yields with Green's formula

$$\begin{aligned} -(s\lambda_*(s))' \|\sqrt{\kappa}v\|_{\Omega}^2 &= s\lambda_*(\kappa v'|v)_{\Omega} - (d\Delta v'|v)_{\Omega} \\ &= (v'|s\lambda_*\kappa v - d\Delta v)_{\Omega} + ([d\partial_{\nu}v']|v)_{\Sigma} - (v'|[d\partial_{\nu}v])_{\Sigma} \\ &= (\delta/l)([d\partial_{\nu}v']|[d\partial_{\nu}v])_{\Sigma} + (\lambda_*l\rho' + \lambda_*'l\rho|\sigma A_{\Sigma}\rho)_{\Sigma} \\ &\quad - (\delta/l)([d\partial_{\nu}v]|[d\partial_{\nu}v'])_{\Sigma} - (\lambda_*l\rho|\sigma A_{\Sigma}\rho')_{\Sigma} \\ &= \lambda_*'l\sigma(A_{\Sigma}\rho|\rho)_{\Sigma}; \end{aligned}$$

hence

$$\lambda_* \|\sqrt{\kappa}v\|_{\Omega}^2 + \lambda_*' \{ \|\sqrt{s\kappa}v\|_{\Omega}^2 + l\sigma(A_{\Sigma}\rho|\rho)_{\Sigma} \} = 0.$$

Employing (4.1) once more (with κ replaced by $s\kappa$), we obtain

$$\lambda_*'(s) = \lambda_*^2(s) \|\sqrt{\kappa}v\|_{\Omega}^2 / \{ \|\sqrt{d}\nabla v\|_{\Omega}^2 + (\delta/l)\|[d\partial_{\nu}v]\|_{\Sigma}^2 \},$$

which yields $\lambda_*'(s) > 0$ for $s \neq s_0$. If $s = s_0$, then we have already established in (b) that $\lambda_*'(s_0) > 0$, and this shows the assertion of Theorem 1.3(c).

(d) If the stability condition (5.1) is violated, then we can conclude from the identity (4.4), where κ is now replaced with $s\kappa$ and λ is replaced with λ_* , that

$$\left(\lambda_*(s)l\delta|\Sigma| + \frac{l^2|\Sigma|^2}{s(\kappa|1)_{\Omega}} - \frac{l\sigma|\Sigma|}{R^2} \right) \leq 0,$$

which shows that $\lambda_*(s) \leq \frac{\sigma}{\delta R^2} (1 - 1/\zeta(s))$ provided that $\delta > 0$, i.e., if kinetic undercooling is present.

(e) To show that $\lambda_*(s) \rightarrow \infty$ as $s \rightarrow \infty$ in case $\delta = 0$, we employ the estimate in Proposition 5.1(b). We first observe that due to the fact that $\lambda_*(s)$ is increasing in s , there exists a number $s_1 > s_0$ such that $s\lambda_*(s) \geq \lambda_0$ for $s \geq s_1$, where λ_0 is the number occurring in Proposition 5.1(b). It then follows from the relation $\lambda_* T_{s\lambda_*} \rho + \sigma A_\Sigma \rho = 0$ that

$$(6.16) \quad \begin{aligned} \sigma^2 \|A_\Sigma \rho\|_\Sigma^2 &= \sigma^2 \|A_\Sigma \rho_0\|_\Sigma^2 + (\sigma^2 |\Sigma|/R^2) \bar{\rho}^2 = \lambda_*^2 \|T_{s\lambda_*} \rho\|_\Sigma^2 \\ &\leq M_1 \|\rho\|_\Sigma^2 \lambda_*(s)/s, \quad s \geq s_1, \end{aligned}$$

where we write $\rho = \rho_0 + \bar{\rho}$ with $(\rho_0|1)_\Sigma = 0$. Multiplying the eigenvalue problem (5.6) with $(s\lambda_* v - (d/\kappa)\Delta v)$ and using the divergence theorem and (6.16), we get

$$\begin{aligned} (s\lambda_*)^2 \|\sqrt{\kappa}v\|_\Omega^2 + 2s\lambda_* \|\sqrt{d}\nabla v\|_\Omega^2 + \|(d/\sqrt{\kappa})\Delta v\|_\Omega^2 + 2s\lambda_*^2 l\sigma(\rho_0|A_\Sigma \rho_0)_\Sigma \\ = (2s\lambda_*^2 l\sigma|\Sigma|/R^2) \bar{\rho}^2 \leq M_2 \|\rho\|_\Sigma^2 \lambda_*^3, \quad s \geq s_1. \end{aligned}$$

The relation $\lambda_* l\rho = [d\partial_\nu v]$ and the inequality above yield

$$\begin{aligned} \lambda_*^2 \|\rho\|_\Sigma^2 &= (1/l)^2 \|[d\partial_\nu v]\|_\Sigma^2 \leq C \|v\|_{W^{3/2+\varepsilon}(\Omega \setminus \Sigma)}^2 \leq C \|v\|_\Omega^{(1/2-\varepsilon)} \|v\|_{W^{3/2+\varepsilon}(\Omega \setminus \Sigma)}^{3/2+\varepsilon} \\ &\leq C \left\{ \|v\|_\Omega^2 + \|v\|_\Omega^{(1/2-\varepsilon)} \|\Delta v\|_\Omega^{3/2+\varepsilon} \right\} \leq C \|\rho\|_\Sigma^2 \lambda_*^3 \left\{ \frac{1}{(s\lambda_*)^2} + \frac{1}{(s\lambda_*)^{1/2-\varepsilon}} \right\} \\ &\leq C \|\rho\|_\Sigma^2 \lambda_*^3 / (s\lambda_*)^{1/2-\varepsilon} = C \|\rho\|_\Sigma^2 \lambda_*^{5/2+\varepsilon} / s^{1/2-\varepsilon} \end{aligned}$$

for $s \geq s_1$, where C is a generic constant that may change from line to line. Dividing by λ_*^2 and by $\|\rho\|_\Sigma^2$ implies

$$\lambda_*(s) \geq cs^{(1-2\varepsilon)/(1+2\varepsilon)}, \quad s \geq s_1.$$

Thus we can conclude that $\liminf_{s \rightarrow \infty} \lambda_*(s)/s^\theta = \infty$ for each $\theta < 1$. \square

Acknowledgments. The authors would like to thank Rico Zacher for helpful discussions. We thank the anonymous referees for carefully reading the manuscript.

REFERENCES

[1] B. V. BAZALII, *Stefan problem for the Laplace equation with regard to the curvature of the free boundary*, Ukrainian Math. J., 49 (1997), pp. 1465–1484.
 [2] G. CAGINALP, *An analysis of a phase field model of a free boundary*, Arch. Rational Mech. Anal., 92 (1986), pp. 205–245.
 [3] B. CHALMERS, *Principles of Solidification*, Krieger, Huntington, NY, 1977.
 [4] X. CHEN, *The Hele-Shaw problem and area-preserving curve-shortening motion*, Arch. Rational Mech. Anal., 123 (1993), pp. 117–151.
 [5] X. CHEN AND F. REITICH, *Local existence and uniqueness of solutions of the Stefan problem with surface tension and kinetic undercooling*, J. Math. Anal. Appl., 164 (1992), pp. 350–362.
 [6] X. CHEN, J. HONG, AND F. YI, *Existence, uniqueness, and regularity of classical solutions of the Mullins-Sekerka problem*, Comm. Partial Differential Equations, 21 (1996), pp. 1705–1727.
 [7] R. DENK, J. PRÜSS, AND R. ZACHER, *Maximal L_p -regularity of Parabolic Problems with Boundary Conditions of Relaxation Type*, submitted.
 [8] J. ESCHER AND G. SIMONETT, *On Hele-Shaw models with surface tension*, Math. Res. Lett., 3 (1996), pp. 467–474.
 [9] J. ESCHER AND G. SIMONETT, *Classical solutions for the quasi-stationary Stefan problem with surface tension*, in Differential Equations, Asymptotic Analysis, and Mathematical Physics (Potsdam, 1996), Math. Res. 100, Akademie Verlag, Berlin, 1997, pp. 98–104.

- [10] J. ESCHER AND G. SIMONETT, *Classical solutions for Hele-Shaw models with surface tension*, Adv. Differential Equations, 2 (1997), pp. 619–642.
- [11] J. ESCHER AND G. SIMONETT, *A center manifold analysis for the Mullins-Sekerka model*, J. Differential Equations, 143 (1998), pp. 267–292.
- [12] J. ESCHER, J. PRÜSS, AND G. SIMONETT, *On the Stefan problem with surface tension*, in Elliptic and Parabolic Problems (Rolduc/Gaeta, 2001), World Scientific, River Edge, NJ, 2002, pp. 377–388.
- [13] J. ESCHER, J. PRÜSS, AND G. SIMONETT, *Analytic solutions for a Stefan problem with Gibbs-Thomson correction*, J. Reine Angew. Math., 563 (2003), pp. 1–52.
- [14] A. FRIEDMAN AND F. REITICH, *The Stefan problem with small surface tension*, Trans. Amer. Math. Soc., 328 (1991), pp. 465–515.
- [15] A. FRIEDMAN AND F. REITICH, *Nonlinear stability of a quasi-static Stefan problem with surface tension: A continuation approach*, Ann. Sc. Norm. Super. Pisa Cl. Sci. (4), 30 (2001), pp. 341–403.
- [16] M. E. GURTIN, *On the two-phase problem with interfacial energy and entropy*, Arch. Rational Mech. Anal., 96 (1986), pp. 199–241.
- [17] M. E. GURTIN, *Multiphase thermomechanics with interfacial structure*, Arch. Rational Mech. Anal., 104 (1988), pp. 195–221.
- [18] E. I. HANZAWA, *Classical solutions of the Stefan problem*, Tôhoku Math. J., 33 (1981), pp. 297–335.
- [19] P. HARTMAN, *Crystal Growth: An Introduction*, North-Holland, Amsterdam, 1973.
- [20] C. KNEISEL, *Über das Stefan-Problem mit Oberflächenspannung und thermischer Unterkühlung*, Ph.D. thesis, Leibniz Universität, Hannover, Germany, 2007.
- [21] J. S. LANGER, *Instabilities and pattern formation in crystal growth*, Rev. Mod. Phys., 52 (1980), pp. 1–28.
- [22] S. LUCKHAUS, *Solutions for the two-dimensional Stefan problem with the Gibbs-Thomson law for melting temperature*, European J. Appl. Math., 1 (1990), pp. 101–111.
- [23] A. M. MEIRMANOV, *The Stefan problem with surface tension in the three dimensional case with spherical symmetry: Non-existence of the classical solution*, European J. Appl. Math., 5 (1994), pp. 1–20.
- [24] W. W. MULLINS, *Thermodynamic equilibrium of a crystal sphere in a fluid*, J. Chem. Phys., 81 (1984), pp. 1436–1442.
- [25] W. W. MULLINS AND R. F. SEKERKA, *Stability of a planar interface during solidification of a dilute binary alloy*, J. Appl. Phys., 35 (1964), pp. 444–451.
- [26] J. PRÜSS AND G. SIMONETT, *Smooth Solutions for Two-phase Stefan Problems with Surface Tension*, in preparation.
- [27] J. PRÜSS, G. SIMONETT, AND R. ZACHER, *Qualitative Behavior of Stefan Problems with Surface Tension*, in preparation.
- [28] E. RADKEVITCH, *The Gibbs-Thompson correction and conditions for the existence of a classical solution of the modified Stefan problem*, Dokl. Akad. Nauk SSSR, 316 (1991), pp. 1311–1315; translation in Soviet Math. Dokl., 43 (1991), pp. 274–278.
- [29] E. RADKEVITCH, *Conditions for the existence of a classical solution of a modified Stefan problem (the Gibbs-Thomson law)*, Mat. Sb., 183 (1992), pp. 77–101; translation in Russian Acad. Sci. Sb. Math., 75 (1993), pp. 221–246.
- [30] H. TRIEBEL, *Höhere Analysis*, Hochschulbücher für Mathematik, Band 76, VEB Deutscher Verlag der Wissenschaften, Berlin, 1972 (in German).
- [31] A. VISINTIN, *Models for supercooling and superheating effects*, in Free Boundary Problems: Applications and Theory, Vol. III, Pitman Res. Notes Math. 120, Pitman, Boston, 1985, pp. 200–207.
- [32] A. VISINTIN, *Models of Phase Transitions*, Progr. Nonlinear Differential Equations Appl. 28, Birkhäuser Boston, Boston, 1996.
- [33] W. YU, *A quasisteady Stefan problem with curvature correction and kinetic undercooling*, J. Partial Differential Equations, 9 (1996), pp. 55–70.

KHOKHLOV–ZABOLOTSKAYA–KUZNETSOV-TYPE EQUATION: NONLINEAR ACOUSTICS IN HETEROGENEOUS MEDIA*

I. KOSTIN[†] AND G. PANASENKO[†]

Abstract. The Khokhlov–Zabolotskaya–Kuznetsov model is a PDE describing the wave profile of an acoustic beam in a nonlinear medium. The paper treats the existence and uniqueness questions for this equation in the case of nonconstant coefficients. For the case of rapidly oscillating coefficients the asymptotic behavior of the solution is studied.

Key words. acoustics, nonlinear, homogenization

AMS subject classifications. 35Q35, 35B27

DOI. 10.1137/060674272

1. Introduction. The so-called Khokhlov–Zabolotskaya–Kuznetsov (KZK) equation [8, 13] belongs to the set of nonlinear acoustics models, such as the well-known Riemann wave equation (or the nonlinear transfer equation), the Burgers equation, the Korteweg–de Vries equation, the Khokhlov–Zabolotskaya equation (see [2, 8, 10, 11, 14]), the Zakharov–Kuznetsov equation [5, 6], and the Rudenko–Sukhorukov equation [7, 9, 12]. These models are derived from the linear or nonlinear wave equation for the acoustic pressure, usually under the hypothesis of small variations of this pressure. More precisely the KZK equation has the form

$$(1.1) \quad \alpha u_{z\tau} = (f(u_\tau))_\tau + \beta u_{\tau\tau\tau} + \gamma u_\tau + \Delta_x u,$$

where $u_\tau = u_\tau(z, x, \tau)$ is the acoustic pressure, $(z, x) \in \mathbb{R} \times \mathbb{R}^d$, $d = 1, 2$, are space variables, and τ is the retarded time.

The nonlinear function f in the KZK equation is quadratic, i.e., $f(s) = \theta s^2$, although for the description of space-limited beams subject to the diffraction and self-action effects it can be taken as cubic: $f(s) = \theta s^3$ (see [13]).

In the real physical setting, f may have a more complicated shape. On the other hand, all of these models are derived under the assumption of small oscillations of the pressure, and so one can always consider f as quadratic for $|s| \leq s^*$ and anything different for $|s| > s^*$ with some finite s^* . If $|u_\tau|$ is smaller than s^* , then the two models (with $f(s) = \theta s^2 \forall s$ and $f(s) = \theta s^2$ for $|s| \leq s^*$) coincide. The advantage of this modified shape of f is that (as it will be proved below) we can get the global existence theorem as soon as f has a bounded derivative.

These arguments motivate us to consider the “KZK-type equation,” that is, (1.1) with a nonlinearity f admitting a bounded derivative. Let us emphasize that this shape gives a more convenient physical description than the classical quadratic shape.

Another particularity of the model we consider in the present paper is that the coefficients are rapidly oscillating functions of z . This corresponds to the heterogeneous (stratified in the direction of the axis z) acoustic media. In this case (1.1) takes

*Received by the editors November 7, 2006; accepted for publication (in revised form) December 17, 2007; published electronically July 3, 2008.

<http://www.siam.org/journals/sima/40-2/67427.html>

[†]Laboratory of Mathematics of the University of Saint-Etienne (LaMUSE), Université Jean Monnet, 23 rue P. Michelon, 42023 Saint-Etienne, France (kostin@free.fr, grigory.panasenko@univ-st-etienne.fr).

the form

$$u_{zt} = r(z, z/\epsilon)(\phi(u_\tau))_\tau + b(z, z/\epsilon)u_{\tau\tau\tau} + g(z, z/\epsilon)u_\tau + \Delta_x u,$$

where ϵ is a small parameter representing the ratio between the microscale and the macroscale, while the coefficients r, b, g are oscillating with respect to the second variable. In particular, they may be 1-periodic functions of this variable. This feature complicates the problem, although it allows us to apply the homogenization method (see [1]) to obtain the homogenized model. Its solution is close to the one of the initial problem.

It seems that so far there was not any publication on the existence and uniqueness of the solution for the KZK (or KZK-type) equation, although the authors discovered that independently and simultaneously these questions (as well as the derivation of the KZK equation from the Navier–Stokes model) in the case of constant coefficients were studied by Bardos and Rozanova [3, 4]. They consider the KZK equation in the whole space ($x \in \mathbb{R}^d$) in the case of constant coefficients. In the present study we shall consider the varying (and even rapidly oscillating) coefficients of the KZK-type equation set for $x \in \omega$, where $\omega \subset \mathbb{R}^d$ is a bounded domain. The boundedness of f' will ensure the global existence, while in [3, 4] it is proved for small initial data only. In the case of stratified media we homogenize the KZK-type equation and prove the closeness of the solutions of the homogenized and initial models.

The main results of the present paper were announced by the authors in [15].

2. Problem setting. Let ω be a bounded open domain in \mathbb{R}^d with a smooth boundary $\partial\omega$. For a real function $u = u(z, x, \tau)$ of the variables $z \in [0, Z]$, $\tau \in \mathbb{R}$, and $x \in \omega$, consider the following PDE:

$$(2.1) \quad \alpha u_{z\tau} = (f(u_\tau))_\tau + \beta u_{\tau\tau\tau} + \gamma u_\tau + \Delta u,$$

where the Laplace operator Δ (as well as ∇ below) derives with respect to the variable $x \in \omega$. The positive coefficients α, β , and γ are functions of z and x . The nonlinearity f may depend on z and x as well, but neither depends on τ .

In the mathematical setting, the variable z is considered as the evolutionary variable, while x and τ shall be called space variables. Let us impose the following boundary conditions on u :

$$(2.2) \quad u \text{ is } 2\pi\text{-periodic with respect to } \tau,$$

$$(2.3) \quad \nu \cdot \nabla u = 0 \quad \text{on } \partial\omega,$$

where ν denotes the unit normal vector of ω . It is clear that these conditions are not sufficient to ensure the uniqueness. Indeed, if the coefficients are constant, then any function u depending only on z solves (2.1)–(2.3). This liberty is eliminated by the additional orthogonality condition

$$(2.4) \quad \int_0^{2\pi} u(z, x, \tau) d\tau = 0 \quad \forall x, \forall z.$$

Finally, the evolutionary problem (2.1)–(2.4) requires an initial condition

$$(2.5) \quad u(0, x, \tau) = u_0(x, \tau).$$

To avoid heavy notation, in the estimates below we shall denote by M any large positive constant depending only on $\omega, f, \alpha, \beta, \gamma$, and Z . By “large” we mean that the estimate stays true with M replaced by any larger value.

Here is the list of assumptions under which the existence and uniqueness result for problem (2.1)–(2.5) will be established.

- (ω) $\omega \subset \mathbb{R}^d$ is an open bounded set. Its boundary $\partial\omega$ is twice continuously differentiable.
- (f1) $f(s) = f(s, z, x)$ is continuously differentiable with respect to s , and its partial derivative $f' \equiv f_s(s, z, x)$ is uniformly bounded on $\mathbb{R} \times \mathbb{R}_+ \times \omega$.
- (f2) $f(0, z, x) = 0$ for all $z \in \mathbb{R}_+$ and $x \in \omega$. As is easy to see, this is not a constraint (otherwise, one can always replace $f(s, z, x)$ by $f(s, z, x) - f(0, z, x)$, leaving the equation intact). Together with the previous assumption this one implies that $|f(s, z, x)| \leq M|s|$. Moreover, we assume that f is continuously differentiable with respect to z and its derivative f_z satisfies $|f_z(s, z, x)| \leq M|s|$.
- (f3) Denote by F the primitive of f with respect to the first argument such that $F(0, z, x) = 0$. Clearly, $|F(s, z, x)| \leq M|s|^2$. We shall assume that F is differentiable with respect to z and $|F_z(v, z, x)| \leq M|v|^2$.
- ($\alpha\beta\gamma$) $\alpha = \alpha(z, x)$, $\beta = \beta(z, x)$, $\gamma = \gamma(z, x)$ are bounded functions defined on $\mathbb{R}_+ \times \omega$. The functions α and β are positive and uniformly bounded away from zero. The derivatives α_z , α_x , β_z , and β_x are assumed to be bounded functions.

These conditions are assumed to hold throughout sections 2–5.

The above list has to be completed by an assumption concerning the regularity of u_0 . It will be given along with the definition of a solution.

3. Notation and preliminaries. Let us introduce some functional spaces. First consider the spaces of functions of the variables x and τ . Set $\Omega = \omega \times (0, 2\pi)$. The scalar product and the norm in $L_2(\Omega)$ are denoted by $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$, respectively.

Let D_Ω be the set of infinitely smooth real functions $v = v(x, \tau)$ defined on Ω and satisfying

$$(3.1) \quad \int_0^{2\pi} v(x, \tau) d\tau = 0 \quad \forall x \in \bar{\omega} \quad (\text{orthogonality}),$$

$$(3.2) \quad \frac{\partial^k}{\partial \tau^k} v(x, 0) = \frac{\partial^k}{\partial \tau^k} v(x, 2\pi) \quad \forall x \in \bar{\omega}, \quad \forall k = 0, 1, \dots \quad (\text{periodicity}),$$

$$(3.3) \quad \nu(x) \cdot \nabla v(x, \tau) = 0 \quad \forall x \in \partial\omega, \quad \forall \tau \in [0, 2\pi] \quad (\text{Neumann condition}).$$

Denote by X_0 the closure of D_Ω in the $L_2(\Omega)$ -norm. Given a function $v \in X_0$, we shall denote by \tilde{v} its (unique in X_0) primitive with respect to τ , that is,

$$\tilde{v}_\tau = v \quad \text{and} \quad \int_0^{2\pi} \tilde{v} d\tau = 0.$$

The Poincaré inequality yields

$$\|\tilde{v}\| \leq \|v\| \leq \|v_\tau\| \leq \|v_{\tau\tau}\|$$

for all $v \in D_\Omega$. Note also that the mappings $v \mapsto v_\tau$ and $v \mapsto \tilde{v}$ are both antisymmetric, that is,

$$\langle v, v_\tau \rangle = 0 \quad \text{and} \quad \langle v, \tilde{v} \rangle = 0$$

for all v in X_0 for which these expressions make sense.

The closures of D_Ω with respect to the norms

$$\begin{aligned} \|v\|_1^2 &= \|v_\tau\|^2 + \|\nabla v\|^2, \\ \|v\|_2^2 &= \|v_{\tau\tau}\|^2 + \|\Delta \tilde{v}\|^2 \end{aligned}$$

are denoted by X_1 and X_2 , respectively. Recall that, as well as the Laplace operator, the gradient is assumed to derive only with respect to $x \in \omega$. The continuous embeddings $X_2 \subset X_1 \subset X_0$ take place, the first one following from the estimate

$$2\|\nabla v\|^2 = 2\langle v_\tau, \Delta \tilde{v} \rangle \leq \|v_\tau\|^2 + \|\Delta \tilde{v}\|^2 \leq \|v\|_2^2,$$

while the second is the classical Poincaré inequality. Moreover, by the Rellich theorem, the embedding $X_1 \subset X_0$ is compact, and therefore the embedding $X_2 \subset X_0$ is also compact.

Now let us consider spaces of functions of three variables. For a given positive Z , set $Q = (0, Z) \times \omega \times (0, 2\pi)$. We shall use the triple bars to denote the $L_2(Q)$ -norm:

$$\|u\|^2 = \int_0^Z \|u(z, \cdot, \cdot)\|^2 dz.$$

Let D_Q be the set of infinitely smooth real functions $u = u(z, x, \tau)$ defined on \bar{Q} such that $u(z, \cdot, \cdot) \in D_\Omega$ for all $z \in [0, Z]$. The closure of D_Q in $L_2(Q)$ is denoted U_0 .

We shall seek for weak solutions of problem (2.1)–(2.5) in the space U_1 , which is defined as the closure of D_Q in the norm

$$\|u\|_1^2 = \|u_{\tau\tau}\|^2 + \|\nabla u_\tau\|^2,$$

while U_2 stands for the space of strong solutions and is defined as the closure of D_Q with respect to the norm

$$\|u\|_2^2 = \|u_{z\tau}\|^2 + \|u_{\tau\tau\tau}\|^2 + \|\Delta u\|^2.$$

The continuous embedding $U_2 \subset U_1$ is proved similarly to the case of $X_2 \subset X_1$ and implies the inequality

$$(3.4) \quad \|u_{z\tau}\|^2 + \|u_{\tau\tau}\|^2 + \|\nabla u_\tau\|^2 \leq 2\|u\|_2^2,$$

which we shall need in what follows.

4. Strong solutions: Definition and existence. If $u \in U_2$, then (2.1) writes down for u as an equality of functions in U_0 . Note also that the functions of U_2 belong to $C([0, Z], X_0)$, so that the initial condition $u(0, \cdot, \cdot) = u_0$ can be understood as an identity in X_0 . These two observations justify the following definition.

DEFINITION 4.1. *For given $u_0 \in X_2$ and $Z > 0$, a function $u \in U_2$ is called a strong solution to problem (2.1)–(2.5) on $[0, Z]$ if (2.1) holds almost everywhere in Q and $u(0, \cdot, \cdot) = u_0$ almost everywhere in Ω .*

From now on we shall consider the real functions of $(z, x, \tau) \in Q$ as X_0 -valued functions of the variable z and write $u(z)$ instead of $u(z, x, \tau)$.

PROPOSITION 4.2. *For any $u_0 \in X_2$ and any $Z > 0$, problem (2.1)–(2.5) admits a strong solution u satisfying*

$$\sup_z \left[\|u_{\tau\tau}\|^2 + \|u_z\|^2 \right] + \|u_{\tau\tau\tau}\|^2 + \|u_{z\tau}\|^2 + \|\Delta u\|^2 \leq M \|u_{0\tau\tau}\|^2 + M \|\Delta \tilde{u}_0\|^2.$$

Proof. Denote by ψ_k the eigenfunctions of the following spectral problem:

$$-\psi_{\tau\tau} - \Delta \psi = \lambda \psi, \quad \psi \in X_2.$$

The functions ψ_k form an orthogonal basis in X_0 as well as in X_1 and X_2 . Let P_m stand for the orthogonal projection onto the subspace of X_0 spanned by ψ_1, \dots, ψ_m . The above spectral problem can be solved by the separation of variables, so that P_m commutes not only with $-\partial^2/\partial\tau^2 - \Delta$ but also with $-\partial^2/\partial\tau^2$ and $-\Delta$. The Galerkin approximations for problem (2.1)–(2.5) are defined as continuously differentiable functions $u^m : [0, Z] \rightarrow P_m X_0$ satisfying

$$(4.1) \quad P_m \alpha u_{z\tau}^m = P_m (f(u_\tau^m))_\tau + P_m \beta u_{\tau\tau\tau}^m + P_m \gamma u_\tau^m + \Delta u^m,$$

$$(4.2) \quad u^m(0) = P_m u_0,$$

which is equivalent to a Cauchy problem for a system of m ODEs.

First let us prove that the sequence u^m is bounded in U_2 . Multiply (4.1) by $u_{\tau\tau}^m$, and observe that $\langle \Delta u^m, u_{\tau\tau}^m \rangle = \langle \nabla u_\tau^m, \nabla u_{\tau\tau}^m \rangle = 0$ to obtain

$$-\langle \alpha u_{z\tau\tau}^m, u_{\tau\tau}^m \rangle = \langle f'(u_\tau^m) u_{\tau\tau}^m, u_{\tau\tau}^m \rangle + \langle \beta u_{\tau\tau\tau}^m, u_{\tau\tau}^m \rangle - \langle \gamma u_{\tau\tau}^m, u_{\tau\tau}^m \rangle.$$

After some obvious transformations using assumptions (f1), (f2), and $(\alpha\beta\gamma)$ as well as the Cauchy–Young inequality, the latter relation becomes

$$\frac{d}{dz} \langle \alpha u_{\tau\tau}^m, u_{\tau\tau}^m \rangle + M^{-1} \|u_{\tau\tau}^m\|^2 \leq M \|u_{\tau\tau}^m\|^2,$$

and therefore the Gronwall lemma along with the equivalence of the norms $\langle \alpha u_{\tau\tau}^m, u_{\tau\tau}^m \rangle$ and $\|u_{\tau\tau}^m\|^2$ implies that

$$(4.3) \quad \sup_z \|u_{\tau\tau}^m\|^2 + \|u_{\tau\tau\tau}^m\|^2 \leq M \|u_{\tau\tau}^m(0)\|^2 \leq M \|u_{0\tau\tau}\|^2.$$

In particular, estimate (4.3) guarantees the existence of the functions u^m satisfying (4.1)–(4.2).

In order to obtain a bound on $\|u_{z\tau}^m\|$, take the primitive of (4.1) with respect to τ , and then derive the result in z . This yields the relation

$$(4.4) \quad \begin{aligned} & P_m \alpha_z u_z^m + P_m \alpha u_{zz}^m \\ &= P_m (f(u_\tau^m))_z + P_m \beta_z u_{\tau\tau}^m + P_m \beta u_{z\tau\tau}^m + P_m \gamma_z u^m + P_m \gamma u_z^m + \Delta \tilde{u}_z^m + c, \end{aligned}$$

where the integration constant c may depend on z and x but not on τ . Let us see now what happens when this identity is multiplied by u_z^m . The last two terms on the right will disappear. On the left, keep the second term and half of the first one:

$$\frac{1}{2} \langle \alpha_z u_z^m, u_z^m \rangle + \langle \alpha u_{zz}^m, u_z^m \rangle = \frac{1}{2} \frac{d}{dz} \langle \alpha u_z^m, u_z^m \rangle.$$

The other half of $\langle \alpha_z u_z^m, u_z^m \rangle$ goes to the right and is estimated by $\|u_z^m\|^2$. The third term on the right goes to the left and, after the integration by parts, admits the following lower bound:

$$-\langle \beta u_{z\tau\tau}^m, u_z^m \rangle = \langle \beta u_{z\tau}^m, u_{z\tau}^m \rangle \geq M^{-1} \|u_{z\tau}^m\|^2.$$

As for the first term on the right of (4.4), use the Cauchy–Young inequality and assumption (f2) to obtain the estimate

$$\begin{aligned} \langle (f(u_\tau^m))_z, u_z^m \rangle &= \langle f'(u_\tau^m) u_{z\tau}^m, u_z^m \rangle + \langle f_z(u_\tau^m), u_z^m \rangle \\ &\leq \delta \|u_{z\tau}^m\|^2 + \frac{M}{\delta} \|u_z^m\|^2 + \|u_z^m\|^2 + M \|u_\tau^m\|^2, \end{aligned}$$

which is valid for any positive δ . Quite similarly, for the second term on the right, we have

$$\langle \beta_z u_{\tau\tau}^m, u_z^m \rangle = -\langle \beta_z u_\tau^m, u_{z\tau}^m \rangle \leq \delta \|u_{z\tau}^m\|^2 + \frac{M}{\delta} \|u_\tau^m\|^2.$$

The treatment of the remaining terms of (4.4) is evident. By putting the pieces together and choosing δ sufficiently small, we finally arrive at

$$\frac{d}{dz} \langle \alpha u_z^m, u_z^m \rangle + M^{-1} \|u_{z\tau}^m\|^2 \leq M \|u_z^m\|^2 + M \|u_\tau^m\|^2.$$

As in the case of estimate (4.3), apply the Gronwall lemma to obtain

$$(4.5) \quad \sup_z \|u_z^m\|^2 + \|u_{z\tau}^m\|^2 \leq M \|u_z^m(0)\|^2 + M \|u_\tau^m\|^2.$$

On the other hand, the multiplication of (4.1) by $-\tilde{u}_z^m$, after the appropriate integration by parts, yields

$$\langle \alpha u_z^m, u_z^m \rangle = \langle f(u_\tau^m), u_z^m \rangle + \langle \beta u_{\tau\tau}^m, u_z^m \rangle + \langle \gamma u^m, u_z^m \rangle + \langle \Delta \tilde{u}^m, u_z^m \rangle.$$

By applying the Cauchy–Young inequality in each of the terms on the right, we have

$$\|u_z^m\|^2 \leq M \|u_{\tau\tau}^m\|^2 + M \|\Delta \tilde{u}^m\|^2,$$

which allows us to eliminate $\|u_z^m(0)\|^2$ from the right-hand side of (4.5). By taking (4.3) into account, we can rewrite estimate (4.5) as

$$(4.6) \quad \sup_z \|u_z^m\|^2 + \|u_{z\tau}^m\|^2 \leq M \|u_{0\tau\tau}\|^2 + M \|\Delta \tilde{u}_0\|^2.$$

The missing bound on $\|\Delta u^m\|$ can be obtained directly from (4.1) with the help of estimates (4.3) and (4.6):

$$(4.7) \quad \begin{aligned} \|\Delta u^m\|^2 &\leq M \|u_{z\tau}^m\|^2 + M \|u_{\tau\tau\tau}^m\|^2 + M \|u_{z\tau}^m\|^2 + M \|u_\tau^m\|^2 \\ &\leq M \|u_{0\tau\tau}\|^2 + M \|\Delta \tilde{u}_0\|^2. \end{aligned}$$

Note that, according to the convention of section 3, we use the symbol M to denote *any* constant, although it is clear that the values of M in the first and in the second line of (4.7) cannot be the same.

By putting together the estimates (4.3), (4.6), and (4.7) we obtain

$$(4.8) \quad \begin{aligned} \sup_z \left[\|u_{\tau\tau}^m\|^2 + \|u_z^m\|^2 \right] + \|u_{\tau\tau\tau}^m\|^2 + \|u_{z\tau}^m\|^2 + \|\Delta u^m\|^2 \\ \leq M \|u_{0\tau\tau}\|^2 + M \|\Delta \tilde{u}_0\|^2. \end{aligned}$$

Thus the sequence u^m is bounded in U_2 and therefore contains a subsequence u^{m_j} such that

$$u_{\tau\tau\tau}^{m_j} \xrightarrow{U_0} u_{\tau\tau\tau}, \quad u_\tau^{m_j} \xrightarrow{U_0} u_\tau, \quad u_{z\tau}^{m_j} \xrightarrow{U_0} u_{z\tau}, \quad \Delta u^{m_j} \xrightarrow{U_0} \Delta u,$$

where $u \in U_2$. In order to show the convergence in the nonlinear term, notice that by (3.4) the sequence u_τ^m is bounded in $H^1(Q)$, and hence $u_\tau^{m_j} \rightarrow u_\tau$ strongly in U_0 .

By the Lipschitzness of f we also have that $f(u_\tau^{m_j}) \rightarrow f(u_\tau)$ in $L_2(Q)$. Moreover, the sequence $(f(u_\tau^{m_j}))_\tau = f'(u_\tau^{m_j})u_{\tau\tau}^{m_j}$ is bounded in U_0 , and therefore

$$(f(u_\tau^{m_j}))_\tau \xrightarrow{U_0} (f(u_\tau))_\tau.$$

Thus we can pass to the U_0 -weak limit in (4.1) and obtain

$$\alpha u_{z\tau} = (f(u_\tau))_\tau + \beta u_{\tau\tau\tau} + \gamma u_\tau + \Delta u$$

as an identity in U_0 . The compact embedding $U_2 \subset C([0, Z], X_0)$ implies the convergence $u^{m_j}(z) \rightarrow u(z)$ in X_0 for each z . In particular, $u(0) = u_0$. Finally, pass to the limit as $m \rightarrow \infty$ in estimate (4.8) to complete the proof. \square

5. Strong solutions: Further properties. The next two statements establish further properties of strong solutions and prepare the proof of the existence result for the weak solutions.

PROPOSITION 5.1. *Each strong solution u of problem (2.1)–(2.5) satisfies*

$$(5.1) \quad \sup_z \left[\|u_\tau\|^2 + \|\nabla u\|^2 \right] + \|u_{\tau\tau}\|^2 + \|\nabla u_\tau\|^2 \leq M \left[\|u_{0\tau}\|^2 + \|\nabla u_0\|^2 \right].$$

Proof. Multiply the equation

$$\alpha u_{z\tau} = (f(u_\tau))_\tau + \beta u_{\tau\tau\tau} + \gamma u_\tau + \Delta u$$

by u_τ , and integrate the result over Ω . The first and the last terms on the right will disappear. After obvious transformations in the remaining terms we obtain for almost all $z \in (0, Z)$ the inequality

$$(5.2) \quad \frac{d}{dz} \langle \alpha u_\tau, u_\tau \rangle + M^{-1} \|u_{\tau\tau}\|^2 \leq M \|u_\tau\|^2.$$

The multiplication of the equation by $-u_z$ suppresses the left-hand side and thus yields

$$(5.3) \quad 0 = \langle f(u_\tau), u_{z\tau} \rangle + \langle \beta u_{\tau\tau\tau}, u_{z\tau} \rangle - \langle \gamma u_\tau, u_z \rangle + \frac{1}{2} \frac{d}{dz} \|\nabla u\|^2.$$

The first term on the right of (5.3) can be transformed as follows:

$$(5.4) \quad \langle f(u_\tau), u_{z\tau} \rangle = \frac{d}{dz} \langle F(u_\tau), 1 \rangle - \langle F_z(u_\tau), 1 \rangle \geq \frac{d}{dz} \langle F(u_\tau), 1 \rangle - M \|u_\tau\|^2$$

(here we have used assumption (f3)). In order to eliminate the term $\langle \beta u_{\tau\tau\tau}, u_{z\tau} \rangle$ from relation (5.3), multiply the equation by $\alpha^{-1} \beta u_{\tau\tau}$. This time the second and the third terms on the right vanish, and thus we have

$$(5.5) \quad \begin{aligned} \langle \beta u_{\tau\tau\tau}, u_{z\tau} \rangle &= \langle \alpha^{-1} \beta f'(u_\tau) u_{\tau\tau\tau}, u_{\tau\tau\tau} \rangle + \langle \nabla(\alpha^{-1} \beta u_\tau), \nabla u_\tau \rangle \\ &\geq -M \|u_{\tau\tau}\|^2 + M^{-1} \|\nabla u_\tau\|^2 - M \|\nabla u_\tau\| \|u_\tau\| \\ &\geq -M \|u_{\tau\tau}\|^2 + M^{-1} \|\nabla u_\tau\|^2. \end{aligned}$$

Recall that, according to the convention about the use of M , its value may change from line to line.

Similarly, to get rid of $\langle \gamma u_\tau, u_z \rangle$ in (5.3), the proper multiplier is $\alpha^{-1}\gamma u$. It follows that

$$(5.6) \quad \begin{aligned} -\langle \gamma u_\tau, u_z \rangle &= -\langle \alpha^{-1}\gamma f(u_\tau), u_\tau \rangle - \langle \nabla(\alpha^{-1}\gamma u), \nabla u \rangle \\ &\geq -M\|u_\tau\|^2 - M\|\nabla u\|^2. \end{aligned}$$

By summing up relations (5.3)–(5.6), we arrive at

$$(5.7) \quad \frac{d}{dz} \left[\|\nabla u\|^2 + 2\langle F(u_\tau), 1 \rangle \right] + M^{-1}\|\nabla u_\tau\|^2 \leq M\|u_{\tau\tau}\|^2 + M\|\nabla u\|^2.$$

Choose a positive μ sufficiently large so that (5.2) and (5.7) imply that

$$(5.8) \quad \begin{aligned} \frac{d}{dz} \left[\mu\langle \alpha u_\tau, u_\tau \rangle + \|\nabla u\|^2 + 2\langle F(u_\tau), 1 \rangle \right] + M^{-1}\|u_{\tau\tau}\|^2 + M^{-1}\|\nabla u_\tau\|^2 \\ \leq M\|u_\tau\|^2 + M\|\nabla u\|^2. \end{aligned}$$

Make μ even larger, if necessary, to ensure that

$$\mu\langle \alpha u_\tau, u_\tau \rangle + 2\langle F(u_\tau), 1 \rangle \geq \|u_\tau\|^2.$$

Now the Gronwall lemma applied to (5.8) yields the desired estimate (5.1). \square

PROPOSITION 5.2. *Let u^1 and u^2 be two strong solutions of problem (2.1)–(2.5) corresponding to the initial conditions $u^1(0) = u_0^1$ and $u^2(0) = u_0^2$, respectively. Then*

$$(5.9) \quad \sup_z \|u_\tau^1 - u_\tau^2\|^2 + \|u_{\tau\tau}^1 - u_{\tau\tau}^2\|^2 \leq M\|u_{0\tau}^1 - u_{0\tau}^2\|^2.$$

In particular, this estimate proves the uniqueness of the strong solution.

Proof. The difference $w = u^1 - u^2$ satisfies the equation

$$(5.10) \quad \alpha w_{z\tau} = (\bar{f}w_\tau)_\tau + \beta w_{\tau\tau\tau} + \gamma w_\tau + \Delta w,$$

where $\bar{f} = (f(u_\tau^1) - f(u_\tau^2))/(u_\tau^1 - u_\tau^2)$ is a bounded function. Multiply it by w_τ , and integrate the result over Ω . The last term on the right of (5.10) disappears, and we obtain

$$\begin{aligned} \frac{1}{2} \frac{d}{dz} \langle \alpha w_\tau, w_\tau \rangle + \langle \beta w_{\tau\tau\tau}, w_{\tau\tau} \rangle &= \frac{1}{2} \langle \alpha_z w_\tau, w_\tau \rangle - \langle \bar{f}w_\tau, w_{\tau\tau} \rangle + \langle \gamma w_\tau, w_\tau \rangle \\ &\leq M\|w_\tau\| \|w_{\tau\tau}\| + M\|w_\tau\|^2 \\ &\leq \frac{1}{2} \langle \beta w_{\tau\tau\tau}, w_{\tau\tau} \rangle + M\|w_\tau\|^2. \end{aligned}$$

Apply the Gronwall lemma to complete the proof. \square

Remark. Inequality (5.9) is obtained from (5.10) with the help of an argument similar to the one used in the first step of the proof of Proposition 5.1 (estimate (5.2)). However, the equivalent of (5.1) does not appear to be possible for (5.10) because of the special treatment required for the nonlinear term (see (5.4)).

6. Weak solutions. We have already qualified U_1 as the space of weak solutions. In order to introduce the weak solutions we also need the space of test functions. Let \hat{Y} be the closure of D_Q with respect to the $H^1(Q)$ -norm. Denote by Y the set of all functions $\eta \in \hat{Y}$ such that $\eta(Z) = 0$. The latter condition makes sense because of the continuous (and even compact) embedding $\hat{Y} \subset C([0, Z], X_0)$.

DEFINITION 6.1. For a given $u_0 \in X_1$, a function $u \in U_1$ is called a weak solution of problem (2.1)–(2.5) if the identity

$$(6.1) \quad \begin{aligned} & -\langle \alpha(0)u_{0\tau}, \eta(0) \rangle - \int_0^Z \langle u_\tau, (\alpha\eta)_z \rangle dz \\ & = \int_0^Z \left[-\langle f(u_\tau), \eta_\tau \rangle - \langle \beta u_{\tau\tau}, \eta_\tau \rangle + \langle \gamma u_\tau, \eta \rangle - \langle \nabla u, \nabla \eta \rangle \right] dz \end{aligned}$$

holds for any $\eta \in Y$.

By multiplying (2.1) by η and integrating by parts where necessary, one easily shows that each strong solution is also a weak one.

PROPOSITION 6.2. For any $u_0 \in X_1$ problem (2.1)–(2.5) admits a weak solution u satisfying

$$(6.2) \quad \sup_z \left[\|u_\tau\|^2 + \|\nabla u\|^2 \right] + \|u_{\tau\tau}\|^2 + \|\nabla u_\tau\|^2 \leq M \left[\|u_{0\tau}\|^2 + \|\nabla u_0\|^2 \right].$$

Proof. Let u_0^m , $m = 1, 2, \dots$, be a sequence in X_2 such that $u_0^m \rightarrow u_0$ in X_1 as $m \rightarrow \infty$. For each m denote by u^m the strong solution corresponding to the initial condition $u^m(0) = u_0^m$. It is also a weak solution, that is,

$$(6.3) \quad \begin{aligned} & -\langle \alpha(0)u_{0\tau}^m, \eta(0) \rangle - \int_0^Z \langle u_\tau^m, (\alpha\eta)_z \rangle dz \\ & = \int_0^Z \left[-\langle f(u_\tau^m), \eta_\tau \rangle - \langle \beta u_{\tau\tau}^m, \eta_\tau \rangle + \langle \gamma u_\tau^m, \eta \rangle - \langle \nabla u^m, \nabla \eta \rangle \right] dz \end{aligned}$$

for any $\eta \in Y$. By Proposition 5.2,

$$\|u_{\tau\tau}^m - u_{\tau\tau}^n\|^2 \leq M \|u_{0\tau}^m - u_{0\tau}^n\|^2 \xrightarrow{m, n \rightarrow \infty} 0,$$

so that there exists an element $u \in U_0$ such that $u_{\tau\tau} \in U_0$ and $u_{\tau\tau}^m \rightarrow u_{\tau\tau}$ in U_0 as $m \rightarrow \infty$. This also implies that $u_\tau^m \rightarrow u_\tau$ in U_0 and therefore $f(u_\tau^m) \rightarrow f(u_\tau)$ in $L_2(Q)$ as well. Finally notice that, by Proposition 5.1,

$$(6.4) \quad \sup_z \left[\|u_\tau^m\|^2 + \|\nabla u^m\|^2 \right] + \|u_{\tau\tau}^m\|^2 + \|\nabla u_\tau^m\|^2 \leq M \left[\|u_{0\tau}^m\|^2 + \|\nabla u_0^m\|^2 \right],$$

so that the sequence ∇u^m is bounded in U_0 , and therefore ∇u^m converges to ∇u weakly in U_0 . Thus we can pass to the limit in (6.3) to obtain (6.1).

In order to establish estimate (6.2), first observe that the sequence ∇u_τ^m is bounded in U_0 and therefore contains a subsequence weakly convergent to ∇u_τ in U_0 . Similarly, the sequences u_τ^m and ∇u^m contain subsequences $*$ -weakly convergent, respectively, to u_τ and ∇u in $L_\infty([0, Z], X_0)$. One can therefore justify the limit as $m \rightarrow \infty$ in (6.4) and thus obtain (6.2). \square

PROPOSITION 6.3. The weak solution is unique.

Proof. Let u^1 and u^2 be two weak solutions corresponding to the same initial condition $u^1(0) = u^2(0) = u_0 \in X_1$. Their difference $w = u^1 - u^2$ satisfies the identity

$$(6.5) \quad \begin{aligned} & - \int_0^Z \langle w_\tau, (\alpha\eta)_z \rangle dz \\ & = \int_0^Z \left[-\langle \bar{f}w_\tau, \eta_\tau \rangle - \langle \beta w_{\tau\tau}, \eta_\tau \rangle + \langle \gamma w_\tau, \eta \rangle - \langle \nabla w, \nabla \eta \rangle \right] dz \end{aligned}$$

for any $\eta \in Y$. As well as in the proof of Proposition 5.2, here \bar{f} stands for the bounded function $(f(u_\tau^1) - f(u_\tau^2))/(u_\tau^1 - u_\tau^2)$.

Formally, the estimate we are going to obtain for w is quite simple and resembles much the one of Proposition 5.2. However, the direct approach of the proof of Proposition 5.2 is not valid any more since it requires the manipulation of objects which no longer exist in the case of the weak solutions. In particular, the weak solutions have no derivatives with respect to z .

Let the functions ψ_k and the projections P_m be those defined in the proof of Proposition 4.2. Given a function $\phi \in C_0^\infty([0, Z])$, set $\eta = \alpha^{-1}\phi\tilde{\psi}_k$ in (6.5). After some obvious transformations this yields

$$\left| \int_0^Z \phi_z \langle w, \psi_k \rangle dz \right| \leq C \int_0^Z |\phi| dz \quad \forall \phi \in C_0^\infty([0, Z]),$$

with some positive C depending on w and k but independent of ϕ . By the Riesz theorem, this proves the existence of a function $\theta \in L_\infty([0, Z])$ such that

$$\int_0^Z \phi_z \langle w, \psi_k \rangle dz = - \int_0^Z \phi \theta dz \quad \forall \phi \in C_0^\infty([0, Z]),$$

and thus the function $\langle w, \psi_k \rangle$ admits a derivative $\langle w, \psi_k \rangle_z \in L_\infty([0, Z])$ understood in the sense of distributions. Therefore, for any m , the function $w^m \equiv P_m w$ also admits a bounded derivative with respect to z .

Fix positive numbers $s < Z$, $\delta < s$, and set

$$\sigma(z) = \begin{cases} 1 & \text{if } z \leq s - \delta, \\ (s - z)/\delta & \text{if } s - \delta < z \leq s, \\ 0 & \text{if } z > s. \end{cases}$$

By taking into account the existence of a bounded derivative w_z^m , it is easy to see that the function $\sigma\alpha^{-1}P_m(\alpha w_\tau^m)$ belongs to the class Y . Inject it in (6.5) instead of η . The integrand at the left becomes

$$\begin{aligned} -\langle w_\tau, (\alpha\eta)_z \rangle &= -\langle w_\tau, P_m(\sigma\alpha w_\tau^m)_z \rangle = -\langle w_\tau^m, (\sigma\alpha w_\tau^m)_z \rangle \\ &= -\frac{1}{2}\sigma_z \langle \alpha w_\tau^m, w_\tau^m \rangle - \frac{\sigma}{2} \langle \alpha_z w_\tau^m, w_\tau^m \rangle - \frac{1}{2} \frac{d}{dz} [\sigma \langle \alpha w_\tau^m, w_\tau^m \rangle]. \end{aligned}$$

By integrating over $[0, Z]$ and recalling the definition of σ , one obtains for the left-hand side of (6.5) the following lower bound:

$$\begin{aligned} & - \int_0^Z \langle w_\tau, (\alpha\eta)_z \rangle dz \\ & \geq \frac{1}{2\delta} \int_{s-\delta}^s \langle \alpha w_\tau^m, w_\tau^m \rangle dz - M \int_0^s \|w_\tau^m\|^2 dz + \frac{1}{2} \langle \alpha(0)w_\tau^m(0), w_\tau^m(0) \rangle \\ & \geq \frac{1}{2\delta} \int_{s-\delta}^s \langle \alpha w_\tau^m, w_\tau^m \rangle dz - M \int_0^s \|w_\tau^m\|^2 dz. \end{aligned}$$

As for the right-hand side of (6.5), the first and the third terms of the integrand admit the estimate

$$-\langle \bar{f}w_\tau, \eta_\tau \rangle + \langle \gamma w_\tau, \eta \rangle \leq \sigma M \|w_\tau\| \|w_\tau^m\|,$$

while the other two terms will be kept intact for the moment. By summing up and passing to the limit as $\delta \rightarrow \infty$, we have

$$\begin{aligned} & \frac{1}{2} \langle \alpha(s)w_\tau^m(s), w_\tau^m(s) \rangle \\ & \leq \int_0^s \left[M \|w_\tau\| \|w_{\tau\tau}^m\| - \langle \beta w_{\tau\tau}, \alpha^{-1} P_m(\alpha w_{\tau\tau}^m) \rangle - \langle \nabla w, \nabla \alpha^{-1} P_m(\alpha w_\tau^m) \rangle \right] dz. \end{aligned}$$

The left-hand side of this relation admits no limit as $m \rightarrow \infty$, so first let us integrate it once more with respect to s :

$$\begin{aligned} & \frac{1}{2} \int_0^\xi \langle \alpha(s)w_\tau^m(s), w_\tau^m(s) \rangle ds \\ & \leq \int_0^\xi \int_0^s \left[M \|w_\tau\| \|w_{\tau\tau}^m\| - \langle \beta w_{\tau\tau}, \alpha^{-1} P_m(\alpha w_{\tau\tau}^m) \rangle - \langle \nabla w, \nabla \alpha^{-1} P_m(\alpha w_\tau^m) \rangle \right] dz ds. \end{aligned}$$

Recall that $w \in U_1$ and therefore $w^m \rightarrow w$ in U_1 . Thus one can pass to the limit as $m \rightarrow \infty$ in all terms to obtain

$$\frac{1}{2} \int_0^\xi \langle \alpha(s)w_\tau(s), w_\tau(s) \rangle ds \leq \int_0^\xi \int_0^s \left[M \|w_\tau\| \|w_{\tau\tau}\| - \langle \beta w_{\tau\tau}, w_{\tau\tau} \rangle \right] dz ds.$$

With the help of the Cauchy–Young inequality this can be rewritten as

$$\int_0^\xi \langle \alpha w_\tau, w_\tau \rangle dz \leq M \int_0^\xi \int_0^s \langle \alpha w_\tau, w_\tau \rangle dz ds.$$

It remains to apply the Gronwall lemma to obtain $w = 0$. □

7. Stratified media. The present section is devoted to the study of the asymptotic behavior of the solutions to problem (2.1)–(2.5) as its coefficients rapidly oscillate. This study will be undertaken under the following additional assumptions:

- ($\beta\gamma$) The coefficients α , β , and γ depend only on z . It is clear that in this case the replacement of z by the new variable $\tilde{z} = \int_0^z \alpha(s) ds$ will bring the equation to the form with $\alpha = 1$. Thus we assume that $\alpha = 1$ throughout the section. The regularity of β and γ is the same as above; that is, β , β_z , and γ are bounded, while β is also positive and bounded away from zero.
- ($f4$) The nonlinear function f is of the form $f(s, z, x) = \rho(z)\phi(s)$, where ρ is bounded with a bounded derivative, while ϕ satisfies the same assumptions as f in (f1)–(f3), that is, ϕ admits a bounded continuous derivative and $\phi(0) = 0$. We shall denote by Φ the primitive of ϕ satisfying $\Phi(0) = 0$.

Thus the equation we study takes the form

$$(7.1) \quad u_{z\tau} = \rho(\phi(u_\tau))_\tau + \beta u_{\tau\tau\tau} + \gamma u_\tau + \Delta u.$$

The existence and uniqueness results of sections 4–6 remain valid.

In order to introduce the oscillating character of the coefficients, let us assume that they are of the form

$$\rho(z) = r(z, z/\epsilon), \quad \beta(z) = b(z, z/\epsilon), \quad \gamma(z) = g(z, z/\epsilon),$$

where $\epsilon < 1$ is a small parameter. The last assumption we introduce is the following.

($\rho\beta\gamma$) The functions $r(z, \zeta)$, $b(z, \zeta)$, and $g(z, \zeta)$ are bounded as well as the derivatives r_z , r_ζ , b_z , and b_ζ . The function b is positive and bounded away from zero. Finally, the functions $r(z, z/\epsilon)$, $b(z, z/\epsilon)$, and $g(z, z/\epsilon)$ admit limits as $\epsilon \rightarrow 0$ in the following weak sense: There exists $q > 1/2$ such that for any $Z > 0$ we have

$$\begin{aligned} \int_0^Z (r(z, z/\epsilon) - \bar{\rho}(z)) dz &= O(\epsilon^q), \\ \int_0^Z (b(z, z/\epsilon) - \bar{\beta}(z)) dz &= O(\epsilon^q), \\ \int_0^Z (g(z, z/\epsilon) - \bar{\gamma}(z)) dz &= O(\epsilon^q), \end{aligned}$$

where the functions $\bar{\rho}(z)$, $\bar{\beta}(z)$, and $\bar{\gamma}(z)$ as well as the derivatives $\bar{\rho}_z(z)$, and $\bar{\beta}_z(z)$ are bounded, while $\bar{\rho}(z)$ and $\bar{\beta}(z)$ are also positive and bounded away from zero.

Note that if r , b , and g are Lipschitzian with respect to the first argument and periodic with respect to the second one, then the above convergence conditions hold with $q = 1$ (see [1, 7, 9]).

The asymptotic study of the solutions to (7.1) as $\epsilon \rightarrow 0$ require a few estimates. The convention concerning the use of the symbol M introduced in section 2 is still valid with the following remark: the constants M are independent of ϵ .

The first estimate we need is similar to the one of Proposition 5.1.

PROPOSITION 7.1. *Let u be the strong solution of (7.1) corresponding to the initial condition $u(0) = u_0 \in X_2$. Then*

$$\begin{aligned} \sup_z \|\nabla u\|^2 + \frac{1}{\epsilon} \sup_z \|u_\tau\|^2 + M^{-1} \int_0^Z \|u_{\tau\tau}\|^2 dz + M^{-1} \int_0^Z \|\nabla u_\tau\|^2 dz \\ \leq M \|\nabla u_0\|^2 + \frac{M}{\epsilon} \|u_{0\tau}\|^2. \end{aligned}$$

Proof. Multiply (7.1) by u_τ to obtain

$$(7.2) \quad \frac{1}{2} \frac{d}{dz} \|u_\tau\|^2 = -\langle \beta u_{\tau\tau}, u_{\tau\tau} \rangle + \langle \gamma u_\tau, u_\tau \rangle \leq -M^{-1} \|u_{\tau\tau}\|^2 + M \|u_\tau\|^2.$$

The multiplication of (7.1) by u_z yields

$$(7.3) \quad \frac{1}{2} \frac{d}{dz} \|\nabla u\|^2 + \frac{d}{dz} \langle \rho \Phi(u_\tau), 1 \rangle = -\langle \beta u_{\tau\tau}, u_{\tau z} \rangle + \langle \rho_z \Phi(u_\tau), 1 \rangle + \langle \gamma u_\tau, u_z \rangle \\ \leq -\langle \beta u_{\tau\tau}, u_{\tau z} \rangle + M \epsilon^{-1} \|u_\tau\|^2 + \langle \gamma u_\tau, u_z \rangle.$$

In order to eliminate the first term on the right-hand side of (7.3), multiply (7.1) by $\beta u_{\tau\tau}$:

$$(7.4) \quad \langle \beta u_{\tau\tau}, u_{\tau z} \rangle = \langle \rho \beta \phi'(u_\tau) u_{\tau\tau}, u_{\tau\tau} \rangle + \langle \beta \nabla u_\tau, \nabla u_\tau \rangle \\ \geq -M \|u_{\tau\tau}\|^2 + M^{-1} \|\nabla u_\tau\|^2.$$

Similarly, after the multiplication of (7.1) by γu , we have

$$(7.5) \quad \langle \gamma u_\tau, u_z \rangle \leq M \|u_\tau\|^2 + M \|\nabla u\|^2.$$

By combining (7.3)–(7.5), we have

$$(7.6) \quad \begin{aligned} & \frac{1}{2} \frac{d}{dz} \|\nabla u\|^2 + \frac{d}{dz} \langle \rho \Phi(u_\tau), 1 \rangle + M^{-1} \|\nabla u_\tau\|^2 \\ & \leq M \|u_{\tau\tau}\|^2 + M \|\nabla u\|^2 + M \epsilon^{-1} \|u_\tau\|^2. \end{aligned}$$

For a given positive μ , set

$$E = \frac{1}{2} \|\nabla u\|^2 + \langle \rho \Phi(u_\tau), 1 \rangle + \frac{\mu}{\epsilon} \|u_\tau\|^2.$$

If the parameter μ is sufficiently large, then

$$E \geq \frac{1}{2} \|\nabla u\|^2 + \frac{\mu}{2\epsilon} \|u_\tau\|^2.$$

Fix μ even larger, if necessary, so that the inequalities (7.2) and (7.6) imply that

$$\frac{dE}{dz} + M^{-1} \|u_{\tau\tau}\|^2 + M^{-1} \|\nabla u_\tau\|^2 \leq ME.$$

It remains to apply the Gronwall lemma to obtain the claimed estimate. \square

The next estimate requires some new notation. The operator $I - \Delta : L_2(\omega) \rightarrow L_2(\omega)$ defined on $H^2(\omega)$ with the Neumann boundary conditions is self-adjoint and positive definite. Denote it by A^2 . Thus A is its square root, which is also positive definite. Clearly, we have $M^{-1} \|A^{-1}u\| \leq \|u\| \leq M \|Au\|$. Note also that

$$\|Au\|^2 = \langle A^2u, u \rangle = \langle u - \Delta u, u \rangle = \|u\|^2 + \|\nabla u\|^2.$$

On X_1 , the norm $\|Au\|$ is equivalent to $\|\nabla u\|$.

PROPOSITION 7.2. *Let u be the strong solution of (7.1) corresponding to the initial condition $u(0) = u_0 \in X_2$. Then*

$$(7.7) \quad \sup_z \|A^{-1}u_{\tau\tau}\|^2 + \int_0^Z \|A^{-1}u_{z\tau}\|^2 dz \leq M \|\nabla u_0\|^2 + \frac{M}{\epsilon} \|u_{0\tau}\|^2 + \|A^{-1}u_{0\tau\tau}\|^2.$$

Proof. The multiplication of (7.1) by $A^{-2}u_{z\tau}$ yields

$$\begin{aligned} & \|A^{-1}u_{z\tau}\|^2 + \beta \frac{d}{dz} \|A^{-1}u_{\tau\tau}\|^2 \\ & \leq M \|u_{\tau\tau}\| \|A^{-2}u_{z\tau}\| + M \|u_\tau\| \|A^{-2}u_{z\tau}\| + M \|A^{-1}\Delta u\| \|A^{-1}u_{z\tau}\| \\ & \leq M \|u_{\tau\tau}\| \|A^{-1}u_{z\tau}\| + M \|\nabla u\| \|A^{-1}u_{z\tau}\|, \end{aligned}$$

which implies, after the division by β , that

$$\beta^{-1} \|A^{-1}u_{z\tau}\|^2 + \frac{d}{dz} \|A^{-1}u_{\tau\tau}\|^2 \leq M \|u_{\tau\tau}\|^2 + M \|\nabla u\|^2.$$

Integrate with respect to z to obtain

$$\sup_z \|A^{-1}u_{\tau\tau}\|^2 + \int_0^Z \|A^{-1}u_{z\tau}\|^2 dz \leq M \int_0^Z [\|u_{\tau\tau}\|^2 + \|\nabla u\|^2] dz + \|A^{-1}u_{0\tau\tau}\|^2.$$

Finally, make use of Proposition 7.1 to complete the proof. \square

We can now state and prove the main result of this section.

THEOREM 7.3. *Let u be the strong solution of (7.1) corresponding to the initial condition $u(0) = u_0 \in X_2$. Denote by v the strong solution of the equation*

$$(7.8) \quad v_{z\tau} = \bar{\rho}(\phi(v_\tau))_\tau + \bar{\beta}v_{\tau\tau\tau} + \bar{\gamma}v_\tau + \Delta v$$

satisfying the same initial condition as u . Then

$$\sup_z \|u - v\|^2 \leq M\epsilon^{q-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right].$$

Proof. It is clear that v is subject to the same type of estimates as u but with the coefficients independent of ϵ . In particular,

$$(7.9) \quad \sup_z \|\nabla v\|^2 + \sup_z \|v_\tau\|^2 + M^{-1} \int_0^Z \|v_{\tau\tau}\|^2 dz + M^{-1} \int_0^Z \|\nabla v_\tau\|^2 dz \\ \leq M\|\nabla u_0\|^2 + M\|u_{0\tau}\|^2.$$

The difference $w = u - v$ satisfies the relation

$$(7.10) \quad w_{z\tau} = \rho[\phi(u_\tau) - \phi(v_\tau)]_\tau + \beta w_{\tau\tau\tau} + \gamma w_\tau + \Delta w \\ - (\bar{\rho} - \rho)(\phi(v_\tau))_\tau + (\bar{\beta} - \beta)v_{\tau\tau\tau} + (\bar{\gamma} - \gamma)v_\tau.$$

Multiply it by $-\tilde{w}$ to obtain

$$(7.11) \quad \frac{1}{2} \frac{d}{dz} \|w\|^2 = \langle \rho[\phi(u_\tau) - \phi(v_\tau)], w \rangle - \langle \beta w_\tau, w_\tau \rangle + \langle \gamma w, w \rangle \\ + (\bar{\rho} - \rho)P + (\bar{\beta} - \beta)Q + (\bar{\gamma} - \gamma)R \\ \leq M\|w\|^2 + (\bar{\rho} - \rho)P + (\bar{\beta} - \beta)Q + (\bar{\gamma} - \gamma)R,$$

where P , Q , and R stand for the residual terms:

$$P = \langle \phi(v_\tau), u - v \rangle, \quad Q = \langle v_{\tau\tau}, u - v \rangle, \quad R = \langle v, u - v \rangle.$$

The integration of (7.11) over $z \in [0, s]$ yields

$$\frac{1}{2} \|w(s)\|^2 \leq M \int_0^s \|w\|^2 dz + \int_0^s [(\bar{\alpha} - \alpha)P + (\bar{\beta} - \beta)Q + (\bar{\gamma} - \gamma)R] dz.$$

By integrating by parts in the second term on the right and making use of assumption $(\rho\beta\gamma)$, we arrive at

$$(7.12) \quad \frac{1}{2} \|w(s)\|^2 \leq M \int_0^s \|w\|^2 dz + M\epsilon^q \int_0^s (|P_z| + |Q_z| + |R_z|) dz.$$

The expressions P_z , Q_z , and R_z need to be estimated separately. The most difficult estimate is required by the term $|\langle v_{\tau\tau}, u_z \rangle|$ appearing in $|Q_z|$. We have

$$(7.13) \quad |\langle v_{\tau\tau}, u_z \rangle| = |\langle (I - \Delta)v_\tau, A^{-2}u_{z\tau} \rangle| \\ \leq |\langle v_\tau, A^{-2}u_{z\tau} \rangle| + |\langle \nabla v_\tau, \nabla A^{-2}u_{z\tau} \rangle| \\ \leq \|A^{-1}v_\tau\| \|A^{-1}u_{z\tau}\| + \|\nabla v_\tau\| \|\nabla A^{-2}u_{z\tau}\|.$$

On the other hand,

$$\|\nabla A^{-2}u_{z\tau}\|^2 = -\langle \Delta A^{-2}u_{z\tau}, A^{-2}u_{z\tau} \rangle = \langle (A^2 - I)A^{-2}u_{z\tau}, A^{-2}u_{z\tau} \rangle \\ = \langle A^{-1}u_{z\tau}, A^{-1}u_{z\tau} \rangle - \langle A^{-2}u_{z\tau}, A^{-2}u_{z\tau} \rangle \leq \|A^{-1}u_{z\tau}\|^2.$$

Thus (7.13) becomes

$$|\langle v_{\tau\tau}, u_z \rangle| \leq (\|v_\tau\| + \|\nabla v_\tau\|) \|A^{-1}u_{z\tau}\|.$$

Now with the help of (7.7) and (7.9) we obtain

$$\begin{aligned} \int_0^s |\langle v_{\tau\tau}, u_z \rangle| dz &\leq M \left[\int_0^s (\|v_\tau\|^2 + \|\nabla v_\tau\|^2) dz \right]^{1/2} \left[\int_0^s \|A^{-1}u_{z\tau}\|^2 dz \right]^{1/2} \\ (7.14) \quad &\leq M \left[\|\nabla u_0\|^2 + \|u_{0\tau}\|^2 \right]^{1/2} \left[\|\nabla u_0\|^2 + \frac{1}{\epsilon} \|u_{0\tau}\|^2 + \|A^{-1}u_{0\tau\tau}\|^2 \right]^{1/2} \\ &\leq M\epsilon^{-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right]. \end{aligned}$$

The latter inequality follows from the boundedness of the operator A^{-1} and the fact that $\|\nabla u_0\|^2 = \langle \nabla u_0, \nabla u_0 \rangle = -\langle \nabla \tilde{u}_0, \nabla u_{0\tau} \rangle = \langle \Delta \tilde{u}_0, u_{0\tau} \rangle \leq \frac{1}{2}(\|\Delta \tilde{u}_0\|^2 + \|u_{0\tau}\|^2)$. Quite similarly to (7.14) one establishes the same upper bound on the term $|\langle v_{\tau\tau}, v_z \rangle|$, and therefore

$$\int_0^s |\langle v_{\tau\tau}, u_z - v_z \rangle| dz \leq M\epsilon^{-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right].$$

Now let us note that

$$\begin{aligned} (7.15) \quad \int_0^s |Q_z| dz &\leq \int_0^s |\langle v_{\tau\tau}, u - v \rangle_z| dz \\ &\leq \int_0^s |\langle v_{\tau\tau}, u_z - v_z \rangle| dz + \int_0^s |\langle v_z, u_{\tau\tau} \rangle| dz + \int_0^s |\langle v_z, v_{\tau\tau} \rangle| dz. \end{aligned}$$

The term involving u_z having already been estimated, let us treat the second term on the right of (7.15). Recall that v is a strong solution to (7.8) whose coefficients are independent of ϵ . Therefore, by Proposition 4.2,

$$(7.16) \quad \sup_z \left[\|v_{\tau\tau}\|^2 + \|v_z\|^2 \right] + \|v_{\tau\tau\tau}\|^2 + \|v_{z\tau}\|^2 + \|\Delta v\|^2 \leq M \|u_{0\tau\tau}\|^2 + M \|\Delta\tilde{u}_0\|^2.$$

Along with Proposition 7.1 the latter inequality implies that

$$\begin{aligned} \int_0^s |\langle v_z, u_{\tau\tau} \rangle| dz &\leq \left[\int_0^s \|v_z\|^2 dz \right]^{1/2} \left[\int_0^s \|u_{\tau\tau}\|^2 dz \right]^{1/2} \\ &\leq M \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right]^{1/2} \times \epsilon^{-1/2} \left[\|\nabla u_0\|^2 + \|u_{0\tau}\|^2 \right]^{1/2}. \end{aligned}$$

As has been mentioned above, $\|\nabla u_0\|^2 \leq \frac{1}{2}(\|\Delta\tilde{u}_0\|^2 + \|u_{0\tau}\|^2)$ and $\|u_{0\tau}\|^2 \leq \|u_{0\tau\tau}\|^2$. So we obtain

$$\int_0^s |\langle v_z, u_{\tau\tau} \rangle| dz \leq M\epsilon^{-1/2} (\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2).$$

The last term on the right of (7.15) admits the same estimate, and we finally obtain

$$\int_0^s |Q_z| dz \leq M\epsilon^{-1/2} (\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2).$$

The estimate on $|P_z|$ is simpler. We have

$$P_z = \langle \phi'(v_\tau)v_{z\tau}, u - v \rangle + \langle \phi(v_\tau), u_z - v_z \rangle.$$

Estimate (7.16) together with Proposition 7.1 yields

$$(7.17) \quad \int_0^s |P_z| dz \leq M\epsilon^{-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right] + \int_0^s |\langle \phi(v_\tau), u_z \rangle| dz.$$

The last term on the right should be treated with the help of Proposition 7.2. First note that

$$\begin{aligned} |\langle \phi(v_\tau), u_z \rangle| &= |\langle A\phi(v_\tau), A^{-1}u_z \rangle| \leq \|A\phi(v_\tau)\| \|A^{-1}u_z\| \\ &\leq M \|\nabla\phi(v_\tau)\| \|A^{-1}u_z\| = M \|\phi'(v_\tau)\nabla v_\tau\| \|A^{-1}u_z\| \\ &\leq M \|\nabla v_\tau\| \|A^{-1}u_z\| \leq M \left[\|\Delta v\|^2 + \|v_{\tau\tau}\|^2 \right]^{1/2} \|A^{-1}u_z\|. \end{aligned}$$

Now the integration over $z \in [0, s]$ and the use of Proposition 7.2 along with estimate (7.16) yield

$$\int_0^s |\langle \phi(v_\tau), u_z \rangle| dz \leq M\epsilon^{-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right].$$

We omit the details concerning the treatment of the term involving R_z in the second integral on the right of (7.12). It is quite similar to the case of Q_z since $v = \tilde{v}_{\tau\tau}$. Thus we finally obtain

$$(7.18) \quad \frac{1}{2} \|w(s)\|^2 \leq M \int_0^s \|w\|^2 dz + M\epsilon^{q-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right].$$

By the Gronwall lemma we have

$$\int_0^s \|w\|^2 dz \leq M\epsilon^{q-1/2} \left[\|u_{0\tau\tau}\|^2 + \|\Delta\tilde{u}_0\|^2 \right].$$

Reinject the latter inequality into (7.18) to complete the proof. \square

8. Concluding remarks. The obtained global existence and uniqueness results for (2.1) for the case of nonlinearity f with bounded derivative show that this formulation may have some mathematical advantages in comparison with the classical KZK equation (see [3, 4]), although for small variations of the acoustic pressure (which is one of the assumptions of the physical model used to derive the equation) these two models coincide.

The introduction of the varying coefficients allows us to model the acoustic beam propagation in a stratified heterogeneous media (for example, in the atmosphere), while the asymptotic analysis by means of the homogenization techniques simplifies reducing it to the case of constant coefficients, when in some situations an analytical solution is possible.

REFERENCES

- [1] N. S. BAKHVALOV AND G. P. PANASENKO, *Homogenization: Averaging Processes in Periodic Media*, Nauka, Moscow, 1984 (in Russian); Kluwer, Dordrecht, 1989, (in English).
- [2] N. S. BAKHVALOV, YA. M. ZHILEIKIN, AND E. A. ZABOLOTSKAYA, *Nonlinear Theory of Sound Beams*, Nauka, Moscow (in Russian); American Institute of Physics, New York, 1987 (in English).
- [3] C. BARDOS AND A. ROZANOVA, *KZK equation*, in Proceedings of the International Crimean Autumn Mathematical School-Symposium, Laspi-Batiliman, Ukraine, 2004.
- [4] C. BARDOS AND A. ROZANOVA, *The Khokhlov-Zabolotskaya-Kuznetsov equation*, C. R. Acad. Sci. Paris Sér. I, 344 (2007), pp. 337–342.
- [5] A. V. FAMINSKY, *Cauchy problem for Zakharov-Kuznetsov equation*, Differ. Equ., 31 (1995), pp. 1070–1081 (in Russian).
- [6] A. V. FAMINSKY, *Cauchy problem for quasi-linear equations of odd order*, Math USSR Sbornik, 180 (1989), pp. 1183–1210 (in Russian).
- [7] E. A. LAPSHIN AND G. P. PANASENKO, *Homogenization of the equations of high frequency nonlinear acoustics*, C. R. Acad. Sci. Paris Sér. I, 325 (1997), pp. 931–936.
- [8] B. K. NOVIKOV, O. V. RUDENKO, AND V. I. TIMOSHENKO, *Nonlinear Underwater Acoustics*, American Institute of Physics, New York, 1987.
- [9] E. A. LAPSHIN AND G. P. PANASENKO, *Homogenization of high frequency nonlinear acoustics equations*, Appl. Anal., 74 (2000), pp. 311–331.
- [10] O. V. RUDENKO, *Nonlinear sawtooth-shaped waves*, Physics-Uspekhi, 38 (1995), pp. 965–989.
- [11] O. V. RUDENKO AND S. I. SOLUYAN, *Theoretical Foundations of Nonlinear Acoustics*, Plenum, New York, 1977.
- [12] O. V. RUDENKO, A. K. SUKHORUKOVA, AND A. P. SUKHORUKOV, *Equations of high frequency nonlinear acoustics of heterogeneous media*, Acoustic J., 40 (1994), pp. 290–294 (in Russian).
- [13] A. P. SARVAZYAN, O. V. RUDENKO, S. D. SWANSON, J. B. FOWLKES, AND S. Y. EMELIANOV, *Shear wave elasticity imaging: A new ultrasonic technology of medical diagnostics*, Ultrasound in Med. and Biol., 24 (1999), pp. 1419–1435.
- [14] O. A. VASIL'eva, A. A. KARABUTOV, E. A. LAPSHIN, AND O. V. RUDENKO, *Interaction of One-Dimensional Waves in Dispersion-Free Media*, Moscow University, Moscow, 1983, (in Russian).
- [15] I. KOSTIN AND G. PANASENKO, *Khokhlov-Zabolotskaya-Kuznetsov type equation: Nonlinear acoustic in heterogeneous media*, C. R. Mécanique, 334 (2006), pp. 220–224.

EXISTENCE OF WEAK SOLUTIONS FOR THE UNSTEADY INTERACTION OF A VISCOUS FLUID WITH AN ELASTIC PLATE*

CÉLINE GRANDMONT†

Abstract. We consider a three-dimensional viscous incompressible fluid governed by the Navier–Stokes equations, interacting with an elastic plate located on one part of the fluid boundary. We do not neglect the deformation of the fluid domain which consequently depends on the displacement of the structure. The purpose of this work is to study the solutions of this unsteady fluid-structure interaction problem, as the coefficient modeling the viscoelasticity (resp., the rotatory inertia) of the plate tends to zero. As a consequence, we obtain the existence of at least one weak solution for the limit problem (Navier–Stokes equation coupled with a plate in flexion) as long as the structure does not touch the bottom of the fluid cavity.

Key words. fluid-structure interaction, weak solutions, plate equations

AMS subject classifications. 74F10, 35Q30, 73K70, 76D03, 35Q35, 35D05

DOI. 10.1137/070699196

1. Introduction. Many physical phenomena deal with a fluid interacting with a moving or deformable structure. These kinds of problems have a lot of important applications, for instance, in areolasticity, biomechanics, hydroelasticity, sedimentation, etc. From the mathematical point of view they have been studied extensively over the past few years. Here we consider a viscous incompressible three-dimensional (3D) fluid described by the Navier–Stokes equations interacting with a two-dimensional elastic plate in flexion.

One already knows that if a viscous term or the rotatory inertia are taken into account in the plate equations, there exists at least one weak solution to this problem (see [4]). From the mechanical point of view, adding a viscous term is a way to introduce dissipation in the plate model, and, from the mathematical point of view, this is a way to regularize the structure velocity. Here the dissipation coming from the fluid enables us to control the space high frequencies of the structure velocity and to pass to the limit in the coupled system as the additional viscous plate coefficient tends to zero and, thus, obtain the existence of weak solutions of the limit problem. This limit system can also be obtained as the limit of the plate-Navier–Stokes system (with a regularized initial plate velocity) as the coefficient of the rotatory inertia tends to zero. In most of the previous studies the structure velocity is quite regular because of the model or because of the presence of a regularization operator in the equations. The existing results are concerned mainly with rigid body motions [5], [10], [11], [13], [16], [17], [18], [21], [22], [24], [25] or with the motion of a structure described by a finite number of modal functions [12] or a structure with additional “viscous” terms [2], [4], [8]. Recently, a significant breakthrough has been made by Coutand and Shkoller. In [6], [7] they prove the existence, locally in time, of a unique regular solution (assuming that the data are smooth enough and satisfy suitable compatibility conditions) for the Navier–Stokes equations coupled with linearized elasticity or quasi-linear elasticity. These are the only existence results where the full 3D elasticity is

*Received by the editors August 3, 2007; accepted for publication (in revised form) March 26, 2008; published electronically July 18, 2008.

<http://www.siam.org/journals/sima/40-2/69919.html>

†REO project, INRIA Rocquencourt, B.P. 105, 78153 Le Chesnay Cedex, France (celine.grandmont@inria.fr).

considered and that don't require additional viscous terms. Nevertheless, despite these new important results, the case of fluid-plate or fluid-shell interaction problems is not, as far as we know, solved. Here, by assuming that the in-plane motions can be neglected and taking advantage of the transverse-only motions of the plate and of the fact that the plate equations enable some regularity of the boundary of the fluid domain, we prove the existence of weak solutions for a fluid-plate interaction problem. Note that the same proof does not apply for general plate or shell models. Yet, even if we consider here a rather simple structure model, this is, to our knowledge, the first existence result of weak solutions in this direction. Note, moreover, that no compatibility assumptions are required and the existence result holds as long as the plate does not touch the bottom of the fluid cavity. Finally, the same results hold for a 2D viscous flow interacting with a one-dimensional membrane.

In the first section we give the equations of the fluid-structure problem for which we derive a priori energy estimates. Next, definitions and properties of the energy spaces are detailed, and the weak formulation of the problem is given. In particular, we build suitable extensions of the fluid test functions and liftings of the structure test functions. In the second section, we state our main results, the third section being devoted to the derivation of compactness properties that enable us to pass to the limit in the equations as the "viscous" plate coefficient tends to zero (section 4).

1.1. Presentation of the problem. We assume that the fluid fills a three-dimensional cavity and interacts with a thin elastic structure, located on a part of the boundary of the fluid, the other part being rigid. For the sake of simplicity, we assume that, in the reference state, the elastic part of the fluid boundary is $\omega \times \{1\}$, where ω denotes a Lipschitz domain in \mathbb{R}^2 . In the initial state the fluid occupies the domain Ω_{η_I} :

$$\Omega_{\eta_I} = \{(x, y, z) \in \mathbb{R}^3, (x, y) \in \omega, 0 < z < 1 + \eta_I(x, y)\},$$

where η_I is a given initial displacement of the elastic part. The rigid part of $\partial\Omega_{\eta_I}$ is denoted by Γ_0 . Note that we could also have considered the case of a fluid between two elastic plates or the case where $\omega \times \{1\}$ is a part of a smooth fluid domain boundary and obtained the same kind of results. We model the deformable part of the boundary by the classical linear plate theory for transverse motions. We take its edge to be clamped. Note that the equations describing the transverse motion and the in-plane motions of a plate are decoupled. Here we ignore in-plane motions. From a mechanical point of view, this assumption is justified (at least for small enough deformations) since the membrane stiffness of a plate is much larger than its flexion stiffness. We denote by $\eta_\varepsilon(t, x, y)$ the vertical displacement with respect to the rest configuration. The subscript ε underlines the dependence of the solution with respect to the parameter $\varepsilon \geq 0$, which measures the "viscosity" of the plate (or the rotatory inertia). Then the equations describing the evolution of the transversal displacement η_ε ($\eta_\varepsilon = \eta_\varepsilon(t, x, y) \in \mathbb{R}$) are

$$(1) \quad \begin{cases} \partial_{tt}\eta_\varepsilon + \Delta^2\eta_\varepsilon + \varepsilon\Delta^2\partial_t\eta_\varepsilon = g + (T_f^\varepsilon)_3 \text{ in } \omega, \\ \eta_\varepsilon = \frac{\partial\eta_\varepsilon}{\partial n} = 0 \text{ on } \partial\omega, \\ \eta_\varepsilon(0) = \eta_I, \quad \partial_t\eta_\varepsilon(0) = \dot{\eta}_I, \end{cases}$$

where g denotes the given body force on the plate and T_f^ε the surface force exerted by the fluid on the structure. The definition of T_f^ε will be made precise later on. Instead

of the additional viscosity term, we could have added $-\varepsilon\Delta\partial_{tt}\eta_\varepsilon$, which models the inertia of rotation. The domain occupied by the fluid at time t is denoted by $\Omega_{\eta_\varepsilon}(t)$:

$$(2) \quad \Omega_{\eta_\varepsilon}(t) = \{(x, y, z) \in \mathbb{R}^3, (x, y) \in \omega, 0 < z < 1 + \eta_\varepsilon(t, x, y)\}.$$

The classical forms of the governing equations for the fluid are

$$(3) \quad \begin{cases} \partial_t \mathbf{u}_\varepsilon + (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon - \nu \Delta \mathbf{u}_\varepsilon + \nabla p_\varepsilon = \mathbf{f} & \text{in } \Omega_{\eta_\varepsilon}(t), \\ \operatorname{div} \mathbf{u}_\varepsilon = 0 & \text{in } \Omega_{\eta_\varepsilon}(t), \\ \mathbf{u}_\varepsilon(t, \cdot) = \mathbf{0} & \text{on } \Gamma_0, \\ \mathbf{u}_\varepsilon(0, \cdot) = \mathbf{u}_I & \text{in } \Omega_{\eta_I}, \end{cases}$$

where \mathbf{u}_ε denotes the fluid velocity and p_ε the pressure field. The body force \mathbf{f} and the initial velocity \mathbf{u}_I are given.

Since the fluid is viscous, it adheres to the plate, and thus the velocities coincide (in a sense to be defined) at the interface. This is written, since we assume that the plate motion is vertical:

$$(4) \quad \mathbf{u}_\varepsilon(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, \partial_t \eta_\varepsilon(t, x, y))^T, \quad (x, y) \in \omega.$$

This condition, together with the incompressibility of the fluid, leads to

$$(5) \quad \int_\omega \partial_t \eta_\varepsilon = 0.$$

Condition (5) states that the global volume of the cavity is preserved. The surface force T_f^ε exerted by the fluid on the plate can be defined by

$$(6) \quad \int_\omega T_f^\varepsilon \cdot \bar{\mathbf{v}} = \int_{\partial\Omega_{\eta_\varepsilon}(t) \setminus \Gamma_0} (-2\nu D(\mathbf{u}_\varepsilon) \cdot \mathbf{n}_t^\varepsilon + p_\varepsilon \mathbf{n}_t^\varepsilon) \cdot \mathbf{v} \quad \forall \mathbf{v},$$

where $D(\mathbf{u}_\varepsilon) = (\nabla \mathbf{u}_\varepsilon + (\nabla \mathbf{u}_\varepsilon)^T)/2$ is the strain tensor, \mathbf{n}_t^ε denotes the outer unit normal at $\partial\Omega_{\eta_\varepsilon}(t) \setminus \Gamma_0$ ($\mathbf{n}_t^\varepsilon = \frac{1}{\sqrt{1+(\partial_x \eta_\varepsilon)^2+(\partial_y \eta_\varepsilon)^2}}(-\partial_x \eta_\varepsilon, -\partial_y \eta_\varepsilon, 1)^T$), and $\bar{\mathbf{v}}(t, x, y) = \mathbf{v}(t, x, y, 1 + \eta_\varepsilon(t, x, y)) \quad \forall (x, y) \in \omega$. Note here that the pressure p_ε is not defined up to a constant but is uniquely defined. Its average value ensures the global volume conservation of the fluid cavity. This average is in fact the Lagrange multiplier associated with the compatibility condition (5). Note that if Neumann boundary conditions had been imposed on Γ_0 , then the plate displacement should not verify an additional ‘‘volume-preserving’’ constraint. As noted in [4], the third component of T_f^ε can be rewritten thanks to

$$2(D(\mathbf{u}_\varepsilon) \cdot \mathbf{n}_t^\varepsilon)_3 = (\nabla \mathbf{u}_\varepsilon \cdot \mathbf{n}_t^\varepsilon)_3.$$

This simplification comes from the fact that the displacement at the fluid-structure interface is only transverse and from the incompressibility of the fluid.

Thus $\forall \mathbf{v}$, such that $v_i(t, x, 1 + \eta_\varepsilon(t, x, y)) = 0, i = 1, 2, (t, x, y) \in (0, T) \times \omega$, we have

$$(7) \quad \int_\omega T_f^\varepsilon \cdot \hat{\mathbf{v}} = \int_{\Gamma_{\eta_\varepsilon}(t)} (-\nu \nabla \mathbf{u}_\varepsilon \cdot \mathbf{n}_t^\varepsilon + p_\varepsilon \mathbf{n}_t^\varepsilon)_3 \quad v_3.$$

Note moreover that the two first components of T_f^ε correspond to the Lagrange multiplier associated to the vertical velocity constrain that the fluid has to verify at the interface.

1.2. A priori estimates. In this subsection we recall the a priori estimates satisfied by any solution, assuming that it is smooth enough. We multiply the Navier–Stokes equations by \mathbf{u}_ε and integrate over $\Omega_{\eta_\varepsilon}(t)$, and we multiply the plate equations by $\partial_t \eta_\varepsilon$, integrate over ω , and add these two contributions. After integrations by parts and taking into account the coupling conditions (equality of the velocities (4) and the definition of T_f^ε (7)), we obtain the following energy equality:

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon|^2 + 2\nu \int_{\Omega_{\eta_\varepsilon}(t)} |\nabla(\mathbf{u}_\varepsilon)|^2 \\
 (8) \quad & + \frac{1}{2} \frac{d}{dt} \int_\omega (\partial_t \eta_\varepsilon)^2 + \frac{1}{2} \frac{d}{dt} \int_\omega (\Delta \eta_\varepsilon)^2 + \varepsilon \int_\omega (\Delta \partial_t \eta_\varepsilon)^2 \\
 & = \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot \mathbf{u}_\varepsilon + \int_\omega g \partial_t \eta_\varepsilon.
 \end{aligned}$$

Hence, using Cauchy–Schwarz and Young’s inequalities and Gronwall’s lemma:

$$\begin{aligned}
 & \frac{1}{2} \|\mathbf{u}_\varepsilon(t, \cdot)\|_{L^2(\Omega_{\eta_\varepsilon}(t))}^2 + 2\nu \int_0^t \|\nabla(\mathbf{u}_\varepsilon)(s, \cdot)\|_{L^2(\Omega_{\eta_\varepsilon}(s))}^2 ds \\
 (9) \quad & + \frac{1}{2} \|\partial_t \eta_\varepsilon(t, \cdot)\|_{L^2(\omega)}^2 + \frac{1}{2} \|\Delta \eta_\varepsilon(t, \cdot)\|_{L^2(\omega)}^2 + \varepsilon \int_0^t \|\Delta \partial_t \eta_\varepsilon(s, \cdot)\|_{L^2(\omega)}^2 ds \\
 & \leq e^t \left(\frac{1}{2} \|\mathbf{u}_I\|_{L^2(\Omega_{\eta_I})}^2 + \frac{1}{2} \|\dot{\eta}_I\|_{L^2(\omega)}^2 + \frac{1}{2} \|\Delta \eta_I\|_{L^2(\omega)}^2 \right) \\
 & + \frac{1}{2} \int_0^t \exp(t-s) \left(\|\mathbf{f}(s, \cdot)\|_{L^2(\Omega_{\eta_\varepsilon}(s))}^2 + \|g(s, \cdot)\|_{L^2(\omega)}^2 \right) ds.
 \end{aligned}$$

Thus, assuming that $\mathbf{f} \in L^2(0, T; L^2(\mathbb{R}^3))$, $g \in L^2(0, T; L^2(\omega))$, $\mathbf{u}_I \in L^2(\Omega_{\eta_I})$, $\eta_I \in H_0^2(\omega)$, and $\dot{\eta}_I \in L^2(\omega)$,

$$(10) \quad \mathbf{u}_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^\infty(0, T; L^2(\Omega_{\eta_\varepsilon}(t))),$$

$$(11) \quad \nabla \mathbf{u}_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t))),$$

and

$$(12) \quad \eta_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } W^{1,\infty}(0, T; L^2(\omega)) \cap L^\infty(0, T; H_0^2(\omega)).$$

Moreover, if $\varepsilon > 0$,

$$\partial_t \eta_\varepsilon \in L^2(0, T; H_0^2(\omega)).$$

Consequently the spaces $L^p(0, T; L^2(\Omega_\gamma(t)))$, $L^2(0, T; H^1(\Omega_\gamma(t)))$ need to be defined for γ belonging to $W^{1,\infty}(0, T; L^2(\omega)) \cap L^\infty(0, T; H_0^2(\omega))$. Note that the following continuous injection holds:

$$W^{1,\infty}(0, T; L^2(\omega)) \cap L^\infty(0, T; H^2(\omega)) \hookrightarrow C^{0,1-\theta}([0, T]; H^{2\theta}(\omega))$$

for all $0 < \theta < 1$. In particular,

$$(13) \quad W^{1,\infty}(0, T; L^2(\omega)) \cap L^\infty(0, T; H^2(\omega)) \hookrightarrow C^{0,\mu}([0, T]; C^{0,1/2-\mu}(\bar{\omega}))$$

for all $0 < \mu < 1/2$. The proof of the first injection relies on standard Hilbertian interpolation inequalities (see [20]). The other is deduced from the first one and from Sobolev injections in dimension two (see [3]). Consequently, this displacement regularity does not imply that the fluid domain boundary is Lipschitz, and we have to pay a special attention to the definitions of the functional spaces. We have also to give a sense to the equality of the velocities. Thus we are going to give some technical lemmas, definitions, and properties, most of which can be found in [4].

1.3. Preliminary definitions and properties. We now turn to the definition of some functional spaces. These definitions can be found in [4], but for the sake of completeness we recall them here. Let $T > 0$ and δ belong to $C^0([0, T] \times \bar{\omega})$ such that, for some positive M and α , $M \geq 1 + \delta(t, x, y) \geq \alpha > 0$ for all $(t, x, y) \in [0, T] \times \bar{\omega}$ and such that $\delta = 0$ on $\partial\omega$. The set $\Omega_\delta(t)$ defined by

$$\Omega_\delta(t) = \{(x, y, z) \in \mathbb{R}^3, (x, y) \in \omega, 0 < z < 1 + \delta(t, x, y)\}$$

is an open subset of \mathbb{R}^3 for every $t \in [0, T]$ which is included in $\mathcal{C}_M = \omega \times (0, M)$. Let $\widehat{\Omega}_\delta$ be the open domain of \mathbb{R}^4 defined by

$$\widehat{\Omega}_\delta = \bigcup_{t \in (0, T)} \{t\} \times \Omega_\delta(t).$$

We set $\widehat{\mathcal{C}}_M = (0, T) \times \mathcal{C}_M$. One can define in a standard way the spaces $L^p(\Omega_\delta(t))$, $H^1(\Omega_\delta(t))$, $H_0^1(\Omega_\delta(t))$, for every t , and $L^p(\widehat{\Omega}_\delta)$, $H^1(\widehat{\Omega}_\delta)$, $L^p(\widehat{\mathcal{C}}_M)$, and $H^1(\widehat{\mathcal{C}}_M)$. The space $H_{0, \Gamma_0}^1(\Omega_\delta(t))$ will denote the subspace of $H_{0, \Gamma_0}^1(\Omega_\delta(t))$ of functions of zero trace on $\Gamma_0 = \omega \times \{0\} \cup \partial\omega \times (0, 1)$. We then define:

$$\begin{aligned} L^2(0, T; H^1(\Omega_\delta(t))) &= \left\{ v \in L^2(\widehat{\Omega}_\delta), \nabla v \in L^2(\widehat{\Omega}_\delta) \right\}, \\ L^2(0, T; H_0^1(\Omega_\delta(t))) &= \overline{L^2(0, T; H^1(\Omega_\delta(t)))}^{\mathcal{D}(\widehat{\Omega}_\delta)}, \\ \mathcal{V}_\delta &= \left\{ \mathbf{v} \in C^1(\widehat{\Omega}_\delta), \operatorname{div} \mathbf{v} = 0, \mathbf{v} = \mathbf{0} \text{ on } (0, T) \times \Gamma_0 \right\}, \\ V_\delta &= \overline{\mathcal{V}_\delta}^{L^2(0, T; H^1(\Omega_\delta(t)))}, \end{aligned}$$

and

$$L^\infty(0, T; L^2(\Omega_\delta(t))) = \left\{ v \in L^2(\widehat{\Omega}_\delta), \sup_{t \in (0, T)} \operatorname{ess} \|v\|_{L^2(\Omega_\delta(t))} < +\infty \right\}.$$

Moreover we define

$$V = \left\{ \mathbf{v} \in L^2(0, T; H^1(\mathcal{C}_M)), \operatorname{div} \mathbf{v} = 0, \mathbf{v} = \mathbf{0} \text{ on } (0, T) \times (\Gamma_0 \cup \Gamma_1) \right\},$$

where $\Gamma_1 = \partial\omega \times (1, M)$.

The space V_δ can be characterized as follows:

$$V_\delta = \left\{ \mathbf{v} \in L^2(0, T; H^1(\Omega_\delta(t))), \operatorname{div} \mathbf{v} = 0, \mathbf{v} = \mathbf{0} \text{ on } (0, T) \times \Gamma_0 \right\}.$$

In the case of a Lipschitz or a star-shaped domain independent of time, this follows from standard arguments (see [26] or [15]). In our case it can be proved by using the fact that the domain $\Omega_\delta(t)$ is locally a subgraph. This property will be extensively used in all that follows.

Next we recall various lemmas that explain how the trace on $\partial\Omega_\delta(t) \setminus \Gamma_0$ makes sense, define extension and lifting operators, and explore some properties of the spaces defined above. We omit the proofs whenever they can be found in [4]. Note that these results take advantage of the fact that the fluid domain is a subgraph because the displacement of the interface is only transverse. Let us consider the linear mapping $\gamma_{\delta(t)}: v \mapsto v(x, y, 1 + \delta(t, x, y))$ defined for $v \in C^0(\overline{\Omega_\delta(t)})$.

LEMMA 1. *For every $t \in [0, T]$, the mapping $\gamma_{\delta(t)}$ from $C^1(\overline{\mathcal{C}_M})$ (resp., $C^1(\overline{\Omega_\delta(t)})$) in $C^0(\overline{\omega})$ can be extended by continuity to a linear mapping from $H^1(\mathcal{C}_M)$ (resp., $H^1(\Omega_\delta(t))$) into $L^2(\omega)$.*

COROLLARY 1. *If $v \in L^2(0, T; H^1(\Omega_\delta(t)))$, then $\gamma_{\delta(t)}(v) \in L^2(0, T; L^2(\omega))$.*

Thus, the trace of $v(x, y, 1 + \delta(t, x, y))$ on ω makes sense at least in $L^2(\omega)$. The following lemma makes precise the regularity of $\gamma_{\delta(t)}(v)$ when assuming moreover that δ belongs to $L^\infty(0, T; H_0^2(\omega))$. This additional regularity will play a crucial role in our asymptotic study and will enable us to control the space high frequencies of the structure velocity.

LEMMA 2. *Assuming that $\delta \in C^0([0, T]; C^0(\overline{\omega})) \cap L^\infty(0, T; H_0^2(\omega))$, then, for any $v \in H^1(\Omega_\delta(t))$, $\gamma_{\delta(t)}(v) \in W^{1-1/p, p}(\omega)$, $\forall 1 < p < 2$ and for $\frac{3}{2} \leq p < 2$, $\gamma_{\delta(t)}(v) \in H^{\frac{3p-2}{p}}(\omega)$, for a.e. t . Moreover, if $v \in L^2(0, T; H^1(\Omega_\delta(t)))$, then $\gamma_{\delta(t)}(v) \in L^2(0, T; W^{1-1/p, p}(\omega)) \forall 1 < p < 2$ and $\gamma_{\delta(t)}(v) \in L^2(0, T; H^{\frac{3p-2}{p}}(\omega)) \forall \frac{3}{2} \leq p < 2$.*

Proof. Let us define an auxiliary function by

$$w(t, x, y, s) = v(x, y, 1 + \delta(t, x, y) - s), \quad 0 \leq s \leq \alpha.$$

It is clear that w belongs to $L^2(\mathcal{C}_\alpha) \forall t$. Moreover

$$\begin{aligned} \nabla w(t, x, y, s) &= \begin{pmatrix} \partial_x v(x, y, 1 + \delta(t, x, y) - s) + \partial_x \delta(t, x, y) \partial_z v(x, y, 1 + \delta(t, x, y) - s) \\ \partial_y v(x, y, 1 + \delta(t, x, y) - s) + \partial_y \delta(t, x, y) \partial_z v(x, y, 1 + \delta(t, x, y) - s) \\ -\partial_z v(x, y, 1 + \delta(t, x, y) - s) \end{pmatrix}. \end{aligned}$$

Since $H^1(\omega)$ is continuously imbedded in $L^q(\omega) \forall q < \infty$ and $\partial_i v(x, y, 1 + \delta(t, x, y) - s) \in L^2(\mathcal{C}_\alpha) \forall t$, we deduce that $\nabla w \in L^p(\mathcal{C}_\alpha)$ for all $1 < p < 2$, for a.e. t . Thus $w \in W^{1, p}(\mathcal{C}_\alpha)$, and $\gamma_{\delta(t)}(v) = w|_{s=0} \in W^{1-1/p, p}(\omega) \forall 1 < p < 2$ for a.e. t . Moreover, for a.e. t ,

$$\|\gamma_{\delta(t)}(v)\|_{W^{1-1/p, p}(\omega)} \leq C(\|\delta\|_{C^0([0, T]; C^0(\overline{\omega})) \cap L^\infty(0, T; H_0^2(\omega))}, \|v\|_{H^1(\Omega_\delta(t))}).$$

Furthermore, from Sobolev injections (see [1, Thm. 7.58, p. 218]), we deduce that, for $2 > p \geq \frac{3}{2}$, $\gamma_{\delta(t)}(v) \in H^{\frac{3p-2}{p}}(\omega)$ and

$$(14) \quad \|\gamma_{\delta(t)}(v)\|_{H^{\frac{3p-2}{p}}(\omega)} \leq C(\|\delta\|_{C^0([0, T]; C^0(\overline{\omega})) \cap L^\infty(0, T; H_0^2(\omega))}, \|v\|_{H^1(\Omega_\delta(t))}). \quad \square$$

Now we are going to give a characterization of $H_0^1(\Omega_\delta(t))$ with the help of the mapping $\gamma_{\delta(t)}$. An additional assumption on the boundary displacement δ is needed: δ is assumed to belong to $C^0([0, T]; H^1(\omega))$ (this is not an optimal assumption). \square

LEMMA 3. *Assuming that $\delta \in C^0([0, T]; C^0(\overline{\omega})) \cap H^1(\omega)$, then*

$$H_0^1(\Omega_\delta(t)) = \{v \in H_{0, \Gamma_0}^1(\Omega_\delta(t)), \gamma_{\delta(t)}(v) = 0\}.$$

COROLLARY 2. *If $v \in L^2(0, T; H_{0, \Gamma_0}^1(\Omega_\delta(t)))$ and $\gamma_{\delta(t)}(v) = 0$, for a.e. t , then $v \in L^2(0, T; H_0^1(\Omega_\delta(t)))$, and the converse is also true.*

We now state a lemma that enables us to extend a function $\mathbf{v} \in V_\delta$ such that $\gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T$, $b \in L^2(0, T; H^1(\omega))$, the extension belonging to V .

LEMMA 4. *We assume that $\delta \in C^0([0, T]; C^0(\bar{\omega}) \cap H_0^1(\omega))$. Let $\mathbf{v} \in V_\delta$ such that, for a.e. t , $\gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T$, $b \in L^2(0, T; H_0^1(\omega))$. The function defined by*

$$(15) \quad \bar{\mathbf{v}} = \begin{cases} \mathbf{v} & \text{in } \widehat{\Omega}_\delta \\ (0, 0, b)^T & \text{in } \widehat{\mathcal{C}}_M \setminus \widehat{\Omega}_\delta \end{cases}$$

belongs to V , and

$$\|\bar{\mathbf{v}}\|_V \leq C(\|\mathbf{v}\|_{V_\delta} + \|b\|_{L^2(0, T; H^1(\omega))}),$$

where C depends only on M .

Remark 1. If $\mathbf{v} \in L^\infty(0, T; L^2(\Omega_\delta(t)))$ and $b \in L^\infty(0, T; L^2(\omega))$, then $\bar{\mathbf{v}} \in L^\infty(0, T; L^2(\mathcal{C}_M))$.

Next we build different lifting operators.

LEMMA 5. *For every $\phi \in H_0^1(\omega)$ there exists $w \in H_{0, \Gamma_0}^1(\Omega_\delta(t))$ such that*

$$\gamma_{\delta(t)}(w) = \phi \quad \text{and} \quad \|w\|_{H^1(\Omega_\delta(t))} \leq C_\alpha \|\phi\|_{H^1(\omega)}.$$

For every $b \in H_0^1(\omega)$ such that $\int_\omega b = 0$ there exists \mathbf{v} such that $\gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T$, $\text{div}(\mathbf{v}) = 0$, and

$$\|\mathbf{v}\|_{H_{0, \Gamma_0}^1(\Omega_\delta(t))} \leq C_\alpha \|b\|_{H_0^1(\omega)}.$$

Proof. Indeed

$$w = \begin{cases} \phi & \text{in } \Omega_\delta(t) \setminus \mathcal{C}_\alpha, \\ \mathcal{R}(\frac{z}{\alpha}\phi) & \text{in } \mathcal{C}_\alpha, \end{cases}$$

where \mathcal{R} is a continuous lifting operator from $H^{1/2}(\partial\mathcal{C}_\alpha)$ into $H^1(\mathcal{C}_\alpha)$ and $\mathcal{C}_\alpha = \omega \times (0, \alpha)$ verifies the desired properties. Moreover, if we consider $b \in H_0^1(\omega)$ such that $\int_\omega b = 0$, then $\tilde{\mathcal{R}}$ is a continuous lifting operator from $H^{1/2}(\partial\mathcal{C}_\alpha)$ into $H^1(\mathcal{C}_\alpha)$ such that $\text{div}(\tilde{\mathcal{R}}\mathbf{v}) = 0$,

$$(16) \quad \mathbf{v} = \begin{cases} (0, 0, b)^T & \text{in } \Omega_\delta(t) \setminus \mathcal{C}_\alpha, \\ \tilde{\mathcal{R}}(0, 0, \frac{z}{\alpha}b)^T & \text{in } \mathcal{C}_\alpha, \end{cases}$$

is divergence-free and belongs to $H_{0, \Gamma_0}^1(\Omega_\delta(t))$. Furthermore, we have for a.e. t

$$\|\mathbf{v}\|_{H_{0, \Gamma_0}^1(\Omega_\delta(t))} \leq C_\alpha \|b\|_{H_0^1(\omega)}.$$

Consequently (16) defines a continuous linear lifting from $\{b \in H_0^1(\omega), \text{ such that } (\text{s.t.}) \int_\omega b = 0\}$ into $\{\mathbf{v} \in H_{0, \Gamma_0}^1(\Omega_\delta(t)), \text{ s.t. } \text{div}(\tilde{\mathcal{R}}\mathbf{v}) = 0\}$. \square

Remark 2. Thanks to the previous lemma the space

$$\{\mathbf{v} \in V_\delta, \gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T \text{ for a.e. } t, b \in L^2(0, T; H_0^1(\omega))\}$$

is equal to the sum of the two following spaces:

$$\overline{\{\mathbf{v} \in \mathcal{D}(\widehat{\Omega}_\delta), \text{div } \mathbf{v} = 0\}}^{L^2(0, T; H^1(\Omega_\delta(t)))}$$

and

$$\left\{ \mathbf{v}, \mathbf{v} = \begin{cases} (0, 0, b)^T & \text{in } \Omega_\delta(t) \setminus \mathcal{C}_\alpha \\ \tilde{\mathcal{R}}(0, 0, \frac{z}{\alpha}b)^T & \text{in } \mathcal{C}_\alpha \end{cases} \text{ for a.e. } t, b \in L^2(0, T; H_0^1(\omega)), \int_\omega b = 0 \right\}.$$

We also need to build a “lifting” operator of $(0, 0, b)^T$ for any b that belongs only to $H^s(\omega), 0 \leq s < \frac{1}{2}$ since the structure velocity $\partial_t \eta_\varepsilon$ will be bounded, uniformly in ε , only in $L^2(0, T; \dot{H}^s(\omega)) \forall 0 \leq s < \frac{1}{2}$ and not in $L^2(0, T; H_0^1(\omega))$ (see Lemma 2).

LEMMA 6. For all $b \in L^2(0, T; \dot{H}^s(\omega)), 0 \leq s < \frac{1}{2}$ such that $\int_\omega b = 0$, there exists a lifting operator \mathcal{R}_α satisfying $\gamma_{\delta(t)}(\mathcal{R}_\alpha(b)) = (0, 0, b)^T$ and $\text{div}(\mathcal{R}_\alpha(b)) = 0$ and for a.e. t

$$(17) \quad \|\mathcal{R}_\alpha(b)\|_{H^s(\mathcal{C}_M)} \leq C\|b\|_{H^s(\omega)} \quad \forall 0 \leq s < \frac{1}{2}.$$

Proof. A rather naive construction will satisfy the desired properties. Let us define \mathcal{R}_α by

$$(18) \quad \mathcal{R}_\alpha(b) = \begin{cases} (0, 0, b)^T & \text{for } z \geq \alpha \\ (0, 0, \frac{z}{\alpha}b)^T + \mathbf{w}_\alpha & \text{in } \mathcal{C}_\alpha \end{cases} \text{ for a.e. } t,$$

where \mathbf{w}_α is such that $\text{div}(\mathbf{w}_\alpha) = b$ and $\mathbf{w}_\alpha \in H_0^1(\mathcal{C}_\alpha), \|\mathbf{w}_\alpha\|_{H_0^1(\mathcal{C}_\alpha)} \leq C\|b\|_{L^2(\omega)}$, for a.e. t . Such a function exists (see, for instance, [14]). It is easy to see that $\mathcal{R}_\alpha(b)$ is divergence-free. Moreover, \mathcal{R}_α is linear continuous from $L^2(0, T, L^2(\omega))$ (resp., $L^2(0, T, H_0^1(\omega))$) into $L^2(0, T, L^2(\mathcal{C}_M))$ (resp., $L^2(0, T, H_{0,\Gamma_0}^1(\mathcal{C}_M))$), and thus, by interpolation (see [19]), \mathcal{R}_α is linear continuous from $L^2(0, T, H^s(\omega))$ into $L^2(0, T, H^s(\mathcal{C}_M)) \forall 0 \leq s < \frac{1}{2}$. Consequently, (17) holds true. \square

Remark 3. The trace $\gamma_{\delta(t)}(\mathcal{R}_\alpha(b))$ makes sense for any $b \in L^2(\omega)$ since $\mathcal{R}_\alpha(b)$ is regular enough with respect to z .

We end this section by mentioning that Korn’s and Poincaré’s inequalities hold in the considered spaces.

LEMMA 7. For all \mathbf{u} and \mathbf{v} in

$$\{ \mathbf{v} \in V_\delta, \exists b \in L^2(0, T; H_0^1(\omega)), \gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T \text{ for a.e. } t \},$$

we have

$$(19) \quad 2 \int_{\Omega_\delta(t)} D(\mathbf{u}) : D(\mathbf{v}) = \int_{\Omega_\delta(t)} \nabla \mathbf{u} : \nabla \mathbf{v} \text{ for a.e. } t,$$

and consequently the following Korn’s “equality” holds:

$$(20) \quad \sqrt{2}\|D(\mathbf{u})\|_{L^2(\Omega_\delta(t))} = \|\nabla \mathbf{u}\|_{L^2(\Omega_\delta(t))} \text{ for a.e. } t.$$

LEMMA 8. Let $v \in H_{0,\Gamma_0}^1(\Omega_\delta(t))$, and then

$$\|v\|_{L^2(\Omega_\delta(t))} \leq M\|\nabla v\|_{L^2(\Omega_\delta(t))}.$$

1.4. Weak formulation. Let $\eta_I \in H_0^2(\omega), (\mathbf{u}_I, \dot{\eta}_I) \in L^2(\Omega_{\eta_I}) \times L^2(\omega)$, such that

$$(21) \quad \begin{aligned} \min_{\bar{\omega}}(1 + \eta_I) &> 0, \\ \text{div } \mathbf{u}_I &= 0 \text{ in } \Omega_{\eta_I}, \\ \mathbf{u}_I \cdot \mathbf{n} &= 0 \text{ on } \Gamma_0, \\ \mathbf{u}_I(x, y, 1 + \eta_I(x, y)) \cdot \mathbf{n}_0 &= (0, 0, \dot{\eta}_I(x, y))^T \cdot \mathbf{n}_0 \text{ on } \omega, \\ \int_\omega \dot{\eta}_I(x, y) &= 0, \end{aligned}$$

where \mathbf{n}_0 denotes the unit normal to the initial position of the plate. We refer to [4], where one proves that the normal trace $\mathbf{u}_I(x, y, 1 + \eta_I(x, y)) \cdot \mathbf{n}_0, (x, y) \in \omega$ makes sense for $\mathbf{u}_I \in L^2(\Omega_{\eta_I})$, with $\eta_I \in H_0^2(\omega)$. We shall say that $(\mathbf{u}_\varepsilon, \eta_\varepsilon)$ is a weak solution of the considered model on $(0, T)$ if it satisfies the following problem that will be denoted $(\mathcal{P}_\varepsilon)$:

- $\mathbf{u}_\varepsilon \in V_{\eta_\varepsilon} \cap L^\infty(0, T; L^2(\Omega_{\eta_\varepsilon}(t))), \eta_\varepsilon \in W^{1,\infty}(0, T; L^2(\omega)) \cap L^\infty(0, T; H_0^2(\omega)),$
- for $\varepsilon > 0, \partial_t \eta_\varepsilon \in L^2(0, T; H_0^2(\omega)),$
- $\mathbf{u}_\varepsilon(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, \partial_t \eta_\varepsilon(t, x, y))^T$ for a.e. $(t, x, y) \in (0, T) \times \omega,$
- for all $(\mathbf{f}, b) \in (V_{\eta_\varepsilon} \cap H^1(\widehat{\Omega}_{\eta_\varepsilon})) \times (L^2(0, T; H_0^2(\omega)) \times H^1(0, T; L^2(\omega)))$ such that $(\mathbf{f}, b)(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T,$ for a.e. $(t, x, y) \in (0, T) \times \omega,$ we have for a.e. t

$$\begin{aligned}
 (22) \quad & \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon(t) \cdot (\mathbf{f}) - \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} \mathbf{u}_\varepsilon \cdot \partial_t \mathbf{f} + \nu \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} \nabla \mathbf{u}_\varepsilon : \nabla \mathbf{f} \\
 & + \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon \cdot \mathbf{f} - \int_0^t \int_\omega (\partial_t \eta_\varepsilon)^2 b + \int_\omega \partial_t \eta_\varepsilon(t) b(t) \\
 & - \int_0^t \int_\omega \partial_t \eta_\varepsilon \partial_t b + \int_0^t \int_\omega \Delta \eta_\varepsilon \Delta b + \varepsilon \int_0^t \int_\omega \Delta \partial_t \eta_\varepsilon \Delta b \\
 & = \int_0^t \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot \mathbf{u}_I + \int_0^t \int_\omega g b + \int_{\Omega_{\eta_I}} \mathbf{u}_I \cdot (0) + \int_\omega \dot{\eta}_I b(0).
 \end{aligned}$$

In what follows, a solution of (\mathcal{P}_0) will be denoted by (\mathbf{u}, η) (instead of (\mathbf{u}_0, η_0)).

Remark 4. The test functions depends on the solution and thus, for $\varepsilon > 0,$ on $\varepsilon.$

Remark 5. The trace at time $t = 0$ of $\mathbf{f} \in H^1(\widehat{\Omega}_{\eta_\varepsilon})$ such that $(\mathbf{f}, b)(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T$ for a.e. $(t, x, y) \in (0, T) \times \omega,$ with $b \in H^1((0, T) \times \omega),$ is well-defined and makes sense at least in $L^2(\Omega_{\eta_I}).$ Indeed we can prove by density arguments ($C^1(\widehat{\Omega}_{\eta_\varepsilon})$ is dense in $H^1(\widehat{\Omega}_{\eta_\varepsilon})$ since the domain is a continuous subgraph; see, for instance, [1, Thm. 2, p. 54]) that

$$\int_{\Omega_{\eta_\varepsilon}(0)} |\mathbf{f}(0)|^2 = 2 \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \psi \partial_t \mathbf{f} + \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\partial_t \psi|^2 + \int_0^T \int_\omega |b|^2 \psi \partial_t \eta_\varepsilon,$$

where ψ belongs to $\mathcal{D}([0, T])$ and satisfies $\psi(0) = 1.$ The right-hand side of this equality makes sense for any $\mathbf{f} \in H^1(\widehat{\Omega}_{\eta_\varepsilon})$ such that $(\mathbf{f}, b)(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T$ for a.e. $(t, x, y) \in (0, T) \times \omega,$ with $b \in H^1((0, T) \times \omega),$ since $\partial_t \eta_\varepsilon \in L^\infty(0, T; L^2(\omega))$ and $b \in L^4((0, T) \times \omega).$

2. Main result. We make the following hypotheses on the data (bulk forces and initial data):

$$(23) \quad \begin{aligned}
 & \mathbf{f} \in L_{loc}^2((0, +\infty) \times \mathbb{R}^2), \quad g \in L_{loc}^2((0, +\infty) \times \omega), \\
 & \mathbf{u}_I \in L^2(\Omega_{\eta_I}), \quad \dot{\eta}_I \in L^2(\omega), \quad \eta_I \in H_0^2(\omega),
 \end{aligned}$$

and we assume moreover that conditions (21) are satisfied.

First we recall that, for $\varepsilon > 0,$ there exists at least one weak solution provided that the plate does not touch the bottom of the fluid cavity, in other words, as long as $\min_{(x,y) \in \overline{\omega}} (1 + \eta_\varepsilon(t, x, y)) > 0.$ The proof of the following theorem can be found in [4].

THEOREM 1. *Let ε be strictly positive. Under assumptions (21) and (23) and if $\min_{(x,y) \in \overline{\omega}} 1 + \eta_I(x, y) > 0,$ there exist $T_\varepsilon^* \in (0, +\infty]$ and a weak solution $(\mathbf{u}_\varepsilon, \eta_\varepsilon)$ of*

$(\mathcal{P}_\varepsilon)$ on $[0, T]$, $T < T_\varepsilon^*$. This solution satisfies an energy estimate for all $T < T_\varepsilon^*$:

$$\begin{aligned} & \|\mathbf{u}_\varepsilon\|_{L^\infty(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} + \|\nabla \mathbf{u}_\varepsilon\|_{L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} \\ & + \|\partial_t \eta_\varepsilon\|_{L^\infty(0, T; L^2(\omega))} + \|\Delta \eta_\varepsilon\|_{L^\infty(0, T; L^2(\omega))} + \sqrt{\varepsilon} \|\Delta \partial_t \eta_\varepsilon\|_{L^2(0, T; L^2(\omega))} \\ & \leq C(T, \|\mathbf{u}_I\|_{L^2(\Omega_{\eta_I})}, \|\mathbf{f}\|_{L^2((0, T) \times \mathbb{R}^2)}, \|g\|_{L^2((0, T) \times \omega)}, \|\eta_I\|_{H_0^2(\omega)}, \|\dot{\eta}_I\|_{L^2(\omega)}), \end{aligned}$$

where $C > 0$ is nondecreasing with respect to its arguments. Moreover, we have the following alternative:

- either $T_\varepsilon^* = +\infty$
- or $\lim_{t \rightarrow T_\varepsilon^*} \min_{\bar{\omega}} (1 + \eta_\varepsilon) = 0$.

Since \mathbf{u}_ε is bounded in $L^2(0, T; H^1(\Omega_{\eta_\varepsilon}(t)))$ uniformly in ε and thanks to the regularity of the moving elastic boundary, the trace $\gamma_{\eta_\varepsilon(t)}(\mathbf{u}_\varepsilon)$ is bounded in $L^2(0, T; W^{1-1/p, p}(\omega)) \forall 1 < p < 2$ and in $L^2(0, T; H^s(\omega)) \forall 0 \leq s < \frac{1}{2}$ uniformly in ε (see Lemma 2 and (14)). Thus, thanks to the equality of the velocities (4),

$$(24) \quad \partial_t \eta_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^2(0, T; W^{1-1/p, p}(\omega)) \forall 1 < p < 2,$$

and

$$(25) \quad \partial_t \eta_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^2(0, T; H^s(\omega)) \forall 0 \leq s < \frac{1}{2}.$$

This will be one of the key arguments for the derivation of the compactness properties of the sequence $(\mathbf{u}_\varepsilon, \eta_\varepsilon)$. In particular thanks to (25) we can control, uniformly in ε , the space high frequencies of $\partial_t \eta_\varepsilon$ in $L^2(0, T; L^2(\omega))$.

In all that follows, the characteristic function of $\widehat{\Omega}_{\eta_\varepsilon}$ will be denoted by ρ_ε , and $\rho_\varepsilon v$ will denote the function equal to v in $\widehat{\Omega}_{\eta_\varepsilon}$ and zero elsewhere. Moreover, we choose M large enough such that, for all $\varepsilon > 0$, $1 + \eta_\varepsilon(t, x, y) \leq M \forall (t, x, y) \in [0, T] \times \bar{\omega}$, which is made possible by (12) and (13). We set

$$\bar{\mathbf{u}}_\varepsilon = \begin{cases} \mathbf{u}_\varepsilon & \text{in } \widehat{\Omega}_{\eta_\varepsilon}, \\ (0, 0, \partial_t \eta_\varepsilon)^T & \text{in } \widehat{C}_M \setminus \widehat{\Omega}_{\eta_\varepsilon}. \end{cases}$$

In all that follows, for any \mathbf{v} in $L^2(0, T; H_{0, \Gamma_0}^1(\Omega_{\eta_\varepsilon}(t)))$ such that $\gamma_{\delta(t)}(\mathbf{v}) = (0, 0, b)^T$, $b \in L^2(0, T; H_0^1(\omega))$, $\bar{\mathbf{v}}$ is defined by (15).

The main results of the present paper are the following.

PROPOSITION 1. *The sequence $(T_\varepsilon^*)_{\varepsilon > 0}$ is bounded from below away from zero, and the following convergences (up to the extractions of subsequences) hold as ε goes to zero:*

$$(26) \quad \begin{array}{llll} \eta_\varepsilon & \rightarrow & \eta & \text{strongly in } C^0([0, T]; C^0(\bar{\omega})), \\ \eta_\varepsilon & \rightharpoonup & \eta & \text{weakly in } L^2(0, T; H_0^2(\omega)), \\ \partial_t \eta_\varepsilon & \rightarrow & \partial_t \eta & \text{strongly in } L^2(0, T; L^2(\omega)), \\ \rho_\varepsilon \mathbf{u}_\varepsilon & \rightarrow & \rho \mathbf{u} & \text{strongly in } L^2(0, T; L^2(C_M)), \\ \bar{\mathbf{u}}_\varepsilon & \rightarrow & \bar{\mathbf{u}} & \text{strongly in } L^2(0, T; L^2(C_M)), \\ \rho_\varepsilon \nabla \mathbf{u}_\varepsilon & \rightharpoonup & \rho \nabla \mathbf{u} & \text{weakly in } L^2(0, T; L^2(C_M)), \end{array}$$

where $T > 0$ is a lower bound of T_ε^* independent of ε . Moreover

$$\gamma_{\eta(t)}(\mathbf{u}) = (0, 0, \partial_t \eta)^T.$$

This enables us to pass to the limit in (22) as ε tends to zero and thus obtain the following theorem.

THEOREM 2. *Under assumptions (21) and (23) and if $\min_{(x,y) \in \bar{\omega}} 1 + \eta_I(x, y) > 0$, there exist $T^* \in (0, +\infty]$ and a weak solution (\mathbf{u}, η) of (\mathcal{P}_0) on $[0, T]$, $T < T^*$. This solution satisfies energy estimates for all $T < T^*$:*

$$(27) \quad \begin{aligned} & \|\mathbf{u}\|_{L^\infty(0,T;L^2(\Omega_\eta(t)))} + \|\nabla \mathbf{u}\|_{L^2(0,T;L^2(\Omega_\eta(t)))} \\ & \quad + \|\partial_t \eta\|_{L^\infty(0,T;L^2(\omega))} + \|\eta\|_{L^\infty(0,T;H_0^2\omega)} \\ & \leq C(T, \|\mathbf{u}_I\|_{L^2(\Omega_{\eta_I})}, \|\mathbf{f}\|_{L^2((0,T) \times \mathbb{R}^2)}, \|g\|_{L^2((0,T) \times \omega)}, \|\eta_I\|_{H_0^2(\omega)}, \|\dot{\eta}_I\|_{L^2(\omega)}), \end{aligned}$$

where $C > 0$ is nondecreasing with respect to its arguments. The following alternatives are satisfied:

- either $T^* = +\infty$
- or $\lim_{t \rightarrow T^*} \min_{\bar{\omega}} (1 + \eta) = 0$.

3. Proof of Proposition 1. First we prove that Proposition 1 holds true. We have to verify that T_ε^* is bounded from below independently of ε and obtain compactness properties on $(\mathbf{u}_\varepsilon, \partial_t \eta_\varepsilon)$ in order to prove the desired strong convergences that will enable us to pass to the limit in $(\mathcal{P}_\varepsilon)$ as ε goes to zero.

Lower bound of T_ε^ .* For $\varepsilon > 0$ the solution η_ε is bounded uniformly in ε in $L^\infty(0, T; H_0^2(\omega)) \cap W^{1,\infty}(0, T; L^2(\omega))$ for all $T < T_\varepsilon^*$. Thus from (13), η_ε is bounded uniformly in ε in $C^{0,\mu}([0, T]; C^0(\bar{\omega}))$, $0 < \mu < \frac{1}{2}$. Consequently

$$1 + \eta_\varepsilon(t, x, y) \geq (1 + \eta_I(x, y)) - Ct^\mu \quad \forall (t, x, y) \in [0, T_\varepsilon^*) \times \bar{\omega},$$

where C depends only on the data of the problem. Thus T_ε^* is bounded from below by a time independent of ε . Let T be such that

$$\forall \varepsilon > 0, \quad \min_{(t,x,y) \in [0,T] \times \bar{\omega}} (1 + \eta_\varepsilon(t, x, y)) \geq \alpha > 0,$$

where α is chosen such that $\min_{(x,y) \in \bar{\omega}} (1 + \eta_I(x, y)) \geq 2\alpha > 0$.

Convergences of the sequence $(\mathbf{u}_\varepsilon, \eta_\varepsilon)$. From (12) and the compact injection (13) we deduce easily the first two convergences announced in Proposition 1. Next we prove strong convergence properties for the fluid and the structure velocities. The solution $(\mathbf{u}_\varepsilon, \eta_\varepsilon)_{\varepsilon > 0}$ we build verifies estimate (9) and (25). Furthermore, since \mathbf{u}_ε is bounded uniformly in ε in $L^2(0, T; H^1(\Omega_{\eta_\varepsilon}(t)))$, it is easy to verify that \mathbf{w}_ε , defined by $\mathbf{w}_\varepsilon(t, x, y, z) = \mathbf{u}_\varepsilon(t, x, y, z(1 + \eta_\varepsilon(t, x, y)))$, is bounded uniformly in ε in $L^2(0, T; W^{1,p}(\mathcal{C}_1)) \forall 1 < p < 2$. This implies, thanks to Sobolev injections (see [1, Thm. 7.58, p. 218]) that \mathbf{w}_ε is uniformly bounded in $L^2(0, T; H^\theta(\mathcal{C}_1))$ for any $\theta < 1$. Moreover $\partial_t \eta_\varepsilon$ is uniformly bounded in $L^2(0, T; H^s(\omega)) \forall 0 \leq s < \frac{1}{2}$. Consequently, $\mathbf{w}_\varepsilon - \mathcal{R}_\alpha(\partial_t \eta_\varepsilon)$ is uniformly bounded in $L^2(0, T; H^s(\mathcal{C}_1))$ for any $0 \leq s < \frac{1}{2}$ and its extension by zero for $z \geq 1$ is uniformly bounded in $L^2(0, T; H^s(\mathcal{C}_L))$ for any $s < \frac{1}{2}$, $L \geq 1$ (see [19]). Thus if we extend \mathbf{w}_ε by $(0, 0, \partial_t \eta_\varepsilon)^T$ for $z \geq 1$, this extension is uniformly bounded in $L^2(0, T; H^s(\mathcal{C}_L)) \forall 0 \leq s < \frac{1}{2}, L \geq 1$. Consequently, since the change of variables $\phi_\varepsilon(t, x, y, z) = (x, y, z(1 + \eta_\varepsilon(t, x, y)))^T$ is in $L^\infty(0, T; C^{0,\beta}(\mathcal{C}_L)) \forall \beta < 1$ as well as its inverse, it is easy to verify that

$$(28) \quad \bar{\mathbf{u}}_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^2(0, T; H^{s'}(\mathcal{C}_M)) \forall 0 \leq s' < s, \forall s < \frac{1}{2}.$$

Moreover thanks to the Sobolev injections, \mathbf{w}_ε is bounded uniformly in ε in $L^2(0, T; L^q(\mathcal{C}_1)) \forall q < 6$ and $\partial_t \eta_\varepsilon$ is uniformly bounded in $L^2(0, T; L^r(\omega)) \forall r < 4$; thus

$$(29) \quad \bar{\mathbf{u}}_\varepsilon \text{ is bounded, uniformly in } \varepsilon, \text{ in } L^2(0, T; L^r(\mathcal{C}_M)) \forall r < 4.$$

Nevertheless, these bounds are not sufficient to obtain the desired strong convergences. We are going to use the following lemma that characterizes the compact sets of $L^p(0, T; X)$, where X is a Banach space (see [23]).

LEMMA 9. *Let X be a Banach space and $F \hookrightarrow L^q(0, T; X)$, with $1 \leq q < \infty$. Then F is a relatively compact set of $L^q(0, T; X)$ if and only if*

- (i) $\{\int_{t_1}^{t_2} f(t)dt, f \in F\}$ is relatively compact in $X \forall t_1 < t_2 < T$;
- (ii) $\|f(t+h) - f(t)\|_{L^q(0, T; X)} \rightarrow 0$ as h goes to zero, uniformly with respect to f in F .

We are going to apply Lemma 9 to $F = (\bar{\mathbf{u}}_\varepsilon, \partial_t \eta_\varepsilon)_{\varepsilon > 0}$, $q = 2$, and $X = L^2(\mathcal{C}_M) \times L^2(\omega)$. The first point (i) is clearly satisfied thanks to (25) and (28), and we have to verify the second point. Given any $h > 0$, we denote that $g^-(t, \cdot) = g(t - h, \cdot)$ and $g^+(t, \cdot) = g(t + h, \cdot)$. The assertion (ii) is a consequence of the following lemma.

LEMMA 10. *Let $T > 0$ such that $\min_{[0, T] \times \bar{\omega}}(1 + \eta_\varepsilon) \geq \alpha > 0$. We have $\forall \beta > 0, \exists h_0 > 0$, s.t. $\forall \varepsilon > 0, \forall h \leq h_0$*

$$(30) \quad \int_0^T \int_{\mathcal{C}_M} \rho_\varepsilon |\bar{\mathbf{u}}_\varepsilon - \bar{\mathbf{u}}_\varepsilon^-|^2 + \int_0^T \int_\omega (\partial_t \eta_\varepsilon - \partial_t \eta_\varepsilon^-)^2 \leq \beta$$

and

$$(31) \quad \int_0^T \int_{\mathcal{C}_M} |\rho_\varepsilon \bar{\mathbf{u}}_\varepsilon - \rho_\varepsilon^- \bar{\mathbf{u}}_\varepsilon^-|^2 \leq \beta,$$

with η_ε extended by η_I for $t < 0$ and $\bar{\mathbf{u}}_\varepsilon$ and $\partial_t \eta_\varepsilon$ extended by 0 for $t < 0$ and where ρ_ε denotes the characteristic function of $\hat{\Omega}_{\eta_\varepsilon}$.

Proof. We first show that (30) implies (31). Indeed:

$$|\rho_\varepsilon \bar{\mathbf{u}}_\varepsilon - \rho_\varepsilon^- \bar{\mathbf{u}}_\varepsilon^-|^2 \leq C (\rho_\varepsilon |\bar{\mathbf{u}}_\varepsilon - \bar{\mathbf{u}}_\varepsilon^-|^2 + |\rho_\varepsilon - \rho_\varepsilon^-| |\bar{\mathbf{u}}_\varepsilon^-|^2).$$

The estimate of the first contribution comes from (30). For the second contribution we use the fact that $\bar{\mathbf{u}}_\varepsilon$ is bounded uniformly in ε in $L^2(0, T; L^3(\mathcal{C}_M))$ (see (29)):

$$\left| \int_0^T \int_{\mathcal{C}_M} |\rho_\varepsilon - \rho_\varepsilon^-| |\bar{\mathbf{u}}_\varepsilon|^2 \right| \leq \int_0^T \|\rho_\varepsilon - \rho_\varepsilon^-\|_{L^3(\mathcal{C}_M)} \|\bar{\mathbf{u}}_\varepsilon\|_{L^3(\mathcal{C}_M)}^2 \leq C \int_0^T \|\rho_\varepsilon - \rho_\varepsilon^-\|_{L^3(\mathcal{C}_M)}.$$

Remember now that $\partial_t \eta_\varepsilon$ is bounded in $L^\infty(0, T; L^2(\omega))$ uniformly in ε , and thus

$$\begin{aligned} \int_{\mathcal{C}_M} |\rho_\varepsilon - \rho_\varepsilon^-|^3 &= \int_\omega |\eta_\varepsilon - (\eta_\varepsilon)^-| \\ &= \int_\omega \left| \int_{t-h}^t \partial_t \eta_\varepsilon(s) ds \right| \\ &\leq \int_\omega \int_{t-h}^t |\partial_t \eta_\varepsilon(s)| ds \\ &\leq Ch. \end{aligned}$$

It leads to

$$(32) \quad \left| \int_0^T \int_{\mathcal{C}_M} (\rho_\varepsilon - \rho_\varepsilon^-) |\bar{\mathbf{u}}_\varepsilon|^2 \right| \leq Ch^{\frac{1}{3}}.$$

This shows that (30) implies (31).

To prove (30) we are going to make a suitable choice for the test functions in the weak formulation satisfied by \mathbf{u}_ε and η_ε :

$$\begin{aligned}
 (33) \quad & \int_{\Omega_{\eta_\varepsilon}(T)} \mathbf{u}_\varepsilon(T) \cdot (T) - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot \partial_t + \nu \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \nabla \mathbf{u}_\varepsilon : \nabla \\
 & + \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon \cdot - \int_0^T \int_\omega (\partial_t \eta_\varepsilon)^2 b + \int_\omega \partial_t \eta_\varepsilon(T) b(T) \\
 & - \int_0^T \int_\omega \partial_t \eta_\varepsilon \partial_t b + \int_0^T \int_\omega \Delta \eta_\varepsilon \Delta b + \varepsilon \int_0^T \int_\omega \Delta \partial_t \eta_\varepsilon \Delta b \\
 & = \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot + \int_0^T \int_\omega g b + \int_{\Omega_{\eta_I}} \mathbf{u}_I \cdot (0) + \int_\omega \dot{\eta}_I b(0),
 \end{aligned}$$

for all $(\cdot, \cdot) \in (V_{\eta_\varepsilon} \cap H^1(\widehat{\Omega}_{\eta_\varepsilon})) \times (L^2(0, T; H_0^2(\omega)) \cap H^1(0, T; L^2(\omega)))$,

s.t. $(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T$ for a.e. $(t, x, y) \in (0, T) \times \omega$.

We are going to study separately the low frequencies and the high frequencies of $\partial_t \eta_\varepsilon$ and take advantage of the fact that $\partial_t \eta_\varepsilon$ is bounded in $L^2(0, T; H^s(\omega)) \forall 0 \leq s < \frac{1}{2}$ uniformly in ε (see (25)). This implies that we can control, uniformly in ε , the space high frequencies of $\partial_t \eta_\varepsilon$ in $L^2(0, T; L^2(\omega))$.

Definition of admissible test functions. First we introduce a basis of $H_0^2(\omega) \cap L_0^2(\omega)$ by taking eigenfunctions $(\xi_i)_{i \in \mathbb{N}}$ defined by:

$$\begin{cases} \int_\omega \Delta \xi_i \Delta b = \lambda_i \int_\omega \xi_i b & \forall b \in H_0^2(\omega) \text{ s.t. } \int_\omega b = 0, \\ \xi_i \in H_0^2(\omega), \int_\omega \xi_i = 0, \end{cases}$$

with $(\lambda_i)_{i \in \mathbb{N}}$ the sequence of increasing eigenvalues: $\lambda_i > 0, \lambda_i \rightarrow +\infty$. We choose $(\xi_i)_{i \in \mathbb{N}}$ orthonormal in $L^2(\omega)$. We denote by d^{N_0} the L^2 -projection on the finite-dimensional space $\text{span}(\xi_i)_{0 \leq i \leq N_0}$ of any function d and by d^{hf, N_0} the difference $d - d^{N_0}$. Thanks to the choice of the ξ_i , the L^2 -projection on the finite-dimensional space $\text{span}(\xi_i)_{0 \leq i \leq N_0}$ is stable in the L^2 -norm as well as in the H_0^2 -norm. In what follows, we will use the fact that the following property, obtained by Hilbertian interpolation, holds true:

$$(34) \quad \forall d \in H^s(\omega), 0 \leq s < \frac{1}{2}, \quad \|d^{hf, N_0}\|_{L^2(\omega)} \leq \lambda_{N_0}^{-\frac{s}{2}} \|d\|_{H^s(\omega)}.$$

Next, for $\sigma > 1$ we define \mathbf{v}_σ by

$$\mathbf{v}_\sigma(x, y, z) = (\sigma v_1(x, y, \sigma z), \sigma v_2(x, y, \sigma z), v_3(x, y, \sigma z)).$$

If \mathbf{v} is divergence-free, \mathbf{v}_σ is also divergence-free.

We want now to define admissible test functions. We set

$$b_\varepsilon = \int_{t-h}^t \partial_t \eta_\varepsilon^{N_0}(s) ds$$

and

$$\varepsilon = \int_{t-h}^t \left(\overline{(\mathbf{u}_\varepsilon - \mathcal{R}_\alpha(\partial_t \eta_\varepsilon))^\lambda} \right)_\sigma (s) ds + \int_{t-h}^t \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})}(s) ds,$$

where the extension $\mathbf{v} \mapsto \bar{\mathbf{v}}$ is defined by (15) and where \mathcal{R}_α is the lifting operator defined at Lemma 6. Moreover a space regularization of $\mathbf{v}_\varepsilon = \mathbf{u}_\varepsilon - \mathcal{R}_\alpha(\partial_t \eta_\varepsilon)$, denoted by $\mathbf{v}_\varepsilon^\lambda$, has been introduced in order to have $\mathbf{v}_\varepsilon^\lambda$ uniformly bounded in $H^1(0, T; H^1(\mathcal{C}_M))$. It verifies $\operatorname{div}(\mathbf{v}_\varepsilon^\lambda) = 0$, $\mathbf{v}_\varepsilon^\lambda \in L^2(0, T; H_0^1(\Omega_{\eta_\varepsilon}(t)))$, and

$$(35) \quad \begin{aligned} & \|\mathbf{v}_\varepsilon - \mathbf{v}_\varepsilon^\lambda\|_{L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} \longrightarrow 0, \text{ uniformly in } \varepsilon, \text{ as } \lambda \text{ goes to zero,} \\ & \|\mathbf{v}_\varepsilon^\lambda\|_{L^2(0, T; H^1(\Omega_{\eta_\varepsilon}(t)))} \leq C\lambda. \end{aligned}$$

Note that the construction of $\mathbf{v}_\varepsilon^\lambda$ relies on the fact that η_ε converges uniformly to η and that the plate does not touch the bottom of the fluid cavity. Moreover, the uniform convergence of $(\mathbf{v}_\varepsilon^\lambda)_\lambda$ as $\lambda \rightarrow 0$ in $L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))$ is made possible since $\bar{\mathbf{v}}_\varepsilon$ is uniformly bounded in $L^2(0, T; H^{s'}(\mathcal{C}_M))$, $0 < s' < 1/2$ thanks to (17), (25), and (28). Note that, with this choice,

$$\begin{aligned} & \|\partial_t \eta_\varepsilon^{N_0}\|_{L^\infty(0, T; L^2(\omega))} \leq C, \quad \|\partial_t \eta_\varepsilon^{N_0}\|_{L^2(0, T; H^s(\omega))} \leq C, \\ & \|b_\varepsilon\|_{W^{1, \infty}(0, T; L^2(\omega))} \leq C, \quad \|b_\varepsilon\|_{H^1(0, T; H^s(\omega))} \leq C, \\ & \|\overline{\mathbf{v}_\varepsilon^\lambda}\|_{L^\infty(0, T; L^2(\mathcal{C}_M))} \leq C, \quad \|\overline{\mathbf{v}_\varepsilon^\lambda}\|_{L^2(0, T; H^{s'}(\mathcal{C}_M))} \leq C \quad \forall s' < s < 1/2, \end{aligned}$$

and

$$\|\partial_t \eta_\varepsilon^{N_0}\|_{L^\infty(0, T; H^2(\omega))} \leq C_{N_0}, \quad \|b_\varepsilon\|_{W^{1, \infty}(0, T; H^2(\omega))} \leq C_{N_0}, \quad \|\overline{\mathbf{v}_\varepsilon^\lambda}\|_{L^2(0, T; H^1(\mathcal{C}_M))} \leq C_\lambda,$$

where C denotes and will denote in all that follows a strictly positive constant that depends only on the data and not on ε and N_0 , and C_{N_0} (resp., C_λ) denotes and will denote a strictly positive constant that depends on the data and not on ε but may depend on N_0 (resp., λ). The integer N_0 (resp., the real λ) will be fixed later on and will be large enough (resp., small enough). Then for well chosen σ , $(\eta_\varepsilon, b_\varepsilon)$ are admissible tests functions. Indeed, η_ε is divergence-free thanks to the definitions of the lifting operator \mathcal{R}_α , the extension operator $\mathbf{v} \mapsto \bar{\mathbf{v}}$, the operator $\mathbf{v} \mapsto \mathbf{v}_\sigma$, and the regularization $\mathbf{v} \mapsto \mathbf{v}^\lambda$. Moreover η_ε belongs to $H^1(0, T; H^1(\mathcal{C}_M))$. The function b_ε belongs to $H^1(0, T; H_0^2(\omega))$. Both of them are bounded in the previous spaces independently of ε but not of N_0 and λ . Moreover since

$$\|\eta_\varepsilon\|_{L^\infty(0, T; H_0^2(\omega)) \cap W^{1, \infty}(0, T; L^2(\omega))} \leq C,$$

and remembering the imbedding (13), we have

$$\|\eta_\varepsilon - \eta_\varepsilon^-\|_{C^0([0, T] \times \bar{\omega})} \leq Ch^\mu, \quad 0 < \mu < \frac{1}{2}.$$

Thus if σ is such that $\sigma \geq 1 + \frac{2C}{\alpha} h^\mu$, we have

$$\eta_\varepsilon(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = \left(0, \int_{t-h}^t \partial_t \eta_\varepsilon^{N_0}(s, x, y) ds \right)^T \text{ on } \omega.$$

In what follows we choose $\sigma = 1 + \frac{2C}{\alpha} h^\mu$, $0 < \mu < \frac{1}{2}$. Hence, with the choice of test

functions that we made, (33) is written:

$$\begin{aligned}
(36) \quad & - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot ((\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma - (\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma^-) \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})}^-) \\
& + \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon \cdot \varepsilon + \nu \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \nabla \mathbf{u}_\varepsilon : \nabla \varepsilon \\
& + \int_{\Omega_{\eta_\varepsilon}(T)} \mathbf{u}_\varepsilon(T) \cdot \varepsilon(T) - \int_0^T \int_\omega \partial_t \eta_\varepsilon \partial_t (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) \\
& - \int_0^T \int_\omega (\partial_t \eta_\varepsilon)^2 (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) + \int_0^T \int_\omega \Delta \eta_\varepsilon \Delta (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) \\
& + \varepsilon \int_0^T \int_\omega \Delta \partial_t \eta_\varepsilon \Delta (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) + \int_\omega \partial_t \eta_\varepsilon(T) (\eta_\varepsilon^{N_0}(T) - \eta_\varepsilon^{N_0}(T-h)) \\
& = \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot \varepsilon + \int_0^T \int_\omega g (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-).
\end{aligned}$$

The two first terms can be written:

$$\begin{aligned}
(37) \quad & - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot ((\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma - (\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma^-) - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})}^-) \\
& = - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-) - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})}^-) \\
& \quad - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (((\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma - \overline{\mathbf{v}_\varepsilon^\lambda}) - ((\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma^- - (\overline{\mathbf{v}_\varepsilon^\lambda})^-)) \\
& \quad - \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (((\overline{\mathbf{v}_\varepsilon^\lambda}) - \overline{\mathbf{v}_\varepsilon}) - ((\overline{\mathbf{v}_\varepsilon^\lambda})^- - (\overline{\mathbf{v}_\varepsilon})^-)).
\end{aligned}$$

We set $I_1 = \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-)$.

$$\begin{aligned}
I_1 &= -\frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon|^2 + \frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\overline{\mathbf{u}_\varepsilon}^-|^2 - \frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-|^2 \\
&= -\frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon|^2 + \frac{1}{2} \int_0^{T-h} \int_{\Omega_{\eta_\varepsilon}(t+h)} |\overline{\mathbf{u}_\varepsilon}^-|^2 - \frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-|^2 \\
&= \frac{1}{2} \int_0^{T-h} \int_{C_M} (\rho_\varepsilon^+ - \rho_\varepsilon) |\overline{\mathbf{u}_\varepsilon}^-|^2 - \frac{1}{2} \int_{T-h}^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon|^2 - \frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-|^2.
\end{aligned}$$

The same argument we use to prove (32) leads to

$$\int_0^T \int_{C_M} |\rho_\varepsilon^+ - \rho_\varepsilon| |\overline{\mathbf{u}_\varepsilon}^-|^2 \leq Ch^{\frac{1}{3}}.$$

This yields

$$(38) \quad I_1 \leq Ch^{\frac{1}{3}} - \frac{1}{2} \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\overline{\mathbf{u}_\varepsilon} - \overline{\mathbf{u}_\varepsilon}^-|^2.$$

For the second term of (37) $\int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})})$ we have, by taking into account the energy estimate (9),

$$\begin{aligned} & \left| \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})}) \right| \\ & \leq C \|\mathbf{u}_\varepsilon\|_{L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} \|\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})}\|_{L^2(0, T; L^2(C_M))} \\ & \leq C (\|\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})\|_{L^2(0, T; L^2(C_\alpha))} + \|\partial_t \eta_\varepsilon^{hf, N_0}\|_{L^2(0, T; L^2(\omega))}). \end{aligned}$$

Thanks to the properties satisfied by \mathcal{R}_α and in particular (17) for $s = 0$ we have

$$(39) \quad \|\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})\|_{L^2(0, T; L^2(C_\alpha))} \leq C \|\partial_t \eta_\varepsilon^{hf, N_0}\|_{L^2(0, T; L^2(\omega))},$$

and thus we obtain

$$\left| \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})}) \right| \leq C \|\partial_t \eta_\varepsilon^{hf, N_0}\|_{L^2(0, T; L^2(\omega))}.$$

The definition of $\partial_t \eta_\varepsilon^{hf, N_0}$ and the fact that $\partial_t \eta_\varepsilon$ is bounded in $L^2(0, T; H^s(\omega))$, $s < 1/2$ independently of ε , imply, remembering (34), that

$$(40) \quad \left| \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (\overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})} - \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{hf, N_0})}) \right| \leq C \|\partial_t \eta_\varepsilon^{hf, N_0}\|_{L^2(0, T; L^2(\omega))} \leq C \lambda_{N_0}^{-\frac{\sigma}{2}}.$$

We have also for the third term of (37):

$$(41) \quad \begin{aligned} & \left| \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (((\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma - \overline{\mathbf{v}_\varepsilon^\lambda}) - ((\overline{\mathbf{v}_\varepsilon^\lambda})^- - (\overline{\mathbf{v}_\varepsilon^\lambda})^-)) \right| \\ & \leq 2 \|\mathbf{u}_\varepsilon\|_{L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} \|(\overline{\mathbf{v}_\varepsilon^\lambda})_\sigma - \overline{\mathbf{v}_\varepsilon^\lambda}\|_{L^2(0, T; L^2(\Omega_{\eta_\varepsilon}(t)))} \\ & \leq (\sigma - 1) \|\overline{\mathbf{v}_\varepsilon^\lambda}\|_{L^2(0, T; H^1(\Omega_{\eta_\varepsilon}(t)))} \\ & \leq C_\lambda (\sigma - 1) \leq C_\lambda h^\mu, \quad 0 < \mu < \frac{1}{2}. \end{aligned}$$

Finally, due to the properties of the space regularization $\mathbf{v} \mapsto \mathbf{v}^\lambda$, the last term of the right-hand side of (37) goes to zero, uniformly in ε , as λ goes to zero. That is written: $\forall \beta > 0$ there exists $\lambda_0 > 0$ such that $\forall \lambda < \lambda_0$

$$(42) \quad \left| \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon \cdot (((\overline{\mathbf{v}_\varepsilon^\lambda}) - \overline{\mathbf{v}_\varepsilon}) - ((\overline{\mathbf{v}_\varepsilon^\lambda})^- - (\overline{\mathbf{v}_\varepsilon})^-)) \right| \leq C\beta \quad \forall \varepsilon.$$

Now we take care of the convective term

$$I_2 = \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon \cdot \left(\int_{t-h}^t \varepsilon \right),$$

with $\eta_\varepsilon = \overline{(\mathbf{v}_\varepsilon^\lambda)_\sigma} + \overline{\mathcal{R}_\alpha(\partial_t \eta_\varepsilon^{N_0})}$. We have

$$\begin{aligned}
 |I_2| &\leq \frac{1}{2} \int_0^T \|\mathbf{u}_\varepsilon\|_{L^4(\Omega_{\eta_\varepsilon}(t))} \|\nabla \mathbf{u}_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))} \int_{t-h}^t \|\eta_\varepsilon\|_{L^4(\Omega_{\eta_\varepsilon}(t))} \\
 (43) \quad &\leq \frac{\sqrt{h}}{2} \int_0^T \|\mathbf{u}_\varepsilon\|_{L^4(\Omega_{\eta_\varepsilon}(t))} \|\nabla \mathbf{u}_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))} \left(\int_{t-h}^t \|\eta_\varepsilon\|_{L^4(\Omega_{\eta_\varepsilon}(t))}^2 \right)^{\frac{1}{2}} \\
 &\leq C_{N_0, \lambda} \sqrt{h}.
 \end{aligned}$$

The next term to consider is $I_3 = \nu \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \nabla \mathbf{u}_\varepsilon : \nabla (\int_{t-h}^t \eta_\varepsilon)$.

$$\begin{aligned}
 |I_3| &\leq \nu \int_0^T \|\nabla \mathbf{u}_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))} \int_{t-h}^t \|\nabla \eta_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))} \\
 (44) \quad &\leq \nu \sqrt{h} \int_0^T \|\nabla \mathbf{u}_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))} \left(\int_{t-h}^t \|\nabla \eta_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(t))}^2 \right)^{\frac{1}{2}} \\
 &\leq C_{N_0, \lambda} \sqrt{h}.
 \end{aligned}$$

The term $I_4 = \int_{\Omega_{\eta_\varepsilon}(T)} \mathbf{u}_\varepsilon(T) \cdot (\int_{T-h}^T \eta_\varepsilon)$ can be estimated as follows:

$$\begin{aligned}
 |I_4| &\leq \|\mathbf{u}_\varepsilon(T)\|_{L^2(\Omega_{\eta_\varepsilon}(T))} \int_{T-h}^T \|\eta_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(T))} \\
 (45) \quad &\leq \sqrt{h} \|\mathbf{u}_\varepsilon(T)\|_{L^2(\Omega_{\eta_\varepsilon}(T))} \left(\int_{T-h}^T \|\eta_\varepsilon\|_{L^2(\Omega_{\eta_\varepsilon}(T))}^2 \right)^{\frac{1}{2}} \\
 &\leq C \sqrt{h}.
 \end{aligned}$$

We set $I_5 = -\int_0^T \int_\omega \partial_t \eta_\varepsilon \partial_t (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-)$.
We have thanks to the definition of $\eta_\varepsilon^{N_0}$

$$\begin{aligned}
 (46) \quad I_5 &= -\frac{1}{2} \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0})^2 + \frac{1}{2} \int_0^T \int_\omega ((\partial_t \eta_\varepsilon^{N_0})^-)^2 - \frac{1}{2} \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0} - (\partial_t \eta_\varepsilon^{N_0})^-)^2 \\
 &= -\frac{1}{2} \int_{T-h}^T \int_\omega (\partial_t \eta_\varepsilon^{N_0})^2 - \frac{1}{2} \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0} - (\partial_t \eta_\varepsilon^{N_0})^-)^2 \\
 &\leq -\frac{1}{2} \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0} - (\partial_t \eta_\varepsilon^{N_0})^-)^2.
 \end{aligned}$$

For the next term we have

$$\begin{aligned}
 |I_6| &= \left| \frac{1}{2} \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0})^2 (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) \right| \\
 &\leq \frac{1}{2} \int_0^T \|\partial_t \eta_\varepsilon^{N_0}\|_{L^3(\omega)}^2 \|\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-\|_{L^3(\omega)},
 \end{aligned}$$

but, taking into account the continuous imbedding of $H^s(\omega), 0 \leq s < 1/2$, in $L^3(\omega)$ and the fact that $\partial_t \eta_\varepsilon^{N_0}$ is bounded uniformly in ε in $L^2(0, T, H^s(\omega)) \forall 0 \leq s < 1/2$,

$$\|\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-\|_{L^3(\omega)} \leq \int_{t-h}^t \|\partial_t \eta_\varepsilon^{N_0}\|_{L^3(\omega)} \leq C \sqrt{h}.$$

Consequently

$$(47) \quad |I_6| \leq C\sqrt{h}.$$

The next term to consider is $I_7 = \int_0^T \int_\omega \Delta \eta_\varepsilon^{N_0} \Delta (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-)$. It can be estimated as follows:

$$(48) \quad \begin{aligned} |I_7| &\leq \int_0^T \|\Delta \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \int_{t-h}^t \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \\ &\leq \sqrt{h} \int_0^T \|\Delta \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \left(\int_{t-h}^t \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)}^2 \right)^{\frac{1}{2}} \\ &\leq C_{N_0} \sqrt{h}. \end{aligned}$$

The additional viscous term gives for $\varepsilon \leq 1$:

$$(49) \quad \begin{aligned} |I_8| &= \varepsilon \left| \int_0^T \int_\omega \Delta \partial_t \eta_\varepsilon^{N_0} \Delta (\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-) \right| \\ &\leq \int_0^T \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \int_{t-h}^t \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \\ &\leq \sqrt{h} \int_0^T \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \left(\int_{t-h}^t \|\Delta \partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)}^2 \right)^{\frac{1}{2}} \\ &\leq C_{N_0} \sqrt{h}. \end{aligned}$$

We set $I_9 = \int_\omega \partial_t \eta_\varepsilon^{N_0}(T) (\eta_\varepsilon(T) - \eta_\varepsilon^{N_0}(T-h))$.

$$(50) \quad \begin{aligned} |I_9| &= \|\partial_t \eta_\varepsilon(T)\|_{L^2(\omega)} \int_{T-h}^T \|\partial_t \eta_\varepsilon^{N_0}\|_{L^2(\omega)} \\ &\leq Ch. \end{aligned}$$

Next $I_{10} = \int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot (\int_{t-h}^t \varepsilon)$ and $I_{11} = \int_0^T \int_\omega g(\eta_\varepsilon^{N_0} - (\eta_\varepsilon^{N_0})^-)$ can be estimated respectively by $C\sqrt{h}$ and Ch .

These last two estimates and estimates (38)–(50) yield, for all $\beta > 0$, for λ small enough

$$\int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon - \bar{\mathbf{u}}_\varepsilon^-|^2 + \int_0^T \int_\omega (\partial_t \eta_\varepsilon^{N_0} - (\partial_t \eta_\varepsilon^{N_0})^-)^2 \leq C_{N_0, \lambda} h^{\frac{1}{3}} + C \lambda_{N_0}^{-\frac{s}{2}} + C\beta \quad \forall \varepsilon < \varepsilon_0,$$

with $s < \frac{1}{2}$. Moreover,

$$\int_0^T \int_\omega (\partial_t \eta_\varepsilon^{hf, N_0} - (\partial_t \eta_\varepsilon^{hf, N_0})^-)^2 \leq C \lambda_{N_0}^{-\frac{s}{2}}.$$

Then for all $\beta > 0$, if λ is chosen small enough, we have

$$\int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon - \bar{\mathbf{u}}_\varepsilon^-|^2 + \int_0^T \int_\omega (\partial_t \eta_\varepsilon - (\partial_t \eta_\varepsilon)^-)^2 \leq C_{N_0, \lambda} h^{\frac{1}{3}} + C \lambda_{N_0}^{-\frac{s}{2}} + C\beta \quad \forall \varepsilon < \varepsilon_0.$$

Thus, by choosing N_0 large enough and λ small enough, we obtain that $\exists h_0 > 0$ such that $\forall h \leq h_0$

$$\int_0^T \int_{\Omega_{\eta_\varepsilon}(t)} |\mathbf{u}_\varepsilon - \bar{\mathbf{u}}_\varepsilon^-|^2 + \int_0^T \int_\omega (\partial_t \eta_\varepsilon - (\partial_t \eta_\varepsilon)^-)^2 \leq C\beta \quad \forall \varepsilon < \varepsilon_0.$$

This proves Lemma 10. \square

Thanks to Lemma 9, we obtain that $\partial_t \eta_\varepsilon$ is compact in $L^2(0, T; L^2(\omega))$ and that $\bar{\mathbf{u}}_\varepsilon$ is compact in $L^2(0, T; L^2(\mathcal{C}_M))$.

We now want to verify the convergences announced in Proposition 1 and to verify that the equality between the structure velocity and the fluid velocity at the interface holds in the limit.

Let $T > 0$ such that $\inf_\varepsilon \min_{[0, T] \times \bar{\omega}} (1 + \eta_\varepsilon) \geq \alpha > 0$. We will denote any subsequence of $(\eta_\varepsilon, \bar{\mathbf{u}}_\varepsilon)$ by $(\eta_\varepsilon, \bar{\mathbf{u}}_\varepsilon)$. Thanks to the energy estimate and to the compactness properties that have just been derived, we have, denoting by (η, \mathbf{u}) the limit of a subsequence of $(\eta_\varepsilon, \bar{\mathbf{u}}_\varepsilon)$, the following convergences as ε goes to zero:

$$\begin{aligned} \eta_\varepsilon &\rightarrow \eta && \text{uniformly in } C^0([0, T]; C^0(\bar{\omega})), \\ \eta_\varepsilon &\rightharpoonup \eta && \text{weakly in } L^2(0, T; H_0^2(\omega)), \\ \partial_t \eta_\varepsilon &\rightarrow \partial_t \eta && \text{strongly in } L^2(0, T; L^2(\omega)), \\ \partial_t \eta_\varepsilon &\rightharpoonup \partial_t \eta && \text{weakly in } L^2(0, T; W^{1-1/p, p}(\omega)) \quad \forall 1 < p < 2, \\ \partial_t \eta_\varepsilon &\rightharpoonup \partial_t \eta && \text{weakly in } L^2(0, T; H^s(\omega)) \quad \forall 0 \leq s < 1/2, \\ \bar{\mathbf{u}}_\varepsilon &\rightarrow \underline{\mathbf{u}} && \text{strongly in } L^2(0, T; L^2(\mathcal{C}_M)), \\ \rho_\varepsilon \mathbf{u}_\varepsilon &\rightarrow \rho \underline{\mathbf{u}} && \text{strongly in } L^2(0, T; L^2(\mathcal{C}_M)). \end{aligned}$$

Moreover $\rho_\varepsilon \nabla \mathbf{u}_\varepsilon$ tends to some \mathbf{z} weakly in $L^2(0, T; L^2(\mathcal{C}_M))$ as ε goes to zero. It is easy to verify, since $\eta_\varepsilon \rightarrow \eta$ in $C^0([0, T] \times \bar{\omega})$, that $\mathbf{z} = 0$ in $\widehat{\mathcal{C}}_M \setminus \widehat{\Omega}_\eta$ and $\mathbf{z}|_{\widehat{\Omega}_\eta} = \nabla(\underline{\mathbf{u}}|_{\widehat{\Omega}_\eta})$. Thus

$$(51) \quad \rho_\varepsilon \nabla \mathbf{u}_\varepsilon \rightharpoonup \rho \nabla(\underline{\mathbf{u}}|_{\widehat{\Omega}_\eta}) \quad \text{in } L^2(0, T; L^2(\mathcal{C}_M)).$$

Note also that $\underline{\mathbf{u}} = (0, 0, \partial_t \eta)^T$ in $\widehat{\mathcal{C}}_M \setminus \widehat{\Omega}_\eta$, and thus $\underline{\mathbf{u}} = \bar{\mathbf{u}}$ by setting $\underline{\mathbf{u}} = \underline{\mathbf{u}}|_{\widehat{\Omega}_\eta}$.

Next we take care of the equality

$$\mathbf{u}_\varepsilon(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, \partial_t \eta_\varepsilon(t, x, y))^T$$

on $(0, T) \times \omega$. The right-hand side converges to $(0, 0, \partial_t \eta)^T$ strongly in $L^2(0, T; L^2(\omega))$. The left-hand side is the trace of the function $\mathbf{w}_\varepsilon(t, x, y, z) = \mathbf{u}_\varepsilon(t, x, y, z(1 + \eta_\varepsilon(t, x, y)))$ on $z = 1$, and \mathbf{w}_ε converges strongly in $L^2(0, T; L^2(\mathcal{C}_1))$ and weakly in $L^2(0, T; W^{1, p}(\mathcal{C}_1)) \quad \forall 1 < p < 2$ to $\mathbf{u}(t, x, y, z(1 + \eta(t, x, y)))$. Hence by the continuity of the trace mapping on $z = 1$, we have, for a.e. t ,

$$\mathbf{u}(t, x, y, 1 + \eta(t, x, y)) = (0, 0, \partial_t \eta(t, x, y))^T \quad \text{on } \omega.$$

This ends the proof of Proposition 1.

4. Passage to the limit—Proof of Theorem 2. Next we pass to the limit in the weak formulation:

$$\begin{aligned} & \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{u}_\varepsilon(t) \cdot \boldsymbol{\varphi}_\varepsilon(t) - \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} \mathbf{u}_\varepsilon \cdot \partial_t \boldsymbol{\varphi}_\varepsilon + \nu \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} \nabla \mathbf{u}_\varepsilon : \nabla \boldsymbol{\varphi}_\varepsilon \\ & + \int_0^t \int_{\Omega_{\eta_\varepsilon}(s)} (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u}_\varepsilon \cdot \boldsymbol{\varphi}_\varepsilon - \int_0^t \int_\omega (\partial_t \eta_\varepsilon)^2 b + \int_\omega \partial_t \eta_\varepsilon(t) b(t) \\ & - \int_0^t \int_\omega \partial_t \eta_\varepsilon \partial_t b + \varepsilon \int_0^t \int_\omega \Delta \partial_t \eta_\varepsilon \Delta b + \int_0^t \int_\omega \Delta \eta_\varepsilon \Delta b \\ & = \int_0^t \int_{\Omega_{\eta_\varepsilon}(t)} \mathbf{f} \cdot \boldsymbol{\varphi}_\varepsilon + \int_0^t \int_\omega g b + \int_{\Omega(0)} \mathbf{u}_I \cdot \boldsymbol{\varphi}_\varepsilon(0) + \int_\omega \dot{\eta}_I b(0) \end{aligned}$$

for a.e. t and for all $(\boldsymbol{\varphi}_\varepsilon, b) \in (V_{\eta_\varepsilon} \cap H^1(\widehat{\Omega}_{\eta_\varepsilon})) \times (L^2(0, T; H_0^2(\omega)) \times H^1(0, T; L^2(\omega)))$ such that $\boldsymbol{\varphi}_\varepsilon(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T, (t, x, y) \in [0, T] \times \omega$.

The fluid test functions should depend on ε . However, it is sufficient to consider a dense family of test functions, and it can be chosen independent of ε and admissible for any ε small enough.

First we consider test functions of the form $(\boldsymbol{\varphi}^0, 0)$ such that $\boldsymbol{\varphi}^0$ belongs to $\mathcal{D}(\cup_{t \in [0, T]} \{t\} \times \Omega_\eta(t))$ and $\operatorname{div} \boldsymbol{\varphi}^0 = 0$. These test functions satisfy the property that $\boldsymbol{\varphi}^0(t, \cdot) \in \mathcal{D}(\Omega_\eta(t)) \forall t$. For ε small enough, since η_ε converges uniformly to η , $\boldsymbol{\varphi}^0 \in \mathcal{D}(\cup_{t \in [0, T]} \{t\} \times \Omega_{\eta_\varepsilon}(t))$.

The second pair of test functions we consider is $(\boldsymbol{\varphi}^1, b)$, where b belongs to $L^2(0, T; H_0^2(\omega)) \cap H^1(0, T; L^2(\omega))$, with $\int_\omega b = 0$ and for a.e. t

$$\boldsymbol{\varphi}^1 = \begin{cases} (0, 0, b)^T \text{ in } \widehat{\mathcal{C}}_M \setminus \overline{\mathcal{C}}_\alpha, \\ \mathcal{R}(0, 0, \frac{z}{\alpha} b)^T \text{ in } \widehat{\mathcal{C}}_\alpha, \end{cases}$$

where $\widehat{\mathcal{C}}_\alpha = (0, T) \times \mathcal{C}_\alpha$ and \mathcal{R} is a linear lifting operator from $\{\mathbf{w} \in H^{\frac{1}{2}}(\partial \mathcal{C}_\alpha); \int_{\partial \mathcal{C}_\alpha} \mathbf{w} \cdot \mathbf{n} = 0\}$ onto $\{\mathbf{v} \in H^1(\mathcal{C}_\alpha); \operatorname{div}(\mathbf{v}) = 0\}$. We have easily that $\boldsymbol{\varphi}^1$ belongs to $L^2(0, T; H^1(\widehat{\mathcal{C}}_M))$ and $\operatorname{div}(\boldsymbol{\varphi}^1) = 0$. Moreover since $\min_{[0, T] \times \overline{\omega}} (1 + \eta_\varepsilon) \geq \alpha$, $\boldsymbol{\varphi}^1(t, x, y, 1 + \eta_\varepsilon(t, x, y)) = (0, 0, b(t, x, y))^T$ on $(0, T) \times (\omega) \forall \varepsilon$. Furthermore we can choose the linear operator \mathcal{R} such that

$$(52) \quad \left\| \mathcal{R} \left(0, 0, \frac{zb}{\alpha} \right)^T \right\|_{L^2(\omega \times (0, \alpha))} \leq C \|b\|_{L^2(\omega)}.$$

It can be done by solving a Stokes problem in $\omega \times (0, \alpha)$. Indeed this type of inequality can be obtained thanks to a transposition argument and relies on a $H^2 \times H^1$ regularity result for the Stokes problem, which is true here since $\omega \times (0, \alpha)$ is a convex set (see [9] for the regularity result for the Stokes problem). Thus if \mathcal{R} is chosen such that (52) holds, we deduce that, for a.e. $t \in I \subset (0, T)$, and for h small enough

$$\left\| \mathcal{R} \left(0, 0, \frac{zb}{\alpha} \right)^T (t) - \mathcal{R} \left(0, 0, \frac{zb}{\alpha} \right)^T (t + h) \right\|_{L^2(\omega \times (0, \alpha))} \leq C \|b(t) - b(t + h)\|_{L^2(\omega)}.$$

Since $\partial_t b \in L^2(0, T; L^2(\omega))$ this implies that $\partial_t \mathcal{R}(0, 0, \frac{zb}{\alpha})^T \in L^2(0, T; L^2(\mathcal{C}_\alpha))$ and that $\partial_t \boldsymbol{\varphi}^1 \in L^2(0, T; L^2(\widehat{\mathcal{C}}_M))$. Consequently $(\boldsymbol{\varphi}^1, b)$ is a pair of admissible test functions for all ε .

With both types of test functions, it is easy to pass to the limit in the weak formulation as ε goes to zero. Since the considered family of test functions is dense in the set of functions $(\mathbf{u}, b) \in (V_\eta \cap H^1(\widehat{\Omega}_\eta)) \times (L^2(0, T; H_0^2(\omega)) \times H^1(0, T; L^2(\omega)))$ such that $(t, x, y, 1 + \eta(t, x, y)) = (0, 0, b(t, x, y))^T, (t, x, y) \in [0, T] \times \omega$, we obtain the existence of one weak solution on $(0, T)$ of (22) that moreover satisfies the energy estimate (27) by passing to the limit as ε tends to zero in (9).

Eventually, we show that we can extend the solution, as long as we have $\min_{[0, T] \times \overline{\omega}} (1 + \eta) > 0$. We do exactly as in [4], but for the sake of completeness we reproduce the proof here. We build an increasing sequence of times $(T_k)_{k \geq 1}$ as follows. First we choose a time $T_1 > 0$ such that there exists a weak solution up to T_1 , with $m_1 = \min_{[0, T_1] \times \overline{\omega}} (1 + \eta) > 0$. Possibly changing slightly T_1 , we may moreover assume that $\eta(T_1) \in H_0^2(\omega), \partial_t \eta(T_1) \in L^2(\omega)$, and $\mathbf{u}(T_1) \in L^2(\Omega_\eta(T_1))$ (since this is true for almost each time).

Now let $k \geq 1$, and assume that we have built a solution up to some time T_k , with $m_k = \min_{[0, T_k] \times \overline{\omega}} (1 + \eta) > 0$. Our construction allows us to build an extension of our solution on some time interval starting from T_k . Thanks to the energy estimate (27) (see also (9)), we have for $s \geq T_k$ for any $0 < \lambda < \frac{3}{4}$

$$(53) \quad 1 + \eta(s) \geq 1 + \eta(T_k) - (s - T_k)^\lambda C(T_k, s) \geq m_k - (s - T_k)^\lambda C(T_k, s),$$

with

$$C(T_k, s) = \tilde{C} \left(\|\mathbf{u}(T_k)\|_{L^2(\Omega_{\eta T})}, \|\eta(T_k)\|_{H_0^2(\omega)}, \|\partial_t \eta(T_k)\|_{L_0^2(\omega)}, \int_{T_k}^s \exp(s - u) (\|\mathbf{f}\|_{L^2(\Omega_\eta(u))}(u) + \|g\|_{L^2(\omega)}(u)) du \right),$$

where \tilde{C} is positive and nondecreasing with respect to its arguments, and $C(T_k, s) \leq C(0, s)$. This a priori estimate shows that if we let

$$\tau_k = \min\{1, (m_k/2C(T_k, T_k + 1))^{\frac{1}{\lambda}}\},$$

we can build a solution starting from $\mathbf{u}(T_k)$ and $\eta(T_k), \partial_t \eta(T_k)$ up to the time $T_k + \tau_k$ (this corresponds to choosing $\alpha = m_k/2$ in the construction of the solution). The time T_{k+1} is chosen close to $T_k + \tau_k$ (in $[T_k + \tau_k/2, T_k + \tau_k]$), in order to have also $\eta(T_k) \in H_0^2(\omega), \partial_t \eta(T_{k+1}) \in L^2(\omega)$, and $\mathbf{u}(T_{k+1}) \in L^2(\Omega_\eta(T_{k+1}))$.

If the sequence $(T_k)_{k \geq 1}$ is infinite, we let $T^* = \sup_k T_k$. If $T < +\infty$, it must be that $m = \min_{[0, T^*] \times \overline{\omega}} (1 + \eta) = 0$. Otherwise, we have $m_k \geq m$ for all k , and hence $\tau_k \geq \min\{1, (m/2C(0, T^*))^{\frac{1}{\lambda}}\} > 0$. But $T_{k+1} - T_k \geq \tau_k/2$ and goes to zero, a contradiction. This achieves the proof of the theorem.

5. Conclusion. We have proved the existence of at least one weak solution for a three-dimensional fluid-plate interaction problem without any (artificial) viscosity of the structure.

REFERENCES

[1] R. A. ADAMS, *Sobolev Spaces*, Pure Appl. Math. 65, Academic Press, New York, 1975.
 [2] M. BOULAKIA, *Existence of weak solutions for the motion of an elastic structure in an incompressible viscous fluid*, C. R. Math. Acad. Sci. Paris, 336 (2003), pp. 985–990.
 [3] H. BREZIS, *Analyse Fonctionnelle, Théorie et applications [Theory and applications]*, Masson, Paris, 1983.

- [4] A. CHAMBOLLE, B. DESJARDINS, M. J. ESTEBAN, AND C. GRANDMONT, *Existence of weak solutions for the unsteady interaction of a viscous fluid with an elastic plate*, J. Math. Fluid Mech., 7 (2005), pp. 368–404.
- [5] C. CONCA, J. H. SAN MARTÍN, AND M. TUCSNAK, *Existence of solutions for the equations modelling the motion of a rigid body in a viscous fluid*, Comm. Partial Differential Equations, 25 (2000), pp. 1019–1042.
- [6] D. COUTAND AND S. SHKOLLER, *Motion of an elastic solid inside an incompressible viscous fluid*, Arch. Ration. Mech. Anal., 176 (2005), pp. 25–102.
- [7] D. COUTAND AND S. SHKOLLER, *The interaction between quasilinear elastodynamics and the Navier-Stokes equations*, Arch. Ration. Mech. Anal., 179 (2006), pp. 303–352.
- [8] H. B. DA VEIGA, *On the existence of strong solutions to a coupled fluid-structure evolution problem*, J. Math. Fluid Mech., 6 (2004), pp. 21–52.
- [9] M. DAUGE, *Stationary Stokes and Navier-Stokes systems on two- or three-dimensional domains with corners, Part I. Linearized equations*, SIAM J. Math. Anal., 20 (1989), pp. 74–97.
- [10] B. DESJARDINS AND M. J. ESTEBAN, *Existence of weak solutions for the motion of rigid bodies in a viscous fluid*, Arch. Ration. Mech. Anal., 146 (1999), pp. 59–71.
- [11] B. DESJARDINS AND M. J. ESTEBAN, *On weak solutions for fluid-rigid structure interaction: Compressible and incompressible models*, Comm. Partial Differential Equations, 25 (2000), pp. 1399–1413.
- [12] B. DESJARDINS, M. J. ESTEBAN, C. GRANDMONT, AND P. LE TALLEC, *Weak solutions for a fluid–structure interaction model*, Rev. Mat. Complut., 14 (2001), pp. 523–538.
- [13] E. FEIREISL, *On the motion of rigid bodies in a viscous compressible fluid*, Arch. Ration. Mech. Anal., 167 (2003), pp. 281–308.
- [14] G. P. GALDI, *On the motion of a rigid body in a viscous liquid: A mathematical analysis with applications*, in Handbook of Mathematical Fluid Dynamics, Vol. I, North-Holland, Amsterdam, 2002, pp. 653–791.
- [15] V. GIRAULT AND P.-A. RAVIART, *Finite Element Methods for Navier-Stokes Equations, Theory and Algorithms*, Springer-Verlag, Berlin, 1986.
- [16] C. GRANDMONT AND Y. MADAY, *Existence for an unsteady fluid-structure interaction problem*, M2AN Math. Model. Numer. Anal., 34 (2000), pp. 609–636.
- [17] M. D. GUNZBURGER, H.-C. LEE, AND G. A. SEREGIN, *Global existence of weak solutions for viscous incompressible flows around a moving rigid body in three dimensions*, J. Math. Fluid Mech., 2 (2000), pp. 219–266.
- [18] K.-H. HOFFMANN AND V. N. STAROVOITOV, *On a motion of a solid body in a viscous fluid. Two-dimensional case*, Adv. Math. Sci. Appl., 9 (1999), pp. 633–648.
- [19] J.-L. LIONS AND E. MAGENES, *Non-homogeneous Boundary Value Problems and Applications, Vol. I*, Grundlehren Math. Wiss. 181, Springer-Verlag, New York, 1972 (in English).
- [20] J.-L. LIONS AND E. MAGENES, *Non-homogeneous Boundary Value Problems and Applications, Vol. I*, Grundlehren Math. Wiss. 181, Springer-Verlag, New York, 1972 (in English).
- [21] J. A. SAN MARTÍN, V. STAROVOITOV, AND M. TUCSNAK, *Global weak solutions for the two-dimensional motion of several rigid bodies in an incompressible viscous fluid*, Arch. Ration. Mech. Anal., 161 (2002), pp. 113–147.
- [22] D. SERRE, *Chute libre d’un solide dans un fluide visqueux incompressible, Existence*, Japan J. Appl. Math., 4 (1987), pp. 99–110.
- [23] J. SIMON, *Compact sets in the space $L^p(0, T; B)$* , Ann. Mat. Pura Appl., 146 (1987), pp. 65–96.
- [24] T. TAKAHASHI, *Analysis of strong solutions for the equations modeling the motion of a rigid-fluid system in a bounded domain*, Adv. Differential Equations, 8 (2003), pp. 1499–1532.
- [25] T. TAKAHASHI AND M. TUCSNAK, *Global strong solutions for the two dimensional motion of a rigid body in a viscous fluid*, J. Math. Fluid Mech., 6 (2004), pp. 53–77.
- [26] R. TEMAM, *Navier-Stokes Equations, Theory and Numerical Analysis*, Stud. Math. Appl. 2, North-Holland, Amsterdam, 1977.

TRANSMISSION EIGENVALUES*

LASSI PÄIVÄRINTA[†] AND JOHN SYLVESTER[‡]

Abstract. The scattering of a time-harmonic plane wave in an inhomogeneous medium is modeled by the scattering problem for the Helmholtz equation. A transmission eigenvalue is a wavenumber at which the scattering operator has a nontrivial kernel or cokernel. Because many sampling methods for locating scatterers succeed only at wavenumbers that are not transmission eigenvalues, they have been studied for some time. Nevertheless, the existence of transmission eigenvalues has previously been proved only for radial scatterers. In this paper, we prove existence for scatterers without radial symmetry.

Key words. inverse scattering, Helmholtz equation, inverse problems, transmission eigenvalues

AMS subject classifications. 81U40, 35P25

DOI. 10.1137/070697525

1. Introduction. The scattering of a time-harmonic plane wave in an inhomogeneous medium is modeled by the scattering problem for the Helmholtz equation. The *total wave* u satisfies the perturbed Helmholtz equation

$$(1.1) \quad (\Delta + k^2(1 + m))u = 0 \quad \text{in } \mathbb{R}^n.$$

The function $m(x)$ denotes the perturbation of the index of refraction from the constant background medium; i.e., $n^2(x) = 1 + m(x)$. We insist that $-1 < m(x)$ be compactly supported and bounded. The relative (far field) scattering operator, s^+ , compares the asymptotics of solutions of the free Helmholtz equation to those of (1.1). Both the linear sampling method and the factorization method use the range of this operator to find the support of the scatterer m . These methods are known to succeed at wavenumbers k for which the range of that operator is dense among all far field patterns (i.e., dense in $L^2(S^{n-1})$). If there exists a bounded domain D that contains the support of $m(x)$, and the wavenumber k is not a *transmission eigenvalue* as defined below, then the range of the scattering operator is dense [5].

DEFINITION 1. A wavenumber k is called a *transmission eigenvalue* if there exists a nontrivial pair (v, w) solving

$$(1.2) \quad \Delta w + k^2 n^2(x)w = 0 \quad \text{in } D,$$

$$(1.3) \quad \Delta v + k^2 v = 0 \quad \text{in } D,$$

$$(1.4) \quad w = v, \quad \frac{\partial w}{\partial \nu} = \frac{\partial v}{\partial \nu} \quad \text{on } \partial D.$$

If D is not smooth enough, we replace (1.4) with the condition that $u - v \in H_0^2(D)$. Under the conditions that $m > 0$ or $m < 0$ on its support, it has been shown that the set of *transmission eigenvalues* is at most discrete [4], [12], but existence has

*Received by the editors July 17, 2007; accepted for publication (in revised form) April 4, 2008; published electronically July 18, 2008.

<http://www.siam.org/journals/sima/40-2/69752.html>

[†]Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland (ljp@rni.helsinki.fi). This author's research was supported by a grant from the Academy of Finland.

[‡]Department of Mathematics, University of Washington, Seattle, WA 98195 (sylvest@u.washington.edu). This author's research was supported by ONR grant N00014-05-1-0716 and NSF grant DMS-0355455.

been established only for m , which depends only on the radius [6]. Under certain conditions, knowledge of the transmission eigenvalues uniquely determines a radial scatterer [9], [10]. For nonradial scatterers, transmission eigenvalues have also been used to infer simple properties of the scatterer [3].

Under the hypothesis that the infimum of $|m|$ is large enough, we prove existence of transmission eigenvalues, as well as upper and lower bounds on the first transmission eigenvalue. The existence and upper bounds are new; the lower bounds are results from [7] and [3].

In [7], we showed that the following three conditions were equivalent. Notice that (1.5) differs from (1.3) in that the condition below requires that v solve the free Helmholtz equation in all of \mathbb{R}^n rather than just in D . Such v which can be represented as superpositions of plane waves with L^2 densities are called Herglotz wave-functions.

1. There exists a nontrivial pair (v, w) solving

$$(1.5) \quad \begin{aligned} \Delta w + k^2 n^2(x)w &= 0 && \text{in } D, \\ \Delta v + k^2 v &= 0 && \text{in } \mathbb{R}^n, \\ w = v, \quad \frac{\partial w}{\partial \nu} &= \frac{\partial v}{\partial \nu} && \text{on } \partial D. \end{aligned}$$

2. There exists a nontrivial $\mu^0 \in \ker s^+$.
3. There exists a nontrivial $\mu^0 \in \text{coker } s^+$.

In the case that v is a Herglotz wave-function, its asymptotic expansion (its far field) belongs to both the kernel and the cokernel of the far field scattering operator. We will show below that, for scatterers supported in a compact set D , the far field scattering operator has a natural extension, and that *transmission eigenvalues* are exactly the wavenumbers for which this natural extension has a kernel or cokernel.

2. The Helmholtz equation and the scattering operator. The scattering operator relates the solutions of (1.1) to solutions of the free Helmholtz equation in all of \mathbb{R}^n :

$$(2.1) \quad (\Delta + k^2) u^0 = 0 \quad \text{in } \mathbb{R}^n.$$

We refer to solutions of (2.1) with finite B^* -norm, defined by

$$\|u^0\|_{B^*} = \sup_{R>0} \frac{1}{\sqrt{R}} \|u^0\|_{L^2(B_R)},$$

as *incident waves* or *free waves*. An *outgoing wave* is a solution to the Helmholtz equation with a compactly supported source f

$$(2.2) \quad (\Delta + k^2) v^+ = f \quad \text{in } \mathbb{R}^n$$

that satisfies the Sommerfeld radiation condition

$$(2.3) \quad \lim_{r \rightarrow \infty} r^{\frac{n-1}{2}} \left(\frac{\partial v^+}{\partial r} - ikv^+ \right) = 0$$

or, equivalently (for $k > 0$), a *limiting absorption principle*

$$(2.4) \quad v^+ = \lim_{\varepsilon \downarrow 0} v_\varepsilon^+,$$

where v_ε^+ is the unique solution to (2.2) with $k^2 \in \mathbb{R}$ replaced by $k^2 + i\varepsilon$ (see, e.g., section 4 of [1]). We could also define an outgoing wave as a solution to

$$(2.5) \quad (\Delta + k^2(1 + m)) w^+ = g$$

with a compactly supported g , and satisfying (2.3) or (2.4). Because m is compactly supported, the definition based on (2.5) and that based on (2.2) coincide. That is, an outgoing solution v^+ to (2.2) is also an outgoing solution w^+ to (2.5) with $g = f - mv^+$.

Existence and uniqueness of outgoing solutions to (2.5) were proved by Agmon in weighted L^2 spaces¹ [1]. Theorem 2 below is a special case of results in [2], and parts of Theorems 3 and 4 are special cases of results in [2] and section 14 of [8].

Define B to be the completion of $C_0^\infty(\mathbb{R}^n)$ in the B -norm

$$\|f\|_B = \|f\|_{L^2(|x| \in [0,1])} + \sum_{j=1}^\infty \frac{1}{\sqrt{2^j}} \|f\|_{L^2(|x| \in [2^j, 2^{j+1}])}.$$

THEOREM 2. *For every compactly supported g , there exists a unique outgoing solution to (2.5), with*

$$\|w^+\|_{B^*} \leq C \|g\|_B,$$

where the constant C depends on m and k .

Because compactly supported functions are dense in B , the correspondence in Theorem 2 defines a bounded map

$$G_m^+ : B \longrightarrow B^*$$

mapping $g \in B$ to w^+ . In the rest of the paper, whenever we refer to waves, we mean subspaces of B^* :

1. B^0 is the subspace of incident waves, i.e., solutions to (2.1).
2. B^m is the subspace of total waves, i.e., solutions to (1.1).
3. B^+ is the subspace of outgoing waves, the range of G_0^+ .²

Both B^0 and B^m are closed in the B^* topology. One way to see this is to note that B^* convergence implies convergence in the sense of tempered distributions so that any u in the closure of B^0 or B^m must satisfy (2.1) or (1.1), respectively, in the sense of distributions. As Schwartz class functions are dense in B , the equations are satisfied in the B^* sense as well. The plane waves, $e^{ik \cdot \Theta \cdot x}$, are not in B^0 . We shall note in Theorem 4 below that B^0 consists of the Herglotz wave-functions, solutions to (2.1) which have square integrable far fields. The far fields of the plane waves are Dirac deltas. The subspace B^+ is not closed in the B^* topology. In particular, every function in B^+ is in $H_{loc}^2(\mathbb{R}^n)$ and satisfies the radiation condition (2.3), and every compactly supported function in $H_{loc}^2(\mathbb{R}^n)$ belongs to B^+ . Because we have defined B^+ as the range of G_0^+ , which is injective, B^+ is a Banach space with norm $\|G_0^+ f\|_{B^+} := \|f\|_B$.

A straightforward consequence of Theorem 2 is the correspondence between incident and total waves.

¹ $\|f\|_{L_\delta^2} = \|(1 + |x|^2)^{\frac{\delta}{2}} f\|_{L^2}$.

²This is also the range of G_m^+ for any bounded compactly supported m .

THEOREM 3. *Every total wave has a unique decomposition into an incident wave plus a scattered wave, and every incident wave has a unique decomposition as a total wave minus a scattered wave:*

$$(2.6) \quad v^m = v^0 + v^+,$$

$$(2.7) \quad u^0 = u^m - u^+.$$

Moreover, the scattering map \mathcal{S} , defined as $u^0 \mapsto u^m$, is an isomorphism from B^0 onto B^m .

Proof. We prove the second assertion first. Any u_0 that solves (2.1) also solves

$$(\Delta + k^2(1 + m)) u^0 = k^2 m u^0.$$

Let u^+ be the unique outgoing solution to (2.5) with $g = -k^2 m u^0$. Note that

$$\|u^+\|_{B^*} \leq C_1 \|g\|_B \leq k^2 C_1 \|m u^0\|_B \leq k^2 C_2 \|u^0\|_{B^*},$$

where both constants depend on an upper bound for m and the size of its support.

Defining $u^m = u^0 + u^+$ and noting that it satisfies (1.1) gives decomposition (2.7) and the estimate

$$(2.8) \quad \|u^m\|_{B^*} \leq C_3 \|u^0\|_{B^*}.$$

If $u^0 = w^m - w^+$ is another such decomposition, then w^+ must also satisfy (2.5) with $g = -k^2 m u^0$, but (2.5) has a unique outgoing solution, so $w^+ = u^+$ and $w^m = u^m$.

Similarly, any v^m solving (1.1) is a solution to

$$(\Delta + k^2) v^m = -k^2 m v^m.$$

Let v^+ be the unique outgoing solution to (2.2) with $f = -k^2 m v^m$, and set $v^0 = v^m - v^+$. Uniqueness follows as in the paragraph above, as does the estimate

$$(2.9) \quad \|u^0\|_{B^*} \leq C_3 \|u^m\|_{B^*}.$$

The existence and uniqueness of the two decompositions (2.6) and (2.7) along with the estimates (2.8) and (2.9) justify the last statement in the theorem—that the scattering map is an isomorphism. \square

In order to see the relationship between the *scattering operator* we have defined above and the scattering operator defined on far fields, we need to discuss asymptotics.

THEOREM 4. *Let $u^0 \in B^0$, $u^+ \in B^+$, and $u^m \in B^m$; then, in spherical coordinates $x = r\Theta$ for large r ,*

$$(2.10) \quad u^0 \sim \mu^0(\Theta) \frac{e^{ikr}}{(ikr)^{\frac{n-1}{2}}} + \mu^0(-\Theta) \frac{e^{-ikr}}{(-ikr)^{\frac{n-1}{2}}},$$

$$(2.11) \quad u^+ \sim \mu^+(\Theta) \frac{e^{ikr}}{(ikr)^{\frac{n-1}{2}}},$$

$$(2.12) \quad u^m \sim (\mu^m(\Theta) + \gamma(\Theta)) \frac{e^{ikr}}{(ikr)^{\frac{n-1}{2}}} + \mu^m(-\Theta) \frac{e^{-ikr}}{(-ikr)^{\frac{n-1}{2}}}.$$

Moreover, the mappings

$$\begin{aligned} b^0 &: B^0 \longrightarrow L^2(S^{n-1}), \\ b^m &: B^m \longrightarrow L^2(S^{n-1}) \end{aligned}$$

defined by $u^0 \mapsto \mu^0$ and by $u^m \mapsto \mu^m$ are isomorphisms. The mapping

$$b^+ : B^+ \longrightarrow L^2(S^{n-1})$$

defined by $u^+ \mapsto \mu^+$ is surjective.³ Every compactly supported function in $H^2(\mathbb{R}^n)$ belongs to its kernel and any function in its kernel is compactly supported (Rellich’s lemma).

Sketch of proof. The operator $(b^0)^{-1}$ is known as the Herglotz operator. We can start with any $\mu_0 \in L^2(S^{n-1})$ and define

$$u^0 = \mathcal{H}\mu^0 = \int_{S^{n-1}} e^{ik\Theta \cdot x} \mu^0(\Theta) dS_\Theta.$$

Noting that any $u^0 \in B^0$ is the inverse Fourier transform of a distribution supported on the sphere $|\xi|^2 - k^2 = 0$ shows that u^0 must have this form, but it requires an estimate [2] to see that $\mu^0 \in L^2$. A stationary phase calculation shows that $\mathcal{H}\mu^0$ has the asymptotics (2.10) when μ^0 is smooth. Again we refer the reader to [2] for the estimate that $\|u\|_{B^*} \leq C\|\mu^0\|_{L^2(S^{n-1})}$.

A similar Fourier transform calculation combined with the limiting absorption principle, or a calculation of the asymptotics of the outgoing Green’s function, gives (2.11). Alternatively, we may note that $u^+(x, k) - u^+(x, -k)$ belongs to B^0 and deduce (2.11) from (2.10).⁴ The surjectivity of b^+ then follows from the surjectivity of b^0 .

Because of the decomposition $u^m = u^0 + u^+$ in Theorem 3, (2.12) follows from (2.10) and (2.11). Rellich’s lemma and unique continuation imply that b^0 and b^m are injective, as well as the final statement in the theorem. \square

We refer to the large r asymptotics as the *far fields* of the corresponding waves; e.g., the far field of u^0 is μ^0 , the far field of u^+ is μ^+ , and the far field of u^m is μ^m . We use (2.12) to define the far field (relative scattering) operator

$$s^+ : L^2(S^{n-1}) \longrightarrow L^2(S^{n-1})$$

by

$$(2.13) \quad \mu^m \mapsto \gamma.$$

The wave scattering operator \mathcal{S} and the far field operator s^+ are closely related.

LEMMA 5. *Let u^0 and w^0 belong to B^0 , with far fields μ^0 and ω^0 , respectively. Then*

$$\int_{\mathbb{R}^n} \overline{u^0} m \mathcal{S} w^0 = \frac{-2i}{k^n} \int_{S^{n-1}} \overline{\mu^0} s^+ \omega^0.$$

³The surjectivity of b^+ , b^0 , and b^m is perhaps the main reason for replacing the L^2_δ spaces in [1] with the Besov spaces, B and B^* , of [2].

⁴Changing the sign of k reverses the sign of ϵ in the limiting absorption principle and changes the sign of the second term in the Sommerfeld radiation condition, thus specifying the unique *incoming*, rather than the *outgoing*, solution.

Proof.

$$\int_{\mathbb{R}^n} \overline{u^0} k^2 m \mathcal{S} w^0 = \int_{\mathbb{R}^n} \overline{u^0} k^2 m w^m,$$

where $w^m = w^0 + w^+$, as in (2.6),

$$\begin{aligned} (2.14) \quad &= - \int_{\mathbb{R}^n} \overline{u^0} (\Delta + k^2) w^m \\ &= - \int_{\mathbb{R}^n} \overline{u^0} (\Delta + k^2) w^+ \\ &= - \lim_{R \rightarrow \infty} \int_{|x| < R} \overline{u^0} (\Delta + k^2) w^+ \\ &= \lim_{R \rightarrow \infty} \int_{|x|=R} \frac{\partial \overline{u^0}}{\partial \nu} w^+ - \overline{u^0} \frac{\partial w^+}{\partial \nu}. \end{aligned}$$

Making use of the asymptotics in (2.11) and (2.12) gives

$$(2.15) \quad = \frac{-2i}{k^{n-2}} \int_{S^{n-1}} \overline{\mu^0} \omega^+,$$

where ω^+ denotes the far field of w^+ ,

$$= \frac{-2i}{k^{n-2}} \int_{S^{n-1}} \overline{\mu^0} s^+ \omega^0. \quad \square$$

A consequence of Lemma 5 is a natural definition of the relative scattering operator, which does not explicitly use asymptotics.

THEOREM 6. *If we define*

$$\mathcal{S}^+ : B^0 \longrightarrow B^{0*}$$

by

$$(2.16) \quad w^0 \xrightarrow{\mathcal{S}^+} \frac{-k^n}{2i} \langle m \mathcal{S} w^0, \cdot \rangle,$$

then

$$(2.17) \quad s^+ = \mathcal{H}^* \mathcal{S}^+ \mathcal{H}.$$

Proof.

$$\begin{aligned} &\langle \mathcal{H}^* \mathcal{S}^+ \mathcal{H} \omega^0, \overline{\mu^0} \rangle \\ &= \langle \mathcal{S}^+ \mathcal{H} \omega^0, \overline{\mathcal{H} \mu^0} \rangle \\ &= \frac{-k^n}{2i} \int_{\mathbb{R}^n} \overline{\mathcal{H} \mu^0} m \mathcal{S} \mathcal{H} \omega^0 \\ &= \frac{-k^n}{2i} \int_{\mathbb{R}^n} \overline{u^0} m \mathcal{S} w^0 \\ &= \int_{S^{n-1}} \overline{\mu^0} s^+ \omega^0 \\ &= \langle s^+ \omega^0, \overline{\mu^0} \rangle. \quad \square \end{aligned}$$

Remark 7. Because $B^0 \subset L^2_{-\delta}$ for any $\delta > \frac{1}{2}$, any $l \in B^{0*}$ has a (nonunique) extension to an element of $L^2_{\delta} = L^2_{-\delta}^*$, so elements of B^{0*} can be represented as functions (and called sources).

3. A generalized scattering operator. We describe an incident wave $\mathcal{H}\alpha$ as illuminating the scatterer m . If we use the far field operator s^+ , the illumination must always come from the *sphere at infinity*. Many useful sources of illumination are generated by sources outside the scatterer. The waves generated by such sources are never incident waves, although they can be approximated by incident waves on certain compact sets. Solutions to the transmission eigenvalue problem (1.2)–(1.4) are not incident waves, so they do not have a direct interpretation in terms of the far field scattering operator. They do, however, span exactly the kernel of the scattering operator we will define below.

If $m \in L^\infty$ is supported in a bounded domain D , Theorem 2 tells us that we can find a unique $u^+ \in B^+$ solving

$$(\Delta + k^2(1 + m)) u^+ = k^2 m u^0$$

for any $u^0 \in L^2(D)$. It follows that $u^m = u^0 + u^+ \in L^2(D)$. Thus the scattering operator \mathcal{S} has a natural extension,

$$\mathcal{S}_D : B_D^0 \longrightarrow B_D^m,$$

where we use the definitions

$$\begin{aligned} B_D^0 &= \{w \in L^2(D) \mid (\Delta + k^2) w = 0 \text{ in } D\}, \\ B_D^m &= \{w \in L^2(D) \mid (\Delta + k^2(1 + m)) w = 0 \text{ in } D\}. \end{aligned}$$

The relative scattering operator \mathcal{S}^+ has a similar extension:

$$\begin{aligned} \mathcal{S}_D^+ : B_D^0 &\longrightarrow B_D^{0*}, \\ w^0 &\xrightarrow{\mathcal{S}_D^+} \frac{k^n}{2i} \langle m \mathcal{S}_D w^0, \cdot \rangle. \end{aligned}$$

The scattering and relative scattering operators \mathcal{S}_D and \mathcal{S}_D^+ are extensions of \mathcal{S} and \mathcal{S}^+ in the sense that, for $u^0, w^0 \in B^0$, then $\mathcal{S}_D u^0$ is the restriction of $\mathcal{S} u^0$ to D and $\langle \mathcal{S}_D^+ u^0, w^0 \rangle = \langle \mathcal{S}^+ u^0, w^0 \rangle$.

4. The interior transmission problem. Let D be a bounded domain and $\text{supp } m \subset D$. We will use the notation

$$\begin{aligned} P^0 &= (\Delta + k^2), \\ P^m &= (\Delta + k^2(1 + m)), \\ H^k(D) &= \{u \in L^2(D) \mid D^\alpha u \in L^2(D) \quad \forall |\alpha| \leq k\}, \\ H_0^k(D) &= \text{the completion of } C_0^\infty(D) \text{ in } H^k(D). \end{aligned}$$

DEFINITION 8. We say that a wavenumber k is a D -transmission eigenvalue of $m \in L^\infty(D)$ if any of the equivalent conditions in Theorem 9 below are satisfied.

THEOREM 9. The following are equivalent:

1. There exist nontrivial $u^0 \in B_D^0$ and $u^m \in B_D^m$ with $u^0 - u^m \in H_0^2(D)$.⁵

⁵This is a restatement of (1.2)–(1.4). The condition that $u^0 \in B_D^0$ is (1.3), $u^m \in B_D^m$ is (1.2), and $u^0 - u^m \in H_0^2(D)$ is (1.4).

2. There exists nontrivial $u^m \in B_D^m$ such that the unique outgoing solution u^+ to

$$(4.1) \quad P^0 u^+ = -k^2 m u^m$$

belongs to $H_0^2(D)$.

3. There exist nontrivial $u^m \in B_D^m$ and some $v \in H_0^2(D)$ satisfying (4.1).
 4. There exists nontrivial $u^0 \in B_D^0$ such that the unique outgoing solution u^+ to

$$(4.2) \quad P^m u^+ = -k^2 m u^0$$

belongs to $H_0^2(D)$.

5. There exist nontrivial $u^0 \in B_D^0$ and some $v \in H_0^2(D)$ satisfying (4.2).
 6. There exists nontrivial $u^0 \in \ker \mathcal{S}_D^+$.
 7. There exists nontrivial $u^0 \in \text{coker } \mathcal{S}_D^+$.

Proof. We first show that items 1–5 are equivalent.

Condition 2 obviously implies 3, but any $H_0^2(D)$ solution, v , to (4.1) extended to be zero in $\mathbb{R}^n \setminus D$ is outgoing. Since the outgoing solution to (4.1) is unique, $v = u^+$, so 3 implies 2.

Similarly, 4 obviously implies 5, but uniqueness of the outgoing solution to (4.2) implies that any $H_0^2(D)$ solution to (4.2), extended to be zero outside D , must be u^+ , so 5 implies 4.

Because Theorem 3 gives a unique decomposition,

$$(4.3) \quad u^m = u^0 + u^+,$$

the unique outgoing solution to (4.1) is also the unique outgoing solution to (4.2). Thus 4 and 2 are equivalent.

The same decomposition shows that $u^+ = u^m - u^0$, so the left-hand side is in $H_0^2(D)$ if and only if the right-hand side is; hence 1 is equivalent to 2.

The equivalence of items 6 and 2 is based on a calculation. Let u^0 and w_0 belong to B_D^0 ,

$$(4.4) \quad \begin{aligned} \langle \overline{w^0}, \mathcal{S}_D^+ u^0 \rangle &= -\frac{k^n}{2i} \langle \overline{w^0}, m \mathcal{S}_D u^0 \rangle \\ &= -\frac{k^n}{2i} \int_D \overline{w^0} m u^m, \end{aligned}$$

where u^m in the line above and u^+ in the line below are those uniquely related to u^0 by (4.3) and Theorem 3,

$$(4.5) \quad = -\frac{k^{n-2}}{2i} \int_D \overline{w^0} P^0 u^+$$

$$(4.6) \quad = \frac{k^{n-2}}{2i} \int_{\partial D} \frac{\partial \overline{w^0}}{\partial \nu} u^+ - \overline{w^0} \frac{\partial u^+}{\partial \nu}.$$

The right-hand side of (4.6) is clearly zero for every w^0 if $u^+ \in H_0^2(D)$, so 2 implies 6. To see the converse, choose $w^0 = \mathcal{H}w^0$ (i.e., $w^0 \in B^0$ with asymptotics as in (2.10)). Every $u^+ \in B^+$ has asymptotics as in (2.11), so we may continue the previous calculation,

$$\begin{aligned} \langle \overline{w^0}, \mathcal{S}_D^+ u^0 \rangle &= \frac{k^{n-2}}{2i} \lim_{R \rightarrow \infty} \int_{B_R} \frac{\partial \overline{w^0}}{\partial \nu} u^+ - \overline{w^0} \frac{\partial u^+}{\partial \nu} \\ &= \int_{S^{n-1}} \overline{\omega^0} \mu^+. \end{aligned}$$

We conclude that, if the left-hand side vanishes for every w^0 , so does the right-hand side for every ω^0 , so $\mu^+ \equiv 0$. Now Rellich’s lemma and unique continuation tell us that $u^+ \in H_0^2(D)$, so 6 implies 2.

Verifying the equivalence of items 7 and 2 requires a similar computation,

$$\begin{aligned} \frac{2i}{k^{n-2}} \langle u^0, \mathcal{S}_D^+ w^0 \rangle &= - \int_D u^0 k^2 m w^m \\ &= \int_D u^0 P^0 w^m \\ &= \int_{\partial D} \frac{\partial u^0}{\partial \nu} w^m - u^0 \frac{\partial w^m}{\partial \nu} \\ &= \int_{\partial D} \frac{\partial u^0}{\partial \nu} w^m - u^0 \frac{\partial w^m}{\partial \nu} - \int_{\partial D} \frac{\partial u^m}{\partial \nu} w^m - u^m \frac{\partial w^m}{\partial \nu}, \end{aligned}$$

because the second integral on the right is always zero. Combining the two terms gives

$$(4.7) \quad \frac{2i}{k^{n-2}} \langle u^0, \mathcal{S}_D^+ w^0 \rangle = - \int_{\partial D} \frac{\partial u^+}{\partial \nu} w^m - u^+ \frac{\partial w^m}{\partial \nu}.$$

Now 2 implies that the right-hand side of (4.7) is zero, so the left-hand side is zero for every w^0 , which implies 7. If we choose $w^0 \in B^0$ and continue the calculation,

$$(4.8) \quad = \int_{S^{n-1}} \mu^+ \omega^m.$$

Item 7 implies that the integral in (4.8) vanishes for every ω^m , so $\mu^+ \equiv 0$. Rellich’s lemma and unique continuation then guarantee that $u^+ \in H_0^2(D)$, which implies item 2. \square

Note that if $\text{supp } m \subset \tilde{D} \subset D$, then $\mathcal{S}_{\tilde{D}}$ is an extension of \mathcal{S}_D . The smaller we make D , the larger we make B_D^0 , the domain of the operator \mathcal{S}_D . Therefore, if k is a D -transmission eigenvalue of m , then k is also a \tilde{D} -transmission eigenvalue.

5. Existence of transmission eigenvalues. In this section we restrict our attention to the case that $D = \text{supp } m$. We assume further that m is bounded away from zero in D . The theorem below was first proved in [12].

THEOREM 10. *If $|m| > \delta > 0$ in D , then k is a D -transmission eigenvalue if and only if there exists $u^+ \in H_0^2(D)$ satisfying*

$$(5.1) \quad P^m \frac{1}{m} P^0 u^+ = 0.$$

Proof. We show that (5.1) is equivalent to item 2 in Theorem 9. If $u^+ \in H_0^2(D)$ satisfies

$$(5.2) \quad P^0 u^+ = -k^2 m u^m,$$

then

$$(5.3) \quad \frac{1}{m} P^0 u^+ = -k^2 u^m,$$

and

$$P^m \frac{1}{m} P^0 u^+ = 0.$$

To see the reverse implication, suppose that $u^+ \in H_0^2(D)$ satisfies (5.1) (recall that any $u^+ \in H_0^2(D)$ is outgoing), and define u^m so that (5.3) holds. It is a consequence of (5.1) that $u^m \in B_D^m$ so that (5.2) implies (4.1). \square

Theorem 10 tells us that k is a D -transmission eigenvalue whenever the operator $P^m \frac{1}{m} P^0$ has a kernel in $H_0^2(D)$. We will investigate the existence of this kernel (5.1) by examining the spectrum of the operator as k^2 changes. We will make use of several equivalent formulas for $P^m \frac{1}{m} P^0$ which we list below. We will let $\tau = k^2$.

$$\begin{aligned}
 P^m \frac{1}{m} P^0 &= P^0 \frac{1}{m} P^m \\
 &= \Delta \frac{1}{m} \Delta + \tau \left(\Delta \frac{1}{m} + \left(1 + \frac{1}{m} \right) \Delta \right) + \tau^2 \left(1 + \frac{1}{m} \right) \\
 &= (\Delta + \tau) \frac{1}{m} (\Delta + \tau) + \tau (\Delta + \tau) \\
 (5.4) \quad &= (\Delta + \tau(1 + m)) \frac{1}{m} (\Delta + \tau(1 + m)) - \tau (\Delta + \tau(1 + m)).
 \end{aligned}$$

The following lemma asserts that $P^m \frac{1}{m} P^0$, with the appropriate domain, defines a semibounded self-adjoint operator on $L^2(D)$.

LEMMA 11. For $\tau \geq 0$, t_τ , defined by

$$(5.5) \quad t_\tau(u) = \int_D \frac{1}{m} |(\Delta + \tau) u|^2 - \tau \int_D |\text{grad } u|^2 + \tau^2 \int_D |u|^2$$

with form domain $H_0^2(D)$, is a densely defined, closed, semibounded quadratic form on $L^2(D)$. T_τ , The unique densely defined self-adjoint operator associated to t_τ , T_τ , is equal to $P^m \frac{1}{m} P^0$ on its domain

$$(5.6) \quad D(T_\tau) = \left\{ u \in H_0^2(D) \mid \frac{1}{m} (\Delta + \tau) u \in H^2(D) \right\}.$$

Proof. We state without proof that $H_0^2(D)$ is dense in $L^2(D)$. To see that t_τ is semibounded, we write

$$\begin{aligned}
 t_\tau(u) &= \int_D \frac{1}{m} |(\Delta + \tau) u|^2 + \tau \int_D \bar{u} ((\Delta + \tau) u) \\
 &\geq \frac{1}{\sup(m)} \| |(\Delta + \tau) u|^2 - \tau \| (\Delta + \tau) u \| \|u\| \\
 &\geq \left(\frac{1}{\sup(m)} - \tau \varepsilon \right) \| |(\Delta + \tau) u|^2 - \frac{\tau}{\varepsilon} \|u\|^2 \\
 (5.7) \quad &\geq \frac{1}{2 \sup(m)} \| |(\Delta + \tau) u|^2 - 2 \sup(m) \tau^2 \|u\|^2
 \end{aligned}$$

after choosing $\varepsilon = \frac{1}{2\tau \sup(m)}$. We record for later use the consequence of (5.7) that

$$(5.8) \quad \| |(\Delta + \tau) u|^2 \leq 2 \sup(m) t_\tau(u) + (2 \sup(m) \tau)^2 \|u\|^2.$$

Every densely defined semibounded quadratic form defines a unique self-adjoint operator [11, page 278], T_τ , with domain the set of $u \in H_0^2(D)$ such that there is an $f \in L^2(D)$ with

$$t(v, u) = (v, f)$$

for all $v \in H_0^2(D)$, where $t(v, u)$ is the bilinear form

$$\begin{aligned} t(v, u) &= \int_D \overline{((\Delta + \tau)v \frac{1}{m}(\Delta + \tau)u + \tau \bar{v}(\Delta + \tau)u)} \\ &= \int_D \bar{v} \left((\Delta + \tau) \frac{1}{m} (\Delta + \tau) u \right) + \tau \int_D \bar{v} (\Delta + \tau) u \\ &= (v, T_\tau u), \end{aligned}$$

where the second and third equalities hold for all $u \in D(T_\tau)$ and illustrate that $T_\tau = P^m \frac{1}{m} P^0$ with $D(T_\tau)$, as asserted in (5.6). \square

LEMMA 12. T_τ has discrete spectrum which depends continuously on τ .

Proof. Because T_τ is semibounded and $H_0^2(D)$ is compactly embedded in $L^2(D)$, T_τ has a compact resolvent and therefore discrete spectrum. For m 's that are not smooth, the domains of T_τ may depend on τ . We give a direct proof of the continuity of the eigenvalues. We shall show below that, for all positive real σ and τ ,

$$(5.9) \quad t_\sigma(u) \leq (1 + M|\sigma - \tau|) t_\tau(u) + M(\tau^2 + \sigma + 1)|\sigma - \tau| \|u\|^2,$$

where the constant M depends only on m . We recall the min-max characterization of the eigenvalues of a self-adjoint operator defined by a quadratic form [13, p. 71]

$$(5.10) \quad \lambda^n = \max_{W \in \mathcal{W}_n} \min_{\substack{u \in W \\ \|u\|=1}} q(u),$$

where \mathcal{W}_n denotes the codimension n subspaces of the form domain of q . An immediate consequence of (5.10) is that inequalities between quadratic forms imply the same inequalities for their ordered eigenvalues so that (5.9) implies

$$\lambda_\sigma^n \leq (1 + M|\sigma - \tau|) \lambda_\tau^n + M(\tau^2 + \sigma + 1)|\sigma - \tau|$$

and, consequently,

$$(5.11) \quad \lambda_\sigma^n - \lambda_\tau^n \leq |\sigma - \tau| M (\lambda_\tau^n + (\tau^2 + \sigma + 1)).$$

Because we may interchange σ and τ ,

$$(5.12) \quad |\lambda_\sigma^n - \lambda_\tau^n| \leq |\sigma - \tau| M (\max(\lambda_\tau^n, \lambda_\sigma^n) + 2(\tau^2 + \sigma^2 + 1)).$$

First fix σ in (5.11), and set $\tau = 0$ to conclude that each λ_σ^n varies only over a compact set when σ varies over a compact set. Thus the maximum, $\max(\lambda_\tau^n, \lambda_\sigma^n)$, is bounded for σ and τ on compact sets, and therefore (5.12) proves continuity of the eigenvalues.

It remains only to prove (5.9). We begin by writing

$$\begin{aligned} t_\sigma(u) - t_\tau(u) &= (\sigma - \tau) \int \left(\bar{u} \frac{1}{m} (\Delta + \tau) u + \overline{(\Delta + \tau) u} \frac{1}{m} u + \tau^2 |u|^2 \right) \\ &\quad + (\sigma - \tau)^2 \int \left(\bar{u} (\Delta + \tau) u + \frac{1}{m} |u|^2 \right) \\ &\leq |\sigma - \tau| (1 + |\sigma - \tau|) M (\|u\| \|(\Delta + \tau) u\| + \tau^2 \|u\|^2), \end{aligned}$$

where M depends only on $\frac{1}{m}$. For any $\varepsilon > 0$

$$\leq |\sigma - \tau| (1 + |\sigma - \tau|) M \left(\varepsilon \|(\Delta + \tau) u\|^2 + \left(\tau^2 + \frac{1}{\varepsilon} \right) \|u\|^2 \right).$$

We make use of (5.8) to obtain with a different M

$$\begin{aligned} &\leq |\sigma - \tau|(1 + |\sigma - \tau|)M \left(\varepsilon t_\tau(u) + \left((\varepsilon + 1)\tau^2 + \frac{1}{\varepsilon} \right) \|u\|^2 \right) \\ &\leq |\sigma - \tau|M (t_\tau(u) + (\tau^2 + \sigma + 1)\|u\|^2) \end{aligned}$$

after choosing $\frac{1}{\varepsilon} = 1 + |\sigma - \tau|$. \square

LEMMA 13. *If*

$$(5.13) \quad \text{sign}(m) \inf_{u \in H_0^2(D)} \frac{t_\tau(u)}{\|u\|^2} > 0,$$

then τ is not a transmission eigenvalue. If there exists $u \in H_0^2(D)$ such that

$$(5.14) \quad \text{sign}(m) \frac{t_\tau(u)}{\|u\|^2} \leq 0,$$

then there is a transmission eigenvalue $\tau^* \in [0, \tau]$.

Proof. The hypothesis (5.13) implies that the spectrum of T_τ is strictly positive or strictly negative; hence it has no kernel.

The hypothesis (5.14) implies that $\text{sign}(m)T_\tau$ has at least one nonpositive eigenvalue. But $\text{sign}(m)T_0$ is easily seen to be positive definite, so the lowest eigenvalue, which is a continuous function of τ , must have passed through zero for some $\tau^* \in [0, \tau]$. \square

We will use a simple modification of Lemma 13 to show the existence of more than one transmission eigenvalue. We define the multiplicity of a transmission eigenvalue τ_* to be the multiplicity of 0 as an eigenvalue of T_{τ_*} .

LEMMA 14. *If there exists a $\tau > 0$ and a p -dimensional subspace $V^p \in H_0^2(D)$ such that*

$$\text{sign}(m) \frac{t_\tau(u)}{\|u\|^2} \leq 0$$

for all $u \in V^p$, then there are p -transmission eigenvalues, counting multiplicity, in $[0, \tau]$.

Proof. The hypothesis guarantees that t_τ has p negative eigenvalues, counting multiplicity. The continuity of the spectrum implies that each of those eigenvalues must pass through zero as τ^* decreases from τ to 0. Each time an eigenvalue passes through 0, the dimension of the negative definite subspace, V^p , decreases by the multiplicity of the zero eigenvalue, so the sum of the multiplicities of the transmission eigenvalues between 0 and τ must be at least p . \square

We will need a few simple inequalities to prove the theorems to follow. We collect them in the lemma below.

LEMMA 15.

$$(5.15) \quad \lambda_0(D) = \inf_{u \in H_0^1(D)} \frac{\int_D |\text{grad } u|^2}{\int_D |u|^2} = \inf_{u \in H_0^2(D)} \frac{\int_D |\text{grad } u|^2}{\int_D |u|^2} > 0,$$

$$(5.16) \quad \mu_0(D) = \inf_{u \in H_0^2(D)} \frac{\int_D |\Delta u|^2}{\int_D |u|^2} \geq \inf_{u \in H_0^1 \cap H^2} \frac{\int_D |\Delta u|^2}{\int_D |u|^2} = \lambda_0(D)^2.$$

If $u \in H_0^2(D)$,

$$(5.17) \quad \lambda_0(D) \leq \frac{\int_D |\text{grad } u|^2}{\int_D |u|^2} \leq \frac{(\int_D |\Delta u|^2)^{\frac{1}{2}}}{(\int_D |u|^2)^{\frac{1}{2}}}.$$

Proof. The first equality in (5.15) is the Rayleigh–Ritz characterization of the first Dirichlet eigenvalue. The second follows because $H_0^2(D)$ is dense in $H_0^1(D)$. The first equality in (5.16) is the Rayleigh–Ritz characterization of the lowest eigenvalue of the biharmonic operator with Dirichlet boundary conditions, the lowest eigenvalue of the *clamped plate*. The inequality holds because $H_0^2(D) \subset H_0^1(D) \cap H^2(D)$, so the first infimum must be larger. The second infimum is exactly the Rayleigh–Ritz characterization for the lowest eigenvalue of the Dirichlet Laplacian squared,⁶ which is the square of the first Dirichlet eigenvalue, which proves the final equality in (5.16). The first inequality in (5.17) follows from the meaning of infimum and the second from integration by parts and the Cauchy–Schwarz inequality. \square

THEOREM 16. *Suppose that $m > -1$ is a constant. If*

$$(5.18) \quad \tau \leq \min \left(1, \frac{1}{m+1} \right) \lambda_0(D),$$

then τ is not a transmission eigenvalue. If

$$(5.19) \quad \frac{(1 + \frac{m}{2})^2}{1 + m} \geq \frac{\mu_p}{\lambda_0^2} \geq 1,$$

where $\mu_p(D)$ is the $(p+1)$ st clamped plate eigenvalue, then there are $p+1$ transmission eigenvalues τ with

$$(5.20) \quad \tau \leq \left(\frac{m+2}{m+1} \right) \frac{\lambda_0(D)}{2}.$$

Proof. It follows from (5.4) that

$$\begin{aligned} mt_\tau(u) &= \|(\Delta + \tau(1+m))u\|^2 - m\tau \int_D \bar{u}(\Delta + \tau(1+m))u \\ &> m\tau [\|\text{grad } u\|^2 - (1+m)\tau \|u\|^2] \\ &\geq m\tau [\lambda_0(D) - (1+m)\tau] \|u\|^2, \end{aligned}$$

which shows that, for $0 < m$ and $\tau \leq \frac{\lambda_0(D)}{1+m}$, t_τ is positive definite and therefore that τ is not a transmission eigenvalue. If $m < 0$, we express t_τ as in (5.5),

$$\begin{aligned} mt_\tau(u) &= \|(\Delta + \tau)u\|^2 - m\tau \|\text{grad } u\|^2 + m\tau^2 \|u\|^2 \\ &\geq (-m)\tau \|u\|^2 (\lambda_0(D) - \tau), \end{aligned}$$

which shows that t_τ is positive definite as long as $\tau < \lambda_0(D)$ and finishes the proof of the assertion that τ satisfying (5.18) is not a transmission eigenvalue.

⁶Functions in the domain of the square of the Dirichlet Laplacian must satisfy a second boundary condition, $\Delta u|_{\partial D} = 0$. Analogous to the case of the Neumann Laplacian, this is a free boundary condition which does not appear explicitly in the definition of the form domain and therefore does not appear explicitly in (5.16).

To prove the existence of transmission eigenvalues, we will use Lemma 14. Restricting our attention to the sphere, $\|u\|^2 = 1$, we may write

$$\begin{aligned} mt_\tau(u) &= (m + 1)\tau^2 - 2\left(1 + \frac{m}{2}\right)\|\text{grad } u\|^2\tau + \|\Delta u\|^2 \\ &\leq (m + 1)\tau^2 - 2\left(1 + \frac{m}{2}\right)\lambda_0\tau + \|\Delta u\|^2. \end{aligned}$$

We choose $\tau = \frac{1+\frac{m}{2}}{1+m}\lambda_0$ to obtain

$$\leq -\frac{\left(1 + \frac{m}{2}\right)^2}{1 + m}\lambda_0^2 + \|\Delta u\|^2$$

and restrict u to the eigenspace associated with the lowest $p + 1$ clamped plate eigenvalues so that

$$mt_\tau(u) \leq -\frac{\left(1 + \frac{m}{2}\right)^2}{1 + m}\lambda_0^2 + \mu_p.$$

Our hypothesis (5.19) is that this quantity is negative, so the conclusion (5.20) follows from Lemma 14. \square

THEOREM 17. *Suppose that $m \in L^\infty(D)$. If*

$$\tau \leq \min\left(1, \frac{1}{\sup(m) + 1}\right)\lambda_0(D),$$

then τ is not a transmission eigenvalue. If $m > 0$ and

$$(5.21) \quad \inf(m) \geq 4\frac{\mu_p^{\frac{1}{2}}}{\lambda_0} + \frac{\mu_p}{\lambda_0^2},$$

then there are $p + 1$ transmission eigenvalues τ with

$$\tau \leq \frac{\lambda_0(D)}{2} \left(\frac{\inf(m) - 2\frac{\mu_p^{\frac{1}{2}}}{\lambda_0}}{\inf(m) + 1} \right).$$

Proof. For $m > 0$,

$$\begin{aligned} t_\tau(u) &= \int \frac{1}{m} |(\Delta + \tau(1 + m))u|^2 - \tau \int \bar{u}(\Delta + \tau(1 + m))u \\ &\geq \tau\|\text{grad } u\|^2 - \tau^2 \int (1 + m)|u|^2 \\ &\geq \tau\|u\|^2 (\lambda_0 - \tau(1 + \sup m)), \end{aligned}$$

which shows that t_τ is positive definite if $\tau < \frac{\lambda_0}{1+\inf m}$ and therefore that τ is not a transmission eigenvalue. For $m < 0$

$$\begin{aligned} -t_\tau(u) &= \int \frac{1}{|m|} |(\Delta + \tau)u|^2 - \tau \int \bar{u}(\Delta + \tau)u \\ &> \tau\|\text{grad } u\|^2 - \tau^2\|u\|^2 \\ &\geq \tau\|u\|^2 (\lambda_0 - \tau) \end{aligned}$$

so that $-t_\tau$ is positive definite if $\tau < \lambda_0$, completing the proof of the first assertion.

To prove existence, we write

$$t_\tau(u) = \tau^2 \int \left(1 + \frac{1}{m}\right) |u|^2 - \tau \left(\|\text{grad } u\|^2 + \int \frac{1}{m} (\bar{u}\Delta u + u\Delta\bar{u}) \right) + \int \frac{1}{m} |\Delta u|^2.$$

We restrict our attention to functions u with $\|u\|^2 = 1$ and write $S = \sup(\frac{1}{m})$ to see that

$$t_\tau(u) \leq \tau^2(1 + S) - \tau (\|\text{grad } u\|^2 - 2S\|\Delta u\|) + S\|\Delta u\|^2.$$

Restricting to V^p gives

$$t_\tau(u) \leq \tau^2(1 + S) - \tau \left(\lambda_0 - 2S\mu_p^{\frac{1}{2}} \right) + S\mu_p.$$

We minimize the sum of the first two terms by choosing $\tau = \frac{\lambda_0 - 2S\mu_p^{\frac{1}{2}}}{2(1+S)}$ to obtain

$$\leq -\frac{\left(\lambda_0 - 2S\mu_p^{\frac{1}{2}}\right)^2}{4(1+S)} + S\mu_p.$$

If we set $A = \frac{\lambda_0}{\mu_p^{\frac{1}{2}}}$, t_τ restricted to V^p is nonpositive if

$$\begin{aligned} (A - 2S)^2 - 4S(1 + S) &\geq 0, \\ A^2 - 4(A + 1)S &\geq 0, \\ \frac{A^2}{4(A + 1)} &\geq S, \end{aligned}$$

which is equivalent to (5.21). \square

6. Conclusions. Under the hypothesis that the perturbation of the index of refraction is large enough, we have shown the existence of D -transmission eigenvalues, given upper and lower bounds for their locations, and identified the corresponding solutions to the transmission eigenvalue problem with the kernel of a scattering operator.

The upper and lower bounds for the transmission eigenvalues depend on the lowest eigenvalues of the Dirichlet Laplacian and the Dirichlet bi-Laplacian (the *clamped plate* operator). These bounds show that the lowest transmission eigenvalue increases as the L^∞ -norm of m , or the size of its support, decreases. All of our bounds depend on D only through these eigenvalues and on m only through its infimum or supremum. Thus our estimates for the transmission eigenvalues, k^2 , scale with dilations just like the Dirichlet eigenvalues, as the reciprocal of the area of D .

In the Born or weak scattering approximation, there are no transmission eigenvalues [7] when m is strictly positive or strictly negative. Our results are consistent with this, but we do not know if there is a threshold below which there are no transmission eigenvalues, or if the lowest transmission eigenvalue simply goes to infinity as m decreases to zero.⁷

In the radial case, there are infinitely many transmission eigenvalues, so it is reasonable to expect the same result here, but we have no results in this direction.

In summary, there are many questions remaining, some of which may be accessible by a further analysis of the quadratic forms of the operators introduced here. Because

⁷See [3] for an inequality relating the lowest transmission eigenvalue (if it exists) to the supremum of m and the diameter of D .

array imaging techniques are making the scattering operator, and hence its kernel, possible to measure, these questions are becoming increasingly relevant.

REFERENCES

- [1] S. AGMON, *Spectral properties of Schrödinger operators and scattering theory*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 2 (1975), pp. 151–218.
- [2] S. AGMON AND L. HÖRMANDER, *Asymptotic properties of solutions of differential equations with simple characteristics*, J. Anal. Math., 30 (1976), pp. 1–38.
- [3] F. CAKONI, D. COLTON, AND P. MONK, *On the use of transmission eigenvalues to estimate the index of refraction from far field data*, Inverse Problems, 23 (2007), pp. 507–522.
- [4] D. COLTON, A. KIRSCH, AND L. PÄIVÄRINTA, *Far-field patterns for acoustic waves in an inhomogeneous medium*, SIAM J. Math. Anal., 20 (1989), pp. 1472–1483.
- [5] D. COLTON AND R. KRESS, *Inverse Acoustic and Electromagnetic Scattering Theory*, 2nd ed., Appl. Math. Sci. 93, Springer-Verlag, Berlin, 1998.
- [6] D. COLTON AND P. MONK, *The inverse scattering problem for acoustic waves in an inhomogeneous medium*, in Inverse Problems in Partial Differential Equations (Arcata, CA, 1989), D. Colton, R. Ewing, and W. Rundell, eds., SIAM, Philadelphia, 1990, pp. 73–84.
- [7] D. COLTON, L. PÄIVÄRINTA, AND J. SYLVESTER, *The interior transmission problem*, Inverse Probl. Imaging, 1 (2007), pp. 13–28.
- [8] L. HÖRMANDER, *The Analysis of Linear Partial Differential Operators. II. Differential Operators with Constant Coefficients*, Classics in Mathematics, Springer-Verlag, Berlin, 2005.
- [9] J. R. McLAUGHLIN AND P. L. POLYAKOV, *On the uniqueness of a spherically symmetric speed of sound from transmission eigenvalues*, J. Differential Equations, 107 (1994), pp. 351–382.
- [10] J. R. McLAUGHLIN, P. L. POLYAKOV, AND P. E. SACKS, *Reconstruction of a spherically symmetric speed of sound*, SIAM J. Appl. Math., 54 (1994), pp. 1203–1223.
- [11] M. REED AND B. SIMON, *Methods of Modern Mathematical Physics. I. Functional Analysis*, 2nd ed., Academic Press, New York, 1980.
- [12] B. P. RYNNE AND B. D. SLEEMAN, *The interior transmission problem and inverse scattering from inhomogeneous media*, SIAM J. Math. Anal., 22 (1991), pp. 1755–1762.
- [13] B. SIMON, *Quantum Mechanics for Hamiltonians Defined as Quadratic Forms*, Princeton Series in Physics, Princeton University Press, Princeton, NJ, 1971.

PRESSURELESS EULER/EULER–POISSON SYSTEMS VIA ADHESION DYNAMICS AND SCALAR CONSERVATION LAWS*

TRUYEN NGUYEN[†] AND ADRIAN TUDORASCU[‡]

Abstract. The “sticky particles” model at the discrete level is employed to obtain global solutions for a class of systems of conservation laws among which lie the pressureless Euler and the pressureless attractive/repulsive Euler–Poisson system with zero background charge. We consider the case of finite, nonnegative initial Borel measures with finite second-order moment, along with continuous initial velocities of at most quadratic growth and finite energy. We prove the time regularity of the solution for the pressureless Euler system and obtain that the velocity satisfies the Oleinik entropy condition, which leads to a partial result on uniqueness. Our approach is motivated by earlier work of Brenier and Grenier, who showed that one-dimensional conservation laws with special initial conditions and fluxes are appropriate for studying the pressureless Euler system.

Key words. pressureless Euler, Euler–Poisson system, sticky particles, scalar conservation laws, Wasserstein distance, adhesion dynamics

AMS subject classifications. 35L65, 35L67, 82C40

DOI. 10.1137/070704459

1. Introduction. Let $\alpha, \beta \in \mathbb{R}$ and consider the system

$$(1) \quad \begin{cases} \partial_t \rho + \partial_x(\rho v) = 0, \\ \partial_t(\rho v) + \partial_x(\rho v^2) = \rho(\alpha \partial_x \Phi + \beta) \quad \text{in } \mathbb{R} \times (0, T), \\ \partial_{xx}^2 \Phi = \rho. \end{cases}$$

If $\alpha = \beta = 0$, then (1) describes the pressureless Euler system in spatial dimension one. The most commonly known form of the pressureless, attractive/repulsive Euler–Poisson system with zero background charge is also obtained from (1) by taking $\alpha = \pm 1$ and $\beta = 0$. In this paper, we are concerned with global existence of solutions for the initial value problem. Unlike the Euler with pressure case, the natural environment for the evolution is the space of nonnegative Borel measures on the real line. We consider the case of finite total mass, which we normalize to unity. The pressureless Euler ($\alpha = \beta = 0$) problem was studied using different techniques in [5], [6], [7], [9], [23], [15], [16], [17], [18], [20]. We point out that in these papers, generally, the velocity is taken to be at least bounded on the support of the initial measure. It appears that [23] and [20] are the only references (to our knowledge) which allow for unbounded velocities. Also, [23] is remarkable for dealing with the gravitational term as well ($\alpha = -1, \beta = 0$). In spite of that, the serious limitation of [23] is the assumptions that the initial velocity be sublinear growth and the initial mass distribution ρ_0 be either discrete or absolutely continuous with respect to the Lebesgue measure. Our main contribution is proving existence of global solutions for (1) if ρ_0 is just in $\mathcal{P}_2(\mathbb{R})$

*Received by the editors October 3, 2007; accepted for publication (in revised form) April 16, 2008; published electronically August 1, 2008. Both authors were supported in part by the School of Mathematics, Georgia Institute of Technology.

<http://www.siam.org/journals/sima/40-2/70445.html>

[†]Department of Theoretical and Applied Mathematics, University of Akron, Akron, OH 44325 (tnguyen@uakron.edu). This author gratefully acknowledges the postdoctoral support provided by NSF grant DMS-03-54729.

[‡]School of Mathematics, Georgia Institute of Technology, Atlanta, GA 30332 (adriant@math.gatech.edu).

and v_0 is continuous of at most quadratic growth and finite energy. As a consequence of an important result from [13], we also manage to show that the solution we obtain for Euler pressureless satisfies the Oleinik entropy condition which was conjectured in [7] and [23]. Note that similar constraints to ours on ρ_0 and v_0 were anticipated by Shnirelman [20] for the pressureless Euler system. However, our approach is more general than that of Shnirelman which cannot easily accommodate the α and β terms on the right-hand side of the momentum equation. The solutions constructed by Shnirelman were expressed in the form of a variational problem, which was shown in [3] to be equivalent to the variational principle considered in [23].

During the last decade, significant progress has been achieved in the study of partial differential equations in the context of optimal mass transportation. Much of the work on parabolic, dissipative equations was synthesized and placed in a very general setting in [2]. Much more recent and much less explored is the study of Hamiltonian systems in this context [1], [14]. The connection with the pressureless Euler system in arbitrary dimension was discovered by Benamou and Brenier [4], who showed that this system describes the geodesics in the Wasserstein space $\mathcal{P}_2(\mathbb{R})$. By definition, $\mathcal{P}_p(\mathbb{R})$ is the set of all Borel probabilities on \mathbb{R} with finite p -order moment. The set $\mathcal{P}_2(\mathbb{R})$ is endowed with the quadratic Wasserstein metric defined by

$$W_2^2(\mu, \nu) := \min_{\gamma} \int_{\mathbb{R}^2} |x - y|^2 d\gamma(x, y),$$

where the infimum is taken among all probabilities γ on the product space \mathbb{R}^2 with marginals μ, ν . The theory of absolutely continuous curves in $\mathcal{P}_2(\mathbb{R})$ [2] asserts the existence of velocities satisfying the conservation of mass equation in (1), regarded as a continuity equation. The left-hand side of the momentum equation can also be interpreted as the acceleration along the curve. We shall discuss these interesting connections at the end of this paper. A different version of the pressureless Euler-Poisson system was analyzed in [13] in the context of optimal mass transportation. The focus was on the two-point boundary problem, and existence and uniqueness for solutions as action-minimizing paths in $\mathcal{P}_2(\mathbb{R})$ were obtained.

We shall need the following assumptions:

- (H1) *The initial distribution of mass $\rho_0 \in \mathcal{P}_2(\mathbb{R})$.*
- (H2) *There exists $0 \leq \Lambda < +\infty$ such that*

$$v_0 \in C(\mathbb{R}) \cap L^2(\rho_0) \text{ and } |v_0(x)| \leq \Lambda(1 + x^2) \text{ for all } x \in \mathbb{R}.$$

The main objective is the following result.

THEOREM 1.1. *The initial-value problem for (1) admits a global weak solution in the sense of distributions if (H1) and (H2) hold.*

Two independent papers that appeared in 1996 and 1998 used adhesion dynamics to obtain global solutions for (1) in the $\alpha = \beta = 0$ case [7], [23] and in the $\alpha = -1, \beta = 0$ case [23]. Not only are we able to deal with the more general (1), but we also establish our results under less restrictive conditions on the initial mass distribution and velocity. In [7] the initial ρ_0 is compactly supported, while the initial velocity v_0 is continuous and bounded. These assumptions are relaxed in [23], e.g., $\sup_{|x| \leq R} |v_0(x)|/R \rightarrow 0$ as $R \rightarrow +\infty$, while $\text{spt}(\rho_0)$ may be unbounded, in which case $\int_0^x y d\rho_0(y) \rightarrow +\infty$ as $|x| \rightarrow +\infty$. As opposed to [7], however, [23] makes the extra assumptions that ρ_0 be either discrete or absolutely continuous with respect to the Lebesgue measure, in which case $\rho_0 > 0$ on $\text{spt}(\rho_0)$.

More recent work [16] treats the case $\alpha = \beta = 0$ for nonnegative Radon measures ρ_0 (not necessarily of finite total mass) and velocities $v_0 \in L^\infty(\rho_0)$. This paper is also remarkable in that it gives a necessary and sufficient condition for uniqueness of the Oleinik entropy solution—the initial weak continuity of the energy. As the example showing the necessity of this condition involves infinite mass initial measures, it remains unclear whether that is really needed in the finite mass case. Another paper that deals with possibly discontinuous (but bounded) initial velocities is [9], where the solution is produced constructively.

In [5] we find the concept of duality solutions, which is based on earlier work by the same authors. Existence and uniqueness are obtained under the assumption of atom-free initial density and bounded and continuous initial velocity. Boudin [6] obtains global existence of smooth solutions when the initial data has some higher regularity and is bounded away from zero and infinity (and thus of infinite total mass). Interestingly, the initial velocity does not have to be nondecreasing in order to rule out formation of singularities in finite time. We will consider this issue in section 4.3.

Whereas [7] and [23] use different approaches, they are still closely connected in principle. The fundamental underlying assumption for the discrete dynamics is the “sticky particle” hypothesis. The idea goes back to Zeldovich [24] and can be briefly described as follows. If m_i , $i = 1, n$ is a discrete system of masses initially located at $-\infty < x_1 < \dots < x_n < +\infty$ and moving with initial velocities v_i , $i = 1, n$, then one makes the assumption that the velocities remain constant while there is no collision. At the collision of a group of particles, the particles stick together and the initial velocity of the newly formed particle is given by the conservation of momentum. In [23] the authors successfully implemented a version adapted to the case $\alpha = -1$, $\beta = 0$. Instead of the constant speed intercollisional motion, we now assume uniformly accelerated motion between collisions. The acceleration is of gravitational nature and is proportional to the difference between the total mass to the left and the total mass to the right. Thus, in both cases, the trajectory of the i th particle before collision is given by

$$x_i(t) = x_i + tv_i + \frac{t^2}{2}a_n^i,$$

where

$$a_n^i = \begin{cases} 0 & \text{if } \alpha = \beta = 0, \\ \frac{1}{2} \left(\sum_{j < i} m_j - \sum_{j > i} m_j \right) & \text{if } \alpha = -1, \beta = 0. \end{cases}$$

(Here we convene that $m_0 = m_{n+1} = 0$.) If the masses m_j , $i \leq j \leq k$, collide at time $t_0 > 0$, then conservation of momentum yields

$$v_i(t_0+) = \frac{\sum_{j=i}^k m_j v_j(t_0-)}{\sum_{j=i}^k m_j}.$$

Of course, only finitely many collisions can occur; therefore, the evolution of the system is completely determined by the above assumptions.

Next we briefly describe the technique employed in [23], whose approach does not distinguish between the discrete and absolutely continuous cases. Here the problem is attacked from a “continuation of characteristics” point of view. When shocks occur, i.e., when the map $\phi_t(y) := y + tv_0(y) + t^2 a_0(y)/2$ is no longer invertible, one needs to

redefine $\phi_t(y)$ in such a way that it remains nondecreasing and $\phi_{t\#}\rho_0 =: \rho_t$ satisfies the equation in a weak sense. This redefinition uses the so-called generalized variational principle which comes from the intuition provided by the discrete case (see [23] for details).

A more elegant approach [7], in our opinion, makes use of standard results on approximations for scalar conservation laws. It is applied to the Euler pressureless system ($\alpha = \beta = 0$) and relies on the fact that the distribution function M of ρ satisfies an autonomous scalar conservation law. Another advantage lies in the fact that the solution for the continuous problem is obtained from the discrete ones via approximation theory for scalar conservation laws. We shall adopt this point of view and prove the more general Theorem 1.1 by an appropriate adaptation of Brenier and Grenier’s method.

The plan is as follows: we use (H1) to produce a sequence of discrete probabilities

$$\rho_0^n := \sum_{i=1}^n m_i^{(n)} \delta_{x_i^{(n)}} \rightarrow \rho_0 \text{ as } n \rightarrow \infty$$

in the 2-Wasserstein distance. We shall, in fact, prove that this sequence may be taken such that $\int_{\mathbb{R}} \zeta d\rho_0^n$ is uniformly bounded for some superquadratic growth function $\zeta : [0, \infty) \rightarrow [0, \infty)$. Denote by M_0^n the right continuous distribution function of ρ_0^n and let

$$(2) \quad a_n(m) := \alpha \left(\sum_{j=1}^i m_j - \frac{1}{2} m_i \right) + \beta$$

whenever $M_0^n(x_i^{(n)} -) \leq m < M_0^n(x_i^{(n)})$. We then define flux functions $\tilde{F}_n : [0, T] \times [0, 1] \rightarrow \mathbb{R}$ as

$$(3) \quad \tilde{F}_n(t, m) := \int_0^m f_n(\omega) d\omega + t \int_0^m a_n(\omega) d\omega,$$

where $N_0^n \# \chi_{(0,1)} = \rho_0^n$ optimally and $f_n := v_0 \circ N_0^n$ (the map N_0^n is taken to be right continuous). By adhesion dynamics we construct the unique entropy solution M_n for the first-order problem

$$(4) \quad \partial_t M + \partial_x [\tilde{F}_n(t, M)] = 0, \quad M(0, \cdot) = M_0^n.$$

We then use (H2) to show that, for an appropriate choice of approximating initial data, the sequence M_n will converge in some sense to the unique entropy solution for

$$(5) \quad \partial_t M + \partial_x [\tilde{F}(t, M)] = 0, \quad M(0, \cdot) = M_0,$$

where $\tilde{F} : [0, T] \times [0, 1] \rightarrow \mathbb{R}$ is given by

$$(6) \quad \tilde{F}(t, m) := \int_0^m f(\omega) d\omega + t \int_0^m a(\omega) d\omega = \int_0^m f(\omega) d\omega + t(\alpha m^2/2 + \beta m)$$

for $f := v_0 \circ N_0$ and $a(m) := \alpha m + \beta$ for $m \in [0, 1]$. Here $N_0 := M_0^{-1}$ (generalized inverse) [12] is the right continuous optimal map pushing $\chi_{(0,1)}$ forward to ρ_0 . The solution M of (5) will produce the solution $\rho := \partial_x M$ and $v\rho := \partial_x [\tilde{F}(t, M)]$ for (1) via a generalization of a result due to Volpert [22] on BV calculus. The last section is

dedicated to the $\alpha = 0 = \beta$ case. We give an explicit formula for v in terms of M , and, more importantly, we prove that our solution satisfies the Oleinik entropy condition. Some qualitative properties of the solution are discussed, e.g., time regularity and the impact of the initial velocity on the occurrence of spatial singularities. We finish with a partial result on uniqueness; i.e., we show that the energy of our solution for the pressureless Euler system is weakly continuous initially, which, along with the Oleinik entropy condition, leads to uniqueness in the case of bounded initial velocities [16].

2. Elements of one-dimensional BV calculus.

2.1. A BV chain rule in dimension one. To prove a chain rule for BV functions, we need the following lemma.

LEMMA 2.1. *Let μ be a Borel probability measure on \mathbb{R} and M be its right continuous distribution function. Write*

$$\mu = \sum_{j \in J} m_j \delta_{x_j} + \rho,$$

where $\{x_j\}_{j \in J}$ is the set (at most countable) of discontinuities of M and $m_j := \mu(\{x_j\})$. If ρ is nonzero, then we have

$$(7) \quad M_{\#}\rho = \chi_{U^c} \text{ for } U := \bigcup_{j \in J} (M(x_j-), M(x_j)).$$

Proof. We first observe that the balance of mass is satisfied. Since M is monotone nondecreasing, (7) is equivalent to the fact that M is the optimal map pushing ρ forward to χ_{U^c} . That is what we prove next. It is well known [12] that, since both measures are atom-free, the optimal map is given by $G^{-1} \circ F$, where F, G are the right continuous distribution functions of ρ and χ_{U^c} , respectively, and G^{-1} is the generalized inverse of G , given by $G^{-1}(y) = \inf\{m \in [0, 1] : G(m) > y\}$. We shall show that $M(x) = G^{-1} \circ F(x)$ for ρ -a.e. $x \in \text{spt}(\rho)$, i.e., $G^{-1}(F(x)) = \inf\{m \in [0, 1] : G(m) > F(x)\} = M(x)$. Note that

$$F(x) = M(x) - \sum_{x_j \leq x} m_j \text{ and } G(m) = m - \sum_{M(x_j) \leq m} m_j \text{ if } m \in U^c.$$

Thus, if $m \in U^c$,

$$G(m) - F(x) = \begin{cases} m - M(x) - \sum_{M(x) < M(x_j) \leq m} m_j & \text{if } m > M(x), \\ 0 & \text{if } m = M(x), \\ m - M(x) + \sum_{m < M(x_j) \leq M(x)} m_j & \text{if } m < M(x). \end{cases}$$

Since $G(M(x)) - F(x) = 0$ and G is a right continuous, nondecreasing function, all we need to prove is that there does not exist a nondegenerate interval $[M(x), M(x) + \epsilon]$ such that $G(m) = F(x)$ for all $m \in [M(x), M(x) + \epsilon]$. Suppose such an interval exists. We know M is right continuous at x ; therefore, if $x < z < x + \delta(\epsilon)$, we have that $M(x) \leq M(z) \leq M(x) + \epsilon$. Thus, as $M(z) \in U^c$, we infer that $M(z) - M(x) - \sum_{x < x_j \leq z} m_j = 0$, i.e., $\rho([x, z]) = 0$. This means $x \in \partial \text{spt}(\rho)$, and, by taking $[x, z_x]$, the maximal interval for which $\rho([x, z]) = 0$, we see that $(x, z_x) \cap (y, z_y) = \emptyset$ for any $x \neq y$ in $\partial \text{spt}(\rho)$ for which these nondegenerate intervals exist. This means that there are at most countably many such points, and, since ρ is atom-free, it follows that the ρ -measure of this set of points is zero. The lemma is thus proved. \square

The following theorem is fundamental. The first part of (i) can be trivially obtained from [22]; see, e.g., [5] for the exact formula on the derivative. The novelty of our result appears in (ii) as we go from the Lipschitz to the $W^{1,p}$ case. However, for the reader's convenience we sketch an elementary proof for (i) as well.

THEOREM 2.2. *Let μ be a Borel probability measure on \mathbb{R} and let M be its right continuous distribution function.*

(i) *Assume $f \in W^{1,\infty}(0,1)$. Then, $f \circ M \in BV(\mathbb{R})$ with distributional derivative $g\mu$, where*

$$(8) \quad g(x) = \begin{cases} \overline{f' \circ M}(x) & \text{if } \mu(\{x\}) = 0, \\ \frac{f \circ M(x) - f \circ M(x-)}{\mu(\{x\})} & \text{if } \mu(\{x\}) \neq 0, \end{cases}$$

μ -a.e. for some function $\overline{f' \circ M} \in L^\infty(\mathbb{R})$. Furthermore, suppose f' is defined unambiguously on $(0,1)$, and there exists a bounded $C^1[0,1]$ sequence f_n such that $f_n \rightarrow f$ uniformly and $f'_n \rightarrow f'$ everywhere on $(0,1)$. Then $\overline{f' \circ M} \equiv f' \circ M$ in the μ -a.e. sense, or, equivalently,

$$(9) \quad g(x) = \int_0^1 f'((1-s)M(x-) + sM(x)) ds \text{ for } \mu\text{-a.e. } x \in \mathbb{R}.$$

(ii) *Let $1 \leq p < +\infty$ and assume $f \in W^{1,p}(0,1) \cap C^1(0,1)$ and g is given by (8), with $\overline{f' \circ M} \equiv f' \circ M$ in the μ -a.e. sense. Then $g \in L^p(\mu)$ with $\|g\|_{L^p(\mu)} \leq \|f'\|_{L^p(0,1)}$ and $f \circ M \in BV(\mathbb{R})$ with distributional derivative $g\mu$.*

Now take $f \in W^{1,p}(0,1)$. Suppose f' is defined unambiguously on $(0,1)$, and there exists a $W^{1,p}(0,1) \cap C^1(0,1)$ sequence f_n such that $f_n \rightarrow f$ in $W^{1,p}$ and $f'_n \rightarrow f'$ everywhere on $(0,1)$. Then we still have the same result with $\overline{f' \circ M} \equiv f' \circ M$ in the μ -a.e. sense.

Proof. W.l.o.g. we may assume that μ is supported in some bounded interval I . Also, we shall first assume $f \in C^1[0,1]$ to show that (8) is valid with $\overline{f' \circ M} \equiv f' \circ M$. Now consider $\varphi \in C_c^\infty(I)$. We need to show that

$$(10) \quad - \int_I \varphi'(x) f \circ M(x) dx = \int_I \varphi(x) g(x) d\mu(x).$$

If μ has finitely many atoms, then the validity of (10) can be checked by direct computation. Indeed, since M is piecewise continuous and bounded, we may approximate it uniformly by nondecreasing piecewise $W^{1,1}$ functions M_ϵ (may take M_ϵ piecewise linear and continuous on each continuity interval for M , such that M and M_ϵ agree at the endpoints). Thus, $\mu_\epsilon := M'_\epsilon - M' = \mu$ weakly as nonnegative, bounded measures. The chain rule for Sobolev functions [8] applies piecewise and yields (10) for μ_ϵ . Then we pass to the limit to obtain the result for μ . Thus, let us assume $D := \{x_1, x_2, \dots, x_n, \dots\}$ is the infinite set of all atoms of μ and write

$$\mu = \sum_{i=1}^\infty m_i \delta_{x_i} + \rho, \text{ where } m_i > 0, i = 1, 2, \dots,$$

and ρ is an atom-free nonnegative Borel measure of total mass $1 - \sum_{i=1}^\infty m_i \geq 0$. We shall call the atomic measure the singular part while ρ shall be called the regular part (although it may not be absolutely continuous with respect to the Lebesgue

measure). Consider now the sequence of measures μ_n given by $\mu_n = \sum_{i=1}^n m_i \delta_{x_i} + \rho$. Of course, $\mu_n \rightharpoonup \mu$ weakly \star as measures. Since M_n , the right continuous distribution function of μ_n , has only finitely many discontinuities, (8) holds for μ_n , as proved above. It is easy to see that $M_n \rightarrow M$ Lebesgue a.e.; thus, the continuity of f along with its boundedness gives the convergence of the left-hand side of (10) by dominated convergence. Therefore, $f \circ M_n \in BV(\mathbb{R})$ with distributional derivative $g_n \mu_n$, where

$$g_n(x) := \begin{cases} f' \circ M_n(x) & \text{if } \mu_n(\{x\}) = 0, \\ \frac{f \circ M_n(x) - f \circ M_n(x-)}{\mu_n(\{x\})} & \text{if } \mu_n(\{x\}) \neq 0. \end{cases}$$

Now we write

$$\int_{\mathbb{R}} \varphi(x) g_n(x) d\mu_n(x) = \sum_{i=1}^n \varphi(x_i) [f(M_n(x_i)) - f(M_n(x_i-))] + \int_{\mathbb{R}} \varphi(x) f' \circ M_n(x) d\rho(x).$$

By the continuity of f' we obtain the convergence of the second term on the right-hand side. Thus, it remains to prove that

$$\sum_{i=1}^n \varphi(x_i) [f(M_n(x_i)) - f(M_n(x_i-))] \rightarrow \sum_{i=1}^{\infty} \varphi(x_i) [f(M(x_i)) - f(M(x_i-))],$$

which can be obtained after some calculations as a consequence of the convergence of the series $\sum m_i$ if $f \in C^2[0, 1]$, and then for $C^1[0, 1]$ functions by approximation. If f is $W^{1,\infty}(0, 1)$, then we conclude by taking a sequence of functions in $C^1[0, 1]$ such that $f_n \rightarrow f$ uniformly and f'_n are uniformly bounded. Indeed, one has

$$-\int_{\mathbb{R}} \varphi' f_n \circ M dx = \int_{\mathbb{R}} \varphi f'_n \circ M d\rho + \int_{\mathbb{R}} \varphi(x) \frac{f_n(M(x)) - f_n(M(x-))}{M(x) - M(x-)} d\mu_s(x),$$

where $\mu_s = \sum m_j \delta_{x_j}$ is the singular part of μ . Note that the ratio in the second term of the right-hand side converges uniformly to $[f(M(x)) - f(M(x-))]/[M(x) - M(x-)]$ on the support of μ_s ; thus the uniform bound on f'_n ensures the convergence of the integral by dominated convergence. Since the left-hand side is trivially convergent, it follows that

$$\int_{\mathbb{R}} \varphi(x) f'_n \circ M(x) d\rho(x) \text{ converges as } n \rightarrow \infty,$$

which, along with the uniform bound on f'_n , yields the convergence of $\frac{f'_n \circ M}{M}$ in the L^∞ weak \star topology. We also deduce that the limit, denoted by $\overline{f' \circ M}$, is μ -a.e. independent of the chosen sequence f_n . The second statement from (i) easily follows (by dominated convergence) from the fact that $g_n \rightarrow g$ everywhere as an L^∞ bounded sequence.

To prove the first part of (ii) we truncate f' by

$$(11) \quad f'_n(x) := \begin{cases} -n & \text{if } f'(x) < -n, \\ f'(x) & \text{if } |f'(x)| \leq n, \\ n & \text{if } f'(x) > n \end{cases}$$

and let f_n be the antiderivative of f'_n vanishing at zero. Note that $|f'_n| \leq |f'|$ on $(0, 1)$, which implies that f_n is uniformly bounded with respect to n in $L^p(0, 1)$. First,

since $f_n \in C^1[0, 1]$, we infer, according to (i), that $f_n \circ M$ is BV with distributional derivative $g_n\mu$, where

$$(12) \quad g_n(x) := \begin{cases} f'_n \circ M(x) & \text{if } \mu(\{x\}) = 0, \\ \frac{f_n \circ M(x) - f_n \circ M(x-)}{\mu(\{x\})} & \text{if } \mu(\{x\}) \neq 0. \end{cases}$$

This (see (9) for the equivalent integral expression) together with the fact that $|f'_n| \leq |f'|$ implies

$$|g_n(x)| \leq \int_0^1 |f'((1-s)M(x-) + sM(x))| ds =: h(x) \text{ for } x \in I.$$

Since $f'_n \rightarrow f'$, we infer that $g_n \rightarrow g$ pointwise. Assume that $h \in L^p(\mu)$. Then, we may pass to the limit in the right-hand side of

$$-\int_I \varphi'(x) f_n \circ M(x) dx = \int_I \varphi(x) g_n(x) d\mu(x).$$

The left-hand side converges to the appropriate quantity because $f_n \rightarrow f$ uniformly and f_n are uniformly bounded.

Thus, we are done if we can prove that $h \in L^p(\mu)$. For this let $\{x_j\}_{j \in J}$ be the set of discontinuities of M and $m_j := \mu(\{x_j\})$. Now consider, for $s \in [0, 1]$, the sum

$$\sum_{i=1}^n m_i |f'((1-s)M(x_i-) + sM(x_i))|^p = \sum_{i=1}^n m_i |f'(M(x_i-) + sm_i)|^p.$$

By monotone convergence

$$\sum_{i=1}^n \int_0^1 m_i |f'(M(x_i-) + sm_i)|^p ds \rightarrow \int_0^1 \sum_{i=1}^\infty m_i |f'(M(x_i-) + sm_i)|^p ds,$$

which, after obvious linear changes of variables, is equivalent to

$$\int_{\cup_{i=1}^n [M(x_i-), M(x_i)]} |f'(m)|^p dm \rightarrow \int_0^1 \sum_{i=1}^\infty m_i |f'(M(x_i-) + sm_i)|^p ds.$$

Again, by monotone convergence the left-hand side converges to $\int_U |f'(m)|^p dm$, where $U := \cup_{j \in J} (M(x_j-), M(x_j))$. Thus, since $f' \in L^p(0, 1)$ and $U \subset [0, 1]$, it follows that

$$\int_0^1 \sum_{i=1}^\infty m_i |f'(M(x_i-) + sm_i)|^p ds = \int_U |f'(m)|^p dm \leq \|f'\|_{L^p(0,1)}^p,$$

so we can apply Fubini's theorem to obtain

$$(13) \quad \int_I \int_0^1 |f'((1-s)M(x-) + sM(x))|^p ds d\mu_s(x) = \int_U |f'(m)|^p dm < +\infty,$$

where μ_s denotes, as before, the singular part of μ (unrelated to the integration variable s). If $\mu = \mu_s$, then we are done. Else, by using Lemma 2.1, we obtain

$$\int_I |f' \circ M(x)|^p d\rho(x) = \int_{U^c} |f'(m)|^p dm,$$

which, combined with (13), yields

$$(14) \quad \int_I \int_0^1 |f'((1-s)M(x-) + sM(x))|^p ds d\mu(x) = \|f'\|_{L^p(0,1)}^p.$$

Therefore, $h \in L^p(\mu)$ with $\|h\|_{L^p(\mu)} \leq \|f'\|_{L^p(0,1)}$ (with equality if $p = 1$). Since $|g| \leq h$ pointwise, the proof of the first part of (ii) is concluded.

Consequently,

$$(15) \quad - \int_{\mathbb{R}} \varphi' f_n \circ M dx = \int_{\mathbb{R}} \varphi g_n d\mu \text{ for all positive integers } n,$$

where g_n is defined by (12) for the approximating sequence f_n considered in the second part of (ii). We have, just as in deducing (14), that

$$\begin{aligned} \int_I |g_n(x) - g_m(x)|^p d\mu(x) &\leq \int_I \int_0^1 |(f'_n - f'_m)((1-s)M(x-) + sM(x))|^p ds d\mu(x) \\ &= \|f'_n - f'_m\|_{L^p(0,1)}^p \end{aligned}$$

for all natural $m, n \geq 1$. Thus, $\{g_n\}$ is convergent in $L^p(\mu)$, and, due to the hypothesis of *everywhere* convergence of f'_n to f' , we obtain $g_n \rightarrow g$ in $L^p(\mu)$, where

$$g(x) := \int_0^1 f'((1-s)M(x-) + sM(x)) ds.$$

Passing to the limit in (15) concludes our proof. \square

The following corollary holds for any $1 \leq p \leq +\infty$.

COROLLARY 2.3. *Let $f \in L^p(0, 1)$ be right continuous (thus unambiguously defined everywhere in $(0, 1)$), and take F to be its antiderivative vanishing at zero. If M is the cumulative distribution function of some Borel probability measure μ on \mathbb{R} , then $F \circ M \in BV(\mathbb{R})$ with distributional derivative $g\mu$, where*

$$(16) \quad g(x) := \begin{cases} f \circ M(x) & \text{if } \mu(\{x\}) = 0, \\ \frac{F \circ M(x) - F \circ M(x-)}{\mu(\{x\})} & \text{if } \mu(\{x\}) \neq 0, \end{cases}$$

in the μ -a.e. sense. Furthermore, $g \in L^p(\mu)$ and $\|g\|_{L^p(\mu)} \leq \|f\|_{L^p(0,1)}$.

The proof is an immediate consequence of Theorem 2.2. Indeed, we take (upon extending f by zero outside $(0, 1)$) the function $f_n := \eta^n * f$, where η^n is obtained from the standard mollifier supported in $[-1/n, 1/n]$ by shifting it to the left by $1/n$. The classic properties of mollification still hold, e.g., $f_n \rightarrow f$ in $L^p(0, 1)$ (if $p \neq +\infty$) and $f_n \in C^\infty[0, 1]$. However, the interesting feature of these “shifted mollifiers” is that the right continuity of f on $(0, 1)$ is enough to easily prove that $f_n \rightarrow f$ *everywhere* in $(0, 1)$. Thus, if we take $F_n(m) := \int_0^m f_n(\omega) d\omega$, we are within the hypotheses of Theorem 2.2.

Remark 2.4. To understand the relevance of this result, observe that the functions f_n, f defined in the introduction are right continuous on $(0, 1)$. Indeed, that comes as a consequence of the continuity of v_0 on \mathbb{R} and the right continuity of N_0^n, N_0 on $(0, 1)$ (as generalized inverses of nondecreasing, right continuous functions [21]).

2.2. A two-dimensional extension. Now let us assume $M : [0, T] \times \mathbb{R} \rightarrow [0, 1]$ for some $0 < T < +\infty$, such that $M \in BV([0, T] \times \mathbb{R})$ and $M(t, \cdot)$ is a right continuous probability distribution function for Lebesgue a.e. $t \in (0, T)$. The following lemma is easy to check as a consequence of Fubini's theorem.

LEMMA 2.5. *Let ∇M be the vector-valued measure given by the BV gradient $\nabla = (\partial_t, \partial_x)$ of M . Then, its x -component $d\partial_x M$ admits the decomposition $d\partial_x M(t, \cdot)dt$.*

We shall use this to prove the main result of this section, which we state below. We denote by C_r the space of right continuous functions.

THEOREM 2.6. *Consider $M : [0, T] \times \mathbb{R} \rightarrow [0, 1]$ as above, and let $f \in L^2(0, 1) \cap C_r(0, 1)$ and F be its antiderivative vanishing at zero. Assume further that*

$$(17) \quad \partial_t M + \partial_x [F \circ M] = 0 \text{ in } \mathcal{D}'([0, T] \times \mathbb{R}).$$

Then $F \circ M \in BV([0, T] \times \mathbb{R})$ and $\nabla [F \circ M] = g\nabla M$ for some $g \in L^1(|\nabla M|)$.

Proof. We may consider w.l.o.g. the case $f \in C(0, 1)$. Else, we simply use Corollary 2.3 instead of Theorem 2.2. Let us start by taking the truncations f_n for f as in (11). The corresponding F_n 's are in $C^1[0, 1]$, $F_n \rightarrow F$ uniformly and $f_n \rightarrow f$ in $L^2(0, 1)$. According to [22], $F_n \circ M \in BV([0, T] \times \mathbb{R})$ and there exist $\bar{f}_n \in L^\infty$ such that $\nabla [F_n \circ M] = \bar{f}_n \nabla M$ as vector-valued measures. This means that

$$(18) \quad - \int_0^T \int_{\mathbb{R}} F_n(M) \nabla \cdot \phi dx dt = \int_0^T \int_{\mathbb{R}} \bar{f}_n \phi_1 d\partial_t M + \int_0^T \int_{\mathbb{R}} \bar{f}_n \phi_2 d\partial_x M$$

for all $\phi := (\phi_1, \phi_2) \in C_c^\infty((0, T) \times \mathbb{R}; \mathbb{R}^2)$. We now use Lemma 2.5 and (17) to infer

$$(19) \quad - \int_0^T \int_{\mathbb{R}} F_n(M) \nabla \cdot \phi dx dt = \int_0^T \int_{\mathbb{R}} (-g\phi_1 + \phi_2) \bar{f}_n d\partial_x M(t, \cdot) dt,$$

where, according to Theorem 2.2 (ii),

$$g(t, \cdot) := \int_0^1 f((1-s)M(t, \cdot -) + sM(t, \cdot)) ds \in L^2(\partial_x M(t, \cdot))$$

for Lebesgue a.e. $t \in (0, T)$. Also, Theorem 2.2 (ii) ensures that

$$\|g(t, \cdot)\|_{L^2(\partial_x M(t, \cdot))} \leq \|f\|_{L^2(0,1)} \text{ uniformly with respect to } t.$$

Note that (19) implies

$$(20) \quad - \int_0^T \int_{\mathbb{R}} F_n(M) \partial_x \phi_1 dx dt = \int_0^T \int_{\mathbb{R}} \bar{f}_n \phi_1 d\partial_x M(t, \cdot) dt$$

for all $\phi_1 \in C_c^\infty((0, T) \times \mathbb{R})$. For Lebesgue a.e. $t \in (0, T)$ one has (Theorem 2.2 (i))

$$\partial_x [F_n \circ M(t, \cdot)] = g_n(t, \cdot) \partial_x M(t, \cdot) \text{ for } g_n(t, x) := \int_0^1 f_n((1-s)M(t, x-) + sM(t, x)) ds$$

and

$$\|g_n(t, \cdot)\|_{L^2(\partial_x M(t, \cdot))} \leq \|f\|_{L^2(0,1)} \text{ uniformly in } t \text{ and } n.$$

Along with (20), this implies $\bar{f}_n \equiv g_n$ in the $d\partial_x M(t, \cdot)dt$ -a.e. sense. Therefore, (19) is equivalent to

$$(21) \quad - \int_0^T \int_{\mathbb{R}} F_n(M) \nabla \cdot \phi dx dt = \int_0^T \int_{\mathbb{R}} (-g\phi_1 + \phi_2) g_n d\partial_x M(t, \cdot) dt.$$

Furthermore, due to the pointwise convergence and the uniform (in n and t) $L^2(\partial_x M)$ -bounds we obtain $g_n \rightarrow g$ in $L^2(d\partial_x M(t, \cdot)dt)$. By the uniform (and bounded) convergence of F_n to F we can also pass to the limit in the left-hand side. Thus,

$$-\int_0^T \int_{\mathbb{R}} F(M) \nabla \cdot \phi dx dt = \int_0^T \int_{\mathbb{R}} (-g^2 \phi_1 + g \phi_2) d\partial_x M(t, \cdot) dt,$$

which, after using Lemma 2.5 and (17) once more, leads to

$$-\int_0^T \int_{\mathbb{R}} F(M) \nabla \cdot \phi dx dt = \int_0^T \int_{\mathbb{R}} g \phi_1 d\partial_t M + \int_0^T \int_{\mathbb{R}} g \phi_2 d\partial_x M.$$

Since $g \in L^2(\partial_x M) \subset L^1(\partial_x M)$ and $g \in L^1(|\partial_t M|)$, we obtain the result. \square

We are going to need a slightly different result which we now state as a consequence. It may be proved by retracing the proof of Theorem 2.6.

COROLLARY 2.7. *Assume now that*

$$\partial_t M + \partial_x [\tilde{F}(t, M)] = 0 \text{ in } \mathcal{D}'([0, T] \times \mathbb{R})$$

for some time-linear perturbation $\tilde{F}(t, \cdot) := F + t\Psi$, $\Psi \in C^1[0, 1]$. Then $F \circ M \in BV([0, T] \times \mathbb{R})$ and

$$\nabla[\tilde{F}(t, M)] = (g + t\psi)\nabla M + \mathbf{V}$$

for g from Theorem 2.6, and the vector field $\mathbf{V} := (\Psi \circ M, 0)$. Here,

$$\psi(t, x) := \int_0^1 \Psi'((1-s)M(t, x-) + sM(t, x)) ds.$$

3. Pressureless Euler/Euler–Poisson systems via scalar conservation laws.

3.1. Convergence of measures and their distribution functions. Let $1 \leq p < \infty$, $\mu \in \mathcal{P}_p(\mathbb{R})$, and $v \in L^p(\mu) \cap C(\mathbb{R})$. Then by the de la Vallée-Poussin lemma which can be found in [10], there exists a nonnegative, convex, increasing function $\zeta \in C^1([0, +\infty))$ satisfying $\zeta(0) = 0$ and $\frac{\zeta(t)}{t} \uparrow +\infty$ as $t \rightarrow +\infty$ such that

$$\int_{\mathbb{R}} \zeta(|x|^p + |v(x)|^p) d\mu(x) < +\infty.$$

It is also well known that (see [11], for example) there exist a probability space (Ω, Σ, P) and a sequence of independent random variables $\xi_i : \Omega \rightarrow \mathbb{R}$ such that

$$\xi_{i\#} P = \mu.$$

Now for each positive integer n and each $\omega \in \Omega$, define

$$\mu^{n, \omega} := \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i(\omega)}.$$

Then by the strong law of large numbers and the separability of $C_c(\mathbb{R})$ we have for P -a.e. $\omega \in \Omega$

$$\int_{\mathbb{R}} f(x) d\mu^{n, \omega}(x) = \frac{1}{n} \sum_{i=1}^n f(\xi_i(\omega)) \rightarrow \mathbb{E}(f \circ \xi_1) = \int_{\mathbb{R}} f(x) d\mu(x)$$

for all functions $f \in C_c(\mathbb{R})$. Consequently, $\mu^{n,\omega}$ converges narrowly to μ . Thus, by using the strong law of large numbers again, we also have that for P -a.e. $\omega \in \Omega$

$$\int_{\mathbb{R}} \zeta(|x|^p) d\mu^{n,\omega}(x) \rightarrow \int_{\mathbb{R}} \zeta(|x|^p) d\mu(x), \quad \int_{\mathbb{R}} \zeta(|v|^p) d\mu^{n,\omega}(x) \rightarrow \int_{\mathbb{R}} \zeta(|v|^p) d\mu(x).$$

The last fact together with the properties of ζ yields for $k > 0$ large enough

$$\frac{\zeta(k^p)}{k^p} \int_{\{x:|x|\geq k\}} |x|^p d\mu^{n,\omega} \leq \int_{\{x:|x|\geq k\}} \frac{\zeta(|x|^p)}{|x|^p} |x|^p d\mu^{n,\omega} \leq \int_{\mathbb{R}} \zeta(|x|^p) d\mu^{n,\omega} \leq C(\omega).$$

Therefore, $\{\mu^{n,\omega}\}$ has uniformly integrable p -moments. Consequently, by Proposition 7.1.5 in [2] we obtain that for P -a.e. $\omega \in \Omega$, $\mu^{n,\omega} \rightarrow \mu$ in W_p and

$$\int_{\mathbb{R}} \zeta(|x|^p) d\mu^{n,\omega}(x) \rightarrow \int_{\mathbb{R}} \zeta(|x|^p) d\mu(x), \quad \int_{\mathbb{R}} \zeta(|v|^p) d\mu^{n,\omega}(x) \rightarrow \int_{\mathbb{R}} \zeta(|v|^p) d\mu(x).$$

LEMMA 3.1. *Suppose $1 \leq p < \infty$ and $\{\mu_n\}$ is a sequence of measures in $\mathcal{P}_p(\mathbb{R})$ converging to $\mu \in \mathcal{P}_p(\mathbb{R})$ in W_p . When $p > 1$ we assume further that $\int_{\mathbb{R}} \zeta(|x|^p) d\mu_n(x)$ is uniformly bounded in n for some nonnegative, convex, increasing function $\zeta \in C^1([0, +\infty))$ satisfying $\zeta(0) = 0$ and $\frac{\zeta(t)}{t} \uparrow +\infty$ as $t \rightarrow +\infty$. Then we have the following:*

(i) *The L^1 norm of $M_n - M$ on \mathbb{R} goes to zero, where*

$$M(x) := \mu((-\infty, x]).$$

(ii) *For any nondecreasing C^1 function B on $[0, 1]$ such that $B(0) = 0, B(1) = 1$,*

$$\partial_x(B(M_n)) \rightarrow \partial_x(B(M)) \quad \text{in } W_p.$$

Proof. From Theorem 6.0.2 in [2], we obtain

$$W_1(\mu_n, \mu) = \int_0^1 |M_n^{-1}(s) - M^{-1}(s)| ds,$$

where the generalized inverse M^{-1} of M is defined by

$$M^{-1}(s) := \inf \{x \in \mathbb{R} : M(x) > s\}, \quad s \in [0, 1].$$

Then by using Fubini's theorem we get $W_1(\mu_n, \mu) = \|M_n - M\|_{L^1(\mathbb{R})}$. This together with the fact that $W_1(\mu_n, \mu) \leq W_p(\mu_n, \mu)$ gives (i). By (i) and the Lipschitz condition on B we clearly have $\|B(M_n) - B(M)\|_{L^1(\mathbb{R})} \rightarrow 0$. But as

$$W_1(\partial_x(B(M_n)), \partial_x(B(M))) = \|B(M_n) - B(M)\|_{L^1(\mathbb{R})},$$

we infer in particular that $\partial_x(B(M_n)) \rightarrow \partial_x(B(M))$ narrowly.

We have from Theorem 2.2

$$\int_{\mathbb{R}} \zeta(|x|^p) d\partial_x(B(M_n)) = \int_{\mathbb{R}} \zeta(|x|^p) g_n(x) d\mu_n(x),$$

where $g_n(x) = \int_0^1 B'(sM_n(x) + (1-s)M_n(x-)) ds$. Therefore,

$$\int_{\mathbb{R}} \zeta(|x|^p) d\partial_x(B(M_n)) \leq \|B'\|_{\infty} \int_{\mathbb{R}} \zeta(|x|^p) d\mu_n(x) \leq C \|B'\|_{\infty}.$$

It follows from this and the properties of the function ζ that for any $k > 0$ large enough,

$$\int_{\{x:|x|\geq k\}} |x|^p d\partial_x(B(M_n)) \leq \frac{k^p}{\zeta(k^p)} C \|B'\|_\infty.$$

That is, the sequence of probability measures $\partial_x(B(M_n))$ has uniformly integrable p -moments. Consequently, we can conclude from Proposition 7.1.5 in [2] that

$$\partial_x(B(M_n)) \rightarrow \partial_x(B(M))$$

in W_p , as desired. \square

3.2. Convergence of the discrete problem. Following [7], we take a discrete probability measure

$$\rho_0^n := \sum_{j=1}^n m_j \delta_{x_j}, \quad x_1 < x_2 < \dots < x_n,$$

and define

$$(22) \quad \rho^n(t, x) := \sum_{j=1}^n m_j \delta_{x_j(t)},$$

where the characteristics are given by

$$(23) \quad x_j(t) = x_j + tv_j + \frac{t^2}{2} a_n(M_n(t, x_j(t)-)).$$

Here, as in the introduction, we have

$$a_n(M_n(t, x_j(t)-)) = \left[\alpha \left(\sum_{i=1}^j m_i - \frac{1}{2} m_j \right) + \beta \right].$$

We impose the adhesion dynamics at collision (see the introduction) and consider

$$(24) \quad M_n(t, x) := \sum_{j=1}^n m_j H(x - x_j(t)),$$

where H is the right continuous Heaviside function. Since M_n is piecewise constant, we need only show that M_n solves (4) by checking the Rankine–Hugoniot jump conditions across the shocks $x = x_j(t)$, i.e.,

$$(25) \quad \dot{x}_j(t) = \frac{\tilde{F}_n(t, M_n(t, x_j(t))) - \tilde{F}_n(t, M_n(t, x_j(t)-))}{M_n(t, x_j(t)) - M_n(t, x_j(t)-)}, \quad j = 1, \dots, n.$$

Assume the masses m_i for $j_0 \leq i \leq j_1$, including m_j , are all amassed at time t . Then, exactly as in [7], we have

$$v_j(t) = \frac{F_n(M_n(t, x_j(t))) - F_n(M_n(t, x_j(t)-))}{M_n(t, x_j(t)) - M_n(t, x_j(t)-)}, \quad j = 1, \dots, n,$$

where F_n is defined by

$$F_n(m) := \int_0^m f_n(\omega) d\omega \quad \text{for } m \in [0, 1].$$

Observe that $\tilde{F}_n(t, m) = F_n(m) + t \int_0^m a_n(\omega) d\omega$ by the definition of \tilde{F}_n in the introduction. The integrand below is constantly a_n^j on each interval of the form $[M_n(t, x_j(t)-), M_n(t, x_j(t))]$; thus,

$$a_n(M_n(t, x_j(t)-)) \sum_{i=j_0}^{j_1} m_i = \int_{M_n(t, x_j(t)-)}^{M_n(t, x_j(t))} a_n(\omega) d\omega$$

implies (25). To prove that M_n is an entropy solution for (4), we check the entropy inequality

$$(26) \quad \dot{x}_j(t) \leq \frac{\tilde{F}_n(t, X) - \tilde{F}_n(t, M_n(t, x_j(t)-))}{X - M_n(t, x_j(t)-)},$$

where $X = \sum_{i \leq k} m_i$, for some $j_0 \leq k \leq j_1$. The inequality

$$v_j(t) \leq \frac{F_n(X) - F_n(M_n(t, x_j(t)-))}{X - M_n(t, x_j(t)-)}$$

is justified in [7] as a consequence of the *barycentric lemma* (which simply formulates the fact that if two groups of particles collide, then the averaged velocity of the group to the left decreases). Then we see that

$$a_n(M_n(t, x_j(t)-)) \sum_{i=j_0}^k m_i = \int_{M_n(t, x_j(t)-)}^X a_n(\omega) d\omega,$$

which, together with the previous inequality, yields (26). We have just sketched the proof of the following proposition.

PROPOSITION 3.2. *The function M_n given by (24) is the entropy solution of the problem (4).*

Next we want to show that M_n converges in some sense to an entropy solution of (5). The following proposition applies if the initial approximating measures are of the form

$$\rho_0^n := \frac{1}{n} \sum_{i=1}^n \delta_{x_i^{(n)}}.$$

Note that section 3.1 shows that we may consider such approximations.

PROPOSITION 3.3. *Let $\rho_0 \in \mathcal{P}_2(\mathbb{R})$ and consider a sequence of discrete probabilities ρ_0^n as above such that $W_2(\rho_0, \rho_0^n) \rightarrow 0$ and*

$$(27) \quad \int_{\mathbb{R}} \zeta(v_0^2) d\rho_0^n \rightarrow \int_{\mathbb{R}} \zeta(v_0^2) d\rho_0 < +\infty, \quad \int_{\mathbb{R}} \zeta(x^2) d\rho_0^n \rightarrow \int_{\mathbb{R}} \zeta(x^2) d\rho_0 < +\infty$$

for some nonnegative, convex, increasing function $\zeta \in C^1([0, +\infty))$ satisfying $\zeta(0) = 0$ and $\frac{\zeta(t)}{t} \uparrow +\infty$ as $t \rightarrow +\infty$. Let M_0^n and M_0 be the right continuous cumulative distribution functions of ρ_0^n and ρ_0 , respectively. Consider, as above, the entropy

solutions M_n of (4) for all n . Then there exists a Borel function $M : [0, \infty) \times \mathbb{R} \rightarrow [0, 1]$ such that, for any given $T > 0$,

$$(28) \quad \max_{0 \leq t \leq T} W_2(\partial_x M_n(t, \cdot), \partial_x M(t, \cdot)) \rightarrow 0 \text{ and } \max_{0 \leq t \leq T} \|M_n(t, \cdot) - M(t, \cdot)\|_{L^1(\mathbb{R})} \rightarrow 0.$$

Moreover, M is the entropy solution of the problem (5).

Proof. First recall that

$$\rho^n(t, x) := \frac{1}{n} \sum_{j=1}^n \delta_{x_j^{(n)}(t)},$$

where

$$x_j^{(n)}(t) = x_j^{(n)} + tv_j^{(n)} + \frac{t^2}{2} \left(\alpha \frac{2j-1}{2n} + \beta \right).$$

Let $\zeta_t(x) := \zeta(\frac{1}{\kappa}x)$, where κ is some positive constant depending only on t which will be determined later. Then, by the properties of ζ , we have

$$\begin{aligned} & \int_{\mathbb{R}} \zeta_t(|x|^2) d\rho_n(t, x) \\ &= \frac{1}{n} \sum_{j=1}^n \zeta \left(\frac{1}{\kappa} |x_j^{(n)}(t)|^2 \right) \leq \frac{1}{n} \sum_{j=1}^n \zeta \left(\frac{3}{\kappa} \left[|x_j^{(n)}|^2 + t^2 |v_j^{(n)}|^2 + \frac{(|\alpha| + |\beta|)^2 t^4}{4} \right] \right) \\ &\leq \frac{1}{3} \zeta \left(\frac{9(|\alpha| + |\beta|)^2 t^4}{4\kappa} \right) + \frac{1}{3n} \sum_{j=1}^n \zeta \left(\frac{9}{\kappa} |x_j^{(n)}|^2 \right) + \frac{1}{3n} \sum_{j=1}^n \zeta \left(\frac{9t^2}{\kappa} |v_j^{(n)}|^2 \right). \end{aligned}$$

Hence, by choosing $\kappa = \max\{9, 9t^2\}$ and using again the fact that ζ is increasing, we obtain

$$\int_{\mathbb{R}} \zeta_t(|x|^2) d\rho_n(t, x) \leq C_t + \frac{1}{3} \int_{\mathbb{R}} \zeta(|x|^2) d\rho_0^n(t, x) + \frac{1}{3} \int_{\mathbb{R}} \zeta(|v_0|^2) d\rho_0^n(t, x) \leq C_t$$

uniformly in n (the last inequality is due to (27)). Since ζ_t has superlinear growth (it is just an argument-rescaled version of ζ), there exists $\rho(t, \cdot) \in \mathcal{P}_2(\mathbb{R})$ such that, up to a subsequence that may depend on t , $W_2(\rho^n(t, \cdot), \rho(t, \cdot)) \rightarrow 0$ as $n \rightarrow \infty$. By a standard diagonal argument we can choose a subsequence independent of t satisfying $W_2(\rho^n(t, \cdot), \rho(t, \cdot)) \rightarrow 0$ for all $t \in [0, \infty) \cap \mathbb{Q}$. In order to see that this conclusion also holds for all t in $[0, \infty)$, we are going to show that the paths $\rho^n(t, \cdot)$ are uniformly locally Lipschitz in t . Indeed, let $T > 0$ and $t, s \in [0, T]$ be arbitrary. Then, since

$$|x_j^{(n)}(t) - x_j^{(n)}(s)|^2 \leq C|t - s|^2 \left(|v_j^{(n)}|^2 + T^2 \right),$$

we have

$$\begin{aligned} W_2^2(\rho^n(t, \cdot), \rho^n(s, \cdot)) &\leq \frac{1}{n} \sum_{j=1}^n |x_j^{(n)}(t) - x_j^{(n)}(s)|^2 \leq C|t - s|^2 \left(\frac{1}{n} \sum_{j=1}^n |v_j^{(n)}|^2 + T^2 \right) \\ &= C|t - s|^2 \left(\int_{\mathbb{R}} |v_0(x)|^2 d\rho_0^n + T^2 \right) \leq C|t - s|^2 \quad \text{uniformly in } n. \end{aligned}$$

Using this uniformly Lipschitz property, we can conclude that, in fact,

$$(29) \quad W_2(\rho^n(t, \cdot), \rho(t, \cdot)) \rightarrow 0 \text{ as } n \rightarrow \infty, \text{ uniformly in } t \in [0, T],$$

which yields (28) with $M(t, x) := \rho(t, (-\infty, x])$.

Next we show that M is a solution of (5). By the assumptions and Lemma 3.1, we have

$$\int_{\mathbb{R}} |M_0^n - M_0| dx \rightarrow 0$$

and

$$\partial_x(B(M_0^n)) \rightarrow \partial_x(B(M_0)) \text{ in } W_2$$

for any nondecreasing C^1 function B on $[0, 1]$ satisfying $B(0) = 0$ and $B(1) = 1$. But as v_0 satisfies the assumption (H2), we obtain

$$\int_{\mathbb{R}} v_0(x) d\partial_x(B(M_0^n)) \rightarrow \int_{\mathbb{R}} v_0(x) d\partial_x(B(M_0)).$$

By the definitions of f_n and f , this means that

$$\int_0^1 f_n(m) B'(m) dm \rightarrow \int_0^1 f(m) B'(m) dm$$

for all nondecreasing C^1 functions B on $[0, 1]$ satisfying $B(0) = 0$ and $B(1) = 1$. Indeed,

$$\begin{aligned} \int_{\mathbb{R}} v_0(x) d\partial_x(B(M_0^n))(x) &= \sum_{i=1}^n v_0(x_i^{(n)}) [B(M_0^n(x_i^{(n)})) - B(M_0^n(x_i^{(n)}-))] \\ &= \sum_{i=1}^n \int_{M_0^n(x_i^{(n)}-)}^{M_0^n(x_i^{(n)})} v_0 \circ N_0^n(m) B'(m) dm = \int_0^1 f_n(m) B'(m) dm. \end{aligned}$$

For the continuous version we use (7) to conclude in a similar way. It then follows that

$$\int_0^1 f_n(m) g(m) dm \rightarrow \int_0^1 f(m) g(m) dm \text{ for all functions } g \in C([0, 1]).$$

This together with the fact

$$\int_0^1 |f_n(m)| dm = \int_{\mathbb{R}} |v_0(x)| d\rho_0^n \leq C \text{ uniformly in } n$$

yields $f_n \rightarrow f$ weakly in L^1 . By using the uniform convergence of \tilde{F}_n to \tilde{F} as continuous and bounded functions, we deduce that M is a solution of the problem (5). Now, for a fixed $t \geq 0$, let U be any C^1 function on $[0, 1]$. Define

$$\begin{aligned} \tilde{F}_{n,U}(t, m) &:= \int_0^m [f_n(\omega) + ta_n(\omega)] U'(\omega) d\omega, \\ \tilde{F}_U(t, m) &:= \int_0^m [f(\omega) + ta(\omega)] U'(\omega) d\omega. \end{aligned}$$

Then we use the uniform convergence of a_n to α and β and the facts that $f_n \rightarrow f$ weakly in L^1 , and

$$(30) \quad \int_0^1 \zeta(|f_n|) dm = \int_{\mathbb{R}} \zeta(|v_0|) d\rho_0^n \leq C$$

to deduce that $\tilde{F}_{n,U}$ converges to \tilde{F}_U uniformly in $(t, m) \in [0, T] \times [0, 1]$. (Notice that we have used the assumption (H2) to derive (30), the equi-integrability of the sequence $\{f_n\}$ in $L^1([0, 1])$.) Therefore, as in [7], we conclude that M is, in fact, an entropy solution for (5). \square

3.3. The existence result. We finally have all the necessary tools to prove Theorem 1.1. Note, however, that we leave out the proof of the fact that the initial conditions are satisfied. We will show that in Proposition 4.9.

Proof of Theorem 1.1. Since M is a solution of the problem (5) by Proposition 3.3, we can use Corollary 2.7 to conclude that $\partial_t M + v \partial_x M = 0$, where $v(t, \cdot) := g(t, \cdot) + t\psi(t, \cdot)$ is well defined $\partial_x M(t, \cdot) =: \rho(t, \cdot)$ -a.e. By differentiation in the sense of distributions, we obtain the first equation in (1). Then, Corollary 2.7 also gives

$$\begin{aligned} \partial_t(\rho v) &= \partial_t[\partial_x[\tilde{F}(t, M)]] = \partial_x[\partial_t[\tilde{F}(t, M)]] \\ &= \partial_x[v\partial_t M + \Psi(M)] = \partial_x(-v^2\rho) + \psi\rho. \end{aligned}$$

Note that, in our case, $\Psi(m) = \alpha m^2/2 + \beta m$. Thus,

$$\psi(t, x) = \beta + \alpha \int_0^1 [(1-s)M(t, x-) + sM(t, x)] ds = \beta + \alpha \left[M(t, x) - \frac{1}{2}\rho(t, \{x\}) \right].$$

Since $\partial_x M(t, \cdot) = \rho(t, \cdot)$ and $\rho(t, \cdot)$ has at most countably many atoms, we may take $\Phi(t, x) = \int_{-\infty}^x M(t, y) dy$ to conclude. \square

Remark 3.4. Note that the term $M(t, x) - \frac{1}{2}\rho(t, \{x\})$ is precisely the *barycentric projection* [2] onto $\rho(t, \cdot)$ of the *optimal coupling* between $\chi_{(0,1)}$ and $\rho(t, \cdot)$. It differs from the projection considered in [13] by an additive factor of 0.5 due to the fact that, instead of $\chi_{(0,1)}$, the reference measure in [13] was $\chi_{(-0.5,0.5)}$.

4. Time regularity, entropy condition, and shocks. In this section we discuss some qualitative properties of the “sticky particle” solution. The family of absolutely continuous curves in $\mathcal{P}_2(\mathbb{R})$ is central to our approach. Thus, recall that $(\mathcal{P}_2(\mathbb{R}), W_2)$ is a Polish space on which we define absolutely continuous curves by saying that $[0, T] \ni t \rightarrow \mu_t \in \mathcal{P}_2(\mathbb{R})$ lies in $AC^2(0, T; \mathcal{P}_2(\mathbb{R}))$ provided that there exists $f \in L^2(0, T)$ such that $W_2(\mu_t, \mu_{t+h}) \leq \int_t^{t+h} f(s) ds$ for all $0 < t < t+h < T$.

4.1. The solution path is locally Lipschitz. First we show that our solution path is locally Lipschitz.

PROPOSITION 4.1. *The solution path $t \rightarrow \rho(t, \cdot)$ satisfies the following:*

(i) *For any $0 < T < +\infty$, we have*

$$W_2(\rho(t, \cdot), \rho(s, \cdot)) \leq C_T |t - s| \quad \text{for all } t, s \in [0, T].$$

(ii) *The energy is nonincreasing, i.e.,*

$$\int_{\mathbb{R}} |v(t, x)|^2 d\rho(t, x) \leq \int_{\mathbb{R}} |v_0(x)|^2 d\rho_0(x) \quad \text{for all } t \geq 0.$$

Proof. Recall that $\rho^n(t, \cdot)$ are uniformly Lipschitz in n and t . Thus, by the triangle inequality,

$$W_2(\rho(t, \cdot), \rho(s, \cdot)) \leq W_2(\rho(t, \cdot), \rho^n(t, \cdot)) + C|t - s| + W_2(\rho^n(s, \cdot), \rho(s, \cdot)).$$

Therefore, by letting n go to infinity and using Proposition 3.3, we obtain (i). Also as $v_n(t, \cdot)\rho^n(t, \cdot) \rightarrow v(t, \cdot)\rho(t, \cdot)$ weakly for each t , we deduce that

$$\begin{aligned} \int_{\mathbb{R}} |v(t, x)|^2 d\rho(t, x) &\leq \liminf_{n \rightarrow \infty} \int_{\mathbb{R}} |v_n(t, x)|^2 d\rho^n(t, x) \leq \liminf_{n \rightarrow \infty} \int_{\mathbb{R}} |v_0(x)|^2 d\rho_0^n(x) \\ &= \int_{\mathbb{R}} |v_0(x)|^2 d\rho_0(x). \quad \square \end{aligned}$$

Remark 4.2. As a consequence of Proposition 4.1, we have $\rho \in AC^2(0, T; \mathcal{P}_2(\mathbb{R}))$.

We now recall a result, slightly modified, proved in [13].

PROPOSITION 4.3. *Suppose $\sigma \in AC^2(0, T; \mathcal{P}_2(\mathbb{R}))$. Let v be the velocity of minimal norm associated to σ and $N(t, \cdot) : (0, 1) \rightarrow \mathbb{R}$ be the monotone nondecreasing map such that $N(t, \cdot) \# \chi_{(0,1)} = \sigma(t, \cdot)$. For each t , modifying $N(t, \cdot)$ on a countable subset of $(0, 1)$ if necessary, we may assume w.l.o.g. that $N(t, \cdot)$ is right continuous. Then, $N \in H^1(0, T; L^2(0, 1))$ and*

$$(31) \quad \dot{N}(t, x) = v(t, N(t, x))$$

for \mathcal{L}^2 -a.e. $(t, x) \in (0, T) \times (0, 1)$.

Note that the *minimal norm* assumption is, in fact, redundant. Indeed, we prove the following lemma.

LEMMA 4.4. *Consider the path $t \rightarrow \mu \in AC^2(0, T; \mathcal{P}_2(\mathbb{R}))$ for some $0 < T < +\infty$. Then the velocity defined in [2] (called “of minimal norm”) is the unique velocity along the curve μ in the following sense: if*

$$\partial_t \mu + \partial_x(\mu v_i) = 0 \text{ in } \mathcal{D}'([0, T] \times \mathbb{R}), \quad i = 1, 2,$$

for some v_i Borel measurable in (t, x) such that $v_i(t, \cdot) \in L^2(\mu(t, \cdot))$ for Lebesgue a.e. $t \in (0, T)$, then for Lebesgue a.e. $t \in (0, T)$ we have $v_1(t, \cdot) \equiv v_2(t, \cdot)$ in the $\mu(t, \cdot)$ -a.e. sense.

Proof. By subtraction and by taking test functions $\varphi(t, x) = \xi(t)\zeta(x)$, the equations above readily yield

$$\int_{\mathbb{R}} u(t, x)\zeta'(x)d\mu(t, x) = 0 \text{ for a.e. } t \in (0, T) \text{ and any } \zeta \in C_c^1(\mathbb{R}),$$

where $u := v_1 - v_2$. Fix $\varepsilon > 0$ and $\phi \in C_c(\mathbb{R})$. If $\phi = 0$ on $[R, +\infty)$, consider, for each natural number $n > R$, the function

$$(32) \quad \Phi_n(x) := \begin{cases} \int_{-\infty}^x \phi(y)dy & \text{if } x < n, \\ \omega(x - n) & \text{if } n \leq x \leq n + 1, \\ 0 & \text{if } x > n + 1, \end{cases}$$

where $\omega \in C^1[0, 1]$ such that $\omega(0) = \int_{-\infty}^R \phi(y)dy$, $\omega(1) = 0$, and $\omega'(0) = 0 = \omega'(1)$. Clearly, $\Phi_n \in C_c^1(\mathbb{R})$. Thus,

$$\int_{\mathbb{R}} u(t, x)\phi(x)d\mu(t, x) + \int_n^{n+1} u(x, t)\omega'(x - n)d\mu(t, x) = 0 \text{ for a.e. } t \in (0, T).$$

We have $|\omega'(x - n)| \leq \|\omega'\|_{L^\infty(0,1)} =: C$ for all $n > R$ and all $x \in (n, n + 1)$. Since $\mu(t, \cdot)$ is a Borel probability for Lebesgue a.e. $t \in (0, T)$, we conclude that for such t we have

$$\left| \int_{\mathbb{R}} u(t, x)\phi(x)d\mu(t, x) \right| \leq C\|u(t, \cdot)\|_{L^2(\mu(t, \cdot))}\mu(t, [n, n + 1])^{1/2} \leq \varepsilon$$

if n is sufficiently large. Due to the arbitrariness of ε and ϕ , the proof is concluded. \square

Remark 4.5. Note that, in fact, we have just proved that $\{\varphi' : \varphi \in C_c^\infty(\mathbb{R})\}$ is dense in $L^2(\mu)$, even though μ may not necessarily have finite p -order moment (for any $p > 0$). Also, as a consequence, the tangent space $\mathcal{T}_\mu\mathcal{P}_2(\mathbb{R})$ [2] is the whole $L^2(\mu)$. This property was brought to our attention by W. Gangbo.

4.2. Recovery of the entropy condition. We shall now prove that the solution we obtained for (1) in the $\alpha = 0 = \beta$ case satisfies the *Oleinik entropy condition*.

THEOREM 4.6. *For Lebesgue a.e. $t \in (0, T)$ we have*

$$(33) \quad v(t, x_2) - v(t, x_1) \leq \frac{1}{t}(x_2 - x_1) \text{ for } \rho(t, \cdot)\text{-a.e. } x_1 \leq x_2.$$

Proof. Let us look back at the discrete problem. Note that

$$N_n(t, \omega) := x_j(t) \text{ whenever } M_n(t, x_j(t)-) \leq \omega < M_n(t, x_j(t))$$

is the optimal map such that $N_n(t, \cdot) \# \chi_{(0,1)} = \rho^n(t, \cdot)$. It is known [7], [15] that the discrete problem satisfies the Oleinik entropy condition, i.e.,

$$t[v_n(t, x_{i_2}(t)) - v_n(t, x_{i_1}(t))] \leq x_{i_2}(t) - x_{i_1}(t) \text{ whenever } i_1 \leq i_2.$$

Since $v_n(t, x_j(t)) = \dot{x}_j(t)$ away from collision times, we infer that the map $t \rightarrow [x_{i_2}(t) - x_{i_1}(t)]/t$ is piecewise nonincreasing. But this map is continuous, so it is globally nonincreasing, and, due to the definition of N_n , it follows that

$$(34) \quad t \rightarrow \frac{1}{t}[N_n(t, \omega_2) - N_n(t, \omega_1)]$$

is nonincreasing in $(0, T)$ for all $\omega_1 \leq \omega_2 \in (0, 1)$. Now let $\Delta := \{\omega = (\omega_1, \omega_2) \in [0, 1]^2 : \omega_1 \leq \omega_2\}$ and $S_n(t, \omega) = N_n(t, \omega_2) - N_n(t, \omega_1)$ be defined on $[0, T] \times \Delta$. It is easy to see that $S_n \in H^1(0, T; L^2(\Delta))$ and that (34) implies

$$(35) \quad \int_0^T \int_\Delta \frac{S_n(t, \omega)}{t} \partial_t \varphi(t, \omega) d\omega dt \geq 0 \text{ for all nonnegative } \varphi \in C_c^\infty((0, T) \times \Delta).$$

On the other hand, due to (28), we infer that

$$\|N_n - N\|_{L^1((0,T) \times (0,1))} = \int_0^T W_1(\rho^n(t, \cdot), \rho(t, \cdot)) dt \rightarrow 0,$$

where $N(t, \cdot)$ is the optimal map such that $N(t, \cdot) \# \chi_{(0,1)} = \rho(t, \cdot)$. In particular, $(S_n/t) \rightarrow (S/t)$ in $\mathcal{D}'((0, T) \times \Delta)$, where $S(t, \omega) = N(t, \omega_2) - N(t, \omega_1)$ if $\omega \in \Delta$. Thus, (35) implies

$$\int_0^T \int_\Delta \frac{S(t, \omega)}{t} \partial_t \varphi(t, \omega) d\omega dt \geq 0 \text{ for all nonnegative } \varphi \in C_c^\infty((0, T) \times \Delta).$$

Therefore, $\partial_t[S/t] \leq 0$ in the distributional sense, which implies $t\dot{S} - S \leq 0$ in the \mathcal{L}^3 -a.e. sense, i.e.,

$$\dot{N}(t, \omega_2) - \dot{N}(t, \omega_1) \leq \frac{1}{t} [N(t, \omega_2) - N(t, \omega_1)] \text{ for Lebesgue a.e. } (t, \omega) \in (0, T) \times \Delta.$$

Consequently, (31) yields

$$v(t, N(t, \omega_2)) - v(t, N(t, \omega_1)) \leq \frac{1}{t} [N(t, \omega_2) - N(t, \omega_1)] \text{ for } \mathcal{L}^3\text{-a.e. } (t, \omega) \in (0, T) \times \Delta,$$

which, due to $N(t, \cdot) \# \chi_{(0,1)} = \rho(t, \cdot)$, implies (33). \square

4.3. Formation of shocks. Now let us assume that ρ_0 is atom-free and take N_0 to be the optimal map pushing $\chi_{(0,1)}$ forward to ρ_0 .

PROPOSITION 4.7. *Let $T := \sup\{t \in [0, \infty) : \text{id} + tv_0 \text{ is nondecreasing on } \text{spt}(\rho_0)\}$. If $T = 0$, then the solution develops atomic singularities instantaneously. If $0 < T < +\infty$, then the solution remains nonatomic before T and develops atomic singularities instantaneously after T . If $T = +\infty$, then the solution is atom-free at all times.*

Proof. Let us treat the case $0 < T < +\infty$. Indeed, it will become clear that the other two cases can be handled almost identically. Note that, due to the definition of T and the fact that ρ_0 is nonatomic, $\text{id} + tv_0$ is (strictly) increasing on $\text{spt}(\rho_0)$ for $t \in [0, T)$. Thus, $N_0 + tv_0 \circ N_0$ is increasing on $(0, 1)$ for $t \in [0, T)$. It follows that

$$\bar{M}(t, \cdot) := (N_0 + tv_0 \circ N_0)^{-1} = M_0 \circ (\text{id} + tv_0)^{-1}$$

is the entropy solution (given by characteristics) of

$$\partial_t \bar{M} + \partial_x [F(\bar{M})] = 0, \quad \bar{M}(0, \cdot) = M_0,$$

where $F' = v_0 \circ N_0$, $F(0) = 0$. Due to uniqueness of the entropy solution, we get $\bar{M} \equiv M$. Thus, $\rho(t, \cdot) = (\text{id} + tv_0) \# \rho_0$ for $t \in [0, T]$ (this is precisely the geodesic connecting ρ_0 and ρ_T in $\mathcal{P}_2(\mathbb{R})$). Again, since $\text{id} + tv_0$ is (strictly) increasing on $\text{spt}(\rho_0)$ for $t \in [0, T)$, we deduce that $\rho(t, \cdot)$ has no atoms if $t \in [0, T)$. To conclude, we argue by contradiction and suppose that $\rho(t, \cdot)$ is atom-free on $[0, T + \epsilon]$ for some $\epsilon > 0$. Thus, $M(t, \cdot)$ is continuous, and so $v(t, \cdot) = v_0 \circ N_0 \circ M(t, \cdot)$ for all $t \in [0, T + \epsilon]$. Since $M(t, \cdot)$ is now the optimal map pushing $\rho(t, \cdot)$ forward to $\chi_{(0,1)}$, we infer that

$$v(t, N(t, m)) = v_0 \circ N_0 \circ M(t, N(t, m)) = v_0 \circ N_0(m) \text{ a.e. } m \in (0, 1),$$

which, in light of (31), yields $\dot{N}(t, m) = v_0 \circ N_0(m)$ for a.e. $(t, m) \in (0, T + \epsilon) \times (0, 1)$. Thus, as an $H^1(0, T + \epsilon; L^2(0, 1))$ map, $N(t, \cdot) = N_0 + tv_0 \circ N_0$. In particular, we obtain that $\text{id} + tv_0$ is nondecreasing on the support of ρ_0 for all $t \in [0, T + \epsilon]$, which contradicts the definition of T . Therefore, the solution becomes atomic instantaneously after T . \square

Remark 4.8. What happens at $t = T$ depends on whether or not $\text{id} + Tv_0$ has “flat spots” on $\text{spt}(\rho_0)$. It is easy to construct examples illustrating that each of these situations may occur.

4.4. Continuity of the energy and a remark on uniqueness.

PROPOSITION 4.9. *Suppose (ρ_0, v_0) satisfies the conditions (H1) and (H2). Let (ρ, v) be the weak solution to the system (1) given in the proof of Theorem 1.1 in subsection 3.3. Then (ρ, v) has the following property:*

$$(36) \quad \lim_{t \rightarrow 0^+} \int_{\mathbb{R}} v(t, x) \varphi(x) d\rho(t, x) = \int_{\mathbb{R}} v_0(x) \varphi(x) d\rho_0(x) \text{ for all } \varphi \in C_b(\mathbb{R}),$$

which shows that the initial condition for the velocity is satisfied. Moreover, we have

$$(37) \quad \lim_{t \rightarrow 0^+} \int_{\mathbb{R}} v^2(t, x) \varphi(x) d\rho(t, x) = \int_{\mathbb{R}} v_0^2(x) \varphi(x) d\rho_0(x) \quad \text{for all } \varphi \in C_b(\mathbb{R}).$$

Proof. Due to $\|M(t, \cdot) - M_0\|_{L^1(\mathbb{R})} = W_1(\rho(t, \cdot), \rho_0) \rightarrow 0$ and the BV calculus, we have for any function $\varphi \in C_c^\infty(\mathbb{R})$

$$(38) \quad \begin{aligned} \lim_{t \rightarrow 0^+} \int_{\mathbb{R}} v(t, x) \varphi(x) d\rho(t, x) &= - \lim_{t \rightarrow 0^+} \int_{\mathbb{R}} \varphi'(x) F(M(t, x)) dx \\ &= - \int_{\mathbb{R}} \varphi'(x) F(M_0(x)) dx = \int_{\mathbb{R}} \varphi(x) g(x) d\rho_0(x), \end{aligned}$$

where $g(x)$ is given by

$$g(x) = \int_0^1 v_0 \circ N_0((1 - s)M_0(x-) + sM_0(x)) ds, \quad x \in \mathbb{R}.$$

We claim that $N_0((1 - s)M_0(x-) + sM_0(x)) = x$ for ρ_0 -a.e. $x \in \mathbb{R}$. Indeed, if a point x satisfies $\rho_0(\{x\}) \neq 0$, then M_0 has a jump at x . Therefore, it is clear in this case that

$$N_0((1 - s)M_0(x-) + sM_0(x)) = \inf \{z : M_0(z) > (1 - s)M_0(x-) + sM_0(x)\} = x.$$

As M_0 has at most countably many “flat spots,” the claim shall be proved if we can show that $g(x) = v_0(x)$ whenever x satisfies $x \in \text{spt}(\rho_0)$, $\rho_0(\{x\}) = 0$ and x is not one of the endpoints of some flat spot of M_0 . To see this is true, let x be a point with such properties. Then there exists $\epsilon > 0$ such that the function M_0 is strictly increasing on $(x - \epsilon, x + \epsilon)$. It follows that

$$N_0((1 - s)M_0(x-) + sM_0(x)) = N_0(M_0(x)) = x.$$

Thus, we obtain the claim, which in turn yields $g(x) = v_0(x)$ for ρ_0 -a.e. $x \in \mathbb{R}$. Hence, by combining this with (38), we get

$$\lim_{t \rightarrow 0^+} \int_{\mathbb{R}} v(t, x) \varphi(x) d\rho(t, x) = \int_{\mathbb{R}} v_0(x) \varphi(x) d\rho_0(x) \quad \text{for all } \varphi \in C_c^\infty(\mathbb{R}).$$

By a simple approximation and the fact that $\|v_0\|_{L^2(\rho_0)}$ is finite, this gives (36). As a consequence of (36) and the fact that the energy is nonincreasing which was proved in Proposition 4.1, we obtain (see, e.g., Theorem 5.4.4 in [2])

$$\lim_{t \rightarrow 0^+} \int_{\mathbb{R}} v^2(t, x) d\rho(t, x) = \int_{\mathbb{R}} v_0^2(x) d\rho_0(x).$$

In addition, we have $W_2(\rho(t, \cdot), \rho_0) \leq Ct$ from Proposition 4.1. Therefore, by using Proposition 7.1.5 and Theorem 5.4.4 in [2] we can conclude that (37) holds. \square

We end the paper by the following remark on the uniqueness of the solution.

Remark 4.10. If we assume, in addition, that v_0 is bounded, then the weak solution (ρ, v) constructed above is the unique weak solution satisfying the entropy condition in Theorem 4.6 and the weak convergence of $v^2(t, \cdot) \rho(t, \cdot)$ to $v_0^2 \rho_0$ (such a weak solution is called an entropy solution for pressureless gases [16]). This fact follows directly from Theorem 4.6, Proposition 4.9, and the uniqueness result in [16]. Thus, if $v_0 \in C_b(\mathbb{R})$, we have proved that any entropy solution to the pressureless gases system can be obtained as a weak limit of sticky particles. We note that a similar result was obtained by Bouchut and James in [5] for duality solutions.

Acknowledgments. The authors would like to thank W. Gangbo, R. Pan, and A. Swiech for numerous and fruitful discussions. We are also grateful to both anonymous referees for their helpful comments and suggestions.

REFERENCES

- [1] L. AMBROSIO AND W. GANGBO, *Hamiltonian ODEs in the Wasserstein spaces of probability measures*, Comm. Pure Appl. Math., 61 (2008), pp. 18–53.
- [2] L. AMBROSIO, N. GIGLI, AND G. SAVARÉ, *Gradient Flows in Metric Spaces and the Wasserstein Spaces of Probability Measures*, Lectures Math. ETH Zürich, Birkhäuser, Basel, 2005.
- [3] A. ANDRIEVSKII, S. GURBATOV, AND A. SOBOLEVSKII, *Ballistic aggregation in symmetric and nonsymmetric flows*, J. Exp. Theor. Phys., 104 (2007), pp. 887–896.
- [4] J. D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem*, Numer. Math., 84 (2000), pp. 375–393.
- [5] F. BOUCHUT AND F. JAMES, *Duality solutions for pressureless gases, monotone scalar conservation laws, and uniqueness*, Comm. Partial Differential Equations, 24 (1999), pp. 2173–2189.
- [6] L. BOUDIN, *A solution with bounded expansion rate to the model of viscous pressureless gases*, SIAM J. Math. Anal., 32 (2000), pp. 172–193.
- [7] Y. BRENIER AND E. GRENIER, *Sticky particles and scalar conservation laws*, SIAM J. Numer. Anal., 35 (1998), pp. 2317–2328.
- [8] H. BREZIS, *Analyse fonctionnelle; théorie et applications*, Collection Mathématiques Appliquées pour la Maîtrise, Masson, Paris, 1983.
- [9] S. CHENG, J. LI, AND T. ZHANG, *Explicit construction of measure solutions of Cauchy problem for transportation equations*, Sci. China Ser. A, 40 (1997), pp. 1287–1299.
- [10] C. DELLACHERIE AND P. A. MEYER, *Probabilités et potentiel*, Hermann, Paris, 1975.
- [11] R. M. DUDLEY, *Real Analysis and Probability*, Cambridge Stud. Adv. Math. 74, Cambridge University Press, Cambridge, UK, 2002.
- [12] M. FRÉCHET, *Sur la distance de deux lois de probabilité*, C. R. Acad. Sci. Paris, 244 (1957), pp. 689–692.
- [13] W. GANGBO, T. NGUYEN, AND A. TUDORASCU, *Euler-Poisson systems as action-minimizing paths in the Wasserstein space*, Arch. Ration. Mech. Anal., to appear.
- [14] W. GANGBO AND T. PACINI, *Infinite dimensional Hamiltonian systems in terms of the Wasserstein distance*, in progress.
- [15] E. GRENIER, *Existence globale pour le système des gaz sans pression*, C. R. Acad. Sci. Paris Sér. I Math., 321 (1995), pp. 171–174.
- [16] F. HUANG AND Z. WANG, *Well posedness for pressureless flow*, Comm. Math. Phys., 222 (2001), pp. 117–146.
- [17] P. MARTIN AND J. PIASECKI, *One-dimensional ballistic aggregation: Rigorous long-time estimates*, J. Statist. Phys., 76 (1994), pp. 447–476.
- [18] P. MARTIN AND J. PIASECKI, *Aggregation dynamics in a self-gravitating one-dimensional gas*, J. Statist. Phys., 84 (1996), pp. 837–857.
- [19] G. MONGE, *Mémoire sur la théorie des déblais et de remblais*, Histoire de l’Académie Royale des Sciences de Paris, avec les Mémoires de Mathématique et de Physique pour la même année, 1781, pp. 666–704.
- [20] A. SHNIRELMAN, *On the principle of the shortest way in the dynamics of systems with constraints*, in Global Analysis—Studies and Applications II, Lecture Notes in Math. 1214, Springer-Verlag, Berlin, 1986, pp. 117–130.
- [21] C. VILLANI, *Topics in optimal transportation*, Grad. Stud. Math. 58, AMS, Providence, RI, 2003.
- [22] A. VOLPERT, *Spaces BV and quasilinear equations*, Mat. Sb. (N.S.), 73 (1967), pp. 255–302.
- [23] E. WEINAN, Y. RYKOV, AND Y. SINAI, *Generalized variational principles, global weak solutions and behavior with random initial data for systems of conservation laws arising in adhesion particle dynamics*, Comm. Math. Phys., 177 (1996), pp. 349–380.
- [24] YA. B. ZELDOVICH, *Gravitational instability: An approximate theory for large density perturbations*, Astro. & Astrophys., 5 (1970), pp. 84–89.

SPREADING SPEEDS AND TRAVELING WAVES FOR NONMONOTONE INTEGRODIFFERENCE EQUATIONS*

SZE-BI HSU[†] AND XIAO-QIANG ZHAO[‡]

Abstract. The spreading speeds and traveling waves are established for a class of nonmonotone discrete-time integrodifference equation models. It is shown that the spreading speed is linearly determinate and coincides with the minimal wave speed of traveling waves.

Key words. integrodifference equations, nonmonotone integral operators, spreading speeds, linear determinacy, traveling waves

AMS subject classifications. 37L15, 39A11, 92D25

DOI. 10.1137/070703016

1. Introduction. The invasion speed is a fundamental characteristic of biological invasions, since it describes the speed at which the geographic range of the population expands; see, e.g., [6, 8, 9, 15] and references therein. Aronson and Weinberger [1, 2] first introduced the concept of the asymptotic speed of spread (in short, spreading speed) for reaction-diffusion equations and showed that it coincides with the minimal wave speed for traveling waves under appropriate assumptions. Weinberger [20] and Lui [13] established the theory of spreading speeds and monostable traveling waves for monotone (order-preserving) operators. This theory has been greatly developed recently in [21, 10, 11, 12] to monotone semiflows so that it can be applied to various discrete- and continuous-time evolution equations admitting the comparison principle.

It is known that many discrete- and continuous-time population models with spatial structure are not monotone. For example, scalar discrete-time integrodifference equations with nonmonotone growth functions, and predator-prey type reaction-diffusion systems are among such models. The spreading speeds were obtained for some nonmonotone continuous-time integral equations and time-delayed reaction-diffusion models in [17, 19], and a general result on the nonexistence of traveling waves was also given in [19, Theorem 3.5]. The existence of monostable traveling waves was established for several classes of nonmonotone time-delayed reaction-diffusion equations in [22, 4, 16, 14]. For certain types of nonmonotone discrete-time integrodifference equation models, nonmonotone traveling waves and even traveling cycles were observed in [7] by numerical simulations. In [7, 9, 15, 8], the monotone linear systems, resulting from the linearization of the nonmonotone discrete-time models at zero, were used to estimate spreading speeds. It is worthy to find sufficient conditions under which the spreading speed is linearly determinate for these nonmonotone systems.

*Received by the editors September 16, 2007; accepted for publication (in revised form) April 22, 2008; published electronically August 1, 2008.

<http://www.siam.org/journals/sima/40-2/70301.html>

[†]Department of Mathematics, National Tsing Hua University, Hsinchu, Taiwan, Republic of China (sbhsu@math.nthu.edu.tw). Research supported in part by the National Council of Science of Republic of China.

[‡]Department of Mathematics and Statistics, Memorial University of Newfoundland, St. John's, NL A1C 5S7, Canada (xzhao@math.mun.ca). Research supported in part by the NSERC of Canada and the MITACS of Canada.

The purpose of our current paper is to study the spreading speeds and traveling waves for nonmonotone discrete-time systems. As a starting point, we consider scalar integrodifference equations with nonmonotone growth functions. The key techniques are to sandwich the given growth function in between two appropriate nondecreasing functions (for spreading speeds) and to construct a closed and convex subset in an appropriate Banach space (for traveling waves). Consequently, we obtain a set of sufficient conditions for the existence of the spreading speed, and the existence and nonexistence of traveling waves. It turns out that the spreading speed is linearly determinate and coincides with the minimal wave speed of traveling waves for this class of nonmonotone discrete-time integrodifference equation population models.

The rest of this paper is organized as follows. In section 2, we first present a general result for monotone integrodifference equations, then we establish the spreading speed c^* by the comparison method and a fluctuation type argument. In section 3, we use the Schauder fixed point theorem to obtain the existence of traveling waves with the wave speed $c > c^*$. The property of the spreading speed is employed to prove the asymptotic property of the wave profile at $+\infty$ and the nonexistence of traveling waves with $c < c^*$. A limiting argument gives the existence of the traveling wave with the wave speed c^* . Section 4 is aimed at the applications of the main results to three types of growth functions arising from population biology.

2. Spreading speeds. Let \mathcal{C} be the space of all bounded and continuous functions from \mathbb{R} to \mathbb{R} equipped with the compact open topology. For a given number $r > 0$, let $\mathcal{C}_r := \{\phi \in \mathcal{C} : 0 \leq \phi(x) \leq r, \forall x \in \mathbb{R}\}$.

Let $k(x)$ be a nonnegative Lebesgue measurable function on \mathbb{R} . Throughout this paper, we assume that the kernel $k(x)$ has the following property:

- (K) $\int_{\mathbb{R}} k(y)dy = 1, k(-y) = k(y), \forall y \in \mathbb{R}$, and $\int_{\mathbb{R}} e^{-\alpha y}k(y)dy < \infty, \forall \alpha \in [0, \Delta)$, where $\Delta > 0$ is the abscissa of convergence and it may be infinity.

We consider a discrete-time integrodifference equation

$$(2.1) \quad u_{n+1}(x) = \int_{\mathbb{R}} h(u_n(y))k(x - y)dy, \quad x \in \mathbb{R}, n \geq 0$$

with $u_0 \in \mathcal{C}$. Assume that there exists $\beta > 0$ such that

- (H1) $h \in C([0, \beta], [0, \beta]), h(0) = 0, h'(0) > 1, h(\beta) = \beta$, and there is $L_0 > 0$ such that $|h(u_1) - h(u_2)| \leq L_0|u_1 - u_2|, \forall u_1, u_2 \in [0, \beta]$.
- (H2) $u < h(u) \leq h'(0)u, \forall u \in (0, \beta)$, and $h(u)$ is nondecreasing in $u \in [0, \beta]$.

Let $U(x)$ be a continuous function on \mathbb{R} . We say $U(x + cn)$ is a traveling wave solution of (2.1) with the wave speed c if $u_n(x) = U(x + cn), \forall n \geq 0$, satisfies (2.1), and $U(x + cn)$ connects 0 to β if $U(-\infty) = 0$ and $U(+\infty) = \beta$. It is easy to see that $U(x + cn)$ is a traveling wave solution of (2.1) if and only if

$$U(\xi) = \int_{\mathbb{R}} h(U(\xi - c - y))k(y)dy, \quad \forall \xi \in \mathbb{R}.$$

Define

$$(2.2) \quad c_h^* = \inf_{\mu \in (0, \Delta)} \frac{\ln(h'(0) \int_{\mathbb{R}} e^{-\mu y}k(y)dy)}{\mu}.$$

The following result is essentially due to Weinberger [20], and shows that c_h^* is not only the spreading speed but also the minimal wave speed of monotone traveling waves for system (2.1).

THEOREM 2.1. *Let (H1) and (H2) hold. Then the following statements are valid:*

(i) *For any $u_0 \in \mathcal{C}_\beta$ with compact support, the solution of (2.1) satisfies*

$$\lim_{n \rightarrow \infty, |x| \geq cn} u_n(x) = 0, \forall c > c_h^*.$$

(ii) *For any $u_0 \in \mathcal{C}_\beta \setminus \{0\}$, the solution of (2.1) satisfies*

$$\lim_{n \rightarrow \infty, |x| \leq cn} u_n(x) = \beta, \forall c \in (0, c_h^*).$$

(iii) *For any $c \geq c_h^*$, (2.1) has a traveling wave $U(x + cn)$ connecting 0 to β such that $U(x)$ is nondecreasing in x , and for any $c \in (0, c_h^*)$, (2.1) has no traveling wave $U(x + cn)$ connecting 0 to β .*

Proof. The existence of the spreading speed, together with the formula (2.2), and traveling waves is a straightforward consequence of [20, Theorems 6.1–6.6] in the case where $\Delta = +\infty$, and [11, Theorem 2.11 and Theorem 2.15] with $\tau = 0$ in the case where $\Delta < +\infty$. Indeed, the proof of [11, Theorem 3.10] implies that [11, Theorem 3.10] with $\inf_{\mu > 0} \Phi(\mu)$ replaced by $\inf_{0 < \mu < \Delta} \Phi(\mu)$ is still valid, provided that (C5) holds for all $\mu_1, \mu_2 \in (-\Delta, \Delta)$ and $\Phi(\mu)$ assumes its minimum value at $\mu^* \in (0, \Delta)$. Thus, the formula (2.2) also holds in the case where $\Delta < +\infty$. By [11, Theorem 3.5], it follows that the number $r = r_\sigma$ in [11, Theorem 2.15] can be chosen to be independent of $\sigma > 0$. For any $u_0 \in \mathcal{C}_\beta \setminus \{0\}$, it is easy to show that there exists an integer $n_0 \geq 1$ such that $u_{n_0}(x) > 0$ for x in an interval of length greater $2r$. Taking $u_{n_0}(x)$ as a new initial data, we see that conclusion (ii) holds. The nonexistence of traveling waves is implied by conclusion (ii) (see also [11, Theorem 4.1]). \square

Next we consider the discrete-time integrodifference equation

$$(2.3) \quad u_{n+1}(x) = \int_{\mathbb{R}} f(u_n(y))k(x - y)dy, \quad x \in \mathbb{R}, n \geq 0$$

with $u_0 \in \mathcal{C}$. Assume that there exists $b > 0$ such that

(F1) $f \in C([0, b], [0, b])$, $f(0) = 0$, $f'(0) > 1$, and there is $L > 0$ such that $|f(u_1) - f(u_2)| \leq L|u_1 - u_2|$, $\forall u_1, u_2 \in [0, b]$.

(F2) $f(u) \leq f'(0)u$, $\forall u \in [0, b]$, and there is $u^* \in (0, b]$ such that $f(u^*) = u^*$, $f(u) > u$, $\forall u \in (0, u^*)$, and $0 < f(u) < u$, $\forall u \in (u^*, b]$.

Define

$$f^+(u) = \max_{0 \leq v \leq u} f(v), \quad f^-(u) = \min_{u \leq v \leq b} f(v), \quad \forall u \in [0, b].$$

It then follows that

$$f^-(u) \leq f(u) \leq f^+(u), \forall u \in [0, b],$$

that both f^+ and f^- are nondecreasing and Lipschitz continuous, with the Lipschitz constant L , on $[0, b]$, and that there exists $\delta_0 \in (0, b]$ such that $f^\pm(u) = f(u)$, $\forall u \in [0, \delta_0]$. Let u_\pm^* be such that $f^\pm(u_\pm^*) = u_\pm^*$. Then $0 < u_-^* \leq u^* \leq u_+^* \leq b$.

To obtain the upward convergence as stated in Theorem 2.1 (ii), we need to impose one of the following two additional conditions on f .

(C1) $u^* = b$ and $f(u)$ is nondecreasing in $u \in [b - \epsilon_0, b]$ for some $\epsilon_0 \in (0, b)$.

(C2) $\frac{f(u)}{u}$ is strictly decreasing for $u \in (0, b]$, and $f(u)$ has the property (P) that for any $v, w \in (0, b]$ satisfying $v \leq u^* \leq w$, $v \geq f(w)$ and $w \leq f(v)$, we have $v = w$.

Motivated by the proofs of [17, Theorems 2.9 and 2.12], we have the following observation.

LEMMA 2.1. *Either of the following two conditions is sufficient for the property (P) in condition (C2) to hold:*

(P1) *$uf(u)$ is strictly increasing for $u \in (0, b]$.*

(P2) *$f(u)$ is nonincreasing for $u \in [u^*, b]$, and $\frac{f^2(u)}{u}$ is strictly decreasing for $u \in (0, u^*]$.*

Proof. Let $v, w \in (0, b]$ be given such that $v \leq u^* \leq w$, $v \geq f(w)$, and $w \leq f(v)$. In the case where (P1) holds, since $vf(v) \geq f(w)f(v) \geq wf(w)$, it follows that $v \geq w$, and hence $v = w$. In the case where (P2) holds, we have $v \geq f(w) \geq f(f(v)) = f^2(v)$, and hence, $\frac{f^2(v)}{v} \leq 1 = \frac{f^2(u^*)}{u^*}$. It then follows that $v \geq u^*$, and hence, $v = u^*$. Since $u^* \leq w \leq f(v) = f(u^*) = u^*$, we further have $w = u^* = v$. \square

Now we are in a position to prove the main result of this section.

THEOREM 2.2. *Let (F1) and (F2) hold and c_f^* be defined as in (2.2) with $h = f$. Then the following statements are valid:*

(1) *For any $u_0 \in \mathcal{C}_{u_+^*}$ with compact support, the solution of (2.3) satisfies*

$$\lim_{n \rightarrow \infty, |x| \geq cn} u_n(x) = 0, \quad \forall c > c_f^*.$$

(2) *For any $u_0 \in \mathcal{C}_{u_+^*} \setminus \{0\}$, the solution of (2.3) satisfies*

$$u_-^* \leq \liminf_{n \rightarrow \infty, |x| \leq cn} u_n(x) \leq \limsup_{n \rightarrow \infty, |x| \leq cn} u_n(x) \leq u_+^*, \quad \forall c \in (0, c_f^*).$$

(3) *If, in addition, either (C1) or (C2) holds, then for any $u_0 \in \mathcal{C}_{u_+^*} \setminus \{0\}$, the solution of (2.3) satisfies*

$$\lim_{n \rightarrow \infty, |x| \leq cn} u_n(x) = u^*, \quad \forall c \in (0, c_f^*).$$

Proof. For convenience, let $c^* = c_f^*$. Define

$$Q(\phi)(x) = \int_{\mathbb{R}} f(\phi(x-y))k(y)dy, \quad Q^\pm(\phi)(x) = \int_{\mathbb{R}} f^\pm(\phi(x-y))k(y)dy.$$

Clearly, Q^\pm is monotone (order preserving) on \mathcal{C}_b and

$$Q^-(\phi) \leq Q(\phi) \leq Q^+(\phi), \quad \forall \phi \in \mathcal{C}_b.$$

By Theorem 2.1, it follows that c^* is the spreading speed for the discrete-time system $u_{n+1} = Q^\pm(u_n)$ on $\mathcal{C}_{u_\pm^*}$.

Case 1. For a given $\phi \in \mathcal{C}_{u_+^*}$ with compact support, let $u_n = Q^n(\phi)$, $u_n^+ = (Q^+)^n(\phi)$, $\forall n \geq 0$. By the comparison principle (see, e.g., [20, Proposition 4.1]), we have

$$0 \leq u_n(x) \leq u_n^+(x), \quad \forall x \in \mathbb{R}, n \geq 0.$$

For any $c > c^*$, Theorem 2.1 (i) implies that $\lim_{n \rightarrow \infty, |x| \geq cn} u_n^+(x) = 0$, and hence $\lim_{n \rightarrow \infty, |x| \geq cn} u_n(x) = 0$.

Case 2. For a given $\phi \in \mathcal{C}_{u_+^*} \setminus \{0\}$, define $\psi(x) = \min(\phi(x), u_-^*)$. Then $\psi \in \mathcal{C}_{u_\pm^*} \setminus \{0\}$. Let $u_n^- = (Q^-)^n(\psi)$, $\forall n \geq 0$. Since $\psi \leq \phi$, it follows from the comparison principle that

$$0 \leq u_n^-(x) \leq u_n(x) \leq u_n^+(x), \quad \forall x \in \mathbb{R}, n \geq 0.$$

For any $c \in (0, c^*)$, Theorem 2.1 (ii) implies that $\lim_{n \rightarrow \infty, |x| \leq cn} u_n^\pm(x) = u_\pm^*$. Thus, we have

$$u_-^* \leq \liminf_{n \rightarrow \infty, |x| \leq cn} u_n(x) \leq \limsup_{n \rightarrow \infty, |x| \leq cn} u_n(x) \leq u_+^*.$$

Case 3. In the case where (C1) holds, we see that $u^* \leq u_+^* \leq b = u^*$. Further, it follows from the definition of f^- that $f^-(u) = f(u), \forall u \in [b - \epsilon_0, b]$, and hence, $u_-^* = u_+^* = u^*$. Since $u_-^* = u_+^* = u^*$, the upward convergence in case (3) follows from the conclusion in case (2).

In the case where (C2) holds, we use similar arguments as in the proof of [17, Lemma 3.10] (see also the proof of [19, Theorem 2.5]). For any $(v, w) \in [0, b]^2$, let

$$(2.4) \quad g(v, w) = \begin{cases} \min\{f(u) : v \leq u \leq w\}, & \text{if } v \leq w, \\ \max\{f(u) : w \leq u \leq v\}, & \text{if } w \leq v. \end{cases}$$

Then $g(v, w)$ is nondecreasing in $v \in [0, b]$ and nonincreasing in $w \in [0, b]$. Moreover, $f(u) = g(u, u)$, and $g(v, w)$ is continuous in $(v, w) \in [0, b]^2$ (see [18, section 2]). For a given $u_0 \in \mathcal{C}_{u_+^*} \setminus \{0\}$, let $u_n = Q^n(u_0), \forall n \geq 0$. Then we have

$$u_{n+1}(x) = \int_{\mathbb{R}} g(u_n(x-y), u_n(x-y))k(y)dy, \quad n \geq 0.$$

For any $\beta \in (0, c^*)$, we define

$$U_*(\beta) := \liminf_{n \rightarrow \infty, |x| \leq \beta n} u_n(x), \quad U^*(\beta) := \limsup_{n \rightarrow \infty, |x| \leq \beta n} u_n(x).$$

Let $c \in (0, c^*)$ be given. We fix a number $\gamma \in (c, c^*)$ and define

$$V_*(c, \gamma) = \inf_{c < \beta < \gamma} U_*(\beta), \quad V^*(c, \gamma) = \sup_{c < \beta < \gamma} U^*(\beta).$$

It then follows that

$$V_*(c, \gamma) \leq U_*(\beta) \leq U^*(\beta) \leq V^*(c, \gamma), \quad \forall \beta \in [c, \gamma].$$

By the conclusion in case 2, we have

$$0 < u_-^* \leq U_*(\beta) \leq U^*(\beta) \leq u_+^*, \quad \forall \beta \in (0, c^*),$$

and hence,

$$0 < u_-^* \leq V_*(c, \gamma) \leq V^*(c, \gamma) \leq u_+^* \leq b.$$

For any $\beta \in (c, \gamma)$, we choose two sequences $n_j \rightarrow \infty$ and $x_j \in \mathbb{R}$ with $|x_j| \leq \beta n_j$ such that $\lim_{j \rightarrow \infty} u_{n_j}(x_j) = U_*(\beta)$. It is easy to see that

$$U_*(\gamma) \leq \liminf_{j \rightarrow \infty} u_{n_j-1}(x_j - y) \leq \limsup_{j \rightarrow \infty} u_{n_j-1}(x_j - y) \leq U^*(\gamma), \quad \forall y \in \mathbb{R}.$$

Since $u_{n_j}(x_j) = \int_{\mathbb{R}} g(u_{n_j-1}(x_j - y), u_{n_j-1}(x_j - y))k(y)dy$, it follows from Fatou's lemma that

$$U_*(\beta) \geq \int_{\mathbb{R}} \liminf_{j \rightarrow \infty} g(u_{n_j-1}(x_j - y), u_{n_j-1}(x_j - y))k(y)dy,$$

and hence

$$\begin{aligned} U_*(\beta) &\geq \int_{\mathbb{R}} g(U_*(\gamma), U^*(\gamma))k(y)dy \\ &= g(U_*(\gamma), U^*(\gamma)) \geq g(V_*(c, \gamma), V^*(c, \gamma)). \end{aligned}$$

Similarly, we have

$$U^*(\beta) \leq g(U^*(\gamma), U_*(\gamma)) \leq g(V^*(c, \gamma), V_*(c, \gamma)).$$

Thus,

$$(2.5) \quad V_*(c, \gamma) \geq g(V_*(c, \gamma), V^*(c, \gamma)), \quad V^*(c, \gamma) \leq g(V^*(c, \gamma), V_*(c, \gamma)).$$

By the definition of function g , we can find $v, w \in [V_*(c, \gamma), V^*(c, \gamma)] \subset (0, b]$ such that

$$g(V^*(c, \gamma), V_*(c, \gamma)) = f(v) \quad \text{and} \quad g(V_*(c, \gamma), V^*(c, \gamma)) = f(w).$$

It then follows from (2.5) that

$$(2.6) \quad f(w) \leq V_*(c, \gamma) \leq v, \quad w \leq V^*(c, \gamma) \leq f(v),$$

and hence,

$$\frac{f(w)}{w} \leq 1 = \frac{f(u^*)}{u^*} \leq \frac{f(v)}{v}.$$

This, together with the strict monotonicity of $\frac{f(u)}{u}$ on $(0, b]$, implies that $v \leq u^* \leq w$. By (2.6) and the property (P), we obtain $v = w$. It then follows from (2.6) that $V_*(c, \gamma) \geq f(w) = f(v) \geq V^*(c, \gamma)$, and hence, $0 < V_*(c, \gamma) = V^*(c, \gamma) \leq b$. From (2.5), we have

$$V_*(c, \gamma) \geq g(V_*(c, \gamma), V_*(c, \gamma)) \geq V^*(c, \gamma),$$

and hence, $0 < V_*(c, \gamma) = f(V_*(c, \gamma))$. By the uniqueness of the positive fixed point of f in $[0, b]$, it follows that $V_*(c, \gamma) = u^*$. Consequently,

$$u^* = V_*(c, \gamma) \leq U_*(c) \leq U^*(c) \leq V^*(c, \gamma) = u^*,$$

which implies that $\lim_{n \rightarrow \infty, |x| \leq cn} u_n(x) = u^*$ for any $c \in (0, c^*)$. \square

Remark 2.1. Under the assumption that $\int_{\mathbb{R}^m} k(y)dy = 1, k(x) = k(y), \forall x, y \in \mathbb{R}^m$ with $|x| = |y|$, Theorem 2.2 is still valid if we replace $\int_{\mathbb{R}}$ with $\int_{\mathbb{R}^m}$.

3. Traveling waves. In this section, we establish the existence and nonexistence of traveling waves for systems (2.3) by appealing to the Schauder fixed point theorem and the property of the spreading speed.

For a given $\lambda > 0$, let

$$X_\lambda := \{ \phi \in C(\mathbb{R}, \mathbb{R}) : \sup_{\xi \in \mathbb{R}} |\phi(\xi)|e^{-\lambda\xi} < +\infty \}$$

and $\|\phi\|_\lambda = \sup_{\xi \in \mathbb{R}} |\phi(\xi)|e^{-\lambda\xi}$. It then follows that $(X_\lambda, \|\cdot\|_\lambda)$ is a Banach space.

Define

$$\Phi(\lambda) = \frac{\ln(f'(0) \int_{\mathbb{R}} e^{-\lambda y} k(y) dy)}{\lambda}, \quad \forall \lambda \in (0, \Delta),$$

and

$$K(c, \lambda) = f'(0) e^{-c\lambda} \int_{\mathbb{R}} e^{-\lambda y} k(y) dy, \quad \forall c \in \mathbb{R}_+, \lambda \in (0, \Delta).$$

By [11, Lemma 3.8], it follows that for any $c > c_f^*$, there exist $0 < \lambda_1 = \lambda_1(c) < \lambda_2 = \lambda_2(c) < \Delta$ such that $\Phi(\lambda_1) = c$ and $\Phi(\lambda) < c, \forall \lambda \in (\lambda_1, \lambda_2)$. Thus, we have

$$K(c, \lambda_1) = 1, \quad K(c, \lambda) < 1, \quad \forall \lambda \in (\lambda_1, \lambda_2).$$

Note that if $f''(0)$ exists, then $f(u) \geq f'(0)u - au^2, \forall u \in [0, \delta]$, for appropriate $a > 0$ and $\delta > 0$. To obtain the existence of traveling waves, we impose the following weaker condition on f (cf. [3]).

(F3) There exist real numbers $\delta^* \in (0, \min(\delta_0, u_+^*)]$, $\sigma > 1$ and $a > 0$ such that $f(u) \geq f'(0)u - au^\sigma, \forall u \in [0, \delta^*]$.

THEOREM 3.1. Let (F1)–(F3) hold. Then the following statements are valid:

- (1) For any $c \in (0, c_f^*)$, (2.3) has no traveling wave $U(x + cn)$ with $U \in \mathcal{C}_{u_+^*} \setminus \{0\}$ and $U(-\infty) = 0$.
- (2) For any $c > c_f^*$, (2.3) has a traveling wave $U(x + cn)$ such that $U \in \mathcal{C}_{u_+^*} \setminus \{0\}$, $U(-\infty) = 0$, and

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

If, in addition, either (C1) or (C2) holds, then $U(+\infty) = u^*$.

Proof. Case 1. Assume, by contradiction, that for some $c_0 \in (0, c_f^*)$, (2.3) has a traveling wave $u_n(x) := U(x + c_0n)$ with $U \in \mathcal{C}_{u_+^*} \setminus \{0\}$ and $U(-\infty) = 0$. By Theorem 2.2 (2), there holds

$$\liminf_{n \rightarrow \infty, |x| \leq cn} u_n(x) \geq u_-^* > 0, \quad \forall c \in (0, c_f^*).$$

Choose $\tilde{c} \in (c_0, c_f^*)$ and let $x = -\tilde{c}n$. Then $\liminf_{n \rightarrow \infty} u_n(-\tilde{c}n) = \liminf_{n \rightarrow \infty} U((c_0 - \tilde{c})n) > 0$, but $\lim_{n \rightarrow \infty} U((c_0 - \tilde{c})n) = U(-\infty) = 0$, a contradiction.

Case 2. Let $c > c_f^*$ be given. Define a mapping $T : \mathcal{C}_b \rightarrow \mathcal{C}_b$ by

$$(3.1) \quad T(\phi)(\xi) = \int_{\mathbb{R}} f(\phi(\xi - c - y))k(y)dy, \quad \forall \xi \in \mathbb{R}, \phi \in \mathcal{C}_b.$$

Let T^\pm be defined as in (3.1) with f replaced by f^\pm . It then follows that T^\pm is nondecreasing with respect to the pointwise ordering on \mathcal{C}_b , and that

$$T^-(\phi) \leq T(\phi) \leq T^+(\phi), \quad \forall \phi \in \mathcal{C}_b.$$

Following [3], we define

$$\phi^+(\xi) := \min\{u_+^* e^{\lambda_1 \xi}, u_+^*\}, \quad \forall \xi \in \mathbb{R}.$$

Since $f^+(u)$ is nondecreasing in u and $\phi^+(\xi) \leq u_+^*, \forall \xi \in \mathbb{R}$, we obtain

$$(3.2) \quad T^+(\phi^+)(\xi) \leq \int_{\mathbb{R}} f^+(u_+^*)k(y)dy = f^+(u_+^*) = u_+^*, \quad \forall \xi \in \mathbb{R}.$$

Note that $f^+(u) \leq f'(0)u, \forall u \in [0, b]$, and $\phi^+(\xi) \leq u_+^* e^{\lambda_1 \xi}, \forall \xi \in \mathbb{R}$. Thus, we further have

$$\begin{aligned}
 T^+(\phi^+)(\xi) &\leq \int_{\mathbb{R}} f'(0)\phi^+(\xi - c - y)k(y)dy \\
 &\leq f'(0) \int_{\mathbb{R}} u_+^* e^{\lambda_1(\xi - c - y)}k(y)dy \\
 (3.3) \qquad &= u_+^* e^{\lambda_1 \xi} K(c, \lambda_1) = u_+^* e^{\lambda_1 \xi}, \forall \xi \in \mathbb{R}.
 \end{aligned}$$

By (3.2), (3.3), and the definition of ϕ^+ , it then follows that $T^+(\phi^+) \leq \phi^+$. We first fix a sufficiently small $\epsilon^* \in (0, \lambda_1(\sigma - 1)]$ such that $\lambda_1 + \epsilon^* < \lambda_2$, and hence $K(c, \lambda_1 + \epsilon^*) < 1$. We then choose a sufficiently large number $M \geq 1$ such that

$$(3.4) \qquad \left(1 + \frac{a(u_+^*)^\sigma}{f'(0)\delta^* M}\right) K(c, \lambda_1 + \epsilon^*) < 1.$$

Following [19], we define

$$\phi^-(\xi) := \max\{0, \delta^*(1 - Me^{\epsilon^* \xi})e^{\lambda_1 \xi}\}, \quad \forall \xi \in \mathbb{R}.$$

Let $\xi_0 := -\frac{\ln M}{\epsilon^*}$. Then we have

$$\phi^-(\xi) = 0, \forall \xi \geq \xi_0, \quad \phi^-(\xi) = \delta^* e^{\lambda_1 \xi} - \delta^* M e^{(\lambda_1 + \epsilon^*)\xi}, \forall \xi \leq \xi_0.$$

Since $\delta^* \leq u_+^*, \xi_0 \leq 0$, and $0 < \epsilon^* \leq \lambda_1(\sigma - 1)$, it is easy to see that

$$0 \leq \phi^-(\xi) \leq \phi^+(\xi), \quad (\phi^-(\xi))^\sigma \leq (u_+^*)^\sigma e^{(\lambda_1 + \epsilon^*)\xi}, \quad \forall \xi \in \mathbb{R}.$$

Clearly, we have

$$(3.5) \qquad T^-(\phi^-)(\xi) \geq 0, \quad \forall \xi \in \mathbb{R}.$$

Since $\phi^-(\xi) \geq \delta^* e^{\lambda_1 \xi} - \delta^* M e^{(\lambda_1 + \epsilon^*)\xi}, \forall \xi \in \mathbb{R}$, it follows from (3.4) that

$$\begin{aligned}
 f'(0)\phi^-(\xi) - a(\phi^-(\xi))^\sigma &\geq f'(0)\delta^* e^{\lambda_1 \xi} - f'(0)\delta^* M e^{(\lambda_1 + \epsilon^*)\xi} - a(u_+^*)^\sigma e^{(\lambda_1 + \epsilon^*)\xi} \\
 &= f'(0)\delta^* e^{\lambda_1 \xi} - f'(0)\delta^* M e^{(\lambda_1 + \epsilon^*)\xi} \left(1 + \frac{a(u_+^*)^\sigma}{f'(0)\delta^* M}\right) \\
 &\geq f'(0)\delta^* e^{\lambda_1 \xi} - \frac{f'(0)\delta^* M}{K(c, \lambda_1 + \epsilon^*)} e^{(\lambda_1 + \epsilon^*)\xi}, \quad \forall \xi \in \mathbb{R}.
 \end{aligned}$$

In view of (F3) and the fact that $f^-(u) = f(u), \forall u \in [0, \delta_0]$, we then have

$$\begin{aligned}
 T^-(\phi^-)(\xi) &\geq \int_{\mathbb{R}} (f'(0)\phi^-(\xi - c - y) - a(\phi^-(\xi - c - y))^\sigma) k(y)dy \\
 &\geq \int_{\mathbb{R}} \left(f'(0)\delta^* e^{\lambda_1(\xi - c - y)} - \frac{f'(0)\delta^* M}{K(c, \lambda_1 + \epsilon^*)} e^{(\lambda_1 + \epsilon^*)(\xi - c - y)}\right) k(y)dy \\
 &= \delta^* e^{\lambda_1 \xi} K(c, \lambda_1) - \frac{\delta^* M}{K(c, \lambda_1 + \epsilon^*)} e^{(\lambda_1 + \epsilon^*)\xi} K(c, \lambda_1 + \epsilon^*) \\
 (3.6) \qquad &= \delta^* e^{\lambda_1 \xi} - \delta^* M e^{(\lambda_1 + \epsilon^*)\xi}, \quad \forall \xi \in \mathbb{R}.
 \end{aligned}$$

By (3.5), (3.6), and the definition of ϕ^- , it follows that $T^-(\phi^-) \geq \phi^-$.

Now we fix a number $\lambda \in (0, \lambda_1)$. It is easy to see that both ϕ^- and ϕ^+ are elements in X_λ . Thus, the set

$$Y := \{\phi \in X_\lambda : \phi^-(\xi) \leq \phi(\xi) \leq \phi^+(\xi), \forall \xi \in \mathbb{R}\}$$

is a nonempty, closed, and convex subset of X_λ . For any $\phi \in Y$, we have

$$\phi^- \leq T^-(\phi^-) \leq T^-(\phi) \leq T(\phi) \leq T^+(\phi) \leq T^+(\phi^+) \leq \phi^+,$$

and hence $T(Y) \subset Y$. For any $\phi, \psi \in Y$, there holds

$$\begin{aligned} \|T(\phi) - T(\psi)\|_\lambda &= \sup_{\xi \in \mathbb{R}} |T(\phi)(\xi) - T(\psi)(\xi)| e^{-\lambda \xi} \\ &\leq L \cdot \sup_{\xi \in \mathbb{R}} \int_{\mathbb{R}} |\phi(\xi - c - y) - \psi(\xi - c - y)| e^{-\lambda \xi} k(y) dy \\ &\leq L \|\phi - \psi\|_\lambda \int_{\mathbb{R}} e^{-\lambda(c+y)} k(y) dy \\ &= \left(L e^{-\lambda c} \int_{\mathbb{R}} e^{-\lambda y} k(y) dy \right) \|\phi - \psi\|_\lambda. \end{aligned}$$

This implies that $T : Y \rightarrow Y$ is continuous. We further show that $T(Y)$ is precompact in X_λ . For any $\phi \in Y, \xi_1, \xi_2 \in \mathbb{R}$, we have

$$\begin{aligned} |T(\phi)(\xi_1) - T(\phi)(\xi_2)| &= \left| \int_{\mathbb{R}} f(\phi(z)) (k(\xi_1 - c - z) - k(\xi_2 - c - z)) dz \right| \\ &\leq b \int_{\mathbb{R}} |k(\xi_1 - c - z) - k(\xi_2 - c - z)| dz \\ &= b \int_{\mathbb{R}} |k(\xi_1 - \xi_2 + y) - k(y)| dy \\ &= b \cdot g(\xi_1 - \xi_2), \end{aligned}$$

where $g(\xi) = \int_{\mathbb{R}} |k(\xi + y) - k(y)| dy, \forall \xi \in \mathbb{R}$. Since $\lim_{\xi \rightarrow 0} g(\xi) = 0$, it follows that the family of functions $\{T(\phi)(\xi) : \phi \in Y\}$ is uniformly bounded and equicontinuous in $\xi \in \mathbb{R}$. Thus, for any given sequence $\{\psi_n\}_{n \geq 1}$ in $T(Y)$, there exist $n_k \rightarrow \infty$ and $\psi \in C(\mathbb{R}, \mathbb{R})$ such that $\lim_{k \rightarrow \infty} \psi_{n_k}(\xi) = \psi(\xi)$ uniformly for ξ in any compact subset of \mathbb{R} . Since $\phi^-(\xi) \leq \psi_{n_k}(\xi) \leq \phi^+(\xi), \forall \xi \in \mathbb{R}$, we have $\phi^-(\xi) \leq \psi(\xi) \leq \phi^+(\xi), \forall \xi \in \mathbb{R}$, and hence, $\psi \in Y$. Note that

$$\lim_{\xi \rightarrow +\infty} (\phi^+(\xi) - \phi^-(\xi)) e^{-\lambda \xi} = 0$$

and

$$\lim_{\xi \rightarrow -\infty} (\phi^+(\xi) - \phi^-(\xi)) e^{-\lambda \xi} = 0.$$

Therefore, for any $\epsilon > 0$, there exists $B > 0$ such that

$$0 \leq (\phi^+(\xi) - \phi^-(\xi)) e^{-\lambda \xi} < \epsilon, \quad \forall |\xi| \geq B.$$

Since $\lim_{k \rightarrow \infty} (\psi_{n_k}(\xi) - \psi(\xi)) e^{-\lambda \xi} = 0$ uniformly for $\xi \in [-B, B]$, there exists an integer $N > 0$ such that

$$|\psi_{n_k}(\xi) - \psi(\xi)| e^{-\lambda \xi} < \epsilon, \quad \forall \xi \in [-B, B], k \geq N.$$

It then follows that

$$\|\psi_{n_k} - \psi\|_\lambda = \sup_{\xi \in \mathbb{R}} |\psi_{n_k}(\xi) - \psi(\xi)|e^{-\lambda\xi} \leq \epsilon, \quad \forall k \geq N.$$

This implies that $\lim_{k \rightarrow \infty} \psi_{n_k} = \psi$ in X_λ . By the Schauder fixed point theorem, there exists $U \in Y$ such that $U = T(U)$, and hence, $U(x + cn)$ is a traveling wave of (2.3). Since $\phi^-(\xi) \leq U(\xi) \leq \phi^+(\xi), \forall \xi \in \mathbb{R}$, we have $U(-\infty) = 0$ and $U \in \mathcal{C}_{u_+^*} \setminus \{0\}$.

Let $u_n(x) := U(x + cn), \forall n \geq 0$, and fix a number $\bar{c} \in (0, c_f^*)$. By Theorem 2.2 (2), it follows that

$$0 < u_-^* \leq \liminf_{n \rightarrow \infty, |x| \leq \bar{c}n} u_n(x) \leq \limsup_{n \rightarrow \infty, |x| \leq \bar{c}n} u_n(x) \leq u_+^*,$$

and hence,

$$u_-^* \leq \liminf_{n \rightarrow \infty} u_n(-\gamma n) \leq \limsup_{n \rightarrow \infty} u_n(-\gamma n) \leq u_+^*$$

uniformly for $\gamma \in [0, \bar{c}]$. This implies that

$$u_-^* \leq \liminf_{n \rightarrow \infty} U(sn) \leq \limsup_{n \rightarrow \infty} U(sn) \leq u_+^*$$

uniformly for $s \in [c - \bar{c}, c]$. Let

$$a_n = n(c - \bar{c}), \quad b_n = nc, \quad \forall n \geq 1.$$

Thus, there exists $N_0 > 0$ such that $a_{n+1} - b_n < 0, \forall n \geq N_0$, and hence,

$$\cup_{n \geq m} [a_n, b_n] = [a_m, +\infty), \quad \forall m \geq N_0.$$

It then follows that

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

If, in addition, either (C1) or (C2) holds, then Theorem 2.2 (3) implies that

$$\lim_{n \rightarrow \infty, |x| \leq \bar{c}n} u_n(x) = u^*, \quad \forall \bar{c} \in (0, c_f^*).$$

By the same arguments as above, we further have $U(+\infty) = u^*$. \square

THEOREM 3.2. *Let (F1)–(F3) hold. Then (2.3) has a traveling wave $U(x + c_f^*n)$ such that $U \in \mathcal{C}_{u_+^*} \setminus \{0, u^*\}$ and*

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

If, in addition, either (C1) or (C2) holds, then $U(+\infty) = u^$.*

Proof. Choose a sequence $\{c_j\}_{j \geq 1} \subset (c_f^*, +\infty)$ such that $\lim_{j \rightarrow \infty} c_j = c_f^*$. By Theorem 3.1 (2), it follows that (2.3) has a traveling wave $U_j(x + c_jn)$ such that $U_j \in \mathcal{C}_{u_+^*} \setminus \{0\}, U_j(-\infty) = 0$, and

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U_j(\xi) \leq \limsup_{\xi \rightarrow +\infty} U_j(\xi) \leq u_+^*.$$

Without loss of generality, we assume that $U_j(0) = \frac{1}{2}u_-^* > 0, \forall j \geq 1$. Note that

$$(3.7) \quad U_j(\xi) = \int_{\mathbb{R}} f(U_j(\xi - c_j - y))k(y)dy, \quad \forall \xi \in \mathbb{R}, j \geq 1.$$

It follows that

$$|U_j(\xi_1) - U_j(\xi_2)| \leq b \cdot g(\xi_1 - \xi_2), \quad \forall \xi_1, \xi_2 \in \mathbb{R}, j \geq 1,$$

where $g(\xi)$ is defined as in the proof of Theorem 3.1. Then the family of functions $\{U_j(\xi) : j \geq 1\}$ is uniformly bounded and equicontinuous in $\xi \in \mathbb{R}$. Thus, there exist $j_k \rightarrow +\infty$ and $U \in C(\mathbb{R}, \mathbb{R})$ such that $\lim_{k \rightarrow \infty} U_{j_k}(\xi) = U(\xi)$ uniformly for ξ in any compact subset of \mathbb{R} . Clearly, $U \in \mathcal{C}_{u_+^*}$ and $U(0) = \frac{1}{2}u_-^*$. Letting $j = j_k \rightarrow +\infty$ in (3.7) and using the dominated convergence theorem, we obtain

$$U(\xi) = \int_{\mathbb{R}} f(U(\xi - c_f^* - y))k(y)dy, \quad \forall \xi \in \mathbb{R},$$

and hence, $u_n(x) := U(x + c_f^*n)$ is a traveling wave of (2.3). As in the proof of Theorem 3.1 (2), we see that Theorem 2.2 (2) and (3) imply the asymptotic behavior of $U(\xi)$ as $\xi \rightarrow +\infty$. \square

Compared with Theorem 3.1 (2), we expect that $U(-\infty) = 0$ in Theorem 3.2. However, we are not able to prove it at this moment since the limiting function $U(\xi)$ may not be nondecreasing on \mathbb{R} .

4. Examples. In this section, we present illustrative examples by choosing three types of growth functions from population biology.

First, we consider the logistic type function $f(u) = ru(1 - \frac{u}{K}), r > 0, K > 0$. Clearly, $f'(0) = r, \max_{u \in [0, K]} f(u) = f(K/2) = \frac{rK}{4}, \frac{f(u)}{u} = r(1 - \frac{u}{K})$ is strictly decreasing on $(0, K]$, and $u^* := K(1 - \frac{1}{r})$ is the unique positive fixed point of f on $[0, K]$. Assume that $1 < r < 4$ so that we have $f'(0) > 1$ and $f((0, K]) \subset (0, K]$. It is easy to verify that $f(u)$ is strictly increasing on $[0, u^*]$ if $r \in (0, 2]$. In the case where $r \in (1, 2]$, we choose $b =: u^*$, and hence, $u_-^* = u_+^* = u^*$. In the case $r \in (2, 4)$, we choose $b =: \frac{rK}{4}$, and hence, $u_+^* = b, u_-^* = f(b) = \frac{r^2K(4-r)}{16}$. Note that

$$\frac{f^2(u)}{u} = \frac{r^2}{K^3} (K^2(K - u) - ru(K - u)^2).$$

It then follows that $f(u)$ satisfies the property (P2) if $r \in (2, 3]$. By Theorems 2.1 and 2.2 and Theorems 3.1 and 3.2, we have the following result.

EXAMPLE 4.1. *Let $f(u) = ru(1 - \frac{u}{K})$ with $K > 0$ and $r \in (1, 4)$, b, u_+^* and u_-^* be defined as above, and c_f^* be defined as in (2.2) with $h = f$. Then the following statements are valid:*

- (i) c_f^* is the spreading speed of (2.3) in the sense that both conclusions (1) and (2) in Theorem 2.2 hold. Further, the conclusion (3) in Theorem 2.2 holds in the case where $r \in (1, 3]$.
- (ii) For any $c \in (0, c_f^*)$, (2.3) has no traveling wave $U(x + cn)$ with $U \in \mathcal{C}_b \setminus \{0\}$ and $U(-\infty) = 0$, and for any $c \geq c_f^*$, (2.3) has a traveling wave $U(x + cn)$ with $U \in \mathcal{C}_{u_+^*} \setminus \{0, u^*\}$ and

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

Further, $U(+\infty) = u^*$ in the case where $r \in (1, 3]$. If $r \in (1, 2]$, then $U(-\infty) = 0$ and $U(\xi)$ is nondecreasing in ξ for all $c \geq c_f^*$. If $r \in (2, 3]$, then $U(-\infty) = 0$ for all $c > c_f^*$.

In Example 4.1, we can also verify that for any $r > 0$, $uf(u)$ is strictly increasing on $[0, 2K/3]$ and strictly decreasing on $[2K/3, +\infty)$. It then follows that $f(u)$ satisfies the property (P1) if $r \in (2, 8/3]$, but does not satisfy the property (P1) if $r \in (8/3, 4)$. So we chose to use the property (P2) to obtain the upward convergence as stated in Theorem 2.2 (3) for r in a larger interval $(1, 3]$. By taking u_0 as a constant function in integrodifference equation (2.3), we see that the upward convergence implies that u^* is a globally attractive fixed point for the map f on $(0, u_+^*]$. Note that $|f'(u^*)| = |2 - r| > 1, \forall r \in (3, 4)$. By [5, Theorem 3.8], it then follows that for any $r \in (3, 4)$, u^* is a unstable fixed point of f . This implies that the upward convergence does not hold for any $r \in (3, 4)$. Thus, the interval $(1, 3]$ for parameter r is optimal for the upward convergence.

Among other things, Kot [7] observed numerically four types of discrete-time traveling waves for the integrodifference equation (2.3) with $f(u) = (1 + r_0)u - r_0u^2$ and $k(x) = 3e^{-6|x|}$: a simple monotone traveling for $r_0 = 0.9$ ([7, Figure 7]); a traveling wave with damped spatial oscillations for $r_0 = 1.9$ ([7, Figure 8]); a traveling two-cycle for $r_0 = 2.2$ ([7, Figures 9a,b]); and a traveling four-cycle for $r_0 = 2.5$ ([7, Figures 10a-d]). Clearly, this system is a special case of Example 4.1 with $r = 1 + r_0$ and $K = (1 + r_0)/r_0$. It is easy to see that our analytic results in Example 4.1 are consistent with these numerical simulations. Note that there is an increasing sequence of parameter values $r_1 = 3 < r_2 \approx 3.449 < r_3 \approx 3.544 < r_4 \approx 3.564 < \dots$ at which the logistic map $f(u) = ru(1 - \frac{u}{K})$ repeatedly undergoes a period-doubling bifurcation (see, e.g., [5, section 3.5]). Further, when $r \approx 3.839$, f has a unique asymptotically stable periodic orbit of minimal period 3. It is a challenging problem to prove the existence of a traveling three-cycle for the integrodifference equation (2.3) associated with the logistic map when $r \approx 3.839$.

Next, we consider the Ricker type function $f(u) = que^{-pu}$, $q > 1, p > 0$. Clearly, $f'(0) = q$ and $\frac{f(u)}{u} = qe^{-pu}$ is strictly decreasing on $(0, +\infty)$. It is easy to see that $u^* := \frac{\ln q}{p}$ is the unique positive fixed point of f on $[0, +\infty)$, that $\max_{u \in [0, +\infty)} f(u) = f(1/p) = \frac{q}{pe}$, and that $f(u)$ is strictly increasing on $[0, u^*]$ if $q \in (1, e]$. In the case where $q \in (1, e]$, we choose $b := u^*$, and hence, $u_-^* = u_+^* = u^*$. In the case where $q > e$, we choose $b := \frac{q}{pe}$, and hence, $u_+^* = b, u_-^* = f(b) = \frac{q^2}{pe}e^{-q/e}$. Note that

$$\frac{f^2(u)}{u} = q^2e^{-p(u+que^{-pu})}.$$

An elementary analysis shows that $f(u)$ satisfies the property (P2) if $q \in (e, e^2]$. By Theorems 2.1 and 2.2 and Theorems 3.1 and 3.2, we have the following result.

EXAMPLE 4.2. Let $f(u) = que^{-pu}$ with $q > 1$ and $p > 0$, b, u_+^* and u_-^* be defined as above, and c_f^* be defined as in (2.2) with $h = f$. Then the following statements are valid:

- (i) c_f^* is the spreading speed of (2.3) in the sense that both conclusions (1) and (2) in Theorem 2.2 hold. Further, conclusion (3) in Theorem 2.2 holds in the case where $q \in (1, e^2]$.
- (ii) For any $c \in (0, c_f^*)$, (2.3) has no traveling wave $U(x + cn)$ with $U \in C_b \setminus \{0\}$ and $U(-\infty) = 0$, and for any $c \geq c_f^*$, (2.3) has a traveling wave $U(x + cn)$

with $U \in \mathcal{C}_{u_+^*} \setminus \{0, u^*\}$ and

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

Further, $U(+\infty) = u^*$ in the case where $q \in (1, e^2]$. If $q \in (1, e]$, then $U(-\infty) = 0$ and $U(\xi)$ is nondecreasing in ξ for all $c \geq c_f^*$. If $q \in (e, e^2]$, then $U(-\infty) = 0$ for all $c > c_f^*$.

In Example 4.2, for any $q > e^2$, we have $|f'(u^*)| = |1 - \ln q| > 1$, and hence, u^* is an unstable fixed point of f . As discussed in Example 4.1, it follows that the interval $(1, e^2]$ for parameter q is optimal for the upward convergence.

Finally, we consider the generalized Beverton–Holt-type function $f(u) = \frac{pu}{q+u^m}$, $m > 0$, and $p > q > 0$. Clearly, $f'(0) = p/q$ and $\frac{f(u)}{u} = \frac{p}{q+u^m}$ is strictly decreasing on $(0, +\infty)$. It is easy to see that $u^* := (p - q)^{\frac{1}{m}}$ is the unique positive fixed point of f on $[0, +\infty)$, and that $f(u)$ is strictly increasing on $[0, +\infty)$ in the case where $m \in (0, 1]$. In the case where $m > 1$, we have

$$\max_{u \in [0, +\infty)} f(u) = f(\bar{u}) = \frac{p(m-1)\bar{u}}{qm}, \quad \bar{u} := \left(\frac{q}{m-1}\right)^{\frac{1}{m}}.$$

By elementary analysis, it follows that $f(u)$ is strictly increasing on $[0, u^*]$ if $m \in (1, p/(p - q)]$, that $uf(u)$ is strictly increasing on $[0, +\infty)$ if $m \in (0, 2]$, and that $uf(u)$ is strictly increasing on $[0, (2q/(m - 2))^{\frac{1}{m}}]$ if $m > 2$. Define $b := u^*$ if $m \in (0, p/(p - q)]$ and $b := f(\bar{u})$ if $m > p/(p - q)$. It then follows that $u_-^* = u_+^* = u^*$ in the case where $m \in (1, p/(p - q)]$, and

$$u_+^* = b, \quad u_-^* = f(b) = \frac{p^2(m-1)\bar{u}}{q^2m + \frac{p^m(m-1)^{m-1}q}{(qm)^{m-1}}},$$

in the case where $m > p/(p - q)$. Note that $f(u)$ satisfies the property (P1) if either $m \in (0, 2]$, or $m > \max(2, p/(p - q))$ and $f(\bar{u}) \leq (2q/(m - 2))^{\frac{1}{m}}$. By Theorems 2.1 and 2.2 and Theorems 3.1 and 3.2, we have the following result.

EXAMPLE 4.3. Let $f(u) = \frac{pu}{q+u^m}$ with $m > 0$ and $p > q > 0$, b, u_+^* and u_-^* be defined as above, and c_f^* be defined as in (2.2) with $h = f$. Then the following statements are valid:

- (i) c_f^* is the spreading speed of (2.3) in the sense that both conclusions (1) and (2) in Theorem 2.2 hold. Further, conclusion (3) in Theorem 2.2 holds in the case where either $m \in (0, \max(2, p/(p - q))]$, or $m > \max(2, p/(p - q))$ and $f(\bar{u}) \leq (2q/(m - 2))^{\frac{1}{m}}$.
- (ii) For any $c \in (0, c_f^*)$, (2.3) has no traveling wave $U(x + cn)$ with $U \in \mathcal{C}_b \setminus \{0\}$ and $U(-\infty) = 0$, and for any $c \geq c_f^*$, (2.3) has a traveling wave $U(x + cn)$ with $U \in \mathcal{C}_{u_+^*} \setminus \{0, u^*\}$ and

$$u_-^* \leq \liminf_{\xi \rightarrow +\infty} U(\xi) \leq \limsup_{\xi \rightarrow +\infty} U(\xi) \leq u_+^*.$$

Further, $U(+\infty) = u^*$ in the case where either $m \in (0, \max(2, p/(p - q))]$, or $m > \max(2, p/(p - q))$ and $f(\bar{u}) \leq (2q/(m - 2))^{\frac{1}{m}}$. If $m \in (0, p/(p - q)]$, then $U(-\infty) = 0$ and $U(\xi)$ is nondecreasing in ξ for all $c \geq c_f^*$. If either $p/(p - q) < m \leq 2$, or $m > \max(2, p/(p - q))$ and $f(\bar{u}) \leq (2q/(m - 2))^{\frac{1}{m}}$, then $U(-\infty) = 0$ for all $c > c_f^*$.

Acknowledgments. We are grateful to two anonymous referees for their careful reading and helpful suggestions which led to an improvement of our original manuscript. Xiao-Qiang Zhao would like to thank the National Center for Theoretical Science, Tsing Hua University, Taiwan for its kind hospitality during his visit there.

REFERENCES

- [1] D. G. ARONSON AND H. F. WEINBERGER, *Nonlinear diffusion in population genetics, combustion, and nerve pulse propagation*, in Partial Differential Equations and Related Topics, J. A. Goldstein, ed., Lecture Notes in Mathematics Ser. 446, Springer-Verlag, Berlin, 1975, pp. 5–49.
- [2] D. G. ARONSON AND H. F. WEINBERGER, *Multidimensional nonlinear diffusion arising in population dynamics*, Adv. Math., 30 (1978), pp. 33–76.
- [3] O. DIEKMANN, *Thresholds and traveling waves for the geographical spread of infection*, J. Math. Biol., 6 (1978), pp. 109–130.
- [4] T. FARIA, W. HUANG, AND J. WU, *Traveling waves for delayed reaction-diffusion equations with global response*, Proc. Roy. Soc. Lond. Ser. A., 462 (2006), pp. 229–261.
- [5] J. K. HALE AND H. KOCAK, *Dynamics and Bifurcations*, Springer-Verlag, New York, 1991.
- [6] A. HASTINGS, K. CUDDINGTON, K. F. DAVIES, C. J. DUGAW, S. ELMENDORF, A. FREESTONE, S. HARRISON, M. HOLLAND, J. LAMBRINOS, U. MALVADKAR, B. A. MELBOURNE, K. MOORE, C. TAYLOR, AND D. THOMSON, *The spatial spread of invasions: New developments in theory and evidence*, Ecology Lett., 8 (2005), pp. 91–101.
- [7] M. KOT, *Discrete-time traveling waves: Ecological examples*, J. Math. Biol., 30 (1992), pp. 413–436.
- [8] J. M. LEVINE, E. PACHEPSKY, B. E. KENDALL, S. G. YELENIK, AND J. H. R. LAMBERS, *Plant-soil feedbacks and invasive spread*, Ecology Lett., 9 (2006), pp. 1005–1014.
- [9] M. KOT, M. A. LEWIS, AND P. VAN DEN DRIESSCHE, *Dispersal data and the spread of invading organisms*, Ecology, 77 (1996), pp. 2027–2042.
- [10] B. LI, H. F. WEINBERGER, AND M. A. LEWIS, *Spreading speeds as slowest wave speeds for cooperative systems*, Math. Biosci., 196 (2005), pp. 82–89.
- [11] X. LIANG AND X.-Q. ZHAO, *Asymptotic speeds of spread and traveling waves for monotone semiflows with applications*, Commun. Pure Appl. Math., 60 (2007), pp. 1–40. Erratum: 61 (2008), pp. 137–138.
- [12] X. LIANG, Y. YI, AND X.-Q. ZHAO, *Spreading speeds and traveling waves for periodic evolution systems*, J. Differential Equations, 231 (2006), pp. 57–77.
- [13] R. LUI, *Biological growth and spread modeled by systems of recursions, I. Mathematical theory*, Math. Biosci., 93 (1989), pp. 269–295.
- [14] S. MA, *Traveling waves for non-local delayed diffusion equations via auxiliary equations*, J. Differential Equations, 237 (2007), pp. 259–277.
- [15] M. NEUBERT AND H. CASWELL, *Demography and dispersal: Calculation and sensitivity analysis of invasion speed for structured populations*, Ecology, 81 (2000), pp. 1613–1628.
- [16] C. OU AND J. WU, *Persistence of wavefronts in delayed nonlocal reaction-diffusion equations*, J. Differential Equations, 235 (2007), pp. 219–261.
- [17] H. R. THIEME, *Density-dependent regulation of spatially distributed populations and their asymptotic speed of spread*, J. Math. Biol., 8 (1979), pp. 173–187.
- [18] H. R. THIEME, *On a class of Hammerstein integral equations*, Manuscripta Math., 29 (1979), pp. 49–84.
- [19] H. R. THIEME AND X.-Q. ZHAO, *Asymptotic speeds of spread and traveling waves for integral equations and delayed reaction-diffusion models*, J. Differential Equations, 195 (2003), pp. 430–470.
- [20] H. F. WEINBERGER, *Long-time behavior of a class of biological models*, SIAM J. Math. Anal., 13 (1982), pp. 353–396.
- [21] H. F. WEINBERGER, M. A. LEWIS, AND B. LI, *Analysis of linear determinacy for spread in cooperative models*, J. Math. Biol., 45 (2002), pp. 183–218.
- [22] J. WU AND X. ZOU, *Traveling wave fronts of reaction-diffusion systems with delay*, J. Dynam. Differential Equations, 13 (2001), pp. 651–687.

ESTIMATES ON THE HAUSDORFF DIMENSION OF THE RUPTURE SET OF A THIN FILM*

KAI-SENG CHOU[†] AND SHI-ZHONG DU[†]

Abstract. Upper bounds on the Hausdorff dimensions of the rupture set of a weak solution of the thin film equation in space-time and in space slices are derived. Finite time rupture is shown to occur for a class of thin films obeying the power law with power in $(0, 1/2)$ under periodic boundary conditions.

Key words. Hausdorff measure, Hausdorff dimension, finite time rupture, thin film equation, partial regularity of a suitable weak solution

AMS subject classifications. 35Q35, 76A20, 35B35, 93D20

DOI. 10.1137/070685348

1. Introduction. In recent years the thin film type equations

$$(1.1) \quad h_t + (h^n h_{xxx})_x = 0, \quad n > 0, \quad h \geq 0,$$

and

$$(1.2) \quad h_t + (h^n (h_{xx} + f(h)))_x = 0, \quad n > 0, \quad h \geq 0,$$

have attracted much attention. These equations, where h describes the height of an axisymmetric thin film in motion, arise from the lubrication approximation of the Navier–Stokes equations when n is equal to 3 under no-slip boundary conditions at the lower boundary. It also describes the motion of fluid in a Hele–Shaw cell when n is equal to 1. The f -term in the second equation describes the presence of other physical effects such as gravity, the van der Waals forces, and thermocapillary effect. It is called a power law if

$$(1.3) \quad f(z) = \begin{cases} B \frac{z^q}{q}, & q \neq 0, \\ B \log z, & q = 0, \end{cases}$$

for some positive B . One may consult the surveys Myers [M] and Oron, Davis, and Bankoff [ODB] for background on these equations and various models. Here we are concerned with the analytic aspects of these equations. To study them one needs to impose initial and boundary conditions so that they become well-posed at least for a short time. In the literature several boundary conditions such as Neumann-type conditions, pressure boundary conditions, and periodic boundary conditions have been used. Throughout this paper we shall employ the periodic boundary conditions. Thus given a nonnegative, periodic function h_0 , we would like to consider the solution to (1.1) starting at h_0 , which is of the same period as h_0 . Observing that the equation is a parabolic equation of fourth order when h is positive, it follows from parabolic theory that, for any positive, sufficiently regular initial datum, the equation admits a classical solution for small time, and it continues to exist as long as the solution

*Received by the editors March 14, 2007; accepted for publication (in revised form) March 11, 2008; published electronically August 20, 2008. This research was supported by an Earmarked Grant for Research, Hong Kong.

<http://www.siam.org/journals/sima/40-2/68534.html>

[†]Department of Mathematics, The Chinese University of Hong Kong, Hong Kong (kschou@math.cuhk.edu.hk, szdu@math.cuhk.edu.hk).

is bounded away from zero. However, unlike second order parabolic equations where positivity is ensured by the maximum principle in an *a priori* way, it is not clear at all whether the film will remain positive for all time, or it will touch down at 0 at some finite time. The film is said to rupture when the latter occurs, and the equation loses its parabolicity. The mathematical foundation of (1.1) was laid down in Bernis-Friedman [BF] where they discovered the following remarkable fact: When $n \geq 4$ equation (1.1) preserves positivity. (Later the requirement $n \geq 4$ was relaxed to $n \geq 3.5$ in Bertozzi et al. [BBDK].) Using this fact, Bernis and Friedman was able to construct a nonnegative, global weak solution of (1.1) for some n less than 4 for nonnegative, H^1 -data. After the works of Beretta, Bertsch, and Dal Passo [BBDP], Bertozzi-Pugh [BP1] and [BP2], nowadays it is known that nonnegative, global weak solutions exist for (1.1) and (1.2) for any positive n when f is in $C^1([0, \infty))$ and grows slower less than a cubic power in h at ∞ .

A natural question is, For n in $(0, 3.5)$, will rupture occur in finite time for some positive initial data? Observing that, when n is equal to 0, finite time rupture occurs for easily constructed initial data, it would not be surprising that finite time rupture occurs for positive, small n . In fact, under the boundary conditions $h = h_{xx} = 1$ at endpoints, such property is established for (1.1) when n is less than 1/2 in [BBDP]. However, for the periodic conditions this remains open. Numerical studies show that there is a critical number $n^* \in (1, 2)$ so that finite time rupture occurs for $n < n^*$, and positivity is preserved for $n > n^*$. Likewise, as for (1.2) under the power law $q < 3$ there should be a critical number $n_c(q)$, independent of B and the period, which separates finite time rupture and the preservation of positivity. Numerical results in Goldstein, Pesci, and Shelly [GPS] show that $1 \leq n_c(1) \leq 3$, and those in Laugesen-Pugh [LP2] show that $1.8 < n_c(0.5) < 1.85$ and $1.65 < n_c(2.5) < 1.6625$. Theoretical upper and lower bounds on $n_c(q)$ can be found in Chou-Kwong [CK] for negative q . In this paper we shall show that finite time rupture occurs for (1.2), where f is a power law (1.3) ($q \in (1, 3)$) under the periodic boundary conditions when n is less than 1/2; see Corollary 3.1 for a precise statement. We shall not consider (1.2) when f grows faster than a cubic power. In addition to the possibility of rupture, it is believed that, in this case the solution may also blow up in finite time. We refer to [BP1] for a conjecture in this direction, and [BP3] and Slepčev-Pugh [SP] for more recent progress.

In this paper we approach the problem of finite time rupture via estimating the possible size of the rupture set. When rupture really occurs, surely our results yield information on the size of the rupture set. On the other hand, even if rupture does not occur, one may still build on these results further criteria for the absence of rupture. A similar situation can be found in the study of the regularity of the weak solution of the three-dimensional Navier–Stokes equations. One may consult Escauriaza, Seregin, and Šverák [ESS] for how the partial regularity estimates in Caffarelli, Kohn, and Nirenberg [CKN] are used in establishing regularity criteria.

Let us review how one excludes finite time rupture for thin films when $n \geq 3.5$. Recall that the weak solution constructed in [BF], [BBDP], [BP1], and [BP2] has many regularity properties including the following two: For any positive initial data, the weak solution of (1.1) or (1.2) ($1 < q < 3$) satisfies, for any $0 < t < T$,

$$(1.4) \quad \int h_x^2(t) + \int_0^T \int h^n [(h_{xx} + f(h))_x]^2 \leq C_1,$$

and

$$(1.5) \quad \int h(t)^{3/2-n} \leq C_2$$

for some constants C_1 and C_2 depending on the initial data. It follows from the conservation of area

$$\int h(t) = \int h(0), \quad \text{all } t > 0,$$

and (1.4) that the solution h has a uniformly bounded $C^{1/2}$ -norm. At a rupture point (x_0, t) the solution is dominated by a constant multiple of $|x - x_0|^{1/2}$. Putting this estimate into (1.5) yields a contradiction when $n \geq 3.5$, so no finite rupture can occur when $n \geq 3.5$. Even when n is less than 3.5, (1.5) tells us that the rupture set has null measure at every time as long as $n > 1.5$. Our first result gives a more precise estimate.

THEOREM 1.1. *Let h be a suitable weak solution of (1.1) or (1.2), where $n \in (1.5, 3.5)$ and f is a power law with $q > 1 - n/2$. Then for any $t > 0$, the $(7 - 2n)/(2n + 1)$ -Hausdorff measure of the set $\{x : h(x, t) = 0\}$ is finite, so its Hausdorff dimension cannot exceed $(7 - 2n)/(2n + 1)$.*

Since weak solutions to (1.1) or (1.2) may not be unique, we use the terminology “a suitable weak solution” to refer to a certain weak solution constructed by the methods in [BF], [BBDP], and [BP1] and [BP2]. We shall review its definition in the next section and sketch a construction of these solutions in Appendix A. The proof of our results including Theorem 1.1 and those stated below involves the use of (1.5) in an essential way, so it does not work when n goes below 1.5. We point out that our bound on the Hausdorff dimension $(7 - 2n)/(2n + 1)$ is always less than 1 and is sharp in the sense that it tends to 0 and 1 as n tends to 3.5 and 1.5, respectively.

The proof of Theorem 1.1 follows rather easily from a covering argument coupling with properties (1.4) and (1.5). However, for the proof of the following size estimate in space-time, our main result, we need to make use of the further regularity properties of the suitable weak solution.

THEOREM 1.2. *Let h be a suitable weak solution of (1.1) or (1.2), where $n \in (1.5, 3.5)$ and f is a power law with $q > 1 - n/2$. The parabolic Hausdorff dimension of the rupture set $\{(x, t) : h(x, t) = 0\}$ cannot exceed $2n/(2n - 3)$ for $n \in [2, 3.5)$ and $6/n + 1$ for $n \in (1.5, 2)$.*

The parabolic Hausdorff measure, to be discussed in Appendix B, is defined through covering a set by cylinders with different weights in space and time. Here we take the parabolic fourth order weights, so the parabolic Hausdorff dimension is equal to 5 for every open set in space-time. Again we observe that our estimate on the dimension of the rupture set lies between 1.75 and 5 and tends to 1.75 and 5 as n tends to 3.5 and 1.5, respectively.

Finally, we give an estimate on the rupture times.

THEOREM 1.3. *Let h be a suitable weak solution of (1.1) or (1.2), where $n \in (2, 3.5)$ and f is a power law with $q > 1 - n/2$. The Hausdorff dimension of the set $\{t > 0 : h(x, t) = 0 \text{ for some } x\}$ cannot exceed $1 - 2(n - 2)^2/(8 - n)$.*

Again this estimate becomes sharp as n tends to 3.5.

Theorems 1.1–1.3 are contained in the results in sections 2 and 3, where we also consider more general f . In section 4 we establish the result on finite time rupture mentioned above.

The problem of rupture for thin films is closely related to the formation of singularities for solutions of (1.1) or (1.2). Indeed, let us assume that the solution remains to be C^2 at the rupture point (x_0, t) . Then it is bounded by a constant multiple of $|x - x_0|^2$ near x_0 . Putting this into (1.5) yields a contradiction unless $n \leq 2$. It demonstrates that the second derivatives of the solution blow up at any rupture point

when $n > 2$, so the rupture set coincides with the singular set of any weak solution of (1.1) or (1.2) in which (1.4) and (1.5) hold. Our proofs of Theorems 1.1, 1.2, and 1.3 make use of some ideas in known works in the partial regularity of solutions for evolution equations. There are quite a number of works on this topic in recent years. We wish to mention the works of Caffarelli–Kohn–Nirenberg [CKN] (later simplified substantially by Lin [L]) on the three-dimensional Navier–Stokes equations, Struwe [S], and Chen–Struwe [CS] on the harmonic heat flows, Lin–Liu [LL] on the motion of liquid crystals, and Chou–Du–Zheng [CDZ] on the semilinear heat equation with supercritical growth.

In Jiang–Lin [JL] estimates for the Hausdorff dimension of the rupture set of a steady state of a multidimensional thin film type equation are present.

Although we are mainly concerned with thin films whose initial data are not too degenerate in the sense that they have finite entropy (1.5), there are significant results for weak solutions starting from nonnegative initial data that could vanish in an open set. Our results do not apply to these solutions. By restricting to the literature on the one-dimensional case (axisymmetric films) only, we point out Bernis [B1], [B2], and Hulshof–Shishkov [HS], where results on the finite speed of the propagation of zero are established for (1.1). That such property continues to hold for (1.2) is explained in [BP2]. Moreover, a waiting time phenomenon which asserts that the support of a thin film does not increase in a small period of time is established for (1.1) ($n \in (0, 3)$) in Dal Passo–Giacomelli–Grun [DPGG]. For $n > 4$, it is known in [BBDP] that the support remains the same all the time.

2. A size estimate on the rupture set in space-time. First we give the definition of a suitable weak solution of the thin film type equation under the periodic boundary conditions. Consider the problem

$$(2.1) \quad \begin{cases} h_t + [h^n(h_{xx} + f(h))]_x = 0, & (x, t) \in (-L/2, L/2) \times (0, \infty), \\ h \text{ is of period } L \text{ for each } t > 0, \end{cases}$$

where n is positive and f is a measurable function in $[0, \infty)$ which belongs to $C^{2,\alpha}(0, \infty)$, $\alpha \in (0, 1)$. Let S_L denote the circle obtained by identifying the endpoints of $(-L/2, L/2)$. A nonnegative function h in

$$C^{\frac{1}{2}, \frac{1}{8}}_{x,t}(S_L \times (0, \infty)) \cap L^\infty(0, \infty; H^1(S_L)) \cap L^2_{loc}(0, \infty; H^2(S_L))$$

is called a suitable weak solution of (2.1) if it fulfills the following requirements. For every $T > 0$,

$$(H_1) \quad h^{\frac{n}{2}}(h_{xx} + f(h))_x \in L^2(Q_T), \quad Q_T = S_L \times (0, T), \text{ and}$$

$$\iint_{Q_T} h \Phi_t = - \iint_{Q_T} h^n (h_{xx} + f(h))_x \Phi_x$$

for all $\Phi \in H^1(Q_T)$ which are compactly supported in Q_T .

(H₂) The following quantities

$$\begin{aligned} K_1 &= \sup_{[0, T]} \int h_x^2, \\ K_2 &= \iint_{Q_T} \left(h^n |(h_{xx} + f(h))_x|^2 + h^n h_{xxx}^2 \right), \\ K_3 &= \sup_{[0, T]} \int h^{\frac{3}{2}-n}, \end{aligned}$$

and

$$K_4 = \iint_{Q_T} \left(h^{-\frac{1}{2}+\delta} h_{xx}^2 + h^{-\frac{5}{2}+\delta} h_x^4 \right) \quad \forall \delta \in (0, \min\{1.5, n - 1.5\})$$

are finite, where K_1-K_4 depend on T , and K_4 also depends on δ .

(H_3) There exists $\Gamma \subseteq [0, T]$ of full measure such that, for every $t_0 \in \Gamma$, the following local energy inequality and local entropy inequality hold: For $\forall t > t_0, \forall \phi \in C^2(S_L)$,

$$\begin{aligned} & \frac{1}{2} \int h_x^2(t) \phi^4 - \int F(h(t)) \phi^4 + \int_{t_0}^t \int h^n |(h_{xx} + f(h))_x|^2 \phi^4 \\ & \leq - \int_{t_0}^t \int h^n (h_{xx} + f(h))_x (h_{xx} + f(h)) (\phi^4)_x \\ & \quad - \int_{t_0}^t \int h^n (h_{xx} + f(h))_x [h_{xx}(\phi^4)_x + h_x(\phi^4)_{xx}] \\ (2.2) \quad & + \frac{1}{2} \int h_x^2(t_0) \phi^4 - \int F(h(t_0)) \phi^4, \end{aligned}$$

and

$$\begin{aligned} & \int h^{\frac{3}{2}-n+\delta}(t) \phi^4 + \sigma \int_{t_0}^t \int \left(h^{-\frac{1}{2}+\delta} h_{xx}^2 \phi^4 + h^{-\frac{5}{2}+\delta} h_x^4 \phi^4 \right) \\ & \leq C \int_{t_0}^t \int \left[h^{\frac{3}{2}+\delta} (\phi_x^4 + \phi^2 \phi_{xx}^2) + f^2(h) h^{-\frac{1}{2}+\delta} \phi^4 \right] \\ (2.3) \quad & + C \int h(t_0) \phi^4 + \int h^{\frac{3}{2}-n+\delta}(t_0) \phi^4, \quad \delta \in (0, \min\{1.5, n - 1.5\}), \end{aligned}$$

where F is a primitive of f and σ and C are positive constants depending on δ .

From the expression for K_4 we see that it is implicitly assumed $n > 1.5$ in the definition of a suitable weak solution. On the other hand, when $n \geq 3.5$, K_1 and K_3 together imply that a suitable weak solution is always positive; hence by parabolic regularity, it is a classical solution. Consequently we will focus on the range $(1.5, 3.5)$. To state our existence result for these solutions we impose the following assumptions on f : For some constants $C_1 - C_3$,

$$(2.4) \quad |f(z)| \leq C_1(1 + z^p), \quad z \geq 1 \quad \text{for some } 0 \leq p < 3,$$

$$(2.5) \quad |f'(z)| \leq C_2, \quad z \in (0, 1],$$

or,

$$(2.6) \quad |f(z)| + z|f'(z)| \leq C_3 z^q, \quad z \in (0, 1] \quad \text{for some } q > 1 - \frac{n}{2}.$$

Equation (2.4) restricts the growth of f at infinity. When $f(z)$ grows faster than z^3 , it is conjectured that some solutions of (2.1) blow up in finite time. So growth restriction like (2.4) is needed to ensure the existence of a global, weak solution. According to the behavior of f at 0 we impose (2.5) and (2.6), respectively. For power laws the former applies when $q \geq 1$, and the latter applies when $q \in (1 - n/2, 1)$.

THEOREM 2.1. *Consider problem (2.1), where $n \in (1.5, 3.5)$, (2.4) and (2.5) or (2.6) hold. Then for any nonnegative $h_0 \in H^1(S_L)$ with finite $\|h_0^{3/2-n}\|_{L^1}$, (2.1) admits a suitable weak solution h satisfying $\|h(t) - h_0\|_{L^2}$ tends to 0 as $t \downarrow 0$.*

An outline of the proof of this theorem can be found in Appendix A. We remark that this theorem is still valid for $p = 3$ in (2.4) provided that $\|h_0\|_{L^1}$ is sufficiently small.

Weak solutions for (1.1) have been constructed in [BF], [BBDP], and [BP1] and for (1.2) in [BP2]. In particular, in Theorem 3.3 of [BP2] the existence of a weak solution essentially satisfying (H₁) and (H₂) is proved. A slight difference is that our Theorem 2.1 applies to the condition $q > 1 - n/2$, which could be negative when $n > 2$, while q is considered to be positive in the previous works. In (H₃) we also need the local energy and entropy inequalities. The inequalities of this type were derived in [B1], [B2], and [BP2] to study the finite speed of the propagation of the weak solution. Unlike previous works where the initial time is always taken to be zero, here we need these inequalities to be valid starting at initial times which form a subset of the time axis of full measure.

We begin with an estimate on the size of the rupture set in a spatial slice.

THEOREM 2.2. *Let h be a nonnegative function in $L^\infty(0, T; H^1(S_L))$ satisfying, for some $n \in (1.5, 3.5)$,*

$$\sup_{0 < t < T} \int h_x^2 < \infty,$$

and

$$\sup_{0 < t < T} \int h^{3/2-n} < \infty.$$

Then for each $T > 0$, the $((7 - 2n)/(2n + 1))$ -Hausdorff measure of the set $\mathcal{R}_T = \{x : h(x, T) = 0\}$ is finite.

Proof. Let h be a function satisfying the conditions in the theorem. We claim that

$$\int_{B_R} h_x^2(T) \leq R^{\gamma_0}, \quad \gamma_0 = \frac{7 - 2n}{2n + 1}$$

and

$$\int_{B_R} h^{\frac{3}{2}-n}(T) \leq \delta_0 R^{\gamma_0}, \quad \delta_0 = \frac{2}{(1 + \gamma_0)(\frac{3}{2} - n) + 2}$$

hold for all small B_R centered at x_0 , then $h(x_0, T) > 0$. Assuming on the contrary that (x_0, T) is a rupture point, we would have

$$\begin{aligned} h(x, T) &= h(x, T) - h(x_0, T) \\ &\leq \left| \int_{x_0}^x 1 \right|^{1/2} \left(\int_{x_0}^x h_x^2 \right)^{1/2} \\ &\leq |x - x_0|^{\frac{1}{2}(1+\gamma_0)} \end{aligned}$$

for all x sufficiently close to x_0 , but then

$$\begin{aligned} \delta_0 R^{\gamma_0} &\geq \int_{B_R} h^{\frac{3}{2}-n}(T) \\ &\geq \frac{2R^{\frac{1}{2}(1+\gamma_0)(\frac{3}{2}-n)+1}}{\frac{1}{2}(1 + \gamma_0)(\frac{3}{2} - n) + 1}. \end{aligned}$$

As $\gamma_0 = \frac{1}{2}(1 + \gamma_0)(\frac{3}{2} - n) + 1$, contradiction holds.

Henceforth, at a rupture point (x_0, T) there exists $\{R_j\}$ (depending on x_0) $\downarrow 0$ such that

$$\int_{B_{R_j}} (h_x^2 + h^{\frac{3}{2}-n})(T) > \delta_0 R_j^{\gamma_0}.$$

For each $\delta > 0$, the collection \mathcal{C}_δ of all of these balls B_{R_j} , with radii less than δ forms a covering of the rupture set \mathcal{R}_T . By a version of the Vitali covering theorem (see Appendix B) we can select from \mathcal{C}_δ a mutually disjoint, countable subcollection \mathcal{C}'_δ such that

$$\mathcal{R}_T \subseteq \bigcup \{B_{5R_j}(x_j) : B_{R_j}(x_j) \in \mathcal{C}'_\delta\}.$$

Therefore, according to the definition of the Hausdorff measure Federer [Fe]

$$\begin{aligned} \mathcal{H}_\delta^{\gamma_0}(\mathcal{R}_T) &\leq \alpha(\gamma_0) \sum (5R_j)^{\gamma_0} \\ &\leq \frac{5^{\gamma_0} \alpha(\gamma_0)}{\delta_0} \sum \int_{B_{R_j}(x_j)} (h_x^2 + h^{\frac{3}{2}-n})(T) \\ &\leq \frac{5^{\gamma_0} \alpha(\gamma_0)}{\delta_0} \int (h_x^2 + h^{\frac{3}{2}-n})(T), \quad \alpha(s) = \frac{\pi^s/2}{\Gamma(s/2 + 1)}, \end{aligned}$$

which leads to the desired conclusion after taking $\delta \rightarrow 0$. \square

Clearly Theorem 1.1 is a special case of Theorem 2.2. Here we have used little beyond the finiteness of K_1 and K_3 . Although the proof of Theorem 2.2 is quite simple, as one will see, the same idea works in the size estimate of the rupture set in space-time. In order to obtain such estimates, we need to make use of the finiteness of the integrals in space-time such as the quantities K_2 and K_4 . We have the following.

THEOREM 2.3. *Consider any suitable weak solution of (2.1), where $n \in (2, 3.5)$ and either (2.5) or (2.6) hold. The parabolic Hausdorff dimension of the rupture set of this solution cannot exceed $2n/(2n - 3)$.*

The proof of this theorem depends on Lemmas 2.1-2.4 below. For the ease of presentation we shall prove the theorem first.

Proof. Let (x_0, T) be a rupture point of the solution h . We claim that, for every $\gamma_0 \in (\frac{7+2\delta}{4}, 4 + 2\delta)$, there exists a sequence $\{R_j\}$, $R_j \downarrow 0$, such that

$$\iint_{Q_{R_j}} h^n \left| (h_{xx} + f(h))_x \right|^2 > R_j^{\gamma_0},$$

or

$$\iint_{Q_{R_j}} h^{-\frac{5}{2}+\delta} h_x^4 > R_j^{\gamma_0}$$

holds. For, were this not true, it means that (2.7) and (2.14) hold for all sufficiently small R . By Lemma 2.4, there exists some $t_1 \in \Gamma \cap [T - R^4, T]$ such that

$$\int_{B_{R/2}} h_x^2(t_1) \leq CR^\gamma \quad \forall \gamma \in \left(0, \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}\right).$$

Then by Lemma 2.3,

$$\int_{B_{R/4}} h_x^2 \leq CR^\gamma \quad \forall t \in [t_1, T].$$

Lemma 2.2 further gives

$$\sup_{B_{R/4} \times [t_1, T]} h \leq CR^{\frac{1+\gamma}{2}},$$

which means that

$$h(x, T) \leq C|x - x_0|^{\frac{1+\gamma}{2}}$$

for all x sufficiently close to x_0 . However, on the other hand,

$$\begin{aligned} \infty &> \int h^{\frac{3}{2}-n}(T) \\ &\geq C \int |x - x_0|^{\frac{1+\gamma}{2}(\frac{3}{2}-n)} = \infty, \end{aligned}$$

if $(\frac{1+\gamma}{2})(\frac{3}{2} - n) \leq -1$, i.e., $\gamma \geq (7 - 2n)/(2n - 3)$. A contradiction will be drawn if we can verify that

$$\frac{7 - 2n}{2n - 3} < \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}, \quad \text{i.e., } \gamma_0 > \frac{2n + 2\delta}{2n - 3}.$$

A direct computation shows that this is true if $n > 2$.

Thus, at a rupture point (x_0, T) , there exists a sequence of cylinders $\{Q_{R_j}\}$, $R_j \downarrow 0$, satisfying

$$\iint_{Q_{R_j}} \left[h^n \left| (h_{xx} + f(h))_x \right|^2 + h^{-\frac{5}{2}+\delta} h_x^4 \right] > R_j^{\gamma_0}.$$

The collection of all of these cylinders forms a covering of the rupture set in space-time, and we can replace the arguments in the second half of the proof of Theorem 2.2 involving the Hausdorff measure by the parabolic Hausdorff measure (see Appendix B) to show that the γ_0 -parabolic Hausdorff measure of the rupture set is finite for all $\gamma_0 > (2n + 2\delta)/(2n - 3)$. By letting $\delta \rightarrow 0$ we conclude that its parabolic Hausdorff dimension cannot exceed $2n/(2n - 3)$. The proof of Theorem 2.3 is completed. \square

LEMMA 2.1. *Let χ be a strictly increasing, continuous function in $[0, \infty)$, with $\chi(0) > 0$. Define a sequence $\{\alpha_n\}$ by $\alpha_{n+1} = \chi(\alpha_n)$, $\alpha_1 = 0$. Then $\{\alpha_n\}$ is increasing and*

- (a) *it diverges to ∞ if $y = \chi(x)$ has no intersection with $y = x$, or*
- (b) *it converges to some α^* which is the x -coordinate of the first intersection between $y = \chi(x)$ and $y = x$.*

Consequently, for any α (less than α^ if α^* exists), there exists some N such that $\alpha_N \geq \alpha$.*

The proof of this lemma is elementary and is omitted.

LEMMA 2.2. *Let h be a suitable weak solution of (2.1), and (x_0, T) a rupture point of h . Suppose that*

$$(2.7) \quad \iint_{Q_R} h^n \left| (h_{xx} + f(h))_x \right|^2 \leq R^{\gamma_0}, \quad Q_R = (x_0 - R, x_0 + R) \times (T - R^4, T),$$

and

$$(2.8) \quad \int_{B_R} h_x^2 \leq C_4 R^\gamma \quad \forall t \in [t_1, T], \quad B_R = (x_0 - R, x_0 + R),$$

hold for some R, γ_0, γ, t_1 , and C_4 , where $t_1 \in [T - R^4, T]$, $R^4 < T$, and γ satisfies

$$(2.9) \quad \gamma < \frac{2\gamma_0 + n}{2 - n} \text{ for } n \in (1.5, 2), \text{ and } < \infty \text{ for } n \in [2, 3.5).$$

Then for every $\tau_0 \in (0, 1)$, there exists a constant C depending on n, γ_0, γ, C_4 , and τ_0 such that

$$\sup_{B_{\tau_0 R} \times [t_1, T]} h \leq CR^{\frac{1+\gamma}{2}}.$$

Proof. By (H_1) in the definition of a suitable weak solution,

$$\iint h\Phi_t = - \iint h^n (h_{xx} + f(h))_x \Phi_x \quad \forall \Phi \text{ test function.}$$

We choose $\Phi = \xi(x)\theta_\delta(t)$, where ξ is compactly supported in B_R , and $\theta_\delta(t) = \int_0^t \eta(t)dt$, with η given by

$$\eta = \begin{cases} -\frac{1}{2\delta}, & t \in [s - \delta, s + \delta], \\ \frac{1}{2\delta}, & t \in [T - \delta, T + \delta], \\ 0, & \text{otherwise,} \end{cases}$$

and s is a fixed number in $[t_1, T]$. Letting $\delta \downarrow 0$, as $\theta_\delta \leq 1$, we obtain

$$(2.10) \quad \left| \int_{B_R} (h(x, T) - h(x, t))\xi(x) \right| \leq \int_t^T \int_{B_R} h^n |(h_{xx} + f(h))_x \xi_x|,$$

where we have changed notation from s to t . Letting $\tau \in (\tau_0, 1)$ to be chosen later, we pick ξ so that $\xi = 1$ in $B_{\tau R}$ and $|\xi_x| \leq 1/R(1 - \tau)$ in B_R . Using (2.7),

$$\begin{aligned} \int_t^T \int_{B_R} h^n |(h_{xx} + f(h))_x \xi_x| &\leq \sqrt{\int_t^T \int_{B_R} h^n (h_{xx} + f(h))_x^2} \sqrt{\int_t^T \int_{B_R} h^n \xi_x^2} \\ &\leq \frac{1}{(1 - \tau)} R^{\frac{\gamma_0+3}{2}} \sup_{Q_R} h^{\frac{n}{2}}. \end{aligned}$$

It follows from (2.10) that

$$\left| \int_{B_R} (h(x, T) - h(x, t))\xi(x) \right| \leq \frac{1}{2\tau(1 - \tau)} R^{\frac{\gamma_0+1}{2}} \sup_{Q_R} h^{\frac{n}{2}},$$

where we have set $\int_{B_R} F \equiv \frac{\int_{B_R} F}{\int_{B_R} \xi}$. By (2.8) and (2.10),

$$\begin{aligned} h(x, t) &\leq \left| \int_{B_R} (h(x, t) - h(y, t))\xi(y)dy \right| + \left| \int_{B_R} (h(y, t) - h(y, T))\xi(y)dy \right| \\ &\quad + \left| \int_{B_R} (h(y, T) - h(x_0, T))\xi(y)dy \right| \\ (2.11) \quad &\leq 2\sqrt{C_4}R^{\frac{1+\gamma}{2}} + \frac{1}{2\tau(1 - \tau)} R^{\frac{1+\gamma_0}{2}} \sup_{Q_R} h^{n/2} \quad \forall x \in B_{\tau R}. \end{aligned}$$

Hence in the case

$$(2.12) \quad \sup_{Q_R} h \leq CR^\alpha \quad \text{for some } \alpha < \frac{1+\gamma}{2},$$

(2.11) implies that

$$(2.13) \quad \begin{aligned} \sup_{Q'_R} h &\leq C \left(R^{\frac{1+\gamma}{2}} + R^{\frac{1+\gamma_0+n\alpha}{2}} \right) \\ &\leq CR^{\chi(\alpha)}, \quad Q'_R = B_{\tau R} \times [t_1, T], \end{aligned}$$

provided that $\chi(\alpha) \equiv (1 + \gamma_0 + n\alpha)/2 < (1 + \gamma)/2$. Starting at $\alpha = 0$, we bootstrap between (2.12) and (2.13) to obtain

$$\sup_{B_{\tau k R} \times [t_1, T]} h \leq C_k R^{\chi(\alpha_k)}$$

after k times of iterations as long as $\chi(\alpha_k) < (1 + \gamma)/2$. When $n \geq 2$, $y = \chi(x)$ does not intersect $y = x$, so there exists an N such that $\chi(\alpha_N) \geq (1 + \gamma)/2$. Then the lemma follows by applying the Holder continuity of h in space. When $n \in (1.5, 2)$, the same conclusion holds if $(\gamma + 1)/2 < \alpha_*$, the intersection of $y = \chi(x)$, and $y = x$. As α_* is given by $\frac{(1+\gamma_0)/2}{1-n/2} = \frac{1+\gamma_0}{2-n}$, $(\gamma + 1)/2 < \alpha_*$ is equivalent to (2.9). \square

LEMMA 2.3. *Let h be a suitable weak solution of (2.1), where $n \in (1.5, 3.5)$ and (2.5) holds, and let (x_0, T) be a rupture point of h . Suppose that (2.7),*

$$(2.14) \quad \iint_{Q_R} h^{-\frac{5}{2}+\delta} h_x^4 \leq R^{\gamma_0}$$

and

$$(2.15) \quad \int_{B_R} h_x^2(t_1) \leq C_5 R^\gamma$$

hold for some $R, \gamma_0, \gamma, t_1 \in \Gamma \cap [T - R^4, T]$, C_5 , where γ and γ_0 also satisfy (2.9),

$$(2.16) \quad \gamma_0 > \frac{7 - 2n + 2\delta}{4}, \quad \delta \in (0, n - 1.5),$$

and

$$(2.17) \quad 0 \leq \gamma < 3.$$

Then for every $\tau_0 \in (0, 1)$, there exists a constant C depending on $n, \gamma_0, \gamma, \delta, C_5$, and τ_0 such that

$$\int_{B_{\tau_0 R}} h_x^2 \leq CR^\gamma \quad \forall t \in [t_1, T].$$

When (2.6) instead of (2.5) holds in the above assumptions, the same conclusion holds when (2.17) is replaced by

$$(2.17)' \quad \gamma < \frac{8 - n}{n}.$$

Proof. Let ϕ be a cut-off function supported in B_R , $\phi = 1$ in $B_{\tau R}$, where $\tau \in (\tau_0, 1)$ is to be chosen later. We plug ϕ into the local energy inequality (2.2) and use the Cauchy–Schwarz inequality to get

$$\begin{aligned} & \left| \iint h^n (h_{xx} + f(h))_x (h_{xx} + f(h)) (4\phi^3 \phi_x) \right| \\ & \leq \frac{1}{4} \iint h^n |(h_{xx} + f(h))_x|^2 \phi^4 + C \iint h^n (h_{xx} + f(h))^2 \phi^2 \phi_x^2 \end{aligned}$$

and

$$\begin{aligned} & \left| \iint h^n (h_{xx} + f(h))_x [4h_{xx} \phi^3 \phi_x + h_x (4\phi^3 \phi_{xx} + 12\phi^2 \phi_x^2)] \right| \\ & \leq \frac{1}{4} \iint h^n |(h_{xx} + f(h))_x|^2 \phi^4 + C \iint [h^n (h_{xx}^2 \phi^2 \phi_x^2 + h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2))]. \end{aligned}$$

It follows that

$$\begin{aligned} \frac{1}{2} \int_{B_R} h_x^2 \phi^4 & \leq \int_{B_R} F(h) \phi^4 + \frac{1}{2} \int_{B_R} h_x^2(t_1) \phi^4 - \int_{B_R} F(h(t_1)) \phi^4 \\ & \quad - \frac{1}{2} \int_{t_1}^t \int_{B_R} h^n |(h_{xx} + f(h))_x|^2 \phi^4 + C \left[\int_{t_1}^t \int_{B_R} h^n (h_{xx} + f(h))^2 \phi^2 \phi_x^2 \right. \\ & \quad \left. + \int_{t_1}^t \int_{B_R} h^n (h_{xx}^2 \phi^2 \phi_x^2 + h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2)) \right] \\ & \leq \int_{B_R} F(h) \phi^4 + \frac{1}{2} \int_{B_R} h_x^2(t_1) \phi^4 \\ & \quad - \int_{B_R} F(h(t_1)) \phi^4 - \frac{1}{4} \int_{t_1}^t \int_{B_R} h^n |(h_{xx} + f(h))_x|^2 \phi^4 \\ & \quad + C \left[\int_{t_1}^t \int_{B_R} h^n f'^2(h) h_x^2 \phi^4 + \int_{t_1}^t \int_{B_R} h^n f^2(h) \phi^2 \phi_x^2 \right. \\ & \quad \left. + \int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^n h_{xx}^2 \phi^2 \phi_x^2 \right] \\ & \leq \int_{B_R} F(h) \phi^4 + \frac{1}{2} \int_{B_R} h_x^2(t_1) \phi^4 \\ & \quad - \int_{B_R} F(h(t_1)) \phi^4 - \frac{1}{4} \int_{t_1}^t \int_{B_R} h^n |(h_{xx} + f(h))_x|^2 \phi^4 \\ & \quad + C \left[\int_{t_1}^t \int_{B_R} h^n f'^2(h) h_x^2 \phi^4 + \int_{t_1}^t \int_{B_R} h^n f^2(h) \phi^2 \phi_x^2 \right. \\ (2.18) \quad & \left. + \int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^{n-2} h_x^4 \phi^2 \phi_x^2 \right]. \end{aligned}$$

Notice that in the last inequality we have used

$$\begin{aligned} \int_{t_1}^t \int_{B_R} h^n h_{xx}^2 \phi^2 \phi_x^2 &= - \int_{t_1}^t \int_{B_R} h^n h_x h_{xxx} \phi^2 \phi_x^2 - n \int_{t_1}^t \int_{B_R} h^{n-1} h_x^2 h_{xx} \phi^2 \phi_x^2 \\ &\quad - \int_{t_1}^t \int_{B_R} h^n h_x h_{xx} (2\phi \phi_x^3 + 2\phi^2 \phi_x \phi_{xx}) \\ &\leq \varepsilon \int_{t_1}^t \int_{B_R} h^n h_{xxx}^2 \phi^4 + \varepsilon \int_{t_1}^t \int_{B_R} h^n h_{xx}^2 \phi^2 \phi_x^2 \\ &\quad + C_\varepsilon \left[\int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^{n-2} h_x^4 \phi^2 \phi_x^2 \right]. \end{aligned}$$

Therefore, under (2.5), when $f \in C^1[0, \infty)$, we have

$$\begin{aligned} \frac{1}{2} \int_{B_R} h_x^2 \phi^4 &\leq \int_{B_R} F(h) \phi^4 + \frac{1}{2} \int_{B_R} h_x^2(t_1) \phi^4 - \int_{B_R} F(h(t_1)) \phi^4 \\ &\quad + C \left[\int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi^4 + \phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^{n-2} h_x^4 \phi^2 \phi_x^2 \right. \\ (2.19) \quad &\quad \left. + \int_{t_1}^t \int_{B_R} h^n \phi^2 \phi_x^2 \right], \end{aligned}$$

since a suitable weak solution is bounded. Here the constant C depends on f . Suppose that

$$(2.20) \quad \int_{B_R} h_x^2 \leq CR^\alpha \quad \forall t \in [t_1, T] \quad \text{for some } \alpha < \gamma.$$

We claim that

$$(2.21) \quad \int_{B_{\tau R}} h_x^2 \leq C'R^{\chi(\alpha)} \quad \forall t \in [t_1, T],$$

where

$$\chi(\alpha) = \min \left\{ \frac{3}{2} + \frac{\alpha}{2}, \frac{n}{2} + \left(1 + \frac{n}{2}\right) \alpha, \gamma_0 + \frac{n - \frac{7}{2} - \delta}{2} + \frac{n + \frac{1}{2} - \delta}{2} \alpha \right\}$$

as long as $\alpha < \gamma$. For, by Lemma 2.2 and (2.20), we have

$$\sup_{B_{\tau R} \times [t_1, T]} h \leq CR^{\frac{1+\alpha}{2}}.$$

Using this estimate we control the integrals on the right of (2.19) as follows.

$$\begin{aligned} \int_{B_R} F(h(t)) \phi^4 &\leq CR^{\frac{3}{2} + \frac{\alpha}{2}}, \\ \int_{B_R} h_x^2(t_1) \phi^4 &\leq C_5 R^\gamma, \quad (\text{from (2.15)}) \\ \int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi^4 + \phi_x^4 + \phi^2 \phi_{xx}^2) &\leq CR^{\alpha + \frac{n}{2}(1+\alpha)}, \\ \int_{t_1}^t \int_{B_R} h^{n-2} h_x^4 \phi^2 \phi_x^2 &\leq CR^{-2} \sup_{Q_R} h^{n+\frac{1}{2}-\delta} \int_{t_1}^t \int_{B_R} h^{-\frac{5}{2}+\delta} h_x^4 \\ &\leq CR^{\gamma_0 + \frac{n-\frac{7}{2}-\delta}{2} + \frac{n+\frac{1}{2}-\delta}{2} \alpha}, \\ \int_{t_1}^t \int_{B_R} h^n \phi^2 \phi_x^2 &\leq CR^{3 + \frac{n}{2}(1+\alpha)} \leq CR^{\alpha + \frac{n}{2}(1+\alpha)}, \end{aligned}$$

so (2.21) follows after putting these estimates back to (2.19). As (2.20) holds for $\alpha = 0$, we can bootstrap between (2.20) and (2.21) to increase the power in R in (1.20). Observe that the function χ is piecewise linear with slopes equal to $1/2, 1+n/2$, or $(2n+1-2\delta)/4$. As $(2n+1-2\delta)/4 > 1$ for $\delta < n-3/2$, χ can possibly intersect $y = x$ at $\alpha^* = 3$. The lemma now follows from Lemma 2.1 after taking $\tau^N = \tau_0$.

In case f satisfies (2.6), the local energy inequality reads as

$$\begin{aligned} \frac{1}{2} \int_{B_R} h_x^2 \phi^4 &\leq \frac{B}{q(q+1)} \int_{B_R} h^{q+1} \phi^4 + \frac{1}{2} \int_{B_R} h_x^2(t_1) \phi^4 - \frac{B}{q(q+1)} \int_{B_R} h^{q+1}(t_1) \phi^4 \\ &+ C \left(\int_{t_1}^t \int_{B_R} h^n h_x^2 (\phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^{n-2} h_x^4 \phi^2 \phi_x^2 \right. \\ &\left. + \int_{t_1}^t \int_{B_R} h^{n+2q} \phi^2 \phi_x^2 + \int_{t_1}^t \int_{B_R} h^{n+2q-2} h_x^2 \phi^4 \right). \end{aligned}$$

Under (2.20), the integral terms involving q can be estimated as follows:

$$\begin{aligned} \int_{B_R} h^{q+1} \phi^4 &\leq CR \sup_{B_R} h^{q+1} \\ &\leq CR^{\frac{q+3}{2} + \frac{4-n}{4} \alpha} \\ &\leq CR^{\frac{8-n}{4} + \frac{4-n}{4} \alpha}, \quad (\because q > 1 - n/2) \\ \int_{t_1}^t \int_{B_R} h^{n+2q} \phi^2 \phi_x^2 &\leq CR^3 \sup_{Q_R} h^{n+2q} \\ &\leq CR^{4+\alpha} \end{aligned}$$

and

$$\begin{aligned} \int_{t_1}^t \int_{B_R} h^{n+2q-2} h_x^2 \phi^4 &\leq CR^4 \left(\sup_{[t_1, t]} \int_{B_R} h_x^2 \right) \sup_{Q_R} h^{n+2q-2} \\ &\leq CR^{4+\alpha}, \quad (\because q > 1 - n/2). \end{aligned}$$

So (2.21) holds after χ is replaced by the function χ' :

$$\begin{aligned} \chi'(\alpha) = \min \left\{ \frac{8-n}{4} + \frac{4-n}{4} \alpha, 4 + \alpha, \frac{n}{2} + \left(1 + \frac{n}{2}\right) \alpha, \right. \\ \left. \gamma_0 + \frac{n - \frac{7}{2} - \delta}{2} + \frac{n + \frac{1}{2} - \delta}{2} \alpha \right\}. \end{aligned}$$

Arguing as before, we can show that Lemma 2.3 holds when $q > 1 - n/2$ and (2.17) is replaced by (2.17)'. \square

LEMMA 2.4. *Let h be a suitable weak solution of (2.1), where $n \in (1.5, 3.5)$ and (2.5) holds, and let (x_0, T) be a rupture point of h . Suppose (2.7), (2.14), (2.16), and (2.17) hold. Then there exists $t_1 \in \Gamma \cap [T - R^4, T]$ such that*

$$\int_{B_{R/2}} h_x^2(t_1) \leq CR^\gamma \quad \forall \gamma \in \left(0, \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}\right),$$

if

$$(2.22) \quad \frac{7 + 2\delta}{4} < \gamma_0 \leq 4 + 2\delta.$$

When (2.5) is replaced by (2.6) and (2.17) by (2.17)', respectively, the same conclusion holds when (2.22) is replaced by

$$(2.22)' \quad \frac{7 + 2\delta}{4} < \gamma_0 < 1 + \frac{6 + 4\delta}{n}.$$

Proof. We only prove this lemma when f satisfies (2.5), for the other case can be proved in a similar way. From (2.14) there exists $t_1 \in \Gamma$ such that

$$\int_{B_R} h^{-\frac{5}{2} + \delta} h_x^4(t_1) \leq R^{\gamma_0 - 4}.$$

Thus

$$(2.23) \quad \begin{aligned} \int_{B_R} h_x^2(t_1) &\leq \sqrt{\int_{B_R} h^{-\frac{5}{2} + \delta}(t_1) h_x^4(t_1)} \sqrt{\int_{B_R} h^{\frac{5}{2} - \delta}(t_1)} \\ &\leq R^{\frac{\gamma_0 - 3}{2}} \sup_{B_R} h^{\frac{5 - 2\delta}{4}}(t_1). \end{aligned}$$

Suppose that

$$(2.24) \quad \int_{B_R} h_x^2(t_1) \leq CR^\alpha$$

for some α satisfying $\alpha < \gamma$. We claim that, for any fixed $\tau_0 \in (0, 1)$,

$$(2.25) \quad \int_{B_{\tau_0 R}} h_x^2(t_1) \leq C' R^{\chi(\alpha)},$$

where $\chi(\alpha) = \frac{\gamma_0 - 3}{2} + \frac{5 - 2\delta}{8}(1 + \alpha)$ as long as $\chi(\alpha) < \gamma$. For, using (2.24) and the fact that $\alpha < \gamma$, we can apply Lemma 2.3 to infer that

$$\int_{B_{\frac{1+\tau_0}{2} R}} h_x^2 \leq C' R^\alpha \quad \forall t \in [t_1, T].$$

By Lemma 2.2,

$$\sup_{B_{\tau_0 R} \times [t_1, T]} h \leq CR^{\frac{1+\alpha}{2}}.$$

Plugging this estimate into (2.23) yields (2.25). Now, using the fact that (2.24) holds for $\alpha = 0$, bootstrapping between (2.24) and (2.25) increases the power in R . Noting that the intersection of χ and $y = x$ is at $\alpha^* = (4\gamma_0 - 7 - 2\delta)/(3 + 2\delta)$, we conclude that the lemma holds if $\gamma < (4\gamma_0 - 7 - 2\delta)/(3 + 2\delta)$. Note that in (2.17) it is required that $\gamma < 3$, and this holds if $(4\gamma_0 - 7 - 2\delta)/(3 + 2\delta) < 3$, that is, $\gamma_0 \leq 4 + 2\delta$. \square

An estimate on the Hausdorff dimension of the rupture set is also available in the range (1.5, 2]. We need to work a bit harder for it though.

THEOREM 2.4. *Let h be a suitable weak solution of (2.1), where $n \in (1.5, 2]$ and f satisfies either (2.5) or (2.6). The Hausdorff dimension of the rupture set of h cannot exceed $1 + 6/n$.*

Proof. Let us consider the case where f satisfies (2.6) only. We claim that at a rupture point (x_0, T) there exists $\{R_j\} \downarrow 0$ such that

$$\iint_{Q_{R_j}} h^n |(h_{xx} + f(h))_x|^2 > R_j^{\gamma_0},$$

$$\iint_{Q_{R_j}} h^{-\frac{5}{2}+\delta} h_x^4 > R_j^{\gamma_0},$$

or

$$\iint_{Q_{R_j}} h^{\frac{3}{2}-n+\delta} > R_j^{\gamma_0+\lambda}$$

hold when γ_0 and δ satisfy (2.22)' and λ is a fixed positive number. For, were this not true, that means (2.7), (2.14), and

$$(2.26) \quad \iint_{Q_R} h^{\frac{3}{2}-n+\delta} \leq R^{\gamma_0+\lambda}$$

hold for all sufficiently small R . By Lemma 2.4,

$$\int_{B_{R/2}} h_x^2(t_1) \leq CR^\gamma \quad \text{for } \gamma < \min \left\{ \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}, \frac{8 - n}{n} \right\} = \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}.$$

By Lemma 2.3 we further infer that

$$\int_{B_{R/4}} h_x^2(t) \leq C'R^\gamma \quad \forall t \in [t_1, T].$$

Consequently,

$$h(x, T) \leq C|x - x_0|^{\frac{1+\gamma}{2}} \quad \forall x \text{ sufficiently close to } x_0.$$

On the other hand, (2.26) implies that there exists some $t_2 \in \Gamma$ such that

$$\int_{B_R} h^{\frac{3}{2}-n+\delta}(t_2) \leq CR^{\gamma_0-4+\lambda}.$$

By Lemma 2.5,

$$\int_{B_{R/2}} h^{\frac{3}{2}-n+\delta}(t) \leq C'R^{\gamma_0-4+\lambda} \quad \forall t \in [t_2, T].$$

So

$$\begin{aligned} 4R^{\gamma_0-4+\lambda} &\geq \int_{B_{R/2}} h^{\frac{3}{2}-n+\delta} \\ &\geq C \int_{B_{R/2}} |x - x_0|^{\frac{1+\gamma}{2}(\frac{3}{2}-n+\delta)} \\ &\geq CR^{\frac{1+\gamma}{2}(\frac{3}{2}-n+\delta)+1} \end{aligned}$$

if

$$\left(\frac{1+\gamma}{2}\right)\left(\frac{3}{2}-n+\delta\right)+1 > 0.$$

To derive a contradiction we choose

$$\gamma_0 = 1 + \frac{6}{n}$$

and

$$\gamma = \frac{4\gamma_0 - 7}{3} - \sigma,$$

where σ is chosen to be small as δ is small and such that

$$\gamma < \frac{4\gamma_0 - 7 - 2\delta}{3 + 2\delta}$$

and

$$\frac{1+\gamma}{2}\left(\frac{3}{2}-n+\delta\right)+1 > 0.$$

As $\delta, \sigma \rightarrow 0$, and γ tends to $(4\gamma_0 - 7)/3$ and

$$\gamma_0 - 4 + \lambda > \gamma_0 - 4 = \frac{1+\gamma}{2}\left(\frac{3}{2}-n\right)+1,$$

the contradiction holds.

Now a covering argument as in the proof of Theorem 2.3 finishes the job. \square

LEMMA 2.5. *Let h be a suitable weak solution (2.1), where $n \in (1.5, 2]$ and either (2.5) or (2.6) hold. There exists a small $R_0 > 0$ such that if*

$$\int_{B_R} h^{\frac{3}{2}-n+\delta}(t_1) \leq C_6 R^\gamma, \quad 0 < \delta < n - \frac{3}{2} \text{ small},$$

for some $\gamma \in [0, 1]$, $R < R_0$ and $t_1 \in \Gamma \cap [T - R^4, T]$, then for every $\tau_0 \in (0, 1)$,

$$\int_{B_{\tau_0 R}} h^{\frac{3}{2}-n+\delta}(t) \leq C R^\gamma \quad \forall t \in [t_1, T],$$

where C depends on n, q, δ, C_6 , and τ_0 .

Proof. We focus on the case where (2.6) holds, since the other case can be handled similarly. By the local entropy inequality (2.3) $\forall t \geq t_1$,

$$\begin{aligned} \int_{B_R} h^{\frac{3}{2}-n+\delta} \phi^4 &\leq -\sigma \left[\int_{t_1}^t \int_{B_R} h^{-\frac{1}{2}+\delta} h_{xx}^2 \phi^4 + \int_{t_1}^t \int_{B_R} h^{-\frac{5}{2}+\delta} h_x^4 \phi^4 \right] \\ &\quad + C \left[\int_{t_1}^t \int_{B_R} h^{\frac{3}{2}+\delta} (\phi_x^4 + \phi^2 \phi_{xx}^2) + \int_{t_1}^t \int_{B_R} h^{2q-\frac{1}{2}+\delta} \phi^4 \right] \\ &\quad + \int_{B_R} h^{\frac{3}{2}-n+\delta}(t_1), \end{aligned}$$

where ϕ is supported in B_R . Using $q \geq 1 - n/2$, we have

$$\int_{t_1}^t \int_{B_R} h^{\frac{3}{2}+\delta} (\phi_x^4 + \phi^2 \phi_{xx}^2) \leq CR,$$

$$\int_{t_1}^t \int_{B_R} h^{2q-\frac{1}{2}+\delta} \phi^4 \leq C(t-t_1) \sup_{[t_1,t]} \int_{B_R} h^{\frac{3}{2}-n+\delta} \phi^4.$$

Therefore,

$$\sup_{[t_1,t]} \int_{B_R} h^{\frac{3}{2}-n+\delta} \phi^4 \leq C \left[(t-t_1) \sup_{[t_1,t]} \int_{B_R} h^{\frac{3}{2}-n+\delta} + R^\gamma \right].$$

By taking t_2 so that $C(t_2 - t_1) = 1/2$, and $\phi \equiv 1$ on $B_{\tau R}$, we get

$$\sup_{[t_1,t_2]} \int_{B_{\tau R}} h^{\frac{3}{2}-n+\delta} \leq CR^\gamma.$$

Now we may use t_2 as the initial time to extend the estimate beyond t_2 . After finitely many times extensions the lemma follows. \square

3. An estimate on rupture times. In this section we establish an estimate on the Hausdorff dimension of the rupture times. An instant t is a rupture time for a suitable weak solution u of (2.1) if $u(x, t) = 0$ for some x . As our derivation does not involve the use of the local energy or entropy inequalities, the condition (H_3) in the definition of a suitable weak solution is not needed. In other words, the main result, Theorem 3.1, holds for any weak solution of (2.1) satisfying (H_1) and (H_2) only.

We start with the following interpolating inequalities controlling $\sup u^{-1}$ in terms of integrals involving the derivatives of the function u and the integrals of its negative powers.

LEMMA 3.1. *Let u be a nonnegative H^1 -function of period L . Then for any $p > 2$, there exists a constant C depending on p and L such that*

$$\sup u^{-1} \leq C \left[\left(\int u^{-p} \right)^{\frac{1}{p}} + \left(\int u^{-p} \right)^{\frac{1}{p-2}} \left(\int u_x^2 \right)^{\frac{1}{p-2}} \right].$$

Proof. We have, for $\alpha > 0$,

$$u^{-\alpha}(x) \leq u^{-\alpha}(y) + \alpha \left| \int u^{-\alpha-1} u_x \right|$$

$$\leq u^{-\alpha}(y) + \alpha \left(\int u^{-2\alpha-2} \right)^{\frac{1}{2}} \left(\int u_x^2 \right)^{\frac{1}{2}}.$$

Integrating in y yields

$$u^{-\alpha}(x) \leq \frac{1}{L} \int u^{-\alpha} + \alpha \left(\int u^{-2\alpha-2} \right)^{\frac{1}{2}} \left(\int u_x^2 \right)^{\frac{1}{2}}.$$

Letting $2\alpha = p - 2$, we get

$$u^{-1}(x) \leq C \left[\left(\int u^{-\frac{p+2}{2}} \right)^{\frac{2}{p-2}} + \left(\int u^{-p} \right)^{\frac{1}{p-2}} \left(\int u_x^2 \right)^{\frac{1}{p-2}} \right]$$

$$\leq C \left[\left(\int u^{-p} \right)^{\frac{1}{p}} + \left(\int u^{-p} \right)^{\frac{1}{p-2}} \left(\int u_x^2 \right)^{\frac{1}{p-2}} \right]. \quad \square$$

LEMMA 3.2. *Let u be a nonnegative H^2 -function of period L . Then for $p > 2/3$, there exists a constant C depending on p and L such that*

$$\sup u^{-1} \leq C \left[\left(\int u^{-p} \right)^{\frac{3}{3p-2}} \left(\int u_{xx}^2 + \beta \right)^{\frac{1}{3p-2}} + \left(\int u_{xx}^2 + \beta \right)^{-\frac{1}{2}} \right] \quad \forall \beta \geq 0.$$

Proof. It suffices to establish this inequality by taking u to be positive and C^2 . The general case follows from approximation. We assume $u(0) = \inf u$ and set

$$E = \int u^{-p} \text{ and } F = \int u_{xx}^2 + \beta.$$

Then by using Taylor's formula

$$u(x) = u(0) + u_x(0)x + \int_0^x u_{xx}(t)(x-t)dt,$$

we have, for all x ,

$$u(x) \leq u(0) + \frac{2}{3} F^{\frac{1}{2}} x^{\frac{3}{2}} \quad \forall x \in [0, L].$$

So

$$\int_0^L \frac{dx}{\left(u(0) + \frac{2}{3} F^{\frac{1}{2}} x^{\frac{3}{2}} \right)^p} \leq \int u^{-p} \leq E,$$

which implies that

$$\int_0^L \frac{dx}{\left(x + 2u(0)^{\frac{2}{3}} F^{-\frac{1}{3}} \right)^{\frac{3p}{2}}} \leq CEF^{\frac{p}{2}}.$$

After an integration over $(0, L/2)$, we have

$$\left(2u(0)^{\frac{2}{3}} F^{-\frac{1}{3}} \right)^{1-\frac{3p}{2}} \leq C(EF^{\frac{p}{2}} + 1),$$

i.e.,

$$\inf u^{-1} \leq C \left(E^{\frac{3}{3p-2}} F^{\frac{1}{3p-2}} + F^{-\frac{1}{2}} \right). \quad \square$$

To apply this lemma to our weak solution, we take advantage of the finite quantities K_4 and K_3 . A direct calculation shows that $u = h^{\frac{3+2\delta}{4}}$ satisfies

$$\int u_{xx}^2 \leq C_\delta \left(\int h^{-\frac{1}{2}+\delta} h_{xx}^2 + \int h^{-\frac{5}{2}+\delta} h_x^4 \right).$$

Now, using Lemma 3.2 ($\beta = 1$) we have, for any $p' > 2/3$,

$$(3.1) \quad \sup h^{-\frac{3+2\delta}{4}} \leq C \left[\left(\int h^{-p' \frac{3+2\delta}{4}} \right)^{3/(3p'-2)} \left(\int h^{-\frac{1}{2}+\delta} h_{xx}^2 + \int h^{-\frac{5}{2}+\delta} h_x^4 \right)^{1/(3p'-2)} + 1 \right].$$

COROLLARY 3.1. *Let h be a suitable weak solution of (2.1) for $n \in (2, 3.5)$. Then*

$$K_5 \equiv \int_0^T \left(\int h^{-p} \right)^{\frac{3(n-2)-\delta}{p+\frac{3}{2}-n}} < \infty \quad \forall p > n - 3/2, \delta \in (0, 3(n-2)).$$

Proof. When $\delta \in (0, 3(n-2))$, we can take $p' = (4n-6)/(3+2\delta) > 2/3$ so that $p'(3+2\delta)/4 = n - 3/2$ in (3.1) to get, for $p > n - 3/2$,

$$\begin{aligned} \int h^{-p} &\leq \sup(h^{-1})^{p+\frac{3}{2}-n} \int h^{\frac{3}{2}-n} \\ &\leq CK_3 \left[K_3^{\frac{1}{n-2-\frac{\delta}{3}}} \left(\int h^{-\frac{1}{2}+\delta} h_{xx}^2 + \int h^{-\frac{5}{2}+\delta} h_x^4 + 1 \right)^{\frac{1}{3(n-2)-\delta}} + 1 \right]^{p+\frac{3}{2}-n}. \end{aligned}$$

Raise both sides to the power $[3(n-2)-\delta]/(p+\frac{3}{2}-n)$ and then integrate over time. The corollary follows from the boundedness of K_3 and K_4 . \square

Now, we have the following result which contains Theorem 1.3.

THEOREM 3.1. *Let h be a suitable weak solution of (2.1), where $n \in (2, 3.5)$ and either (2.5) or (2.6) holds. The Hausdorff dimension of the set of rupture times cannot exceed $1 - 2(n-2)^2/(8-n)$.*

Once again we point out that the theorem holds for any weak solution of (2.1) satisfying (H₁) and (H₂) only.

Proof. Let n be given as above and p, δ satisfy

$$p \in (2, n], \quad p + \frac{3}{2} > n, \quad \delta \in (0, 3(n-2)).$$

We claim that there exists a small ε_0 depending on n, p, δ such that whenever

$$\int_{T-R}^T \left(\int h^{-p} \right)^\beta \leq \varepsilon_0 R^\gamma \quad \forall \text{ small } R,$$

where

$$\beta = \frac{3(n-2)-\delta}{p+\frac{3}{2}-n}, \quad \gamma = 1 - \beta \left(\frac{p-2}{8-n} \right),$$

then $h(x_0, T) > 0$. To prove this claim we first observe that by the mean-value theorem there is some $t_1 \in [T-R, T]$ such that

$$\int h^{-p}(\cdot, t_1) \leq 2\varepsilon_0^{\frac{1}{\beta}} R^{\frac{\gamma-1}{\beta}}.$$

By Lemma 3.1 and $K_1 < \infty$, h is positive at t_1 . By continuity h is positive for all t close to t_1 , in particular, it is a classical solution near t_1 . Let

$$t^* = \sup \{ t_2 \in (t_1, T] : h \text{ is positive in } [t_1, t_2] \}.$$

We shall show that $t^* = T$. As for $t \in (t_1, t^*)$,

$$\begin{aligned} \frac{d}{dt} \int h^{-p} &= -p(p+1) \int h^{n-p-2} h_{xx}^2 - p(p+1)(n-p-2) \int h^{n-p-3} h_x^2 h_{xx} \\ (3.2) \quad &+ p(p+1) \int h^{n-p-2} f'(h) h_x^2. \end{aligned}$$

We estimate the integrals on the right as follows: For $p \in (2, n)$,

$$\begin{aligned} \left| \int h^{n-p-3} h_x^2 h_{xx} \right| &= \left| \frac{n-p-3}{3} \int h^{n-p-4} h_x^4 \right| \\ &\leq C \int h^{p-n} \left| \left(h^{\frac{n-p}{2}} \right)_x \right|^4 \\ &\leq C \sup (h^{-1})^{n-p} \int \left| \left(h^{\frac{n-p}{2}} \right)_x \right|^4 \\ &\leq C \sup (h^{-1})^{n-p} \left(\int \left(h^{\frac{n-p}{2}} \right)_x^2 \right)^{3/2} \left(\int \left(h^{\frac{n-p}{2}} \right)_{xx}^2 \right)^{1/2} \\ &\leq \varepsilon \int \left(h^{\frac{n-p}{2}} \right)_{xx}^2 + C \sup (h^{-1})^{2(n-p)} \left(\int \left(h^{\frac{n-p}{2}} \right)_x^2 \right)^3 \\ &\leq C\varepsilon \left(\int h^{n-p-2} h_{xx}^2 + \int h^{n-p-4} h_x^4 \right) + C \sup (h^{-1})^{6+p-n} \left(\int h_x^2 \right)^3, \end{aligned}$$

which means, for small ε ,

$$\int h^{n-p-4} h_x^4 \leq C \left[\varepsilon \int h^{n-p-2} h_{xx}^2 + \sup (h^{-1})^{6+p-n} \right].$$

By Lemma 3.1 we further have

$$(3.3) \quad \int h^{n-p-4} h_x^4 \leq C \left[\varepsilon \int h^{n-p-2} h_{xx}^2 + C \left(1 + \int h^{-p} \right)^{\frac{6+p-n}{p-2}} \right].$$

By modifying the above arguments slightly we see that (3.3) also holds when $p = n$. For the third term on the right-hand side of (3.2) we have

$$\begin{aligned} \int h^{n-p-2} f'(h) h_x^2 &\leq C \left(\int h^{n-p-2} h_x^2 + \int h^{n-p+q-3} h_x^2 \right) \\ &\leq C \left(\int h^{n-p-4} h_x^4 + \int h^{n-p} + \int h^{n-p+2q-2} \right) \\ (3.4) \quad &\leq C \left[\int h^{n-p-4} h_x^4 + \left(\int h^{-p} \right)^{\frac{n-p+2q-2}{-p}} \right]. \end{aligned}$$

Putting (3.3) and (3.4) back to (3.2) yields, after taking ε small,

$$\frac{d}{dt} \int h^{-p} \leq C \left(1 + \int h^{-p} \right)^\mu,$$

where $\mu = \max \left\{ \frac{6+p-n}{p-2}, \frac{p-n-2q+2}{p} \right\} = \frac{6+p-n}{p-2}$ (because $n \leq 8$). By comparing $\int h^{-p}$ with the solution Y of

$$\frac{dY}{dt} = C(1+Y)^\mu, \quad Y(t_1) = 2\varepsilon_0^{\frac{1}{\beta}} R^{\frac{\gamma-1}{\beta}},$$

we see that Y exists in $[t_1, T]$, when $\varepsilon_0^{1/\beta}$ is sufficiently small and $\int h^{-p} \leq Y(t)$ as long as both functions exist. So $\int h^{-p}$ is finite in $[t_1, T]$. It follows from Lemma 3.1 that h is positive up to T , that is, $t^* = T$.

We fix an ε_0 satisfying the property in the above claim. It follows that T is a rupture time if and only if there exists $R_j \downarrow 0$ such that

$$\int_{T-R_j}^T \left(\int h^{-p} \right)^\beta > \varepsilon_0 R_j^\gamma$$

for each p, n, δ in the range. Applying a covering argument to the time interval $(0, \infty)$ one can show that, using the finiteness of K_5 , the γ -Hausdorff measure of the set of rupture times is finite for all $\gamma = 1 - \beta(p-2)/(8-n)$. Letting $\delta \rightarrow 0$ and then taking $p = n$, we conclude the proof of this theorem. \square

4. A result on finite time rupture. For any positive, smooth, periodic function h_0 , it is known that

$$(4.1) \quad \begin{cases} h_t + [h^n(h_{xx} + f(h))]_x = 0, \\ h(0) = h_0 \end{cases}$$

admits a unique, classical solution in $[0, \omega)$, which is of the same period in x and is maximal in the sense that whenever ω is finite,

$$\inf_{t \uparrow \omega} h(t) = 0 \quad \text{or} \quad \sup_{t \uparrow \omega} h(t) = \infty.$$

A proof of this fact, which can be found in [BF], is based on parabolic regularity theory (Eidelman [Ei] or Friedman [F1]) when there is no f -term in (4.1). Nevertheless, an examination of their arguments shows that it still works with the presence of the f -term as long as f is in $C^{2,\alpha}((0, \infty))$ for some $\alpha \in (0, 1)$. In the following we denote by $X(L, c_0)$ the subset of $H^1(S_L)$ consisting of all of the positive functions with area c_0 . As the area is conserved by the equation in (4.1), the classical solution of (4.1) forms a flow in $X(L, c_0)$. For our result the following condition on f will be imposed:

$$(4.2) \quad f(0) = 0, \quad |f'(z)| \leq C_1 z^{q-1}, \quad z \in (0, 1] \quad \text{for some } q \geq \frac{1-n}{2}.$$

THEOREM 4.1. *Let $n \in (0, 1/2)$ be in (4.1), and $f \in C([0, \infty)) \cap C^{2,\alpha}((0, \infty))$ satisfy (4.2). Suppose that there exists a bounded open set \mathcal{U} in $X(L, c_0)$ satisfying the following properties:*

(P₁) *The closure of \mathcal{U} does not contain any positive steady state, fully supported droplet with zero contact angle, or an array of identical droplets with zero contact angle of (4.1); and*

(P₂) *Any maximal solution h starting at some positive smooth h_0 in \mathcal{U} stays inside \mathcal{U} .*

Then ω is finite, that is, h ruptures in finite time.

From this theorem we deduce the following result on finite time rupture.

COROLLARY 4.1. *Consider (4.1), where $n \in (0, 1/2)$ and f is a power law, with q satisfying $(1-n)/2 \leq q < 3$. Then for any positive c_0 , there exists some L_0 depending on q and c_0 such that, for any $L \geq L_0$, there are solutions of (4.1) in $X(L, c_0)$ which rupture in finite time.*

Steady states arise in the study of long-time behavior of the solutions of (4.1), and the definition of a steady state is motivated from the energy dissipation relation. Indeed, a nonnegative function h in the closure of $X(L, c_0)$ is called a steady state of (4.1) if, on each component of the positivity set of h , $h_{xx} + f(h) = c$ holds for some

constant c (which could be different on components). Besides the constant states c_0/L that always exist in any $X(L, c_0)$, there could be other steady states in the same $X(L, c_0)$. Among others, a droplet is a steady state that is positive in some $(-a, a)$, $a \leq L/2$ and equal to 0 in $[-L/2, L/2] \setminus (-a, a)$ after a suitable translation. It is of full support if a is equal to $L/2$. The droplet has zero contact angle if its derivatives vanish at $\pm a$. A configuration of droplets is a steady state, where the positivity set $\{x \in (-L/2, L/2) : h(x) > 0\}$ is a union of at least two disjoint, open subintervals of $(-L/2, L/2)$. An array of identical droplets is a configuration of droplets in which, after a suitable translation, there is a partition of $(-L/2, L/2)$ into subintervals with equal length so that the restrictions of the steady state to these subintervals are congruent droplets with full support.

Criteria for the existence and multiplicity of positive steady states and droplets in a given $X(L, c_0)$ for general f and for power laws can be found in [LP1], and the stability of positive steady states, especially those for power laws, is studied systematically in [LP2], [LP3], and [LP4].

Proof of Theorem 4.1. We let $h(t)$ be a solution of (4.1) starting inside \mathcal{U} . Under (P_2) , h cannot blow up in any time, so either it ruptures in finite time or it is a positive global solution. Assuming that the latter holds we will derive a contradiction.

As h stays inside \mathcal{U} , it is uniformly bounded in H^1 -norm. From the energy dissipation relation

$$\mathcal{E}(h) + \int_0^T \int h^n \left| (h_{xx} + f(h))_x \right|^2 = \mathcal{E}(h_0),$$

where the energy is given by

$$\mathcal{E}(h) = \frac{1}{2} \int h_x^2 - \int F(h), \quad F'(z) = f(z), \quad F(0) = 0,$$

we know that

$$(4.3) \quad \int_0^\infty \int \left| (h_{xx} + f(h))_x \right|^2 < \infty.$$

On the other hand, the entropy dissipation relation is given by

$$\frac{d}{dt} \int h^{2-n} = -(2-n)(1-n) \int (h_{xx}^2 - f(h)h_{xx}).$$

As h is uniformly bounded for all time and $n < 1$, we can find a constant C_2 so that

$$(4.4) \quad \int_T^{T+1} \int (h_{xx}^2 - f(h)h_{xx}) \leq C_2$$

holds for all $T > 0$. From (4.3) and (4.4) it is easy to find $\{t_j\} \rightarrow \infty$ and a constant C_3 such that, for $h_j = h(t_j)$,

$$(4.5) \quad \int h_j^n (h_{jxx} + f(h_j))_x^2 \rightarrow 0$$

and

$$(4.6) \quad \int h_{jxx}^2 \leq C_3.$$

By passing to a subsequence if necessary, we may further assume that

$$\begin{aligned} h_j &\rightarrow \hat{h} \text{ in } C^{1,\frac{1}{2}}(S_L), \\ h_j &\rightarrow \hat{h} \text{ in } H^2(S_L), \\ h_j &\rightarrow \hat{h} \text{ in } C_{loc}^{2,\frac{1}{2}}(\{\hat{h} > 0\}), \end{aligned}$$

where \hat{h} is a steady state of (4.1). By (P_1) \hat{h} cannot be positive, so after a suitable translation, we may assume that there is a component (a, b) of the positivity set of h in $(-L/2, L/2)$ on which

$$\hat{h}'' + f(\hat{h}) = c.$$

(4.6) implies that $\hat{h}'(a) = \hat{h}'(b) = 0$. By integrating this equation we obtain

$$\frac{\hat{h}_x^2}{2} + F(\hat{h}) = c\hat{h}.$$

As $f(0) = 0$, the constant c must be positive unless the derivatives of \hat{h} vanish at the endpoints. But the derivative cannot be equal to 0, for the vanishing of both \hat{h} and \hat{h}_x at the endpoints would force \hat{h} to vanish identically by the uniqueness of the solution to the ordinary differential equation satisfied by \hat{h} . Therefore c is positive. By (P_1) there are two possibilities: First, there exists some $\delta_0 > 0$ such that $\hat{h} = 0$ in $[a - \delta_0, a]$ or $[b, b + \delta_0]$. Second, there exist two subintervals (a, b) and (b, b_1) in $(-L/2, L/2)$ such that $\hat{h}_{xx} + f(\hat{h}) = c_1, x \in (a, b)$ and $\hat{h}_{xx} + f(\hat{h}) = c_2, x \in (b, b_1)$, where c_1 and c_2 are different positive constants. Moreover, $\hat{h}_x^-(b) = \hat{h}_x^+(b) = 0$ holds.

Let us treat the first case first. For simplicity we take b to be 0 and so $\hat{h} = 0$ in $[0, \delta_0]$. For each $0 < \delta \leq \delta_0$, we rescale h_j and \hat{h} by

$$h_j(x) = \delta^2 H_j \left(\frac{x}{\delta} \right), \quad \hat{h}(x) = \delta^2 \hat{H} \left(\frac{x}{\delta} \right).$$

Then, for each fixed δ ,

$$\begin{aligned} H_j &\rightarrow \hat{H} && \text{in } C^{1,\frac{1}{2}}([-1, 1]) \cap C_{loc}^{2,\frac{1}{2}}((-1, 1)), \\ \hat{H}(y), \hat{H}_y(y) &\rightarrow 0 && \text{as } y \uparrow 0, \\ \hat{H}_{yy}(y) &\rightarrow c && \text{as } y \uparrow 0. \end{aligned}$$

We claim that there exists a constant C_* independent of δ such that

$$(4.7) \quad \int_{-1}^1 H_j^n H_{jyyy}^2 \geq C_* > 0 \quad \forall j \text{ large.}$$

To prove (4.7) we first fix $\sigma_0 \in (-1, 0)$ such that

$$0 > H_{jy}(\sigma_0) \geq -\frac{c}{32}$$

and

$$H_{jyy}(\sigma_0) \geq \frac{c}{2}$$

hold (σ_0 may depend on j). Next we claim that, for each large j , there exists some $\xi \in [\sigma_0, 1]$ such that

$$H_{jyy}(\xi) = \frac{c}{4}.$$

For, if not, then $H_{jyy} > c/4$ in $[\sigma_0, 1]$. By Taylor expansion at σ_0 , there exists some $z \in (\sigma_0, 1)$ such that

$$\begin{aligned} H_j(1) &= H_j(\sigma_0) + H_{jy}(\sigma_0)(1 - \sigma_0) + \frac{H_{jyy}(z)}{2}(1 - \sigma_0)^2 \\ &\geq -\frac{c}{32}(1 - \sigma_0) + \frac{c}{8}(1 - \sigma_0)^2 \\ &\geq \frac{c}{16}, \end{aligned}$$

contradicting the fact that $H_j(1)$ tends to 0 as $j \rightarrow \infty$. Thus $H_{jyy} = c/4$ somewhere in $(\sigma_0, 1)$. Let

$$\sigma_1 = \inf \{ \sigma > \sigma_0 : H_{jyy}(\sigma) = c/4 \} < 1.$$

We consider two subcases separately:

- (a) H_{jy} does not change sign in (σ_0, σ_1) , and
- (b) $H_{jy} = 0$ somewhere in (σ_0, σ_1) .

In subcase (a), $H_{jy} < 0$ in (σ_0, σ_1) because $H_{jy}(\sigma_0) < 0$. By Taylor expansion at σ_1 , for $y \in (\sigma_0, \sigma_1)$,

$$\begin{aligned} H_j(y) &= H_j(\sigma_1) + H_{jy}(\sigma_1)(y - \sigma_1) + \frac{H_{jyy}(z)}{2}(y - \sigma_1)^2 \\ &\geq H_j(\sigma_1) + \frac{c}{8}(y - \sigma_1)^2 \\ (4.8) \quad &\geq \frac{c}{8}(y - \sigma_1)^2 \quad \forall j \text{ large.} \end{aligned}$$

On the other hand,

$$\begin{aligned} \int_{\sigma_0}^{\sigma_1} |H_{jyyy}| &\geq |H_{jyy}(\sigma_1) - H_{jyy}(\sigma_0)| \\ &\geq \frac{c}{4}. \end{aligned}$$

So,

$$\begin{aligned} \int_{\sigma_0}^{\sigma_1} H_j^n |H_{jyyy}|^2 \int_{\sigma_0}^{\sigma_1} H_j^{-n} &\geq \left(\int_{\sigma_0}^{\sigma_1} |H_{jyyy}| \right)^2 \\ &\geq \frac{c^2}{16}. \end{aligned}$$

By (4.8),

$$\int_{\sigma_0}^{\sigma_1} H_j^{-n} \leq \left(\frac{8}{c} \right)^n \int_{\sigma_0}^{\sigma_1} \frac{1}{(y - \sigma_1)^{2n}} = \left(\frac{8}{c} \right)^n \frac{2}{1 - 2n}.$$

It follows that

$$\int_{-1}^1 H_j^n |H_{jyyy}|^2 \geq C_* \equiv \frac{c^{n+2}(1 - 2n)}{2^{3n+5}},$$

that is, (4.7) holds.

In subcase (b), let $\sigma_2 \in (\sigma_0, \sigma_1)$ be a point at which $H_{jy}(\sigma_2) = 0$. By Taylor expansion at σ_2 , for $y \in (\sigma_0, \sigma_1)$,

$$\begin{aligned} H_j(y) &= H_j(\sigma_2) + \frac{H_{jyy}(z)}{2}(y - \sigma_2)^2 \\ &\geq \frac{c}{8}(y - \sigma_2)^2, \end{aligned}$$

and the above argument still works to show that (4.7) holds for the same C_* .

Now, we scale (4.7) back to get

$$(4.9) \quad \begin{aligned} \int_{-\delta}^{\delta} h_j^n h_{jxxx}^2 &= \delta^{2n-1} \int_{-1}^1 H_j^n H_{jyyy}^2 \\ &\geq C_* \delta^{2n-1}. \end{aligned}$$

However, on the other hand, for small δ ,

$$\begin{aligned} \int_{-\delta}^{\delta} h_j^n h_{jxxx}^2 &\leq \int_{-\delta}^{\delta} h_j^n (h_{jxxx} + f'(h_j)h_{jx})^2 + 4 \int_{-\delta}^{\delta} f'(h_j)^2 h_{jx}^2 \\ &\leq \int_{-1}^1 h_j^n (h_{jxxx} + f'(h_j)h_{jx})^2 + C \int_{-\delta}^{\delta} (1 + h_j^{n+2q-2}) h_{jx}^2. \end{aligned}$$

As

$$\begin{aligned} \int h_j^{n+2q-2} h_{jx}^2 &= \frac{-1}{n + 2q - 1} \int h_j^{n+2q-1} h_{jxx} \\ &\leq C \left(\int h_j^{2(n+2q-1)} \right)^{1/2} \left(\int h_{jxx}^2 \right)^{1/2}, \end{aligned}$$

by (4.2), (4.5), and (4.6) we conclude that

$$\int_{-\delta}^{\delta} h_j^n h_{jxxx}^2$$

is bounded uniformly in j , but this is in conflict with (4.9) as $n < 1/2$ and δ could be arbitrarily small, so the classical solution cannot be a global one, contradiction holds.

The proof for the second case is similar. Without loss of generality let us assume that $c_1 > c_2$. Again we take $b = 0$ and fix some small δ_0 so that \hat{h} is decreasing in $(-\delta_0, 0)$ and increasing in $(0, \delta_0)$. For each $\delta < \delta_0$ we use the same rescalings as before to define H_j and \hat{H} and claim that (4.7) holds. For, we pick σ_0 in $(-1, 0)$ such that $H_{jyy}(\sigma_0) = c_1 - \varepsilon_0$, where ε_0 is a fixed number less than $(c_1 - c_2)/2$. As before one can show that there is some $\xi \in (\sigma_0, 1)$ such that $H_{jyy}(\xi) = (c_1 + c_2)/2$. So we can define

$$\sigma_1 = \inf\{\sigma > \sigma_0 : H_{jyy}(\sigma) = (c_1 + c_2)/2\} < 1.$$

We consider the two subcases as before. We have in subcase (a) $H_{jyy} \geq \frac{c_1+c_2}{4}(y-\sigma_1)^2$ and in subcase (b) $H_{jyy} \geq \frac{c_1+c_2}{4}(y-\sigma_2)^2$ for $y \in (\sigma_0, \sigma_1)$. In both cases

$$\int_{\sigma_0}^{\sigma_1} |H_{jyyy}| \geq \left| \frac{c_1 - c_2}{2} - \varepsilon_0 \right| > 0.$$

Now the same arguments as before would show that (4.7) holds, so the same contradiction can be drawn. \square

Proof of Corollary 4.1. First of all, from the interpolation inequality [F2] for functions in $H^1(S_L)$:

$$\int u^4 \leq C \left(\int |u| \right)^2 \left(\int u_x^2 + \int u^2 \right),$$

it is easy to see that, for $q \in (-1, 3)$, the energy $\mathcal{E}(h)$ satisfies

$$\int h_x^2 \leq C\mathcal{E}(h)$$

for some positive constant C . By a standard argument based on the direct method, the minimization problem

$$m = \inf\{\mathcal{E}(h) : h \in X(L, c_0)\}$$

admits a minimizer in the closure of $X(L, c_0)$, which is a steady state of (4.1).

Let us examine the candidates for the minimizer. First, according to Theorem 10 in [LP3], the constant state c_0/L becomes energy unstable whenever L satisfies $L > L_1$, where L_1 satisfies $L_1^{3-q}c_0^{q-1} = 4\pi^2$, and hence it cannot be a minimizer. Next, the classification results in [LP1] shows that, for any pair (L, c_0) , there are, up to translations, at most finitely many positive steady states with period $L/j, j \geq 2$ in $X(L, c_0)$, but they cannot be minimizers as they are also energy unstable due to Theorem 1 of [LP2]. Furthermore, the classification in section 5 of [LP1] asserts that positive steady states with minimal period and fully supported droplets with zero contact angle do not exist if $L > L_2$ for some L_2 , depending on c_0 and q . From the same source we know that there are, up to translations, at most finitely many arrays of identical droplets with zero angle in $X(L, c_0)$. We claim that they too cannot be minimizers. For, consider an array of identical droplets h , which vanishes at $-L/2$ and $L/2$. Define h^* to be the function equal to h on $(-L/2, L/2)$ and zero outside this interval. Then the rearrangement of h^* gives a function which is positive inside and zero outside $(-L/2, L/2)$. Extending the restriction of h^* on $(-L/2, L/2)$ as an L -periodic function h^{**} . From the properties of rearrangement h^{**} belongs to the closure of $X(L, c_0)$, with energy strictly lower than that of h ; see, for instance, [LiL]. Therefore h cannot be a minimizer. Since the energy is invariant under translation, we conclude that no array of identical droplets can be a minimizer of the energy.

As a result, for any L greater than L_1 and L_2 , if we denote by μ the minimum of the energies of all of the possible positive steady states, fully supported droplets with zero contact angle and arrays of identical droplets with zero contact angle in $X(L, c_0)$, then $m < \mu$. Now the corollary follows from Theorem 4.1 by taking

$$\mathcal{U} = \{h \in X(L, c_0) : \mathcal{E}(h) < (m + \mu)/2\}. \quad \square$$

Remark 4.1. A sharp version of the interpolation inequality used in the above proof is the Nagy's inequality [SP]:

$$\int u^4 \leq \frac{9}{4\pi^2} \left(\int |u| \right)^2 \left(\int u_x^2 \right),$$

which is valid for all $u \in H^1(\mathbf{R})$. By modifying this inequality one can show that the following inequality holds for all $u \in H^1(S_L)$:

$$\int u^4 \leq C \left(\int |u| \right)^2 \left(\int u_x^2 + \int u^2 \right)$$

for any constant $C > 9/4\pi^2$. Using this result, one can show that Corollary 4.1 continues to hold when $q = 3$ provided c_0 satisfies

$$3Bc_0^2 < 8\pi^2.$$

Appendix A. Suitable weak solutions. Here we give a sketch of the proof of Theorem 2.1. We shall treat problem (2.1), where $n \in (1.5, 3.5)$ in (2.1) and f satisfies (2.4) and (2.6). The proof when f satisfies (2.5) instead of (2.6) can be handled in a similar and simpler way.

Step 1. (2.1) admits a positive, classical solution for all time if $n \geq 4$ and $f \in C^{2,\alpha}([0, \infty))$.

This is achieved by considering the approximation problem ($\varepsilon > 0$)

$$\begin{cases} h_t + [(|h|^n + \varepsilon)(h_{xx} + f(h))]_x = 0, \\ h(0) = h_0 > 0 \text{ in } C^\infty(S_L). \end{cases}$$

One first shows that this problem admits a global classical solution h_ε . In fact, using the interpolation inequality on $H^1(S_L)$ [F2],

$$\int h^4 \leq C \left(\int h \right)^2 \left(\int h_x^2 + \int h^2 \right)$$

for some positive constant C in the energy dissipation relation

$$\begin{aligned} \frac{1}{2} \int h_{\varepsilon x}^2 - \int F(h_\varepsilon) + \int_0^t \int (|h|^n + \varepsilon) |(h_{\varepsilon xx} + f(h_\varepsilon))_x|^2 \\ = \frac{1}{2} \int h_{0x}^2 - \int F(h_0), \end{aligned}$$

where F is a primitive of f , with the help of the conservation of area

$$\int h = c_0 \equiv \int h_0,$$

and the growth condition $p < 3$ in (1.4) or $p = 3$ when c_0 is small one obtains a uniform H^1 -bound on h_ε . Using the entropy inequality one further shows that, although h_ε may change sign, its limit of as $\varepsilon \downarrow 0$ through a subsequence is a positive function h which solves (2.1). We refer to [BF] and [BP2] for details.

Step 2. (2.1) admits a positive, classical solution for all time if $n \geq 4$ and $f \in C^{2,\alpha}((0, \infty))$ satisfies (2.4) ($p < 3$) and (2.6) ($q > 1 - n/2$).

Consider the approximation problem ($\delta > 0$)

$$\begin{cases} h_t + [h^n (h_{xx} + f_\delta(h))]_x = 0, & f_\delta(z) = f(z + \delta), \\ h(0) = h_0 > 0 \text{ in } C^\infty(S_L). \end{cases}$$

By Step 1, this problem admits a positive, classical solution h_δ , which is uniformly bounded in H^1 -norm again by (2.4) ($p < 3$). Further we compute the change rate of the entropy $\int h_\delta^{2-n}$ as follows:

$$\begin{aligned} \frac{d}{dt} \int h_\delta^{2-n} &= (2-n)(1-n) \int h_{\delta x} (h_{\delta xx} + f_\delta(h_\delta))_x \\ &\leq -\frac{1}{2}(2-n)(1-n) \int h_{\delta xx}^2 + C \int f_\delta^2(h_\delta) \\ &\leq -\frac{1}{2}(2-n)(1-n) \int h_{\delta xx}^2 + C \left(\int h_\delta^{2q} + 1 \right) \\ &\leq C \left(\int h_\delta^{2-n} + 1 \right) \end{aligned}$$

as $q \geq 1 - n/2$. It follows that $\|h_\delta^{2-n}\|_{L^1}$ is uniformly bounded on every interval $[0, T]$. By passing to limit and using $n \geq 4$, h_δ subconverges to a positive, classical solution of (2.1), with initial datum h_0 .

We remark that the assertions in Steps 1 and 2 hold when the coefficient h^n is replaced by some positive, smooth $a(h)$, which behaves like h^n near 0.

Step 3. Equation (2.1) admits a suitable weak solution if $n \in (1.5, 4)$, (2.4) and (2.6) hold for $0 \leq p < 3$ and $q > 1 - n/2$.

Consider the approximation problem ($\varepsilon > 0$)

$$\begin{cases} h_t + \left[a_\varepsilon(h)(h_{xx} + f(h))_x \right]_x = 0, & a_\varepsilon(z) = \frac{z^4}{z^{4-n} + \varepsilon}, \\ h(0) = h_0 > 0 \text{ in } C^\infty(S_L). \end{cases}$$

As the coefficients behave like h^4 near 0, Step 2 guarantees that this problem admits a positive, classical solution h_ε for all $\varepsilon \leq 1$. Moreover, these solutions satisfy

$$(A.1) \quad \sup_{[0, T]} \int h_{\varepsilon x}^2 \leq C_1$$

and

$$(A.2) \quad \iint_{Q_T} a_\varepsilon(h_\varepsilon) \left| (h_{\varepsilon xx} + f(h_\varepsilon))_x \right|^2 \leq C_2.$$

It follows that

$$(A.3) \quad \|h_\varepsilon\|_{C_{x,t}^{\frac{1}{2}, \frac{1}{8}}(Q_T)} \leq C_3$$

(we refer to [BF] for the deduction of (A.3) from (A.1) and (A.2)), and, by (A.1), (A.2), and (2.6) ($q > 1 - n/2$),

$$\begin{aligned} \iint_{Q_T} a_\varepsilon(h_\varepsilon) h_{\varepsilon xxx}^2 &\leq 2 \iint_{Q_T} a_\varepsilon(h_\varepsilon) |(h_{\varepsilon xx} + f(h_\varepsilon))_x|^2 \\ &\quad + 2 \iint_{Q_T} a_\varepsilon(h_\varepsilon) f'^2(h_\varepsilon) h_{\varepsilon x}^2 \\ &\leq C \left(1 + \iint_{Q_T} h_\varepsilon^{n+2q-2} h_{\varepsilon x}^2 \right) \\ (A.4) \quad &\leq C \left(1 + \iint_{Q_T} h_{\varepsilon x}^2 \right) \leq C_4. \end{aligned}$$

Define

$$G_{s,\varepsilon}(z) = \int_z^M \int_y^M \frac{\tau^s}{a_\varepsilon(\tau)} d\tau dy,$$

where M is a fixed number greater than $\sup_{[0,T]} h_\varepsilon$ for all $\varepsilon \leq 1$. Then we have

$$(A.5) \quad \begin{aligned} \frac{d}{dt} \int G_{s,\varepsilon}(h_\varepsilon) &= - \int \left(h_\varepsilon^s h_{\varepsilon xx}^2 + s \int h_\varepsilon^{s-2} h_{\varepsilon x}^2 h_{\varepsilon xx} \right) \\ &\quad - \int f(h_\varepsilon) (h_\varepsilon^s h_{\varepsilon xx} + s h_\varepsilon^{s-1} h_{\varepsilon x}^2). \end{aligned}$$

We estimate the terms on the left of (A.5) as follows: For $s \in (-1/2, 1)$, there exists a positive number σ depending on s such that

$$\int h_\varepsilon^s h_{\varepsilon xx}^2 + s \int h_\varepsilon^{s-1} h_{\varepsilon x}^2 h_{\varepsilon xx} \geq \sigma_s \left(\int h_\varepsilon^s h_{\varepsilon xx}^2 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \right),$$

and, for $s = -1/2$,

$$\int h_\varepsilon^s h_{\varepsilon xx}^2 + s \int h_\varepsilon^{s-1} h_{\varepsilon x}^2 h_{\varepsilon xx} \geq 0.$$

Therefore, as $q \geq 1 - n/2$,

$$\begin{aligned} &\left| \int f(h_\varepsilon) (h_\varepsilon^s h_{\varepsilon xx} + s h_\varepsilon^{s-1} h_{\varepsilon x}^2) \right| \\ &\leq \delta \left(\int h_\varepsilon^s h_{\varepsilon xx}^2 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \right) + C_\delta \int h_\varepsilon^s f^2(h_\varepsilon) \\ &\leq \delta \left(\int h_\varepsilon^s h_{\varepsilon xx}^2 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \right) + C_\delta \left(\int h_\varepsilon^{s-n+2} + 1 \right) \\ &\leq \delta \left(\int h_\varepsilon^s h_{\varepsilon xx}^2 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \right) + C_\delta \left(1 + \int h_\varepsilon + \int G_{s,\varepsilon}(h_\varepsilon) \right). \end{aligned}$$

In the last step we have used the relation

$$(A.6) \quad \begin{aligned} G_{s,\varepsilon}(z) &= \int_z^M \int_y^M \tau^{s-n} d\tau dy + \varepsilon \int_z^M \int_y^M \tau^{s-4} d\tau dy \\ &= \frac{1}{(s-n+1)(s-n+2)} z^{s-n+2} - \frac{M^{s-n+1}}{s-n+1} z + \frac{M^{s-n+2}}{s-n+2} \\ &\quad + \varepsilon \left[\frac{1}{(s-3)(s-2)} z^{s-2} - \frac{M^{s-3}}{s-3} z + \frac{M^{s-2}}{s-2} \right]. \end{aligned}$$

Putting this identity back to (A.5) yields, for $s \in (-1/2, \min\{1, n - 2\})$,

$$(A.7) \quad \sup_{[0, T]} \left[\int h_\varepsilon^{s-n+2} + \varepsilon \int h_\varepsilon^{s-2} \right] + \sigma \iint_{Q_T} (h_\varepsilon^s h_{\varepsilon xx}^2 + h_\varepsilon^{s-2} h_{\varepsilon x}^4) \leq C_5.$$

Next, we fix $s' > -1/2$ such that $-\frac{1}{2} - s' + 2q > -n + 2$ still holds. By (A.5) and (A.7),

$$\begin{aligned} \frac{d}{dt} \int G_{-\frac{1}{2}, \varepsilon}(h_\varepsilon) &\leq C \int (h_\varepsilon^q + 1)(h_\varepsilon^{-1/2} h_{\varepsilon xx} + s h_\varepsilon^{-3/2} h_{\varepsilon x}^2) \\ &\leq C \int (h_\varepsilon^{s'} h_{\varepsilon xx}^2 + h_\varepsilon^{s'-2} h_{\varepsilon x}^4) + C \left(\int h_\varepsilon^{-1-s'+2q} + 1 \right) \\ &\leq C \int (h_\varepsilon^{s'} h_{\varepsilon xx}^2 + h_\varepsilon^{s'-2} h_{\varepsilon x}^4) + C \left(\int G_{-\frac{1}{2}, \varepsilon}(h_\varepsilon) + \int h_\varepsilon + 1 \right), \end{aligned}$$

so

$$(A.8) \quad \sup_{[0, T]} \left[\int h_\varepsilon^{\frac{3}{2}-n} + \varepsilon \int h_\varepsilon^{-\frac{5}{2}} \right] \leq C_6.$$

To let ε go to 0, first observe that, by (A.2) and (A.4),

$$(A.9) \quad \sqrt{a_\varepsilon(h_\varepsilon)} h_{\varepsilon xxx} \rightharpoonup h^{\frac{n}{2}} h_{xxx} \quad \text{in } L^2(Q_T),$$

and

$$(A.10) \quad \sqrt{a_\varepsilon(h_\varepsilon)} (h_{\varepsilon xx} + f(h_\varepsilon))_x \rightharpoonup h^{\frac{n}{2}} (h_{xx} + f(h))_x \quad \text{in } L^2(Q_T).$$

By parabolic regularity, for all $s > -1/2$ (since h_ε are uniformly bounded in ε , (A.7) holds for all $s > -1/2$),

$$\begin{aligned} h_\varepsilon^{\frac{s}{2}} h_{\varepsilon xx} &\rightarrow h^{\frac{s}{2}} h_{xx}, \\ h_\varepsilon^{\frac{s-2}{4}} h_{\varepsilon x} &\rightarrow h^{\frac{s-2}{4}} h_x, \quad \text{uniformly in any compact subset of } \{h > 0\}. \end{aligned}$$

Taking $s' \in (-1/2, s)$ and using (A.7), we have

$$\begin{aligned} \iint_{Q_T} |h_\varepsilon^{\frac{s}{2}} h_{\varepsilon xx} - h^{\frac{s}{2}} h_{xx}|^2 &= \left(\iint_{\{h \geq \delta\}} + \iint_{\{h < \delta\}} \right) |h_\varepsilon^{\frac{s}{2}} h_{\varepsilon xx} - h^{\frac{s}{2}} h_{xx}|^2 \\ &\leq o(\varepsilon) + C \delta^{s-s'} \rightarrow 0, \quad \text{as } \varepsilon \rightarrow 0 \text{ and } \delta \rightarrow 0, \end{aligned}$$

so

$$(A.11) \quad h_\varepsilon^{\frac{s}{2}} h_{\varepsilon xx} \rightarrow h^{\frac{s}{2}} h_{xx} \quad \text{in } L^2(Q_T).$$

Similarly, we have

$$(A.12) \quad h_\varepsilon^{\frac{s-2}{4}} h_{\varepsilon x} \rightarrow h^{\frac{s-2}{4}} h_x \quad \text{in } L^4(Q_T).$$

Now, using (A.9)–(A.12) to pass to limit we obtain a subsequence of h_ε converging uniformly to some h in Q_T . From (A.1)–(A.4), (A.7), and (A.8) we see that (H_1) and (H_2) hold. It remains to establish (H_3) .

LEMMA A.1. *There exists a set Γ of full measure in $[0, T]$ such that, for a sequence $\varepsilon_j \rightarrow 0$,*

$$h_{\varepsilon_j}(x, t) \rightarrow h(x, t) \quad \text{in } H^1(S_L) \quad \forall t \in \Gamma,$$

and

$$\int G_{s, \varepsilon_j}(h_{\varepsilon_j})(t) \phi^4 \rightarrow \int G_s(h)(t) \phi^4 \quad \forall t \in [0, T] \quad \forall \phi \in C^\infty(S_L), \quad s \in (-1/2, n-2).$$

Proof. It follows from (A.12) (taking $s = 2$) that

$$\iint_{Q_T} |h_{\varepsilon x} - h_x|^2 \rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0.$$

Therefore, there exists a subset $\Gamma \subset [0, T]$ with full measure such that, for a subsequence $\varepsilon = \varepsilon_j$,

$$\int |h_{\varepsilon x} - h_x|^2(t) \rightarrow 0 \quad \forall t \in \Gamma.$$

To prove the second statement, observe that each $t \in [0, T]$, K_3 implies that $h(x, t) > 0$ a.e. in S_L . So it follows (A.6), (A.7), the pointwise convergence

$$G_{s, \varepsilon}(h_\varepsilon)(t) \rightarrow G_s(h)(t) \quad \text{a.e. in } S_L,$$

and the dominated convergence theorem. \square

Now we can prove (H₃). First, by a direct computation we have

$$\begin{aligned} & \frac{1}{2} \int h_{\varepsilon x}^2 \phi^4 - \int F(h_\varepsilon) \phi^4 + \int_{t_0}^t \int a_\varepsilon(h_\varepsilon) \left| (h_{\varepsilon x x} + f(h_\varepsilon))_x \right|^2 \phi^4 \\ &= - \int_{t_0}^t \int a_\varepsilon(h_\varepsilon) (h_{\varepsilon x x} + f(h_\varepsilon))_x (h_{\varepsilon x x} + f(h_\varepsilon)) (\phi^4)_x \\ & \quad - \int_{t_0}^t \int a_\varepsilon(h_\varepsilon) (h_{\varepsilon x x} + f(h_\varepsilon))_x \left[h_{\varepsilon x x} (\phi^4)_x + h_{\varepsilon x} (\phi^4)_{x x} \right] \\ & \quad + \frac{1}{2} \int h_{\varepsilon x}^2(t_0) \phi^4 - \int F(h_\varepsilon(t_0)) \phi^4, \end{aligned} \tag{A.13}$$

and, for $s = -\frac{1}{2} + \delta \in (-\frac{1}{2}, \min\{1, n-2\})$,

$$\begin{aligned} & \frac{d}{dt} \int G_{s, \varepsilon}(h_\varepsilon) \phi^4 = \int G'_{s, \varepsilon}(h_\varepsilon) h_{\varepsilon t} \phi^4 \\ &= - \left[\int h_\varepsilon^s h_{\varepsilon x x}^2 \phi^4 + s \int h_\varepsilon^{s-1} h_{\varepsilon x}^2 h_{\varepsilon x x} \phi^4 \right] - \int h_\varepsilon^s h_{\varepsilon x} h_{\varepsilon x x} (\phi^4)_x \\ & \quad - \int f(h_\varepsilon) \left[h_\varepsilon^s h_{\varepsilon x x} \phi^4 + s h_\varepsilon^{s-1} h_{\varepsilon x}^2 \phi^4 + h_\varepsilon^s h_{\varepsilon x} (\phi^4)_x \right] \\ & \quad - \int \left[a'_\varepsilon(h_\varepsilon) G'_{s, \varepsilon}(h_\varepsilon) + a_\varepsilon(h_\varepsilon) G''_{s, \varepsilon}(h_\varepsilon) \right] h_{\varepsilon x} \left[h_{\varepsilon x x} + f(h_\varepsilon) \right] (\phi^4)_x \\ & \quad - \int a_\varepsilon(h_\varepsilon) G'_{s, \varepsilon}(h_\varepsilon) \left[h_{\varepsilon x x} + f(h_\varepsilon) \right] (\phi^4)_{x x} \\ & \leq -\sigma \left[\int h_\varepsilon^s h_{\varepsilon x x}^2 \phi^4 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \phi^4 \right] \\ & \quad + C \left[\int h_\varepsilon^{s+2} (\phi_x^4 + \phi^2 \phi_{x x}^2) + \int h_\varepsilon^s f^2(h_\varepsilon) \phi^4 \right], \end{aligned} \tag{A.14}$$

where the constants σ and C are positive constants. Note that, in the last inequality of (A.14), we have used

$$\begin{aligned} a_\varepsilon(z)G_{s,\varepsilon}(z) &= z^s, \\ |a'_\varepsilon(z)G'_{s,\varepsilon}(z)| &\leq C_7 z^s, \\ |a_\varepsilon(z)G'_{s,\varepsilon}(z)| &\leq C_8 z^{s+1}, \end{aligned}$$

and, for $-1/2 < s < 1$,

$$\begin{aligned} &\int h_\varepsilon^s h_{\varepsilon xx}^2 \phi^4 + s \int h_\varepsilon^{s-1} h_{\varepsilon x}^2 h_{\varepsilon xx} \phi^4 \\ &\geq \sigma_1 \left[\int h_\varepsilon^s h_{\varepsilon xx}^2 \phi^4 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \phi^4 \right] - C_1 \left| \int h_\varepsilon^{s-1} h_{\varepsilon x}^3 (\phi^4)_x \right| \\ &\geq \sigma_2 \left[\int h_\varepsilon^s h_{\varepsilon xx}^2 \phi^4 + \int h_\varepsilon^{s-2} h_{\varepsilon x}^4 \phi^4 \right] - C_2 \int h_\varepsilon^{s+2} \phi_x^4. \end{aligned}$$

Using (A.9)–(A.12) we can pass limit in (A.13) to obtain the local energy inequality. Similarly, by passing limit in the integral form of (A.14) using also Lemma A.1, we see that the local entropy inequality also holds.

Finally, by an approximation argument the initial data can be taken to be non-negative H^1 -functions, with finite $\|h_0^{3/2-n}\|_{L^1}$. One may consult [BF] for the proof of $\lim_{t \rightarrow 0} h(t) = h_0$ in H^1 . We have finished our outline of the proof of Theorem 2.1.

Appendix B. The parabolic Hausdorff dimension. Same as the Hausdorff measures, fourth order parabolic Hausdorff measures can be defined by a general construction due to Carathéodory; see [Fe]. Here we state their definitions and list some of their basic properties.

Let s be nonnegative, and $\delta > 0$. For any set X in \mathbf{R}^2 , we let \mathcal{C} be a collection of cylinders Q_α of the form $\{(x, t) : |x - x_0| < \rho_\alpha, \text{ and } |t - t_0| < \rho_\alpha^4\}$ ($\rho_\alpha > 0$) whose union contains X . Then let

$$\mathcal{P}_\delta^s(X) \equiv \inf \left\{ \sum \rho_\alpha^s : \text{All collections } \mathcal{C} \text{ that cover } X \text{ with } \rho_\alpha < \delta \right\},$$

and define the s -parabolic Hausdorff measure to be

$$\mathcal{P}^s(X) = \sup_{\delta \rightarrow 0} \mathcal{P}_\delta^s(X).$$

Then \mathcal{P}_δ^s and \mathcal{P}^s are outer measures for which all Borel subsets are measurable. We define the s -parabolic Hausdorff dimension to be the infimum of $\{s : \mathcal{P}^s(X) = 0.\}$ It is clear that $\mathcal{P}^t(X) = 0$ if $\mathcal{P}^s(X) < \infty$ whenever $t > s$, and $\mathcal{P}^5(X)$ is positive for any bounded, open X .

We state the following version of the Vitali covering theorem; see section 6 in [CKN] for a proof.

THEOREM B.1. *Let \mathcal{C} be a collection of parabolic cylinders described as above. Assuming that there is a uniform bound on the diameters of these cylinders, we can find a countable subcollection consisting of mutually disjoint cylinders $Q_j = \{(x, t) : |x - x_0| < \rho_j, |t - t_0| < \rho_j^4\}$ so that the union of the enlarged cylinders $Q'_j = \{(x, t) : |x - x_0| < 5\rho_j, |t - t_0| < 625\rho_j^4\}$ covers \mathcal{C} .*

Acknowledgments. We are grateful to the referees whose suggestions improve the presentation of this work.

REFERENCES

- [BBDP] E. BERETTA, M. BERTSCH, AND R. DAL PASSO, *Nonnegative solutions of a fourth order nonlinear degenerate parabolic equation*, Arch. Rat. Mech. Anal., 129 (1995), pp. 175–200.
- [B1] F. BERNIS, *Finite speed of propagation and continuity of the interface for thin viscous flows*, Adv. Differential Equations, 1 (1996), pp. 337–368.
- [B2] F. BERNIS, *Finite speed of propagation for thin viscous flows when $2 \leq n < 3$* , C. R. Acad. Sci. Paris Sér. I Math., 322 (1996), pp. 1169–1174.
- [BF] F. BERNIS AND A. FRIEDMAN, *Higher order nonlinear degenerate parabolic equations*, J. Diff. Eq., 83 (1990), pp. 179–206.
- [BBDK] A.L. BERTOZZI, M.P. BRENNER, T.F. DUPONT, AND L.P. KADANOFF, *Singularities and similarities in interface flow*, in Trends and Perspectives in Applied Mathematics, v. 100, L. Sirovich, ed., Springer-Verlag, New York, 1994, pp. 155–208.
- [BP1] A.L. BERTOZZI AND M.C. PUGH, *The lubrication approximation for their viscous films: Regularity and long time behavior of weak solutions*, Comm. Pure Appl. Math., 49 (1996), pp. 85–123.
- [BP2] A.L. BERTOZZI AND M.C. PUGH, *Long-wave instabilities and saturation in thin film equations*, Comm. Pure Appl. Math., 51 (1998), pp. 625–661.
- [BP3] A.L. BERTOZZI AND M.C. PUGH, *Finite-time blow-up of solutions of some long-wave unstable thin film equations*, Indiana Univ. Math. J., 49 (2000), pp. 1323–1366.
- [CKN] L. CAFFARELLI, R. KOHN, AND L. NIRENBERG, *Partial regularity of suitable weak solutions of the Navier-Stokes equations*, Comm. Pure Appl. Math., 35 (1982), pp. 771–831.
- [CS] Y.M. CHEN AND M. STRUWE, *Existence and partial regularity results for the heat flow for harmonic maps*, Math. Z., 201 (1989), pp. 83–103.
- [CDZ] K.S. CHOU, S.Z. DU, AND G.F. ZHENG, *On partial regularity of the borderline solution of semilinear parabolic problems*, Calc. Var. Partial Differential Equations, 30 (2007), pp. 251–275.
- [CK] K.S. CHOU AND Y.C. KWONG, *Finite time rupture for thin films under Van Der Waals forces*, Nonlinearity, 20 (2007), pp. 299–317.
- [DPGG] R. DAL PASSO, L. GIACOMELLI, AND G. GRÜN, *A waiting time phenomenon for thin film equations*, Ann. Sc. Norm. Super. Pisa Cl. Sci., 30 (2001), pp. 437–463.
- [E] J. EGGERS, *Nonlinear dynamics and breakup of free-surface flows*, Rev. Modern Phys., 69 (1997), pp. 865–929.
- [Ei] S.D. EIDELMAN, *Parabolic Systems*, North-Holland, Amsterdam, 1969.
- [ESS] L. ESCAURIAGA, G.A. SEREGIN, AND V. ŠVERÁK, *$L_{3,\infty}$ -solutions of Navier-Stokes equations and backward uniqueness*, Uspekhi Mat. Nauk, 58 (2003), pp. 3–44 (in Russian); Russian Math. Surveys, 58 (2003), pp. 211–250 (in English).
- [Fe] H. FEDERER, *Geometric Measure Theory*, Springer-Verlag, New York, 1969.
- [F1] A. FRIEDMAN, *Interior estimates for parabolic systems of partial differential equations*, J. Math. Mech., 7 (1958), pp. 393–418.
- [F2] A. FRIEDMAN, *Partial Differential Equations*, Holt, Rinehart and Winston, New York, 1969.
- [GPS] R.E. GOLDSTEIN, A.I. PESCI, AND M.J. SHELLY, *Instabilities and singularities in Hele-Shaw flow*, Phy. Fluids, 10 (1998), pp. 2701–2723.
- [HS] J. HULSHOF AND A. SHISHKOV, *The thin film equation with $2 \leq n < 3$: Finite speed of propagation in terms of the L^1 -norm*, Adv. Differential Equations, 3 (1998), pp. 625–642.
- [JL] H.Q. JIANG AND F.H. LIN, *Zero set of Sobolev functions with negative power of integrability*, Chinese Ann. Math. Ser. B, 25 (2004), pp. 65–72.
- [LP1] R.S. LAUGESEN AND M.C. PUGH, *Properties of steady states for thin film equations*, European J. Appl. Math., 11 (2000), pp. 293–351.
- [LP2] R.S. LAUGESEN AND M.C. PUGH, *Linear stability of steady states for thin film and Cahn-Hilliard type equations*, Arch. Ration. Mech. Anal., 154 (2000), pp. 3–51.
- [LP3] R.S. LAUGESEN AND M.C. PUGH, *Energy levels of steady states for thin film type equations*, J. Diff. Eq., 182 (2002), pp. 377–415.
- [LP4] R.S. LAUGESEN AND M.C. PUGH, *Heteroclinic orbits, mobility parameters and stability for thin film type equations*, Electron. J. Diff. Eq., 95 (2002), p. 29.

- [LiL] E.H. LIEB AND M. LOSS, *Analysis*, 2nd ed., American Mathematical Society, Providence, RI, 2001.
- [L] F.H. LIN, *A new proof of the Caffarelli-Kohn-Nirenberg theorem*, *Comm. Pure Appl. Math.*, 51 (1998), pp. 241–257.
- [LL] F.H. LIN AND C. LIU, *Partial regularity of the dynamic system modeling the flow of liquid crystals*, *Discrete Contin. Dyn. Syst.*, 2 (1996), pp. 1–22.
- [M] T.G. MYERS, *Thin films with high surface tension*, *SIAM Rev.*, 40 (1998), pp. 441–462.
- [ODB] A. ORON, S.H. DAVIS, AND S.G. BANKOFF, *Long-scale evolution of thin liquid films*, *Rev. Modern Phys.*, 69 (1997), pp. 931–980.
- [SP] D. SLEPČEV AND M.C. PUGH, *Self-similar blow-up of unstable thin film equations*, *Indiana Univ. Math. J.*, 54 (2005), pp. 1697–1738.
- [S] M. STRUWE, *On the evolution of harmonic maps in higher dimensions*, *J. Differential Geom.*, 28 (1988), pp. 485–502.

THE STABILITY OF THE NORMAL STATE OF SUPERCONDUCTORS IN THE PRESENCE OF ELECTRIC CURRENTS*

Y. ALMOG[†]

Abstract. The stability of the normal state of superconductors in the presence of electric currents is studied in the large domain limit. The model being used is the time-dependent Ginzburg–Landau model, in the absence of an applied magnetic field, and with the effect of the induced magnetic field being neglected. We find that if the current is nowhere perpendicular to the boundary, or if the minimal current on the boundary, at points where it is perpendicular to it, is greater than the critical current in the one-dimensional case, then the normal state is stable. We also prove some short-time instability when the current is both perpendicular to the boundary and smaller than the one-dimensional critical current.

Key words. superconductivity, Ginzburg–Landau, electric current, non–self-adjoint

AMS subject classifications. 82D55, 35P10, 35P15

DOI. 10.1137/070699755

1. Introduction. It is well known that when a superconductor is placed at a temperature lower than the critical one, it loses its electrical resistivity. This means that current can flow through a superconducting sample with a vanishingly small voltage drop. If one raises the current above a certain critical level, superconductivity will be destroyed and the material would revert to the normal state, even if the temperature is kept fixed below the critical one.

The reverse experiment can also be considered. One can flow a strong current through the sample which would set it in the normal state. Then, if we lower the current, there is a critical current where the sample would abruptly become purely superconducting. Though the two experiments substantially differ from each other from a theoretical point of view, hysteresis was not experimentally observed in the current-voltage characteristics of the sample [8, 14, 15].

We consider here the second experiment. To this end we must analyze the stability of the normal state. The model we use in this work is the time-dependent Ginzburg–Landau model [7, 4], presented here in a dimensionless form as follows:

$$(1.1a) \quad \frac{\partial \psi}{\partial t} + i\phi\psi = (\nabla - iA)^2 \psi + \psi(1 - |\psi|^2) \quad \text{in } \Omega,$$

$$(1.1b) \quad -\kappa^2 \nabla \times (\nabla \times A) - \sigma \left(\frac{\partial A}{\partial t} + \nabla \phi \right) = \frac{i}{2} (\bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}) + |\psi|^2 A \quad \text{in } \Omega,$$

$$(1.1c) \quad \psi = 0 \quad \text{on } \partial\Omega_c,$$

$$(1.1d) \quad (i\nabla + A)\psi \cdot \nu = 0 \quad \text{on } \partial\Omega_i.$$

*Received by the editors August 9, 2007; accepted for publication (in revised form) April 8, 2008; published electronically August 20, 2008. The author was supported by NSF grant DMS 0604467 and by the Summer Stipend Program at LSU.

<http://www.siam.org/journals/sima/40-2/69975.html>

[†]Department of Mathematics, Louisiana State University, Baton Rouge, LA 70803 (almog@math.lsu.edu).

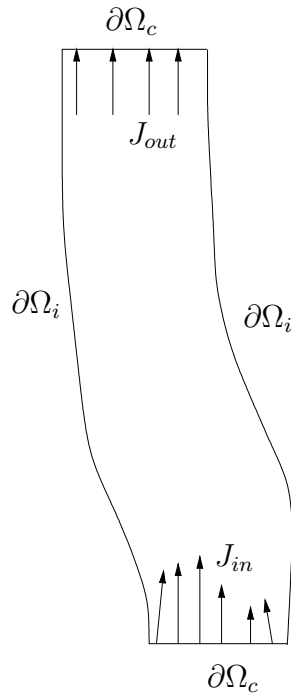


FIG. 1. Typical superconducting sample. The arrows denote the direction of the current flow (J_{in} for the inlet and J_{out} for the outlet).

In (1.1) ψ is the superconducting order parameter, so that $|\psi|$ represents the number density of superconducting electrons. Superconductors for which $|\psi| = 1$ are said to be wholly superconducting, and those for which $\psi = 0$ are said to be at the normal state. A is the magnetic vector potential and ϕ is the electric scalar potential. The constant σ is a measure of the normal conductivity of the superconducting material so that $-\sigma\nabla\phi$ is the normal current, and κ is the Ginzburg–Landau parameter. Length has been scaled with the coherence length ξ , which is the natural length-scale for variations in ψ . The domain $\Omega \subset\subset \mathbb{R}^n$ ($n = 1, 2, 3$), where the superconducting sample resides, is smooth and has an interface, denoted by $\partial\Omega_c$, with a conducting metal which is at the normal state. The rest of the boundary, denoted by $\partial\Omega_i$, is adjacent to an insulator. We allow nonsmoothness of $\partial\Omega$ in the sense that $\partial\Omega_c$ and $\partial\Omega_i$ are required to be perpendicular to each other in order to include cylindrical-like domains. Figure 1 presents a typical two-dimensional sample, where the current flows into the sample from one part of $\partial\Omega_c$ and exits from another part, disconnected from the first one. Most wires would fall into the above class of domains.

Equations (1.1) are gauge invariant in the sense that they are invariant under transformations of the form

$$A \rightarrow A + \nabla\omega, \quad \phi \rightarrow \phi - \frac{\partial\omega}{\partial t}, \quad \psi \rightarrow \psi e^{i\kappa\omega}.$$

Note that none of the important physical properties, i.e., $|\psi|$, the magnetic field $H = \nabla \times A$, and the electric field $E = -\partial A/\partial t - \nabla\phi$, are altered by the above transformation.

To obtain a well-posed problem one must add to (1.1) initial conditions, and the equations satisfied by A and ϕ outside Ω , that is, the Maxwell equations. Continuity

of the tangential components of A and $\nabla \times A$ through $\partial\Omega$ and some conditions on $\nabla \times A$ at infinity should be required as well. Since these details are irrelevant in the context of the present contribution, we omit them. Interested readers may be able to find them in [4, 6].

We consider here the stability of the normal state. If $\psi \equiv 0$, we obtain that the steady state solution must satisfy

$$\nabla \times H = -\frac{\sigma}{\kappa^2} \nabla \phi,$$

and hence ϕ is harmonic in Ω . To obtain ϕ we thus need to solve the following problem:

$$(1.2) \quad \begin{cases} \Delta \phi = 0 & \text{in } \Omega, \\ \phi = \phi_0(x) & \text{on } \partial\Omega_c, \\ \frac{\partial \phi}{\partial \nu} = 0 & \text{on } \partial\Omega_i. \end{cases}$$

We note that instead of prescribing the potential on $\partial\Omega_c$ we can prescribe the current in the normal direction to the boundary $J_n(x) = -\sigma \partial\phi / \partial\nu$. For simplicity we assume that $\nabla\phi \neq 0$ everywhere in Ω . This can easily be achieved; for instance, for the samples described in Figure 1, if $\inf \phi$ on one connected component of $\partial\Omega_c$ is greater than the $\sup \phi$ on the other connected component, then $\nabla\phi$ never vanishes.

Once (1.2) is solved, one can solve for the magnetic field. Here we need to solve a problem in \mathbb{R}^n ,

$$(1.3) \quad \begin{cases} \nabla \times H = -\frac{\sigma}{\kappa^2} \nabla \phi = \frac{1}{\kappa^2} J & \text{in } \Omega, \\ \nabla \times H = 0 & \text{in } \mathbb{R}^n \setminus \Omega, \\ H \rightarrow h_{ex} & \text{as } |x| \rightarrow \infty, \end{cases}$$

with H continuous across $\partial\Omega$. We assume zero applied magnetic field ($h_{ex} = 0$). To simplify our problem further we shall assume $A = H = 0$. This assumption can be justified in the case where $\kappa \gg 1$ in view of (1.3). However, since we intend to consider large domains, one must assume that $\kappa \gg \text{diam}\Omega$. In real-world coordinates this means that our domain size must be much larger than the coherence length ξ but also much smaller than the penetration depth λ , which is the length-scale characterizing variations in H ($\kappa = \lambda/\xi$). While this assumption significantly limits the validity of our results, it has been made very often by physicists [15, 8] and is reasonable to adopt as a starting point.

Once the above assumption is adopted one obtains

$$\begin{cases} \frac{\partial \psi}{\partial t} - \Delta \psi + i\phi\psi - \psi(1 - |\psi|^2) = 0 & \text{in } \Omega, \\ \psi = 0 & \text{on } \partial\Omega_c, \\ \frac{\partial \psi}{\partial \nu} = 0 & \text{on } \partial\Omega_i, \\ \psi(x, 0) = \psi_0(x) & \text{in } \Omega. \end{cases}$$

It should be noted that in [6] Du and Gray prove, within the framework of a more general case, convergence in the limit $\kappa \rightarrow \infty$ of (1.1) to a different limit problem where the magnetic field is not negligible. The domain size considered there is, however, much larger than in our case, as it is comparable with the penetration depth.

Linearizing the above near the normal state, we obtain

$$(1.4) \quad \begin{cases} \frac{\partial \psi}{\partial t} - \Delta \psi + i\phi \psi - \psi = 0 & \text{in } \Omega, \\ \psi = 0 & \text{on } \partial\Omega_c, \\ \frac{\partial \psi}{\partial \nu} = 0 & \text{on } \partial\Omega_i, \\ \psi(x, 0) = \psi_0(x) & \text{in } \Omega. \end{cases}$$

The above problem has been analyzed by physicists in one-dimensional settings. Ivlev and Kopnin review these results in [8]. In these settings we have $\phi = Jx + \mu$, where J (denoting the current) and μ are constants. Previous results include the closed form solution of (1.4) in \mathbb{R} for any initial condition. In \mathbb{R}_+ , the first critical current J_c for which a steady state solution ($\partial|\psi|/\partial t = 0$) exists is found. Then weakly nonlinear analysis is performed, where it is shown that [9] the bifurcation taking place near $J = J_c$ is subcritical (i.e., unstable).

We consider (1.4) in three-dimensional settings. While all the results are stated for three-dimensional objects, they are equally valid for two-dimensional objects as well. We deal with (1.4) in the large domain limit; i.e., we consider a domain Ω_R which is obtained from Ω via the transformation

$$(1.5) \quad x \rightarrow Rx.$$

The portions of the boundaries $\partial\Omega_R^c$ and $\partial\Omega_R^i$ are similarly obtained from $\partial\Omega_c$ and $\partial\Omega_i$, respectively. To keep $\nabla\phi$ unaltered we consider also potentials of the form

$$\phi_R = R\phi(x/R).$$

Thus, we consider the following problem:

$$(1.6) \quad \begin{cases} \frac{\partial \psi_R}{\partial t} - \Delta \psi_R + i\phi_R \psi_R - \psi_R = 0 & \text{in } \Omega_R, \\ \psi_R = 0 & \text{on } \partial\Omega_R^c, \\ \frac{\partial \psi_R}{\partial \nu} = 0 & \text{on } \partial\Omega_R^i, \\ \psi_R(x, 0) = \psi_0(x) & \text{in } \Omega_R. \end{cases}$$

Our main result is the following.

THEOREM 1.1. *Let ϕ satisfy (1.2). Let $\partial\Omega_c^n$ denote the portion of $\partial\Omega_c$ on which $\nabla\phi$ is perpendicular to the boundary. Let further*

$$J_m = \min_{x \in \partial\Omega_c^n} \left| \frac{\partial\phi}{\partial\nu} \right|,$$

and let J_c denote the critical current for the problem in \mathbb{R}_+ (which is precisely defined in (2.19)). Suppose further that $|J| > 0$ everywhere on $\partial\Omega_c$. Then, if $J_m > J_c$, or if $\partial\Omega_c^n$ is empty, there exists $R_0 > 0$ such that $\psi_R \equiv 0$ is a stable solution of (1.6) in the sense of $L^2(\Omega_R)$ for all $R > R_0$.

Furthermore, if $J_m < J_c$, there exists $\psi_0 \in L^2(\Omega_R)$ and $T_R > 0$ such that

$$\liminf_{R \rightarrow \infty} \frac{T_R}{\ln R} > 0,$$

and such that the solution ψ_R of (1.6) has the following property:

$$(1.7) \quad t < T_R \Rightarrow \|\psi_R\|_{L^2(\Omega_R)} > \frac{1}{2} \|\psi_0\|_{L^2(\Omega_R)} e^{\beta t},$$

where

$$\beta = 1 - (J/J_c)^{2/3}.$$

In the next section we review and enhance the results in [8] in the one-dimensional case. In section 3 we extend some of the results in section 2 to unbounded three-dimensional domains. In section 4 we provide the proof of the theorem. Finally, in the last section we highlight some possible directions for future research. The appendix provides a technical result needed in section 2.

2. One-dimensional problems. In this section we consider (1.4) in two different one-dimensional settings: on \mathbb{R} and on \mathbb{R}_+ . The solution of these simple problems would provide us with some important intuition on the solution of (1.6) in three-dimensional bounded domains in the large domain limit. In both cases we shall assume $\phi = Jx$, i.e., that the current is uniform and equals J throughout the sample.

2.1. Infinite one-dimensional domain. Here we consider the problem

$$(2.1) \quad \frac{\partial \psi}{\partial t} - \psi'' - \psi + iJx\psi = 0.$$

We consider here only the case $J > 0$. Otherwise, if $J < 0$, we can consider the complex conjugate of (2.1). Applying the coordinate transformation

$$(2.2) \quad x \rightarrow J^{1/3}x; \quad t \rightarrow J^{2/3}t,$$

we obtain the problem

$$\frac{\partial \psi}{\partial t} + \mathcal{L}\psi = \lambda_J \psi,$$

where

$$(2.3) \quad \mathcal{L}\psi = -\psi'' + ix\psi,$$

and $\lambda_J = J^{-2/3}$.

We first focus our interest on the spectrum of the operator $\mathcal{L} : \mathcal{D}_{\mathbb{R}}(\mathcal{L}) \rightarrow L^2(\mathbb{R}, \mathbb{C})$, where $\mathcal{D}_{\mathbb{R}}(\mathcal{L})$ is the dense subset of $L^2(\mathbb{R}, \mathbb{C})$ defined as

$$\mathcal{D}_{\mathbb{R}}(\mathcal{L}) = \{u \in L^2(\mathbb{R}, \mathbb{C}) \mid -u'' + ixu \in L^2(\mathbb{R}, \mathbb{C})\}.$$

LEMMA 2.1. *The operator $\mathcal{L} - \lambda I$ is invertible for all $\lambda \in \mathbb{C}$. (Thus $\sigma(\mathcal{L}) = \phi$.)*

Proof. It is sufficient to consider here $\lambda \in \mathbb{R}$. Otherwise, if $\lambda = \lambda_r + i\lambda_i$, we apply the transformation $x \rightarrow x - \lambda_i$.

Though it is not necessary, we first prove injectivity of $\mathcal{L} - \lambda I$. We shall later make use of a similar argument in three dimensions. To the problem

$$\mathcal{L}u = \lambda u$$

we apply the Fourier transform

$$(2.4) \quad \hat{u}(\omega) = \mathcal{F}(u) = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-i\omega x} u(x) dx$$

to obtain

$$\omega^2 \hat{u} - \frac{\partial \hat{u}}{\partial \omega} = \lambda \hat{u}.$$

Since the above equation doesn't possess any nontrivial solutions with bounded $L^2(\mathbb{R})$ norm, $\mathcal{L} - \lambda I$ must be injective.

To find the spectrum of \mathcal{L} we construct the Green's function of $\mathcal{L} - \lambda I$. Let

$$(2.5) \quad \begin{cases} w_1(x, \lambda) = A_i(e^{i\pi/6}x + \lambda e^{i2\pi/3}), \\ w_2(x, \lambda) = A_i(-e^{-i\pi/6}x + \lambda e^{-i2\pi/3}), \end{cases}$$

where A_i denotes Airy's function [1] (cf. (A.2) in Appendix A for the asymptotic behavior of A_i). It is easy to show that w_1, w_2 constitute a fundamental set of solutions to $(\mathcal{L} - \lambda I)u = 0$ [13]. We can now write the Green's function in the form

$$G(x, \xi, \lambda) = \begin{cases} \frac{w_2(\xi, \lambda)}{W(w_1, w_2)} w_1(x, \lambda), & x > \xi, \\ \frac{w_1(\xi, \lambda)}{W(w_1, w_2)} w_2(x, \lambda), & x < \xi, \end{cases}$$

where $W(w_1, w_2)$ denote the Wronskian (which is a constant that clearly doesn't vanish since w_1 and w_2 are linearly independent).

Using the asymptotic behavior of Airy's functions [1], we can show, using the same procedure applied in Appendix A to the semi-infinite case, that $G(\cdot, \lambda) \in L^2(\mathbb{R}^2, \mathbb{C})$ for all $\lambda \in \mathbb{R}$. The lemma is proved. \square

The spectrum of the operator $\mathcal{L} : \mathcal{D}_{\mathbb{R}}(\mathcal{L}) \rightarrow L^2(\mathbb{R}, \mathbb{C})$ can therefore teach us very little on the stability of the normal state in this case. The following lemma proves its stability in a direct manner.

LEMMA 2.2. *The trivial solution of (2.1), $\psi \equiv 0$, is globally stable in L^2 ; i.e., if $f(x) \in L^2(\mathbb{R}, \mathbb{C})$, then the solution of (2.1), $\psi(x, t)$, satisfying the initial condition $\psi(x, 0) = f(x)$ satisfies $\psi(x, t) \xrightarrow[t \rightarrow \infty]{} 0$ in $L^2(\mathbb{R}, \mathbb{C})$.*

Proof. As long as $\psi(x, t) \in L^2(\mathbb{R}, \mathbb{C})$, we can apply to (2.1) the Fourier transform (2.4). We obtain

$$\frac{\partial \hat{\psi}}{\partial t} + \omega^2 \hat{\psi} - \frac{\partial \hat{\psi}}{\partial \omega} - \lambda_J \hat{\psi} = 0.$$

The unique solution to the above problem is given by

$$(2.6) \quad \hat{\psi}(\omega, t) = \hat{f}(\omega) \exp \left\{ -\omega^2 t - \omega t^2 - \frac{1}{3} t^3 + \lambda_J t \right\},$$

in which $\hat{f}(\omega)$ denotes the Fourier transform of f . Integrating the modulus square of the above over \mathbb{R} with respect to ω , we obtain

$$\|\psi(\cdot, t)\|_2 \leq C \|f\|_2 \exp \left\{ -\frac{1}{12} t^3 + \lambda_J t \right\},$$

where $\|\cdot\|_2$ denotes the $L^2(\mathbb{R}, \mathbb{C})$ norm. \square

The above superlinear convergence is in accordance with the result proved in the previous lemma by which the spectrum of \mathcal{L} is empty. We note that in [8] the inverse transform of (2.6) is obtained, but no decay proof is given.

2.2. Semi-infinite one-dimensional domain. We now consider (2.1) on \mathbb{R}_+ . We concentrate here on the Dirichlet boundary condition $\psi(0) = 0$; however, the same analysis applies to Neumann and mixed boundary conditions as well.

We start by proving the following result on the spectrum of the operator $\mathcal{L} : \mathcal{D}_{\mathbb{R}_+}(\mathcal{L}) \rightarrow L^2(\mathbb{R}, \mathbb{C})$, where

$$\mathcal{D}_{\mathbb{R}_+}(\mathcal{L}) = \{u \in L^2(\mathbb{R}_+, \mathbb{C}) \mid -u'' + ixu \in L^2(\mathbb{R}_+, \mathbb{C}), u \in H_0^1(\mathbb{R}_+, \mathbb{C})\}.$$

LEMMA 2.3.

1. There exists a sequence of eigenvalues $\{\lambda_n\}_{n=1}^\infty$ and eigenfunctions of \mathcal{L} , with unity norm, $\{u_n\}_{n=1}^\infty \subset \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$, i.e.,

$$\mathcal{L}u_n = \lambda_n u_n.$$

2. We have

$$(2.7) \quad m = \max_{n \in \mathbb{N}} \Re \lambda_n > 0.$$

3. $\text{span}\{u_n\}_{n=1}^\infty = L^2(\mathbb{R}_+, \mathbb{C})$.

4. Suppose that $u, v \in L^2(\mathbb{R}_+, \mathbb{C})$ can be represented in the form

$$v = \sum_{n=1}^\infty \alpha_n u_n; \quad w = \sum_{n=1}^\infty \beta_n u_n,$$

where the convergence is in the L^2 sense. Let then

$$(2.8) \quad \langle v, w \rangle_U = \sum_{n=1}^\infty \alpha_n \bar{\beta}_n.$$

Then,

$$(2.9) \quad \inf_{v \in \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})} \Re \langle v, \mathcal{L}v \rangle_U \geq m \|v\|_U^2,$$

where $\|\cdot\|_U$ denotes the norm induced by (2.8).

Proof. Let

$$z = -ix + \lambda.$$

Let $u(x, \lambda) \in \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$ denote an eigenfunction of \mathcal{L} , i.e., $\mathcal{L}u = \lambda u$. Let $v(z, \lambda) = u(x, \lambda)$. We have

$$(2.10) \quad \begin{cases} \frac{\partial^2 v}{\partial z^2} - zv = 0, & z \in \mathbb{C}, \\ v(\lambda, \lambda) = 0. \end{cases}$$

Since $u \in L^2(\mathbb{R}_+, \mathbb{C})$, v must be subdominant (i.e., it decays exponentially fast [13]) in the sector

$$S_1 : -\pi < \arg z < -\frac{\pi}{3}.$$

The decaying solution of (2.10) in S_1 is given by (cf. [13])

$$v = A_i(e^{2\pi i/3} z).$$

Since the zeros of Airy's functions are eigenvalues of the self-adjoint operator $d^2/dx^2 - x$ in $\mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$, they must all be real. Let $\{\mu_n\}_{n=1}^\infty \subset \mathbb{R}$ denote the zeroes of Airy's function on the real axis. By the maximum principle they must all be strictly negative. We arrange them so that $\mu_n \downarrow -\infty$. As every eigenvalue of \mathcal{L} must satisfy $v(\lambda, \lambda) = 0$, the set $\{\lambda_n\}_{n=1}^\infty$, where

$$(2.11) \quad \lambda_n = e^{-i2\pi/3}\mu_n,$$

contains all the eigenvalues of \mathcal{L} . Since $\mu_1 < 0$ we have that

$$m = \Re\lambda_1 = -\frac{\mu_1}{2} > 0.$$

The set $\{\tilde{u}_n\}_{n=1}^\infty$ of eigenfunctions of \mathcal{L} in $\mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$ is given by

$$(2.12) \quad \tilde{u}_n = A_i(e^{i2\pi/3}(-ix + \lambda_n)) = A_i(e^{i\pi/6}x + \mu_n) \quad \forall n \in \mathbb{N}.$$

We then set

$$(2.13) \quad u_n = \frac{\tilde{u}_n}{\|\tilde{u}_n\|_{L^2(\mathbb{R}_+)}}.$$

To prove that $\{u_n\}_{n=1}^\infty$ is complete in $L^2(\mathbb{R}, \mathbb{C})$ we consider the resolvent $\mathcal{L}_\lambda^{-1} = (\mathcal{L} - \lambda I)^{-1}$ (which is also the modified resolvent of \mathcal{L}^{-1}). We have

$$(2.14a) \quad \mathcal{L}_\lambda^{-1}f = \int_0^\infty \tilde{G}(x, \xi, \lambda)f(\xi)d\xi,$$

in which

$$(2.14b) \quad \tilde{G}(x, \xi, \lambda) = \begin{cases} \frac{\tilde{w}_2(\xi, \lambda)}{W(w_1, \tilde{w}_2)}w_1(x, \lambda), & x > \xi, \\ \frac{w_1(\xi, \lambda)}{W(w_1, \tilde{w}_2)}\tilde{w}_2(x, \lambda), & x < \xi, \end{cases}$$

where

$$(2.15) \quad \tilde{w}_2(x, \lambda) = \frac{w_1(0, \lambda)}{w_2(0, \lambda)}w_1(x, \lambda) - w_2(x, \lambda),$$

and w_1 and w_2 are given in (2.5). In Appendix A, we prove that $\tilde{G} \in L^2(\mathbb{R}_+ \times \mathbb{R}_+)$, and that

$$(2.16) \quad \|\tilde{G}(\cdot, \cdot, \lambda)\|_{L^2(\mathbb{R}_+ \times \mathbb{R}_+)} \leq e^{M|\lambda|^{3/2}},$$

as long as $\lambda \notin \{\lambda_n\}_{n=1}^\infty$.

Let $u \in \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$. We now multiply $\mathcal{L}u$ by $e^{i\theta}\bar{u}$ and integrate over \mathbb{R}_+ to obtain

$$\Re \langle e^{i\theta}\mathcal{L}u, u \rangle = \int_0^\infty (\cos\theta|u'|^2 - \sin\theta x|u|^2)dx.$$

By Theorem 12.8 in [2] we have that every direction $e^{i\arg\lambda}$ with

$$\pi/2 < \arg\lambda < 3\pi/2$$

is a direction of minimal growth of the resolvent of $e^{i\theta}\mathcal{L}$ for every $-\pi/2 < \theta < 0$. Consequently, every direction $e^{i \arg \lambda}$ with

$$\pi/2 < \arg \lambda < 2\pi$$

is a direction of minimal growth of \mathcal{L}_λ^{-1} , i.e. (cf. [2]),

$$(2.17) \quad \|\mathcal{L}_\lambda^{-1}\| \sim \mathcal{O}(|\lambda|^{-1}) \quad \pi/2 < \arg \lambda < 2\pi.$$

We now apply the same argument used in the proof of Theorem 16.4 in [2]. Let $f \in L^2(\mathbb{R}_+)$, and let $g \in V^\perp$, where $V = \text{span}\{A_i(e^{i\pi/6}x + \mu_n)\}_{n=1}^\infty = \text{sp}'(\mathcal{L}^{-1})$. Then,

$$F(\lambda) = \langle \mathcal{L}_\lambda^{-1}f, g \rangle$$

is an entire function (cf. [2]) of λ satisfying, by (2.16) and (2.17),

$$\begin{cases} |F(\lambda)| \leq \frac{C}{|\lambda|}, & \pi/2 < \arg \lambda < 2\pi, \\ |F(\lambda)| \leq Ce^{M|\lambda|^{3/2}}, & \lambda \in \mathbb{C}. \end{cases}$$

By Theorem 16.1 in [2] (or the Phragmen–Leindelöf theorem) and Liouville’s theorem, we must have $F(\lambda) \equiv 0$. Hence, $\bar{V} = \text{range}(\mathcal{L}^{-1}) = \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$.

It remains still necessary to prove (2.9). Let $w \in \mathcal{D}_{\mathbb{R}_+}(\mathcal{L})$. Then,

$$w = \sum_{n=1}^\infty \alpha_n u_n$$

and

$$\mathcal{L}w = \sum_{n=1}^\infty \alpha_n \lambda_n u_n.$$

Hence

$$\Re \langle w, \mathcal{L}w \rangle_U = \sum_{n=1}^\infty \Re \lambda_n |\alpha_n|^2 \geq m \|w\|_U^2. \quad \square$$

Similar techniques were used in [12, 10] to prove completeness of the system of eigenfunctions of some nonlinear eigenvalue problems in \mathbb{R} . We note that the set $\{u_n\}_{n=1}^\infty$ is not a basis in the usual sense in Banach spaces. In fact, it has been demonstrated in [5] that the system $\{\tilde{u}_n, \bar{\tilde{u}}_n\}_{n=1}^\infty$, which is a biorthogonal system after we appropriately normalize it, is wild. This means that $\|\tilde{u}_n\|_2$ grows faster than any algebraic rate as $n \rightarrow \infty$.

We now prove the existence of a critical current J_c obtained in [8, 9].

LEMMA 2.4. *Let $\psi(x, t) \in H_0^2(\mathbb{R}_+ \times \mathbb{R}_+, \mathbb{C})$ denote a solution of the equation*

$$(2.18) \quad \frac{\partial \psi}{\partial t} - \frac{\partial^2 \psi}{\partial x^2} - \psi + iJx\psi = 0 \quad \text{in } \mathbb{R}_+ \times \mathbb{R}_+.$$

If

$$(2.19) \quad J > J_c = \left(-\frac{\mu_1}{2}\right)^{-3/2},$$

in which μ_1 is the rightmost zero of Airy's function, then $\|\psi(\cdot, t)\|_U \xrightarrow[t \rightarrow \infty]{} 0$. Otherwise, if $J < J_c$, then $\psi \equiv 0$ is an unstable solution of (2.18).

Proof. We first apply (2.2) to obtain

$$\begin{cases} \frac{\partial \psi}{\partial t} + \mathcal{L}\psi - \lambda_J \psi = 0, \\ \psi(0, t) = 0, \\ \psi(x, 0) = \psi_0(x), \end{cases}$$

where $\|\psi_0\|_U < \infty$. Taking the inner product (2.8) of the above equation with ψ , we obtain

$$\frac{1}{2} \frac{d}{dt} \|\psi\|_U^2 \leq \left(\lambda_J + \frac{\mu_1}{2} \right) \|u\|_U^2.$$

Hence, if (2.19) is satisfied, we have $\lambda_J + \mu_1/2 < 0$, and hence ψ tends to 0 exponentially fast as $t \rightarrow \infty$, in the $\|\cdot\|_U$ sense.

If $J < J_c$, let $\psi(x, 0) = u_1(x)$, where u_1 is given in (2.13). Then

$$\psi(x, t) = u_1(x)e^{(\lambda_J + \mu_1/2)t}.$$

Returning to the original variables by applying the inverse of (2.2), we obtain

$$\psi(x, t) = u_1(J^{1/3}x)e^{[1 - (J/J_c)^{2/3}]t},$$

and hence $\psi \equiv 0$ is unstable. □

3. Unbounded domains in \mathbb{R}^3 . In this section we consider several different problems: in \mathbb{R}^3 , in \mathbb{R}_+^3 , and in a quarter-space. In contrast with the previous section, we consider these problems only to the extent needed in the next section, that is, we analyze the existence of eigenvalues with nonpositive real part of the elliptic operator in the right-hand side of (1.4).

3.1. Eigenfunctions in \mathbb{R}^3 . We consider here (1.4) with $\phi = Jx_1$. This choice summarizes all possible electric potentials with constant gradients, as the problem is invariant to translations and rotations. Thus, we have for every eigenfunction u ,

$$(3.1) \quad -\Delta u - u + iJx_1u = -\lambda u.$$

We shall assume here that $\lambda \in \mathbb{R}$; otherwise we can apply the transformation $x_1 \rightarrow x_1 - \Im \lambda / J$.

It is easy to find all L^2 solutions of (3.1) in \mathbb{R}^3 as follows: apply the Fourier transform (2.4) in the x_2 and x_3 directions (using the respective Fourier coordinates ω_2 and ω_3) to obtain

$$(3.2) \quad \mathcal{L}\hat{u} = \left((1 - \lambda)J^{-2/3} - \omega_2^2 - \omega_3^2 \right) \hat{u},$$

where assuming $J > 0$ we have applied the transformation

$$x \rightarrow J^{1/3}x.$$

(We confine the discussion in what follows to the case $J > 0$. If $J < 0$, we can consider the complex conjugate of (3.2) to obtain a new problem with $J > 0$.)

By Lemma 2.1 we have that $\hat{u} \equiv 0$ is the unique L^2 solution of (3.2). However, for the blow-up arguments employed in the next section we need to obtain the above result for any uniformly bounded solution of (3.1) in \mathbb{R}^3 . This is exactly what the next lemma states.

LEMMA 3.1. *Let u denote a uniformly bounded solution of (3.1) in \mathbb{R}^3 . Then, $u \equiv 0$.*

Proof. We first show that $u(\cdot, x_2, x_3) \in L^2(\mathbb{R}, \mathbb{C})$. Let $\chi_r \in C^\infty(\mathbb{R}_+, [0, 1])$ satisfy

$$(3.3) \quad \chi_r(x) = \begin{cases} 1 & x < r/2, \\ 0 & x > r, \end{cases} \quad |\chi'| \leq \frac{C}{r}.$$

Multiplying (3.1) by $\chi_r^2(|x - x_0|)\bar{u}$ we obtain, taking the real part of identity, that

$$(3.4) \quad \int_{B(x_0, r)} |\nabla(\chi_r u)|^2 \leq \int_{B(x_0, r)} [\chi^2 + |\nabla\chi|^2] |u|^2.$$

Consequently, since u is bounded in $L^\infty(\mathbb{R}^3)$, we have

$$(3.5) \quad \int_{B(x_0, r/2)} |\nabla u|^2 \leq C \quad \forall x_0 \in \mathbb{R}^3.$$

From the imaginary part of the identity, we obtain that

$$\int_{B(x_0, r)} \left(\nabla(\chi_r^2) \cdot \Im(\bar{u} \nabla u) + J_{x_1} \chi_r |u|^2 \right) = 0.$$

Let $x_0 = (x_0^1, x_0^2, x_0^3)$. Then,

$$\int_{B(x_0, r/2)} |x_1| |u|^2 dx_1 \leq C \int_{B(x_0, r)} [|u|^2 + |\nabla u|^2].$$

Consequently, for $|x_0^1| > r$, since u is bounded and in view of (3.5), we have

$$\int_{B(x_0, r/2)} |u|^2 \leq \frac{C}{|x_0^1| - \frac{r}{2}} \leq \frac{C}{|x_0^1|}.$$

Repeating the above steps (from (3.4) to the above inequality) k times we obtain that

$$\int_{B(x_0, r/2^k)} |u|^2 \leq \frac{C_k}{|x_0^1|^k}.$$

Hence,

$$\int_{B(x_0, r/2^k)} |x_1|^2 |u|^2 \leq \frac{C_k}{|x_0^1|^{k-2}},$$

which allows us to apply standard elliptic estimates [3] to obtain that

$$|u| \leq \frac{C_k}{(|x_1| + 1)^k} \quad \forall k \in \mathbb{N}.$$

In view of the above we have that $x_1^k u(x_1, x') \in L^2(\mathbb{R}, \mathbb{C})$ for all fixed $x' \in \mathbb{R}^2$. Thus, one can apply to (3.1) the Fourier transform (2.4) to obtain

$$(3.6) \quad -\Delta_\perp \hat{u} + (\omega^2 - 1 + \lambda)\hat{u} - J \frac{\partial \hat{u}}{\partial \omega} = 0,$$

where

$$-\Delta_{\perp} = \frac{\partial^2}{\partial x_2^2} + \frac{\partial^2}{\partial x_3^2}.$$

Let

$$(3.7) \quad \tilde{x}_r(x) = \begin{cases} 1, & |x| < r, \\ e^{-\frac{1}{r}(|x|-r)}, & |x| > r. \end{cases}$$

Multiplying (3.6) by $\tilde{\chi}_r^2(|x' - x^0|)\bar{u}$ and integrating over \mathbb{R}^2 , we obtain

$$-J \frac{dU_r}{d\omega} + (\omega^2 - 1 + \lambda)U_r = \int_{\mathbb{R}^2} (-|\nabla(\tilde{\chi}_r \hat{u})|^2 + |\nabla \tilde{\chi}_r|^2 |\hat{u}|^2) dx',$$

where

$$U_r(\omega) = \int_{\mathbb{R}^2} \tilde{\chi}_r^2 |\hat{u}|^2 dx'.$$

For the last term, we have

$$\int_{\mathbb{R}^2} |\nabla \tilde{\chi}_r|^2 |\hat{u}|^2 \leq \frac{1}{r^2} U_r.$$

Consequently,

$$-J \frac{dU_r}{d\omega} + (\omega^2 - 1 + \lambda - r^{-2})U_r \leq 0,$$

and therefore, for every $\omega_0 \in \mathbb{R}$ and $\omega > \omega_0$, we have

$$U_r(\omega) \geq U_r(\omega_0) \exp \left\{ \frac{1}{3}(\omega^3 - \omega_0^3) - (1 + \lambda + r^{-2})(\omega - \omega_0) \right\}.$$

Thus, since U_r is positive, it must diverge exponentially fast, unless

$$U_r \equiv 0,$$

from which the lemma easily follows. \square

3.1.1. Eigenfunctions in \mathbb{R}_+^3 : Perpendicular current. Let

$$\mathbb{R}_+^3 = \{(x_1, x_2, x_3) \mid x_1 > 0\}.$$

We consider here solutions of (3.1) in \mathbb{R}_+^3 satisfying a Dirichlet boundary condition on $\partial\mathbb{R}_+^3$. Instead of considering complex eigenvalues, we consider only real ones and treat their imaginary part as part of the electric potential.

LEMMA 3.2. *Let $u \in H^2(\mathbb{R}_+^3)$ denote a uniformly bounded solution of*

$$(3.8) \quad \begin{cases} -\Delta u - u + iJ(x_1 - \mu)u = -\lambda u & \text{in } \mathbb{R}_+^3, \\ u = 0 & \text{on } \partial\mathbb{R}_+^3, \end{cases}$$

with $\lambda \in \mathbb{R}_+$. Then, if $J > J_c$, where J_c is defined as in (2.19), u must vanish identically.

Proof. Let $u_n(x_1)$ be defined as in (2.12). Let further

$$a_n(x_2, x_3) = \int_0^\infty u_n(x_1)u(x_1, x_2, x_3)dx_1 = \langle u_n, \bar{u} \rangle,$$

where the inner product is in the regular L^2 sense. Clearly, a_n is uniformly bounded in \mathbb{R}^2 as $u \in L^\infty(\mathbb{R}^3, \mathbb{C})$ and $u_n \in L^1(\mathbb{R}, \mathbb{C})$.

Applying the transformation

$$(3.9) \quad x \rightarrow J^{1/3}x,$$

multiplying (3.8) by u_n , and integrating over \mathbb{R}_+ with respect to x_1 we obtain, in view of the boundedness of u and the exponential rate of decay of u_n as $x_1 \rightarrow \infty$, that

$$(3.10) \quad -\Delta_\perp a_n + (\lambda_n - \tilde{\lambda}_J - i\mu)a_n = 0,$$

where $\tilde{\lambda}_J = (1-\lambda)J^{-2/3}$ and the definition of λ_n is given as in Lemma 2.3. Multiplying (3.10) by $\tilde{\chi}_r(|x'|)$, we obtain for the real part

$$\left(\frac{|\mu_n|}{2} - \tilde{\lambda}_J\right) \int_{\mathbb{R}^2} |\tilde{\chi}_r|^2 |a_n|^2 = - \int_{\mathbb{R}^2} |\nabla(\tilde{\chi}_r a_n)|^2 + \int_{\mathbb{R}^2} |\nabla \tilde{\chi}_r|^2 |a_n|^2 \leq \frac{C}{r^2} \int_{\mathbb{R}^2} |\tilde{\chi}_r|^2 |a_n|^2.$$

Since $\tilde{\lambda}_J < |\mu_1|/2 \leq |\mu_n|/2$ by our assumption we obtain that for sufficiently large r we must have

$$\int_{\mathbb{R}^2} |\tilde{\chi}_r|^2 |a_n|^2 = 0.$$

Hence, $a_n \equiv 0$ in \mathbb{R}^2 . Since by Lemma 2.3 $\{u_n\}_{n=1}^\infty$ is a basis for $L^2(\mathbb{R}_+, \mathbb{C})$, we must have $u \equiv 0$. \square

3.1.2. Steady solutions in \mathbb{R}_+^3 : Nonperpendicular current. This problem is very similar to the problem in \mathbb{R}^3 . Consider the equation

$$(3.11) \quad \begin{cases} -\Delta u - u + i(J_1 x_1 + J_2 x_2 - \mu)u = -\lambda u & \text{in } \mathbb{R}_+^3, \\ u = 0 & \text{on } \partial\mathbb{R}_+^3, \end{cases}$$

with $J_2 \neq 0$ and $\lambda \in \mathbb{R}_+$. Like the problem in \mathbb{R}^3 , there is no need to consider $\mu \neq 0$ here since the transformation

$$x_2 \rightarrow x_2 + \frac{\mu}{J_2}$$

sets $\mu = 0$ in the transformed problem. Furthermore, we also obtain the following result, which is exactly the same as the result obtained in \mathbb{R}^3 .

LEMMA 3.3. *Let u denote a bounded solution of (3.11) with $J_2 \neq 0$. Then $u \equiv 0$.*

Proof. Consider first the case where $J_1 \neq 0$. We first apply the transformation (3.9) with $J = J_1$ to obtain

$$(3.12) \quad -\Delta_\perp u + \mathcal{L}u - \tilde{\lambda}_{J_1} u + i\gamma x_2 u = 0,$$

where $\tilde{\lambda}_{J_1} = (1-\lambda)J_1^{-2/3}$ and $\gamma = J_2/J_1$. Multiplying (3.12) by u_n and integrating over \mathbb{R}_+ with respect to x_1 , we obtain

$$-\Delta_\perp a_n + (\lambda_n - \tilde{\lambda}_{J_1})a_n + i\gamma x_2 a_n = 0 \quad \text{in } \mathbb{R}^2,$$

where $a_n = \langle u_n, \tilde{u} \rangle$. The above equation cannot have any nontrivial bounded solution in \mathbb{R}^2 ; otherwise it would also be a bounded solution in \mathbb{R}^3 , which by Lemma 3.1 must identically vanish. Consequently, since a_n must be bounded, we must have $a_n \equiv 0$ for all n , from which the lemma easily follows in this case.

Consider now the case where $J_1 = 0$. Here we define

$$\tilde{u}(x_1, x_2, x_3) = \begin{cases} u(x_1, x_2, x_3), & x_1 > 0, \\ -u(-x_1, x_2, x_3), & x_1 < 0. \end{cases}$$

Clearly, \tilde{u} is a bounded weak solution of (3.1) in \mathbb{R}^3 and hence, by Lemma 3.1, $\tilde{u} \equiv 0$. \square

For later reference we shall also need the following lemma.

LEMMA 3.4. *Let u denote a bounded solution of*

$$\begin{cases} -\Delta u - u + iJ_2x_2u = -\lambda u & \text{in } \mathbb{R}_+^3, \\ \frac{\partial u}{\partial x_1} = 0 & \text{on } \partial\mathbb{R}_+^3, \end{cases}$$

with $\lambda \in \mathbb{R}_+$. Then $u \equiv 0$.

Proof. Once again we define, this time an even function,

$$\tilde{u}(x_1, x_2, x_3) = \begin{cases} u(x_1, x_2, x_3), & x_1 > 0, \\ u(-x_1, x_2, x_3), & x_1 < 0. \end{cases}$$

Clearly, \tilde{u} is a bounded weak solution of (3.1) in \mathbb{R}^3 and hence, by Lemma 3.1, $\tilde{u} \equiv 0$. \square

The last result in this section is needed in the next section in order to deal with the interface between $\partial\Omega^i$ and $\partial\Omega_n$ (that are perpendicular by assumption).

LEMMA 3.5. *Let*

$$Q = \{(x_1, x_2, x_3) \in \mathbb{R}_+^3 \mid x_2 > 0\}.$$

Let u denote a bounded solution of

$$\begin{cases} -\Delta u - u + i(J_2x_2 + J_3x_3)u = -\lambda u & \text{in } Q, \\ \frac{\partial u}{\partial x_1}(0, x_2, x_3) = 0, & x_2 > 0, \quad x_3 \in \mathbb{R}, \\ u(x_1, 0, x_3) = 0, & x_1 > 0, \quad x_3 \in \mathbb{R}, \end{cases}$$

with $\lambda \in \mathbb{R}_+$. Then $u \equiv 0$.

Proof. Once again we define an even extension of u ,

$$\tilde{u}(x_1, x_2, x_3) = \begin{cases} u(x_1, x_2, x_3), & x_1 > 0, \\ u(-x_1, x_2, x_3), & x_1 < 0. \end{cases}$$

By Lemma 3.3, $\tilde{u} = 0$. \square

4. Large bounded domains in \mathbb{R}^3 . We consider here (1.6) in the limit $R \rightarrow \infty$ which is the large domain limit. We first show that any eigenfunctions of the elliptic operator in (1.6) must decay exponentially fast, as $R \rightarrow \infty$ away from the boundary. As in the previous section we insert the imaginary part of λ into the electric potential

(and consequently consider a family of potentials) and then confine the discussion to real values of λ .

Before getting into the main discussion we repeat here the definition of ϕ_R from section 1 and list some of its properties. Recall that ϕ is the unique solution of (1.2), and that

$$\phi_R(x) = R\phi(x/R).$$

The following proposition lists some of the properties of ϕ and ϕ_R .

PROPOSITION 4.1. *Let ϕ denote a solution of (1.2) and $\phi_R(x) = R\phi(x/R)$. Let $\mu \in \mathbb{R}$ and $\mu_R = R\mu$. Let further $x_j \in \Omega_{R_j}$ where $R_j \uparrow \infty$, and $\{\mu_j\}_{j=1}^\infty \subset \mathbb{R}$. Then*

- (i) *we have either $|\phi_{R_j}(x_j) - \mu_j|$ is unbounded, or else we must have, up to a subsequence,*

$$\exists(b, J) \in \mathbb{R} \times \mathbb{R}^3 : \|\phi_{R_j} - \mu_j - J \cdot x - b\|_{L^\infty(D_r(x_j))} \rightarrow 0 \quad \forall r > 0,$$

where $D_r(x_j) = \Omega_R \cap B(x_j, r)$.

- (ii) *we let Γ_μ denote the level set $\phi = \mu$. Let $M = \max_{x \in \partial\Omega_c} \phi(x)$ and $m = \min_{x \in \partial\Omega_c} \phi(x)$. If $\mu \notin [m, M]$, then Γ_μ is empty.*
- (iii) *assume that $\partial\Omega_c$ is composed of exactly two connected sets as in Figure 1, and that Ω is diffeomorphic to a cylinder. If $\mu \in [m, M]$, but $\mu \notin \phi(\partial\Omega_c)$, then $\Gamma_\mu \cap \partial\Omega_i$ is a simple closed contour, separating $\partial\Omega$ into two subsets, such that none of them is a subset of $\partial\Omega_i$. Furthermore, $\nabla\phi \neq 0$ on Γ_μ .*

Proof. Let $b_j = \phi_{R_j}(x_j) - \mu_j$ and $J_j = \nabla\phi_{R_j}(x_j)$. To prove (i) we first choose a subsequence such that $(b_j, J_j) \rightarrow (b, J)$. The claim then follows from the Taylor expansion of ϕ_{R_j} near x_j .

The proof of (ii) follows immediately from the maximum principle.

To prove (iii) we notice first that since ϕ is real analytic, Γ_μ must either be closed or intersect $\partial\Omega$. If it is closed, then $\phi \equiv \mu$ inside Γ_μ , and hence also outside Γ_μ (in view of its analyticity), which is clearly a contradiction. Thus Γ_μ must intersect the boundary on $\partial\Omega_i$.

Since ϕ is continuous and since on one of the connected sets we must have $\phi < \mu$ and on the other one $\phi > \mu$, the intersection of $\partial\Omega_i$ with Γ_μ must contain at least one closed contour. This contour separates $\partial\Omega$ into two disjoint subsets $\partial\Omega_+$ and $\partial\Omega_-$, the first of them contains the connected subset of $\partial\Omega_c$; over which $\phi > \mu$, and contains a portion of $\partial\Omega^i$. Moreover, this contour is the boundary of a continuous subset of Γ_μ which we denote by A .

To see this, define cylindrical coordinates (r, θ, z) in Ω , where $\theta \in [-\pi, \pi)$ and $0 \leq r < R(\theta, z)$. Then, for each (r, θ) there exists a finite (as ϕ is real analytic) set $\{z_j\}$ such that $(r, \theta, z_j) \in \Gamma_\mu$. Denote the minimum in this set by z_m . Clearly, $z = z_m(r, \theta)$ is continuous, and thus we can define

$$A = \{(r, \theta, z) \mid z = z_m(r, \theta)\}.$$

We now consider the problem for ϕ in a subdomain of Ω whose boundary consists of A and $\partial\Omega_+$. By the maximum principle and Hopf's lemma we have $\phi > \mu$ on $\partial\Omega_+$ and in the interior. In a similar manner we show that $\phi < \mu$ both on $\partial\Omega_-$ and in the interior of the subdomain surrounded by A and $\partial\Omega_-$. This shows that $A = \Gamma_\mu$. From Hopf's lemma it follows that $\nabla\phi \neq 0$ on Γ_μ . \square

Remark 4.1. While property (i) will be used extensively throughout this section, properties (ii) and (iii) are brought here to provide the reader with some intuition of the behavior of ϕ in the “wire-like” domain presented in Figure 1.

Since (1.6) contains the term $i(\phi_R - \mu_R)\psi_R$ which might be unbounded as $R \rightarrow \infty$, we must provide here the following elliptic estimate.

LEMMA 4.1. *Let u_R denote a solution of*

$$(4.1) \quad \begin{cases} -\Delta u_R - u_R + i(\phi_R - \mu_R)u_R = -\lambda u_R & \text{in } \Omega_R, \\ \frac{\partial \psi_R}{\partial \nu} = 0, & \text{on } \partial\Omega_R^i, \\ \psi_R = 0, & \text{on } \partial\Omega_R^c, \end{cases}$$

in which $\lambda \in \mathbb{R}_+$. Let $D_r(x_0) = \Omega_R \cap B_r(x_0)$, where $x_0 \in \Omega_R$ is chosen such that either

- (i) $D_r(x_0) = B_r(x_0)$, or
- (ii) $x_0 \in \partial\Omega_R$ and either $(\partial D_r(x_0) \cap \partial\Omega_R) \subset \partial\Omega_R^c$ or $(\partial D_r(x_0) \cap \partial\Omega_R) \subset \partial\Omega_R^i$.

Then,

$$(4.2) \quad \exists \tilde{r} > 0 : \|u_R\|_{L^\infty(D_r(x_0))} \leq C_r \|u_R\|_{L^2(D_{2r}(x_0))} \quad \forall x_0 \in \Omega_R, r < \tilde{r},$$

where C_r is independent of R and x_0 .

Proof. Let $\rho_R = |u_R|$. By (4.1) we have that

$$-\Delta \rho_R - \rho_R \leq 0$$

in Ω_R . Let $x_0 \in \Omega_R$. We set U_r to be the solution (we discuss its existence below) of

$$(4.3) \quad \begin{cases} -\Delta U_r - U_r = 0 & \text{in } D_r(x_0), \\ U_r = 0 & \text{on } \partial D_r^c(x_0), \\ \frac{\partial U_r}{\partial \nu} = 0 & \text{on } \partial D_r^i(x_0), \\ U_r = \rho_R & \text{on } \partial D_r^s(x_0), \end{cases}$$

where $\partial D_r^c(x_0) = \partial D_r(x_0) \cap \partial\Omega_R^c$, $\partial D_r^i(x_0) = \partial D_r(x_0) \cap \partial\Omega_R^i$, and $\partial D_r^s(x_0) = \partial D_r(x_0) \setminus (\partial D_r(x_0) \cap \partial\Omega_R)$. Note that either $\partial D_r^c(x_0) = \emptyset$ or $\partial D_r^i(x_0) = \emptyset$. Clearly, there exists \tilde{r} , independent of x_0 and R , such that for every $r < \tilde{r}$, we have

$$(4.4) \quad \inf_{u \in \mathcal{D}} \frac{\text{int}_{D_r(x_0)} |\nabla u|^2}{\text{int}_{D_r(x_0)} |u|^2} > 1,$$

where

$$\mathcal{D} = \{u \in H^1(D_r) \mid u = 0 \text{ on } \partial D_r(x_0) \setminus \partial D_r^i(x_0)\}.$$

For $r < \tilde{r}$ the elliptic operator in (4.3) is invertible, and hence a unique U_r exists.

Let then $V = \rho_R - U_r$. Clearly,

$$(4.5) \quad \begin{cases} -\Delta V - V \leq 0 & \text{in } D_r(x_0), \\ \frac{\partial V}{\partial \nu} = 0 & \text{on } \partial D_r^i(x_0), \\ V = 0 & \text{on } \partial D_r(x_0) \setminus \partial D_r^i(x_0). \end{cases}$$

Let further

$$V_+ = \begin{cases} V, & V \geq 0, \\ 0, & V < 0. \end{cases}$$

Multiplying (4.5) by V_+ and integrating over $D_r(x_0)$, we obtain that

$$\int_{D_r} |\nabla V_+|^2 - |V_+|^2 \leq 0,$$

and since $V \in \mathcal{D}$, we obtain by (4.4) that $V_+ = 0$. Consequently we have

$$(4.6) \quad \rho_R \leq U_r \quad \text{in } D_r(x_0).$$

To complete the proof of (4.2) it is thus necessary to obtain an estimate of $\|U_r\|_{L^\infty(D_r)}$ in terms of $\|\rho_R\|_{L^2(D_r)}$. In case (i) above, since U_r is unique, we can use Theorem 10.5 in [3], together with Sobolev embeddings, to obtain

$$(4.7) \quad \|U_r\|_{L^\infty(D_r(x_0))} \leq C_r \|\rho_R\|_{H^{1/2}(\partial D_r^s)},$$

where C_r is independent of x_0 and R .

In case (ii) $D_r(x_0)$ is diffeomorphic to a hemisphere with radius r . Denote the diffeomorphism by T_{R,x_0} . Note that $T_{R,x_0} \rightarrow I$ as $R \rightarrow \infty$ uniformly in x_0 (as long as the assumption in (ii) holds) in view of the transformation (1.5). Let $B_+ = T_{R,x_0}(D_r)$. On the flat surface of B_+ we have either $\partial U_r / \partial \nu = 0$ or $U_r = 0$. In the first case one can extend U_r evenly to a sphere B , whereas in the second case we use an odd extension for that matter. In both cases U_r satisfies

$$A : \nabla(A^t \nabla U_r) + U_r = 0 \quad \text{in } B,$$

where $A = DT_{R,x_0}$ is smooth in B and satisfies $A \rightarrow I$ uniformly in B as $R \rightarrow \infty$. On ∂B , U_r is equal to either the even extension or the odd extension of ρ_R . Hence, for sufficiently large R , U_r must satisfy (4.7) in case (ii) as well.

Clearly,

$$(4.8) \quad \|\rho_R\|_{H^{1/2}(\partial D_r^s)} \leq \|\rho_R\|_{H^1(D_r(x_0))}.$$

Multiplying (4.1) by $\chi_{2r}^2(|x - x_0|)\bar{u}_R$ and integrating over Ω_R , we obtain from the real part of the identity

$$\int_{D_{2r}(x_0)} |\nabla(\chi_{2r} u_R)|^2 - (|\nabla \chi_{2r}|^2 + (1 - \lambda)\chi_{2r}^2)|u_R|^2 = 0.$$

Hence,

$$(4.9) \quad \int_{D_r(x_0)} |\nabla \rho_R|^2 \leq C_r \int_{D_{2r}(x_0)} |\rho_R|^2.$$

Combining the above with (4.6)–(4.8) yields (4.2). \square

As an immediate conclusion of Lemma 4.1 we prove the following lemma

LEMMA 4.2. *Let u_R denote a solution of (4.1) with $\lambda \in \mathbb{R}_+$. Let X_R denote a maximum point of $|u_R|$ in Ω_R . Then, $|\phi_R(x_R) - \mu_R|$ is bounded as $R \rightarrow \infty$.*

Proof. We first normalize u_R by $|u_R(x_R)|$ so that $\|u_R\|_\infty = 1$. Let $D_r(x_R) = B(x_R, r) \cap \Omega_R$. Multiplying (4.1) by $\chi_r^2(|x - x_R|)\bar{u}_R$ and integrating over Ω_R , we obtain from the imaginary part of the identity

$$\int_{D_r} \nabla(\chi_r^2) \cdot \frac{1}{2i} (\bar{u}_{R,j} \nabla u_R - \nabla \bar{u}_R) + \int_{D_r} \chi_r^2 (\phi_R - \mu_R) |u_R|^2 = 0.$$

Let $b_R = \phi_R(x_R) - \mu_R$. Since $\nabla\phi_R$ is bounded in Ω_R , and since r is fixed, we have

$$\inf_{x \in D_r} |\phi_R - \mu_R| \geq \frac{1}{2} b_R.$$

Thus,

$$b_R \int_{D_{r/2}} |u_R|^2 \leq C_r \int_{D_r} |u_R|^2 + |\nabla u_R|^2.$$

Using (4.9) and the fact that $|u_R| \leq 1$, we obtain that

$$\int_{D_{r/2}} |u_R|^2 \leq \frac{C}{b_R}.$$

By Lemma 4.1 we then have that

$$1 = |u_R(x_R)| \leq \frac{C_r}{b_R},$$

from which the lemma immediately follows. \square

Remark 4.2. Note that if $\phi_R \neq \mu_R$ for all $x \in \Omega_R$, then $u_R \equiv 0$ must be the unique solution of (4.1). To see this multiply (4.1) by \bar{u}_R and integrate over Ω_R to obtain from the imaginary part,

$$\int_{\Omega_R} (\phi_R - \mu_R) |u_R|^2 = 0.$$

Since $\phi_R - \mu_R$ is either positive or negative throughout Ω_R , u_R must vanish everywhere.

Denote the curve in Ω along which we have $\phi = \mu$ by Γ_μ . Denote its image under the mapping (1.5) by Γ_R . By the previous lemma we have that $d(x_R, \Gamma_R)$ is bounded as $R \rightarrow \infty$. We now prove that u_R must decay exponentially fast away from Γ_R .

LEMMA 4.3. *Let u_R denote a solution of (4.1) with $\lambda \in \mathbb{R}_+$. Then, there exists $\alpha > 0$ such that*

$$(4.10) \quad \int_{\Omega_R} |u_R|^2 e^{2\alpha s} \leq C,$$

where $s = d(x, \Gamma_R)$ and C is independent of R .

Recall that by Proposition 4.1, for domains that are diffeomorphic to a cylinder, when $\phi \neq \mu$ for every $x \in \partial\Omega^c$, Γ_μ must be a surface whose boundary is a closed simple contour on $\partial\Omega^i$. We also have $\nabla\phi \neq 0$ on Γ_μ . Note also that by the above remark, if Γ_R is empty, then $u_R \equiv 0$.

Proof. It is convenient to consider here u_R for which $\|u_R\|_{L^2(\Omega_R)} = 1$. Let Ω_β^+ and Ω_β^- be, respectively, defined by

$$\Omega_\beta^+ = \{x \in \Omega_R \mid \phi_R - \mu_R > \beta\},$$

$$\Omega_\beta^- = \{x \in \Omega_R \mid \phi_R - \mu_R < -\beta\}$$

for some $\beta > 0$ which is independent of R .

Let $\eta_\beta^+ \in C^\infty(\Omega_R, [0, 1])$ and η_β^- , respectively, be defined by

$$\eta_\beta^+ = \begin{cases} 1, & x \in \Omega_\beta^+, \\ 0, & x \in \Omega_0^-, \end{cases}$$

and

$$\eta_\beta^- = \begin{cases} 1, & x \in \Omega_\beta^-, \\ 0, & x \in \Omega_0^+. \end{cases}$$

Let \mathcal{C}_β^\pm denote the portion of $\partial\Omega_\beta^\pm$ which is not on $\partial\Omega_R$. As $\nabla\phi_R(x) = \nabla\phi(x/R)$ and since $\nabla\phi$ is bounded in Ω , we have that

$$d(\mathcal{C}_\beta^\pm, \Gamma_\mu) \geq \frac{\beta}{\|\nabla\phi\|_{L^\infty(\Omega)}}.$$

Hence, we can choose η_β^\pm such that $|\nabla\eta_\beta^\pm| < C$. Let

$$D_\beta^+ = \{x \in \Omega_R \mid 0 < \eta_\beta^+ < 1\}.$$

We choose η_β^\pm such that

$$\sup_{x \in D_\beta^+} s \leq 1.$$

Multiplying (4.1) by $(\eta_\beta^+)^2 e^{2\alpha s} \bar{u}_R$ and integrating over Ω_R we obtain, for the imaginary part and the real part, respectively,

(4.11a)

$$\int_{\Omega_R} \nabla((\eta_\beta^+)^2 e^{2\alpha s}) \cdot \frac{1}{2i} (\bar{u}_R \nabla u_R - u_R \nabla \bar{u}_R) + \int_{\Omega_R} (\phi_R - \mu_R) (\eta_\beta^+)^2 |u_R|^2 e^{2\alpha s} = 0,$$

(4.11b)

$$\int_{\Omega_R} (\eta_\beta^+)^2 |\nabla u_R|^2 e^{2\alpha s} = (1 - \lambda) \int_{\Omega_R} (\eta_\beta^+)^2 |u_R|^2 e^{2\alpha s} - \frac{1}{2} \int_{\Omega_R} \nabla((\eta_\beta^+)^2 e^{2\alpha s}) \cdot \nabla |u_R|^2.$$

From the real part, (4.11b), we obtain that for every $\epsilon > 0$ we have

(4.12)

$$(1 - 2\alpha\epsilon) \int_{\Omega_R} (\eta_\beta^+)^2 e^{2\alpha s} |\nabla u_R|^2 \leq \left(1 + \frac{\alpha}{2\epsilon}\right) \int_{\Omega_R} (\eta_\beta^+)^2 e^{2\alpha s} |u_R|^2 + C \int_{D_\beta^+} e^{2\alpha s} |\nabla u_R|^2.$$

From (4.11a), or the imaginary part, we obtain

$$\beta \int_{\Omega_\beta^+} |u_R|^2 e^{2\alpha s} \leq \alpha \int_{\Omega_R} (\eta_\beta^+)^2 e^{2\alpha s} [|u_R|^2 + |\nabla u_R|^2] + C \int_{D_\beta^+} e^{2\alpha s} [|u_R|^2 + |\nabla u_R|^2].$$

Combining the above with (4.12) for $\epsilon = 4\alpha^{-1}$, we obtain

$$(4.13) \quad (\beta - 2\alpha - 4\alpha^3) \int_{\Omega_\beta^+} |u_R|^2 e^{2\alpha s} \leq C \int_{D_\beta^+} e^{2\alpha s} [|u_R|^2 + |\nabla u_R|^2].$$

Multiplying (4.1) by \bar{u}_R and integrating over Ω_R , we obtain (for the real part)

$$\int_{\Omega_R} |\nabla u_R|^2 = (1 - \lambda) \int_{\Omega_R} |u_R|^2 = 1 - \lambda.$$

Consequently, we obtain from (4.13) that for any given α we may choose β to be sufficiently large (but still independent of R) so that $\beta > 2\alpha + 4\alpha^3$, and hence

$$\int_{\Omega_\beta^+} |u_R|^2 e^{2\alpha s} \leq C e^{2\alpha} \int_{D_\beta^+} [|u_R|^2 + |\nabla u_R|^2] \leq 2C e^{2\alpha}. \quad \square$$

Remark 4.3. By Lemma 4.1 it follows that $|u_R|$ decays exponentially fast away from Γ_R also in a pointwise sense.

We now prove that any eigenfunction corresponding to a nonpositive eigenvalue must decay exponentially fast away from the boundary.

LEMMA 4.4. *Let u_R denote a solution of (4.1). Then*

$$(4.14) \quad \exists \alpha > 0 : |u_R| \leq C e^{-\alpha d(x, \partial\Omega)} \quad \forall \mu_R \in \mathbb{R},$$

where α is independent of R , μ_R , and λ . Furthermore, denote by x_R the maximum point of $|u_R|$. Then $d(x_R, \partial\Omega_R)$ is bounded as $R \rightarrow \infty$.

Proof. We apply standard blow-up arguments to prove the lemma. Let

$$\Omega(R, k, s) = \{x \in \Omega_R \mid d(x, \partial\Omega_R) \geq ks\}.$$

We prove the exponential rate of decay by showing first that

$$(4.15) \quad \exists R_0, s_0 : \|u_R\|_{L^\infty(\Omega(R, k+1, s))} \leq \frac{1}{2} \|u_R\|_{L^\infty(\Omega(R, k, s))} \quad \forall s > s_0, R > R_0, k \in \mathbb{N}.$$

Suppose, for contradiction, that (4.15) does not hold. Then, there exist sequences $\{R_j\}_{j=1}^\infty$, $\{s_j\}_{j=1}^\infty$, and $\{k_j\}_{j=1}^\infty$ satisfying $R_j \uparrow \infty$, $s_j \uparrow \infty$, $k_j \in \mathbb{N}$, and

$$(4.16) \quad \|u_{R_j}\|_{L^\infty(\Omega(R_j, k_j+1, s_j))} \geq \frac{1}{2} \|u_{R_j}\|_{L^\infty(\Omega(R_j, k_j, s_j))} \stackrel{def}{=} \frac{1}{2} m_j.$$

Let

$$\tilde{u}_{R_j} = \frac{u_{R_j}}{m_j}.$$

By (4.15) there exists $x_j \in \Omega(R_j, k_j + 1, s_j)$ such that

$$(4.17) \quad |\tilde{u}_{R_j}(x_j)| \geq \frac{1}{2}.$$

For notational convenience we also let $f_j(x) = \tilde{u}_{R_j}(x_j + x)$.

We now distinguish between two different cases.

Case 1.

$$b_j = \inf_{x \in B(x_j, s_j)} |\phi_{R_j} - \mu_{R_j}| \rightarrow \infty$$

up to a subsequence.

Let $\chi_r \in C^\infty(\mathbb{R}_+, [0, 1])$ be defined by (3.3). Since f_j satisfies (4.1) we multiply it by $\chi_r^2(0)\bar{f}_j$ and integrate over $B(0, r)$ to obtain, for the imaginary part,

$$\int_{B(0,r)} \nabla(\chi_r^2) \cdot \frac{1}{2i}(\bar{f}_j \nabla f_j - f_j \nabla \bar{f}_j) + \int_{B(0,r)} \chi_r^2(\phi_{R_j} - \mu_{R_j})|f_j|^2 = 0,$$

yielding

$$(4.18) \quad b_j \int_{B(0,r/2)} |f_j|^2 \leq C_r \int_{B(0,r)} |f_j|^2 + |\nabla f_j|^2$$

for all $r < s_j$. For the real part we obtain that

$$\int_{B(0,r)} |\nabla(\chi_r f_j)|^2 - (|\nabla \chi_r|^2 + (1 - \lambda)\chi_r^2)|f_j|^2 = 0,$$

and hence

$$(4.19) \quad \int_{B(0,r/2)} |\nabla f_j|^2 \leq C \int_{B(0,r)} |f_j|^2.$$

Combining (4.18) with (4.19), we obtain

$$\int_{B(0,r/2)} |f_j|^2 \leq \frac{C_r}{b_j} \int_{B(0,2r)} |f_j|^2.$$

As $|f_j| \leq 1$ we obtain that

$$\int_{B(0,r/2)} |f_j|^2 \rightarrow 0$$

as $j \rightarrow \infty$, and by Lemma 4.1 also that $f_j(0) \rightarrow 0$, a contradiction.

Case 2. $\limsup_{j \rightarrow \infty} b_j < \infty$.

Let $J = |\nabla \phi|$. We choose a coordinate system, where $\nabla \phi(x_j)$ is parallel to the x_1 axis. Then, by Lemma 4.2 we have a subsequence for which

$$\phi_{R_j} - \mu_{R_j} \rightarrow Jx_1 + b$$

uniformly in $B(0, r)$ for all $r > 0$, where b is a constant. Thus, by standard elliptic estimates and Sobolev embeddings, there exists a subsequence $\{f_{j_k}\}_{k=1}^\infty$ such that $f_{j_k} \rightarrow f_\infty$ uniformly on every compact set in \mathbb{R}^3 and such that f_∞ is a bounded solution of

$$-\Delta f_\infty - f_\infty + i(Jx_1 + b)f_\infty = -\lambda f_\infty \quad \text{in } \mathbb{R}^3.$$

By Lemma 3.1 we must have $f_\infty \equiv 0$, a contradiction.

Thus, we have proved (4.15), and hence also (4.14). The boundedness of $d(x_R, \partial\Omega_R)$ follows from (4.14) as well. \square

The following lemma provides the basis for our main stability result.

LEMMA 4.5. *Let $\partial\Omega_n$ denote the subset of $\partial\Omega_c$, where $\nabla \phi$ is perpendicular to $\partial\Omega$. Suppose that either $|\nabla \phi| > J_c$ for all $x \in \partial\Omega_n$ or that $\partial\Omega_n$ is empty. Then, for sufficiently large R , $u_R \equiv 0$ is the unique solution of (4.1) for all $\mu_R \in \mathbb{R}$.*

Proof. Let x_R be the point where u_R obtains its maximum in Ω_R . Let x_0^R denote its projection on $\partial\Omega_R$. Recall that by Lemma 4.4 $|x_R - x_0^R|$ is bounded as $R \rightarrow \infty$. Note that by Lemma 4.2 ϕ_R converges uniformly to a linear function in $D_r(x_0^R)$ for all fixed $r > 0$ as $R \rightarrow \infty$. Suppose first that $x_0^R \in \partial\Omega_i^R$. Following [11], let (t_1, t_2, t_3) denote a local curvilinear coordinate system, whose origin lies at x_0^R , such that $t_3 = d(x, \partial\Omega_R)$ when $x \in \Omega_R$ and such that the t_1 and t_2 curves on $\partial\Omega$ are the lines of curvature. Let further κ_1^R and κ_2^R denote the respective principal curvatures on $\partial\Omega_R$. Clearly, $\kappa_i^R = \kappa_i/R$ ($i = 1, 2$), where κ_i is the corresponding principal curvature on $\partial\Omega$, at $x_0 = x_0^R/R$.

Since $\partial\Omega$ is smooth near x_0^R , this curvilinear coordinate system is properly defined in some neighborhood of x_0^R . Let

$$B^+(0, r) = \{(t_1, t_2, t_3) \in B(0, r) \mid t_3 > 0\}.$$

Then, the above coordinate system is well defined in $B^+(0, \delta R)$ for some $\delta > 0$. We can now present any x in this neighborhood by

$$x = r(t_1, t_2) - t_3\nu,$$

where ν is the outward normal at $(t_1, t_2, 0)$. Let

$$g_{ij}(t_1, t_2) = \frac{\partial r}{\partial t_i} \cdot \frac{\partial r}{\partial t_j}, \quad i, j = 1, 2,$$

and

$$G_{ij} = [1 - \kappa_i t_3/R]g_{ij}, \quad i, j = 1, 2.$$

Since our coordinate system is orthogonal, we have

$$g_{12} = 0.$$

Furthermore, we can scale t_1 and t_2 so that

$$g_{11} = g_{22} = 1 + O(1/R) \quad \text{as } R \rightarrow \infty$$

uniformly in $B^+(0, r)$ for every fixed $r > 0$. Finally, we define

$$\begin{cases} G = \sqrt{G_{11}G_{22}}, \\ \alpha_j = \frac{G}{G_{jj}}, \quad j = 1, 2. \\ \alpha_3 = G, \end{cases}$$

Let $w_R(x) = u_R(x)/|u_R(x_R)|$. In the new coordinates, (4.1) takes the form

$$-\sum_{j=1}^3 \frac{1}{G} \frac{\partial}{\partial t_j} \left(\alpha_j \frac{\partial w_R}{\partial t_j} \right) - w_R + i(\phi_R - \mu_R)w_R = -\lambda w_R$$

in $B^+(0, \delta R)$. Standard elliptic estimates then prove the existence of a sequence $\{w_{R_j}\}_{j=1}^\infty$ such that $w_{R_j} \rightarrow w_\infty$ uniformly on every compact set in \mathbb{R}_+^3 , where here

$$\mathbb{R}_+^3 = \{(t_1, t_2, t_3) \mid t_3 > 0\}.$$

By standard elliptic estimates again we have that w_∞ satisfies the following problem:

$$\begin{cases} -\Delta w_\infty - w_\infty + i(J_1 t_1 + J_2 t_2 + b)w_\infty = -\lambda w_\infty & \text{in } \mathbb{R}_+^3, \\ \frac{\partial w_\infty}{\partial \nu} = 0 & \text{on } \partial \mathbb{R}_+^3. \end{cases}$$

By Lemma 3.4 we have $w_\infty \equiv 0$ in \mathbb{R}_+^3 , a contradiction since $|w_R(x_R)| = 1$ and $|x_R - x_0^R|$ is bounded.

Consider now the case where $x_0^R \in \partial \Omega_c^R$. Following the same procedure as before we obtain that $w_R \rightarrow w_\infty$ uniformly on every compact set in \mathbb{R}_+^3 , where w_∞ must satisfy

$$\begin{cases} -\Delta w_\infty - w_\infty + i(J_1 t_1 + J_2 t_2 + J_3 t_3 + b)w_\infty = -\lambda w_\infty & \text{in } \mathbb{R}_+^3, \\ w_\infty = 0 & \text{on } \partial \mathbb{R}_+^3. \end{cases}$$

If $\partial \Omega_n$ is empty, we have $J_1^2 + J_2^2 > 0$. By Lemma 3.3 we then have $w_\infty \equiv 0$. Otherwise, if $x_0 \in \partial \Omega_n$, we must have $w_\infty \equiv 0$ by Lemma 3.2 since $J_3 > J_c$.

Finally, if x_0 lies on the interface between $\partial \Omega_i$ and $\partial \Omega_c$, we obtain that w_∞ must satisfy, since the two surfaces are perpendicular to each other at the interface, a problem in Q , where

$$Q = \{(t_1, t_2, t_3) \mid t_3 > 0, t_1 > 0\}.$$

We have

$$\begin{cases} -\Delta w_\infty - w_\infty + i(J_1 t_1 + J_2 t_2 + b)w_\infty = 0 & \text{in } Q, \\ \frac{\partial w_\infty}{\partial \nu} = 0, & x \in \partial Q : t_3 = 0, \\ w_\infty = 0, & x \in \partial Q : t_1 = 0. \end{cases}$$

By Lemma 3.5 we have $w_\infty \equiv 0$. (Note that some modification of the local coordinate system is necessary in this case.) \square

Proof of Theorem 1.1. Since the principal part of the differential operator on the left-hand side of (4.1) is the Laplacian, it can be regarded as a perturbation of a self-adjoint operator. Thus, it follows from the discussion below Theorem 15.2 in [2] regarding such perturbations that it has exactly one direction which is not a direction of minimal growth, that is, $\arg \lambda = 0$. Hence, it follows from Theorem 15.1 in [2] that the spectrum of this differential operator must be discrete and that all its eigenvalues must have finite multiplicity. Furthermore, by Theorem 16.5 in [2], the eigenfunctions span $L^2(\Omega_R, \mathbb{C})$. Hence, for $J > J_c$, since all the eigenvalues of the above operator must have positive real part, the normal state must be stable.

Consider now the case when a point $x \in \partial \Omega_c$ exists, where

$$\left| \frac{\partial \phi}{\partial \nu} \right| = |\nabla \phi| = J < J_c.$$

To prove the short-time instability we look at the solution of (1.6), after applying to it the transformation (2.2), with the initial condition

$$\psi(x, 0) = u_1(t_3)\chi_{R^{1/2}}(t_1, t_2)\eta_R(t_3),$$

where (t_1, t_2, t_3) are the above-defined system of local curvilinear coordinates, χ_r is defined in (3.3), and

$$\eta_R(x) = \begin{cases} 1, & x < \frac{1}{2}\delta R, \\ 0, & x > \delta R, \end{cases}$$

is a smooth cutoff function.

Let $\tilde{\beta} = \lambda_J - \lambda_{J_e}$. We write

$$(4.20) \quad \psi_R = v + \psi_0(x)e^{\tilde{\beta}t}$$

to obtain

$$(4.21) \quad \frac{\partial v}{\partial t} - \Delta v - \lambda_J v = f.$$

The precise form of f need not concern us except for the fact that

$$(4.22) \quad \|f\|_2 \leq \frac{C_\alpha}{R^\alpha} \|\psi_0\|_2 e^{\tilde{\beta}t} \quad \forall \alpha < 1.$$

Multiplying (4.21) by \bar{v} and integrating by parts, we obtain for the real part

$$\begin{cases} \frac{\partial \|v\|_2}{\partial t} - \lambda_J \|v\|_2 = \|f\|_2, \\ \|v\|_2(0) = 0. \end{cases}$$

Consequently,

$$\|v\|_2(t) \leq \int_0^t e^{\lambda(t-\tau)} \|f\|_2(\tau) d\tau,$$

and hence,

$$\|v\|_2 \leq \frac{C_\alpha}{R^\alpha} \|\psi_0\|_2 e^{(\lambda_J + \tilde{\beta})t}.$$

Clearly, there exist $T_R \sim \mathcal{O}(\ln R)$, as $R \rightarrow \infty$, such that

$$t < T_R \Rightarrow \frac{C_\alpha}{R^\alpha} e^{\lambda_J t} < \frac{1}{2},$$

and thus

$$t < T_R \Rightarrow \|\psi\|_2 \geq \frac{1}{2} \|\psi_0\|_2 e^{\tilde{\beta}t}.$$

Applying the inverse of (2.2), we obtain (1.7). \square

The above instability result is valid, of course, only for $T < T_R$. Proving long-time instability appears to be a much more difficult problem. The stability proof presented above relies on the convergence of any solution of (4.1) to a solution of (3.8) uniformly on every compact set near the point on the boundary where J is perpendicular to it. This, however, does not prove convergence of the spectrum or even of its bottom. What has been demonstrated is only the upper semicontinuity of the spectrum of the differential operator on the left-hand side of (4.1). Lower semicontinuity of the spectrum appears to be much harder to prove, especially since the operator is not self-adjoint. Nevertheless, it does seem reasonable to conjecture that the solution would continue to grow exponentially fast as $t \rightarrow \infty$, in view of the above short-time instability result. Further research is necessary in order to establish that result.

5. Concluding remarks. In the previous section we proved that the normal state remains stable in the large domain limit, as long as the current on the boundary, at points where it is perpendicular to it, is greater than J_c . If the current is nowhere perpendicular to the boundary, then as long as it doesn't vanish there, the normal state must be stable. We also demonstrate short-time instability when $J_m < J_c$.

In the following we provide a short list of interesting problems that are waiting to be resolved:

1. Proving long-time instability when $J > J_c$. We have elaborated on this matter at the end of the preceding section.
2. Adding the effect of magnetic fields. This magnetic field can be either induced by the electric current (via (1.3)) or else be applied externally (or both). It has been verified experimentally that an induced magnetic field can generate vortices [14] if the current is sufficiently large and the material is close to the wholly superconducting state. However, its effect on the critical current J_c has not been investigated. It is reasonable to believe that J_c would become smaller if we add the effect of the magnetic field.
3. Adding the effect of temperature, since electric currents have the tendency to heat the sample, thereby creating vortices [14]. Incorporating this effect requires modification of (1.1), and the use of different nondimensionalization; otherwise the domain would become temperature-dependent.
4. Proving that the bifurcating branch (at $J = J_c$) is unstable.

Appendix A. The Hilbert–Schmidt norm of \mathcal{L}_λ^{-1} .

LEMMA A.1. *Let \tilde{G} be as given by (2.14b). Then, for any $\lambda \in \mathbb{C} \setminus \{\lambda_n\}_{n=1}^\infty$, where $\{\lambda_n\}_{n=1}^\infty$ are given as in (2.11),*

$$(A.1) \quad \exists C, M > 0 : \|\tilde{G}\| \leq Ce^{M\lambda^{3/2}}.$$

Proof. We first note that

$$W(w_1, \tilde{w}_2)(x, \lambda) = w_2'(x, \lambda)w_1(x, \lambda) - w_1'(x, \lambda)w_2(x, \lambda)$$

is independent of x by Abel's formula. Furthermore, for $y = x + i\lambda$ we have

$$W(w_1, \tilde{w}_2)(x, \lambda) = W(w_1, \tilde{w}_2)(y, 0) = W(w_1, \tilde{w}_2)(0, 0).$$

Therefore, W is independent of both x and λ .

Since W is constant, it follows that \tilde{G} is symmetric, i.e.,

$$\tilde{G}(x, \xi, \lambda) = \tilde{G}(\xi, x, \lambda).$$

Consequently, it suffices to prove that

$$\int_0^\infty d\xi \int_\xi^\infty dx |\tilde{G}(x, \xi)|^2 = C_1 \int_0^\infty d\xi |\tilde{w}_2(\xi)|^2 \int_\xi^\infty dx |\tilde{w}_1(x)|^2 \leq Ce^{M\lambda^{3/2}},$$

where C, C_1 , and M are all independent of λ .

To obtain the above estimate we use asymptotic properties of Airy's functions [1, 13], from which it follows that

$$(A.2) \quad |A_i(z)| \leq \frac{C}{|z|^{1/4}} \left| e^{-\frac{2}{3}z^{3/2}} \right|.$$

Consider then first the domain $\xi > M_0\lambda$ for sufficiently large $M_0 > 0$. We have

$$\int_{\xi}^{\infty} |w_1(x)|^2 dx \leq \int_{\xi}^{\infty} \frac{dx}{|x + i\lambda|^{1/2}} e^{-\beta(x)|x+i\lambda|^{3/2}},$$

where

$$\beta(x) = \frac{4}{3} \cos\left(\frac{3}{2} \arg(x + i\lambda) + \frac{\pi}{4}\right).$$

It is easy to show that

$$(A.3) \quad \begin{cases} |\beta'(x)| \leq C \frac{|\lambda|}{|x+i\lambda|^2}, \\ |\beta''(x)| \leq C \frac{|\lambda|}{|x+i\lambda|^3}. \end{cases}$$

Let then

$$\gamma(x) = \beta(x)|x + i\lambda|^{3/2}.$$

By (A.3) we have

$$(A.4) \quad \begin{cases} |\gamma'(x)| \geq \frac{3}{2}\beta(x)|x + i\lambda|^{1/2} \left[1 - C \frac{|\lambda|}{|x+i\lambda|}\right], \\ |\gamma''(x)| \leq \frac{3}{4}\beta(x)|x + i\lambda|^{-1/2} \left[1 + C \frac{|\lambda|}{|x+i\lambda|}\right]. \end{cases}$$

Using (A.4) and integration by parts, we obtain

$$\begin{aligned} \int_{\xi}^{\infty} e^{-\gamma(x)} dx &\leq \frac{1}{|\gamma'(\xi)|} e^{-\gamma(\xi)} + \int_{\xi}^{\infty} \frac{|\gamma''(x)|}{|\gamma'(x)|} e^{-\gamma(x)} dx \\ &\leq \frac{C}{|\xi + i\lambda|^{1/2}} e^{-\gamma(\xi)} + \frac{C}{|\xi + i\lambda|^{3/2}} \int_{\xi}^{\infty} e^{-\gamma(x)} dx. \end{aligned}$$

Hence,

$$\int_{\xi}^{\infty} e^{-\gamma(x)} dx \leq \frac{C}{|\xi + i\lambda|^{1/2}} e^{-\gamma(\xi)},$$

from which we easily obtain that

$$(A.5) \quad \int_{\xi}^{\infty} |w_1(x)|^2 dx \leq \frac{C}{|\xi + i\lambda|} e^{-\gamma(\xi)}.$$

From the asymptotic behavior of Airy's function (A.2), we obtain again

$$|\tilde{w}_2(\xi)|^2 \leq \frac{C}{|\xi + i\lambda|^{1/2}} e^{\gamma(\xi)}.$$

Thus,

$$\int_{M_0\lambda}^{\infty} d\xi |\tilde{w}_2(\xi)|^2 \int_{\xi}^{\infty} dx |\tilde{w}_1(x)|^2 \leq \frac{C}{\lambda^{1/2}}.$$

To complete the proof we need to bound the norm for $0 < \xi < M_0\lambda$. By (A.2) and (2.14b) we have

$$|\tilde{G}(x, \xi, \lambda)| \leq C \exp \left\{ \frac{2}{3} (M_0 + 1)^{3/2} |\lambda|^{3/2} \right\}.$$

Consequently, from the above and (A.5) we obtain

$$\begin{aligned} \int_0^{M_0\lambda} d\xi |\tilde{w}_2(\xi)|^2 \int_\xi^\infty dx |\tilde{w}_1(x)|^2 &= \int_0^{M_0\lambda} d\xi |\tilde{w}_2(\xi)|^2 \int_\xi^{M_0\lambda} dx |\tilde{w}_1(x)|^2 \\ &+ \int_0^{M_0\lambda} d\xi |\tilde{w}_2(\xi)|^2 \int_{M_0\lambda}^\infty dx |\tilde{w}_1(x)|^2 \leq CM_0^2 \lambda^2 \exp \left\{ \frac{2}{3} (M_0 + 1)^{3/2} |\lambda|^{3/2} \right\}, \end{aligned}$$

from which (A.1) easily follows. \square

Acknowledgments. The author wishes to thank Professor Bernard Helffer and Professor Yehuda Pinchover for their comments.

REFERENCES

- [1] M. ABRAMOWITZ AND I. A. STEGUN, *Handbook of Mathematical Functions*, Dover, Mineola, NY, 1972.
- [2] S. AGMON, *Lectures on Elliptic Boundary Value Problems*, prepared for publication by B. Frank Jones, Jr., with the assistance of George W. Batten, Jr., Van Nostrand Math. Stud. 2, D. Van Nostrand Co., Inc., Princeton, NJ, 1965.
- [3] S. AGMON, A. DOUGLIS, AND L. NIRENBERG, *Estimates near the boundary for solutions of elliptic partial differential equations satisfying general boundary conditions. II*, Comm. Pure Appl. Math., 17 (1964), pp. 35–92.
- [4] S. J. CHAPMAN AND D. R. HERON, *A hierarchy of models for superconducting thin films*, SIAM J. Appl. Math., 63 (2003), pp. 2087–2127.
- [5] E. B. DAVIES, *Will spectral behaviour of anharmonic oscillators*, Bull. London Math. Soc., 32 (2000), pp. 432–438.
- [6] Q. DU AND P. GRAY, *High-kappa limits of the time-dependent Ginzburg–Landau model*, SIAM J. Appl. Math., 56 (1996), pp. 1060–1093.
- [7] L. P. GOR'KOV AND G. M. ÉLIASHBERG, *Generalisation of the Ginzburg–Landau equations for non-stationary problems in the case of alloys with paramagnetic impurities*, Soviet Phys. JETP, 27 (1968), p. 328.
- [8] B. I. IVLEV AND N. B. KOPNIN, *Electric currents and resistive states in thin superconductors*, Adv. in Phys., 33 (1984), pp. 47–114.
- [9] B. I. IVLEV, N. B. KOPNIN, AND L. A. MASLOVA, *Stability of current-carrying states in narrow finite-length superconducting channels*, Soviet Phys. JETP, 56 (1982), pp. 884–890.
- [10] P. T. LAÏ AND D. ROBERT, *Sur un problème aux valeurs propres non linéaire*, Israel J. Math., 36 (1980), pp. 169–186.
- [11] X.-B. PAN, *Surface superconductivity in 3 dimensions*, Trans. Amer. Math. Soc., 356 (2004), pp. 3899–3937.
- [12] D. ROBERT, *Non-linear eigenvalue problems*, Mat. Contemp., 26 (2004), pp. 109–127.
- [13] Y. SIBUYA, *Global Theory of a Second Order Linear Ordinary Differential Equation with a Polynomial Coefficient*, North–Holland, Amsterdam, 1975.
- [14] A. G. SIVAKOV, A. M. GLUKHOV, A. N. OMELYANCHOUK, Y. KOVAL, P. MÜLLER, AND A. V. USTINOV, *Josephson behavior of phase-slip lines in wide superconducting strips*, Phys. Rev. Lett., 91 (2003), 267001.
- [15] D. Y. VODOLAZOV, F. M. PEETERS, L. PIRAUX, S. MATEFI-TEMPFLI, AND S. MICHOTTE, *Current-voltage characteristics of quasi-one-dimensional superconductors: An s-shaped curve in the constant voltage regime*, Phys. Rev. Lett., 91 (2003), 157001.

CARLEMAN AND OBSERVABILITY ESTIMATES FOR STOCHASTIC WAVE EQUATIONS*

XU ZHANG[†]

Abstract. Based on a fundamental identity for stochastic hyperbolic-like operators, we derive in this paper a global Carleman estimate (with singular weight function) for stochastic wave equations. This leads to an observability estimate for stochastic wave equations with nonsmooth lower order terms. Moreover, the observability constant is estimated by means of an explicit function of the norm of the coefficients involved in the equation. An application to the state observation problem for semilinear stochastic wave equations is also given.

Key words. Carleman estimate, singular weight function, observability estimate, stochastic wave equation, state observation problem

AMS subject classifications. Primary, 60H15; Secondary, 93B07, 35B45

DOI. 10.1137/070685786

1. Introduction. Let $T > 0$, $G \subset \mathbb{R}^n$ ($n \in \mathbb{N}$) be a given bounded domain with a C^2 boundary Γ , with Γ_0 a given nonempty open subset of Γ . Put $Q \triangleq (0, T) \times G$, $\Sigma \triangleq (0, T) \times \Gamma$, and $\Sigma_0 \triangleq (0, T) \times \Gamma_0$. Throughout this paper, we will use C to denote a generic positive constant depending only on T , G , and Γ_0 , which may change from line to line.

Let $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, P)$ be a complete filtered probability space on which a one dimensional standard Brownian motion $\{B(t)\}_{t \geq 0}$ is defined. Let H be a Banach space. We denote by $L^2_{\mathcal{F}}(0, T; H)$ the Banach space consisting of all H -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted processes $X(\cdot)$ such that $\mathbb{E}(|X(\cdot)|^2_{L^2(0, T; H)}) < \infty$, with the canonical norm; by $L^\infty_{\mathcal{F}}(0, T; H)$ the Banach space consisting of all H -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted bounded processes; and by $L^2_{\mathcal{F}}(\Omega; C([0, T]; H))$ the Banach space consisting of all H -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted continuous processes $X(\cdot)$ such that $\mathbb{E}(|X(\cdot)|^2_{C([0, T]; H)}) < \infty$, with the canonical norm (similarly, one can define $L^2_{\mathcal{F}}(\Omega; C^k([0, T]; H))$ for any positive integer k).

Let us consider the following stochastic wave equation:

$$(1.1) \quad \begin{cases} dy_t - \Delta y dt = (a_1 y_t + a_2 \cdot \nabla y + a_3 y + f) dt + (a_4 y + g) dB(t) & \text{in } Q, \\ y = 0 & \text{on } \Sigma, \\ y(0) = y_0, \quad y_t(0) = y_1 & \text{in } G, \end{cases}$$

with initial data $(y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, suitable coefficients a_i ($i = 1, 2, 3, 4$), and source terms f and g . Here, $y_t = \frac{dy}{dt}$.

*Received by the editors March 20, 2007; accepted for publication (in revised form) April 3, 2008; published electronically August 27, 2008. This work was supported by the NSF of China under grant 10525105, grant MTM2005-00714 of the Spanish MEC, and the SIMUMAT projet of the CAM (Spain). Part of this work was done when the author visited the Departamento de Matemáticas at Universidad Autónoma de Madrid and the Shanghai Key Laboratory for Contemporary Applied Mathematics at Fudan University.

<http://www.siam.org/journals/sima/40-2/68578.html>

[†]Key Laboratory of Systems and Control, Academy of Mathematics and Systems Sciences, Chinese Academy of Sciences, Beijing 100080, China, and Yangtze Center of Mathematics, Sichuan University, Chengdu 610064, China (xuzhang@amss.ac.cn).

Put

$$(1.2) \quad \mathcal{H}_T \triangleq L^2_{\mathcal{F}}(\Omega; C([0, T]; H^1_0(G))) \cap L^2_{\mathcal{F}}(\Omega; C^1([0, T]; L^2(G))).$$

Clearly, \mathcal{H}_T is a Banach space with the canonical norm. We begin with the following notion.

DEFINITION 1.1. We call $y \in \mathcal{H}_T$ a solution of (1.1) if the following hold:

- (i) $y(0) = y_0$ in G , P -a.s.
- (ii) For any $t \in [0, T]$ and any $\eta \in H^1_0(G)$, it holds that

$$(1.3) \quad \begin{aligned} & \int_G y_t(t, x)\eta(x)dx - \int_G y_t(0, x)\eta(x)dx \\ &= \int_0^t \int_G \left\{ -\nabla y(s, x) \cdot \nabla \eta(x) + \left[a_1(s, x)y_t(s, x) + a_2(s, x) \cdot \nabla y(s, x) \right. \right. \\ & \quad \left. \left. + a_3(s, x)y(s, x) + f(s, x) \right] \eta(x) \right\} dx ds \\ & \quad + \int_0^t \int_G \left[a_4(s, x)y(s, x) + g(s, x) \right] \eta(x) dx dB(s), \quad P\text{-a.s.} \end{aligned}$$

Under some assumptions, for any initial data $(y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H^1_0(G) \times L^2(G))$, one can show that system (1.1) admits one and only one solution $y \in \mathcal{H}_T$ (see Proposition 3.1 in section 3).

The main purpose of this paper is to derive a (partial) boundary observability estimate for system (1.1), i.e., find (if possible) a constant $\mathcal{C}(a_1, a_2, a_3, a_4) > 0$ such that solutions of system (1.1) satisfy

$$(1.4) \quad \begin{aligned} & |(y(T), y_t(T))|_{L^2(\Omega, \mathcal{F}_T, P; H^1_0(G) \times L^2(G))} \\ & \leq \mathcal{C}(a_1, a_2, a_3, a_4) \left[\left| \frac{\partial y}{\partial \nu} \right|_{L^2_{\mathcal{F}}(0, T; L^2(\Gamma_0))} + |f|_{L^2_{\mathcal{F}}(0, T; L^2(G))} + |g|_{L^2_{\mathcal{F}}(0, T; L^2(G))} \right] \\ & \quad \forall (y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H^1_0(G) \times L^2(G)). \end{aligned}$$

As we shall see in the last section of this paper, inequality (1.4) is strongly related to the state observation problem of semilinear stochastic wave equations.

It is well known that the observability estimate is an important tool for the study of stabilization and controllability problems for deterministic PDEs. We refer the reader to [16] for a recent survey in this respect. Although there are numerous references addressing the observability problems for deterministic PDEs, very little is known for the stochastic counterpart, and it remains to be further understood. Indeed, to the best of our knowledge, [1] is the only publication in this direction, which is devoted to the controllability/observability for the stochastic heat equation. As far as we know, nothing is known for the observability estimate on the stochastic wave equation.

Similar to the deterministic setting, we shall use a stochastic version of the global Carleman estimate to establish inequality (1.4). The difficulty in doing this is the very fact that, unlike the deterministic situation, system (1.1), a stochastic wave equation, is *time-irreversible*. Therefore, one cannot simply mimic the usual Carleman inequality for the deterministic wave equations (see [3, 4, 5, 8, 9, 13] and the references cited therein). Rather, in order to overcome this difficulty, instead of the usual smooth

weight function, we need to introduce another singular weight function in this paper to derive the desired Carleman estimate for system (1.1) (see (2.13) in the next section).

On the other hand, the Carleman estimate is itself a fundamental tool for the study of control and inverse problems for deterministic PDEs (see [7, 16] and the references cited therein). Similar to the situation for the observability estimate, although there are numerous references addressing the Carleman estimate for deterministic PDEs, to the best of our knowledge, [1, 12] are the only references for the stochastic counterpart, which are devoted to the stochastic heat equation. It would be quite interesting to extend the deterministic Carleman estimate for other PDEs to the stochastic ones, but there are many things to be done, and some of them seem to be challenging. In this paper, in order to present the key idea in the simplest way, we do not pursue the full technical generality.

The rest of this paper is organized as follows. The main results of this paper are stated in section 2. In section 3, we show some preliminary results. In section 4, we present a crucial identity for a stochastic hyperbolic-like operator. Then, in section 5, we derive pointwise Carleman-type estimates for the stochastic wave operator. Section 6 is devoted to the proof of Theorems 2.1–2.2. Finally, section 7 addresses giving an application of our observability result.

2. Statement of the main results. Fix any $x_0 \in \mathbb{R}^d \setminus \overline{G}$. It is clear that

$$(2.1) \quad 0 < R_0 \triangleq \min_{x \in G} |x - x_0| < R_1 \triangleq \max_{x \in G} |x - x_0|.$$

Put

$$(2.2) \quad \Gamma_0 \triangleq \{x \in \Gamma \mid (x - x_0) \cdot \nu(x) > 0\},$$

where $\nu(x)$ is the unit outward normal vector of G at $x \in \Gamma$.

Assume that

$$(2.3) \quad \begin{aligned} a_1 &\in L^\infty_{\mathcal{F}}(0, T; L^\infty(G)), & a_2 &\in L^\infty_{\mathcal{F}}(0, T; L^\infty(G; \mathbb{R}^n)), \\ a_3 &\in L^\infty_{\mathcal{F}}(0, T; L^n(G)), & a_4 &\in L^\infty_{\mathcal{F}}(0, T; L^\infty(G)) \end{aligned}$$

and that

$$(2.4) \quad f \in L^2_{\mathcal{F}}(0, T; L^2(G)), \quad g \in L^2_{\mathcal{F}}(0, T; L^2(G)).$$

In what follows, we use the following notation:

$$(2.5) \quad \begin{aligned} \mathcal{A}(a_1, a_2, a_3, a_4) &\triangleq |(a_1, a_4)|^2_{L^\infty_{\mathcal{F}}(0, T; (L^\infty(G))^2)} + |a_2|^2_{L^\infty_{\mathcal{F}}(0, T; L^\infty(G; \mathbb{R}^n))} \\ &\quad + |a_3|^2_{L^\infty_{\mathcal{F}}(0, T; L^n(G))}. \end{aligned}$$

We choose a sufficiently small constant $c \in (0, 1)$ so that (recall (2.1) for R_0 and R_1)

$$(2.6) \quad \frac{(4 + 5c)R_0^2}{9c} > R_1^2.$$

In what follows, we take $T (> 2R_1)$ sufficiently large such that

$$(2.7) \quad \frac{4(4 + 5c)R_0^2}{9c} > c^2 T^2 > 4R_1^2.$$

Our observability estimate for system (1.1) is stated as follows.

THEOREM 2.1. *Let (2.3)–(2.4) hold, Γ_0 be given by (2.2), and T satisfy (2.7). Then solutions of system (1.1) satisfy (1.4) with*

$$(2.8) \quad \mathcal{C}(a_1, a_2, a_3, a_4) = Ce^{C\mathcal{A}(a_1, a_2, a_3, a_4)}.$$

Remark 2.1. From the proof of Theorem 2.1, it is easy to see that the conclusion can be slightly strengthened as follows: for any $t \in (0, T]$, solutions of system (1.1) satisfy

$$(2.9) \quad \begin{aligned} & |(y(t), y_t(t))|_{L^2(\Omega, \mathcal{F}_t, P; H_0^1(G) \times L^2(G))} \\ & \leq e^{Ct^{-1}} \mathcal{C}(a_1, a_2, a_3, a_4) \left[\left| \frac{\partial y}{\partial \nu} \right|_{L^2_{\mathcal{F}}(0, T; L^2(\Gamma_0))} + |f|_{L^2_{\mathcal{F}}(0, T; L^2(G))} + |g|_{L^2_{\mathcal{F}}(0, T; L^2(G))} \right] \\ & \quad \forall (y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G)), \end{aligned}$$

where $\mathcal{C}(a_1, a_2, a_3, a_4)$ is given by (2.8).

Remark 2.2. A deterministic version of (1.1) reads

$$(2.10) \quad \begin{cases} w_{tt} - \Delta w = b_1 w_t + b_2 \cdot \nabla w + b_3 w + h & \text{in } Q, \\ w = 0 & \text{on } \Sigma, \\ w(0) = w_0, \quad w_t(0) = w_1 & \text{in } G, \end{cases}$$

where $b_1 \in L^\infty(Q)$, $b_2 \in L^\infty(Q; \mathbb{R}^n)$, $b_3 \in L^\infty(0, T; L^p(G))$ with $p \in [n, \infty]$, and $h \in L^2(Q)$. As a special case of [3, Theorems 2.2 and 2.3] and noting the time reversibility of system (2.10), the following counterpart of Theorem 2.1 holds: If $T > 2R_1$, then solutions of (2.10) satisfy

$$(2.11) \quad \begin{aligned} & |(w(T), w_t(T))|_{H_0^1(G) \times L^2(G)} \\ & \leq Ce^C \left[|b_1|_{L^\infty(Q)}^2 + |b_2|_{L^\infty(Q; \mathbb{R}^n)}^2 + |b_3|_{L^\infty(0, T; L^p(G))}^{\frac{1}{3/2 - n/p}} \right] \left[\left| \frac{\partial w}{\partial \nu} \right|_{L^2(\Sigma_0)} + |h|_{L^2(Q)} \right] \\ & \quad \forall (w_0, w_1) \in H_0^1(G) \times L^2(G). \end{aligned}$$

There are three main differences between Theorem 2.1 and this result. The first is that one can replace the left-hand side of (2.11) by $|(w_0, w_1)|_{H_0^1(G) \times L^2(G)}$. However, due to the time irreversibility of system (1.1), one cannot do the same in the stochastic setting, i.e., replacing the left-hand side of (1.4) by $|(y_0, y_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))}$. The second is that we assume that T in Theorem 2.1 satisfies (2.7), which is usually much more restrictive than that $T > 2R_1$ for the deterministic setting (see Remark 2.4 below for more explanation about this). The third is that the observability constant $\mathcal{C}(a_1, a_2, a_3, a_4)$ in Theorem 2.1 is not as sharp as that in (2.11). Indeed, it is clear that (recall (2.5) for $\mathcal{A}(\cdot, \cdot, \cdot, \cdot)$)

$$|b_1|_{L^\infty(Q)}^2 + |b_2|_{L^\infty(Q; \mathbb{R}^n)}^2 + |b_3|_{L^\infty(0, T; L^p(G))}^{\frac{1}{3/2 - n/p}} \leq \mathcal{A}(b_1, b_2, b_3, 0) + C \quad \forall p \in [n, \infty].$$

Remark 2.3. It is well known that a sharp condition guaranteeing observability inequality (2.11) (at least when b_1, b_2 , and b_3 are time-invariant) is that the triple

(G, Γ_0, T) satisfies the geometric optic condition introduced in [2]. It would be quite interesting to extend this result to the stochastic setting, but this is an open problem.

As mentioned in section 1, in order to prove Theorem 2.1, we need to derive a Carleman estimate (with singular weight function) for system (1.1). For this purpose, for any (large) $\lambda > 0$ and $c \in (0, 1)$ given in (2.6), set

$$(2.12) \quad \ell = \ell(t, x) \triangleq \lambda \left[|x - x_0|^2 - c \left(t - \frac{T}{2} \right)^2 \right], \quad \theta \triangleq e^\ell.$$

Also, for any $\beta > 0$, set

$$(2.13) \quad \Theta = \Theta(t) \triangleq \exp \left\{ -\frac{\beta}{t(T-t)} \right\}, \quad 0 < t < T.$$

It is easy to see that $\Theta(t)$ decays rapidly to 0 as $t \rightarrow 0$ or $t \rightarrow T$. Our Carleman estimate for system (1.1) is stated as follows.

THEOREM 2.2. *Let (2.3)–(2.4) hold, Γ_0 be given by (2.2), and c and T satisfy, respectively, (2.6) and (2.7). Then there exist a constant $\beta > 0$ (which is very small) and a constant*

$$\lambda^* = C \left[1 + \mathcal{A}(a_1, a_2, a_3, a_4) \right]$$

such that solutions of system (1.1) satisfy

$$(2.14) \quad \begin{aligned} & \lambda \mathbb{E} \int_Q \Theta \theta^2 (y_t^2 + |\nabla y|^2 + \lambda^2 y^2) dx dt \\ & \leq C \mathbb{E} \left\{ \lambda \int_{\Sigma_0} \Theta \theta^2 \left| \frac{\partial y}{\partial \nu} \right|^2 d\Sigma_0 + \int_Q \Theta \theta^2 (f^2 + \lambda g^2) dx dt \right\} \\ & \quad \forall (y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G)), \quad \forall \lambda \geq \lambda^*. \end{aligned}$$

Remark 2.4. The restriction on T in (2.7) is a technical condition, which does not seem to be natural. However, this condition plays a key role in Step 3 in the proof of Theorem 5.1, a crucial preliminary for the proof of Theorem 2.2. As in the deterministic setting (recall Remark 2.2), it is reasonable to expect that it should be improved to be $T > 2R_1$. But this is an open problem.

Remark 2.5. We now recall the Carleman estimate for the deterministic wave equation, i.e., system (2.10). Fix a constant $c \in (0, 1)$ so that $\left(\frac{2R_1}{T}\right)^2 < c < \frac{2R_1}{T}$ (which is possible since $T > 2R_1$ for the deterministic situation). By [5, equation (11.12) in the proof of Theorem 5.1] and similar to [3, Theorem 2.3], we conclude that there exists a constant

$$\lambda_* = C \left[1 + |b_1|_{L^\infty(Q)}^2 + |b_2|_{L^\infty(Q; \mathbb{R}^n)}^2 + |b_3|_{L^\infty(0,T; L^p(G))}^{\frac{1}{3/2-n/p}} \right]$$

such that solutions of system (2.10) satisfy

$$(2.15) \quad \begin{aligned} & \lambda \int_Q \theta^2 (\lambda^2 w^2 + w_t^2 + |\nabla w|^2) dx dt \leq C \left[\lambda \int_{\Sigma_0} \theta^2 \left| \frac{\partial w}{\partial \nu} \right|^2 d\Sigma_0 + \int_Q \theta^2 |h|^2 dx dt \right] \\ & \quad \forall (w_0, w_1) \in H_0^1(G) \times L^2(G), \quad \forall \lambda \geq \lambda_*. \end{aligned}$$

The main difference between Theorem 2.2 and this result is, as mentioned before, one has to introduce a singular weight Θ in (2.14).

Remark 2.6. We consider here the simplest case of one dimensional standard Brownian motion. It would be interesting to extend the results in this paper to the case of colored (infinite dimensional) noise, or even with both state- and control-dependent noise. But these remain to be done.

Remark 2.7. It would be quite interesting to study Carleman and observability estimates for backward stochastic wave equations. To the best of our knowledge, this is a challenging problem, and nothing is known in this respect.

3. Preliminaries. In this section, we show some preliminary results that will subsequently be used.

In what follows, for simplicity, we denote $\sum_{i,j=1}^n$ and $\sum_{i=1}^n$ simply by $\sum_{i,j}$ and \sum_i , respectively. Also, we will use the notation $u_i = u_{x_i}$, where x_i is the i th coordinate of a generic point $x = (x_1, \dots, x_n)$ in \mathbb{R}^n . In a similar manner, we use the notation ℓ_i, v_i , etc., for the partial derivatives of ℓ and v with respect to x_i .

First of all, we have the following well-posedness result for system (1.1).

PROPOSITION 3.1. *Under assumptions (2.3)–(2.4), for any $(y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, system (1.1) admits one and only one solution $y \in \mathcal{H}_T$. Moreover (recall (2.8) for $\mathcal{C}(a_1, a_2, a_3, a_4)$),*

$$(3.1) \quad |y|_{\mathcal{H}_T} \leq \mathcal{C}(a_1, a_2, a_3, a_4) \left[|(y_0, y_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))} + |f|_{L^2_{\mathcal{F}}(0, T; L^2(G))} + |g|_{L^2_{\mathcal{F}}(0, T; L^2(G))} \right].$$

Proof. The detailed proof is lengthy but almost standard. Therefore, we give below only a sketch.

Step 1. First of all, following the *parabolic regularization approach* in the proof of [6, Theorem 4.1, p. 159] (for this, one needs to make numerous but small changes), one can show that, for any $(w_0, w_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, $\hat{f} \in L^2_{\mathcal{F}}(0, T; L^2(G))$, and $\hat{g} \in L^2_{\mathcal{F}}(0, T; L^2(G))$, the system

$$(3.2) \quad \begin{cases} dw_t - \Delta w dt = \hat{f} dt + \hat{g} dB(t) & \text{in } Q, \\ w = 0 & \text{on } \Sigma, \\ w(0) = w_0, \quad w_t(0) = w_1 & \text{in } G \end{cases}$$

admits a unique solution $w \in \mathcal{H}_T$. Furthermore, there is a constant $C = C(T) > 0$ such that, for any $\tau \in [0, T]$, it holds that

$$(3.3) \quad |w|_{\mathcal{H}_\tau} \leq C \left[|(w_0, w_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))} + |\hat{f}|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))} + |\hat{g}|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))} \right],$$

where \mathcal{H}_τ is defined similarly to (1.2).

Step 2. Next, we show that system (1.1) admits a local (in time) solution $y \in \mathcal{H}_\tau$ for small $\tau \in (0, T]$. For this purpose, fixing any $y \in \mathcal{H}_\tau$, we solve the following system:

$$(3.4) \quad \begin{cases} dz_t - \Delta z dt = (a_1 y_t + a_2 \cdot \nabla y + a_3 y + f) dt + (a_4 y + g) dB(t) & \text{in } Q, \\ z = 0 & \text{on } \Sigma, \\ z(0) = y_0, \quad z_t(0) = y_1 & \text{in } G. \end{cases}$$

By Step 1, system (3.4) (by viewing, respectively, $a_1y_t + a_2 \cdot \nabla y + a_3y + f$ and $a_4y + g$ as nonhomogeneous terms \hat{f} and \hat{g}) admits a unique solution $z \in \mathcal{H}_\tau$. We define a map \mathcal{F} from \mathcal{H}_τ into itself by

$$\mathcal{F}(y) = z.$$

For any $y^1, y^2 \in \mathcal{H}_\tau$, by means of (3.3), it is easy to verify that

$$|\mathcal{F}(y^1) - \mathcal{F}(y^2)|_{\mathcal{H}_\tau} \leq C\tau|y^1 - y^2|_{\mathcal{H}_\tau}.$$

This shows that the map \mathcal{F} is contractive when τ is small enough. Therefore, \mathcal{F} admits a fixed point $y \in \mathcal{H}_\tau$. Hence, system (1.1) admits a local (in time) solution $y \in \mathcal{H}_\tau$.

Step 3. To conclude the proof of Proposition 3.1, it suffices to show that estimate (3.1) holds.

For any $\tau \in [0, T]$, applying estimate (3.3) to system (1.1) (by viewing, respectively, $a_1y_t + a_2 \cdot \nabla y + a_3y + f$ and $a_4y + g$ as nonhomogeneous terms \hat{f} and \hat{g}), we obtain (recall (2.5) for $\mathcal{A}(a_1, a_2, a_3, a_4)$)

(3.5)

$$\begin{aligned} |y|_{\mathcal{H}_\tau}^2 &\leq C \left[|(y_0, y_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))}^2 \right. \\ &\quad \left. + |a_1y_t + a_2 \cdot \nabla y + a_3y + f|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))}^2 + |a_4y + g|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))}^2 \right] \\ &\leq C \left[|(y_0, y_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))}^2 + |f|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))}^2 + |g|_{L^2_{\mathcal{F}}(0, \tau; L^2(G))}^2 \right. \\ &\quad \left. + \mathcal{A}(a_1, a_2, a_3, a_4) \int_0^\tau |y|_{\mathcal{H}_s}^2 ds \right]. \end{aligned}$$

Now the desired result (3.1) follows from (3.5) and Gronwall’s inequality. \square

Next, we show the following identity.

PROPOSITION 3.2. *Let $\mu = \mu(x) \triangleq (\mu^1, \dots, \mu^n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a vector field of class C^1 and w an $H_{loc}^2(\mathbb{R}^n)$ -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted process such that w_t is an $L_{loc}^2(\mathbb{R}^n)$ -valued semimartingale. Then, for a.e. $x \in \mathbb{R}^n$ and P -a.s. $\omega \in \Omega$, it holds that*

$$\begin{aligned} &-\nabla \cdot [2(\mu \cdot \nabla w)\nabla w + \mu(w_t^2 - |\nabla w|^2)] dt \\ (3.6) \quad &= 2 \left[\mu \cdot \nabla w(dw_t - \Delta w dt) - d(w_t \mu \cdot \nabla w) - \sum_{i,j} w_i w_j \frac{\partial \mu^j}{\partial x_i} dt \right] \\ &+ (|\nabla w|^2 - w_t^2)\nabla \cdot \mu dt. \end{aligned}$$

Proof. Noting that μ is time-independent, it follows that

$$\begin{aligned} 2[\mu \cdot \nabla w \Delta w dt + d(w_t \mu \cdot \nabla w)] &= 2 \left[\sum_{i,j} \mu^j w_j w_{ii} dt + \mu \cdot \nabla w dw_t + w_t \mu \cdot \nabla w_t dt \right] \\ &= 2 \sum_{i,j} \left[(w_i \mu^j w_j)_i - w_i w_j \frac{\partial \mu^j}{\partial x_i} \right] dt - \mu \cdot \nabla |\nabla w|^2 dt + 2\mu \cdot \nabla w dw_t + \mu \cdot \nabla w_t^2 dt \\ &= \nabla \cdot [2(\mu \cdot \nabla w) \nabla w + \mu (w_t^2 - |\nabla w|^2)] dt + 2 \left[\mu \cdot \nabla w dw_t - \sum_{i,j} w_i w_j \frac{\partial \mu^j}{\partial x_i} dt \right] \\ &\quad + (|\nabla w|^2 - w_t^2) \nabla \cdot \mu dt. \end{aligned}$$

This gives (3.6). \square

Finally, we show the following hidden regularity for the solution of system (1.1) (here, by hidden regularity we mean that it follows from the equation rather than from the usual trace theorem in the theory of Sobolev spaces).

PROPOSITION 3.3. *Under assumptions (2.3)–(2.4), for any $(y_0, y_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, the solution of system (1.1) satisfies $\frac{\partial y}{\partial \nu} \in L^2_{\mathcal{F}}(0, T; L^2(\Gamma))$. Moreover,*

(3.7)

$$\begin{aligned} &\left| \frac{\partial y}{\partial \nu} \right|_{L^2_{\mathcal{F}}(0, T; L^2(\Gamma))} \\ &\leq C \left[|(y_0, y_1)|_{L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))} + |f|_{L^2_{\mathcal{F}}(0, T; L^2(G))} + |g|_{L^2_{\mathcal{F}}(0, T; L^2(G))} \right] \\ &\quad \times \exp \left\{ C \left[|(a_1, a_4)|_{L^{\infty}_{\mathcal{F}}(0, T; (L^{\infty}(G))^2)}^2 + |a_2|_{L^{\infty}_{\mathcal{F}}(0, T; L^{\infty}(G; \mathbb{R}^n))}^2 + |a_3|_{L^{\infty}_{\mathcal{F}}(0, T; L^n(G))}^2 \right] \right\}. \end{aligned}$$

Proof. Since $\Gamma \in C^2$, one can find a vector field $\mu_0 = (\mu_0^1, \dots, \mu_0^n) \in C^1(\bar{\Omega}; \mathbb{R}^n)$ such that $\mu_0 = \nu$ on Γ (see [10]). Applying Proposition 3.2 with $\mu = \mu_0$ and $w = y$, by means of Proposition 3.1, following [10, 13], it is not difficult to show $\frac{\partial y}{\partial \nu} \in L^2_{\mathcal{F}}(0, T; L^2(\Gamma))$ and the desired estimate (3.7) (hence we omit the details). \square

4. Identity for a stochastic hyperbolic-like operator. In this section, we show the following fundamental identity for a stochastic hyperbolic-like operator.

THEOREM 4.1. *Let $b^{ij} \in C^1((0, T) \times \mathbb{R}^n)$ satisfy*

$$(4.1) \quad b^{ij} = b^{ji}, \quad i, j = 1, 2, \dots, n,$$

$\ell, \Psi \in C^2((0, T) \times \mathbb{R}^n)$. Assume u is an $H^2_{loc}(\mathbb{R}^n)$ -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted process such that u_t is an $L^2_{loc}(\mathbb{R}^n)$ -valued semimartingale. Set $\theta = e^{\ell}$ and $v = \theta u$. Then, for

a.e. $x \in \mathbb{R}^n$ and P -a.s. $\omega \in \Omega$,

$$\begin{aligned}
 & \theta \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) \left[du_t - \sum_{i,j} (b^{ij} u_i)_j dt \right] \\
 & + \sum_{i,j} \left[\sum_{i',j'} \left(2b^{ij} b^{i'j'} \ell_{i'} v_i v_{j'} - b^{ij} b^{i'j'} \ell_i v_{i'} v_{j'} \right) - 2b^{ij} \ell_t v_i v_t + b^{ij} \ell_i v_t^2 \right. \\
 & \quad \left. + \Psi b^{ij} v_i v - \left(A\ell_i + \frac{\Psi_i}{2} \right) b^{ij} v^2 \right]_j dt \\
 (4.2) \quad & + d \left[\sum_{i,j} b^{ij} \ell_t v_i v_j - 2 \sum_{i,j} b^{ij} \ell_i v_j v_t + \ell_t v_t^2 - \Psi v_t v + \left(A\ell_t + \frac{\Psi_t}{2} \right) v^2 \right] \\
 & = \left\{ \left[\ell_{tt} + \sum_{i,j} (b^{ij} \ell_i)_j - \Psi \right] v_t^2 - 2 \sum_{i,j} [(b^{ij} \ell_j)_t + b^{ij} \ell_{tj}] v_i v_t \right. \\
 & \quad \left. + \sum_{i,j} \left[(b^{ij} \ell_t)_t + \sum_{i',j'} \left(2b^{ij'} (b^{i'j} \ell_{i'})_{j'} - (b^{ij} b^{i'j'} \ell_{i'})_{j'} \right) + \Psi b^{ij} \right] v_i v_j \right. \\
 & \quad \left. + Bv^2 + \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right)^2 \right\} dt + \theta^2 \ell_t (du_t)^2,
 \end{aligned}$$

where $(du_t)^2$ denotes the quadratic variation process of u_t ,

$$(4.3) \quad \begin{cases} A \triangleq (\ell_t^2 - \ell_{tt}) - \sum_{i,j} (b^{ij} \ell_i \ell_j - b_j^{ij} \ell_i - b^{ij} \ell_{ij}) - \Psi, \\ B \triangleq A\Psi + (A\ell_t)_t - \sum_{i,j} (Ab^{ij} \ell_i)_j + \frac{1}{2} \left[\Psi_{tt} - \sum_{i,j} (b^{ij} \Psi_i)_j \right]. \end{cases}$$

Proof. By $v(t, x) = \theta(t, x)u(t, x)$, we have $u_t = \theta^{-1}(v_t - \ell_t v)$ and $u_j = \theta^{-1}(v_j - \ell_j v)$ for $j = 1, 2, \dots, n$. Hence,

$$(4.4) \quad du_t = \theta^{-1} [dv_t - 2\ell_t v_t dt + (\ell_t^2 - \ell_{tt}) v dt].$$

Similarly, by symmetry condition (4.1), one may check that

$$(4.5) \quad \sum_{i,j} (b^{ij} u_i)_j = \theta^{-1} \sum_{i,j} \left[(b^{ij} v_i)_j - 2b^{ij} \ell_i v_j + (b^{ij} \ell_i \ell_j - b_j^{ij} \ell_i - b^{ij} \ell_{ij}) v \right].$$

Therefore, by (4.4)–(4.5) and recalling the definition of A in (4.3), we get

$$\begin{aligned}
 & \theta \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) \left[du_t - \sum_{i,j} (b^{ij} u_i)_j dt \right] \\
 &= \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) \left\{ dv_t - \left[\sum_{i,j} (b^{ij} v_i)_j - Av \right. \right. \\
 &\quad \left. \left. + 2\ell_t v_t - 2 \sum_{i,j} b^{ij} \ell_i v_j - \Psi v \right] dt \right\} \\
 (4.6) \quad &= \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) dv_t \\
 &\quad + \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) \left[- \sum_{i,j} (b^{ij} v_i)_j + Av \right] dt \\
 &\quad + \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right)^2 dt.
 \end{aligned}$$

We now analyze the first two terms in the right-hand side of (4.6).

First, using Itô's formula, we have

$$\begin{aligned}
 (4.7) \quad & \left(-2\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) dv_t \\
 &= d \left[\left(-\ell_t v_t + 2 \sum_{i,j} b^{ij} \ell_i v_j + \Psi v \right) v_t \right] \\
 &\quad - \left[-\ell_{tt} v_t^2 + 2 \sum_{i,j} (b^{ij} \ell_i)_t v_j v_t + 2 \sum_{i,j} b^{ij} \ell_i v_{tj} v_t + \Psi v_t^2 + \Psi_t v v_t \right] dt + \ell_t (dv_t)^2 \\
 &= d \left(-\ell_t v_t^2 + 2 \sum_{i,j} b^{ij} \ell_i v_j v_t + \Psi v v_t - \frac{\Psi_t}{2} v^2 \right) \\
 &\quad + \left\{ - \sum_{i,j} (b^{ij} \ell_i v_t^2)_j + \left[\ell_{tt} + \sum_{i,j} (b^{ij} \ell_i)_j - \Psi \right] v_t^2 - 2 \sum_{i,j} (b^{ij} \ell_j)_t v_i v_t + \frac{\Psi_{tt}}{2} v^2 \right\} dt \\
 &\quad + \theta^2 \ell_t (du_t)^2.
 \end{aligned}$$

Next,

$$\begin{aligned}
 & -2\ell_t v_t \left[-\sum_{i,j} (b^{ij} v_i)_j + Av \right] \\
 (4.8) \quad & = 2 \left[\sum_{i,j} (b^{ij} \ell_t v_i v_t)_j - \sum_{i,j} b^{ij} \ell_{tj} v_i v_t \right] - \sum_{i,j} b^{ij} \ell_t (v_i v_j)_t - A\ell_t (v^2)_t \\
 & = 2 \left[\sum_{i,j} (b^{ij} \ell_t v_i v_t)_j - \sum_{i,j} b^{ij} \ell_{tj} v_i v_t \right] + \sum_{i,j} (b^{ij} \ell_t)_t v_i v_j \\
 & \quad - \left(\sum_{i,j} b^{ij} \ell_t v_i v_j + A\ell_t v^2 \right)_t + (A\ell_t)_t v^2.
 \end{aligned}$$

Further, by means of a direct computation, one may check that

$$\begin{aligned}
 & 2 \sum_{i,j} b^{ij} \ell_i v_j \left[-\sum_{i,j} (b^{ij} v_i)_j + Av \right] \\
 (4.9) \quad & = -\sum_{i,j} \left[\sum_{i',j'} \left(2b^{ij} b^{i'j'} \ell_{i'} v_i v_{j'} - b^{ij} b^{i'j'} \ell_i v_{i'} v_{j'} \right) - Ab^{ij} \ell_i v^2 \right]_j \\
 & \quad + \sum_{i,j,i',j'} \left[2b^{ij'} (b^{i'j} \ell_{i'})_{j'} - (b^{ij} b^{i'j'} \ell_{i'})_{j'} \right] v_i v_j - \sum_{i,j} (Ab^{ij} \ell_i)_j v^2
 \end{aligned}$$

and

$$\begin{aligned}
 (4.10) \quad \Psi v \left[-\sum_{i,j} (b^{ij} v_i)_j + Av \right] & = -\sum_{i,j} \left(\Psi b^{ij} v_i v - \frac{\Psi_i}{2} b^{ij} v^2 \right)_j + \Psi \sum_{i,j} b^{ij} v_i v_j \\
 & \quad + \left[-\frac{1}{2} \sum_{i,j} (b^{ij} \Psi_i)_j + A\Psi \right] v^2.
 \end{aligned}$$

Finally, combining (4.6)–(4.10), we arrive at the desired equality (4.2). \square

5. Pointwise Carleman-type estimates for the stochastic wave operator.

In this section, we show a pointwise Carleman-type estimate (with singular weight) for the stochastic wave operator “ $du_t - \Delta u dt$.”

To begin with, by taking $(b^{ij})_{n \times n} = I$, the identity matrix, and $\theta = e^\ell$ (with ℓ given in (2.12)) in Theorem 4.1, one obtains the following pointwise Carleman-type estimate for the stochastic wave operator.

LEMMA 5.1. *Let $\ell, \Psi \in C^2((0, T) \times \mathbb{R}^n)$ and $k \in \mathbb{R}$. Assume u is an $H^2_{loc}(\mathbb{R}^n)$ -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted process such that u_t is an $L^2_{loc}(\mathbb{R}^n)$ -valued semimartingale.*

Set $v = \theta u$. Then, for a.e. $x \in \mathbb{R}^n$ and P -a.s. $\omega \in \Omega$, it holds that

$$\begin{aligned}
 & \theta(-2\ell_t v_t + 2\nabla\ell \cdot \nabla v + \psi v)(du_t - \Delta u dt) \\
 & + d\left[\ell_t(v_t^2 + |\nabla v|^2) - 2(\nabla\ell) \cdot (\nabla v)v_t - \Psi v v_t + A\ell_t v^2\right] \\
 (5.1) \quad & + \sum_i \left\{2v_i(\nabla\ell) \cdot (\nabla v) - \ell_i|\nabla v|^2 - 2\ell_t v_t v_i + \ell_i v_t^2 + \Psi v v_i - A\ell_i v^2\right\}_i dt \\
 & \geq \left[(1-k)\lambda v_t^2 + (k+3-4c)\lambda|\nabla v|^2 + Bv^2\right. \\
 & \left. + \left(-2\ell_t v_t + 2\nabla\ell \cdot \nabla v + \psi v\right)^2\right] dt + \theta^2 \ell_t (du_t)^2,
 \end{aligned}$$

where

$$(5.2) \quad \begin{cases} \Psi \triangleq (2n - 2c - 1 + k)\lambda, \\ A = 4 \left[c^2 \left(t - \frac{T}{2} \right)^2 - |x - x_0|^2 \right] \lambda^2 + \lambda(4c + 1 - k), \\ B = 4 \left[(4c + 5 - k)|x - x_0|^2 - (8c + 1 - k)c^2 \left(t - \frac{T}{2} \right)^2 \right] \lambda^3 + O(\lambda^2). \end{cases}$$

The desired pointwise Carleman-type estimate (with singular weight function Θ) for the stochastic wave operator reads as follows.

THEOREM 5.1. *Assume u is an $H^2(G)$ -valued $\{\mathcal{F}_t\}_{t \geq 0}$ -adapted process such that u_t is an $L^2(G)$ -valued semimartingale. Let $v = \theta u$ and T satisfy (2.7). Then there exist three constants $\lambda_0 > 0$, $\beta_0 > 0$, and $c_0 > 0$, independent of u , such that, for all $\beta \in (0, \beta_0)$ and $\lambda \geq \lambda_0$ and a.e. $x \in G$ and P -a.s. $\omega \in \Omega$, it holds that*

(5.3)

$$\begin{aligned}
 & \Theta\theta(-2\ell_t v_t + 2\nabla\ell \cdot \nabla v + \psi v)(du_t - \Delta u dt) \\
 & + d\left\{\Theta\left[\ell_t(v_t^2 + |\nabla v|^2) - 2(\nabla\ell) \cdot (\nabla v)v_t - \Psi v v_t + A\ell_t v^2\right]\right\} \\
 & + \sum_i \left\{\Theta\left[2v_i(\nabla\ell) \cdot (\nabla v) - \ell_i|\nabla v|^2 - 2\ell_t v_t v_i + \ell_i v_t^2 + \Psi v v_i - A\ell_i v^2\right]\right\}_i dt \\
 & \geq \left[c_0\lambda\Theta\theta^2(u_t^2 + |\nabla u|^2 + \lambda^2 u^2) + \Theta\left(-2\ell_t v_t + 2\nabla\ell \cdot \nabla v + \psi v\right)^2\right] dt + \Theta\theta^2 \ell_t (du_t)^2,
 \end{aligned}$$

with A and Ψ given by (5.2).

Remark 5.1. The main difference between the pointwise estimates (5.1) and (5.3) is that we introduce a singular “pointwise” weight in (5.3). Another difference between (5.1) and (5.3) is that T is arbitrary in the former estimate, while for the latter one needs to take T to be large enough.

Proof of Theorem 5.1. We borrow some idea from the proof of [15, Theorem 1]. The proof is divided into several steps.

Step 1. We multiply both sides of inequality (5.1) by Θ . Obviously, we have (recall (5.2) for A and Ψ)

$$\begin{aligned}
 & \Theta d \left[\ell_t (v_t^2 + |\nabla v|^2) - 2(\nabla \ell) \cdot (\nabla v) v_t - \Psi v v_t + A \ell_t v^2 \right] \\
 (5.4) \quad & = d \left\{ \Theta \left[\ell_t (v_t^2 + |\nabla v|^2) - 2(\nabla \ell) \cdot (\nabla v) v_t - \Psi v v_t + A \ell_t v^2 \right] \right\} \\
 & \quad - \frac{\beta(T-2t)}{t^2(T-t)^2} \Theta \left[\ell_t (v_t^2 + |\nabla v|^2) - 2(\nabla \ell) \cdot (\nabla v) v_t - \Psi v v_t + A \ell_t v^2 \right] dt.
 \end{aligned}$$

Note that

$$\begin{aligned}
 & \left| -\frac{\beta(T-2t)}{t^2(T-t)^2} \Theta [-2(\nabla \ell) \cdot (\nabla v) v_t - \Psi v v_t] \right| \\
 (5.5) \quad & \leq \frac{\beta|T-2t|}{t^2(T-t)^2} \Theta [2|(\nabla \ell) \cdot (\nabla v) v_t| + |\Psi v v_t|] \\
 & \leq \frac{\beta|T-2t|}{t^2(T-t)^2} \Theta \left[(|\nabla \ell| + 1)v_t^2 + |\nabla \ell| |\nabla v|^2 + \frac{1}{4} \Psi^2 v^2 \right].
 \end{aligned}$$

Thus, by (5.1) and using (5.4)–(5.5), we get

$$\begin{aligned}
 & \Theta \theta (-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v) (du_t - \Delta u dt) \\
 & \quad + d \left\{ \Theta \left[\ell_t (v_t^2 + |\nabla v|^2) - 2(\nabla \ell) \cdot (\nabla v) v_t - \Psi v v_t + A \ell_t v^2 \right] \right\} \\
 & \quad + \sum_i \left\{ \Theta \left[2v_i (\nabla \ell) \cdot (\nabla v) - \ell_i |\nabla v|^2 - 2\ell_t v_t v_i + \ell_i v_t^2 + \Psi v v_i - A \ell_i v^2 \right] \right\}_i dt \\
 (5.6) \quad & \geq \left\{ \Theta(1-k)\lambda v_t^2 + \Theta(k+3-4c)\lambda |\nabla v|^2 + \frac{\beta(T-2t)}{t^2(T-t)^2} \ell_t \Theta (v_t^2 + |\nabla v|^2) \right. \\
 & \quad - \frac{\beta|T-2t|}{t^2(T-t)^2} \Theta [(|\nabla \ell| + 1)v_t^2 + |\nabla \ell| |\nabla v|^2] \\
 & \quad + \left[B + \frac{\beta(T-2t)}{t^2(T-t)^2} \ell_t A - \frac{\beta|T-2t|}{4t^2(T-t)^2} \Psi^2 \right] \Theta v^2 \\
 & \quad \left. + \Theta \left(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v \right)^2 \right\} dt + \Theta \theta^2 \ell_t (du_t)^2,
 \end{aligned}$$

where B is given by (5.2).

Step 2. Recalling that ℓ and Ψ are given, respectively, by (2.12) and (5.2), we get

$$\begin{aligned}
 (5.7) \quad \text{right-hand side of (5.6)} & = \left[\lambda \Theta (F_1 v_t^2 + F_2 |\nabla v|^2) + \lambda^3 \Theta G v^2 \right. \\
 & \quad \left. + \Theta \left(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v \right)^2 \right] dt + \Theta \theta^2 \ell_t (du_t)^2,
 \end{aligned}$$

where

$$(5.8) \quad F_1 \triangleq 1 - k + \frac{c\beta(T-2t)^2}{t^2(T-t)^2} - \frac{\beta|T-2t|}{t^2(T-t)^2} (2|x-x_0| + \lambda^{-1}),$$

$$(5.9) \quad F_2 \triangleq k + 3 - 4c + \frac{c\beta(T-2t)^2}{t^2(T-t)^2} - \frac{2\beta|T-2t||x-x_0|}{t^2(T-t)^2},$$

and

$$(5.10) \quad G \triangleq 4 \left[(4c + 5 - k)|x - x_0|^2 - (8c + 1 - k)c^2 \left(t - \frac{T}{2} \right)^2 \right] + O(\lambda^{-1}) \\ + \frac{\beta|T - 2t|}{t^2(T - t)^2} \left\{ 4c|T - 2t| \left[c^2(t - T/2)^2 - |x - x_0|^2 \right] + O(\lambda^{-1}) \right\}.$$

Step 3. Let us show that F_1 , F_2 , and G are positive when λ is large enough and β is sufficiently small. For this, put

$$F_1^0 \triangleq 1 - k, \quad F_2^0 \triangleq k + 3 - 4c, \\ G^0 \triangleq 4 \left[(4c + 5 - k)|x - x_0|^2 - (8c + 1 - k)c^2 \left(t - \frac{T}{2} \right)^2 \right] + O(\lambda^{-1}),$$

which are, respectively, the nonsingular part of F_1 , F_2 , and G . Similarly, put

$$F_1^1 \triangleq \frac{c\beta(T - 2t)^2}{t^2(T - t)^2} - \frac{\beta|T - 2t|}{t^2(T - t)^2} (2|x - x_0| + \lambda^{-1}), \\ F_2^1 \triangleq \frac{c\beta(T - 2t)^2}{t^2(T - t)^2} - \frac{2\beta|T - 2t||x - x_0|}{t^2(T - t)^2}, \\ G^1 \triangleq \frac{\beta|T - 2t|}{t^2(T - t)^2} \left\{ 4c|T - 2t| \left[c^2(t - T/2)^2 - |x - x_0|^2 \right] + O(\lambda^{-1}) \right\},$$

which are, respectively, the singular part of F_1 , F_2 , and G .

Further, we choose $k = 1 - c$. It is easy to see that both F_1^0 and F_2^0 are positive, and

$$G^0 \geq 4(4 + 5c)R_0^2 - 9c^3T^2 + O(\lambda^{-1}),$$

which, via the first inequality in (2.7), is positive, provided that λ is sufficiently large.

When t is close to 0 or T , i.e., $t \in I_0 \triangleq (0, \delta_0) \cup (T - \delta_0, T)$ for some sufficiently small $\delta_0 \in (0, T/2)$, the dominant terms in F_i ($i = 1, 2$) and G are the singular ones. For $t \in I_0$,

$$F_1^1 \geq \frac{\beta|T - 2t|}{t^2(T - t)^2} [c(T - 2\delta_0) - 2R_1 - \lambda^{-1}] = \frac{\beta|T - 2t|}{t^2(T - t)^2} (cT - 2R_1 - 2c\delta_0 - \lambda^{-1}),$$

which, via the second inequality in (2.7), is positive, provided that both δ_0 and λ^{-1} are sufficiently small. Similarly, for $t \in I_0$, F_2^0 is positive, provided that δ_0 is sufficiently small. Further, for $t \in I_0$,

$$G^1 \geq \frac{\beta|T - 2t|}{t^2(T - t)^2} \left\{ 4c|T - 2\delta_0| \left[c^2(\delta_0 - T/2)^2 - R_1^2 \right] + O(\lambda^{-1}) \right\} \\ \geq \frac{\beta|T - 2t|}{t^2(T - t)^2} \left\{ 4c|T - 2\delta_0| \left[c^2T^2/4 - R_1^2 + c^2\delta_0(\delta_0 - T) \right] + O(\lambda^{-1}) \right\},$$

which, via the second inequality in (2.7), is positive, provided that both δ_0 and λ^{-1} are sufficiently small.

By (5.8)–(5.10), we see that $F_1 = F_1^0 + F_1^1$, $F_2 = F_2^0 + F_2^1$, and $G = G^0 + G^1$. Noting the positivity of F_1^0 , F_2^0 , and G^0 , by the above argument, we see that F_1 , F_2 , and G are positive for $t \in I_0$. For $t \in (0, T) \setminus I_0$, noting again the positivity of F_1^0 , F_2^0 , and G^0 , one can choose $\beta > 0$ small enough such that F_1^1 , F_2^1 , and G^1 are very small so that F_1 , F_2 , and G are positive. Hence (5.6)–(5.7) yield the desired (5.3). This completes the proof of Theorem 5.1. \square

6. Proof of Theorems 2.1–2.2. We are now in a position to prove Theorems 2.1–2.2.

Proof of Theorem 2.2. The key idea is to apply Theorem 5.1. Integrating both sides of (5.3) (with u replaced by y , and $v = \theta y$), using integration by parts, and recalling that $\Theta(t)$ decays exponentially to 0 as $t \rightarrow 0$ or $t \rightarrow T$, noting that $v|_\Sigma = 0$ (and hence $\nabla v = \frac{\partial v}{\partial \nu} \nu$ on Σ), we arrive at

$$\begin{aligned}
 & \mathbb{E} \int_Q \left[c_0 \lambda \Theta \theta^2 (y_t^2 + |\nabla y|^2 + \lambda^2 y^2) + \Theta \left(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v \right)^2 \right] dx dt \\
 (6.1) \quad & \leq \mathbb{E} \int_Q \Theta \theta \left[(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v)(dy_t - \Delta y dt) - \theta \ell_t (dy_t)^2 \right] dx \\
 & \quad + \mathbb{E} \int_\Sigma \Theta \frac{\partial \ell}{\partial \nu} \left| \frac{\partial v}{\partial \nu} \right|^2 d\Gamma dt.
 \end{aligned}$$

By the first equation of system (1.1), we get

$$\begin{aligned}
 (6.2) \quad & \mathbb{E} \int_Q \Theta \theta \left[(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v)(dy_t - \Delta y dt) - \theta \ell_t (dy_t)^2 \right] dx \\
 & = \mathbb{E} \int_Q \Theta \theta \left[(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v)(a_1 y_t + a_2 \cdot \nabla y + a_3 y + f) \right. \\
 & \quad \left. - \theta \ell_t (a_4 y + g)^2 \right] dx dt \\
 & \leq \mathbb{E} \int_Q \Theta \left(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v \right)^2 dx dt \\
 & \quad + C \left\{ \mathbb{E} \int_Q \Theta \theta^2 \left(a_1 y_t + a_2 \cdot \nabla y + a_3 y + f \right)^2 dx dt + \lambda \mathbb{E} \int_Q \Theta \theta^2 (a_4 y + g)^2 dx dt \right\} \\
 & \leq \mathbb{E} \int_Q \Theta \left(-2\ell_t v_t + 2\nabla \ell \cdot \nabla v + \psi v \right)^2 dx dt \\
 & \quad + C \left\{ \mathbb{E} \int_Q \Theta \theta^2 (f^2 + \lambda g^2) dx dt + |a_1|_{L^\infty(0,T;L^\infty(G))}^2 \mathbb{E} \int_Q \Theta \theta^2 y_t^2 dx dt \right. \\
 & \quad + \lambda \left[|a_3|_{L^\infty(0,T;L^n(G))}^2 + |a_4|_{L^\infty(0,T;L^\infty(G))}^2 \right] \mathbb{E} \int_Q \Theta \theta^2 y^2 dx dt \\
 & \quad \left. + \left[|a_2|_{L^\infty(0,T;L^\infty(G;\mathbb{R}^n))}^2 + |a_3|_{L^\infty(0,T;L^n(G))}^2 \right] \mathbb{E} \int_Q \Theta \theta^2 |\nabla y|^2 dx dt \right\}.
 \end{aligned}$$

On the other hand, recalling (2.2), we have

$$\begin{aligned}
 \mathbb{E} \int_{\Sigma} \Theta \left| \frac{\partial \ell}{\partial \nu} \right|^2 d\Gamma dt &= 2\lambda \mathbb{E} \int_{\Sigma} \Theta \theta^2 (x - x_0) \cdot \nu(x) \left| \frac{\partial y}{\partial \nu} \right|^2 d\Gamma dt \\
 (6.3) \quad &\leq 2\lambda \mathbb{E} \int_{\Sigma_0} \Theta \theta^2 (x - x_0) \cdot \nu(x) \left| \frac{\partial y}{\partial \nu} \right|^2 d\Gamma_0 dt \leq C\lambda \mathbb{E} \int_{\Sigma_0} \Theta \theta^2 \left| \frac{\partial y}{\partial \nu} \right|^2 d\Gamma_0 dt.
 \end{aligned}$$

Finally, combining (6.1), (6.2), and (6.3), we arrive at the desired estimate (2.14). This completes the proof of Theorem 2.2. \square

Proof of Theorem 2.1. The proof follows easily from Theorem 2.2 and the energy estimate in Proposition 3.1. We omit the details. \square

7. Application to state observation problem of semilinear stochastic wave equations. This section is devoted to giving an application of the observability inequality in section 2. Fixing two known nonlinear functions $F(\eta, v, \zeta) : \mathbb{R}^1 \times \mathbb{R}^1 \times \mathbb{R}^n \rightarrow \mathbb{R}^1$ and $G(\eta) : \mathbb{R}^1 \rightarrow \mathbb{R}^1$, and unknown initial data $(w_0, w_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, for the semilinear stochastic wave equation

$$(7.1) \quad \begin{cases} dw_t - \Delta w dt = F(w, w_t, \nabla w) dt + G(w) dB(t) & \text{in } Q, \\ w = 0 & \text{on } \Sigma, \\ w(0) = w_0, \quad w_t(0) = w_1 & \text{in } G, \end{cases}$$

we consider the following.

State observation problem: Let $T > 0$ be given. Determine the state $(w(t), w_t(t)) \in L^2(\Omega, \mathcal{F}_t, P; H_0^1(G) \times L^2(G))$ of (7.1) (for $t \in (0, T]$) from the observed boundary data

$$\frac{\partial w}{\partial \nu}(t, x), \quad (t, x) \in \Sigma_0, \quad P\text{-a.s.}$$

We refer the reader to [11, 14] for some studies on the state observation problems for (deterministic) semilinear wave equations. To the best of our knowledge, nothing is published for the stochastic setting. We need the following assumption:

(H) The nonlinear functions $F(\cdot)$ and $G(\cdot)$ satisfy the following:
(1)

$$\begin{aligned}
 |F(\eta_1, v, \zeta) - F(\eta_2, v, \zeta)| + |G(\eta_1) - G(\eta_2)| &\leq L(1 + |\eta_1|^{p-1} + |\eta_2|^{p-1})|\eta_1 - \eta_2| \\
 &\forall \eta_1, \eta_2, v \in \mathbb{R}^1, \zeta \in \mathbb{R}^n
 \end{aligned}$$

with $1 \leq p \leq \frac{n}{n-2}$ if $n \geq 3$; $1 \leq p < \infty$ if $n = 1, 2$;
(2)

$$\begin{aligned}
 |F(\eta, v_1, \zeta_1) - F(\eta, v_2, \zeta_2)| &\leq L(|v_1 - v_2| + |\zeta_1 - \zeta_2|) \\
 &\forall (\eta, v_i, \zeta_i) \in \mathbb{R}^1 \times \mathbb{R}^1 \times \mathbb{R}^n, \quad i = 1, 2, \\
 |F(0, v, \zeta)| &\leq L(|v| + |\zeta|) \quad \forall (v, \zeta) \in \mathbb{R}^1 \times \mathbb{R}^n
 \end{aligned}$$

for some constant $L > 0$;

(3) for any given initial data $(w_0, w_1) \in L^2(\Omega, \mathcal{F}_0, P; H_0^1(G) \times L^2(G))$, (7.1) admits a unique solution $w = w(\cdot; w_0, w_1) \in \mathcal{H}_T$ (the solution of (7.1) is defined similarly to Definition 1.1).

Since we do not introduce any sign condition on the nonlinearities $F(\cdot)$ and $G(\cdot)$, the global existence of a solution of (7.1) is not guaranteed. This is why we need to impose the third assumption in (H). The (global) well-posedness of (7.1) is an interesting but difficult problem, which is beyond the scope of this paper.

It is easy to see that, by the Sobolev embedding theorem, assumption (H) implies that

$$F(z, z_t, \nabla z) \in L^2_{\mathcal{F}}(0, T; L^2(G)), \quad G(z) \in L^2_{\mathcal{F}}(0, T; L^2(G))$$

for any $z \in \mathcal{H}_T$. Thus, thanks to Proposition 3.3, we can define a nonlinear map $\gamma : L^2(\Omega, \mathcal{F}_0, P; H^1_0(G) \times L^2(G)) \rightarrow L^2_{\mathcal{F}}(0, T; L^2(\Gamma_0))$ by

$$\gamma(w_0, w_1) = \frac{\partial w}{\partial \nu} \Big|_{\Sigma_0 \times \Omega},$$

where w is the solution of (7.1). Now we have the following result.

THEOREM 7.1. *Let (H) hold, Γ_0 be given by (2.2), and T satisfy (2.7). Then, for any $(w_0, w_1), (\hat{w}_0, \hat{w}_1) \in L^2(\Omega, \mathcal{F}_0, P; H^1_0(G) \times L^2(G))$ and any $t \in (0, T]$, there exists a constant $\tilde{C} = \tilde{C}(F, G, w_0, \hat{w}_0, w_1, \hat{w}_1) > 0$ such that*

$$\begin{aligned} & |(w(t) - \hat{w}(t), w_t(t) - \hat{w}_t(t))|_{L^2(\Omega, \mathcal{F}_t, P; H^1_0(G) \times L^2(G))} \\ & \leq e^{Ct-1} \tilde{C} |\gamma(w_0, w_1) - \gamma(\hat{w}_0, \hat{w}_1)|_{L^2_{\mathcal{F}}(0, T; L^2(\Gamma_0))}, \end{aligned}$$

where $\hat{w} = \hat{w}(\cdot; \hat{w}_0, \hat{w}_1) \in \mathcal{H}_T$ is the solution of (7.1) with (w_0, w_1) replaced by (\hat{w}_0, \hat{w}_1) .

Remark 7.1. Theorem 7.1 indicates that the state $(w(t), w_t(t)) \in L^2(\Omega, \mathcal{F}_t, P; H^1_0(G) \times L^2(G))$ of (7.1) (for $t \in (0, T]$) can be uniquely determined from the observed boundary data $\frac{\partial w}{\partial \nu}(t, x)|_{\Sigma_0}$, P -a.s., and continuously depend on it.

Proof of Theorem 7.1. Set

$$y = \hat{w} - w.$$

It is easy to see that y is a solution of (1.1) with

$$\begin{cases} a_1 = \int_0^1 \partial_{\eta} F(w + s(\hat{w} - w), w_t, \nabla w) ds, & a_2 = \int_0^1 \partial_{\nu}(\hat{w}, w_t + s(\hat{w}_t - w_t), \nabla w) ds, \\ a_3 = \int_0^1 \partial_{\zeta} F(\hat{w}, \hat{w}_t, \nabla w + s(\nabla \hat{w} - \nabla w)) ds, & a_4 = \int_0^1 \partial_{\eta} G(w + s(\hat{w} - w)) ds. \end{cases}$$

The desired result follows immediately from Remark 2.1. □

Acknowledgments. The author acknowledges Professor Shanjian Tang for stimulating discussion and the anonymous referees for helpful comments.

REFERENCES

[1] V. BARBU, A. RĂSCANU, AND G. TESSITORE, *Carleman estimate and controllability of linear stochastic heat equations*, Appl. Math. Optim., 47 (2003), pp. 97–120.
 [2] C. BARDOS, G. LEBEAU, AND J. RAUCH, *Sharp sufficient conditions for the observation, control, and stabilization of waves from the boundary*, SIAM J. Control Optim., 30 (1992), pp. 1024–1065.

- [3] T. DUYSKAERTS, X. ZHANG, AND E. ZUAZUA, *On the optimality of the observability inequalities for parabolic and hyperbolic systems with potentials*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 25 (2008), pp. 1–41.
- [4] X. FU, *A weighted identity for partial differential operators of second order and its applications*, C. R. Math. Acad. Sci. Paris, 342 (2006), pp. 579–584.
- [5] X. FU, J. YONG, AND X. ZHANG, *Exact controllability for the multidimensional semilinear hyperbolic equations*, SIAM J. Control Optim., 46 (2007), pp. 1578–1614.
- [6] W. GRECKSCH AND C. TUDOR, *Stochastic Evolution Equations: A Hilbert Space Approach*, Math. Res. 85, Akademie-Verlag, Berlin, 1995.
- [7] V. ISAKOV, *Carleman estimates and applications to inverse problems*, Milan J. Math., 72 (2004), pp. 249–271.
- [8] M. KAZEMI AND M. V. KLIBANOV, *Stability estimates for ill-posed Cauchy problem involving hyperbolic equations and inequalities*, Appl. Anal., 50 (1993), pp. 93–102.
- [9] M. M. LAVRENT'EV, V. G. ROMANOV, AND S. P. SHISHAT-SKIĬ, *Ill-Posed Problems of Mathematical Physics and Analysis*, Transl. Math. Monogr. 64, AMS, Providence, RI, 1986.
- [10] J.-L. LIONS, *Contrôlabilité exacte, perturbations et stabilisation de systèmes distribués, Tome 1, Contrôlabilité exacte*, Rech. Math. Appl. 8, Masson, Paris, 1988.
- [11] L. PAN, K. L. TEO, AND X. ZHANG, *State-observation problem for a class of semi-linear hyperbolic systems*, Chinese Ann. Math. Ser. A, 25 (2004), pp. 189–198 (in Chinese); Chinese J. Contemp. Math., 25 (2004), pp. 163–172 (in English).
- [12] S. TANG AND X. ZHANG, *Carleman inequality for backward stochastic parabolic equations with general coefficients*, C. R. Math. Acad. Sci. Paris, 339 (2004), pp. 775–780.
- [13] X. ZHANG, *Explicit observability inequalities for the wave equation with lower order terms by means of Carleman inequalities*, SIAM J. Control Optim., 39 (2000), pp. 812–834.
- [14] X. ZHANG, *Exact Controllability of Semi-linear Distributed Parameter Systems*, Gaodeng Jiaoyu Chubanshe (Higher Education Press), Beijing, China, 2004 (in Chinese).
- [15] X. ZHANG AND E. ZUAZUA, *A sharp observability inequality for Kirchoff plate systems with potentials*, Comput. Appl. Math., 25 (2006), pp. 353–373.
- [16] E. ZUAZUA, *Controllability and observability of partial differential equations: Some results and open problems*, in Handbook of Differential Equations: Evolutionary Equations, Vol. 3, Elsevier Science, Amsterdam, 2006, pp. 527–621.

GLOBAL BEHAVIOR OF SPHERICALLY SYMMETRIC NAVIER–STOKES EQUATIONS WITH DEGENERATE VISCOSITY COEFFICIENTS*

MINGJUN WEI[†], TING ZHANG[†], AND DAOYUAN FANG[†]

Abstract. In this paper, we study a free boundary problem for compressible spherically symmetric Navier–Stokes equations with a gravitational force and degenerate viscosity coefficients. Under certain assumptions that are imposed on the initial data, we obtain the global existence and uniqueness of the weak solution and give some uniform bounds (with respect to time) of the solution. Moreover, we obtain some stabilization rate estimates in L^∞ -norm and weighted H^1 -norm of the solution. The results show that such a system is stable under small perturbations and could be applied to the astrophysics.

Key words. compressible Navier–Stokes equations, density-dependent viscosity, free boundary, asymptotic behavior

AMS subject classifications. 35Q35, 35D05, 76N10

DOI. 10.1137/070681703

1. Introduction. We consider the compressible Navier–Stokes equations with density-dependent viscosity in \mathbb{R}^n ($n \geq 2$), which can be written in Eulerian coordinates as

$$(1.1) \quad \begin{cases} \partial_\tau \rho + \nabla \cdot (\rho \vec{u}) = 0, \\ \partial_\tau (\rho \vec{u}) + \nabla \cdot (\rho \vec{u} \otimes \vec{u}) + \nabla P = \operatorname{div}(\mu(\nabla \vec{u} + \nabla \vec{u}^\top)) + \nabla(\lambda \operatorname{div} \vec{u}) - \rho \vec{f}, \end{cases}$$

in a domain $\{(\vec{\xi}, \tau) | \vec{\xi} \in \Omega_\tau \subset \mathbb{R}^n, \tau > 0\}$, the initial conditions and boundary conditions are

$$(1.2) \quad (\rho, \vec{u})(\vec{\xi}, 0) = (\rho_0, u_0)(\vec{\xi}), \quad \vec{\xi} \in \Omega_0 = \{\vec{\xi} \in \mathbb{R}^n | a < |\vec{\xi}| < b\},$$

$$(1.3) \quad \vec{u}|_{|\vec{\xi}|=a} = 0, \quad \rho|_{\vec{\xi} \in \partial\Omega_\tau \setminus \{|\vec{\xi}|=a\}} = 0,$$

where $0 < a < b < \infty$, $\Omega_\tau = \psi(\Omega_0, \tau)$ and ψ is the flow of \vec{u} :

$$(1.4) \quad \begin{cases} \partial_\tau \psi(\vec{\xi}, \tau) = \vec{u}(\psi(\vec{\xi}, \tau), \tau), \quad \vec{\xi} \in \mathbb{R}^n, \\ \psi(\vec{\xi}, 0) = \vec{\xi}. \end{cases}$$

Here ρ , P , $\vec{u} = (u_1, \dots, u_n)$, and \vec{f} are the density, pressure, velocity, and external force, respectively; $\lambda = \lambda(\rho)$ and $\mu = \mu(\rho)$ are the viscosity coefficients.

*Received by the editors February 3, 2007; accepted for publication (in revised form) May 21, 2008; published electronically September 8, 2008. This work is supported by NSFC 10571158, Zhejiang Provincial NSF of China (Y605076), and China Postdoctoral Science Foundation 20060400335.

<http://www.siam.org/journals/sima/40-3/68170.html>

[†]Department of Mathematics, Zhejiang University, Hangzhou 310027, People’s Republic of China (m.j.wei@126.com, zhangting79@hotmail.com (corresponding author), dyf@zju.edu.cn).

For the initial-boundary value problem (1.1)–(1.3) with the spherically symmetric initial data and external force

$$(\rho, \vec{u})(\vec{\xi}, 0) = \left(\rho_0(r), u_0(r) \frac{\vec{\xi}}{r} \right), \quad \vec{\xi} \in \Omega_0,$$

$$\vec{f} = f(m, r, \tau) \frac{\vec{\xi}}{r}, \quad m(\rho, r) = \int_a^r \rho(s, \tau) s^{n-1} ds, \quad \vec{\xi} \in \Omega_\tau,$$

where $r = |\vec{\xi}| = \sqrt{\xi_1^2 + \dots + \xi_n^2}$, we are looking for spherically symmetric solutions (ρ, \vec{u}) :

$$\rho(\vec{\xi}, \tau) = \rho(r, \tau), \quad \vec{u}(\vec{\xi}, \tau) = u(r, \tau) \frac{\vec{\xi}}{r}, \quad \vec{\xi} \in \Omega_\tau,$$

where $\Omega_\tau = \{\vec{\xi} \in \mathbb{R}^n \mid a < |\vec{\xi}| < b(\tau), b(0) = b, b'(\tau) = u(b(\tau), \tau)\}$. Then $(\rho, u)(r, \tau)$ is determined by

$$(1.5) \quad \begin{cases} \partial_t \rho + \partial_r(\rho u) + \frac{n-1}{r} \rho u = 0, \\ \rho(\partial_t u + u \partial_r u) + \partial_r P = \partial_r [(\lambda + 2\mu)(\partial_r u + \frac{n-1}{r} u)] - 2(n-1) \frac{u}{r} \partial_r \mu - \rho f, \end{cases}$$

where $(r, \tau) \in (a, b(\tau)) \times (0, \infty)$, with the initial data

$$(1.6) \quad (\rho, u)|_{\tau=0} = (\rho_0, u_0)(r), \quad a \leq r \leq b,$$

the boundary conditions

$$(1.7) \quad u|_{r=a} = 0, \quad \rho|_{r=b(\tau)} = 0,$$

where $b(0) = b, b'(\tau) = u(b(\tau), \tau), \tau > 0$.

To simplify the presentation, we consider only the famous polytropic model, i.e., $P(\rho) = A\rho^\gamma$, with $\gamma > 1$ and $A > 0$ being constants. And we assume that the viscosity coefficients μ and λ are proportional to ρ^θ , i.e., $\mu(\rho) = c_1\rho^\theta$ and $\lambda(\rho) = c_2\rho^\theta$, where c_1, c_2 , and θ are three constants.

Additionally, we assume that the external force $f(m, r, \tau)$ satisfies

$$(1.8) \quad f(m, r, \tau) = f_\infty(m, r) + \Delta f(m, r, \tau),$$

for all $m \geq 0, r \geq a$, and $\tau \geq 0$, with

$$(1.9) \quad f_\infty(m, r) = G \frac{M_0 + m}{r^{n-1}}, \quad \Delta f(m, r, \tau) \in C^1(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}_+)$$

$$(1.10) \quad \|\Delta f(\cdot, \cdot, \tau)\|_{L^\infty(\mathbb{R}_+ \times \mathbb{R}_+)} \leq f_1(\tau), \quad \|(\partial_r \Delta f, \partial_\tau \Delta f)(\cdot, \cdot, \tau)\|_{L^\infty(\mathbb{R}_+ \times \mathbb{R}_+)} \leq f_2(\tau),$$

$$(1.11) \quad f_1 \in L^\infty \cap L^2(\mathbb{R}_+), \quad f_2 \in L^2(\mathbb{R}_+),$$

where $\mathbb{R}_+ = [0, \infty)$, $G > 0$ is a constant, $M_0 \geq 0$ is the total mass of the solid core surrounded by the gas, and the perturbation Δf tends to 0 as $\tau \rightarrow \infty$ in some weak sense. If $M_0 = 0$, we ignore the gravitational effect of the solid core. Δf expresses the influence of the outside gravitational force, f_∞ is the precise expression for its own gravitational force and the gravitational force of the solid core, in the astrophysical case (with spherical symmetry).

It is very interesting to study the stabilization problem for the model of the atmosphere of a planet, whose evolution is influenced by the gravitational force. In [22], we considered a simple model that is one-dimensional Navier–Stokes equations. In this paper, we study the stabilization problem of the spherically symmetric model (1.5). We should consider the essential multidimensional model in the future.

Now, we consider the stationary problem, namely

$$(1.12) \quad (P(\rho_\infty))_r = -\rho_\infty f_\infty(m(\rho_\infty, r), r)$$

in an interval $r \in (a, l_\infty)$, with the end l_∞ satisfying

$$(1.13) \quad \rho_\infty(l_\infty) = 0, \int_a^{l_\infty} \rho_\infty r^{n-1} dr = M := \int_a^b \rho_0 r^{n-1} dr.$$

The unknown quantities are the stationary density $\rho_\infty \geq 0$ and free boundary $l_\infty > 0$. If $\gamma > \frac{2n-2}{n}$, from Proposition 2.5, we know that there exists a unique solution (ρ_∞, l_∞) to the stationary system (1.12)–(1.13), satisfying $\rho_\infty(r) \sim (l_\infty^n - r^n)^{\frac{1}{\gamma-1}}$, $(\rho_\infty)_r(r) < 0$, $a < r < l_\infty$, with $l_\infty < +\infty$.

To handle the free boundary problem (1.5)–(1.7), it is convenient to reduce the problem in Eulerian coordinates (r, τ) to the problem in Lagrangian coordinates (x, t) moving with the fluids, via the transformation:

$$(1.14) \quad x = \int_a^r y^{n-1} \rho(y, \tau) dy, \quad t = \tau.$$

Then the fixed boundary $r = a$ and the free boundary $r = b(\tau)$ become

$$x = 0 \quad \text{and} \quad x = \int_a^{b(\tau)} y^{n-1} \rho(y, \tau) dy = \int_a^b y^{n-1} \rho_0(y) dy = M,$$

where M is the total mass initially so that the region $\{(r, \tau) | a \leq r \leq b(\tau), \tau \geq 0\}$ under consideration is transformed into the region $\{(x, t) | 0 \leq x \leq M, t \geq 0\}$, and the function $m(\rho, r)$ becomes x . Under the coordinate transformation (1.14), the equations (1.5)–(1.7) are transformed into

$$(1.15) \quad \begin{cases} \partial_t \rho(x, t) = -\rho^2 \partial_x (r^{n-1} u), \\ \partial_t u(x, t) = r^{n-1} \{ \partial_x [\rho(\lambda + 2\mu) \partial_x (r^{n-1} u) - P] - 2(n-1) \frac{u}{r} \partial_x \mu \} - f(x, r, t), \\ r^n(x, t) = a^n + n \int_0^x \rho^{-1}(y, t) dy, \end{cases}$$

where $(x, t) \in (0, M) \times (0, \infty)$, with the initial data

$$(1.16) \quad (\rho, u)|_{t=0} = (\rho_0, u_0)(x), r|_{t=0} = r_0(x) = \left(a^n + n \int_0^x \rho_0^{-1}(y) dy \right)^{\frac{1}{n}},$$

and the boundary conditions

$$(1.17) \quad u|_{x=0} = 0, \quad \rho|_{x=M} = 0, \quad t > 0.$$

It is standard that if we can solve the problem (1.15)–(1.17), then the free boundary problem (1.1)–(1.3) has a solution.

From (1.12)–(1.13), it is easy to see that $\rho_\infty(x)$ is the solution to the stationary system

$$(1.18) \quad Ar_\infty^{n-1}(\rho_\infty^\gamma)_x = -f_\infty(x, r_\infty), \quad r_\infty^n(x) = a^n + n \int_0^x \rho_\infty^{-1}(y)dy, \quad x \in (0, M),$$

$$(1.19) \quad \rho_\infty(M) = 0.$$

The results in [8, 20] show that the compressible Navier–Stokes system with the constant viscosity coefficient has the singularity at the vacuum. Considering the modified Navier–Stokes system in which the viscosity coefficient depends on the density, Liu, Xin, and Yang [11] proved that such system is local well posed. It is motivated by the physical consideration that, in the derivation of the Navier–Stokes equations from the Boltzmann equation through the Chapman–Enskog expansion to the second order, cf. [6], viscosity is a function of temperature. If we consider the case of isentropic fluids, this dependence is reduced to the dependence on the density function.

Since $n \geq 2$ and the viscosity coefficient μ depends on ρ , the nonlinear term $2(n-1)\frac{1}{r}u\partial_x\mu$ in (1.15)₂ makes the analysis significantly different from the one-dimensional case [11, 15, 19, 21, 22]. When $\mu \geq \underline{\mu} > 0$ and $\rho_0 \geq \underline{\rho} > 0$, authors in [4, 24] obtained the existence, uniqueness, and global behavior of the solution for compressible spherically symmetric Navier–Stokes equations with an external pressure and without the nonlinear term $2(n-1)\frac{1}{r}u\partial_x\mu$. Following the ideas in [24], we can obtain the existence and uniqueness results for the stationary problem in section 2. In this paper, since viscosity coefficients and density will degenerate at the free boundary, we need to use the weighted function $(M-x)$ to control the lower bound of the density in section 3.

Considering the system (1.15)–(1.17) with a general external force, Chen-Zhang obtained the local existence and uniqueness of the solution in [3]. In this paper, when the initial data (ρ_0, u_0, r_0) are close to the stationary state $(\rho_\infty, 0, r_\infty)$, we will obtain some appropriate a priori estimates and prove that the maximum existence time $T^* = \infty$. The difficulty of this problem is to obtain the lower bound of the density ρ . The key ideas are using the classical continuation method and the result of Claim 1. In Claim 1, we want to prove that there is a small positive constant ϵ_1 , such that, for any $T > 0$, if

$$(1.20) \quad I(t) = \|g(\cdot, t) - g_\infty\|_{L^\infty} \leq 2\epsilon_1 \quad \forall t \in [0, T],$$

where $g(x, t) = (M-x)^{-\frac{1}{\gamma}}\rho(x, t)$ and $g_\infty(x) = (M-x)^{-\frac{1}{\gamma}}\rho_\infty(x)$, then

$$(1.21) \quad I(t) \leq \epsilon_1 \quad \forall t \in [0, T].$$

Using the energy method and induction method, we can estimate the weighted L^2 -norm of $g - g_\infty$ in Lemma 3.7. In such a process (see Lemmas 3.5–3.6), we use the weight function $(1+t)^\alpha$ (with $\alpha = -\frac{5}{8}$) to remedy the disadvantage of the nonlinear term $2(n-1)\frac{1}{r}u\partial_x\mu$, and we use the induction method to increase α to $-\epsilon_2$. Then, by the reduction to absurdity, we can finish the proof of Claim 1 in Lemma 3.8.

Our results show that such a system is stable under small perturbations, it does not develop vacuum states or concentration states for all time, and the interface $\partial\Omega_\tau$ propagates with finite speed.

The assumptions on c_1, c_2, θ, γ , and initial data can be stated as follows:

(A1) $\gamma > \frac{2n-2}{n}, \theta \in (0, \gamma-1) \cap (0, \frac{\gamma}{2}]$, c_1 and c_2 satisfy that $c_1 > 0$ and $2c_1 + nc_2 > 0$;

(A2) $N_1(M-x)^{1/\gamma} \leq \rho_0 \leq N_2(M-x)^{1/\gamma}$, with some positive constants $0 < N_1 < N_2$,
and $(M-x)^{1-\frac{\theta}{\gamma}}(\rho_0^\theta)_x \in L^1([0, M])$, $(\rho_0^\gamma)_x \in L^2([0, M])$;

(A3) $u_0 \in L^2([0, M])$, $\rho_0^{\frac{\theta+1}{2}}(u_0)_x \in L^2([0, M])$, $u_0(0) = 0$,

$$(1.22) \quad ((2c_1 + c_2)\rho_0^{\theta+1}(r_0^{n-1}u_0)_x)_x - 2c_1(n-1)\frac{u_0}{r_0}(\rho_0^\theta)_x \in L^2([0, M]).$$

Under the above assumptions (A1)–(A3), we will prove the existence of the global weak solution to the initial-boundary value problem (1.15)–(1.17) in the sense of the following definition.

DEFINITION 1.1. *A pair of functions $(\rho, u, r)(x, t)$ is called a global weak solution to the initial-boundary value problem (1.15)–(1.17) if, for any $T > 0$,*

$$\rho, u \in L^\infty([0, M] \times [0, T]) \cap C^1([0, T]; L^2([0, M])),$$

$$r \in C^1([0, T]; L^\infty([0, M])),$$

$$\rho^{-1}, (r^{n-2}u)_x, (r^{n-1})_x \in L^\infty([0, T]; L^1([0, M])),$$

and

$$\rho^{1+\theta}(r^{n-1}u)_x \in L^\infty([0, M] \times [0, T]) \cap C^{\frac{1}{2}}([0, T]; L^2([0, M])).$$

Furthermore, the following equations hold:

$$\rho_t + \rho^2(r^{n-1}u)_x = 0, \quad \rho(x, 0) = \rho_0(x), \quad a.e.$$

$$r_t = u, \quad r(x, 0) = r_0(x), \quad r^n(x, t) = a^n + n \int_0^x \rho^{-1}(y, t) dy, \quad a.e.$$

$$\int_0^\infty \int_0^M [w\psi_t + (P - \rho(\lambda + 2\mu)(r^{n-1}u)_x)(r^{n-1}\psi)_x + 2(n-1)\mu(r^{n-2}u\psi)_x - f(x, r, t)\psi] dx dt + \int_0^M u_0(x)\psi(x, 0) dx = 0$$

for any test function $\psi(x, t) \in C_0^\infty(\Omega)$, with $\Omega = \{(x, t) \mid 0 < x \leq M, t \geq 0\}$ and

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \int_0^\epsilon u dx = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \int_{1-\epsilon}^1 \rho dx = 0.$$

In what follows, we always use $C(C_i)$ to denote a generic positive constant depending only on γ, θ, f_1, f_2 , and the initial data, independent of the given time T .

We now state the main theorems in this paper.

THEOREM 1.1. *Under the conditions (1.8)–(1.11) and (A1)–(A3), there exists a constant $\epsilon_0 > 0$ such that if*

$$(1.23) \quad \|f_1\|_{L^\infty}^2 + \int_0^\infty (1+t)f_1^2(t) dt \leq \epsilon_0^2$$

and

$$(1.24) \quad \|u_0\|_{L^2} + \|(M-x)^{-\frac{1}{\gamma}}(\rho_0 - \rho_\infty)\|_{L^\infty} \leq \epsilon_0,$$

then the system (1.15)–(1.17) has a unique global weak solution (ρ, u, r) satisfying

$$(1.25) \quad C^{-1}(M-x)^{\frac{1}{\gamma}} \leq \rho(x, t) \leq C(M-x)^{\frac{1}{\gamma}},$$

$$(1.26) \quad r(x, t) \in [a, C],$$

$$(1.27) \quad \int_0^M (M-x)^{1-\frac{\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2 dx \leq C$$

and

$$(1.28) \quad \|u(\cdot, t)\|_{L^\infty} + \|[\rho(r^{n-1}u)_x](\cdot, t)\|_{L^\infty} \leq C$$

for all $t \geq 0$ and $x \in [0, M]$. Furthermore, if $(1+t)^{\frac{2(\gamma+\theta)}{\gamma+\theta+1}} f_1^2(t) \in L^1(\mathbb{R}_+)$, for any $\eta > 0$, we have

$$(1.29) \quad \int_0^M \left\{ u^2 + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \leq C_\eta (1+t)^{-\frac{2(\gamma+\theta)}{\gamma+\theta+1} + \eta},$$

$$(1.30) \quad \int_0^M (\rho^{\theta-1} u^2 + \rho^{\theta+1} u_x^2) dx \leq C_\eta (1+t)^{-\frac{2(\gamma+\theta)}{\gamma+\theta+1} + \eta},$$

$$(1.31) \quad \int_0^M (M-x)^{2-\frac{2\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2 dx \leq C_\eta (1+t)^{-\frac{\gamma+\theta-1}{\gamma+\theta+1} + \eta},$$

$$(1.32) \quad \|\rho^\gamma(\cdot, t) - \rho_\infty^\gamma(\cdot)\|_{L^\infty} \leq C_\eta (1+t)^{-\frac{3\gamma+3\theta-1}{4(\gamma+\theta+1)} + \frac{\eta}{2}},$$

and

$$(1.33) \quad \|\rho^{\frac{\gamma+\theta}{2}}(\cdot, t) - \rho_\infty^{\frac{\gamma+\theta}{2}}(\cdot)\|_{L^\infty}^2 + \|u(\cdot, t)\|_{L^\infty} \leq C_\eta (1+t)^{-\frac{\gamma+\theta}{\gamma+\theta+1} + \frac{\eta}{2}},$$

for all $t \geq 0$, where C_η is a positive constant depending on η .

Remark 1.1. The constant $\epsilon_0 = \epsilon_0(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ is defined in (3.55). There are no smallness assumptions on $\|(M-x)^{1-\frac{\theta}{\gamma}}(\rho_0^\theta)_x\|_{L^1}$ and $\|\rho_0^{\frac{1+\theta}{2}}(u_0)_x\|_{L^2}$.

Remark 1.2. The uniqueness of the solution in Theorems 1.1 or 3.1 means that if (ρ_1, u_1, r_1) and (ρ_2, u_2, r_2) are two solutions to the system (1.15)–(1.17) with the same initial data (ρ_0, u_0, r_0) and satisfy regularity conditions in the theorem, then we have that $(\rho_1, u_1, r_1) = (\rho_2, u_2, r_2)$.

Remark 1.3. In particular, the viscosity of the gas is proportional to the square root of the temperature for the hard sphere model (as pointed out in [15, 21]), and the relation between θ and γ is

$$\theta = \frac{\gamma - 1}{2}.$$

Our condition (A1) covers it. Since the Navier–Stokes system with constant viscosity coefficients has the singularity at the vacuum [8, 20], we assume that $\theta > 0$ in (A1).

Remark 1.4. Considering the no-vacuum system in an exterior domain in \mathbb{R}^3 , Kobayashi and Shibata [9] obtained $\|(\rho - \rho_\infty, u)(\cdot, t)\|_{L^2} \lesssim (1+t)^{-\frac{3}{4}}$ and $\|(\rho - \rho_\infty, u)(\cdot, t)\|_{L^\infty} \lesssim (1+t)^{-\frac{3}{2}}$ when ρ_∞ is a positive constant. Considering the no-vacuum system in \mathbb{R}^n ($n \geq 3$), Ukai, Yang, and Zhao [18] obtained $\|(\rho - \rho_\infty, u)(\cdot, t)\|_{L^2 \cap L^\infty} \leq C(1+t)^{-\frac{n}{4} + \epsilon}$ when ρ_∞ is close to a positive constant. Considering the one-dimensional system with a degenerate viscosity coefficient, we [22] obtained $\|(\rho^\gamma - \rho_\infty^\gamma, u)(\cdot, t)\|_{L^\infty} \leq$

$C(1+t)^{-\frac{1}{2}}$ when the external force f_∞ is close to a positive constant N_0 and the stationary density ρ_∞ is close to $(\frac{N_0(M-x)}{A})^{\frac{1}{\gamma}}$. Since $\frac{\gamma+\theta}{\gamma+\theta+1} > \frac{1}{2}$ and $\frac{3\gamma+3\theta-1}{4(\gamma+\theta+1)} > \frac{1}{2}$ (if $\gamma + \theta > 3$), it is easy to see that our results in Theorem 1.1 are better than that in [22]. Using similar arguments as that in Theorem 1.1, we can also obtain similar results in the one-dimensional case, which are better than that in [22]. For example, the stabilization rate estimate $\|(\rho^\gamma - \rho_\infty^\gamma, u)(\cdot, t)\|_{L^\infty} \leq C(1+t)^{-\frac{1}{2}}$ in [22] can be replaced by $\|(\rho^\gamma - \rho_\infty^\gamma, u)(\cdot, t)\|_{L^\infty} \leq C_\eta(1+t)^{-\frac{\gamma+\theta}{\gamma+\theta+1} + \frac{\eta}{2}}$, $\|\rho^{\frac{\gamma+\theta}{2}}(\cdot, t) - \rho_\infty^{\frac{\gamma+\theta}{2}}(\cdot)\|_{L^\infty} \leq C_\eta(1+t)^{-\frac{\gamma+\theta}{2(\gamma+\theta+1)} + \frac{\eta}{4}}$, $t \geq 0$ for any $\eta > 0$.

Remark 1.5. Since the information of the dimension n mainly appear in the index of the radii r , and we consider only the system with a solid core $r \geq a > 0$ in this paper; our results cannot show the effect of the dimension. In [23], we studied the global behavior of the solution to a similar problem with a positive external pressure and without a solid core, and we obtained the stabilization rate estimates for the solution of exponential type. The admissible range of the parameter $\frac{\lambda}{\mu}$ depends on the dimension n in [23]. We will study the system without a solid core $r \geq 0$ and with degenerate viscosity coefficients in the future, and we will guess that stabilization rate estimates of the solution cannot be better than that in Theorem 1.1.

THEOREM 1.2 (continuous dependence). *For each $i = 1, 2$, let (ρ_i, u_i, r_i) be the solution to the system (1.15)–(1.17) with the initial data $(\rho_{0i}, u_{0i}, r_{0i})$, which satisfies regularity conditions in Theorem 1.1. Then, we have*

$$\begin{aligned} & \int_0^M ((u_1 - u_2)^2 + \rho_1^{1-\theta} \rho_2^{2\theta-4} (\rho_1 - \rho_2)^2 + \rho_1^\theta \rho_2^{-1} (r_1 - r_2)^2) dx \\ & \leq C e^{Ct} \int_0^M ((u_{01} - u_{02})^2 + \rho_{01}^{1-\theta} \rho_{02}^{2\theta-4} (\rho_{01} - \rho_{02})^2 + \rho_{01}^\theta \rho_{02}^{-1} (r_{01} - r_{02})^2) dx \end{aligned}$$

for all $t \geq 0$.

Remark 1.6. Using similar arguments as that in [3], we can easily obtain such continuous dependence of the solution on the initial data and omit the details.

Remark 1.7. If we ignore the influence of self-gravitation, i.e., $f_\infty(m, r) = G \frac{M_0}{r^{n-1}}$, with $M_0 > 0$, then we can also obtain the same results in Theorems 1.1–1.2.

We now briefly review the previous works in this direction. For the related free boundary problem of one-dimensional isentropic fluids with density-dependent viscosity (like $\mu(\rho) = c\rho^\theta$), see [11, 15, 19, 21, 22] and the references therein. For the related stabilization rate estimates of the one-dimensional free boundary problem, see [5, 12, 16, 22], etc. For the spherically symmetric solutions of the Navier–Stokes equations with a free boundary, see [2, 4, 13, 14, 23, 24], etc. Also see Bresch and Desjardins [1], Lions [10], and Vaigant and Kazhikhov [17] for multidimensional isentropic fluids.

The rest of this paper is organized as follows. First, we obtain the existence and the uniqueness of the solution to the stationary problem in section 2. In section 3, we will prove some a priori estimates and extend the local solution in [3] to the global solution in time. In section 4, we obtain the stabilization rate estimates of the solution.

2. The stationary problem. Zlotnik and Ducomet [24] obtained the existence of the positive solution to the stationary problem with a positive external pressure. Using similar arguments as that in [24], we can obtain the following results for the stationary problem without an external pressure. We start with a proof of the existence of a nonnegative solution to the Lagrangian stationary problem.

PROPOSITION 2.1. *If*

$$(2.1) \quad \gamma > \frac{2n-2}{n}$$

or

$$(2.2) \quad \gamma = \frac{2n-2}{n} \quad \text{and} \quad \left(\frac{n\gamma}{\gamma-1} M^{\frac{\gamma-1}{\gamma}} \right)^{\frac{2n-2}{n}} < \frac{G}{A} \left(\frac{M}{2} + M_0 \right),$$

or

$$(2.3) \quad n > 2, \quad 1 < \gamma < \frac{2n-2}{n} \quad \text{and} \quad \delta_6^\gamma \left(a^n + \frac{n\gamma}{\delta_6(\gamma-1)} M^{\frac{\gamma-1}{\gamma}} \right)^{\frac{2n-2}{n}} \leq \frac{G}{A} \left(\frac{M}{2} + M_0 \right),$$

where $\delta_6 = a^{-n} \left(1 - \frac{\gamma n}{2n-2} \right) \frac{2n-2}{\gamma-1} M^{\frac{\gamma-1}{\gamma}}$, then the Lagrangian stationary problem (1.18)–(1.19) has a nonnegative solution $\rho_\infty \in W^{1,\beta}([0, M])$ satisfying $C^{-1}(M-x)^{\frac{1}{\gamma}} \leq \rho_\infty(x) \leq C(M-x)^{\frac{1}{\gamma}}$, where $\beta \in [1, \frac{\gamma}{\gamma-1}]$ is a constant.

Proof. We introduce the nonlinear operator

$$I : K \rightarrow W^{1,\beta}([0, M]),$$

where $K = \{f \in C([0, M]) \mid f \geq 0, \|\frac{(M-x)^{\frac{1}{\gamma}}}{f(x)}\|_{L^\infty} < \infty, \|\frac{f(x)}{(M-x)^{\frac{1}{\gamma}}}\|_{L^\infty} < \infty\}$, by setting

$$I(f)(x) = \left(\frac{\int_x^M G \frac{M_0+y}{r_f^{2n-2}(y)} dy}{A} \right)^{\frac{1}{\gamma}}, \quad \text{with } r_f^n(x) = a^n + n \int_0^x f^{-1}(y) dy, \quad x \in [0, M].$$

We can restate the problem (1.18)–(1.19) as the fixed-point problem

$$(2.4) \quad \rho_\infty = I(\rho_\infty).$$

For all $f \in K_{\delta_1, \delta_2} = \{f \in K \mid \delta_1(M-x)^{\frac{1}{\gamma}} \leq f(x) \leq \delta_2(M-x)^{\frac{1}{\gamma}}\}$, with $0 < \delta_1 \leq \delta_2 < \infty$, we have

$$a^n \leq r_f^n(x) \leq a^n + \frac{n\gamma}{\delta_1(\gamma-1)} M^{\frac{\gamma-1}{\gamma}} := B^n$$

and

$$\left(\frac{G(\frac{M}{2} + M_0)}{AB^{2n-2}} \right)^{\frac{1}{\gamma}} (M-x)^{\frac{1}{\gamma}} \leq I(f)(x) \leq \left(\frac{G(M + M_0)}{Aa^{2n-2}} \right)^{\frac{1}{\gamma}} (M-x)^{\frac{1}{\gamma}}, \quad x \in [0, M].$$

If $\gamma > \frac{2n-2}{n}$, then $I(K_{\delta_3, \delta_4}) \subset K_{\delta_3, \delta_4}$, where $\delta_4 = \left(\frac{G(M+M_0)}{Aa^{2n-2}} \right)^{\frac{1}{\gamma}}$ and δ_3 is a positive constant satisfying $\delta_3^\gamma \left(a^n + \frac{n\gamma}{\delta_3(\gamma-1)} M^{\frac{\gamma-1}{\gamma}} \right)^{\frac{2n-2}{n}} \leq \frac{G}{A} \left(\frac{M}{2} + M_0 \right)$. And one can immediately verify that I is a compact operator on K_{δ_3, δ_4} . Since K_{δ_3, δ_4} is a convex closed bounded nonempty subset of $C([0, M])$, the problem (2.4) has a solution $\rho \in K_{\delta_3, \delta_4}$ by Schauder's fixed-point theorem.

Similarly, if $\gamma = \frac{2n-2}{n}$ and $(\frac{n\gamma}{\gamma-1}M^{\frac{\gamma-1}{\gamma}})^{\frac{2n-2}{n}} < \frac{G}{A}(\frac{M}{2} + M_0)$, then $I(K_{\delta_5, \delta_4}) \subset K_{\delta_5, \delta_4}$, where $\delta_5 = a^{-n}[(\frac{G}{A}(\frac{M}{2} + M_0))^{\frac{n}{2n-2}} - \frac{n\gamma}{\gamma-1}M^{\frac{\gamma-1}{\gamma}}]$, and problem (2.4) has a solution $\rho \in K_{\delta_5, \delta_4}$.

Similarly, if $n > 2$, $1 < \gamma < \frac{2n-2}{n}$ and $\delta_6^\gamma(a^n + \frac{n\gamma}{\delta_6(\gamma-1)}M^{\frac{\gamma-1}{\gamma}})^{\frac{2n-2}{n}} \leq \frac{G}{A}(\frac{M}{2} + M_0)$, then $I(K_{\delta_6, \delta_4}) \subset K_{\delta_6, \delta_4}$, and problem (2.4) has a solution $\rho \in K_{\delta_6, \delta_4}$. \square

Similar to [24], we say a stationary solution $(\rho_\infty, r_\infty^n)$ is *statically stable* if

$$(2.5) \quad \begin{aligned} J[W] &:= \int_0^M (\gamma A \rho_\infty^{1+\gamma} W_x^2 - (2n-2)G(M_0+x)r_\infty^{2-3n}W^2) dx \\ &\geq \delta_7 \int_0^M ((M-x)^{\frac{1+\gamma}{\gamma}} W_x^2 + W^2) dx \end{aligned}$$

for some $\delta_7 > 0$ and all $W \in K_1$,

$$K_1 = \left\{ f \in C([0, M]) \mid f(0) = 0, \|(M-x)^{\frac{1}{\gamma}} f'(x)\|_{L^\infty} < \infty, \left\| \frac{1}{(M-x)^{\frac{1}{\gamma}} f'(x)} \right\|_{L^\infty} < \infty \right\}.$$

Now, the static potential energy takes the following form:

$$(2.6) \quad S[V] = \int_0^M \left(\frac{A}{\gamma-1} (V_x)^{1-\gamma} + \int_{\frac{a^n}{n}}^V G(M_0+x)(nh)^{\frac{2-2n}{n}} dh \right) dx.$$

We call

$$V \in K_2 = \left\{ f \in C([0, M]) \mid f(0) = \frac{a^n}{n}, \|(M-x)^{\frac{1}{\gamma}} f'(x)\|_{L^\infty} < \infty, \left\| \frac{1}{(M-x)^{\frac{1}{\gamma}} f'(x)} \right\|_{L^\infty} < \infty \right\}$$

is a point of *local quadratic minimum* of S if

$$(2.7) \quad S[V+W] - S[V] \geq \delta_8 \int_0^M ((M-x)^{\frac{1+\gamma}{\gamma}} W_x^2 + W^2) dx,$$

for all $W \in K_1$ and $\|(M-x)^{\frac{1}{\gamma}} W_x\|_{L^\infty([0, M])} + \|W\|_{L^\infty} \leq \delta_9$, for some $\delta_8 > 0$ and $\delta_9 > 0$.

PROPOSITION 2.2. *If $\gamma > \frac{2n-2}{n}$ and ρ_∞ is a solution of the problem (1.18)–(1.19) satisfying $\rho_\infty \in W^{1,\beta}([0, M])$ and $C^{-1}(M-x)^{\frac{1}{\gamma}} \leq \rho_\infty(x) \leq C(M-x)^{\frac{1}{\gamma}}$, then we have that (2.5) and (2.7) hold, with $V = V_\infty = \frac{r_\infty^n}{n}$.*

Proof. From $r_\infty \geq a$, $(A\rho_\infty^\gamma)_x = -G\frac{M_0+x}{r_\infty^{\frac{2n-2}{n}}}$ and $(r_\infty^n)_x = n\rho_\infty^{-1}$, using integration by parts, we have

$$\begin{aligned} J[W] &= \int_0^M (\gamma A \rho_\infty^{1+\gamma} W_x^2 + (2n-2)A(\rho_\infty^\gamma)_x r_\infty^{-n} W^2) dx \\ &= \int_0^M (\gamma A \rho_\infty^{1+\gamma} W_x^2 - 2(2n-2)A\rho_\infty^\gamma r_\infty^{-n} W W_x \\ &\quad + n(2n-2)A\rho_\infty^{\gamma-1} r_\infty^{-2n} W^2) dx \text{ for all } W \in K_1. \end{aligned}$$

If $\gamma > \frac{2n-2}{n}$, we have

$$J[W] \geq C^{-1} \int_0^M \left((M-x)^{\frac{1+\gamma}{\gamma}} W_x^2 + (M-x)^{\frac{\gamma-1}{\gamma}} W^2 \right) dx.$$

From $r_\infty \geq a$ and $(A\rho_\infty^\gamma)_x = -G\frac{M_0+x}{r_\infty^{2n-2}}$, using integrating by parts and the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \int_0^M W^2 dx &\leq C \int_0^{\frac{M}{2}} (M-x)^{1-\frac{1}{\gamma}} W^2 dx + C \int_0^M G \frac{M_0+x}{r_\infty^{2n-2}} W^2 dx \\ &= C \int_0^{\frac{M}{2}} (M-x)^{1-\frac{1}{\gamma}} W^2 dx - C \int_0^M A(\rho_\infty^\gamma)_x W^2 dx \\ (2.8) \qquad &\leq C \int_0^M \left((M-x)^{\frac{1+\gamma}{\gamma}} W_x^2 + (M-x)^{\frac{\gamma-1}{\gamma}} W^2 \right) dx, \end{aligned}$$

then we can immediately obtain (2.5), with some $\delta_7 = \delta_7(G, A, M_0, M, a, n, \gamma)$.

Similarly, we obtain

$$\begin{aligned} &S[V_\infty + W] - S[V_\infty] \\ &= \frac{1}{2} \int_0^M \left\{ A[\gamma + O(|(M-x)^{\frac{1}{\gamma}} W_x|)] \rho_\infty^{1+\gamma} W_x^2 \right. \\ &\quad \left. - [2n-2 + O(|W|)] G(M_0+x) r_\infty^{2-3n} W^2 \right\} dx \\ &\geq \frac{\delta_7}{2} \int_0^M \left((M-x)^{\frac{1+\gamma}{\gamma}} W_x^2 + W^2 \right) dx \\ &\quad + \int_0^M \left\{ O(|(M-x)^{\frac{1}{\gamma}} W_x|) \rho_\infty^{1+\gamma} W_x^2 + O(|W|) G(M_0+x) r_\infty^{2-3n} W^2 \right\} dx \end{aligned}$$

for all $W \in K_1$. Here, $O(d)$ means $O(d) \rightarrow 0$ as $d \rightarrow 0$. If $\gamma > \frac{2n-2}{n}$, choosing δ_9 small enough, we can immediately obtain (2.7), with some $(\delta_8, \delta_9)(G, A, M_0, M, a, n, \gamma)$. \square

Using a similar argument as that in Proposition 2.2, we could obtain the following uniqueness result.

PROPOSITION 2.3. *Let ρ_∞ be a solution obtained in Proposition 2.1, and ρ_2 be another solution of the problem (1.18)–(1.19) satisfying $\rho_2 \in W^{1,\beta}([0, M])$ and $C^{-1}(M-x)^{\frac{1}{\gamma}} \leq \rho_2(x) \leq C(M-x)^{\frac{1}{\gamma}}$. If $\gamma > \frac{2n-2}{n}$ and $\|(M-x)^{-\frac{1}{\gamma}}(\rho_\infty - \rho_2)(x)\|_{L^\infty} \leq \delta_{10}$ with a small enough positive constant δ_{10} , then we have that $\rho_\infty(x) = \rho_2(x)$, a.e. $x \in [0, M]$.*

Proof. From (1.18)–(1.19), we have

$$\begin{aligned} A\rho_\infty^\gamma(x) &= \int_x^M G \frac{M_0+y}{r_\infty^{2n-2}} dy, \quad r_\infty^n(x) = a^n + n \int_0^x \rho_\infty^{-1}(y) dy, \\ A\rho_2^\gamma(x) &= \int_x^M G \frac{M_0+y}{r_2^{2n-2}} dy, \quad r_2^n(x) = a^n + n \int_0^x \rho_2^{-1}(y) dy, \end{aligned}$$

and

$$A(\rho_\infty^\gamma - \rho_2^\gamma) = \int_x^M G(M_0+y) (r_\infty^{2-2n} - r_2^{2-2n}) dy.$$

Multiplying the above equality by $(\rho_\infty^{-1} - \rho_2^{-1})$, integrating over $[0, M]$, and using the fact that $\int_0^M n(\rho_\infty^{-1} - \rho_2^{-1})(x) \int_x^M g(y) dy dx = \int_0^M g(r_\infty^n - r_2^n) dx$, we obtain

$$\begin{aligned} 0 &= \int_0^M \left\{ A(\rho_\infty^\gamma - \rho_2^\gamma)(\rho_\infty^{-1} - \rho_2^{-1}) - G \frac{(M_0 + x)}{n} (r_\infty^{2-2n} - r_2^{2-2n})(r_\infty^n - r_2^n) \right\} dx \\ &= \int_0^M \left\{ -A[\gamma + O(|(M-x)^{-\frac{1}{\gamma}}(\rho_\infty - \rho_2)|)] \rho_\infty^{1+\gamma} (\rho_\infty^{-1} - \rho_2^{-1})^2 \right. \\ &\quad \left. - \frac{A}{n^2} [2n - 2 + O(|r_\infty^n - r_2^n|)] (\rho_\infty^\gamma)_x r_\infty^{-n} (r_\infty^n - r_2^n)^2 \right\} dx \\ &\leq -C^{-1} \int_0^M \left((M-x)^{\frac{1+\gamma}{\gamma}} (\rho_\infty^{-1} - \rho_2^{-1})^2 + (r_\infty^n - r_2^n)^2 \right) dx, \end{aligned}$$

when $\gamma > \frac{2n-2}{n}$ and δ_{10} is small enough. Thus, we can immediately obtain that $\rho_\infty = \rho_2$. \square

Now, we shall use the shooting method to prove the uniqueness of the solution $\rho_\infty \in K$.

PROPOSITION 2.4. *Under the assumption (2.1), the Lagrangian stationary problem (1.18)–(1.19) has a unique solution $\rho_\infty \in K$.*

Proof. We consider the Cauchy problem

$$(2.9) \quad (A\rho^\gamma)_x = -G(M_0 + x)(nV)^{\frac{2-2n}{n}}, \quad (V)_x = \rho^{-1}, \quad x \in (0, M),$$

$$(2.10) \quad \rho|_{x=0} = \sigma, \quad V|_{x=0} = \frac{\sigma^n}{n}$$

for the unknown functions $\rho(\sigma, x)$ and $V(\sigma, x)$, where $\sigma > 0$ is the shooting parameter. Thus, for each $\sigma > 0$, using the classical ODE theory, there exists a unique solution to the problem (2.9)–(2.10) satisfying $\rho(\sigma, x) > 0$ for $x \in [0, M_\sigma]$, where either $\rho|_{x=M_\sigma} = 0$ and $M_\sigma \in (0, M)$ or $M_\sigma = M$.

Clearly, if $\rho_\infty \in K$ is a solution to the problem (1.18)–(1.19), then ρ_∞ satisfies (2.9)–(2.10) for some $\sigma_0 > 0$, and $M_{\sigma_0} = M$. We will show that it is possible only for one value of σ . Using similar arguments as that in [7, section V.3], we obtain that $(\partial_\sigma \rho^\gamma, \partial_\sigma V)$ is well defined and satisfies the linear Cauchy problem

$$(2.11) \quad A(\partial_\sigma \rho^\gamma)_x = (2n - 2)G(M_0 + x)(nV)^{\frac{2-3n}{n}} \partial_\sigma V, \quad (\partial_\sigma V)_x = -\frac{1}{\gamma} \rho^{-\gamma-1} \partial_\sigma \rho^\gamma, \quad x \in [0, M_\sigma),$$

$$(2.12) \quad \partial_\sigma \rho^\gamma|_{x=0} = 1, \quad \partial_\sigma V|_{x=0} = 0.$$

It is easy to see that

$$\partial_\sigma \rho^\gamma > 0, \quad (\partial_\sigma V)_x < 0, \quad \text{and} \quad \partial_\sigma V < 0$$

hold on $(0, M_4)$, where either $\partial_\sigma \rho^\gamma|_{x=M_4} = 0$ and $M_4 \in (0, M_\sigma)$ or $M_4 = M_\sigma$. We claim that only $M_4 = M_\sigma$ can occur.

Assume that $M_4 \in (0, M_\sigma)$. Letting $\phi = A\rho^\gamma(\partial_\sigma V)_x + \frac{n}{2n-2} A\partial_\sigma \rho^\gamma(V)_x$, from (2.9) and (2.11), we have

$$\int_0^{M_4} \phi dx = \left\{ A\rho^\gamma \partial_\sigma V + \frac{n}{2n-2} A\partial_\sigma \rho^\gamma V \right\} \Big|_0^{M_4}.$$

By the estimates $\rho(\sigma, M_4) > 0$, $\partial_\sigma \rho^\gamma|_{x=M_4} = 0$, $\partial_\sigma V|_{x=M_4} < 0$ and the initial conditions (2.10) and (2.12), we get

$$\int_0^{M_4} \phi dx < 0.$$

On the other hand, from (2.9) and (2.11), we have

$$\phi = A\rho^{-1}\partial_\sigma\rho^\gamma\left(\frac{n}{2n-2} - \frac{1}{\gamma}\right) > 0, \quad x \in (0, M_4).$$

It is a contradiction.

Thus, we obtain

$$\rho(\sigma, x) > 0, \quad \partial_\sigma \rho(\sigma, x) > 0, \quad x \in (0, M_\sigma),$$

and M_σ is nondecreasing on $\sigma \in (0, \infty)$. Therefore, for each fixed point $x \in [0, M_b)$, the function $\rho(\sigma, x)$ is strictly increasing on $\sigma \geq b$.

If there exists $\sigma_1 \neq \sigma_0$ such that $M_{\sigma_1} = M_{\sigma_0} = M$ and $\rho(\sigma_1, x) \in K$, then there exists $\min\{\sigma_0, \sigma_1\} < \sigma_2 < \max\{\sigma_0, \sigma_1\}$ such that $0 < \|(M-x)^{-\frac{1}{\gamma}}(\rho(\sigma_2, x) - \rho(\sigma_0, x))\|_{L^\infty} \leq \delta_{10}$. From Proposition 2.3, we have that $\rho(\sigma_2, x) = \rho(\sigma_0, x) = \rho_\infty(x)$, which is a contradiction. Thus, we finish the proof of Proposition 2.4. \square

Using the properties of the transformation (1.14) and Propositions 2.1–2.4, we can immediately obtain the following proposition.

PROPOSITION 2.5. *Under the assumption (2.1), the Eulerian stationary problem (1.12)–(1.13) has a unique solution (ρ_∞, l_∞) satisfying $\rho_\infty(r) \sim (l_\infty^n - r^n)^{\frac{1}{\gamma-1}}$, $(\rho_\infty)_r < 0$, $a < r < l_\infty$, with $l_\infty < +\infty$.*

Remark 2.1. The uniqueness of the solution in Proposition 2.5 means that if $(\rho_{\infty 1}, l_{\infty 1})$ and $(\rho_{\infty 2}, l_{\infty 2})$ are two solutions to the Eulerian stationary problem (1.12)–(1.13) with the same total mass M , and satisfy $\rho_{\infty i}(r) \sim (l_{\infty i}^n - r^n)^{\frac{1}{\gamma-1}}$, $i = 1, 2$, then we have that $(\rho_{\infty 1}, l_{\infty 1}) = (\rho_{\infty 2}, l_{\infty 2})$.

3. Global existence. Using similar arguments as that in [3], we obtain the following local existence and uniqueness result and omit the proof.

THEOREM 3.1 (local result). *Under the assumptions in Theorem 1.1, there is a positive constant $T_1 > 0$ such that the free boundary problem (1.15)–(1.17) admits a unique weak solution $(\rho, u, r)(x, t)$ on $[0, M] \times [0, T_1]$ in the sense that*

$$\rho(x, t), u(x, t), r(x, t) \in L^\infty([0, M] \times [0, T_1]) \cap C^1([0, T_1]; L^2([0, M])),$$

$$\rho^{\theta+1} \partial_x (r^{n-1} u) \in L^\infty([0, M] \times [0, T_1]) \cap C^{\frac{1}{2}}([0, T_1]; L^2([0, M])),$$

$$\partial_x r^{n-1}, \partial_x (r^{n-2} u) \in L^\infty([0, T_1], L^1([0, M])),$$

and the following equations hold:

$$\partial_t \rho = -\rho^2 \partial_x (r^{n-1} u), \quad \rho(x, 0) = \rho_0,$$

$$(3.1) \quad \partial_t r(x, t) = u(x, t), \quad r^n(x, t) = a^n + n \int_0^x \rho^{-1}(y, t) dy,$$

$$(3.2) \quad (r^\beta (\rho^\theta)_x)_t = -\frac{\theta r^{1+\beta-n}}{2c_1 + c_2} u_t - \frac{\theta}{2c_1 + c_2} (Ar^\beta (\rho^\gamma)_x + r^{1+\beta-n} f),$$

$$(3.3) \quad (2c_1 + c_2)\rho^{1+\theta}(r^{n-1}u)_x \\ = A\rho^\gamma + 2c_1(n-1)\rho^\theta \frac{u}{r} + \int_x^M \left\{ -\frac{u_t}{r^{n-1}} + 2c_1(n-1)\rho^\theta \left(\frac{u}{r}\right)_x - \frac{f}{r^{n-1}} \right\} dy,$$

for almost all $x \in [0, M]$, any $t \in [0, T_1]$, where $\beta = \frac{2(n-1)c_1\theta}{2c_1+c_2}$,

$$(3.4) \quad \int_0^\infty \int_0^M [u\psi_t + (P - \rho(\lambda + 2\mu)(r^{n-1}u)_x)(r^{n-1}\psi)_x \\ + 2(n-1)\mu(r^{n-2}u\psi)_x - f(x, r, t)\psi] dx dt + \int_0^M u_0(x)\psi(x, 0) dx = 0$$

for any test function $\psi(x, t) \in C_0^\infty((0, M) \times [0, T_1])$. Furthermore, we have

$$(3.5) \quad \frac{N_1}{3}(1-x)^{\frac{1}{\gamma}} \leq \rho(x, t) \leq 3N_2(1-x)^{\frac{1}{\gamma}}, \quad (x, t) \in [0, M] \times [0, T_1],$$

$$(3.6) \quad (M-x)^{-\frac{1}{\gamma}}\rho(x, t) \in C([0, T_1]; L^\infty([0, M])), \\ (M-x)^{\frac{\gamma-\theta}{2\gamma}}(\rho^\theta)_x, \rho_t, u_t \in L^\infty([0, T_1]; L^2([0, M])),$$

and

$$\rho^{\frac{\theta+1}{2}}u_{xt} \in L^2([0, M] \times [0, T_1]), \quad \rho\partial_x u \in L^\infty([0, M] \times [0, T_1]).$$

Remark 3.1. From (1.15)₁, (3.5), and $\rho\partial_x u \in L_{t,x}^\infty$, we have that $(M-x)^{-\frac{1}{\gamma}}\partial_t \rho \in L_{t,x}^\infty$. Thus, (3.6) holds.

Assume the maximum existence time of the weak solution in Theorem 3.1 is T^* . In this section, under the small assumptions on the initial data, we will obtain the following a priori estimates and prove that $T^* = \infty$. In the following, we may assume that $(\rho, u, r)(x, t)$ is suitably smooth since the following estimates are valid for the solutions with the regularities indicated in Theorem 3.1 by using the Friedrichs mollifier.

From (1.9), (1.18), and Proposition 2.1, we could easily obtain the following lemma.

LEMMA 3.1. *Under the assumptions of Theorem 1.1, we have*

$$(3.7) \quad A\rho_\infty^\gamma(x) = \int_x^M G \frac{M_0 + y}{r_\infty^{2n-2}} dy,$$

$$(3.8) \quad C^{-1}(M-x)^{\frac{1}{\gamma}} \leq \rho_\infty \leq C(M-x)^{\frac{1}{\gamma}}, \quad r_\infty(x) \in [a, C],$$

and

$$(3.9) \quad \frac{d}{dx}(A\rho_\infty^\gamma(x)) = -G \frac{M_0 + x}{r_\infty^{2n-2}}, \quad C^{-1} \leq (M-x)^{1-\frac{1}{\gamma}} \frac{d}{dx}\rho_\infty(x) \leq C$$

for all $x \in [0, M]$.

LEMMA 3.2. *Under the assumptions of Theorem 1.1, we have*

$$(3.10) \quad \frac{d}{dt} \int_0^M \left(\frac{1}{2}u^2 + \frac{A\rho^{\gamma-1}}{\gamma-1} + \int_a^r G \frac{M_0 + x}{s^{n-1}} ds \right) dx \\ + \int_0^M \left\{ \left(\frac{2}{n}c_1 + c_2 \right) \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 + \frac{2(n-1)}{n}c_1\rho^{1+\theta} \left(r^{n-1}u_x - \frac{u}{r\rho} \right)^2 \right\} dx \\ = - \int_0^M \Delta f u dx, \quad t \in [0, T^*].$$

Proof. Multiplying (1.15)₂ by u , integrating the resulting equation over $[0, M]$, and using integration by parts and the boundary conditions (1.17), we obtain

$$(3.11) \quad \begin{aligned} & \frac{d}{dt} \int_0^M \frac{1}{2} u^2 dx + \int_0^M \{ (2c_1 + c_2) \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 - 2c_1(n-1) \rho^\theta \partial_x(r^{n-2}u^2) \} dx \\ &= \int_0^M A \rho^\gamma \partial_x(r^{n-1}u) dx - \int_0^M f u dx. \end{aligned}$$

From (1.15), we have

$$(3.12) \quad \int_0^M A \rho^\gamma \partial_x(r^{n-1}u) dx = -\frac{d}{dt} \int_0^M \frac{A}{\gamma-1} \rho^{\gamma-1} dx,$$

$$(3.13) \quad -\int_0^M f u dx = -\frac{d}{dt} \int_0^M \int_a^r G \frac{M_0+x}{s^{n-1}} ds dx - \int_0^M \Delta f u dx,$$

and

$$(3.14) \quad \begin{aligned} & (2c_1 + c_2) \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 - 2c_1(n-1) \rho^\theta \partial_x(r^{n-2}u^2) \\ &= \left(\frac{2}{n} c_1 + c_2 \right) \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 + \frac{2(n-1)}{n} c_1 \rho^{1+\theta} \left(r^{n-1} u_x - \frac{u}{r \rho} \right)^2. \end{aligned}$$

From (3.11)–(3.14), we can immediately obtain (3.10). \square

Now, using the classical continuation method, we will obtain the estimate of $\|(M-x)^{-\frac{1}{\gamma}}(\rho - \rho_\infty)\|_{L^\infty}$.

CLAIM 1. *Under the assumptions of Theorem 1.1, for any $T \in (0, T^*)$, if*

$$(3.15) \quad I(t) = \|g(\cdot, t) - g_\infty\|_{L^\infty} \leq 2\epsilon_1 \quad \forall t \in [0, T],$$

where $g(x, t) = (M-x)^{-\frac{1}{\gamma}} \rho(x, t)$ and $g_\infty(x) = (M-x)^{-\frac{1}{\gamma}} \rho_\infty(x)$, then

$$(3.16) \quad I(t) \leq \epsilon_1 \quad \forall t \in [0, T],$$

where

$$\epsilon_1 = \epsilon_0 + C_{10} \epsilon_0^{\frac{\theta}{\gamma} 2^{-N_4-1}},$$

and positive constants N_4 and C_{10} are defined in Lemmas 3.6 and 3.8, respectively.

Using the results in Lemmas 3.3–3.8, we can finish the proof of Claim 1.

LEMMA 3.3. *Under the assumptions of Theorem 1.1 and (3.15), if ϵ_1 is small enough, then there exists a positive constant $C_1 = C_1(G, A, M_0, M, a, n, \gamma)$ such that*

$$(3.17) \quad C_1^{-1} (M-x)^{\frac{1}{\gamma}} \leq \rho(x, t) \leq C_1 (M-x)^{\frac{1}{\gamma}},$$

$$(3.18) \quad r(x, t) \in [a, C_1]$$

for all $t \in [0, T]$ and $x \in [0, M]$.

Proof. From (1.15)₃, (3.15), and Lemma 3.1, we can easily obtain the estimate (3.17) and (3.18), when $4\epsilon_1 < \min_{x \in [0, M]} g_\infty$. \square

LEMMA 3.4. *Under the assumptions of Lemma 3.3, if ϵ_1 is small enough, then there exists a positive constant $C_2 = C_2(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ such that*

$$(3.19) \quad \int_0^M \left\{ u^2 + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \leq C_2 \epsilon_0^2,$$

and

$$(3.20) \quad \int_0^t \|u(\cdot, s)\|_{L^\infty}^2 ds + \int_0^t \int_0^M (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2)(x, s) dx ds \leq C_2 \epsilon_0^2$$

for all $t \in [0, T]$.

Proof. From (2.6), (3.7), and (3.10), we have

$$(3.21) \quad \begin{aligned} & \frac{d}{dt} \left(\int_0^M \frac{1}{2} u^2 dx + S[V] - S[V_\infty] \right) \\ & + \int_0^M \left\{ \left(\frac{2}{n} c_1 + c_2 \right) \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 + \frac{2(n-1)}{n} c_1 \rho^{1+\theta} \left(r^{n-1} u_x - \frac{u}{r\rho} \right)^2 \right\} dx \\ & = - \int_0^M \Delta f u dx, \end{aligned}$$

where $V_\infty = \frac{r_\infty^n}{n}$ and $V = \frac{r^n}{n}$. From (2.7), (3.17)–(3.18), and Proposition 2.2, we have

$$(3.22) \quad \begin{aligned} & C^{-1} \int_0^M (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 dx \\ & \leq S[V] - S[V_\infty] \leq C \int_0^M (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 dx, \end{aligned}$$

when $\|(M-x)^{\frac{1}{\gamma}}(\rho^{-1} - \rho_\infty^{-1})\|_{L^\infty} + \|\frac{1}{n}(r^n - r_\infty^n)\|_{L^\infty} \leq C_3 \epsilon_1 \leq \delta_9$. From (1.24), (3.17)–(3.18), and (3.21)–(3.22), we obtain

$$(3.23) \quad \begin{aligned} & \int_0^M \left\{ u^2 + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \\ & + \int_0^t \int_0^M \{ \rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2 \} dx ds \\ & \leq C \epsilon_0^2 + C \int_0^t f_1(s) \|u(\cdot, s)\|_{L^\infty} ds. \end{aligned}$$

Since $\theta \in (0, \gamma - 1)$, we obtain

$$(3.24) \quad \begin{aligned} & |u(x, t)| = \left| \int_0^x u_x dy \right| \leq C \left(\int_0^x \rho^{\theta+1} u_x^2 dy \right)^{\frac{1}{2}} \left(\int_0^x \rho^{-\theta-1} dy \right)^{\frac{1}{2}} \\ & \leq C \left(\int_0^x \rho^{\theta+1} u_x^2 dy \right)^{\frac{1}{2}} \left(\int_0^x (M-y)^{-\frac{\theta+1}{\gamma}} dy \right)^{\frac{1}{2}} \leq C \left(\int_0^x \rho^{\theta+1} u_x^2 dy \right)^{\frac{1}{2}} \end{aligned}$$

and

$$(3.25) \quad C \int_0^t f_1(s) \|u(\cdot, s)\|_{L^\infty} ds \leq \frac{1}{2} \int_0^t \int_0^M \rho^{\theta+1} u_x^2 dy ds + C^2 \int_0^t f_1^2 dt.$$

From (1.23) and (3.23)–(3.25), we can immediately obtain (3.19)–(3.20). \square

LEMMA 3.5. *Under the assumptions of Lemma 3.3, if ϵ_1 is small enough, then there exists a positive constant $C_4 = C_4(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ such that*

$$(3.26) \quad (1+t)^\alpha \int_0^M \rho_\infty^{\theta-1} (g - g_\infty)^2 dx + \int_0^t \int_0^M (1+s)^\alpha [\rho_\infty^{\gamma-1} (g - g_\infty)^2 + (r - r_\infty)^2] dx ds \leq C_4 \epsilon_0$$

for all $t \in [0, T]$, where $\alpha = -\frac{5}{8}$.

Proof. Multiplying (1.15)₂ by $(1+t)^\alpha r^{1-n} (\frac{r^n}{n} - \frac{r_\infty^n}{n})$, and integrating over $[0, M]$, using integration by parts and the boundary conditions (1.17), we obtain

$$(3.27) \quad (1+t)^\alpha \int_0^M \left[A(\rho_\infty^\gamma - \rho^\gamma)(\rho^{-1} - \rho_\infty^{-1}) + G(M_0 + x)(r^{2-2n} - r_\infty^{2-2n}) \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right] dx \\ = - (1+t)^\alpha \int_0^M \frac{u_t}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx - (1+t)^\alpha \int_0^M \Delta f r^{1-n} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx \\ + (1+t)^\alpha \int_0^M (2c_1 + c_2) \rho^{1+\theta} \partial_x (r^{n-1} u) (\rho_\infty^{-1} - \rho^{-1}) dx \\ + (1+t)^\alpha \int_0^M 2c_1 (n-1) \rho^\theta \left(\frac{u}{r} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right)_x dx := \sum_{i=1}^4 B_i.$$

We can rewrite the left-hand side (L.H.S.) of (3.27) as follows:

$$\text{L.H.S. of (3.27)} = (1+t)^\alpha \int_0^M \left[A(\gamma + O(\epsilon_1)) \rho_\infty^{\gamma+1} (\rho^{-1} - \rho_\infty^{-1})^2 - (2n - 2 + O(\epsilon_1)) G(M_0 + x) r_\infty^{2-3n} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right)^2 \right] dx.$$

Similar to (2.5), we have

$$(3.28) \quad \text{L.H.S. of (3.27)} \geq C_5 (1+t)^\alpha \int_0^M [\rho_\infty^{\gamma-1} (g - g_\infty)^2 + (r - r_\infty)^2] dx,$$

when $\epsilon_1 \leq \delta_{10}$ is small enough.

Using (3.17)–(3.19), integration by parts, and Hölder’s inequality, we can estimate B_i as follows:

$$(3.29) \quad B_1 = -\frac{d}{dt} \int_0^M (1+t)^\alpha \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx + \alpha (1+t)^{\alpha-1} \int_0^M \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx \\ + (1+t)^\alpha \int_0^M u^2 \left(\frac{1}{n} + \frac{(n-1)r_\infty^n}{nr^n} \right) dx \\ \leq -\frac{d}{dt} \int_0^M (1+t)^\alpha \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx + C \int_0^M u^2 dx + C \epsilon_0^2 (1+t)^{\alpha-1},$$

$$(3.30) \quad B_2 \leq C\epsilon_0(1+t)^\alpha f_1,$$

$$\begin{aligned} B_3 &= -\frac{2c_1+c_2}{\theta} \int_0^M (\rho^\theta)_t(1+t)^\alpha \left(\frac{1}{\rho_\infty} - \frac{1}{\rho}\right) dy \\ &= -\frac{2c_1+c_2}{\theta} \int_0^M \partial_t h(\rho, \rho_\infty)(1+t)^\alpha dx \\ &= -\frac{2c_1+c_2}{\theta} \frac{d}{dt} \int_0^M h(\rho, \rho_\infty)(1+t)^\alpha dx + \frac{\alpha(2c_1+c_2)}{\theta} \int_0^M h(\rho, \rho_\infty)(1+t)^{\alpha-1} dx, \end{aligned}$$

(3.31)

where $h(\rho, \rho_\infty) = \int_{\rho_\infty}^\rho \theta s^{\theta-1} (\frac{1}{\rho_\infty} - \frac{1}{s}) ds \sim \rho_\infty^{\theta-1} (g - g_\infty)^2$ and

$$\begin{aligned} B_4 &\leq C(1+t)^\alpha \int_0^M [\rho^\theta |u_x| + \rho^{\theta-1} |u|] dx \\ &\leq C(1+t)^\alpha \left[\int_0^M (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dx \right]^{\frac{1}{2}} \left[\int_0^M (M-x)^{\frac{\theta-1}{\gamma}} dx \right]^{\frac{1}{2}} \\ (3.32) \quad &\leq C(1+t)^\alpha \left[\int_0^M (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dx \right]^{\frac{1}{2}}, \end{aligned}$$

since $\gamma + \theta - 1 > 0$. From (3.27)–(3.32), we get

$$\begin{aligned} &\frac{d}{dt} \int_0^M (1+t)^\alpha \left\{ \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) + \frac{2c_1+c_2}{\theta} h(\rho, \rho_\infty) \right\} dx \\ &+ C^{-1} \int_0^M (1+t)^{\alpha-1} \rho_\infty^{\theta-1} (g - g_\infty)^2 + (1+t)^\alpha \{ \rho_\infty^{\gamma-1} [(g - g_\infty)^2 + (r - r_\infty)^2] \} dx \\ &\leq C(1+t)^\alpha \left[\int_0^M \rho_\infty^{\theta+1} u_x^2 dx + \|u(\cdot, t)\|_{L^\infty}^2 \right]^{\frac{1}{2}} + C\epsilon_0(1+t)^\alpha f_1 \\ &+ C \int_0^M u^2 dx + C\epsilon_0^2(1+t)^{\alpha-1}. \end{aligned}$$

And using (1.23), (3.17)–(3.20), and Hölder’s inequality, we can immediately obtain (3.26). \square

Let $\epsilon_2 \in (0, \min\{\frac{1}{4}, \frac{\gamma-\theta-1}{\gamma-\theta}, \frac{\gamma-1}{2(3\gamma-1)}\})$ be a constant. Define $\{\beta_j\}$ and $\{\alpha_j\}$ by $\beta_{j+1} = \frac{\beta_j}{2} + \frac{1}{2} - \frac{\epsilon_2}{4}$, $\alpha_j = \frac{\beta_j}{2} - \frac{1}{2} - \frac{\epsilon_2}{4}$ and $\alpha_0 = \alpha = -\frac{5}{8}$, $j = 0, 1, \dots$. Let N_4 be an integer satisfying $\beta_{N_4} \in [1 - \epsilon_2, 1 - \frac{3\epsilon_2}{4})$ and $\alpha_{N_4} \in (-\epsilon_2, -\frac{\epsilon_2}{4})$. It is easy to see that $\beta_0 = -\frac{1}{4} + \frac{\epsilon_2}{2} < 0$, $\alpha_j \in [-\frac{5}{8}, -\frac{\epsilon_2}{4})$ and $\beta_j \in (-\frac{1}{4}, 1 - \frac{3\epsilon_2}{4})$, $j = 0, 1, \dots, N_4$. Then, the following lemma can be proved by induction.

LEMMA 3.6. *Under the assumptions of Lemma 3.3, if ϵ_1 is small enough, then there exists a positive constant $C_7 = C_7(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ such that*

$$(3.33) \quad \int_0^M \left\{ u^2 + (M-x)^{\frac{\gamma-1}{\gamma}} (g - g_\infty)^2 + (r - r_\infty)^2 \right\} dx \leq C_7 \epsilon_0^{2^{1-N_4}} (1+t)^{\epsilon_2-1},$$

$$(3.34) \quad \int_0^t (1+s)^{1-\epsilon_2} \|u(\cdot, s)\|_{L^\infty}^2 ds + \int_0^t \int_0^M (1+s)^{1-\epsilon_2} (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dx ds \leq C_7 \epsilon_0^{2^{1-N_4}},$$

and

$$(3.35) \quad \int_0^M \frac{(M-x)^{\frac{\theta-1}{\gamma}}}{(1+t)^{\epsilon_2}} (g-g_\infty)^2 dx + \int_0^t \int_0^M \frac{\rho_\infty^{\gamma-1} (g-g_\infty)^2 + (r-r_\infty)^2}{(1+s)^{\epsilon_2}} dx ds \leq C_7 \epsilon_0^{2-N_4}$$

for all $t \in [0, T]$.

Proof. The following estimates can be proved by induction:

$$(3.36) \quad \int_0^M \left\{ u^2 + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \leq C_{j,\epsilon_2} \epsilon_0^{2^{1-j}} (1+t)^{-\beta_j},$$

$$(3.37) \quad \int_0^t (1+s)^{\beta_j} \|u(\cdot, s)\|_{L^\infty}^2 ds + \int_0^t \int_0^M (1+s)^{\beta_j} (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2)(x, s) dx ds \leq C_{j,\epsilon_2} \epsilon_0^{2^{1-j}},$$

and

$$(3.38) \quad (1+t)^{\alpha_j} \int_0^M \rho_\infty^{\theta-1} (g-g_\infty)^2 dx + \int_0^t \int_0^M (1+s)^{\alpha_j} [\rho_\infty^{\gamma-1} (g-g_\infty)^2 + (r-r_\infty)^2] dx ds \leq C_{j,\epsilon_2} \epsilon_0^{2^{-j}}$$

for all $t \geq 0$, where C_{j,ϵ_2} is a constant depending on j and ϵ_2 , $j = 0, 1, \dots, N_4$.

From (3.19)–(3.20) and (3.26), we obtain that (3.36)–(3.38) hold with $j = 0$. Now, suppose that (3.36)–(3.38) hold with $j = k \geq 0$. To show (3.36)–(3.37) hold with $j = k + 1$, from (3.21), we have

$$\begin{aligned} & \frac{d}{dt} \left\{ (1+t)^{\beta_{k+1}} \left(\int_0^M \frac{1}{2} u^2 dx + S[V] - S[V_\infty] \right) \right\} \\ & + (1+t)^{\beta_{k+1}} \int_0^M \left(\frac{2}{n} c_1 + c_2 \right) \rho^{1+\theta} [\partial_x (r^{n-1} u)]^2 + \frac{2(n-1)}{n} c_1 \rho^{1+\theta} \left(r^{n-1} u_x - \frac{u}{r\rho} \right)^2 dx \\ & = \beta_{k+1} (1+t)^{\alpha_k} \left(\int_0^M \frac{1}{2} u^2(x, t) dx + S[V] - S[V_\infty] \right) - (1+t)^{\beta_{k+1}} \int_0^M \Delta f u dx. \end{aligned}$$

Integrating the above equality in $[0, t]$, using (1.24), (3.17)–(3.20), (3.22), and (3.38) with $j = k$ and the fact that $\alpha_k < 0$, we obtain

$$(3.39) \quad \begin{aligned} & (1+t)^{\beta_{k+1}} \int_0^M \left\{ u^2(x, t) + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \\ & + \int_0^t \int_0^M (1+s)^{\beta_{k+1}} \{ \rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2 \} dx ds \\ & \leq C \epsilon_0^{2^{-k}} + C \int_0^t (1+s)^{\beta_{k+1}} f_1(s) \|u(\cdot, s)\|_{L^\infty} ds. \end{aligned}$$

From (1.23) and (3.25), we can immediately obtain (3.36)–(3.37) with $j = k + 1$.

To show (3.38) with $j = k + 1$, from (3.27)–(3.28), we have

$$\begin{aligned}
 & (1+t)^{\alpha_{k+1}} \int_0^M \left[A(\rho_\infty^\gamma - \rho^\gamma)(\rho^{-1} - \rho_\infty^{-1}) \right. \\
 & \left. + G(M_0 + x)(r^{2-2n} - r_\infty^{2-2n}) \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right] dx \\
 = & -(1+t)^{\alpha_{k+1}} \int_0^M \frac{u_t}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx - (1+t)^{\alpha_{k+1}} \int_0^M \Delta f r^{1-n} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx \\
 & + (1+t)^{\alpha_{k+1}} \int_0^M (2c_1 + c_2) \rho^{1+\theta} \partial_x (r^{n-1} u) (\rho_\infty^{-1} - \rho^{-1}) dx \\
 (3.40) \quad & + (1+t)^{\alpha_{k+1}} \int_0^M 2c_1(n-1) \rho^\theta \left(\frac{u}{r} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right)_x dx := \sum_{i=1}^4 E_i,
 \end{aligned}$$

and

$$(3.41) \quad \text{L.H.S. of (3.40)} \geq C_8(1+t)^{\alpha_{k+1}} \int_0^M \left[\rho_\infty^{\gamma-1} (g - g_\infty)^2 + (r - r_\infty)^2 \right] dx.$$

Similar to (3.29)–(3.32), applying the estimates (3.17)–(3.19), integration by parts, Hölder’s inequality, and the fact that $\alpha_{k+1} < 0$, we can estimate E_i as follows:

$$(3.42) \quad E_1 \leq -\frac{d}{dt} \int_0^M (1+t)^{\alpha_{k+1}} \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx + C \|u\|_{L_x^\infty}^2 + C \epsilon_0^2 (1+t)^{\alpha_{k+1}-1},$$

$$(3.43) \quad E_2 \leq C \epsilon_0 f_1 (1+t)^{\alpha_{k+1}},$$

$$\begin{aligned}
 (3.44) \quad E_3 = & -\frac{2c_1 + c_2}{\theta} \frac{d}{dt} \int_0^M h(\rho, \rho_\infty) (1+t)^{\alpha_{k+1}} dx \\
 & + \frac{\alpha_{k+1}(2c_1 + c_2)}{\theta} \int_0^M h(\rho, \rho_\infty) (1+t)^{\alpha_{k+1}-1} dx,
 \end{aligned}$$

and

$$(3.45) \quad E_4 \leq C(1+t)^{-\frac{1}{2} - \frac{\epsilon_2}{4}} \left[(1+t)^{\beta_{k+1}} \int_0^M \rho^{\theta+1} u_x^2 dx + \|u(\cdot, t)\|_{L^\infty}^2 \right]^{\frac{1}{2}}.$$

Using (1.23), (3.40)–(3.45), (3.36)–(3.37) with $j = k + 1$ and Hölder’s inequality, we get

$$\begin{aligned}
 & \int_0^M (1+t)^{\alpha_{k+1}} \rho_\infty^{\theta-1} (g - g_\infty)^2 dx + \int_0^t \int_0^M (1+s)^{\alpha_{k+1}-1} \rho_\infty^{\theta-1} (g - g_\infty)^2 dx ds \\
 & + \int_0^t \int_0^M (1+s)^{\alpha_{k+1}} \left\{ \rho_\infty^{\gamma-1} [(g - g_\infty)^2 + (r - r_\infty)^2] \right\} dx ds \\
 \leq & C(1+t)^{\alpha_{k+1}} \int_0^M |u| |r - r_\infty| dx + C \int_0^M |u_0| |r_0 - r_\infty| dx \\
 & + C \int_0^t \left[f_1 (1+s)^{\alpha_{k+1}} + \epsilon_0^2 (1+s)^{\alpha_{k+1}-1} + \|u\|_{L_x^\infty}^2 \right] ds \\
 & + C \int_0^t (1+s)^{-\frac{1}{2} - \frac{\epsilon_2}{4}} \left[(1+t)^{\beta_{k+1}} \left(\int_0^M \rho^{\theta+1} u_x^2 dx + \|u(\cdot, t)\|_{L^\infty}^2 \right) \right]^{\frac{1}{2}} ds \\
 (3.46) \quad & \leq C \epsilon_0^{2-(k+1)}
 \end{aligned}$$

and finish the proof of (3.38) with $j = k + 1$. Thus, we show that (3.36)–(3.38) hold for $j = 0, 1, \dots, N_4$ and obtain (3.33)–(3.35). \square

From Lemma 3.6, we can obtain the following estimate of the weighted L^2 -norm of $g - g_\infty$.

LEMMA 3.7. *Under the assumptions of Lemma 3.3, then there exists a positive constant $C_9 = C_9(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ such that*

$$(3.47) \quad \int_0^M (M - x)^{\frac{\theta-1+(\gamma-\theta)\epsilon_2}{\gamma}} (g - g_\infty)^2 dx \leq C_9 \epsilon_0^{2-N_4}, \quad t \in [0, T].$$

Proof. Using (3.33), (3.35), and Hölder’s inequality, we have

$$\begin{aligned} & \int_0^M (M - x)^{\frac{\theta-1+(\gamma-\theta)\epsilon_2}{\gamma}} (g - g_\infty)^2 dx \\ & \leq C \left[\int_0^M (1+t)^{1-\epsilon_2} (M-x)^{\frac{\gamma-1}{\gamma}} (g - g_\infty)^2 dx \right]^{\epsilon_2} \left[\int_0^M \frac{(M-x)^{\frac{\theta-1}{\gamma}}}{(1+t)^{\epsilon_2}} (g - g_\infty)^2 dx \right]^{1-\epsilon_2} \\ & \leq C \epsilon_0^{2-N_4}. \quad \square \end{aligned}$$

Then, using a similar argument as that in [22], we can finish the proof of Claim 1 in the following lemma.

LEMMA 3.8. *Under the assumptions of Lemma 3.3, if ϵ_0 is small enough, then there exists a positive constant $C_{10} = C_{10}(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$ such that*

$$(3.48) \quad |g(x, t) - g_\infty(x)| \leq C_{10} \epsilon_0^{\frac{\theta}{\gamma} 2^{-N_4-1}} \quad \text{for all } (x, t) \in [0, M] \times [0, T].$$

Proof. From (3.2), for any fixed $x \in [0, M]$, we have

$$(3.49) \quad \begin{aligned} & I_1(x, t) + \frac{\theta}{2c_1 + c_2} \int_0^t A r^\beta(x, \tau) (\rho^\gamma(x, \tau) - \rho_\infty^\gamma(x)) d\tau \\ & = r_0^\beta(x) \rho_0^\theta(x) + I_2(x, t), \quad x \in [0, M], \quad t \in [0, T], \end{aligned}$$

where

$$\begin{aligned} I_1(x, t) &= r_\infty^\beta(x) \rho^\theta(x, t) - (r_\infty^\beta(x) - r^\beta(x, t)) \rho^\theta(x, t) \\ & \quad + \int_x^M \beta [(r^{\beta-n} \rho^{\theta-1})(y, t) - (r_0^{\beta-n} \rho_0^{\theta-1})(y)] dy \\ & \quad - \frac{\theta}{2c_1 + c_2} \int_x^M [(r^{\beta-n+1} u)(y, t) - (r_0^{\beta-n+1} u_0)(y)] dy \\ & \quad + \frac{\theta(\beta - n + 1)}{2c_1 + c_2} \int_0^t \int_x^M r^{\beta-n} u^2 dy d\tau \end{aligned}$$

and

$$\begin{aligned} & I_2(x, t) \\ &= -\frac{\theta A \beta}{2c_1 + c_2} \int_0^t \int_x^M r^{\beta-n} \frac{\rho^\gamma - \rho_\infty^\gamma}{\rho} dy d\tau \\ & \quad + \frac{\theta}{2c_1 + c_2} \int_0^t \int_x^M \{r^\beta G(M_0 + y)(r^{2-2n} - r_\infty^{2-2n}) + r^{\beta-n+1} \Delta f\} dy d\tau. \end{aligned}$$

Using (3.17)–(3.18), (3.47), Hölder’s inequality, and the condition $\epsilon_2 < \frac{\gamma-\theta-1}{\gamma-\theta}$, i.e., $\frac{\theta+1+(\gamma-\theta)\epsilon_2}{\gamma} < 1$, we have

$$\begin{aligned}
 & |(r - r_\infty)(x)| \leq C|r^n - r_\infty^n| \\
 & \leq C \int_0^x |\rho^{-1} - \rho_\infty^{-1}| dy \leq C \int_0^x (M - y)^{-\frac{1}{\gamma}} |g - g_\infty| dy \\
 & \leq C \left(\int_0^x (M - y)^{\frac{\theta-1+(\gamma-\theta)\epsilon_2}{\gamma}} (g - g_\infty)^2 dy \right)^{\frac{1}{2}} \left(\int_0^x (M - y)^{-\frac{\theta+1+(\gamma-\theta)\epsilon_2}{\gamma}} dy \right)^{\frac{1}{2}} \\
 (3.50) \quad & \leq C\epsilon_0^{2^{-N_4-1}}
 \end{aligned}$$

and

$$\begin{aligned}
 (3.51) \quad & \int_x^M |\rho^{\theta-1} - \rho_\infty^{\theta-1}| dy \leq C \int_x^M (M - y)^{\frac{\theta-1}{\gamma}} |g - g_\infty| dy \\
 & \leq C \left(\int_x^M (M - y)^{\frac{\theta-1+(\gamma-\theta)\epsilon_2}{\gamma}} (g - g_\infty)^2 dy \right)^{\frac{1}{2}} \left(\int_x^M (M - y)^{\frac{2\theta}{\gamma} - \frac{\theta+1+(\gamma-\theta)\epsilon_2}{\gamma}} dy \right)^{\frac{1}{2}} \\
 & \leq C\epsilon_0^{2^{-N_4-1}} (M - x)^{\frac{\theta}{\gamma}}.
 \end{aligned}$$

From the fact $\theta \in (0, \frac{\gamma}{2}]$ and the estimate (3.18)–(3.19), we have

$$\begin{aligned}
 & \left| -\frac{\theta}{2c_1 + c_2} \int_x^M [(r^{\beta-n+1}u)(y, t) - (r_0^{\beta-n+1}u_0)(y)] dy \right| \\
 & \leq C(M - x)^{\frac{1}{2}} (\|u\|_{L_x^2} + \|u_0\|_{L_x^2}) \\
 (3.52) \quad & \leq C\epsilon_0^{2^{-N_4-1}} (M - x)^{\frac{\theta}{\gamma}}, \quad x \in [0, M].
 \end{aligned}$$

Thus, from (1.23), (3.17)–(3.20), and (3.50)–(3.52), we obtain

$$(3.53) \quad |I_1(x, t) - r_\infty^\beta \rho^\theta| \leq C_{1,1} (M - x)^{\frac{\theta}{\gamma}} \epsilon_0^{2^{-N_4-1}},$$

and

$$(3.54) \quad |I_2(x, t_1) - I_2(x, t_2)| \leq C_{1,2} (M - x) \epsilon_0^{2^{-N_4-1}} |t_2 - t_1|, \quad x \in [0, M].$$

CLAIM 2. For any fixed $x \in [0, M]$, we have

$$\begin{aligned}
 I_1(x, t) & \geq \min \left\{ I_1(x, 0), r_\infty^\beta \left(\rho_\infty^\gamma - \frac{C_{1,2}}{C_{1,3}} \epsilon_0^{2^{-N_4-1}} (M - x) \right)^{\frac{\theta}{\gamma}} - C_{1,1} \epsilon_0^{2^{-N_4-1}} (M - x)^{\frac{\theta}{\gamma}} \right\} \\
 & := M_{1,1},
 \end{aligned}$$

where $C_{1,3} := \frac{A\theta a^\beta}{2c_1+c_2} \leq \frac{A\theta r^\beta}{2c_1+c_2}$.

Proof of Claim 2. If not, there exists $t_{1,1} > 0$ such that $I_1(x, t_{1,1}) < M_{1,1}$, then we can find $t_{1,2} \in (0, t_{1,1})$ such that $I_1(x, t_{1,2}) = M_{1,1}$ and $I_1(x, t) < M_{1,1}$ for all $t \in (t_{1,2}, t_{1,1})$. From (3.54), we have

$$I_1(x, t_{1,1}) - I_1(x, t_{1,2}) + \frac{A\theta}{2c_1 + c_2} \int_{t_{1,2}}^{t_{1,1}} r^\beta (\rho^\gamma - \rho_\infty^\gamma) \geq -C_{1,2} \epsilon_0^{2^{-N_4-1}} (M - x) (t_{1,1} - t_{1,2}).$$

From (3.53), we have

$$\begin{aligned} \rho^\theta(x, t) &= r_\infty^{-\beta}(I_1(x, t) - (I_1(x, t) - r_\infty^\beta \rho^\theta)) \\ &\leq r_\infty^{-\beta}(M_{1,1} + C_{1,1}\epsilon_0^{2^{-N_4-1}}(M-x)^{\frac{\theta}{\gamma}}) \leq \left(\rho_\infty^\gamma - \frac{C_{1,2}}{C_{1,3}}\epsilon_0^{2^{-N_4-1}}(M-x)\right)^{\frac{\theta}{\gamma}} \end{aligned}$$

and

$$\rho^\gamma \leq \rho_\infty^\gamma - \frac{C_{1,2}}{C_{1,3}}\epsilon_0^{2^{-N_4-1}}(M-x),$$

then $I_1(x, t_{1,1}) \geq I_1(x, t_{1,2})$. It is a contradiction. Thus, Claim 2 holds.

Similarly, we can obtain the following claim.

CLAIM 3. *For any fixed $x \in [0, M]$, we have*

$$\begin{aligned} I_1(x, t) &\leq \max \left\{ I_1(x, 0), r_\infty^\beta \left(\rho_\infty^\gamma + \frac{C_{1,2}}{C_{1,4}}\epsilon_0^{2^{-N_4-1}}(M-x) \right)^{\frac{\theta}{\gamma}} + C_{1,1}\epsilon_0^{2^{-N_4-1}}(M-x)^{\frac{\theta}{\gamma}} \right\} \\ &:= M_{1,2}, \end{aligned}$$

where $C_{1,4}$ is a positive constant satisfying $C_{1,4} \geq \frac{A\theta r^\beta}{2c_1+c_2}$.

From Claims 2 and 3, we have

$$|g(x, t) - g_\infty(x)| \leq C_{1,5}\epsilon_0^{\frac{\theta}{\gamma}2^{-N_4-1}}, \quad x \in [0, M], \quad t \in [0, T],$$

when $\epsilon_0 \leq \delta_{11}$, with a small positive constant $\delta_{11}(G, A, M_0, M, a, n, \gamma, \theta, c_1, c_2)$. □

Using the results in Lemmas 3.3–3.8, we have that there exists ϵ_0 satisfying

$$(3.55) \quad 4\epsilon_1 < \min_{x \in [0, M]} g_\infty, \quad C_3\epsilon_1 \leq \delta_9, \quad \epsilon_1 \leq \delta_{10} \quad \text{and} \quad \epsilon_0 \leq \delta_{11}$$

such that Claim 1 holds with

$$\epsilon_1 = \epsilon_0 + C_{10}\epsilon_0^{\frac{\theta}{\gamma}2^{-N_4-1}}.$$

From (3.6) and Claim 1, using the classical continuation method, we can easily obtain the following lemma.

LEMMA 3.9. *Under the assumptions of Theorem 1.1, we obtain that (3.17)–(3.18), (3.33)–(3.35), (3.48), and*

$$(3.56) \quad |r(x, t) - r_\infty(x)| \leq C_{11}\epsilon_0^{\frac{\theta}{\gamma}2^{-N_4-1}}$$

hold for all $x \in [0, M]$ and $t \in [0, T^*)$.

Proof. Let $\mathcal{A} = \{T \in [0, T^*) \mid I(t) \leq \epsilon_1 \text{ for all } t \in [0, T]\}$. Since $I(0) \leq \epsilon_0 < \epsilon_1$ and $I(t) \in C([0, T^*))$, then there exists a constant $T_0 > 0$ such that $I(t) \leq \epsilon_1$ for all $t \in [0, T_0]$. Thus, \mathcal{A} is not empty and relatively closed in $[0, T^*)$. To show that \mathcal{A} is also relatively open in $[0, T^*)$, and hence the entire interval, it therefore suffices to show that the weaker bound

$$I(t) \leq 2\epsilon_1, \quad \text{for all } t \in [0, T'] \subset [0, T^*),$$

implies $I(t) \leq \epsilon_1$ for all $t \in [0, T']$. From Claim 1, we have that $\mathcal{A} = [0, T^*)$.

Then, from Lemmas 3.3–3.8, we obtain that (3.17)–(3.18), (3.33)–(3.35), (3.48), and (3.56) hold for all $x \in [0, M]$ and $t \in [0, T^*)$. □

We will prove an estimate in weighted $L^2([0, M] \times [0, T^*))$ -norm of the function $g - g_\infty$.

LEMMA 3.10. *Under the assumptions of Theorem 1.1, we obtain*

$$(3.57) \quad \int_0^t \int_0^M (1+s)^{-\epsilon_2} (g - g_\infty)^2 dx ds \leq C,$$

where $t \in [0, T^*)$.

Proof. From (1.15), we have

$$\begin{aligned} A(\rho^\gamma - \rho_\infty^\gamma) &= \int_x^M \left(\frac{u_t}{r^{n-1}} + \frac{\Delta f}{r^{n-1}} \right) dy + \int_x^M G(M_0 + y)(r^{2-2n} - r_\infty^{2-2n}) dy \\ &\quad + (2c_1 + c_2)\rho^{\theta+1}(r^{n-1}u)_x - 2c_1(n-1)\rho^\theta \frac{u}{r} - 2c_1(n-1) \int_x^M \rho^\theta \left(\frac{u}{r} \right)_x dy. \end{aligned}$$

Multiplying the above equality by $(1+t)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma)$, and integrating the resulting equation over $[0, M] \times [0, t]$, we obtain

$$\begin{aligned} (3.58) \quad & \int_0^t \int_0^M A(1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma)^2 dx ds \\ &= \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \frac{u_t}{r^{n-1}} dy dx ds \\ &\quad + \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \frac{\Delta f}{r^{n-1}} dy dx ds \\ &\quad + \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M G(M_0 + y)(r^{2-2n} - r_\infty^{2-2n}) dy dx ds \\ &\quad + \int_0^t \int_0^M (2c_1 + c_2)(1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \rho^{\theta+1} (r^{n-1}u)_x dx ds \\ &\quad - \int_0^t \int_0^M 2c_1(n-1)(1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \rho^\theta \frac{u}{r} dx ds \\ &\quad - \int_0^t \int_0^M 2c_1(n-1)(1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \rho^\theta \left(\frac{u}{r} \right)_x dy dx ds \\ &:= \sum_{i=1}^6 F_i. \end{aligned}$$

Using (1.15), (1.23), Lemma 3.9, integration by parts, and the Cauchy-Schwarz inequality, we can estimate F_i as follows.

$$\begin{aligned} (3.59) \quad F_1 &= \left\{ \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \frac{u}{r^{n-1}} dy dx \right\} \Big|_0^t \\ &\quad + \int_0^t \int_0^M \epsilon_2 (1+s)^{-\epsilon_2-1} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \frac{u}{r^{n-1}} dy dx ds \\ &\quad + \int_0^t \int_0^M \gamma (1+s)^{-\epsilon_2} (M-x)^{-2} \rho^{\gamma+1} (r^{n-1}u)_x \int_x^M \frac{u}{r^{n-1}} dy dx ds \\ &\quad + \int_0^t \int_0^M (n-1)(1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) \int_x^M \frac{u^2}{r^n} dy dx ds \end{aligned}$$

$$\begin{aligned}
 &\leq C\|g - g_\infty\|_{L^\infty_{xt}} \left(\int_0^M u^2 dx \right)^{\frac{1}{2}} \left(\int_0^M (M-x)^{-\frac{1}{2}} dx \right)^{\frac{1}{2}} + C \\
 &\quad + C \left(\int_0^t \int_x^M \frac{\rho_\infty^{\gamma-1}}{(1+s)^{\epsilon_2}} (g - g_\infty)^2 dx ds \right)^{\frac{1}{2}} \\
 &\quad \quad \left(\int_0^t \|u(\cdot, s)\|_{L^\infty}^2 ds \int_0^M (M-x)^{-\frac{\gamma-1}{\gamma}} dx \right)^{\frac{1}{2}} \\
 &\quad + C \left(\int_0^t \int_x^M \rho^{\theta+1} (r^{n-1} u)_x^2 dx ds \right)^{\frac{1}{2}} \left(\int_0^t \|u(\cdot, s)\|_{L^\infty}^2 ds \int_0^M (M-x)^{-\frac{\theta-1}{\gamma}} dx \right)^{\frac{1}{2}} \\
 &\quad + C\|g - g_\infty\|_{L^\infty_{xt}} \int_0^t \|u(\cdot, s)\|_{L^\infty}^2 ds \\
 &\leq C,
 \end{aligned}$$

$$(3.60) \quad F_2 \leq C \left(\int_0^t (1+t)^{-1-\epsilon_2} \right)^{\frac{1}{2}} \left(\int_0^t (1+t) f_1^2 dt \right)^{\frac{1}{2}} \leq C,$$

$$\begin{aligned}
 &|r(x, t) - r_\infty(x)| \\
 &\leq C \int_0^x \rho_\infty^{-1} |g - g_\infty| dy \leq C \left(\int_0^x \rho_\infty^{\gamma-1} (g - g_\infty)^2 dy \right)^{\frac{1}{2}} \left(\int_0^x \rho_\infty^{-\gamma-1} dy \right)^{\frac{1}{2}} \\
 (3.61) \quad &\leq C(M-x)^{-\frac{1}{2\gamma}} \left(\int_0^M \rho_\infty^{\gamma-1} (g - g_\infty)^2 dx \right)^{\frac{1}{2}},
 \end{aligned}$$

$$\begin{aligned}
 F_3 &\leq C \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} |\rho^\gamma - \rho_\infty^\gamma| \\
 &\quad \times \left(\int_0^M \rho_\infty^{\gamma-1} (g - g_\infty)^2 dz \right)^{\frac{1}{2}} \int_x^M (M-y)^{-\frac{1}{2\gamma}} dy dx \\
 &\leq \frac{A}{4} \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma)^2 dx ds \\
 &\quad + C \int_0^t \int_0^M (1+s)^{-\epsilon_2} \rho_\infty^{\gamma-1} (g - g_\infty)^2 dx ds \int_0^M (M-z)^{-\frac{1}{\gamma}} dz \\
 (3.62) \quad &\leq \frac{A}{4} \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma)^2 dx ds + C,
 \end{aligned}$$

$$\begin{aligned}
 F_4 &= -\frac{2c_1 + c_2}{\theta} \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-2} (\rho^\gamma - \rho_\infty^\gamma) (\rho^\theta)_t dx ds \\
 &= -\frac{2c_1 + c_2}{\theta} \left\{ (1+s)^{-\epsilon_2} \int_0^M (M-x)^{-2} \left(\frac{\theta}{\gamma + \theta} \rho^{\gamma+\theta} - \rho_\infty^\gamma \rho^\theta \right) dx \right\} \Big|_0^t \\
 &\quad - \frac{\epsilon_2(2c_1 + c_2)}{\theta} \int_0^t \int_0^M (1+s)^{-\epsilon_2-1} (M-x)^{-2} \left(\frac{\theta}{\gamma + \theta} \rho^{\gamma+\theta} - \rho_\infty^\gamma \rho^\theta \right) dx \\
 &\leq C\|g - g_\infty\|_{L^\infty_{tx}} \int_0^M (M-x)^{\frac{\theta}{\gamma}-1} dx \left(1 + \int_0^t (1+s)^{-1-\epsilon_2} ds \right) + C \\
 (3.63) \quad &\leq C,
 \end{aligned}$$

$$\begin{aligned}
 F_5 &\leq C \int_0^t \int_0^M (1+s)^{-\epsilon_2} \|u(\cdot, t)\|_{L^\infty} (M-x)^{\frac{\theta}{\gamma}-1} dx ds \\
 (3.64) \quad &\leq C \left\{ \int_0^t (1+s)^{-1-\epsilon_2} ds \right\}^{\frac{1}{2}} \left\{ \int_0^t (1+s)^{1-\epsilon_2} \|u(\cdot, t)\|_{L^\infty}^2 ds \right\}^{\frac{1}{2}} \leq C
 \end{aligned}$$

and

$$\begin{aligned}
 (3.65) \quad F_6 &\leq C \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-1} \int_x^M (|\rho^\theta u_x| + |\rho^{\theta-1} u|) dy dx ds \\
 &\leq C \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{-1} \left[\int_x^M (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dy \right]^{\frac{1}{2}} \left[\int_x^M \rho^{\theta-1} dy \right]^{\frac{1}{2}} dx ds \\
 &\leq C \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{\frac{\theta-1}{2\gamma}-\frac{1}{2}} \left[\int_x^M (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dy \right]^{\frac{1}{2}} dx ds \\
 &\leq C \left[\int_0^t \int_0^M (1+s)^{1-\epsilon_2} (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2) dx ds \right]^{\frac{1}{2}} \left[\int_0^t (1+s)^{-1-\epsilon_2} dt \right]^{\frac{1}{2}} \\
 &\leq C.
 \end{aligned}$$

From (3.58)–(3.65), we can immediately get (3.57). \square

LEMMA 3.11. *Under the assumptions of Theorem 1.1, we obtain*

$$\begin{aligned}
 (3.66) \quad &(1+t)^{-\epsilon_2} \int_0^M (M-x)^{1-\frac{\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2(x, t) dx \\
 &+ \int_0^t \int_0^M (1+s)^{-\epsilon_2} (M-x)^{2-\frac{2\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2(x, s) dx ds \leq C
 \end{aligned}$$

for all $t \in [0, T^*)$.

Proof. From (3.2), we have

$$\begin{aligned}
 &\partial_t \left[\frac{\theta}{2c_1 + c_2} r^{1+\beta-n} u + r^\beta (\rho^\theta)_x - r_\infty^\beta (\rho_\infty^\theta)_x \right] \\
 &+ \frac{A\gamma\rho^{\gamma-\theta}}{2c_1 + c_2} \left[\frac{\theta}{2c_1 + c_2} r^{1+\beta-n} u + r^\beta (\rho^\theta)_x - r_\infty^\beta (\rho_\infty^\theta)_x \right] \\
 (3.67) \quad &= \frac{A\gamma\theta}{(2c_1 + c_2)^2} \rho^{\gamma-\theta} r^{1+\beta-n} u + \frac{\theta(1+\beta-n)}{2c_1 + c_2} r^{\beta-n} u^2 - \frac{\theta}{2c_1 + c_2} r^{1+\beta-n} \Delta f \\
 &- \frac{\theta}{2c_1 + c_2} \left(r^\beta \frac{G(M_0 + x)}{r^{2n-2}} + r_\infty^\beta \frac{\rho^{\gamma-\theta}}{\rho_\infty^{\gamma-\theta}} (A\rho_\infty^\gamma)_x \right).
 \end{aligned}$$

Let $H = \frac{\theta}{2c_1+c_2} r^{1+\beta-n} u + r^\beta (\rho^\theta)_x - r_\infty^\beta (\rho_\infty^\theta)_x$. Multiplying (3.67) by $(1+t)^{-\epsilon_2} (M-x)^{1-\frac{\theta}{\gamma}} H$, integrating the resulting equation over $[0, M]$, and using the Cauchy-Schwarz

inequality, we obtain

$$\begin{aligned}
 & \frac{d}{dt} \int_0^M (1+t)^{-\epsilon_2} (M-x)^{1-\frac{\theta}{\gamma}} H^2 dx + C_{13} \int_0^M (1+t)^{-\epsilon_2} (M-x)^{2-\frac{2\theta}{\gamma}} H^2 dx \\
 & \leq C \int_0^M (1+t)^{-\epsilon_2} (M-x)^{1-\frac{\theta}{\gamma}} \left((M-x)^{1-\frac{\theta}{\gamma}} |Hu| + |u^2 H| + |\Delta f H| \right) dx \\
 & \quad + \int_0^M (1+t)^{-\epsilon_2} (M-x)^{1-\frac{\theta}{\gamma}} \left| r^\beta G \frac{M_0+x}{r^{2n-2}} + \frac{\rho^{\gamma-\theta} r_\infty^\beta}{\rho_\infty^{\gamma-\theta}} (A\rho_\infty^\gamma)_x \right| |H| dx \\
 & \leq \frac{C_{13}}{4} \int_0^M (1+t)^{-\epsilon_2} (M-x)^{2-\frac{2\theta}{\gamma}} H^2 dx + C \|u(\cdot, t)\|_{L^\infty}^2 + C \|u(\cdot, t)\|_{L^\infty}^2 \|u(\cdot, t)\|_{L^2}^2 \\
 (3.68) \quad & + C f_1^2 + C \int_0^M (1+t)^{-\epsilon_2} \left| r^\beta G \frac{M_0+x}{r^{2n-2}} + \frac{\rho^{\gamma-\theta} r_\infty^\beta}{\rho_\infty^{\gamma-\theta}} (A\rho_\infty^\gamma)_x \right|^2 dx.
 \end{aligned}$$

Here, we use the estimates (3.17)–(3.18) and the condition $\theta \in (0, \gamma - 1)$. From (3.9) and (3.17)–(3.18), we have

$$\begin{aligned}
 & \int_0^M (1+t)^{-\epsilon_2} \left| r^\beta G \frac{M_0+x}{r^{2n-2}} + \frac{\rho^{\gamma-\theta} r_\infty^\beta}{\rho_\infty^{\gamma-\theta}} (A\rho_\infty^\gamma)_x \right|^2 dx \\
 & = \int_0^M (1+t)^{-\epsilon_2} \left| r^\beta G \frac{M_0+x}{r^{2n-2}} - G \frac{(\rho^{\gamma-\theta} r_\infty^\beta)(M_0+x)}{(\rho_\infty^{\gamma-\theta})(r_\infty^{2n-2})} \right|^2 dx \\
 (3.69) \quad & \leq C \int_0^M (1+t)^{-\epsilon_2} [(r-r_\infty)^2 + (g-g_\infty)^2] dx.
 \end{aligned}$$

From (3.33) and (3.68)–(3.69), we obtain

$$\begin{aligned}
 & \frac{d}{dt} \int_0^M (1+t)^{-\epsilon_2} (M-x)^{1-\frac{\theta}{\gamma}} H^2 dx + \frac{C_{13}}{2} \int_0^M (1+t)^{-\epsilon_2} (M-x)^{2-\frac{2\theta}{\gamma}} H^2 dx \\
 (3.70) \quad & \leq C \int_0^M (1+t)^{-\epsilon_2} ((r-r_\infty)^2 + (g-g_\infty)^2) dx + C \|u(\cdot, t)\|_{L^\infty}^2 + C f_1^2.
 \end{aligned}$$

From (A2), (1.23), (3.33)–(3.35), (3.56)–(3.57), and (3.70), we can immediately obtain (3.66). \square

LEMMA 3.12. *Under the assumptions of Theorem 1.1, we obtain*

$$\begin{aligned}
 (3.71) \quad & (1+t)^{1-\epsilon_2} \int_0^M (\rho^{\theta-1} u^2 + \rho^{\theta+1} u_x^2)(x, t) dx + \int_0^t \int_0^M (1+s)^{1-\epsilon_2} u_t^2(x, s) dx ds \leq C,
 \end{aligned}$$

$$(3.72) \quad \|u(\cdot, t)\|_{L^\infty} \leq C(1+t)^{-\frac{1}{2} + \frac{\epsilon_2}{2}}$$

for all $t \in [0, T^*)$.

Proof. Multiplying (1.15)₂ by $(1+t)^{1-\epsilon_2}u_t$, integrating the resulting equation over $[0, M] \times [0, t]$, and using integration by parts and the boundary conditions (1.17), we obtain

$$\begin{aligned}
 (3.73) \quad & \int_0^t \int_0^M (1+s)^{1-\epsilon_2} u_t^2(x, s) dx ds \\
 &= \int_0^t \int_0^M A(1+s)^{1-\epsilon_2} \rho^\gamma \partial_x(r^{n-1}u_t) dx ds \\
 &\quad - \int_0^t \int_0^M (2c_1 + c_2)(1+s)^{1-\epsilon_2} \rho^{1+\theta} \partial_x(r^{n-1}u) \partial_x(r^{n-1}u_t) dx ds \\
 &\quad + \int_0^t \int_0^M 2c_1(n-1)(1+s)^{1-\epsilon_2} \rho^\theta \partial_x(r^{n-2}uu_t) dx ds - \int_0^t \int_0^M (1+s)^{1-\epsilon_2} f u_t dx ds \\
 &:= \sum_{i=1}^4 H_i.
 \end{aligned}$$

Using (A3), (1.23), (3.17)–(3.18), (3.33)–(3.34), and the Cauchy–Schwarz inequality, we obtain

$$\begin{aligned}
 & H_2 + H_3 \\
 &= \left\{ (1+s)^{1-\epsilon_2} \int_0^M \left[-\frac{2c_1 + c_2}{2} \rho^{1+\theta} [\partial_x(r^{n-1}u)]^2 \right. \right. \\
 &\quad \left. \left. + c_1(n-1) \rho^\theta \partial_x(r^{n-2}u^2) \right] dx \right\} \Big|_0^t \\
 &\quad + \int_0^t \int_0^M \frac{(2c_1 + c_2)(1-\epsilon_2)}{2} (1+s)^{-\epsilon_2} \rho^{1+\theta} (r^{n-1}u)_x^2 dx ds \\
 &\quad - \int_0^t \int_0^M c_1(n-1)(1-\epsilon_2)(1+s)^{-\epsilon_2} \rho^\theta \partial_x(r^{n-2}u^2) dx ds \\
 &\quad + \int_0^t \int_0^M (1+s)^{1-\epsilon_2} \left\{ (2c_1 + c_2)(n-1) \rho^{1+\theta} \partial_x(r^{n-1}u) \partial_x(r^{n-2}u^2) \right. \\
 &\quad - \frac{(2c_1 + c_2)}{2} (1+\theta) \rho^{2+\theta} [\partial_x(r^{n-1}u)]^3 + 2\theta c_1(n-1) \rho^{\theta+1} \frac{u}{r} [\partial_x(r^{n-1}u)]^2 \\
 &\quad - \theta c_1 n(n-1) \rho^\theta \frac{u^2}{r^2} \partial_x(r^{n-1}u) + 2nc_1(n-1)(n-2) \rho^{\theta-1} \frac{u^3}{r^3} \\
 &\quad \left. - 3c_1(n-1)(n-2) \rho^\theta \frac{u^2}{r^2} \partial_x(r^{n-1}u) \right\} dx ds \\
 &\leq C + C \int_0^t (\|u\|_{L_x^\infty} + \|\rho(r^{n-1}u)_x\|_{L_x^\infty}) (1+s)^{1-\epsilon_2} \\
 &\quad \int_0^M [\rho^{1+\theta} (r^{n-1}u)_x^2 + \rho^{\theta-1} u^2] dx ds \\
 (3.74) \quad & - C_{14}(1+t)^{1-\epsilon_2} \int_0^M [\rho^{1+\theta} (r^{n-1}u)_x^2 + \rho^{\theta-1} u^2] dx,
 \end{aligned}$$

$$\begin{aligned}
 H_1 &= \left\{ (1+s)^{1-\epsilon_2} \int_0^M A\rho^\gamma \partial_x(r^{n-1}u) dx \right\} \Big|_0^t \\
 &\quad + \int_0^t \int_0^M A\gamma(1+s)^{1-\epsilon_2} \rho^{\gamma+1} [\partial_x(r^{n-1}u)]^2 dx ds \\
 &\quad - \int_0^t \int_0^M 2A(n-1)(1+s)^{1-\epsilon_2} \rho^\gamma \frac{u}{r} \partial_x(r^{n-1}u) dx ds \\
 &\quad + \int_0^t \int_0^M An(n-1)(1+s)^{1-\epsilon_2} \rho^{\gamma-1} \frac{u^2}{r^2} dx ds \\
 &\quad - \int_0^t \int_0^M A(1-\epsilon_2)(1+s)^{-\epsilon_2} \rho^\gamma \partial_x(r^{n-1}u) dx ds \\
 (3.75) \quad &\leq (1+s)^{1-\epsilon_2} \int_0^M A\rho^\gamma \partial_x(r^{n-1}u) dx + C,
 \end{aligned}$$

and

$$\begin{aligned}
 H_4 &= - \left\{ (1+s)^{1-\epsilon_2} \int_0^M G \frac{u(M_0+x)}{r^{n-1}} dx \right\} \Big|_0^t - \int_0^t \int_0^M (1+s)^{1-\epsilon_2} \Delta f u_t dx ds \\
 &\quad + (1-\epsilon_2) \int_0^t \int_0^M (1+s)^{-\epsilon_2} \frac{Gu(M_0+x)}{r^{n-1}} dx ds \\
 &\quad + \int_0^t \int_0^M (1-n)(1+s)^{1-\epsilon_2} G(M_0+x)r^{-n}u^2 dx ds \\
 (3.76) \quad &\leq -(1+s)^{1-\epsilon_2} \int_0^M G \frac{u(M_0+x)}{r^{n-1}} dx + C + \frac{1}{2} \int_0^t \int_0^M (1+s)^{1-\epsilon_2} u_t^2 dx ds.
 \end{aligned}$$

Using (3.17)–(3.18), (3.24), (3.33), integration by parts, and the Cauchy–Schwarz inequality, we obtain

$$\begin{aligned}
 (3.77) \quad &(1+s)^{1-\epsilon_2} \int_0^M A\rho^\gamma \partial_x(r^{n-1}u) dx - (1+s)^{1-\epsilon_2} \int_0^M G \frac{u(M_0+x)}{r^{n-1}} dx \\
 &= (1+s)^{1-\epsilon_2} \int_0^M (A(\rho^\gamma - \rho_\infty^\gamma) \partial_x(r^{n-1}u) - Gr^{n-1}u(M_0+x)(r^{2-2n} - r_\infty^{2-2n})) dx \\
 &\leq \frac{C_{14}}{4} (1+t)^{1-\epsilon_2} \int_0^M \rho^{1+\theta} (r^{n-1}u)_x^2 dx + C.
 \end{aligned}$$

From (3.73)–(3.77), we can obtain

$$\begin{aligned}
 &(1+t)^{1-\epsilon_2} \int_0^M [\rho^{1+\theta} (r^{n-1}u)_x^2 + \rho^{\theta-1}u^2] dx + \int_0^t \int_0^M (1+s)^{1-\epsilon_2} u_t^2 dx ds \\
 &\leq C + C \int_0^t (\|u\|_{L_x^\infty} + \|\rho(r^{n-1}u)_x\|_{L_x^\infty}) (1+s)^{1-\epsilon_2} \\
 (3.78) \quad &\times \int_0^M [\rho^{1+\theta} (r^{n-1}u)_x^2 + \rho^{\theta-1}u^2] dx ds.
 \end{aligned}$$

From (3.3), we have

$$\rho(r^{n-1}u)_x = \frac{1}{(2c_1 + c_2)\rho^\theta} \left\{ A\rho^\gamma + 2c_1(n-1)\rho^\theta \frac{u}{r} + \int_x^M \left[-\frac{u_t}{r^{n-1}} + 2c_1(n-1)\rho^\theta \left(\frac{u}{r} \right)_x - \frac{f}{r^{n-1}} \right] dy \right\}.$$

Using conditions $\theta \in (0, \frac{\gamma}{2}) \cap (0, \gamma - 1)$, (1.23), estimates (3.17)–(3.20), and Hölder’s inequality, we conclude that

$$(3.79) \quad \|\rho\partial_x(r^{n-1}u)\|_{L^\infty} \leq C + C \left(\|u(\cdot, t)\|_{L^\infty}^2 + \int_0^M [\rho^{1+\theta}(r^{n-1}u)_x^2 + u_t^2] dx \right)^{\frac{1}{2}}.$$

Using (3.24), (3.34), (3.78)–(3.79), and the Cauchy–Schwarz inequality, we can obtain

$$\begin{aligned} & (1+t)^{1-\epsilon_2} \int_0^M [\rho^{1+\theta}(r^{n-1}u)_x^2 + \rho^{\theta-1}u^2] dx + \int_0^t \int_0^M (1+s)^{1-\epsilon_2} u_t^2 dx ds \\ & \leq C + C \int_0^t \left(\int_0^M (1+s)^{1-\epsilon_2} [\rho^{1+\theta}(r^{n-1}u)_x^2 + \rho^{\theta-1}u^2] dx \right)^2 ds. \end{aligned}$$

Using Gronwall’s inequality and the estimate (3.34), we can immediately get (3.71). From (3.17), (3.71), and the fact $\theta \in (0, \gamma - 1)$, we can obtain

$$\begin{aligned} |u(x, t)| & \leq \left| \int_0^x u_x dy \right| \leq C \left(\int_0^x \rho^{\theta+1} u_x^2 dy \right)^{\frac{1}{2}} \left(\int_0^x (M-y)^{-\frac{\theta+1}{\gamma}} dy \right)^{\frac{1}{2}} \\ & \leq C(1+t)^{-\frac{1}{2} + \frac{\epsilon_2}{2}}, \quad (x, t) \in [0, M] \times [0, T^*]. \quad \square \end{aligned}$$

LEMMA 3.13. *Under the assumptions of Theorem 1.1, we obtain*

$$(3.80) \quad \int_0^M u_t^2(x, t) dx + \int_0^t \int_0^M [\rho^{\theta+1} u_{xt}^2 + \rho^{\theta-1} u_t^2] dx ds \leq C_{11},$$

$$(3.81) \quad \|\rho(r^{n-1}u)_x(\cdot, t)\|_{L^\infty} \leq C_{11}$$

for all $t \in [0, T^*]$.

Proof. We differentiate the equation (1.15)₂ with respect to t , multiply it by u_t , and integrate it over $[0, M] \times [0, t]$ using the boundary conditions (1.17), then derive

$$\begin{aligned} (3.82) \quad & \int_0^M \frac{1}{2} u_t^2 dx \\ & = \int_0^M \frac{1}{2} u_t^2(x, 0) dx - \int_0^t \int_0^M \left[(2c_1 + c_2)\rho^{1+\theta} \partial_x(r^{n-1}u) - A\rho^\gamma - 2c_1(n-1)\rho^\theta \frac{u}{r} \right] \\ & \quad \times \partial_x((n-1)r^{n-2}uu_t) dx ds - \int_0^t \int_0^M \partial_t \left[(2c_1 + c_2)\rho^{1+\theta} \partial_x(r^{n-1}u) - A\rho^\gamma \right. \\ & \quad \left. - 2c_1(n-1)\rho^\theta \frac{u}{r} \right] \partial_x(r^{n-1}u_t) dx ds + \int_0^t \int_0^M 2c_1(n-1)\partial_t \left(r^{n-1}\rho^\theta \partial_x \left(\frac{u}{r} \right) \right) u_t dx ds \end{aligned}$$

$$\begin{aligned}
 & - \int_0^t \int_0^M f_t u_t dx ds \\
 & := \sum_{i=1}^5 J_i.
 \end{aligned}$$

From (A2)–(A3), we have

$$\begin{aligned}
 J_1 & \leq C \left(\left\| \left((2c_1 + c_2)\rho_0^{\theta+1}(r_0^{n-1}u_0)_x \right)_x - 2c_1(n-1)\frac{u_0}{r_0}(\rho_0^\theta)_x \right\|_{L^2} \right. \\
 & \quad \left. + \|(\rho_0^\gamma)_x\|_{L^2} + \|f(x, r_0, 0)\|_{L^2} \right)^2 \\
 (3.83) \quad & \leq C.
 \end{aligned}$$

From (3.17)–(3.20) and the Cauchy–Schwarz inequality, we get

$$\begin{aligned}
 & J_3 + J_4 \\
 & = - \int_0^t \int_0^M \left[(2c_1 + c_2)\rho^{1+\theta}(r^{n-1}u_t)_x^2 - 2c_1(n-1)\rho^\theta(r^{n-2}u_t^2)_x \right] dx ds \\
 & \quad + \int_0^t \int_0^M \left\{ (2c_1 + c_2)(1 + \theta)\rho^{\theta+2}[\partial_x(r^{n-1}u)]^2 - (n-1)(2c_1 + c_2)\rho^{1+\theta}\partial_x(r^{n-2}u^2) \right. \\
 & \quad \left. - \gamma\rho^{\gamma+1}\partial_x(r^{n-1}u) - 2c_1(n-1)\theta\rho^{\theta+1}\partial_x(r^{n-1}u)\frac{u}{r} - 2c_1(n-1)\rho^\theta\frac{u^2}{r^2} \right\} \\
 & \quad \times \left[(n-1)\frac{u_t}{r\rho} + r^{n-1}u_{tx} \right] dx ds + 2c_1(n-1) \int_0^t \int_0^M \left\{ (n-1)r^{n-2}u\rho^\theta \left(\frac{u}{r}\right)_x u_t \right. \\
 & \quad \left. - \theta r^{n-1}\rho^{\theta+1}(r^{n-1}u)_x \left(\frac{u}{r}\right)_x u_t - r^{n-1}\rho^\theta \left(\frac{u^2}{r^2}\right)_x u_t \right\} dx ds \\
 & \leq -C_{15} \int_0^t \int_0^M (\rho^{\theta+1}u_{xt}^2 + \rho^{\theta-1}u_t^2) dx ds + C \\
 (3.84) \quad & + C \int_0^t \left(\|u\|_{L^\infty}^2 + \|\rho(r^{n-1}u)_x\|_{L^\infty}^2 \right) \int_0^M [\rho^{\theta+1}u_x^2 + \rho^{\theta-1}u^2] dx ds.
 \end{aligned}$$

From (1.11), (3.17)–(3.20), and the Cauchy–Schwarz inequality, we obtain

$$\begin{aligned}
 J_2 & \leq \frac{C_{15}}{8} \int_0^t \int_0^M (\rho^{\theta-1}u_t^2 + \rho^{\theta+1}u_{xt}^2) dx ds + C \\
 (3.85) \quad & + C \int_0^t \left(\|u\|_{L^\infty}^2 + \|\rho(r^{n-1}u)_x\|_{L^\infty}^2 \right) \int_0^M [\rho^{\theta+1}u_x^2 + \rho^{\theta-1}u^2] dx ds
 \end{aligned}$$

and

$$\begin{aligned}
 J_5 & \leq \frac{C_{15}}{8} \int_0^t \int_0^M \rho^{\theta-1}u_t^2 dx ds \\
 & \quad + C \int_0^t \int_0^M (G(M_0 + x)r^{-n}|u| + |\partial_r \Delta f u| + |\partial_t \Delta f|)^2 \rho^{1-\theta} dx ds \\
 (3.86) \quad & \leq \frac{C_{15}}{8} \int_0^M r^{\alpha-2}u_t^2 dx ds + C.
 \end{aligned}$$

From (3.82)–(3.86), we have

$$\begin{aligned}
 & \int_0^M u_t^2(x, t) dx + \int_0^t \int_0^M [\rho^{\theta+1} u_{xt}^2 + \rho^{\theta-1} u_t^2] dx ds \\
 (3.87) \quad & \leq C + C \int_0^t \left(\|u\|_{L_x^\infty}^2 + \|\rho(r^{n-1}u)_x\|_{L_x^\infty}^2 \right) \int_0^M [\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2] dx ds.
 \end{aligned}$$

From (3.71)–(3.72) and (3.79), we have

$$(3.88) \quad \|\rho \partial_x(r^{n-1}u)\|_{L_x^\infty} \leq C + C \|u_t\|_{L_x^2}.$$

From (3.20), (3.71)–(3.72), and (3.87)–(3.88), we obtain

$$\begin{aligned}
 & \int_0^M u_t^2(x, t) dx + \int_0^t \int_0^M [\rho^{\theta+1} u_{xt}^2 + \rho^{\theta-1} u_t^2] dx ds \\
 & \leq C + C \int_0^t \int_0^M [\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2] dx \|u_t\|_{L_x^2}^2 ds.
 \end{aligned}$$

Using Gronwall’s inequality and the estimate (3.20), we can immediately obtain (3.80)–(3.81). \square

Proof of existence and uniqueness. If $T^* < \infty$, from Lemmas 3.9–3.13, we have, for all $t \in [0, T^*)$,

$$C^{-1}(M - x)^{1/\gamma} \leq \rho(x, t) \leq C(M - x)^{1/\gamma}, \quad x \in [0, M],$$

$$\int_0^M (M - x)^{1-\frac{\theta}{\gamma}} (\rho^\theta)_x^2 dx \leq C, \quad \int_0^M (\rho^\gamma)_x^2 dx \leq C,$$

$$\int_0^M u^2 + (M - x)^{\frac{\theta+1}{\gamma}} u_x^2 dx \leq C, \quad u(0, t) = 0$$

and

$$\int_0^M \left\{ ((2c_1 + c_2)\rho^{\theta+1}(r^{n-1}u)_x)_x - 2c_1(n - 1)\frac{u}{r}\partial_x\rho^\theta \right\}^2 dx \leq C.$$

Thus, from Theorem 3.1, there exists $T_2 > 0$ such that the free boundary problem (1.15)–(1.17) admits a unique weak solution $(\rho_2, u_2, r_2)(x, t)$ on $[0, M] \times [T^* - \frac{T_2}{2}, T^* + \frac{T_2}{2}]$, with initial data $(\rho, u, r)(x, T^* - \frac{T_2}{2})$. Using the uniqueness result in Theorem 3.1, we obtain that

$$(\tilde{\rho}, \tilde{u}, \tilde{r})(x, t) = \begin{cases} (\rho, u, r)(x, t), & t \in [0, T^* - \frac{T_2}{2}], \\ (\rho_2, u_2, r_2)(x, t), & t \in [T^* - \frac{T_2}{2}, T^* + \frac{T_2}{2}], \end{cases}$$

is a solution of the system (1.15)–(1.17), which is in contradiction with the definition of T^* . Thus, we have that $T^* = \infty$. From Lemmas 3.9–3.13, we can show that the global weak solution satisfies the regularity conditions (1.25)–(1.26) and (1.28) in Theorem 1.1.

Remark 3.2. The uniqueness of the solution of Theorem 3.1 is obtained by the energy method. Let (u_i, ρ_i, r_i) , $i = 1, 2$, be two solutions of the system (1.15)–(1.17)

satisfying the regularity conditions in Theorem 1.1. Using similar arguments as that in the uniqueness part in [3], we can obtain, for all $T > 0$,

$$\begin{aligned} & \frac{d}{dt} \int_0^M (w^2 + \rho_1^{1-\theta} \rho_2^{2\theta-4} \varrho^2 + \rho_1^\theta \rho_2^{-1} \mathcal{R}^2) dx \\ & + C^{-1} \int_0^M \rho_1^{1+\theta} \left(\rho_1 r_1^{2n-2} (\partial_x w)^2 + \frac{w^2}{r_1^2 \rho_1} \right) dx \\ & \leq C \int_0^M (w^2 + \rho_1^{1-\theta} \rho_2^{2\theta-4} \varrho^2 + \rho_1^\theta \rho_2^{-1} \mathcal{R}^2) dx, \quad t \in [0, T], \end{aligned}$$

where $(w, \varrho, \mathcal{R}) = (u_1 - u_2, \rho_1 - \rho_2, r_1 - r_2)$. Using Gronwall's inequality, we could obtain that $(u_1, \rho_1, r_1) = (u_2, \rho_2, r_2)$, a.e. $(x, t) \in [0, M] \times [0, T]$.

4. Further decay result.

LEMMA 4.1. *Let ν be a positive constant satisfying $\nu < \min\{1, \frac{2\gamma-2}{\gamma+\theta}\}$. Under the assumptions of Theorem 1.1, we obtain*

$$(4.1) \quad \left\| \rho^{\frac{\gamma+\theta}{2}}(\cdot, t) - \rho_\infty^{\frac{\gamma+\theta}{2}}(\cdot) \right\|_{L^\infty} \leq C(1+t)^{-\frac{1}{4} + \frac{\epsilon_2}{2}},$$

and

$$(4.2) \quad \|r(\cdot, t) - r_\infty(\cdot)\|_{L^\infty} \leq C_\nu(1+t)^{-\frac{1}{4}\nu + \frac{\epsilon_2\nu}{2}}$$

for all $t \geq 0$, where C_ν is a positive constant depending on ν .

Proof. From (3.48), (3.56), and (3.66), we have

$$(4.3) \quad \int_0^M \left(\rho^{\frac{\gamma+\theta}{2}} - \rho_\infty^{\frac{\gamma+\theta}{2}} \right)_x^2 dx \leq C(1+t)^{\epsilon_2}, \quad t \geq 0.$$

Combining (3.33) and the Galiardo–Nirenberg inequality $\|\phi\|_{L^\infty} \leq \|\phi\|_{L^2}^{\frac{1}{2}} \|\phi'\|_{L^2}^{\frac{1}{2}}$, we obtain

$$\left\| \rho^{\frac{\gamma+\theta}{2}} - \rho_\infty^{\frac{\gamma+\theta}{2}} \right\|_{L^\infty} \leq C(1+t)^{-\frac{1}{4} + \frac{\epsilon_2}{2}}, \quad t \geq 0.$$

From (3.17)–(3.18), (3.48), and (4.1), we have

$$\left\| \rho^{\frac{\nu(\gamma+\theta)}{2}}(\cdot, t) - \rho_\infty^{\frac{\nu(\gamma+\theta)}{2}}(\cdot) \right\|_{L^\infty} \leq C \left\| \rho^{\frac{\gamma+\theta}{2}}(\cdot, t) - \rho_\infty^{\frac{\gamma+\theta}{2}}(\cdot) \right\|_{L^\infty}^\nu \leq C(1+t)^{-\frac{\nu}{4} + \frac{\epsilon_2\nu}{2}}$$

and

$$\begin{aligned} \|r(\cdot, t) - r_\infty(\cdot)\|_{L^\infty} & \leq C \int_0^M |\rho^{-1} - \rho_\infty^{-1}| dx \\ & \leq C \left\| \rho^{\frac{\nu(\gamma+\theta)}{2}}(\cdot, t) - \rho_\infty^{\frac{\nu(\gamma+\theta)}{2}}(\cdot) \right\|_{L^\infty} \int_0^M (M-x)^{-\frac{1}{\gamma} - \frac{\nu(\gamma+\theta)}{2\gamma}} dx \\ & \leq C_\nu(1+t)^{-\frac{\nu}{4} + \frac{\epsilon_2\nu}{2}} \end{aligned}$$

for all $t \geq 0$. \square

Using similar arguments as that in Lemmas 3.6, 3.11–3.12, and 4.1 with $\nu = \frac{\gamma-1}{2\gamma}$, we can obtain the following lemma and omit the proof.

LEMMA 4.2. *Under the assumptions of Theorem 1.1, we have*

$$(4.4) \quad \int_0^M (M-x)^{\frac{\theta-1}{\gamma}} (g-g_\infty)^2 dx + \int_0^t \int_0^M [\rho_\infty^{\gamma-1} (g-g_\infty)^2 + (r-r_\infty)^2] dx ds \leq C,$$

$$(4.5) \quad \int_0^M \left\{ u^2(x,t) + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \leq C(1+t)^{-1},$$

$$(4.6) \quad \int_0^t (1+s) \|u(\cdot, s)\|_{L^\infty}^2 ds + \int_0^t \int_0^M (1+s) (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2)(x,s) dx ds \leq C,$$

$$(4.7) \quad \int_0^M (M-x)^{1-\frac{\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2 dx + \int_0^t \int_0^M (M-x)^{2-\frac{2\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2 dx ds \leq C,$$

and

$$(4.8) \quad (1+t) \int_0^M (\rho^{\theta-1} u^2 + \rho^{\theta+1} u_x^2)(x,t) dx + \int_0^t \int_0^M (1+s) u_t^2(x,s) dx ds \leq C$$

for all $t \geq 0$.

Remark 4.1. The key point is as follows: Similar to (3.32), using the estimates (3.18), (3.34), (4.2), and the condition $\epsilon_2 < \frac{\gamma-1}{2(3\gamma-1)}$, we have

$$\begin{aligned} & \int_0^t \int_0^M 2c_1(n-1)\rho^\theta \left(\frac{u}{r} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right)_x dx ds \\ & \leq C \int_0^t \int_0^M \{ |r-r_\infty| (|\rho^\theta u_x| + |\rho^{\theta-1} u|) + \rho^\theta |u(\rho^{-1} - \rho_\infty^{-1})| \} dx ds \\ & \leq C \int_0^t \int_0^M (1+s)^{1-\epsilon_2} (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2)(x,s) dx ds \\ & \quad + C \int_0^t (1+s)^{\epsilon_2-1} \|r(\cdot, t) - r_\infty(\cdot)\|_{L^\infty}^2 \int_0^M (M-x)^{\frac{\theta-1}{\gamma}} dx ds \\ & \quad + C \int_0^t (1+s)^{\epsilon_2-1} \|\rho^{\frac{\nu(\gamma+\theta)}{2}}(\cdot, t) - \rho_\infty^{\frac{\nu(\gamma+\theta)}{2}}(\cdot)\|_{L^\infty}^2 \int_0^M (M-x)^{\frac{\theta+1}{\gamma} - \frac{2}{\gamma} - \frac{\nu(\gamma+\theta)}{\gamma}} dx ds \\ & \leq C + C \int_0^t (1+s)^{\epsilon_2-1-\frac{\nu}{2}+\epsilon_2\nu} ds \leq C, \quad t \geq 0. \end{aligned}$$

Without loss of generality, we assume that $\eta \in (0, \frac{2(\gamma+\theta)}{\gamma+\theta+1})$. Let $\epsilon_4 \in (0, \frac{\gamma+\theta-1}{3(\gamma+\theta)})$ be a constant satisfying $\frac{1-\epsilon_4}{2-\frac{3\gamma+3\theta-1}{2(\gamma+\theta)}+\frac{\epsilon_4}{2}} > \frac{2(\gamma+\theta)}{\gamma+\theta+1} - \eta$. Define $\{\kappa_j\}$ and $\{\eta_j\}$ by $\eta_{j+1} = 1 + \kappa_j$, $\kappa_j = \frac{3\gamma+3\theta-1}{4(\gamma+\theta)}\eta_j - \frac{\epsilon_4}{4}\eta_j - \frac{1}{2} - \frac{\epsilon_4}{2}$, and $\eta_0 = 1$. Let N_5 be a positive integer satisfying $\eta_{N_5} > \frac{2(\gamma+\theta)}{\gamma+\theta+1} - \eta$. It is easy to see that $\eta < 2$ and $\kappa_j < 1$, $j = 0, 1, \dots, N_5$. Using similar arguments as that in Lemma 4.2, applying the induction method, we can obtain the following lemma and omit the proof.

LEMMA 4.3. *Under the assumptions of Theorem 1.1, we have*

$$(4.9) \quad \int_0^M \left\{ u^2(x,t) + (M-x)^{\frac{\gamma-1}{\gamma}} (g-g_\infty)^2 + (r-r_\infty)^2 \right\} dx \leq C_{\eta,j} (1+t)^{-\eta_j},$$

$$(4.10) \quad \int_0^t (1+s)^{\eta_j} \|u(\cdot, s)\|_{L^\infty}^2 ds + \int_0^t \int_0^M (1+s)^{\eta_j} (\rho^{\theta+1} u_x^2 + \rho^{\theta-1} u^2)(x,s) dx ds \leq C_{\eta,j},$$

$$(4.11) \quad \left\| \rho^{\frac{\gamma+\theta}{2}}(\cdot, t) - \rho_\infty^{\frac{\gamma+\theta}{2}}(\cdot) \right\|_{L^\infty} \leq C_{\eta,j}(1+t)^{-\frac{\eta_j}{4}},$$

$$(4.12) \quad \int_0^t \int_0^M (1+s)^{\kappa_j} [\rho_\infty^{\gamma-1}(g-g_\infty)^2 + (r-r_\infty)^2] dx ds \leq C_{\eta,j},$$

$$(4.13) \quad (1+t)^{\eta_j} \int_0^M (\rho^{\theta-1}u^2 + \rho^{\theta+1}u_x^2)(x, t) dx + \int_0^t \int_0^M (1+s)^{\eta_j} u_t^2(x, s) dx ds \leq C_{\eta,j},$$

and

$$(4.14) \quad \|u(\cdot, t)\|_{L^\infty} \leq C_{\eta,j}(1+t)^{-\frac{\eta_j}{2}}$$

for all $t \geq 0$, $j = 0, \dots, N_5$, where $C_{\eta,j}$ is a positive constant depending on η and j .

Remark 4.2. The main difficulty is to show (4.12) with $j = k$, when (4.9)–(4.11) hold with $j = k$. From (3.27)–(3.28), we have

$$\begin{aligned} & \int_0^T \int_0^M (1+t)^{\kappa_k} \left[A(\rho_\infty^\gamma - \rho^\gamma)(\rho^{-1} - \rho_\infty^{-1}) \right. \\ & \quad \left. + G(M_0 + x)(r^{2-2n} - r_\infty^{2-2n}) \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right] dx dt \\ &= - \int_0^T \int_0^M (1+t)^{\kappa_k} \frac{u_t}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx dt \\ & \quad - \int_0^T \int_0^M (1+t)^{\kappa_k} \Delta f r^{1-n} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx dt \\ & \quad + \int_0^T \int_0^M (1+t)^{\kappa_k} (2c_1 + c_2) \rho^{1+\theta} \partial_x (r^{n-1}u) (\rho_\infty^{-1} - \rho^{-1}) dx dt \\ & \quad + \int_0^T \int_0^M (1+t)^{\kappa_k} 2c_1(n-1) \rho^\theta \left(\frac{u}{r} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) \right)_x dx dt \\ (4.15) \quad & := \sum_{i=1}^4 Q_i, \quad T > 0 \end{aligned}$$

and

$$(4.16) \quad \text{L.H.S. of (4.15)} \geq C_{12} \int_0^T \int_0^M (1+t)^{\kappa_k} [\rho_\infty^{\gamma-1}(g-g_\infty)^2 + (r-r_\infty)^2] dx dt.$$

Similar to (3.29)–(3.32), applying the estimates (3.17)–(3.20), (4.4), integration by parts, the Cauchy–Schwarz inequality, and the fact that $\kappa_j = \frac{3\gamma+3\theta-1}{4(\gamma+\theta)}\eta_j - \frac{\epsilon_4}{4}\eta_j - \frac{1}{2} - \frac{\epsilon_4}{2} < \eta_j$, we can estimate Q_i as follows:

$$\begin{aligned} Q_1 &\leq - \int_0^M (1+t)^{\kappa_k} \frac{u}{r^{n-1}} \left(\frac{r^n}{n} - \frac{r_\infty^n}{n} \right) dx \Big|_0^T + C \int_0^T (1+t)^{\kappa_k} \|u\|_{L^\infty}^2 dt \\ & \quad + C \int_0^T \int_0^M (1+t)^{\kappa_k-1} |u| |r - r_\infty| dx dt \\ (4.17) \quad &\leq C, \end{aligned}$$

$$Q_2 \leq \frac{C_{12}}{6} \int_0^T \int_0^M (1+t)^{\kappa_k} (r - r_\infty)^2 dx dt + C \int_0^T f_1^2(1+t)^{\kappa_k} dt$$

$$(4.18) \quad \leq \frac{C_{12}}{6} \int_0^T \int_0^M (1+t)^{\kappa_k} (r-r_\infty)^2 dxdt + C,$$

$$(4.19) \quad \begin{aligned} Q_3 &\leq C \int_0^T \int_0^M (1+t)^{\eta_k} \rho^{1+\theta} (r^{n-1} u_x)_x^2 dxdt \\ &\quad + C \int_0^T \int_0^M (1+t)^{2\kappa_k - \eta_k - \frac{\eta_k \nu_1}{2}} (M-x)^{\frac{\theta-1}{\gamma} - \frac{\nu_1(\gamma+\theta)}{\gamma}} dxdt \\ &\leq C, \end{aligned}$$

where $\nu_1 = \frac{\gamma+\theta-1}{\gamma+\theta} - \epsilon_4$,

$$\|r(\cdot, t) - r_\infty(\cdot)\|_{L^\infty} \leq C(1+t)^{\frac{\nu_1 \eta_k}{4}},$$

and

$$(4.20) \quad \begin{aligned} Q_4 &\leq C \int_0^T \int_0^M (1+t)^{\eta_k} (\rho^{1+\theta} (r^{n-1} u_x)_x^2 + \rho^{\theta-1} u^2) dxdt \\ &\quad + C \int_0^T \int_0^M (1+t)^{2\kappa_k - \eta_k - \frac{\eta_k \nu_1}{2}} (M-x)^{\frac{\theta-1}{\gamma}} dxdt \\ &\leq C. \end{aligned}$$

From (4.15)–(4.20), we finish the proof of (4.12) with $j = k$.

From (4.12) with $j = N_5$, using similar arguments as that in Lemma 3.11, we can obtain the following lemma and omit the proof.

LEMMA 4.4. *Under the assumptions of Theorem 1.1, we have*

$$(4.21) \quad \int_0^M (M-x)^{2-\frac{2\theta}{\gamma}} (\rho^\theta - \rho_\infty^\theta)_x^2 dx \leq C_\eta (1+t)^{\eta - \frac{\gamma+\theta-1}{\gamma+\theta}},$$

and

$$(4.22) \quad \|\rho^\gamma(\cdot, t) - \rho_\infty^\gamma(\cdot)\|_{L^\infty} \leq C(1+t)^{\frac{\eta}{2} - \frac{3\gamma+3\theta-1}{4(\gamma+\theta+1)}}, \quad t \geq 0.$$

Thus, we finish the proof of Theorem 1.1.

REFERENCES

- [1] D. BRESCH AND B. DESJARDINS, *On the construction of approximate solutions for the 2D viscous shallow water model and for compressible Navier-Stokes models*, J. Math. Pure Appl., 86 (2006), pp. 362–368.
- [2] G. Q. CHEN AND M. KRATKA, *Global solutions to the Navier-Stokes equations for compressible heat-conducting flow with symmetry and free boundary*, Comm. Partial Differential Equations, 27 (2002), pp. 907–943.
- [3] P. CHEN AND T. ZHANG, *A vacuum problem for multidimensional compressible Navier-Stokes equations with degenerate viscosity coefficients*, Commun. Pure Appl. Anal., 7 (2008), pp. 987–1016.
- [4] B. DUCOMET AND A. A. ZLOTNIK, *Viscous compressible barotropic symmetric flows with free boundary under general mass force. I. uniform-in-time bounds and stabilization*, Math. Methods Appl. Sci., 28 (2005), pp. 827–863.
- [5] B. DUCOMET AND A. A. ZLOTNIK, *Lyapunov functional method for 1D radiative and reactive viscous gas dynamics*, Arch. Ration. Mech. Anal., 177 (2005), pp. 185–229.
- [6] H. GRAD, *Asymptotic theory of the Boltzmann equation II*, in Rarefied Gas Dynamics, Vol. 1. J. Laurmann, ed., Academic Press, New York, 1963, pp. 26–59.

- [7] P. HARTMAN, *Ordinary Differential Equations*, Wiley, New York, 1964.
- [8] D. HOFF AND D. SERRE, *The failure of continuous dependence on initial data for the Navier-Stokes equations of compressible flow*, SIAM J. Appl. Math., 51 (1991), pp. 887–898.
- [9] T. KOBAYASHI AND Y. SHIBATA, *Decay estimates of solutions for the equations of motion of compressible viscous and heat-conductive gases in an exterior domain in R^3* , Comm. Math. Phys., 200 (1999), pp. 621–659.
- [10] P. L. LIONS, *Mathematical Topics in Fluid Mechanics*, Vol. 1-2, Oxford University Press, New York, 1996, 1998.
- [11] T. P. LIU, Z. P. XIN, AND T. YANG, *Vacuum states of compressible flow*, Discrete Contin. Dyn. Syst., 4 (1998), pp. 1–32.
- [12] A. MATSUMURA AND S. YANAGI, *Uniform boundedness of the solutions for a one-dimensional isentropic model system of a compressible viscous gas*, Comm. Math. Phys., 175 (1996), pp. 259–274.
- [13] Š. MATUŠŮ-NEČASOVÁ, M. OKADA, AND T. MAKINO, *Free boundary problem for the equation of spherically symmetric motion of viscous gas (II)–(III)*, Japan J. Indust. Appl. Math., 12 (1995) pp. 195–203; 14 (1997), pp. 199–213.
- [14] M. OKADA AND T. MAKINO, *Free boundary value problems for the equation of spherically symmetrical motion of viscous gas*, Japan J. Appl. Math., 10 (1993), pp. 219–235.
- [15] M. OKADA, Š. MATUŠŮ-NEČASOVÁ, AND T. MAKINO, *Free boundary problem for the equation of one-dimensional motion of compressible gas with density-dependent viscosity*, Ann. Univ. Ferrara Sez. VII (N.S.), 48 (2002), pp. 1–20.
- [16] I. STRAŠKRABA AND A. A. ZLOTNIK, *Global behavior of 1d-viscous compressible barotropic fluid with a free boundary and large data*, J. Math. Fluid Mech., 5 (2003), pp. 119–143.
- [17] V. A. VAIGANT AND A. V. KAZHIKHOV, *On existence of global solutions to the two-dimensional Navier-Stokes equations for a compressible viscous fluid*, Siberian Math. J., 36 (1995), pp. 1108–1141.
- [18] S. UKAI, T. YANG, AND H. J. ZHAO, *Convergence rate for the compressible Navier-Stokes equations with external force*, J. Hyperbolic Differ. Equ., 3 (2006), pp. 561–574.
- [19] S. W. VONG, T. YANG, AND C. J. ZHU, *Compressible Navier-Stokes equations with degenerate viscosity coefficient and vacuum (II)*, J. Differential Equations, 192 (2003), pp. 475–501.
- [20] Z. P. XIN, *Blow-up of smooth solution to the compressible Navier-Stokes equations with compact density*, Comm. Pure Appl. Math., 51 (1998), pp. 229–240.
- [21] T. YANG AND C. J. ZHU, *Compressible Navier-Stokes equations with degenerate viscosity coefficient and vacuum*, Comm. Math. Phys., 230 (2002), pp. 329–363.
- [22] T. ZHANG AND D. Y. FANG, *Global behavior of compressible Navier-Stokes equations with a degenerate viscosity coefficient*, Arch. Ration. Mech. Anal., 182 (2006), pp. 223–253.
- [23] T. ZHANG AND D. Y. FANG, *Global behavior of spherically symmetric Navier-Stokes equations with density-dependent viscosity*, J. Differential Equations, 236 (2007), pp. 293–341.
- [24] A. A. ZLOTNIK AND B. DUCOMET, *The stabilization rate and stability of viscous compressible barotropic symmetric flows with a free boundary for a general mass force*, Sb. Math., 196 (2005), pp. 1745–1799.

COMPACTLY SUPPORTED SYMMETRIC C^∞ WAVELETS WITH SPECTRAL APPROXIMATION ORDER*

BIN HAN[†] AND ZUOWEI SHEN[‡]

Abstract. In this paper, we obtain symmetric C^∞ real-valued tight wavelet frames in $L_2(\mathbb{R})$ with compact support and the spectral frame approximation order. Furthermore, we present a family of symmetric compactly supported C^∞ orthonormal complex wavelets in $L_2(\mathbb{R})$. A complete analysis of nonstationary tight wavelet frames and orthonormal wavelet bases in $L_2(\mathbb{R})$ is given.

Key words. symmetric tight wavelet frames, spectral frame approximation order, symmetric orthonormal complex wavelets, nonstationary cascade algorithm, nonstationary C^∞ wavelets

AMS subject classifications. 42C40, 41A25, 41A05, 42C05

DOI. 10.1137/060675009

1. Introduction. In this paper, we are interested in symmetric compactly supported C^∞ tight wavelet frames with the spectral frame approximation order. Since it is impossible to achieve all these properties under the framework of stationary tight wavelet frames, it is natural for us to consider nonstationary tight wavelet frames, in particular, nonstationary tight wavelet frames derived from nonstationary multiresolution analysis by the new (nonstationary) unitary extension principle.

We start with a family of 2π -periodic trigonometric polynomials $\widehat{a}_j, j \in \mathbb{N}$, and their associated nonstationary refinable functions (or tempered distributions) $\widehat{\phi}_{j-1}, j \in \mathbb{N}$, defined by

$$(1.1) \quad \widehat{\phi}_{j-1}(\xi) := \widehat{a}_j(\xi/2)\widehat{\phi}_j(\xi/2) = \prod_{n=1}^{\infty} \widehat{a}_{n+j-1}(2^{-n}\xi), \quad \xi \in \mathbb{R}, j \in \mathbb{N},$$

where the 2π -periodic trigonometric polynomials $\widehat{a}_j, j \in \mathbb{N}$, are called refinement *masks*. Here, the Fourier transform \widehat{f} of a function $f \in L_1(\mathbb{R})$ used in this paper is defined to be $\widehat{f}(\xi) := \int_{\mathbb{R}} f(t)e^{-it\xi} dt$ and can be naturally extended to square integrable functions and tempered distributions.

The stationary multiresolution analysis corresponds to the case that all the masks \widehat{a}_j are the same; therefore, all the functions $\widehat{\phi}_j$ are the same, and, in particular, $\widehat{\phi}_0(\xi) = \widehat{a}_1(\xi/2)\widehat{\phi}_1(\xi/2)$. We say that a function $\phi : \mathbb{R} \mapsto \mathbb{C}$ is refinable with a 2π -periodic trigonometric polynomial refinement mask \widehat{a} if $\widehat{\phi}(\xi) = \widehat{a}(\xi/2)\widehat{\phi}(\xi/2)$. The frame generators ψ^ℓ are generally obtained from the refinable function ϕ via $\widehat{\psi}^\ell(\xi) = \widehat{b}^\ell(\xi/2)\widehat{\phi}(\xi/2)$ for some 2π -periodic trigonometric polynomials \widehat{b}^ℓ with some desirable properties.

*Received by the editors November 15, 2006; accepted for publication (in revised form) April 3, 2008; published electronically September 8, 2008. This research was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC Canada) under grant RGP 228051 and in part by several grants from the National University of Singapore.

<http://www.siam.org/journals/sima/40-3/67500.html>

[†]Department of Mathematical and Statistical Sciences, University of Alberta, Edmonton T6G 2G1, AB, Canada (bhan@math.ualberta.ca, <http://www.ualberta.ca/~bhan>).

[‡]Department of Mathematics, National University of Singapore, Singapore (matzuows@nus.edu.sg, <http://www.math.nus.edu.sg/~matzuows>).

A tight wavelet frame in $L_2(\mathbb{R})$ (in the stationary case) is generated by the integer translates and dyadic dilates of a finite set of elements in $L_2(\mathbb{R})$. More precisely, we say that $\{\psi^1, \dots, \psi^L\}$ generates a (normalized) *tight wavelet frame* in $L_2(\mathbb{R})$ if

$$(1.2) \quad \|f\|_{L_2(\mathbb{R})}^2 = \sum_{\ell=1}^L \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} |\langle f, \psi_{j,k}^\ell \rangle|^2 \quad \forall f \in L_2(\mathbb{R}),$$

where $\psi_{j,k}^\ell := 2^{j/2} \psi^\ell(2^j \cdot -k)$ and $\langle f, g \rangle := \int_{\mathbb{R}} f(t) \overline{g(t)} dt$. As a redundant wavelet system, tight frame wavelet systems are easier to design and provide more flexibilities in applications than orthonormal wavelet bases, especially in image inpainting (see [1, 2, 3, 4, 5] for details). Because of this, tight wavelet frames have been extensively studied in the literature; to mention only a few, see [6, 7, 8, 14, 16, 18, 19, 21, 25, 26, 27, 33] and the many references therein.

Tight wavelet frames obtained from refinable functions are of particular interest, due to their associated multiresolution structure (which we refer to as MRA-based) and fast frame algorithms. Constructions of tight wavelet frames from a refinable function can be done by the unitary extension principle in [33]. In fact, many tight wavelet frames have been constructed in [6, 16, 33]. Later, by using the more general oblique extension principle, which is independently developed in [7, 14], more tight wavelet frames with various desirable properties have been obtained in [7, 14, 21, 25, 26, 27] and many other references therein.

For the stationary case, it already has been pointed out in [13] that there does not exist a compactly supported refinable function ϕ with a 2π -periodic trigonometric polynomial refinement mask such that ϕ belongs to $C^\infty(\mathbb{R})$. Hence, it is impossible to obtain MRA-based compactly supported (stationary) tight wavelet frames in $L_2(\mathbb{R})$ whose generators are in $C^\infty(\mathbb{R})$. However, it is shown in [10] that, by using the class of masks for orthonormal refinable functions of [12] whose integer shifts form an orthonormal system, one can obtain a family of nonstationary refinable functions such that every nonstationary refinable function belongs to $C^\infty(\mathbb{R})$ and its integer shifts still form an orthonormal system in $L_2(\mathbb{R})$. For this family of nonstationary refinable functions, a C^∞ nonstationary orthonormal wavelet basis in $L_2(\mathbb{R})$ is derived in [10]. In fact, ideas of generating a class of nonstationary refinable functions in $C^\infty(\mathbb{R})$ from a given family of masks for stationary refinable functions have already been discussed in [15, 35]. One such example is the *up-function* [10, 15, 35] generated from the family of masks for the B-splines. Let $\hat{a}_j(\xi) = 2^{-j}(1 + e^{-i\xi})^j$, $j \in \mathbb{N}$, be the mask for the B-spline of order j and define ϕ_{j-1} , $j \in \mathbb{N}$, as in (1.1). Then all ϕ_{j-1} , $j \in \mathbb{N}$, are compactly supported C^∞ functions. In particular, the function ϕ_0 is supported on $[0, 2]$ (see [10, 15, 35]).

Motivated by the interesting work of Cohen and Dyn [10] and equipped with the pseudosplines (a more general class of refinable functions containing B-splines, interpolatory refinable functions, and Daubechies orthonormal refinable functions in [12] as special cases), together with the idea of the unitary extension principle, we establish the analysis needed here for constructing nonstationary $C^\infty(\mathbb{R})$ tight wavelet frames in $L_2(\mathbb{R})$ with desirable properties, especially the symmetry property, which cannot be achieved by real-valued orthonormal dyadic refinable functions. As we will see, the construction more or less follows the idea of the unitary extension principle for the stationary case, while the main analysis of this paper is somehow different from that of [10]. For example, in the orthonormal wavelet case, the approximation order of the truncated wavelet series in [10] is the same as that of the (nonstationary)

multiresolution analysis, while in the tight wavelet frame case, they are different, even for the stationary case, as shown in [14].

Next, we briefly describe ideas of the construction of tight wavelet frames. Although one of our major objectives of this paper is to use the family of refinement masks for pseudosplines to construct tight wavelet frames and to provide the corresponding analysis, the construction in this paper is given for the general setting.

We start with 2π -periodic measurable functions $\widehat{a}_j, j \in \mathbb{N}$, as a sequence of refinement masks. To make the idea of the unitary extension principle work, it is necessary to require that for every $j \in \mathbb{N}$, the mask \widehat{a}_j should satisfy

$$(1.3) \quad |\widehat{a}_j(\xi)|^2 + |\widehat{a}_j(\xi + \pi)|^2 \leq 1 \quad \text{a.e. } \xi \in \mathbb{R}.$$

Since we are interested only in compactly supported tight wavelet frames, it is natural to start with compactly supported refinable functions ϕ_j , which, in turn, require that the degrees of the trigonometric polynomials \widehat{a}_j do not increase too fast. For a 2π -periodic trigonometric polynomial \widehat{a} , we denote $\text{deg}(\widehat{a})$ the smallest nonnegative integer such that its Fourier coefficients of \widehat{a} vanish outside $[-\text{deg}(\widehat{a}), \text{deg}(\widehat{a})]$. We note that $\text{deg}(\widehat{a})$ defined here is somewhat slightly different from the usual definition of the degree of a trigonometric polynomial; $\text{deg}(\widehat{a})$ here is the minimal integer k such that $[-k, k]$ contains the support of the Fourier coefficients of both \widehat{a} and $\widehat{a}(\cdot)$. For ϕ_0 in (1.1) to be compactly supported, by a simple calculation, it is very natural to require [10] that

$$(1.4) \quad \sum_{j=1}^{\infty} 2^{-j} \text{deg}(\widehat{a}_j) < \infty.$$

With (1.3) and (1.4), under the condition that $\sum_{j=1}^{\infty} |\widehat{a}_j(0) - 1| < \infty$, it can be proved that all the corresponding refinable functions ϕ_{j-1} in (1.1) are well-defined compactly supported functions in $L_2(\mathbb{R})$.

Wavelet functions $\psi_{j-1}^\ell, j \in \mathbb{N}$ and $\ell \in \{1, \dots, \mathcal{J}_j\}$, are obtained from ϕ_j by

$$(1.5) \quad \widehat{\psi_{j-1}^\ell}(\xi) := \widehat{b_j^\ell}(\xi/2)\widehat{\phi_j}(\xi/2), \quad \ell = 1, \dots, \mathcal{J}_j,$$

where \mathcal{J}_j are positive integers and each $\widehat{b_j^\ell}, \ell = 1, \dots, \mathcal{J}_j$, is called a (high-pass) wavelet mask. Denote $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$. We say that $\{\phi_0\} \cup \{\psi_j^\ell : j \in \mathbb{N}_0, \ell = 1, \dots, \mathcal{J}_{j+1}\}$ generates a *nonstationary tight wavelet frame* in $L_2(\mathbb{R})$ if

$$(1.6) \quad \{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;k}^\ell := 2^{j/2}\psi_j^\ell(2^j \cdot -k) : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, \dots, \mathcal{J}_{j+1}\}$$

is a tight frame of $L_2(\mathbb{R})$; that is, the following holds:

$$(1.7) \quad \|f\|_{L_2(\mathbb{R})}^2 = \sum_{k \in \mathbb{Z}} |\langle f, \phi_0(\cdot - k) \rangle|^2 + \sum_{j=0}^{\infty} \sum_{\ell=1}^{\mathcal{J}_{j+1}} \sum_{k \in \mathbb{Z}} |\langle f, \psi_{j;k}^\ell \rangle|^2 \quad \forall f \in L_2(\mathbb{R}).$$

We say that ψ_j^ℓ has ν *vanishing moments* if $\widehat{\psi_j^\ell}^{(n)}(0) = 0$ for all $n = 0, \dots, \nu - 1$, where $\widehat{\psi_j^\ell}^{(n)}$ denotes the n th derivative of $\widehat{\psi_j^\ell}$. It is clear that (1.7) is equivalent to

$$(1.8) \quad f = \sum_{k \in \mathbb{Z}} \langle f, \phi_0(\cdot - k) \rangle \phi_0(\cdot - k) + \sum_{j=0}^{\infty} \sum_{\ell=1}^{\mathcal{J}_{j+1}} \sum_{k \in \mathbb{Z}} \langle f, \psi_{j;k}^\ell \rangle \psi_{j;k}^\ell, \quad f \in L_2(\mathbb{R}).$$

The frame approximation operators $Q_n, n \in \mathbb{N}$, associated with the truncation of the tight wavelet frame in (1.6) at level n , are defined to be

$$(1.9) \quad Q_n(f) := \sum_{k \in \mathbb{Z}} \langle f, \phi_0(\cdot - k) \rangle \phi_0(\cdot - k) + \sum_{j=0}^{n-1} \sum_{\ell=1}^{\mathcal{J}_{j+1}} \sum_{k \in \mathbb{Z}} \langle f, \psi_{j,j,k}^\ell \rangle \psi_{j,j,k}^\ell, \quad f \in L_2(\mathbb{R}).$$

For $\nu \geq 0$, we denote $W_2^\nu(\mathbb{R})$ the Sobolev space of all functions $f \in L_2(\mathbb{R})$ such that

$$(1.10) \quad \|f\|_{W_2^\nu(\mathbb{R})}^2 := \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\hat{f}(\xi)|^2 d\xi < \infty.$$

The Sobolev seminorm $|f|_{W_2^\nu(\mathbb{R})}$ is defined to be

$$(1.11) \quad |f|_{W_2^\nu(\mathbb{R})}^2 := \int_{\mathbb{R}} |\xi|^{2\nu} |\hat{f}(\xi)|^2 d\xi, \quad f \in W_2^\nu(\mathbb{R}).$$

The unitary extension principle provides a sufficient condition on the wavelet masks \hat{a}_j and $b_j^\ell, \ell = 1, \dots, \mathcal{J}_j$, so that, with ψ_{j-1}^ℓ defined in (1.5), the wavelet system in (1.6) forms a tight frame in $L_2(\mathbb{R})$. Altogether, we have the following result on nonstationary tight wavelet frames in $L_2(\mathbb{R})$.

THEOREM 1.1. *Let $\hat{a}_j, j \in \mathbb{N}$, be 2π -periodic trigonometric polynomials with $\hat{a}_j(0) = 1$ for all $j \in \mathbb{N}$. If (1.3) and (1.4) hold, letting ϕ_j and $\psi_{j-1}^\ell, j \in \mathbb{N}$ and $\ell \in \{1, \dots, \mathcal{J}_j\}$, be defined in (1.1) and (1.5), respectively, then the following hold:*

- (i) *All functions $\phi_{j-1}, j \in \mathbb{N}$, are well-defined compactly supported functions in $L_2(\mathbb{R})$.*
- (ii) *If $b_j^\ell, j \in \mathbb{N}$ and $\ell \in \{1, \dots, \mathcal{J}_j\}$, are 2π -periodic trigonometric polynomials satisfying*

$$(1.12) \quad \begin{aligned} |\hat{a}_j(\xi)|^2 + \sum_{\ell=1}^{\mathcal{J}_j} |b_j^\ell(\xi)|^2 &= 1 \quad \text{and} \\ \hat{a}_j(\xi) \overline{\hat{a}_j(\xi + \pi)} + \sum_{\ell=1}^{\mathcal{J}_j} b_j^\ell(\xi) \overline{b_j^\ell(\xi + \pi)} &= 0, \end{aligned}$$

then the wavelet system in (1.6) is a compactly supported tight wavelet frame in $L_2(\mathbb{R})$.

- (iii) *If, in addition to (1.12), we assume that*

$$(1.13) \quad \deg(\hat{a}_j) = O(j^\alpha 2^{\beta j}) \quad \text{as } j \rightarrow \infty \quad \text{for some } \alpha \geq 0, 0 \leq \beta < 1$$

and assume that there exist a positive number $\nu \in \frac{1}{2}\mathbb{N}$ and a positive integer N such that $1 - |\hat{a}_j(\xi)|^2$ has a zero of order 2ν at $\xi = 0$ for all $j \geq N$, that is, for $j \geq N$,

$$(1.14) \quad |\hat{a}_j(\xi)|^2 = 1 + O(|\xi|^{2\nu}), \quad \xi \rightarrow 0,$$

then there exists a positive constant C , independent of f and n , such that

$$(1.15) \quad \begin{aligned} \|f - Q_n(f)\|_{L_2(\mathbb{R})} &\leq C n^{\nu\alpha} 2^{-\nu(1-\beta)n} |f|_{W_2^\nu(\mathbb{R})} \\ &\forall f \in W_2^\nu(\mathbb{R}) \quad \text{and } n \geq N, \end{aligned}$$

where the linear operators Q_n are defined in (1.9).

Item (ii) of Theorem 1.1 is called the unitary extension principle for the nonstationary case. Theorem 1.1 will be proved in section 4. As we shall see in section 4, the main effort there is to prove items (i) and (iii) of Theorem 1.1. To show item (ii) of Theorem 1.1, one needs to show the convergence of the frame series in the right side of (1.8) to the function f in $L_2(\mathbb{R})$. When (1.13) holds, the convergence of the frame series follows from (iii) by observing that the masks in item (ii) satisfy (1.14) for $\nu = 1/2$. Furthermore, a refined analysis establishes the convergence of the frame series even without assuming (1.13).

We further remark that (1.12) guarantees the multiresolution frame decomposition algorithm whose proof can be straightforwardly verified and is more or less known. In fact, Theorem 1.1 generalizes the unitary extension principle from the stationary case in [33] to the general nonstationary case. It is clear that, similar to the stationary case, for every fixed $j \in \mathbb{N}$, in order to construct a set of 2π -periodic trigonometric polynomials \widehat{b}_j^ℓ , $\ell = 1, \dots, \mathcal{J}_j$, derived from the mask \widehat{a}_j so that (1.12) is satisfied, the mask \widehat{a}_j must satisfy (1.3). Hence, (1.3) is a necessary and sufficient condition to make (1.12) hold, as we shall see later in this section.

Nonstationary spline tight wavelet frames using the oblique extension principle developed in [7, 14] have been systematically studied in Chui, He, and Stöckler [8] recently. There, they considered an even more general nonstationary setting; i.e., it is not even shift-invariant at each level. Since the oblique extension principle is a generalization of the unitary extension principle, the proof of [8] might be modified to prove item (ii) of Theorem 1.1. However, this at most leads to the conclusion $Q_n(f) \rightarrow f$ in $L_2(\mathbb{R})$. Our approach of item (ii) is beyond the proof of the tight frame property itself in (1.8). Instead, we analyze the approximation power of the truncated tight wavelet frame series as stated in item (iii) of Theorem 1.1. As a consequence of this analysis, we obtain the tight frame property stated in item (ii) of Theorem 1.1. Finally, we remark that a systematic study of general nonstationary wavelet frames that may not be MRA-based was given in [34]

Following [14], we say that a tight wavelet frame $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;j,k}^\ell : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, \dots, \mathcal{J}_{j+1}\}$ provides *frame approximation order* ν if there exist a positive constant C , independent of f and n , and a positive integer N such that

$$(1.16) \quad \|f - Q_n(f)\|_{L_2(\mathbb{R})} \leq C 2^{-\nu n} \|f\|_{W_2^\nu(\mathbb{R})} \quad \forall f \in W_2^\nu(\mathbb{R}) \quad \text{and} \quad n \geq N.$$

We say that a tight wavelet frame provides *the spectral frame approximation order* if it provides frame approximation order ν for any positive integer ν . Here, we point out that for the frame approximation order discussed in this paper, the constant C in (1.15) of Theorem 1.1 and the constant C in (1.16) of Theorems 1.2 and 1.4 can be explicitly obtained.

In section 4, we shall study when $Q_n(f)$ approaches f with an approximation order ν as $n \rightarrow \infty$. As a consequence, we prove item (ii) of Theorem 1.1 and (1.8) by showing that $Q_n(f) \rightarrow f$ in $L_2(\mathbb{R})$ as $n \rightarrow \infty$ for every $f \in L_2(\mathbb{R})$, provided that the conditions in Theorem 1.1 are satisfied. The approximation order of $Q_n(f)$ was not studied in [10], since it was not needed there. In fact, since only orthonormal wavelet systems were considered in [10], the associated operators Q_n become orthogonal projections and attain the approximation order provided by the nonstationary multiresolution analysis. Therefore, one only needs to understand the conditions under which $Q_n(f) \rightarrow f$ in $L_2(\mathbb{R})$ as $n \rightarrow \infty$ in the orthonormal wavelet case. Nevertheless, our approach here applies to this special case as well and simplifies the conditions given in [10]. The approximation order of $Q_n(f)$ was not studied in [8] either, since it is a

more challenging problem in its more general setting of [8]. For the stationary case, it is evident that (1.13) holds with $\alpha = \beta = 0$ and, consequently, the notion of the frame approximation power in (1.15) agrees with that of the frame approximation order in (1.16). However, we shall present an example of nonstationary tight wavelet frames derived from the up-function (see Theorem 1.3) to demonstrate that (1.15) holds with $\nu = 2$, $\alpha = 1$, and $\beta = 0$, while (1.16) fails for any $\nu > 0$; that is, this particular nonstationary tight wavelet frame has a “weak” frame approximation order 2 in the sense of (1.15), but it does not have any “strong” frame approximation order in the sense of (1.16).

Finally, we note that the 2π -periodic trigonometric polynomial wavelet masks \widehat{b}_j^ℓ , $j \in \mathbb{N}$ and $\ell \in \{1, \dots, \mathcal{J}_j\}$, can be constructed from the masks \widehat{a}_j by many ways provided that the refinement masks $\widehat{a}_j, j \in \mathbb{N}$, satisfy (1.3). Here is one such construction modified from the stationary case of [6] (also cf. [16, 27, 25]). For every $j \in \mathbb{N}$, from the mask \widehat{a}_j with real coefficients and satisfying (1.3), define

$$\begin{aligned}
 \widehat{b}_j^1(\xi) &:= e^{-i\xi} \overline{\widehat{a}_j(\xi + \pi)}, \\
 \widehat{b}_j^2(\xi) &:= 2^{-1} [A_j(\xi) + e^{-i\xi} \overline{A_j(\xi)}], \\
 \widehat{b}_j^3(\xi) &:= 2^{-1} [A_j(\xi) - e^{-i\xi} \overline{A_j(\xi)}],
 \end{aligned}
 \tag{1.17}$$

where A_j is a π -periodic trigonometric polynomial with real coefficients such that

$$|A_j(\xi)|^2 = 1 - |\widehat{a}_j(\xi)|^2 - |\widehat{a}_j(\xi + \pi)|^2.$$

Then, $\widehat{a}_j, \widehat{b}_j^1, \widehat{b}_j^2$, and \widehat{b}_j^3 , $j \in \mathbb{N}$, satisfy (1.12) with $\mathcal{J}_j = 3$. Furthermore, the corresponding wavelets defined by (1.5) using masks in (1.17) are symmetric or antisymmetric whenever ϕ_j is symmetric.

After establishing Theorem 1.1, we focus on constructing nonstationary $C^\infty(\mathbb{R})$ wavelets derived from a family of refinement masks for pseudosplines. Pseudosplines (of type I) were first introduced in [14] and [36] to improve the approximation order of truncated tight wavelet frame series for the tight wavelet frame system obtained by the unitary extension principle. The pseudosplines in [14] are generally not symmetric. The pseudosplines of type II are symmetric and were introduced in [16]. Since we are aiming at constructing symmetric tight wavelet frames, we will use pseudosplines of type II. For positive integers $m, l \in \mathbb{N}$, throughout the paper we denote

$$P_{m,l}(x) := \sum_{j=0}^{l-1} \binom{m+j-1}{j} x^j = \sum_{j=0}^{l-1} \frac{(m+j-1)!}{j!(m-1)!} x^j, \quad x \in \mathbb{R}.
 \tag{1.18}$$

The masks for pseudosplines of type II with order (m, l) [16] are given by

$$\widehat{a}_{m,l}(\xi) := \cos^{2m}(\xi/2) P_{m,l}(\sin^2(\xi/2)), \quad m \in \mathbb{N}, l = 1, \dots, m.
 \tag{1.19}$$

Since it is evident that $\widehat{a}_{m,l}(\xi) \geq 0$ for all $\xi \in \mathbb{R}$, the mask $\widehat{a}_{m,l}^I$, for the pseudospline of type I with order (m, l) introduced in [14] and [36], is obtained by taking the square root of the mask $\widehat{a}_{m,l}$ in (1.19) for the pseudospline of type II with order (m, l) using the Fejér–Riesz lemma such that

$$|\widehat{a}_{m,l}^I(\xi)|^2 = \widehat{a}_{m,l}(\xi), \quad \xi \in \mathbb{R}.
 \tag{1.20}$$

While the pseudosplines of type II and their masks in (1.19) are symmetric, their type I counterparts usually do not have symmetry. For the case $l = 1$, the corresponding refinable pseudosplines are B-splines for both types. For the case $l = m$, the corresponding refinable pseudospline ϕ of type I with mask $\widehat{a_{m,m}^I}$ in (1.20) has orthonormal integer shifts (i.e., $\{\phi(\cdot - k) : k \in \mathbb{Z}\}$ is an orthonormal system in $L_2(\mathbb{R})$), and the corresponding refinable pseudospline ϕ of type II with mask $\widehat{a_{m,m}}$ in (1.19) is interpolatory (i.e., $\phi(0) = 1$ and $\phi(k) = 0$ for all $k \in \mathbb{Z} \setminus \{0\}$). It is easy to verify that the condition in (1.3) is satisfied for all the masks for pseudosplines of type I and type II (e.g., see [14, 16]).

Our construction here employs masks $\widehat{a_{m,l}}$ in (1.19) for pseudosplines of type II, since we are interested in constructing symmetric tight wavelet frames. We have the following result on symmetric C^∞ tight wavelet frames in $L_2(\mathbb{R})$ with compact support and the spectral frame approximation order.

THEOREM 1.2. *Let $\widehat{a_j} := \widehat{a_{m_j,l_j}}$ be defined in (1.19), where $1 \leq l_j \leq m_j$ and m_j ($j \in \mathbb{N}$) are positive integers satisfying*

$$(1.21) \quad \lim_{j \rightarrow \infty} m_j = \infty \quad \text{and} \quad \sum_{j=1}^{\infty} 2^{-j} m_j < \infty.$$

For $j \in \mathbb{N}$, define ϕ_{j-1} as in (1.1) and $\psi_{j-1}^1, \psi_{j-1}^2$, and ψ_{j-1}^3 as in (1.5) with the wavelet masks $\widehat{b_j^1}, \widehat{b_j^2}$, and $\widehat{b_j^3}$ being derived from $\widehat{a_j}$ in (1.17). Then the following hold:

- (1) Each nonstationary refinable function $\phi_j, j \in \mathbb{N}_0$, is a compactly supported C^∞ real-valued function that is symmetric about the origin: $\phi_j(-\cdot) = \phi_j$.
- (2) Each wavelet function $\psi_j^\ell, \ell = 1, 2, 3$ and $j \in \mathbb{N}_0$, is a compactly supported C^∞ function with l_{j+1} vanishing moments and satisfies $\psi_j^\ell(1 - \cdot) = \psi_j^\ell$ for $\ell = 1, 2$ and $\psi_j^3(1 - \cdot) = -\psi_j^3$.
- (3) The system $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j,j,k}^\ell := 2^{j/2} \psi_j^\ell(2^j \cdot - k) : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, 2, 3\}$ is a compactly supported symmetric C^∞ tight wavelet frame in $L_2(\mathbb{R})$.
- (4) If in addition $\liminf_{j \rightarrow \infty} l_j/m_j > 0$, then the tight wavelet frame in item (3) has the spectral frame approximation order.

The simplest choice in Theorem 1.2 is $m_j = l_j = j$ for all $j \in \mathbb{N}$, for which the condition in (1.21) is evidently satisfied and $\liminf_{j \rightarrow \infty} l_j/m_j = 1 > 0$. Therefore, by Theorem 1.2, we have a symmetric C^∞ tight wavelet frame in $L_2(\mathbb{R})$ with compact support and the spectral frame approximation order. Of course, the claims in Theorem 1.2 also hold if one chooses $m \approx \rho_1 j$ and $l_j \approx \rho_2 j$ for all $j \in \mathbb{N}$ with some fixed positive numbers ρ_1 and ρ_2 . In order to have refinable functions $\phi_j, j \in \mathbb{N}$, in (1.1) with support as small as possible, one should choose a sequence $\{m_j\}_{j=1}^\infty$ so that m_j goes to ∞ as slowly as possible. This is one of our motivations to choose a general integer m_j instead of the standard choice $m_j = j$ for our setup. We point out that such a strategy has already been considered by Cohen [9]. We also mention that all the claims in Theorem 1.2 hold, except possibly for the symmetry property, if one chooses the masks $\widehat{a_j} := \widehat{a_{m_j,l_j}^I}$ in (1.20) for the pseudosplines of type I instead of type II in Theorem 1.2.

It is clear that the frame approximation order in (1.16) implies (1.15). For the stationary case, it is evident that (1.13) holds with $\alpha = \beta = 0$ and, consequently, the notion of the frame approximation power in (1.15) agrees with that of the frame approximation order in (1.16). However, as illustrated by the following result, they could be quite different in the case of nonstationary tight wavelet frames.

THEOREM 1.3. *Let $\widehat{a}_j(\xi) := 2^{-j}(1 + e^{-i\xi})^j$, $j \in \mathbb{N}$, be the masks for the up-function; in other words, we take $m_j := j$ and $l_j := 1$ in Theorem 1.2. For $j \in \mathbb{N}$, define ϕ_{j-1} as in (1.1) and $\psi_{j-1}^1, \psi_{j-1}^2$, and ψ_{j-1}^3 as in (1.5) with the wavelet masks $\widehat{b}_j^1, \widehat{b}_j^2$, and \widehat{b}_j^3 being derived from \widehat{a}_j in (1.17). Then the following hold:*

- (i) $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j,j,k}^\ell : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, 2, 3\}$ is a compactly supported symmetric C^∞ tight wavelet frame in $L_2(\mathbb{R})$, and each ψ_j has one vanishing moment.
- (ii) There exists a positive constant C , independent of f and n , such that

$$(1.22) \quad \|f - Q_n(f)\|_{L_2(\mathbb{R})} \leq Cn^2 2^{-2n} |f|_{W_2^2(\mathbb{R})} \quad \forall f \in W_2^2(\mathbb{R}) \quad \text{and} \quad n \geq 2,$$

where the linear operators Q_n are defined in (1.9).

- (iii) The nonstationary tight wavelet frame in (i) does not have any frame approximation order; i.e., for any given $\nu > 0$, there does not exist a positive constant C such that (1.16) is satisfied.

The Daubechies orthogonal masks $\widehat{a}_{j,j}^I$ in (1.20) with real coefficients for the pseudosplines of type I with order (j, j) have been considered in [10] (also see [9] for the general case \widehat{a}_{m_j, m_j}^I) to obtain C^∞ compactly supported (nonstationary) orthonormal refinable functions, from which (nonstationary) orthonormal wavelets with the spectral approximation order are derived in [9, 10]. However, it is well known [13] that such Daubechies orthogonal masks $\widehat{a}_{j,j}^I$, having real coefficients and obtained from $\widehat{a}_{j,j}$ via the Fejér–Riesz lemma in (1.20), are not symmetric (except $j = 1$) and therefore, all the associated nonstationary refinable functions $\phi_j, j \in \mathbb{N}_0$, are not symmetric. One way to achieve symmetry is to split the masks $\widehat{a}_{j,j}$ into masks that are similar to $\widehat{a}_{j,j}^I$ in (1.20) but allow complex-valued coefficients (see [30]). Examples of symmetric orthonormal complex wavelets were first constructed in [30] in this way from Daubechies orthogonal masks of odd orders. Recently, symmetric orthonormal complex-valued wavelets have been systematically studied in Han [23].

Let $P_{j,j}$ be the polynomial defined in (1.18). For an odd integer j , one can always construct [23, Lemma 6 and section 2] two polynomials P_j^r and P_j^i with real coefficients such that

$$(1.23) \quad P_{j,j}(x) = [P_j^r(x)]^2 + [P_j^i(x)]^2, \quad x \in \mathbb{R} \quad \text{with} \quad P_j^r(0) = 1, \quad P_j^i(0) = 0.$$

Now define

$$(1.24) \quad \widehat{a}_j^S(\xi) := e^{i(j-1)\xi/2} 2^{-j} (1 + e^{-i\xi})^j [P_j^r(\sin^2(\xi/2)) + iP_j^i(\sin^2(\xi/2))].$$

It is easy to check [23, Lemma 3] that $|\widehat{a}_j^S(\xi)|^2 = |\widehat{a}_{j,j}^I(\xi)|^2 = \widehat{a}_{j,j}(\xi)$ and the integer shifts of the stationary refinable function associated with the mask \widehat{a}_j^S are orthonormal. Using this family of masks, one can obtain symmetric C^∞ orthonormal complex wavelets with compact support. We summarize the above discussion into the following result.

THEOREM 1.4. *Let $m_j, j \in \mathbb{N}$, be positive odd integers such that (1.21) holds. Take $\widehat{a}_j(\xi) := \widehat{a}_{m_j}^S$, where $\widehat{a}_{m_j}^S$ is defined in (1.24). Define*

$$\widehat{\psi}_{j-1}(2\xi) := e^{-i\xi} \overline{\widehat{a}_j(\xi + \pi)} \widehat{\phi}_j(\xi), \quad j \in \mathbb{N},$$

where $\phi_j, j \in \mathbb{N}_0$, are defined in (1.1). Then the following hold:

- (1) Each refinable function $\phi_j, j \in \mathbb{N}_0$, is a compactly supported C^∞ complex-valued function such that $\phi_j(1 - \cdot) = \phi_j$ and $\{\phi_j(\cdot - k) : k \in \mathbb{Z}\}$ is an orthonormal system in $L_2(\mathbb{R})$.
- (2) Each wavelet function $\psi_j, j \in \mathbb{N}_0$, is a compactly supported C^∞ complex-valued function such that $\psi_j(1 - \cdot) = -\psi_j$ and ψ_j has m_{j+1} vanishing moments.
- (3) $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;k} := 2^{j/2}\psi_j(2^j \cdot - k) : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is a compactly supported symmetric C^∞ orthonormal basis of $L_2(\mathbb{R})$ and has the spectral approximation order.

This paper is organized as follows. In section 2, we shall discuss nonstationary cascade algorithms and some properties of nonstationary refinable functions. In particular, we study the initial functions in a nonstationary cascade algorithm and provide a sufficient condition for the convergence of a nonstationary cascade algorithm in a Sobolev space $W_2^\nu(\mathbb{R})$. As a consequence, we obtain a characterization for nonstationary orthonormal wavelet bases in $L_2(\mathbb{R})$. In section 3, we shall study the frame approximation order of a nonstationary tight wavelet frame. The proofs of Theorems 1.1–1.4 will be given in section 4.

2. Nonstationary cascade algorithms and refinable functions. In this section, we first discuss the existence of L_2 -solutions of nonstationary refinable functions for a general set of masks satisfying (1.3). In fact, this follows from the following result, proved in this section:

$$(2.1) \quad [\widehat{\phi}_j, \widehat{\phi}_j](\xi) := \sum_{k \in \mathbb{Z}} |\widehat{\phi}_j(\xi + 2\pi k)|^2 \leq 1, \quad \text{a.e. } \xi \in \mathbb{R} \quad \forall j \in \mathbb{N}_0$$

for all masks satisfying (1.3), provided that the infinite products in (1.1) exist almost everywhere. The above inequality plays a critical role in our study of nonstationary tight wavelet frames and their frame approximation orders.

The question of when the refinable functions are in Sobolev spaces is discussed next. In fact, we prove it as a consequence of the convergence of the cascade algorithm in various Sobolev spaces when $\widehat{a}_j, j \in \mathbb{N}$, are masks of pseudosplines. The proof is done in the Fourier domain with the initial function whose Fourier transform is the characteristic function of $[-\pi, \pi]$. We then prove that when the cascade algorithm converges for one initial function, it converges for a large class of functions. Although a similar result is well known for the stationary case, it is not straightforward for the case of nonstationary cascade algorithms. However, this result is important in computer aided geometric design, because it results in a compactly supported function in each iteration of the cascade algorithm that generates a curve from a finitely supported sequence of points to approximate the underlying curve. Hence, it is desirable to prove the convergence of a cascade algorithm with a compactly supported initial function instead of an infinitely supported band-limited function in computer aided geometric design [17].

For the nonstationary refinable functions $\phi_j, j \in \mathbb{N}_0$, defined by masks for pseudosplines as in Theorem 1.2, one could use the same techniques developed in [10] to show that the cascade algorithm converges for the special initial function whose Fourier transform is the characteristic function of $[-\pi, \pi]$ that leads to $\phi_j \in C^\infty(\mathbb{R})$ for all $j \in \mathbb{N}_0$. But our discussion on nonstationary cascade algorithms in this section will supplement the results in [10] on nonstationary cascade algorithms. We use the results in [10] whenever they can be directly applied, e.g., Lemma 2.1, and at the same

time develop our own results to achieve our goal with a systematic and comprehensive approach. We also believe that some results (e.g., Theorem 2.4 and Lemmas 2.2 and 2.7) derived in this section have their own value in addition to being used to prove Theorem 2.8 in this section.

2.1. L_2 -solutions. We start with a basic property about the pointwise convergence of the infinite product in (1.1). A sufficient condition for the convergence of the infinite product in (1.1) has been established in the following lemma by Cohen and Dyn in [10, Theorem 2.1].

LEMMA 2.1. *Let $\widehat{a}_j, j \in \mathbb{N}$, be 2π -periodic trigonometric polynomials such that $\sup_{j \in \mathbb{N}} \|\widehat{a}_j\|_{L_\infty(\mathbb{R})} < \infty$. If (1.4) holds and $\sum_{j=1}^\infty |\widehat{a}_j(0) - 1| < \infty$, then the infinite product in (1.1) converges uniformly on every compact set of \mathbb{R} and all $\phi_j, j \in \mathbb{N}_0$, in (1.1) are well-defined compactly supported tempered distributions.*

Next, we consider when $\phi_j \in L_2(\mathbb{R}), j \in \mathbb{N}_0$, provided that the infinite products in (1.1) exist almost everywhere. In order to investigate the frame approximation order of a nonstationary tight wavelet frame, we establish (2.1), which is the following lemma.

LEMMA 2.2. *Let $\widehat{a}_j, j \in \mathbb{N}$, be 2π -periodic measurable functions satisfying (1.3) for each $j \in \mathbb{N}$. Assume that, for every $j \in \mathbb{N}_0, \widehat{\phi}_j(\xi) := \lim_{N \rightarrow \infty} \prod_{n=1}^N \widehat{a}_{n+j}(2^{-n}\xi)$ is well defined for almost every $\xi \in \mathbb{R}$; that is, the infinite product in (1.1) exists for almost every point in \mathbb{R} . Then (2.1) holds, and, consequently, $\phi_j \in L_2(\mathbb{R})$ with $\|\phi_j\|_{L_2(\mathbb{R})} \leq 1$ for every $j \in \mathbb{N}_0$.*

Proof. It suffices to prove the case $j = 0$, since the proof of the general case $j \in \mathbb{N}_0$ is the same. Note that $\widehat{\phi}_0(\xi) = \lim_{n \rightarrow \infty} \prod_{j=1}^n \widehat{a}_j(2^{-j}\xi)$ for almost every $\xi \in \mathbb{R}$. For any fixed positive integer K , we have

$$(2.2) \quad \sum_{k=-K}^K |\widehat{\phi}_0(\xi + 2\pi k)|^2 = \lim_{n \rightarrow \infty} \sum_{k=-K}^K \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 \quad \text{a.e. } \xi \in \mathbb{R}.$$

Let N be the smallest positive integer such that $N > 1 + \log_2 K$. Then we have $[-K, K] \subseteq [-K, 2^N - 1 - K]$. Consequently, for all $n \geq N$, we have $[-K, K] \subseteq [-K, 2^n - 1 - K]$ and

$$(2.3) \quad \sum_{k=-K}^K \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 \leq \sum_{k=-K}^{2^n-1-K} \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2.$$

Let $L_\infty(\mathbb{T}) := \{f \in L_\infty(\mathbb{R}) : f \text{ is } 2\pi\text{-periodic}\}$. The transition operator $T_j : L_\infty(\mathbb{T}) \rightarrow L_\infty(\mathbb{T})$ is defined for each $f \in L_\infty(\mathbb{T})$ as follows:

$$[T_j f](\xi) := |\widehat{a}_j(\xi/2)|^2 f(\xi/2) + |\widehat{a}_j(\xi/2 + \pi)|^2 f(\xi/2 + \pi), \quad \xi \in \mathbb{R}.$$

Observing that $\|[T_j f](\xi)\| \leq [|\widehat{a}_j(\xi/2)|^2 + |\widehat{a}_j(\xi/2 + \pi)|^2] \|f\|_{L_\infty(\mathbb{R})}$, by (1.3), we deduce that

$$(2.4) \quad \|T_j f\|_{L_\infty(\mathbb{R})} \leq \|f\|_{L_\infty(\mathbb{R})} \left\| |\widehat{a}_j(\cdot/2)|^2 + |\widehat{a}_j(\cdot/2 + \pi)|^2 \right\|_{L_\infty(\mathbb{R})} \leq \|f\|_{L_\infty(\mathbb{R})}.$$

By induction on n , we can verify (e.g., [21, Lemma 2.1]) that

$$(2.5) \quad \sum_{k=-K}^{2^n-1-K} \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 = [T_1 T_2 \cdots T_{n-1} T_n 1](\xi), \quad \xi \in \mathbb{R}.$$

Now it follows from (2.3) and (2.4) that, for $n \geq N$ and for almost every $\xi \in \mathbb{R}$,

$$\sum_{k=-K}^K \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 \leq [T_1 T_2 \cdots T_{n-1} T_n 1](\xi) \leq \|T_1 T_2 \cdots T_{n-1} T_n 1\|_{L_\infty(\mathbb{R})} \leq 1.$$

That is, we have

$$\sum_{k=-K}^K \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 \leq 1 \quad \text{a.e. } \xi \in \mathbb{R}, n \geq N.$$

It follows from the above inequality and (2.2) that, for any fixed positive integer K ,

$$\sum_{k=-K}^K |\widehat{\phi}_0(\xi + 2\pi k)|^2 = \lim_{n \rightarrow \infty} \sum_{k=-K}^K \prod_{j=1}^n |\widehat{a}_j(2^{-j}(\xi + 2\pi k))|^2 \leq 1 \quad \text{a.e. } \xi \in \mathbb{R}.$$

Taking $K \rightarrow \infty$ in the above inequality, we conclude that (2.1) is true for $j = 0$.

Since (2.1) implies that $\|\widehat{\phi}_j\|_{L_2(\mathbb{R})}^2 = \int_{-\pi}^\pi [\widehat{\phi}_j, \widehat{\phi}_j](\xi) d\xi \leq \int_{-\pi}^\pi 1 d\xi = 2\pi$, by Plancherel's theorem, it follows that $\|\phi_j\|_{L_2(\mathbb{R})} \leq 1$. \square

As Lemma 2.1 states, the assumption that, for $j \in \mathbb{N}_0$,

$$\widehat{\phi}_j(\xi) := \lim_{N \rightarrow \infty} \prod_{n=1}^N \widehat{a_{n+j}}(2^{-n}\xi)$$

is well defined for almost every $\xi \in \mathbb{R}$, required in Lemma 2.2, is satisfied whenever the conditions $\widehat{a}_j(0) = 1$, $j \in \mathbb{N}$, and (1.4) hold. In other words, Lemma 2.2 says that if the masks \widehat{a}_j , $j \in \mathbb{N}$, satisfy (1.3), (1.4), and $\widehat{a}_j(0) = 1$, then the corresponding nonstationary refinable functions $\phi_j \in L_2(\mathbb{R})$, $j \in \mathbb{N}_0$.

Since the approximation property of ϕ_j , $j \in \mathbb{N}_0$, discussed in this paper depends only on ϕ_j for large enough j , without loss of generality, throughout the paper we shall assume that the normalization condition $\widehat{a}_j(0) = 1$ holds for all $j \in \mathbb{N}$. In fact, if the conclusion in Lemma 2.1 holds, since $\prod_{n=1}^\infty \widehat{a_{n+j}}(0)$ converges and is nonzero for sufficiently large j , then we can replace $\widehat{\phi}_j$ and \widehat{a}_j with $\widehat{\phi}_j/\widehat{\phi}_j(0)$ and $\widehat{a}_j(\xi)/\widehat{a}_j(0)$, respectively.

2.2. Cascade algorithms. A cascade algorithm is often used to study various properties of refinable functions and is closely related to a subdivision scheme in computer aided geometric design for generating smooth curves [10, 15, 20, 24]. For a given sequence of masks $\{\widehat{a}_j\}_{j=1}^\infty$, starting with an initial function $f \in L_2(\mathbb{R})$, one computes a sequence of cascade functions f_n by

$$(2.6) \quad \widehat{f}_n(\xi) := \widehat{f}(2^{-n}\xi) \prod_{j=1}^n \widehat{a}_j(2^{-j}\xi), \quad \xi \in \mathbb{R}, n \in \mathbb{N}.$$

If $\lim_{n \rightarrow \infty} \widehat{f}(2^{-n}\xi) = 1$ and $\widehat{\phi}(\xi) := \lim_{n \rightarrow \infty} \prod_{j=1}^n \widehat{a}_j(2^{-j}\xi)$ exists for almost every $\xi \in \mathbb{R}$, then by (2.6) it is evident that $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \widehat{\phi}(\xi)$ for almost every $\xi \in \mathbb{R}$.

The cascade algorithm is closely related to another algorithm, called a subdivision scheme, which we define next. For a sequence $u : \mathbb{Z} \mapsto \mathbb{C}$, we denote \hat{u} its Fourier series as $\hat{u}(\xi) := \sum_{k \in \mathbb{Z}} u(k)e^{-ik\xi}$. In particular, by δ we denote the *Dirac sequence* on \mathbb{Z} such

that $\delta(0) = 1$ and $\delta(k) = 0$ for $k \in \mathbb{Z} \setminus \{0\}$. That is, $\hat{\delta} = 1$. For a sequence u and a mask a , the subdivision operator S_a maps the sequence u into a new sequence, $S_a u$ on \mathbb{Z} , which is determined by $\widehat{S_a u}(\xi) = 2\hat{a}(\xi)\hat{u}(2\xi)$. In fact, the product $2^n \prod_{j=1}^n \hat{a}_j(2^{n-j}\xi)$ is the Fourier series of the subdivision sequence $S_{a_n} S_{a_{n-1}} \cdots S_{a_2} S_{a_1} \delta$. More precisely, it follows from (2.6) that the cascade sequence $\{f_n\}_{n=1}^\infty$ and the subdivision sequence $\{S_{a_n} S_{a_{n-1}} \cdots S_{a_2} S_{a_1} \delta\}_{n=1}^\infty$ are related by

$$(2.7) \quad f_n = \sum_{k \in \mathbb{Z}} [S_{a_n} S_{a_{n-1}} \cdots S_{a_2} S_{a_1} \delta](k) f(2^n \cdot -k), \quad n \in \mathbb{N}.$$

Recall that $f \in W_2^\nu(\mathbb{R})$ if $\|f\|_{W_2^\nu(\mathbb{R})}^2 := \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\hat{f}(\xi)|^2 d\xi < \infty$. For a sequence of masks $\{\hat{a}_j\}_{j=1}^\infty$ and an initial function $f \in W_2^\nu(\mathbb{R})$, we say that the (nonstationary) cascade algorithm associated with masks $\{\hat{a}_j\}_{j=1}^\infty$ and an initial function f converges in the Sobolev space $W_2^\nu(\mathbb{R})$ if $f_n \in W_2^\nu(\mathbb{R})$ for all $n \in \mathbb{N}$ and the sequence $\{f_n\}_{n=1}^\infty$ is convergent in $W_2^\nu(\mathbb{R})$. Many (but not all) functions in $W_2^\nu(\mathbb{R})$ can serve as an initial function in a cascade algorithm. One popular and natural choice of an initial function f in computer aided geometric design is from the B-spline functions, since they are compactly supported functions of piecewise polynomials. Hence, it is easy to compute the values of the underlying approximating function. However, to analyze the convergence of the cascade algorithm in the frequency domain, the sinc function $f(x) = \frac{\sin(\pi x)}{\pi x}$, that is, $\hat{f} = \chi_{[-\pi, \pi]}$, the characteristic function of the interval $[-\pi, \pi]$, is a more natural choice (see, e.g., [9, 10, 11, 13, 22]). Our analysis will show that the cascade algorithm generated by pseudospline masks with the sinc function being the initial function converges in $W_2^\nu(\mathbb{R})$. To make sure that this cascade algorithm also converges in $W_2^\nu(\mathbb{R})$ when the initial seed is replaced by splines, we prove a more general result as follows: if a cascade algorithm converges in $W_2^\nu(\mathbb{R})$ for one initial seed with stable integer shifts, then it converges in $W_2^\nu(\mathbb{R})$ for a class of initial functions. As we will see, the proof is more technical than the stationary case, because a stationary refinable function is a fixed point of a stationary cascade algorithm, while this is no longer the case for the nonstationary case.

Before proceeding further, let us introduce the following notation. For $\nu \in \mathbb{R}$ and $f \in L_2(\mathbb{R})$, we define

$$(2.8) \quad [\hat{f}, \hat{f}]_\nu(\xi) := \frac{1}{|\xi|^{2\nu}} \sum_{k \in \mathbb{Z}} |\hat{f}(\xi + 2\pi k)|^2 |\xi + 2\pi k|^{2\nu}, \quad \xi \in \mathbb{R},$$

and

$$(2.9) \quad \{\hat{f}, \hat{f}\}_\nu(\xi) := \frac{1}{|\xi|^{2\nu}} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\hat{f}(\xi + 2\pi k)|^2 |\xi + 2\pi k|^{2\nu}, \quad \xi \in \mathbb{R}.$$

Clearly, we have $[\hat{f}, \hat{f}] := [\hat{f}, \hat{f}]_0 = \sum_{k \in \mathbb{Z}} |\hat{f}(\cdot + 2\pi k)|^2$ and

$$(2.10) \quad [\hat{f}, \hat{f}]_\nu(\xi) = \{\hat{f}, \hat{f}\}_\nu(\xi) + |\hat{f}(\xi)|^2.$$

Following [20, 22], we introduce the set \mathcal{F}_ν of initial functions in a cascade algorithm as

$$(2.11) \quad \mathcal{F}_\nu := \left\{ f \in W_2^\nu(\mathbb{R}) : \lim_{n \rightarrow \infty} \hat{f}(2^{-n}\xi) = 1, \quad \lim_{n \rightarrow \infty} \{\hat{f}, \hat{f}\}_\nu(2^{-n}\xi) = 0, \right. \\ \left. \text{a.e. } \xi \in [-\pi, \pi], \quad [\hat{f}, \hat{f}]_\nu \in L_\infty([-\pi, \pi]) \right\}.$$

The following result will be needed later; its proof is rather simple and therefore is omitted.

LEMMA 2.3. *Let $f \in \mathcal{F}_\nu$ for some $\nu \geq 0$. Then $[\hat{f}, \hat{f}](\xi) \leq [\hat{f}, \hat{f}]_\nu(\xi)$ for almost every $\xi \in [-\pi, \pi]$ (consequently, $[\hat{f}, \hat{f}] \in L_\infty(\mathbb{R})$), and*

$$(2.12) \quad \lim_{n \rightarrow \infty} [\hat{f}, \hat{f}](2^{-n}\xi) = \lim_{n \rightarrow \infty} [\hat{f}, \hat{f}]_\nu(2^{-n}\xi) = 1 \quad \text{a.e. } \xi \in \mathbb{R}.$$

We say that the integer shifts of a function $f \in L_2(\mathbb{R})$ are *stable* in $L_2(\mathbb{R})$ if there exists a positive constant C such that

$$(2.13) \quad C^{-1} \leq [\hat{f}, \hat{f}](\xi) \leq C, \quad \text{a.e. } \xi \in \mathbb{R}.$$

Now we state the following result on an initial function with stable integer shifts in a nonstationary cascade algorithm.

THEOREM 2.4. *Let $\hat{a}_j, j \in \mathbb{N}$, be 2π -periodic measurable functions such that*

$$\widehat{f}_\infty(\xi) := \lim_{n \rightarrow \infty} \prod_{j=1}^n \hat{a}_j(2^{-j}\xi)$$

exists for almost every $\xi \in \mathbb{R}$. For a function $f \in \mathcal{F}_\nu$ with $\nu \geq 0$ and stable integer shifts, define $f_n, n \in \mathbb{N}$, by (2.6). Assume that $\{f_n\}_{n=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$. Then it converges to f_∞ in $W_2^\nu(\mathbb{R})$. Furthermore, for every $g \in \mathcal{F}_\nu$, the sequence of functions $g_n, n \in \mathbb{N}$, defined by

$$(2.14) \quad \widehat{g}_n(\xi) := \hat{g}(2^{-n}\xi) \prod_{j=1}^n \hat{a}_j(2^{-j}\xi), \quad \xi \in \mathbb{R}, n \in \mathbb{N},$$

converges to f_∞ in $W_2^\nu(\mathbb{R})$.

Proof. By the definition of f_n in (2.6), we deduce that

$$\begin{aligned} & \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{f}_n(\xi)|^2 d\xi \\ &= \int_{\mathbb{R}} \chi_{[-\pi, \pi]}(2^{-n}\xi) \left([\hat{f}, \hat{f}](2^{-n}\xi) + |\xi|^{2\nu} [\hat{f}, \hat{f}]_\nu(2^{-n}\xi) \right) \prod_{j=1}^n |\hat{a}_j(2^{-j}\xi)|^2 d\xi. \end{aligned}$$

That is, (2.6) implies

$$(2.15) \quad \|f_n\|_{W_2^\nu(\mathbb{R})}^2 := \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{f}_n(\xi)|^2 d\xi = \int_{\mathbb{R}} F_n(\xi) d\xi$$

with

$$F_n(\xi) := \left([\hat{f}, \hat{f}](2^{-n}\xi) + |\xi|^{2\nu} [\hat{f}, \hat{f}]_\nu(2^{-n}\xi) \right) \chi_{[-\pi, \pi]}(2^{-n}\xi) \prod_{j=1}^n |\hat{a}_j(2^{-j}\xi)|^2, \quad \xi \in \mathbb{R}.$$

Similarly, by (2.14), we deduce that

$$(2.16) \quad \|g_n\|_{W_2^\nu(\mathbb{R})}^2 := \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_n(\xi)|^2 d\xi = \int_{\mathbb{R}} G_n(\xi) d\xi$$

with

$$G_n(\xi) := \left([\hat{g}, \hat{g}](2^{-n}\xi) + |\xi|^{2\nu} [\hat{g}, \hat{g}]_\nu(2^{-n}\xi) \right) \chi_{[-\pi, \pi]}(2^{-n}\xi) \prod_{j=1}^n |\hat{a}_j(2^{-j}\xi)|^2.$$

On the one hand, since $g \in \mathcal{F}_\nu$, by the definition of \mathcal{F}_ν in (2.11) and Lemma 2.3, we see that there exists a positive constant C_1 such that $[\hat{g}, \hat{g}](\xi) \leq [\hat{g}, \hat{g}]_\nu(\xi) \leq C_1$ for almost every $\xi \in [-\pi, \pi]$. By Lemma 2.3, it follows from (2.13) that $C^{-1} \leq [\hat{f}, \hat{f}]_\nu(\xi)$ for almost every $\xi \in [-\pi, \pi]$ and

$$[\hat{g}, \hat{g}](\xi) \leq C_1 \leq CC_1[\hat{f}, \hat{f}](\xi) \quad \text{and} \quad [\hat{g}, \hat{g}]_\nu(\xi) \leq C_1 \leq CC_1[\hat{f}, \hat{f}]_\nu(\xi) \quad \text{a.e. } \xi \in [-\pi, \pi].$$

Now it follows from the above inequalities that

$$(2.17) \quad |G_n(\xi)| \leq CC_1 F_n(\xi) \quad \text{a.e. } \xi \in \mathbb{R} \quad \text{and} \quad n \in \mathbb{N}.$$

On the other hand, by $f, g \in \mathcal{F}_\nu$ and Lemma 2.3, since $\lim_{n \rightarrow \infty} \prod_{j=1}^n \hat{a}_j(2^{-j}\xi) = \widehat{f_\infty}(\xi)$ for almost every $\xi \in \mathbb{R}$, we see that $\lim_{n \rightarrow \infty} F_n(\xi) = \lim_{n \rightarrow \infty} G_n(\xi) = (1 + |\xi|^{2\nu}) |\widehat{f_\infty}(\xi)|^2$ for almost every $\xi \in \mathbb{R}$. Since $\{f_n\}_{n=1}^\infty$ is a convergent sequence in $W_2^\nu(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \widehat{f_\infty}(\xi)$ for almost every $\xi \in \mathbb{R}$, we have $f_\infty \in W_2^\nu(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})} = 0$. In particular, we have $\lim_{n \rightarrow \infty} \|f_n\|_{W_2^\nu(\mathbb{R})}^2 = \|f_\infty\|_{W_2^\nu(\mathbb{R})}^2$. By (2.15), we have $\|f_n\|_{W_2^\nu(\mathbb{R})}^2 = \int_{\mathbb{R}} F_n(\xi) d\xi$. Therefore, we conclude that $\lim_{n \rightarrow \infty} \int_{\mathbb{R}} F_n(\xi) d\xi = \|f_\infty\|_{W_2^\nu(\mathbb{R})}^2$. Now by (2.17) and the generalized Lebesgue dominated convergence theorem, it follows from (2.16) and $\lim_{n \rightarrow \infty} G_n(\xi) = (1 + |\xi|^{2\nu}) |\widehat{f_\infty}(\xi)|^2$ for almost every $\xi \in \mathbb{R}$ that

$$(2.18) \quad \lim_{n \rightarrow \infty} \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_n(\xi)|^2 d\xi = \lim_{n \rightarrow \infty} \int_{\mathbb{R}} G_n(\xi) d\xi = \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{f_\infty}(\xi)|^2 d\xi.$$

But we also have

$$(1 + |\xi|^{2\nu}) |\widehat{g}_n(\xi) - \widehat{f_\infty}(\xi)|^2 \leq 2(1 + |\xi|^{2\nu}) \left[|\widehat{g}_n(\xi)|^2 + |\widehat{f_\infty}(\xi)|^2 \right].$$

By (2.18), we have

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} 2(1 + |\xi|^{2\nu}) \left[|\widehat{g}_n(\xi)|^2 + |\widehat{f_\infty}(\xi)|^2 \right] d\xi = 4 \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{f_\infty}(\xi)|^2 d\xi < \infty.$$

Now by the generalized Lebesgue dominated convergence theorem again, we conclude that

$$\begin{aligned} \lim_{n \rightarrow \infty} \|g_n - f_\infty\|_{W_2^\nu(\mathbb{R})}^2 &= \lim_{n \rightarrow \infty} \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_n(\xi) - \widehat{f_\infty}(\xi)|^2 d\xi \\ &= \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) \lim_{n \rightarrow \infty} |\widehat{g}_n(\xi) - \widehat{f_\infty}(\xi)|^2 d\xi = 0, \end{aligned}$$

since $\lim_{n \rightarrow \infty} \widehat{g}_n(\xi) = \widehat{f_\infty}(\xi)$ for almost every $\xi \in \mathbb{R}$. This completes the proof. \square

As shown in the next result, the conditions in (2.11) and (2.13) are not very restrictive. In fact, the sinc function and all B-spline functions belong to \mathcal{F}_ν for some ν .

PROPOSITION 2.5. *If $f(x) := \frac{\sin(\pi x)}{\pi x}$, then $f \in \mathcal{F}_\nu$ for all $\nu \geq 0$ and (2.13) holds, where \mathcal{F}_ν is defined in (2.11). For the B-spline B_m of order m , that is, $\widehat{B_m}(\xi) =$*

$(\frac{1-e^{-i\xi}}{i\xi})^m$, $B_m \in \mathcal{F}_\nu$ and (2.13) holds for all $0 \leq \nu < m - 1/2$. Note that $B_m \notin W_2^{m-1/2}(\mathbb{R})$.

Proof. Letting $\nu \geq 0$, we note that $\hat{f} = \chi_{[-\pi, \pi]}$ and therefore, $f \in W_2^\nu(\mathbb{R})$. For $\xi \in (-\pi, \pi)$, we have $\hat{f}(\xi + 2\pi k) = 0$ for all $k \in \mathbb{Z} \setminus \{0\}$ and, hence,

$$\{\hat{f}, \hat{f}\}_\nu(\xi) := \frac{1}{|\xi|^{2\nu}} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\hat{f}(\xi + 2\pi k)|^2 |\xi + 2\pi k|^{2\nu} = 0, \quad \xi \in (\pi, \pi) \setminus \{0\}.$$

Therefore, it is evident that $\lim_{n \rightarrow \infty} \hat{f}(2^{-n}\xi) = 1$ and $\lim_{n \rightarrow \infty} \{\hat{f}, \hat{f}\}_\nu(2^{-n}\xi) = 0$ for almost every $\xi \in [-\pi, \pi]$. Moreover, by (2.10), we have

$$[\hat{f}, \hat{f}]_\nu(\xi) = |\hat{f}(\xi)|^2 + \{\hat{f}, \hat{f}\}_\nu(\xi) = |\hat{f}(\xi)|^2 = 1, \quad \xi \in (-\pi, \pi) \setminus \{0\}.$$

Hence, $[\hat{f}, \hat{f}]_\nu \in L_\infty([-\pi, \pi])$. Thus, $f \in \mathcal{F}_\nu$. Inequality (2.13) is obviously true by $[\hat{f}, \hat{f}]_\nu(\xi) = 1$ for almost every $\xi \in \mathbb{R}$.

Letting $0 \leq \nu < m - 1/2$, from $\widehat{B_m}(\xi) = (\frac{1-e^{-i\xi}}{i\xi})^m$, we have $\lim_{j \rightarrow \infty} \widehat{B_m}(2^{-j}\xi) = \widehat{B_m}(0) = 1$, $B_m \in W_2^\nu(\mathbb{R})$, and $|\widehat{B_m}(\xi)|^2 = \frac{\sin^{2m}(\xi/2)}{(\xi/2)^{2m}}$. Now a simple computation shows that

$$\{\widehat{B_m}, \widehat{B_m}\}_\nu(\xi) \leq 2^{2(m-\nu)} \sin^{2(m-\nu)}(\xi/2) \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2(m-\nu)}.$$

Since $0 \leq \nu < m - 1/2$, we have $2(m - \nu) > 1$. Hence, for all $\xi \in [-\pi, \pi]$, the series $\sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2(m-\nu)}$ uniformly converges and therefore is uniformly bounded. Now it follows from the above inequality that $\{\widehat{B_m}, \widehat{B_m}\}_\nu(\xi)$ is uniformly bounded on $[-\pi, \pi]$ and

$$\lim_{n \rightarrow \infty} \{\widehat{B_m}, \widehat{B_m}\}_\nu(2^{-n}\xi) \leq \lim_{n \rightarrow \infty} \sin^{2(m-\nu)}(2^{-n-1}\xi) = 0 \quad \forall \xi \in [-\pi, \pi].$$

Now it follows from (2.10) that $[\widehat{B_m}, \widehat{B_m}]_\nu \in L_\infty([-\pi, \pi])$, since $\widehat{B_m} \in L_\infty(\mathbb{R})$ and $\{\widehat{B_m}, \widehat{B_m}\}_\nu \in L_\infty([-\pi, \pi])$. It is well known that (2.13) holds for B_m . \square

The following result provides us a sufficient condition on the convergence of a nonstationary cascade algorithm in a Sobolev space $W_2^\nu(\mathbb{R})$. As we will see, this result is sufficient to study the convergence of nonstationary cascade algorithms with masks for pseudosplines.

PROPOSITION 2.6. *Let \hat{a}_j and \hat{b}_j ($j \in \mathbb{N}$) be 2π -periodic measurable functions such that, for all $j \in \mathbb{N}$,*

$$(2.19) \quad |\hat{a}_j(\xi)| \leq |\hat{b}_j(\xi)| \quad a.e. \quad \xi \in \mathbb{R}.$$

Let $\eta \in W_2^\nu(\mathbb{R})$ such that $\lim_{j \rightarrow \infty} \hat{\eta}(2^{-j}\xi) = 1$ for almost every $\xi \in \mathbb{R}$. Define

$$\hat{f}_n(\xi) := \hat{\eta}(2^{-n}\xi) \prod_{j=1}^n \hat{a}_j(2^{-j}\xi) \quad \text{and} \quad \hat{g}_n(\xi) := \hat{\eta}(2^{-n}\xi) \prod_{j=1}^n \hat{b}_j(2^{-j}\xi), \quad \xi \in \mathbb{R}.$$

Assume that $\hat{f}_\infty(\xi) := \lim_{n \rightarrow \infty} \prod_{j=1}^n \hat{a}_j(2^{-j}\xi)$ and $\hat{g}_\infty(\xi) := \lim_{n \rightarrow \infty} \prod_{j=1}^n \hat{b}_j(2^{-j}\xi)$ are well defined for almost every $\xi \in \mathbb{R}$. Then, $\lim_{n \rightarrow \infty} \|g_n - g_\infty\|_{W_2^\nu(\mathbb{R})} = 0$ implies

$\lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})} = 0$. In particular, suppose that there are a positive integer J and a 2π -periodic measurable function \hat{b} such that

$$(2.20) \quad |\hat{a}_j(\xi)| \leq |\hat{b}(\xi)| \quad \text{a.e. } \xi \in \mathbb{R}, \forall j > J \quad \text{and} \quad \hat{a}_j \in L_\infty(\mathbb{R}), \quad 1 \leq j \leq J.$$

For $n \in \mathbb{N}$, define $\widehat{h}_n(\xi) := \hat{\eta}(2^{-n}\xi) \prod_{j=1}^n \hat{b}(2^{-j}\xi)$. If $\{h_n\}_{n=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$, then f_n converges to f_∞ in $W_2^\nu(\mathbb{R})$; i.e., $\lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})} = 0$.

Proof. The assumption that the functions \widehat{f}_∞ and \widehat{g}_∞ are well defined implies that $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \widehat{f}_\infty(\xi)$ and $\lim_{n \rightarrow \infty} \widehat{g}_n(\xi) = \widehat{g}_\infty(\xi)$ for almost every $\xi \in \mathbb{R}$.

By assumption, $\{g_n\}_{n=1}^\infty$ converges to g_∞ in $W_2^\nu(\mathbb{R})$; together with the fact that $\lim_{n \rightarrow \infty} \widehat{g}_n(\xi) = \widehat{g}_\infty(\xi)$ for almost every $\xi \in \mathbb{R}$, we must have $g_\infty \in W_2^\nu(\mathbb{R})$ and $\lim_{n \rightarrow \infty} \|g_n - g_\infty\|_{W_2^\nu(\mathbb{R})} = 0$. In particular, we have

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_n(\xi)|^2 d\xi = \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_\infty(\xi)|^2 d\xi < \infty.$$

Denote $\eta_n(\xi) := 2(1 + |\xi|^{2\nu}) [|\widehat{g}_n(\xi)|^2 + |\widehat{g}_\infty(\xi)|^2]$. It follows from the above identity that

$$(2.21) \quad \lim_{n \rightarrow \infty} \int_{\mathbb{R}} \eta_n(\xi) d\xi = 4 \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{g}_\infty(\xi)|^2 d\xi < \infty.$$

By (2.19), it follows from the definition of \widehat{g}_n and \widehat{f}_n that $|\widehat{f}_n(\xi)| \leq |\widehat{g}_n(\xi)|$ for almost every $\xi \in \mathbb{R}$. Since we have $\lim_{n \rightarrow \infty} \widehat{f}_n(\xi) = \widehat{f}_\infty(\xi)$ for almost every $\xi \in \mathbb{R}$, we also have $|\widehat{f}_\infty(\xi)| \leq |\widehat{g}_\infty(\xi)|$ for almost every $\xi \in \mathbb{R}$. Consequently, by $g_\infty, g_n \in W_2^\nu(\mathbb{R})$ for all $n \in \mathbb{N}$, we have $f_\infty, f_n \in W_2^\nu(\mathbb{R})$ for all $n \in \mathbb{N}$. Moreover, we have

$$(1 + |\xi|^{2\nu}) |\widehat{f}_n(\xi) - \widehat{f}_\infty(\xi)|^2 \leq 2(1 + |\xi|^{2\nu}) [|\widehat{f}_n(\xi)|^2 + |\widehat{f}_\infty(\xi)|^2] \leq \eta_n(\xi) \quad \text{a.e. } \xi \in \mathbb{R}.$$

By (2.21) and the generalized Lebesgue dominated convergence theorem, we conclude that

$$\begin{aligned} \lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})}^2 &= \lim_{n \rightarrow \infty} \int_{\mathbb{R}} (1 + |\xi|^{2\nu}) |\widehat{f}_n(\xi) - \widehat{f}_\infty(\xi)|^2 d\xi \\ &= \int_{\mathbb{R}} \lim_{n \rightarrow \infty} (1 + |\xi|^{2\nu}) |\widehat{f}_n(\xi) - \widehat{f}_\infty(\xi)|^2 d\xi = 0. \end{aligned}$$

Therefore, $\lim_{n \rightarrow \infty} \|g_n - g_\infty\|_{W_2^\nu(\mathbb{R})} = 0$ implies $\lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})} = 0$.

If (2.20) holds, for $n > J$ we deduce that, for almost every $\xi \in \mathbb{R}$,

$$\begin{aligned} |\widehat{f}_n(\xi)| &= \left[\prod_{j=1}^J |\hat{a}_j(2^{-j}\xi)| \right] |\hat{\eta}(2^{-n}\xi)| \prod_{j=J+1}^n |\hat{a}_j(2^{-j}\xi)| \\ &\leq C |\hat{\eta}(2^{-n}\xi)| \prod_{j=J+1}^n |\hat{b}(2^{-j}\xi)| = C |\widehat{h_{n-J}}(2^{-J}\xi)|, \end{aligned}$$

where $C := \prod_{j=1}^J \|\hat{a}_j\|_{L_\infty(\mathbb{R})} < \infty$. Since $\{h_n\}_{n=1}^\infty$ is convergent in $W_2^\nu(\mathbb{R})$, it is evident that $\{2^J h_{n-J}(2^J \cdot)\}_{n=J+1}^\infty$ is also convergent in $W_2^\nu(\mathbb{R})$. Note that the Fourier transform of $2^J h_{n-J}(2^J \cdot)$ is $\widehat{h_{n-J}}(2^{-J} \cdot)$. Now by the generalized Lebesgue dominated

convergence theorem, we conclude that $\{f_n\}_{n=1}^\infty$ is also convergent in $W_2^\nu(\mathbb{R})$; that is, we have $\lim_{n \rightarrow \infty} \|f_n - f_\infty\|_{W_2^\nu(\mathbb{R})} = 0$. \square

Let $\{\widehat{a}_j\}_{j=1}^\infty$ be a sequence of 2π -periodic measurable functions. Define $\{f_n\}_{n=1}^\infty$ by

$$(2.22) \quad \widehat{f}_n(\xi) := \chi_{[-\pi, \pi]}(2^{-n}\xi) \prod_{j=1}^n \widehat{a}_j(2^{-j}\xi), \quad \xi \in \mathbb{R}, n \in \mathbb{N},$$

where $\chi_{[-\pi, \pi]}$ denotes the characteristic function of the interval $[-\pi, \pi]$. This can be understood as a representation of the nonstationary cascade algorithm associated with the masks $\{\widehat{a}_j\}_{j=1}^\infty$ in the frequency domain. Due to Theorem 2.4, we say that a nonstationary cascade algorithm associated with masks $\{\widehat{a}_j\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$ if the sequence $\{f_n\}_{n=1}^\infty$ in (2.22) converges in $W_2^\nu(\mathbb{R})$. Note that the initial function here in (2.22) is the sinc function $f(x) = \frac{\sin(\pi x)}{\pi x}$ since $\widehat{f}(\xi) = \chi_{[-\pi, \pi]}$. Similarly, we say that a stationary cascade algorithm associated with a mask \widehat{a} converges in $W_2^\nu(\mathbb{R})$ if the cascade algorithm associated with $\{\widehat{a}_j\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$ with $\widehat{a}_j = \widehat{a}$ for all $j \in \mathbb{N}$.

Basically, Proposition 2.6 says that if (2.19) holds and if the nonstationary cascade algorithm associated with $\{\widehat{b}_j\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$, then so does the nonstationary cascade algorithm associated with $\{\widehat{a}_j\}_{j=1}^\infty$. Similarly, if (2.20) holds and the stationary cascade algorithm associated with mask \widehat{b} converges in $W_2^\nu(\mathbb{R})$, then Proposition 2.6 says that the nonstationary cascade algorithm associated with masks $\{\widehat{a}_j\}_{j=1}^\infty$ must converge in $W_2^\nu(\mathbb{R})$.

The convergence of a stationary cascade algorithm associated with a finitely supported mask can be verified easily by calculating the spectrum of the transition operator. Let \widehat{a} be a 2π -periodic trigonometric polynomial with $\widehat{a}(0) = 1$. Write $\widehat{a}(\xi) = (1 + e^{-i\xi})^m \widehat{c}(\xi)$ for some nonnegative integer m and some 2π -periodic trigonometric polynomial $\widehat{c}(\xi)$ with $\widehat{c}(\pi) \neq 0$. Write $|\widehat{c}(\xi)|^2 = \sum_{k=-K}^K c_k e^{-ik\xi}$, where K is some nonnegative integer. Denote $\rho(\widehat{a})$ the spectral radius of the square matrix $(c_{2j-k})_{-K \leq j, k \leq K}$ and define $\nu_2(\widehat{a}) := -1/2 - \log_2 \sqrt{\rho(\widehat{a})}$. It is known ([20, Theorem 4.3 and Proposition 7.2] and [22, Theorem 2.1]) that the stationary cascade algorithm associated with a 2π -periodic trigonometric polynomial mask \widehat{a} converges in $W_2^\nu(\mathbb{R})$ if and only if $\nu_2(\widehat{a}) > \nu$. Moreover, $\phi \in W_2^\nu(\mathbb{R})$ for all $0 \leq \nu < \nu_2(\widehat{a})$, where ϕ is the nontrivial compactly supported refinable function associated with mask \widehat{a} such that $\widehat{\phi}(2\xi) = \widehat{a}(\xi)\widehat{\phi}(\xi)$. Moreover, $\phi \notin W_2^\nu(\mathbb{R})$ for $\nu > \nu_2(\widehat{a})$ whenever the integer shifts of ϕ are stable. See [17, 20, 22, 24, 32, 37] and the many references therein on the convergence of stationary cascade algorithms.

2.3. Convergence of cascade algorithms with pseudospline masks.

We show that the nonstationary cascade algorithm associated with masks for pseudosplines in Theorem 1.2 converges in $W_2^\nu(\mathbb{R})$ for arbitrary $\nu \geq 0$. We first prove the following lemma which is not only used in the proof of the convergence of the nonstationary cascade algorithms associated with pseudospline refinement masks, but also plays an important role in our proof of Theorem 1.2 and the spectral frame approximation order.

LEMMA 2.7. *Let $\widehat{a}_{m,l}$ be the refinement mask for the pseudospline of type II with order (m, l) in (1.19). For positive integers m_1, m_2 , and l_2 such that $1 \leq m_1 \leq m_2$*

and $1 \leq l_2 \leq m_2$, the following inequality holds:

$$(2.23) \quad \begin{aligned} |\widehat{a_{m_2, l_2}}(\xi)| &\leq |\widehat{a_{m_2, l_2}^I}(\xi)| \leq |\widehat{b_{m_1}}(\xi)| \quad \forall \xi \in \mathbb{R} \quad \text{with} \\ \widehat{b_{m_1}}(\xi) &:= 2(1 + e^{-i\xi})^{-1} \widehat{a_{m_1, m_1}^I}(\xi). \end{aligned}$$

Note that $\widehat{b_{m_1}}$ is uniquely determined by $\widehat{a_{m_1, m_1}^I}(\xi) = 2^{-1}(1 + e^{-i\xi})\widehat{b_{m_1}}(\xi)$.

Proof. Since $|\widehat{a_{m, l}}(\xi)| \leq 1$ for all $\xi \in \mathbb{R}$ and $1 \leq l \leq m$, it follows from the relation $|\widehat{a_{m, l}^I}(\xi)|^2 = \widehat{a_{m, l}}(\xi)$ and (1.18) that $|\widehat{a_{m, l}}(\xi)| \leq |\widehat{a_{m, l}^I}(\xi)| \leq |\widehat{a_{m, m}^I}(\xi)|$ for all $\xi \in \mathbb{R}$ and $1 \leq l \leq m$. Now in order to prove (2.23), it suffices to prove that $\widehat{a_{m_2, m_2}}(\xi) = |\widehat{a_{m_2, m_2}^I}(\xi)|^2 \leq |\widehat{b_{m_1}}(\xi)|^2$ for all $\xi \in \mathbb{R}$ and $1 \leq m_1 \leq m_2$.

Setting $x = \sin^2(\xi/2)$, by the definition of $\widehat{a_{m_2, m_2}}$ in (1.19), we see that $\widehat{a_{m_2, m_2}}(\xi) \leq |\widehat{b_{m_1}}(\xi)|^2$ for all $\xi \in \mathbb{R}$ and $1 \leq m_1 \leq m_2$ is equivalent to

$$(2.24) \quad \begin{aligned} (1-x)^{m_2} P_{m_2, m_2}(x) &\leq (1-x)^{m_1-1} P_{m_1, m_1}(x) \\ &\quad \forall x \in [0, 1] \quad \text{and} \quad 1 \leq m_1 \leq m_2. \end{aligned}$$

By the definition of $P_{m, l}$ in (1.18), we deduce that

$$\begin{aligned} (1-x)^{m+1} P_{m+1, m+1}(x) &= (1-x)^{m+1} \sum_{j=0}^m \frac{(m+j)!}{j!m!} x^j \\ &= (1-x)^m \sum_{j=0}^m \frac{(m+j)!}{j!m!} (x^j - x^{j+1}) \\ &= (1-x)^m \left[\sum_{j=0}^m \frac{(m+j)!}{j!m!} x^j - \sum_{j=1}^{m+1} \frac{(m+j-1)!}{(j-1)!m!} x^j \right] \\ &= (1-x)^m \left[1 - \frac{(2m)!}{m!m!} x^{m+1} + \sum_{j=1}^m \left(\frac{(m+j)!}{j!m!} - \frac{(m+j-1)!}{(j-1)!m!} \right) x^j \right] \\ &= (1-x)^m \left[1 - \frac{(2m)!}{m!m!} x^{m+1} + \sum_{j=1}^m \frac{(m+j-1)!}{j!(m-1)!} x^j \right]. \end{aligned}$$

That is, we have

$$\begin{aligned} (1-x)^{m+1} P_{m+1, m+1}(x) &= (1-x)^m \left[\frac{(2m-1)!}{m!(m-1)!} x^m - \frac{(2m)!}{m!m!} x^{m+1} + \sum_{j=0}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^j \right] \\ &= (1-x)^m \left[\frac{(2m-1)!}{m!(m-1)!} x^m (1-2x) + P_{m, m}(x) \right] \\ &= (1-x)^m P_{m, m}(x) - \frac{(2m-1)!}{m!(m-1)!} x^m (1-x)^m (2x-1) \quad \forall m \in \mathbb{N}. \end{aligned}$$

Consequently, for $x \in [1/2, 1]$ and $m_2 \geq m_1$, we have $2x - 1 \geq 0$ and

$$(1-x)^{m_2} P_{m_2, m_2}(x) \leq (1-x)^{m_2-1} P_{m_2-1, m_2-1}(x) \leq \dots \leq (1-x)^{m_1} P_{m_1, m_1}(x) \leq (1-x)^{m_1-1} P_{m_1, m_1}(x).$$

Therefore, (2.24) holds for $x \in [1/2, 1]$.

It remains to prove (2.24) for all $x \in [0, 1/2]$. Since

$$(1-x)^{m_2} P_{m_2, m_2}(x) \leq (1-x)^{m_2} P_{m_2, m_2}(x) + x^{m_2} P_{m_2, m_2}(1-x) = 1,$$

in order to prove (2.24) for all $x \in [0, 1/2]$, now it suffices to show that

$$(2.25) \quad 1 \leq (1-x)^{m_1-1} P_{m_1, m_1}(x) \quad \forall x \in [0, 1/2].$$

Note that

$$\begin{aligned} 1 &= (1-x)(1-x)^{m_1-1} P_{m_1, m_1}(x) + x x^{m_1-1} P_{m_1, m_1}(1-x) \\ &= (1-x)^{m_1-1} P_{m_1, m_1}(x) - x[(1-x)^{m_1-1} P_{m_1, m_1}(x) - x^{m_1-1} P_{m_1, m_1}(1-x)], \end{aligned}$$

from which we have

$$(1-x)^{m_1-1} P_{m_1, m_1}(x) = 1 + x[(1-x)^{m_1-1} P_{m_1, m_1}(x) - x^{m_1-1} P_{m_1, m_1}(1-x)].$$

In order to prove (2.25), by the above identity, it suffices to prove that

$$(2.26) \quad (1-x)^{m_1-1} P_{m_1, m_1}(x) \geq x^{m_1-1} P_{m_1, m_1}(1-x) \quad \forall x \in [0, 1/2].$$

Note that for $x \in [0, 1/2]$, we have $0 \leq x/(1-x) \leq 1$. By the definition of $P_{m_1, m_1}(x) := \sum_{j=0}^{m_1-1} \binom{m_1+j-1}{j} x^j$, we have

$$\begin{aligned} 1 &= \frac{1}{1} \geq \frac{\binom{m_1}{1} x}{\binom{m_1}{1} (1-x)} \geq \frac{\binom{m_1+1}{2} x^2}{\binom{m_1+1}{2} (1-x)^2} \geq \dots \geq \frac{\binom{2m_1-2}{m_1-1} x^{m_1-1}}{\binom{2m_1-2}{m_1-1} (1-x)^{m_1-1}} \\ &= \frac{x^{m_1-1}}{(1-x)^{m_1-1}}, \quad x \in [0, 1/2]. \end{aligned}$$

But for positive numbers a, b, c, d , it is easy to see that $\frac{a}{b} \geq \frac{c}{d}$ implies $\frac{a}{b} \geq \frac{a+c}{b+d} \geq \frac{c}{d}$. Now it follows from the above inequalities that

$$1 \geq \frac{P_{m_1, m_1}(x)}{P_{m_1, m_1}(1-x)} = \frac{\sum_{j=0}^{m_1-1} \binom{m_1+j-1}{j} x^j}{\sum_{j=0}^{m_1-1} \binom{m_1+j-1}{j} (1-x)^j} \geq \frac{x^{m_1-1}}{(1-x)^{m_1-1}}, \quad x \in [0, 1/2],$$

from which we see that (2.26) holds. \square

Next we establish the following result on C^∞ nonstationary refinable functions. In fact, we prove that the nonstationary cascade algorithm associated with the masks for pseudosplines of type I or type II in Theorem 1.2 converges in $W_2^\nu(\mathbb{R})$ for any $\nu \geq 0$.

THEOREM 2.8. *Let \hat{a}_j be the mask for the pseudospline of type I or type II with order (m_j, l_j) , where $1 \leq l_j \leq m_j$ and $j \in \mathbb{N}$ are positive integers such that (1.21) holds. Then, for every $n \in \mathbb{N}_0$, the nonstationary cascade algorithm (2.22)*

associated with $\{\widehat{a_{j+n}}\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$ for any $\nu \geq 0$. Consequently, the nonstationary refinable functions $\phi_j, j \in \mathbb{N}_0$, in (1.1) must be well-defined compactly supported $C^\infty(\mathbb{R})$ functions.

Proof. Since $\widehat{a_j}(\xi) = \widehat{a_{m_j, l_j}}(\xi)$ or $\widehat{a_j}(\xi) = \widehat{a_{m_j, l_j}^I}(\xi)$, it is easy to see that $\deg(\widehat{a_j}) \leq 2m_j$ and $\widehat{a_j}(0) = 1$. Therefore, by our assumption in (1.21), we see that

$$\sum_{j=1}^\infty 2^{-j} \deg(\widehat{a_j}) \leq 2 \sum_{j=1}^\infty 2^{-j} m_j < \infty.$$

Moreover, we have $|\widehat{a_j}(\xi)| \leq 1$ for all $j \in \mathbb{N}$ and $\xi \in \mathbb{R}$. Thus, the condition of Lemma 2.1 is satisfied. By Lemma 2.1, we conclude that $\phi_j, j \in \mathbb{N}_0$, are well-defined compactly supported tempered distributions.

Let $\widehat{a_{j,j}}$ be the refinement mask for the pseudospline of type II with order (j, j) . It was proved by Daubechies [12, 13] that $\lim_{j \rightarrow \infty} \nu_2(\widehat{a_{j,j}}) = \lim_{j \rightarrow \infty} \nu_2(\widehat{a_{j,j}^I}) = \infty$. Hence, there exists a positive integer J such that $\nu_2(\widehat{a_{j,j}^I}) \geq \nu + 2$. By $\lim_{j \rightarrow \infty} m_j = \infty$, there exists a positive integer N such that

$$(2.27) \quad m_j \geq J \quad \text{and} \quad \nu_2(\widehat{a_j}) \geq \nu_2(\widehat{a_{j,j}^I}) \geq \nu + 2 \quad \forall j \geq N.$$

Let \hat{b} be the unique 2π -periodic trigonometric polynomial such that $\widehat{a_{j,j}^I}(\xi) = 2^{-1}(1 + e^{-i\xi})\hat{b}(\xi)$. By the definition of $\nu_2(\hat{b})$ and (2.27), it is straightforward to see that $\nu_2(\hat{b}) = \nu_2(\widehat{a_{j,j}^I}) - 1 \geq \nu + 1 > \nu$. Therefore, the stationary cascade algorithm associated with the mask \hat{b} converges in $W_2^\nu(\mathbb{R})$ (see [20, Theorem 4.3]). On the other hand, by (2.23) of Lemma 2.7, since $m_j \geq J$ for $j \geq N$, we deduce that

$$|\widehat{a_j}(\xi)| \leq |\widehat{a_{m_j, l_j}^I}(\xi)| \leq |\hat{b}(\xi)| \quad \forall \xi \in \mathbb{R}, j \geq N.$$

Since the stationary cascade algorithm associated with the mask \hat{b} converges in $W_2^\nu(\mathbb{R})$, by Proposition 2.6, the nonstationary cascade algorithm associated with masks $\{\widehat{a_j}\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$. Therefore, we have $\phi_0 \in W_2^\nu(\mathbb{R})$ for all $\nu \geq 0$. That is, ϕ_0 is a compactly supported C^∞ function. The same proof works for every ϕ_n and for the nonstationary cascade algorithm associated with masks $\{\widehat{a_{n+j}}\}_{j=1}^\infty$. \square

In computer aided geometric design, it is of interest to consider the convergence of a subdivision scheme and a cascade algorithm in $C^\kappa(\mathbb{R})$, the space of functions with the κ th continuous derivative, instead of the Sobolev space $W_2^\nu(\mathbb{R})$. When a cascade algorithm is implemented in the space domain as given in (2.7), the initial function is often chosen to be a compactly supported function such as spline functions of order ν with $\nu > \kappa + 1/2$. This is indeed true and can be proved easily by the imbedding theorems of Sobolev spaces.

COROLLARY 2.9. *Let $\widehat{a_j} (j \in \mathbb{N})$ be the masks satisfying the conditions given in Theorem 2.8. Let κ be any nonnegative integer, and let f be a compactly supported initial function in $W_2^\nu \cap \mathcal{F}_\nu$ with $\nu > \kappa + 1/2$. Then the nonstationary cascade algorithm defined by (2.7) associated with $\{\widehat{a_{n+j}}\}_{j=1}^\infty$ converges in $C^\kappa(\mathbb{R})$.*

Proof. The convergence of the cascade algorithm in $W_2^\nu(\mathbb{R})$ defined by (2.7) follows from Proposition 2.6 and Theorem 2.8. Since all $f_n (n \in \mathbb{N}_0)$ and ϕ_0 are supported inside some compact set, it follows from the imbedding theorem that the sequence f_n also converges in $C^\kappa(\mathbb{R})$. \square

2.4. Orthogonality. For the stationary case, it is well known that for a compactly supported refinable function ϕ with a 2π -periodic trigonometric polynomial \hat{a} , the integer shifts of ϕ form an orthonormal system if and only if its refinement mask \hat{a} satisfies

$$|\hat{a}(\xi)|^2 + |\hat{a}(\xi + \pi)|^2 = 1 \quad \text{a.e. } \xi \in \mathbb{R}$$

and the corresponding stationary cascade algorithm converges in $L_2(\mathbb{R})$ (that is, $\nu_2(\hat{a}) > 0$; see [20, 22]). Furthermore, if one chooses the wavelet function ψ by $\hat{\psi}(2\xi) := e^{-i\xi} \hat{a}(\xi + \pi) \hat{\phi}(\xi)$, then the wavelet system generated by ψ forms an orthonormal basis in $L_2(\mathbb{R})$. For example, see [11, 12, 13, 20, 22, 24, 31, 32, 37] and the references therein. It turns out that this is also true for the nonstationary case, as a consequence of the proof of Theorem 2.4.

THEOREM 2.10. *Let $\hat{a}_j, j \in \mathbb{N}$, be 2π -periodic measurable functions such that, for every $j \in \mathbb{N}_0$, $\widehat{\phi}_j(\xi) := \lim_{N \rightarrow \infty} \prod_{n=1}^N \widehat{a_{n+j}}(2^{-n}\xi)$ is well defined for almost every $\xi \in \mathbb{R}$. Then the integer shifts of ϕ_j form an orthonormal system in $L_2(\mathbb{R})$ for all $j \in \mathbb{N}$; i.e.,*

$$(2.28) \quad \langle \phi_j(\cdot - k), \phi_j \rangle := \int_{\mathbb{R}} \phi_j(x - k) \overline{\phi_j(x)} dx = \delta(k) \quad \forall k \in \mathbb{Z} \quad \text{and} \quad j \in \mathbb{N}_0,$$

where $\delta(0) = 1$ and $\delta(k) = 0$ for all $k \neq 0$ if and only if

(1) all the masks \hat{a}_j satisfy

$$(2.29) \quad |\widehat{a}_j(\xi)|^2 + |\widehat{a}_j(\xi + \pi)|^2 = 1 \quad \text{a.e. } \xi \in \mathbb{R}, \quad \forall j \in \mathbb{N}_0;$$

(2) the nonstationary cascade algorithm associated with masks $\{a_{n+j}\}_{n=1}^\infty$ converges in $L_2(\mathbb{R})$ for large enough $j \in \mathbb{N}_0$.

Moreover, if (1.4) and (2.28) hold, define $\widehat{\psi}_{j-1}(\xi) := e^{-i\xi/2} \widehat{a}_j(\xi/2 + \pi) \widehat{\phi}_j(\xi/2)$ for $j \in \mathbb{N}$; then $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j,j,k} : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is an orthonormal basis of $L_2(\mathbb{R})$.

Proof. It is known that (2.28) holds for each j if and only if $[\widehat{\phi}_j, \widehat{\phi}_j](\xi) = 1$ almost everywhere $\xi \in \mathbb{R}$.

Assume that (1) and (2) hold. Then $[\widehat{\phi}_j, \widehat{\phi}_j](\xi) = 1$ almost everywhere $x \in \mathbb{R}$ for all $j \in \mathbb{N}_0$ can be proved by an argument similar to that in the stationary case (see [11, 13, 22, 24, 32]). We omit the details here.

The necessity part is proved as follows. If (2.28) holds, then $[\widehat{\phi}_j, \widehat{\phi}_j] = 1$ for all $j \in \mathbb{N}_0$. Now (2.29) can be verified by the same argument as in the stationary case. So, item (1) holds. Next, we prove item (2); that is, the sequence $\{f_N\}_{N=1}^\infty$, defined by

$$\widehat{f}_N(\xi) := \chi_{[-\pi, \pi]}(2^{-N}\xi) \prod_{n=1}^N \widehat{a_{n+j}}(2^{-n}\xi), \quad N \in \mathbb{N},$$

converges in $L_2(\mathbb{R})$ for every $j \in \mathbb{N}_0$. By the relation $\widehat{\phi}_j(\xi) = \widehat{a_{j+1}}(\xi/2) \widehat{\phi}_{j+1}(\xi/2)$, we deduce that

$$(2.30) \quad \widehat{\phi}_j(\xi) = \widehat{\phi_{N+j}}(2^{-N}\xi) \prod_{n=1}^N \widehat{a_{n+j}}(2^{-n}\xi), \quad N \in \mathbb{N} \quad \text{and} \quad j \in \mathbb{N}_0.$$

Applying (2.30) and $[\widehat{\phi}_N, \widehat{\phi}_N] = 1$, one obtains that

$$\|\widehat{\phi}_j\|_{L_2(\mathbb{R})}^2 = \int_{\mathbb{R}} |\widehat{f}_N(\xi)|^2 d\xi, \quad N \in \mathbb{N}.$$

In particular, we have $\lim_{N \rightarrow \infty} \int_{\mathbb{R}} |\widehat{f}_N(\xi)|^2 d\xi = \|\widehat{\phi}_j\|_{L_2(\mathbb{R})}^2 < \infty$. Note that

$$|\widehat{f}_N(\xi) - \widehat{\phi}_j(\xi)|^2 \leq 2 \left[|\widehat{f}_N(\xi)|^2 + |\widehat{\phi}_j(\xi)|^2 \right] =: \eta_N(\xi), \quad \xi \in \mathbb{R},$$

and

$$\lim_{N \rightarrow \infty} \int_{\mathbb{R}} \eta_N(\xi) d\xi = 4 \int_{\mathbb{R}} |\widehat{\phi}_j(\xi)|^2 d\xi < \infty.$$

By $\lim_{N \rightarrow \infty} \widehat{f}_N(\xi) = \widehat{\phi}_j(\xi)$ for almost every $\xi \in \mathbb{R}$ and the generalized Lebesgue dominated convergence theorem, we conclude that $\lim_{N \rightarrow \infty} \|\widehat{f}_N - \widehat{\phi}_j\|_{L_2(\mathbb{R})} = 0$. Thus, for every $j \in \mathbb{N}_0$, the cascade algorithm associated with masks $\{\widehat{a_{n+j}}\}_{n=1}^\infty$ converges in $L_2(\mathbb{R})$.

If (2.28) holds, by the definition of ψ_j and (2.29), then it is easy to check that $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;k} : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is an orthonormal system of $L_2(\mathbb{R})$. By Theorem 1.1, we see that $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;k} : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is a tight frame in $L_2(\mathbb{R})$. Therefore, $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;k} : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is an orthonormal basis of $L_2(\mathbb{R})$. \square

3. The approximation order of the truncated frame series. In this section, we shall study the approximation property of a nonstationary tight wavelet frame, i.e., the frame approximation properties of the operators Q_n in (1.9). The approximation operators Q_n and their approximation order for a given stationary tight wavelet frame have been extensively studied in [14]. The approximation operators Q_n provide a simple approximation scheme for a given tight wavelet frame and have close links to the frame decomposition and reconstruction algorithms for a tight wavelet frame (see, e.g., [14]). Moreover, their approximation order determines the accuracy of the truncation operators and is not necessarily equal to the best approximation order provided by the underlying nonstationary multiresolution analysis.

Since the approximation operators Q_n provide a simple approximation scheme for a given tight wavelet frame, they are often used in various applications. For example, in [1, 2, 3, 4, 5], where the (stationary) tight wavelet frame based algorithms for high/super resolution image reconstructions, image inpainting, and deconvolutions are given, the operators Q_n are used there to approximate the underlying function from a given data set. The interested reader should consult [1, 2, 3, 4, 5] for details.

The operators Q_n are closely related to other operators. For a sequence $\{\phi_n\}_{n=0}^\infty$ of functions in $L_2(\mathbb{R})$, we define the linear operators $P_n(f), n \in \mathbb{N}_0$, by

$$P_n(f) := \sum_{k \in \mathbb{Z}} \langle f, \phi_{n;n,k} \rangle \phi_{n;n,k}, \quad f \in L_2(\mathbb{R}) \quad \text{with} \tag{3.1}$$

$$\phi_{n;n,k} := 2^{n/2} \phi_n(2^n \cdot -k).$$

Similar to the stationary case, by calculation, it is easy to verify that (1.12) implies

$$\begin{aligned}
 (3.2) \quad & \sum_{\ell=1}^{\mathcal{J}_j} \sum_{k \in \mathbb{Z}} \langle f, \psi_{j-1;j-1,k}^\ell \rangle \psi_{j-1;j-1,k}^\ell \\
 &= \sum_{k \in \mathbb{Z}} \langle f, \phi_{j;j,k} \rangle \phi_{j;j,k} - \sum_{k \in \mathbb{Z}} \langle f, \phi_{j-1;j-1,k} \rangle \phi_{j-1;j-1,k}
 \end{aligned}$$

for $f \in L_2(\mathbb{R})$. Consequently, by the definition of the linear operators Q_n in (1.9), it follows from the relation in (3.2) that

$$\begin{aligned}
 Q_n(f) &= \sum_{k \in \mathbb{Z}} \langle f, \phi_0(\cdot - k) \rangle \phi_0(\cdot - k) + \sum_{j=0}^{n-1} \sum_{\ell=1}^{\mathcal{J}_{j+1}} \sum_{k \in \mathbb{Z}} \langle f, \psi_{j;j,k}^\ell \rangle \psi_{j;j,k}^\ell \\
 &= \sum_{k \in \mathbb{Z}} \langle f, \phi_{n;n,k} \rangle \phi_{n;n,k} = P_n(f).
 \end{aligned}$$

That is, if (1.12) holds, then the linear operators Q_n in (1.9) and P_n in (3.1) are the same.

The approximation order of Q_n 's for a stationary tight wavelet frame is investigated in [14] through that of P_n 's, since $Q_n = P_n$ for a tight frame system constructed from the unitary extension principle. The relationship has been studied in [14] for stationary tight wavelet frames between the approximation order of P_n 's and the (best) approximation order provided by the spaces $S_n(\phi_n)$, where $S_n(\phi_n)$ is the smallest closed subspace of $L_2(\mathbb{R})$ generated by the linear span of $\phi_n(2^n \cdot -k)$, $k \in \mathbb{Z}$; that is, $S_n(\phi_n)$ is the same as the smallest closed subspace of $L_2(\mathbb{R})$ containing the truncated tight frame system $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;j,k}^\ell : k \in \mathbb{Z}, 0 \leq j < n, \ell = 1, \dots, \mathcal{J}_{j+1}\}$. It is well known that the approximation order provided by the spaces $S_n(\phi_n)$ is determined by the order of the Strang–Fix conditions satisfied by ϕ_n . However, the approximation order of P_n 's is determined by the order of the zero at the origin of the function $1 - |\widehat{\phi}|^2$, in addition to the order of the Strang–Fix conditions satisfied by ϕ_n . Consequently, the frame approximation order can be (much) smaller than the approximation order provided by the spaces $S_n(\phi_n)$ (see [14] for details). This is also true for the nonstationary case. For example, let $\widehat{a}_j(\xi) := 2^{-j}(1 + e^{-i\xi})^j$, $j \in \mathbb{N}$, be the masks for the up-functions. Then, item (iii) of Theorem 1.3 says that it does not have any “strong” frame approximation order in the sense of (1.16); i.e., for any given $\nu > 0$, there does not exist a positive constant C such that (1.16) is satisfied. But one can check that the corresponding spaces $S_n(\phi_n)$ provide a spectral approximation order.

If the integer shifts of ϕ_n are orthonormal, then the linear operator P_n in (3.1) becomes an orthogonal projection from $L_2(\mathbb{R})$ to $S_n(\phi_n)$. That is, for this case, $P_n(f)$ is the best approximation of $f \in L_2(\mathbb{R})$ in the closed subspace $S_n(\phi_n)$ of $L_2(\mathbb{R})$. This is the reason why the approximation order of Q_n 's is identified with the best approximation order provided by the spaces $S_n(\phi_n)$ in [10], which is simpler to understand, since only orthonormal wavelets are studied in [10].

To summarize our discussion here, the understanding of the approximation order of the approximation operators Q_n for a given tight wavelet frame is necessary, since it is simple and used in applications such as image inpainting, and since, unlike orthonormal wavelets, the approximation order of a truncated tight frame series is

not necessarily the same as the best approximation order provided by the underlying nonstationary multiresolution analysis.

The main result of this section is Theorem 3.2, which is interesting in its own right and is independent of its role in our proofs of some of major parts of Theorems 1.1–1.4.

The following result can be directly obtained by applying Jetter and Zhou [28] and [29, Theorem 2.1].

PROPOSITION 3.1. *Let $\varphi \in L_2(\mathbb{R})$ and $\nu \geq 0$. Define a linear operator P by*

$$P(f) := \sum_{k \in \mathbb{Z}} \langle f, \varphi(\cdot - k) \rangle \varphi(\cdot - k), \quad f \in L_2(\mathbb{R}).$$

Then $\|f - P(f)\|_{L_2(\mathbb{R})} \leq C_\varphi |f|_{W_2^\nu(\mathbb{R})}$ for all $f \in W_2^\nu(\mathbb{R})$ with a positive constant

$$(3.3) \quad C_\varphi := \pi^{-1/2} \sqrt{\max(c_1, c_3) + \max(2c_2, 2c_4 + 1)},$$

provided that there exist positive constants c_1, c_2, c_3, c_4 such that, for almost every $\xi \in [-\pi, \pi]$, the following inequalities hold:

$$(3.4) \quad |1 - |\hat{\varphi}(\xi)|^2|^2 \leq c_1 |\xi|^{2\nu},$$

$$(3.5) \quad \sum_{k \in \mathbb{Z} \setminus \{0\}} |\hat{\varphi}(\xi)|^2 |\hat{\varphi}(\xi + 2\pi k)|^2 \leq c_2 |\xi|^{2\nu},$$

$$(3.6) \quad \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2\nu} |\hat{\varphi}(\xi)|^2 |\hat{\varphi}(\xi + 2\pi k)|^2 \leq c_3,$$

$$(3.7) \quad \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2\nu} \sum_{\ell \in \mathbb{Z} \setminus \{0\}} |\hat{\varphi}(\xi + 2\pi \ell)|^2 |\hat{\varphi}(\xi + 2\pi k)|^2 \leq c_4.$$

Next, we present the following result on the approximation properties of the operators P_n defined in (3.1).

THEOREM 3.2. *Let $\hat{a}_j, j \in \mathbb{N}$, be 2π -periodic measurable functions such that (1.3) holds for all $j \in \mathbb{N}$, and, for every $n \in \mathbb{N}_0$, the function $\hat{\phi}_n(\xi) := \lim_{J \rightarrow \infty} \prod_{j=1}^J \hat{a}_{j+n}(2^{-j}\xi)$ is well defined for almost every $\xi \in \mathbb{R}$. Let $\nu \geq 0$. If, for $n \in \mathbb{N}$,*

$$(3.8) \quad \begin{aligned} |1 - |\hat{\phi}_n(\xi)|^2|^2 &\leq C_{\phi_n} |\xi|^{2\nu} && \text{a.e. } \xi \in [-\pi, \pi], \\ \sum_{k \in \mathbb{Z} \setminus \{0\}} |\hat{\phi}_n(\xi)|^2 |\hat{\phi}_n(\xi + 2\pi k)|^2 &\leq C_{\phi_n} |\xi|^{2\nu} && \text{a.e. } \xi \in [-\pi, \pi], \end{aligned}$$

where C_{ϕ_n} is a constant depending only on ϕ_n , then, for the linear operators P_n in (3.1),

$$(3.9) \quad \|f - P_n(f)\|_{L_2(\mathbb{R})} \leq \max(2, \sqrt{C_{\phi_n}}) 2^{-\nu n} |f|_{W_2^\nu(\mathbb{R})} \quad \forall f \in W_2^\nu(\mathbb{R}) \quad \text{and } n \in \mathbb{N}.$$

In particular, (3.8) is satisfied if

$$(3.10) \quad 1 - |\hat{\phi}_n(\xi)|^2 \leq C_{\phi_n} |\xi|^{2\nu} \quad \text{a.e. } \xi \in [-\pi, \pi].$$

Proof. By (1.3) and Lemma 2.2, we have $\phi_n \in L_2(\mathbb{R})$ for all $n \in \mathbb{N}_0$. For each fixed $n \in \mathbb{N}_0$, we denote $P_{n,0}$ the following linear operator on $L_2(\mathbb{R})$:

$$P_{n,0}(f) := \sum_{k \in \mathbb{Z}} \langle f, \phi_n(\cdot - k) \rangle \phi_n(\cdot - k), \quad f \in L_2(\mathbb{R}).$$

It is apparent [28] that the operators P_n and $P_{n,0}$ are linked through the relation $P_n(f) = [P_{n,0}(f(2^{-n}\cdot))](2^n\cdot)$ and

$$(3.11) \quad \|f - P_n(f)\|_{L_2(\mathbb{R})} = 2^{-n/2} \|f(2^{-n}\cdot) - P_{n,0}(f(2^{-n}\cdot))\|_{L_2(\mathbb{R})}.$$

Since (1.3) holds, by Lemma 2.2, we have

$$[\widehat{\phi}_n, \widehat{\phi}_n](\xi) := \sum_{k \in \mathbb{Z}} |\widehat{\phi}_n(\xi + 2\pi k)|^2 \leq 1 \quad \text{a.e. } \xi \in \mathbb{R}.$$

In particular, we have $|\widehat{\phi}_n(\xi)| \leq 1$ for almost every $\xi \in \mathbb{R}$.

Since $\nu \geq 0$, for $\xi \in [-\pi, \pi]$, it is evident that $|\xi + 2\pi k|^{-2\nu} \leq 1$ for all $k \in \mathbb{Z} \setminus \{0\}$ and

$$\begin{aligned} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2\nu} \sum_{\ell \in \mathbb{Z} \setminus \{0\}} |\widehat{\phi}_n(\xi + 2\pi \ell)|^2 |\widehat{\phi}_n(\xi + 2\pi k)|^2 \\ \leq \sum_{k \in \mathbb{Z}} |\widehat{\phi}_n(\xi + 2\pi k)|^2 \sum_{\ell \in \mathbb{Z}} |\widehat{\phi}_n(\xi + 2\pi \ell)|^2 \leq 1. \end{aligned}$$

Hence, (3.7) holds with $c_4 = 1$ and $\varphi = \phi_n$. Similarly, for $\xi \in [-\pi, \pi]$, we have

$$\sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-2\nu} |\widehat{\phi}_n(\xi)|^2 |\widehat{\phi}_n(\xi + 2\pi k)|^2 \leq \sum_{k \in \mathbb{Z}} |\widehat{\phi}_n(\xi + 2\pi k)|^2 \leq 1.$$

Therefore, (3.6) holds with $c_3 = 1$ and $\varphi = \phi_n$.

By (3.8), we see that (3.4) and (3.5) hold with $c_1 = c_2 = C_{\phi_n}$ and $\varphi = \phi_n$.

Note that $|f(2^{-n}\cdot)|_{W_2^\nu(\mathbb{R})} = 2^{n/2} 2^{-\nu n} |f|_{W_2^\nu(\mathbb{R})}$. Therefore, by Proposition 3.1 and (3.11), we conclude that, for all $f \in W_2^\nu(\mathbb{R})$,

$$(3.12) \quad \begin{aligned} \|f - P_n(f)\|_{L_2(\mathbb{R})} &= 2^{-n/2} \|f(2^{-n}\cdot) - P_{n,0}(f(2^{-n}\cdot))\|_{L_2(\mathbb{R})} \\ &\leq \max(2, \sqrt{C_{\phi_n}}) 2^{-\nu n} |f|_{W_2^\nu(\mathbb{R})}, \end{aligned}$$

since

$$\begin{aligned} \max(c_1, c_3) + \max(2c_2, 2c_4 + 1) &= \max(C_{\phi_n}, 1) + \max(2C_{\phi_n}, 3) \leq \max(3C_{\phi_n}, 6) \\ &\leq \pi \max(C_{\phi_n}, 4). \end{aligned}$$

Therefore, (3.9) is verified.

If (3.10) holds, then

$$|1 - |\widehat{\phi}_n(\xi)|^2|^2 \leq C_{\phi_n}^2 |\xi|^{2\nu} |\xi|^{2\nu} \leq C_{\phi_n} |\xi|^{2\nu}, \quad |\xi| \leq C_{\phi_n}^{-\frac{1}{2\nu}},$$

and by $|\widehat{\phi}_n(\xi)| \leq 1$,

$$|1 - |\widehat{\phi}_n(\xi)|^2|^2 \leq 1 = C_{\phi_n} [C_{\phi_n}^{-\frac{1}{2\nu}}]^{2\nu} \leq C_{\phi_n} |\xi|^{2\nu}, \quad |\xi| \geq C_{\phi_n}^{-\frac{1}{2\nu}}.$$

Also, by $|\widehat{\phi}_n(\xi)|^2 \leq [\widehat{\phi}_n, \widehat{\phi}_n](\xi) \leq 1$ and the above two inequalities, we have

$$\begin{aligned} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\widehat{\phi}_n(\xi)|^2 |\widehat{\phi}_n(\xi + 2\pi k)|^2 &\leq \sum_{k \in \mathbb{Z} \setminus \{0\}} |\widehat{\phi}_n(\xi + 2\pi k)|^2 = [\widehat{\phi}_n, \widehat{\phi}_n](\xi) - |\widehat{\phi}_n(\xi)|^2 \\ &\leq 1 - |\widehat{\phi}_n(\xi)|^2. \end{aligned}$$

Consequently, (3.10) implies (3.8). \square

The behavior of $1 - |\widehat{\phi}_n(\xi)|^2$ near the origin $\xi = 0$ in (3.10) is closely related to that of the masks $1 - |\widehat{a}_j(\xi)|^2$ for all $j \in \mathbb{N}$ near the origin $\xi = 0$. As we will see in the proof of the next section, when only masks are available and the nonstationary refinable functions are not explicitly given, one can use the estimate of $1 - |\widehat{a}_j(\xi)|^2$ near $\xi = 0$ for all $j \in \mathbb{N}$ to obtain the estimate of $1 - |\widehat{\phi}_n(\xi)|^2$ near $\xi = 0$. The following result provides the estimate of $1 - |\widehat{a}_j(\xi)|^2$ near the origin for the masks of pseudosplines, which will be needed later in our proof of the spectral frame approximation order in Theorems 1.2 and 1.4.

LEMMA 3.3. *Let $a_{m,l}^I$ and $\widehat{a_{m,l}}$ be pseudospline masks of types I and II of order (m, l) defined in (1.19) and (1.20), respectively. For any $0 < \rho \leq 1$ and $\nu \geq 0$, there exist a positive integer N and a positive constant C (both of them depend only on ρ and ν), such that, for all $N \leq \rho m < l \leq m$,*

$$(3.13) \quad 0 \leq 1 - |a_{m,l}^I(\xi)|^2 \leq 1 - |\widehat{a_{m,l}}(\xi)|^2 \leq C|\xi|^{2\nu} \quad \forall \xi \in [-\pi, \pi].$$

Proof. Since $|a_{m,l}^I(\xi)|^2 = \widehat{a_{m,l}}(\xi) \leq 1$, we have

$$0 \leq 1 - |a_{m,l}^I(\xi)|^2 \leq 1 - |\widehat{a_{m,l}}(\xi)|^2 = [1 + \widehat{a_{m,l}}(\xi)][1 - \widehat{a_{m,l}}(\xi)] \leq 2[1 - \widehat{a_{m,l}}(\xi)].$$

Setting $x = \sin^2(\xi/2)$, by the definition of $\widehat{a_{m,l}}(\xi)$ in (1.19), we have $\widehat{a_{m,l}}(\xi) = (1 - x)^m P_{m,l}(x)$, where the polynomial $P_{m,l}$ is defined in (1.18). In order to prove (3.13), now it is easy to see that it is equivalent to proving that, for any $0 < \rho \leq 1$ and any positive integer ν , there exist a positive integer N and a positive constant C , all depending only on ρ and ν , such that

$$(3.14) \quad 1 - (1 - x)^m P_{m,l}(x) \leq Cx^\nu \quad \forall x \in [0, 1] \quad \text{and} \quad N \leq \rho m < l \leq m.$$

Since $(1 - x)^m P_{m,m}(x) + x^m P_{m,m}(1 - x) = 1$, we deduce that

$$\begin{aligned} 1 - (1 - x)^m P_{m,l}(x) &= (1 - x)^m P_{m,m}(x) + x^m P_{m,m}(1 - x) - (1 - x)^m P_{m,l}(x) \\ &= x^m P_{m,m}(1 - x) + (1 - x)^m \sum_{j=l}^{m-1} \frac{(m + j - 1)!}{j!(m - 1)!} x^j. \\ &= x^\nu \left[x^{m-\nu} P_{m,m}(1 - x) + (1 - x)^m \sum_{j=l}^{m-1} \frac{(m + j - 1)!}{j!(m - 1)!} x^{j-\nu} \right]. \end{aligned}$$

By Lemma 2.7, (2.24) holds. In particular, replacing x by $1 - x$ in (2.24), we conclude that

$$x^{m-\nu} P_{m,m}(1 - x) \leq x^{N-\nu-1} P_{N,N}(1 - x) \quad \forall x \in [0, 1] \quad \text{and} \quad m \geq N \geq \nu + 1.$$

Therefore, on the one hand, we have

$$(3.15) \quad x^{m-\nu} P_{m,m}(1 - x) \leq C_N := \max_{x \in [0,1]} x^{N-\nu-1} P_{N,N}(1 - x) < \infty$$

$$\forall x \in [0, 1], \quad m \geq N \geq \nu + 1.$$

On the other hand, for $x \in [0, 1]$, we have

$$\begin{aligned} (1-x)^m \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^{j-\nu} &= (1-x)^m x^{l-\nu} \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^{j-l} \\ &\leq (1-x)^m x^{l-\nu} \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} \\ &= (1-x)^m x^{l-\nu} \sum_{j=l}^{m-1} \left[\frac{(m+j)!}{j!m!} - \frac{(m+j-1)!}{(j-1)!m!} \right] \\ &\leq (1-x)^m x^{l-\nu} \frac{(2m-1)!}{(m-1)!m!}. \end{aligned}$$

By Stirling's formula, we have $\lim_{n \rightarrow \infty} \frac{n!}{\sqrt{2\pi n} n^{n+1/2} e^{-n}} = 1$. Thus, when m is large enough and for all $x \in [0, 1]$, we have

$$\frac{(2m-1)!}{(m-1)!m!} = \frac{1}{2} \frac{(2m)!}{m!m!} \leq \frac{1}{\sqrt{2\pi}} \frac{(2m)^{2m+1/2} e^{-2m}}{m^{2m+1} e^{-2m}} \leq 4^m m^{-1/2}.$$

Thus, for large enough m , we have

$$(1-x)^m \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^{j-\nu} \leq (1-x)^m x^{l-\nu} 4^m m^{-1/2} = [4(1-x)x^{(l-\nu)/m}]^m m^{-1/2}.$$

Since $\rho m < l$, we have $l/m > \rho$. Since ν is fixed, when m is large enough, we have $(l-\nu)/m \geq \rho$ and, hence, $x^{(l-\nu)/m} \leq x^\rho$ for all $x \in [0, 1]$. Therefore, we have

$$(1-x)^m \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^{j-\nu} \leq [4(1-x)x^\rho]^m m^{-1/2} \quad \forall x \in [0, 1].$$

Note that $\rho > 0$. The continuous function $(1-x)x^\rho$ has only one critical point $x = \frac{\rho}{1+\rho}$ on the interval $(0, 1)$, and it takes value zero at $x = 0$. Thus, we can choose $0 < \tau \leq 1$ such that $4(1-x)x^\rho \leq 1$ for all $x \in [0, \tau]$; one may prefer to choose such τ as large as possible; in particular, τ may be obtained by solving $4(1-\tau)\tau^\rho = 1$. Therefore, there exists a positive integer N such that, for $N \leq \rho m < l \leq m$,

$$(3.16) \quad (1-x)^m \sum_{j=l}^{m-1} \frac{(m+j-1)!}{j!(m-1)!} x^{j-\nu} \leq [4(1-x)x^\rho]^m m^{-1/2} \leq m^{-1/2} \leq 1$$

$$\forall x \in [0, \tau].$$

Combining (3.15) and (3.16), we conclude that

$$1 - (1-x)^m P_{m,l}(x) \leq x^\nu [C_N + 1] \quad \forall x \in [0, \tau], \quad N \leq \rho m < l \leq m.$$

Observing that $0 \leq (1-x)^m P_{m,l}(x) \leq 1$ for all $x \in [0, 1]$, we deduce that

$$1 - (1-x)^m P_{m,l}(x) \leq 1 = x^{-\nu} x^\nu \leq \tau^{-\nu} x^\nu \quad \forall x \in [\tau, 1], \quad 1 \leq l \leq m, \quad m \in \mathbb{N}.$$

Thus, (3.14) is verified with $C := \max(C_N + 1, \tau^{-\nu})$. This completes the proof. \square

4. Proofs of Theorems 1.1–1.4. In this section, we shall prove Theorems 1.1–1.4. We start with a proof of Theorem 1.1.

Proof of Theorem 1.1. By our assumption in Theorem 1.1, item (i) follows from Lemmas 2.1 and 2.2. Since (1.3) holds, by Lemma 2.2, we have $[\widehat{\phi}_n, \widehat{\phi}_n] \leq 1$ and $\phi_n \in L_2(\mathbb{R})$. Therefore, the linear operators Q_n in (1.9) and P_n in (3.1) are well defined, bounded, and the same (see section 3).

Let us first prove (1.15) in item (iii). In order to do so, in the following we estimate the constants C_{ϕ_n} in (3.10). Denote $\widehat{d}_j(\xi) := |\widehat{a}_j(\xi)|^2$ for $j \in \mathbb{N}$. Since $\widehat{a}_j(0) = 1$, we have $\widehat{d}_j(0) = 1$ for all $j \in \mathbb{N}$. By our assumption and Lemma 2.1, we have $|\widehat{\phi}_n(\xi)|^2 = \prod_{j=1}^{\infty} \widehat{d_{j+n}}(2^{-j}\xi)$ for all $\xi \in \mathbb{R}$. Therefore, we have

$$\begin{aligned} 1 - |\widehat{\phi}_n(\xi)|^2 &= \prod_{j=1}^{\infty} \widehat{d_{j+n}}(0) - \prod_{j=1}^{\infty} \widehat{d_{j+n}}(2^{-j}\xi) \\ &= \sum_{\ell=1}^{\infty} \left[\prod_{j=1}^{\ell-1} \widehat{d_{j+n}}(0) \right] [\widehat{d_{\ell+n}}(0) - \widehat{d_{\ell+n}}(2^{-\ell}\xi)] \left[\prod_{j=\ell+1}^{\infty} \widehat{d_{j+n}}(2^{-j}\xi) \right]. \end{aligned}$$

Since $\widehat{d_{j+n}}(0) = 1$ and $0 \leq \widehat{d_{j+n}}(\xi) \leq 1$ by (1.3), we conclude that

$$(4.1) \quad 0 \leq 1 - |\widehat{\phi}_n(\xi)|^2 \leq \sum_{\ell=1}^{\infty} |\widehat{d_{\ell+n}}(0) - \widehat{d_{\ell+n}}(2^{-\ell}\xi)|, \quad \xi \in \mathbb{R}.$$

Since \widehat{a}_j is a 2π -periodic trigonometric polynomial, by the definition of \widehat{d}_j , we see that \widehat{d}_j is a real-valued C^∞ function. Note that, by our assumption, 2ν is a positive integer. Therefore, for $\xi \in [-\pi, \pi]$, there exists $\zeta_{\xi,j} \in [-\pi, \pi]$ such that

$$(4.2) \quad \widehat{d}_j(\xi) = \widehat{d}_j(0) + \frac{\widehat{d}_j^{(1)}(0)}{1!} \xi + \dots + \frac{\widehat{d}_j^{(2\nu-1)}(0)}{(2\nu-1)!} \xi^{2\nu-1} + \frac{\widehat{d}_j^{(2\nu)}(\zeta_{\xi,j})}{(2\nu)!} \xi^{2\nu}.$$

By Bernstein’s inequality and $0 \leq \widehat{d}_j(\xi) \leq 1$, we have

$$\|\widehat{d}_j^{(2\nu)}\|_{L_\infty(\mathbb{R})} \leq [\deg(\widehat{d}_j)]^{2\nu} \|\widehat{d}_j\|_{L_\infty(\mathbb{R})} \leq [\deg(\widehat{d}_j)]^{2\nu}.$$

By assumption in (1.14), for $j \geq N$, we have $\widehat{d}_j^{(\ell)}(0) = 0$ for all $\ell = 1, \dots, 2\nu - 1$. Therefore, it follows from (4.2) that, for $n \geq N$, $\ell \in \mathbb{N}$, and $\xi \in [-\pi, \pi]$,

$$|\widehat{d_{\ell+n}}(0) - \widehat{d_{\ell+n}}(2^{-\ell}\xi)| \leq \frac{[\deg(\widehat{d_{\ell+n}})]^{2\nu}}{(2\nu)!} |2^{-\ell}\xi|^{2\nu} \leq \frac{1}{(2\nu)!} |\xi|^{2\nu} 2^{-2\nu\ell} [\deg(\widehat{d_{\ell+n}})]^{2\nu}.$$

Therefore, we have

$$\sum_{\ell=1}^{\infty} |\widehat{d_{\ell+n}}(0) - \widehat{d_{\ell+n}}(2^{-\ell}\xi)| \leq \frac{1}{(2\nu)!} |\xi|^{2\nu} \sum_{\ell=1}^{\infty} 2^{-2\nu\ell} [\deg(\widehat{d_{\ell+n}})]^{2\nu}.$$

That is, by (4.1) we see that, for every $n \geq N$, (3.10) holds with

$$(4.3) \quad C_{\phi_n} := \frac{1}{(2\nu)!} \sum_{\ell=1}^{\infty} 2^{-2\nu\ell} [\deg(\widehat{d_{\ell+n}})]^{2\nu}.$$

Now we estimate C_{ϕ_n} using the condition in (1.13). By (1.13), there exists a positive constant C_1 such that

$$(4.4) \quad \deg(\widehat{a}_j) \leq C_1 j^\alpha 2^{\beta j} \quad \forall j \in \mathbb{N}.$$

By the definition of \widehat{d}_j , we have $\deg(\widehat{d}_j) \leq 2 \deg(\widehat{a}_j)$ for all $j \in \mathbb{N}$. Therefore, from (4.4), we deduce that

$$\begin{aligned} 2^{-2\nu\ell} [\deg(\widehat{d}_{\ell+n})]^{2\nu} &\leq 2^{2\nu} C_1^{2\nu} (\ell+n)^{2\nu\alpha} 2^{2\nu\beta(\ell+n)} 2^{-2\nu\ell} \\ &= 2^{2\nu} C_1^{2\nu} n^{2\nu\alpha} 2^{2\nu\beta n} (1 + \ell/n)^{2\nu\alpha} 2^{-2\nu(1-\beta)\ell} \\ &\leq 2^{2\nu} C_1^{2\nu} n^{2\nu\alpha} 2^{2\nu\beta n} (1 + \ell)^{2\nu\alpha} 2^{-2\nu(1-\beta)\ell}. \end{aligned}$$

Consequently, we have the following estimate for the constant C_{ϕ_n} :

$$\begin{aligned} C_{\phi_n} &= \frac{1}{(2\nu)!} \sum_{\ell=1}^{\infty} 2^{-2\nu\ell} [\deg(\widehat{d}_{\ell+n})]^{2\nu} \leq \frac{2^{2\nu} C_1^{2\nu}}{(2\nu)!} n^{2\nu\alpha} 2^{2\nu\beta n} \sum_{\ell=1}^{\infty} (1 + \ell)^{2\nu\alpha} 2^{-2\nu(1-\beta)\ell} \\ &= C_2 n^{2\nu\alpha} 2^{2\nu\beta n}, \end{aligned}$$

where $C_2 := 2^{2\nu} C_1^{2\nu} [(2\nu)!]^{-1} \sum_{\ell=1}^{\infty} (1 + \ell)^{2\nu\alpha} 2^{-2\nu(1-\beta)\ell} < \infty$, since $1 - \beta > 0$ and $\nu > 0$. Since $Q_n = P_n$, by Theorem 3.2, we conclude that

$$\|f - Q_n(f)\|_{L_2(\mathbb{R})} \leq \max(2, \sqrt{C_2}) n^{\nu\alpha} 2^{-\nu(1-\beta)n} |f|_{W_2^\nu(\mathbb{R})} \quad \forall f \in W_2^\nu(\mathbb{R}) \quad \text{and} \quad n \geq N.$$

That is, (1.15) holds with $C := \max(2, \sqrt{C_2}) < \infty$, which is independent of f and n .

Now we prove item (ii). In order to show that $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;j,k}^\ell : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, \dots, J_{j+1}\}$ is a tight frame of $L_2(\mathbb{R})$, since $Q_n = P_n$, it now suffices to show that

$$\lim_{n \rightarrow \infty} \|f - Q_n(f)\|_{L_2(\mathbb{R})} = \lim_{n \rightarrow \infty} \|f - P_n(f)\|_{L_2(\mathbb{R})} = 0 \quad \forall f \in L_2(\mathbb{R}).$$

Since $\widehat{a}_j(0) = 1$ and $\widehat{d}_j(\xi) = |\widehat{a}_j(\xi)|^2$, it is evident that $\widehat{d}_j(0) = 1$. Since \widehat{d}_j is a 2π -periodic trigonometric polynomial, the condition in (1.14) is automatically satisfied with $\nu = 1/2$.

Now from the above proof of item (iii) and by Theorem 3.2, we see that (3.9) holds with the constant C_{ϕ_n} defined in (4.3) and $\nu = 1/2$. More precisely, by Theorem 3.2, for $\nu = 1/2$, we have

$$(4.5) \quad \|f - P_n(f)\|_{L_2(\mathbb{R})}^2 \leq C_n |f|_{W_2^{1/2}(\mathbb{R})}^2 \quad \forall f \in W_2^{1/2}(\mathbb{R})$$

with

$$(4.6) \quad C_n := \max(4, C_{\phi_n}) 2^{-n} \quad \text{and} \quad C_{\phi_n} := \sum_{\ell=1}^{\infty} 2^{-\ell} \deg(\widehat{d}_{\ell+n}).$$

Now we prove that $\lim_{n \rightarrow \infty} C_n = 0$ by showing $\lim_{n \rightarrow \infty} 2^{-n} C_{\phi_n} = 0$. Note that

$$(4.7) \quad 2^{-n} C_{\phi_n} = \sum_{\ell=1}^{\infty} 2^{-(\ell+n)} \deg(\widehat{d}_{\ell+n}) = \sum_{j=n+1}^{\infty} 2^{-j} \deg(\widehat{d}_j).$$

Since $\deg(\widehat{d}_j) \leq 2 \deg(\widehat{a}_j)$, by our assumption in (1.4), we have

$$\sum_{j=1}^{\infty} 2^{-j} \deg(\widehat{d}_j) \leq 2 \sum_{j=1}^{\infty} 2^{-j} \deg(\widehat{a}_j) < \infty.$$

Consequently, by (4.7), we conclude that

$$0 \leq \lim_{n \rightarrow \infty} 2^{-n} C_{\phi_n} \leq \lim_{n \rightarrow \infty} \sum_{j=n+1}^{\infty} 2^{-j} \deg(\widehat{d}_j) = 0.$$

That is, $\lim_{n \rightarrow \infty} 2^{-n} C_{\phi_n} = 0$. Thus, we have $\lim_{n \rightarrow \infty} C_n = \lim_{n \rightarrow \infty} \max(4, C_{\phi_n}) 2^{-n} = 0$. Now from (4.5), we see that

$$\lim_{n \rightarrow \infty} \|f - P_n(f)\|_{L_2(\mathbb{R})}^2 = \lim_{n \rightarrow \infty} C_n \|f\|_{W_2^{1/2}(\mathbb{R})}^2 = 0 \quad \forall f \in W_2^{1/2}(\mathbb{R}).$$

Since $P_n = Q_n$ in (3.1), we conclude that

$$\begin{aligned} \|f\|_{L_2(\mathbb{R})}^2 &= \lim_{n \rightarrow \infty} \langle Q_n(f), f \rangle = \sum_{k \in \mathbb{Z}} |\langle f, \phi_0(\cdot - k) \rangle|^2 + \sum_{j=0}^{\infty} \sum_{\ell=1}^{\mathcal{J}_{j+1}} \sum_{k \in \mathbb{Z}} |\langle f, \psi_{j;j,k}^{\ell} \rangle|^2 \\ &\quad \forall f \in W_2^{1/2}(\mathbb{R}). \end{aligned}$$

Since $W_2^{1/2}(\mathbb{R})$ is dense in $L_2(\mathbb{R})$, (1.2) must hold for all $f \in L_2(\mathbb{R})$. Therefore, $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;j,k}^{\ell} : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, \dots, \mathcal{J}_{j+1}\}$ is a tight frame of $L_2(\mathbb{R})$. \square

Next, we prove Theorem 1.2.

Proof of Theorem 1.2. For item (1), applying Theorem 2.8, we conclude that all $\phi_j, j \in \mathbb{N}_0$, are compactly supported functions in $C^\infty(\mathbb{R})$. Since all the masks \widehat{a}_j are 2π -periodic trigonometric polynomials with real coefficients and are symmetric about the origin, we have $\overline{\widehat{a}_j(\xi)} = \widehat{a}_j(\xi)$. Now by the definition of $\widehat{\phi}_j$ in (1.1), it is straightforward to see that all $\phi_j, j \in \mathbb{N}_0$, are real-valued and $\overline{\widehat{\phi}_j(\xi)} = \widehat{\phi}_j(\xi)$; that is, all $\phi_j, j \in \mathbb{N}_0$, are symmetric about the origin.

By the definition of $\widehat{a_{m,l}^I}$ and $\widehat{a_{m,l}}$, we have

$$(4.8) \quad 1 - |\widehat{a_{m,l}^I}(\xi)|^2 = 1 + O(|\xi|^{2l}) \quad \text{and} \quad 1 - |\widehat{a_{m,l}}(\xi)|^2 = 1 + O(|\xi|^{2l}), \quad \xi \rightarrow 0.$$

By (1.12) and (4.8), we see that $\widehat{b_j^\ell}(\xi) = O(|\xi|^{l_j})$ as $\xi \rightarrow 0$. Therefore, ψ_j^ℓ has l_{j+1} vanishing moments. Thus, item (2) holds.

For item (3), by the definition of $\widehat{b_j^\ell}$ in (1.17), it is straightforward to check that (1.12) holds with $\mathcal{J}_j = 3$ for all $j \in \mathbb{N}$. Now by Theorem 1.1, we see that $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j;j,k}^\ell : j \in \mathbb{N}_0, k \in \mathbb{Z}, \ell = 1, \dots, \mathcal{J}_{j+1}\}$ is a tight wavelet frame in $L_2(\mathbb{R})$.

Now we prove item (4) by using Theorem 3.2 and Lemma 3.3. Let ν be an arbitrary positive integer. Since $\liminf_{j \rightarrow \infty} l_j/m_j > 0$, there exist a positive integer N and $0 < \rho < \liminf_{j \rightarrow \infty} l_j/m_j$ such that $2\nu < N < \rho m_j < l_j \leq m_j$ for all $j \geq N$. Denote $\widehat{d}_j(\xi) := |\widehat{a}_j(\xi)|^2$. By Lemma 3.3, we see that (3.13) holds. That is, there exists a positive constant C , independent of j , such that

$$(4.9) \quad 0 \leq 1 - \widehat{d}_j(\xi) \leq C|\xi|^{2\nu}, \quad \xi \in [-\pi, \pi] \quad \text{and} \quad j \geq N.$$

We now use (4.9) to estimate the constants C_{ϕ_n} in (3.10) of Theorem 3.2. For $n \geq N$ and $\ell \in \mathbb{N}$, since $\widehat{d_{\ell+n}}(0) = 1$, it follows from (4.9) that

$$|\widehat{d_{\ell+n}}(0) - \widehat{d_{\ell+n}}(2^{-\ell}\xi)| = |1 - \widehat{d_{\ell+n}}(2^{-\ell}\xi)| \leq C2^{-2\nu\ell}|\xi|^{2\nu} \quad \forall \xi \in [-\pi, \pi].$$

Now by (4.1), we conclude that

$$1 - |\widehat{\phi}_n(\xi)|^2 \leq C|\xi|^{2\nu} \sum_{\ell=1}^{\infty} 2^{-2\nu\ell}, \quad \xi \in [-\pi, \pi].$$

Therefore, (3.10) holds with

$$C_{\phi_n} := C \sum_{\ell=1}^{\infty} 2^{-2\nu\ell} = \frac{C}{1 - 2^{-2\nu}} < \infty.$$

Consequently, by $Q_n = P_n$ and Theorem 3.2, we conclude that

$$\|f - Q_n(f)\|_{L_2(\mathbb{R})} \leq C_1 2^{-\nu n} \|f\|_{W_2^\nu(\mathbb{R})} \quad \forall f \in W_2^\nu(\mathbb{R}) \quad \text{and} \quad n \geq N,$$

where

$$C_1 := \max(2, \sqrt{C/(1 - 2^{-2\nu})}) < \infty$$

is independent of f and n . Since ν is arbitrary, the tight wavelet frame has the desired spectral frame approximation order. \square

Now we prove Theorem 1.3.

Proof of Theorem 1.3. Item (i) has been proved in Theorem 1.2. It is evident that

$$\deg(\widehat{a}_j) = j \quad \text{and} \quad |\widehat{a}_j(\xi)|^2 = \cos^{2j}(\xi/2) = 1 + O(|\xi|^2), \quad \xi \rightarrow 0.$$

Now it follows from the proof of item (iii) of Theorem 1.1 that there exists a positive constant C_1 such that

$$1 - |\widehat{\phi}_n(\xi)|^2 \leq C_1 n^2 |\xi|^2 \quad \forall \xi \in [-\pi, \pi] \quad \text{and} \quad n \in \mathbb{N}.$$

That is, we conclude that

$$(4.10) \quad \left|1 - |\widehat{\phi}_n(\xi)|^2\right|^2 \leq C_1^2 n^4 |\xi|^4 \quad \forall \xi \in [-\pi, \pi] \quad \text{and} \quad n \in \mathbb{N}.$$

Let B_2 be the B-spline of order 2. Then

$$|\widehat{B}_2(\xi)|^2 = \frac{\sin^4(\xi/2)}{(\xi/2)^4} \quad \text{and} \quad |\widehat{B}_2(2\xi)|^2 = \cos^4(\xi/2) |\widehat{B}_2(\xi)|^2.$$

Since $|\widehat{a}_j(\xi)|^2 = \cos^{2j}(\xi/2) \leq \cos^4(\xi/2)$ for all $\xi \in \mathbb{R}$ and $j \geq 2$, it is evident that $|\widehat{\phi}_n(\xi)|^2 \leq |\widehat{B}_2(\xi)|^2$ for all $\xi \in \mathbb{R}$ and $n \geq 2$. In particular, for $\xi \in [-\pi, \pi]$, we deduce that

$$\begin{aligned} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\widehat{\phi}_n(\xi)|^2 |\widehat{\phi}_n(\xi + 2\pi k)|^2 &\leq \sum_{k \in \mathbb{Z} \setminus \{0\}} |\widehat{B}_2(\xi + 2\pi k)|^2 \\ &= |\xi|^4 \frac{\sin^4(\xi/2)}{(\xi/2)^4} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-4}. \end{aligned}$$

Setting $C_2 := \sup_{\xi \in [-\pi, \pi]} \sum_{k \in \mathbb{Z} \setminus \{0\}} |\xi + 2\pi k|^{-4} < \infty$, we conclude that

$$\sum_{k \in \mathbb{Z} \setminus \{0\}} |\widehat{\phi}_n(\xi)|^2 |\widehat{\phi}_n(\xi + 2\pi k)|^2 \leq C_2 |\xi|^4, \quad \xi \in [-\pi, \pi] \quad \text{and} \quad n \geq 2.$$

Now taking into account (4.10), we see that the two inequalities in (3.8) hold with $\nu = 2$ and $C_{\phi_n} := \max(C_1, C_2)n^4$. Thus, by Theorem 3.2, we see that item (ii) holds.

Now we prove item (iii) using proof by contradiction. Suppose that (1.16) holds for some $\nu > 0$. By [29, Theorem 2.2], we have $|1 - |\widehat{\phi}_n(\xi)|^2|^2 \leq \pi C^2 |\xi|^{2\nu}$ for almost every $\xi \in [-\pi, \pi]$ and $n \geq N$, where C is the positive constant in (1.16). That is, we must have

$$(4.11) \quad 1 - |\widehat{\phi}_j(\xi)|^2 \leq C_3 |\xi|^\nu \quad \forall \xi \in [-\pi, \pi] \quad \text{and} \quad j \geq N,$$

where $C_3 := \sqrt{\pi}C$ is a positive constant independent of j .

By the definition of $\widehat{\phi}_j$ in (1.1) and $\widehat{a}_j(\xi) = 2^{-j}(1 + e^{-i\xi})^j$, it is evident that $\widehat{\phi}_j(\xi) = \widehat{B}_j(\xi)\widehat{\phi}_0(\xi)$, where B_j is the B-spline of order j . Since $|\widehat{\phi}_0(\xi)| \leq 1$, by $\widehat{\phi}_j(\xi) = \widehat{B}_j(\xi)\widehat{\phi}_0(\xi)$, we have $|\widehat{\phi}_j(\xi)| \leq |\widehat{B}_j(\xi)|$, and, therefore, it follows from (4.11) that

$$(4.12) \quad 1 - |\widehat{B}_j(\xi)|^2 \leq 1 - |\widehat{\phi}_n(\xi)|^2 \leq C_3 |\xi|^\nu \quad \forall \xi \in [-\pi, \pi] \quad \text{and} \quad j \geq N.$$

Since $|\widehat{B}_j(\xi)|^2 = \cos^{2j}(\xi/4)|\widehat{B}_j(\xi/2)|^2$ and $|\widehat{B}_j(\xi)| \leq 1$, we have

$$1 - |\widehat{B}_j(4\xi)|^2 = 1 - \cos^{2j}(\xi) + \cos^{2j}(\xi)(1 - |\widehat{B}_j(2\xi)|^2) \geq 1 - \cos^{2j}(\xi).$$

Consequently, (4.12) implies

$$(4.13) \quad \frac{1 - \cos^{2j}(\xi)}{|\xi|^\nu} \leq \frac{1 - |\widehat{B}_j(4\xi)|^2}{|\xi|^\nu} \leq C_4 \quad \forall \xi \in [-\pi/4, \pi/4] \quad \text{and} \quad j \geq N,$$

where $C_4 := 4^\nu C_3 < \infty$. Noting $1 = [\cos^2(\xi) + \sin^2(\xi)]^j \geq \cos^{2j}(\xi) + j \cos^{2j-2}(\xi) \sin^2(\xi)$, by $4\pi^{-2}\xi^2 \leq \sin^2(\xi) \leq |\xi|^2$ for all $\xi \in [-\pi/2, \pi/2]$, we deduce that, for $\xi \in [-\pi/4, \pi/4]$,

$$\begin{aligned} \frac{1 - \cos^{2j}(\xi)}{|\xi|^\nu} &\geq j \frac{\sin^2(\xi) \cos^{2j-2}(\xi)}{|\xi|^\nu} = j \frac{\sin^2(\xi)(1 - \sin^2(\xi))^{j-1}}{|\xi|^\nu} \\ &\geq 4\pi^{-2}j(\xi^2)^{1-\nu/2}(1 - \xi^2)^{j-1}. \end{aligned}$$

Taking $\xi_j := \sqrt{\frac{1-\nu/2}{j-\nu/2}}$, we observe that $\lim_{j \rightarrow \infty} \xi_j = 0$, and it follows from the above inequalities and (4.13) that, for $\xi \in [-\pi/4, \pi/4]$ and sufficiently large j ,

$$\begin{aligned} C_4 &\geq \frac{1 - \cos^{2j}(\xi_j)}{|\xi_j|^\nu} \geq 4\pi^{-2}j(\xi_j^2)^{1-\nu/2}(1 - \xi_j^2)^{j-1} \\ &= 4\pi^{-2}j \frac{(1 - \nu/2)^{1-\nu/2}(j-1)^{j-1}}{(j - \nu/2)^{j-\nu/2}} := C_5 c_j, \end{aligned}$$

where $0 < C_5 := 4\pi^{-2}(1 - \nu/2)^{1-\nu/2} < \infty$ and $c_j := j \frac{(j-1)^{j-1}}{(j-\nu/2)^{j-\nu/2}}$. That is, we must have

$$(4.14) \quad c_j := j \frac{(j-1)^{j-1}}{(j - \nu/2)^{j-\nu/2}} \leq C_4/C_5$$

for all sufficiently large integers j .

By calculation, we have

$$c_j = j \frac{(j-1)^{j-1}}{(j-\nu/2)^{j-\nu/2}} = j \frac{j^{j-1}(1-\frac{1}{j})^{j-1}}{j^{j-\nu/2}(1-\frac{\nu}{2j})^{j-\nu/2}} = j^{\nu/2} \left(1-\frac{1}{j}\right)^{j-1} \left(1-\frac{\nu}{2j}\right)^{-j+\nu/2}.$$

We note that $\lim_{j \rightarrow \infty} (1-\frac{1}{j})^{j-1} = e^{-1}$ and $\lim_{j \rightarrow \infty} (1-\frac{\nu}{2j})^{-j+\nu/2} = e^{\nu/2}$. Hence, by $\nu > 0$, we conclude that $\lim_{j \rightarrow \infty} c_j = \lim_{j \rightarrow \infty} j^{\nu/2} e^{\nu/2-1} = \infty$, which is a contradiction to (4.14). So, the tight wavelet frame does not have any frame approximation order. Now item (iii) is verified. \square

We finish this paper by proving Theorem 1.4.

Proof of Theorem 1.4. For item (1), by a direct calculation [23, Lemma 3], we observe that

$$\begin{aligned} |\widehat{a}_j(\xi)|^2 &= \cos^{2m_j}(\xi/2) ([P_{m_j}^r(\sin^2(\xi/2))]^2 + [P_{m_j}^i(\sin^2(\xi/2))]^2) \\ &= \cos^{2m_j}(\xi/2) P_{m_j, m_j}(\sin^2(\xi/2)) = |\widehat{a_{m_j, m_j}^I}(\xi)|^2. \end{aligned}$$

Hence $|\widehat{a}_j(\xi)|^2 + |\widehat{a}_j(\xi + \pi)|^2 = 1$. By Theorem 2.8 and Proposition 2.6, we see that the nonstationary cascade algorithm associated with $\{\widehat{a}_j\}_{j=1}^\infty$ converges in $W_2^\nu(\mathbb{R})$ for any $\nu \geq 0$ and, therefore, all ϕ_j , $j \in \mathbb{N}_0$, are well-defined compactly supported functions in $C^\infty(\mathbb{R})$. By Theorem 2.10, (2.28) holds, and $\{\phi_0(\cdot - k) : k \in \mathbb{Z}\} \cup \{\psi_{j; j, k} : j \in \mathbb{N}_0, k \in \mathbb{Z}\}$ is an orthonormal basis of $L_2(\mathbb{R})$. By the same proof as in Theorem 1.2 and Lemma 3.3 with $l_j = m_j$, this orthonormal wavelet basis has the spectral approximation order. So, item (3) is verified.

The symmetry $\phi_j(1 - \cdot) = \phi_j$ follows [23, Lemma 2] from the definition of ϕ_j in (1.1) and the symmetry of the masks \widehat{a}_j : $\widehat{a}_j(\xi) = e^{-i\xi} \widehat{a}_j(-\xi)$. Item (2) can be easily verified by (4.8). \square

Acknowledgment. This work was done and completed while BH visited the National University of Singapore in 2006. BH would like to thank the National University of Singapore for their hospitality and support during his visit.

REFERENCES

- [1] J.-F. CAI, R. CHAN, L. SHEN, AND Z. SHEN, *Restoration of chopped and noded images by framelets*, SIAM J. Sci. Comput., 30 (2008), pp. 1205–1227.
- [2] J.-F. CAI, R. H. CHAN, AND Z. SHEN, *A framelet-based image inpainting algorithm*, Appl. Comput. Harmon. Anal., 24 (2008), pp. 131–149.
- [3] A. CHAI AND Z. SHEN, *Deconvolution: A wavelet frame approach*, Numer. Math., 106 (2007), pp. 529–587.
- [4] R. H. CHAN, S. D. RIEMENSCHNEIDER, L. SHEN, AND Z. SHEN, *Tight frame: An efficient way for high-resolution image reconstruction*, Appl. Comput. Harmon. Anal., 17 (2004), pp. 91–115.
- [5] R. H. CHAN, Z. SHEN, AND T. XIA, *A framelet algorithm for enhancing video stills*, Appl. Comput. Harmon. Anal., 23 (2007), pp. 153–170.
- [6] C. K. CHUI AND W. HE, *Compactly supported tight frames associated with refinable functions*, Appl. Comput. Harmon. Anal., 8 (2000), pp. 293–319.
- [7] C. K. CHUI, W. HE, AND J. STÖCKLER, *Compactly supported tight and sibling frames with maximum vanishing moments*, Appl. Comput. Harmon. Anal., 13 (2002), pp. 224–262.
- [8] C. K. CHUI, W. HE, AND J. STÖCKLER, *Nonstationary tight wavelet frames, II: Unbounded interval*, Appl. Comput. Harmon. Anal., 18 (2005), pp. 25–66.
- [9] A. COHEN, *Non-stationary multiscale analysis*, in Wavelets: Theory, Algorithms, and Applications (Taormina, 1993), Wavelet Anal. Appl. 5, Academic Press, San Diego, CA, 1994, pp. 3–12.

- [10] A. COHEN AND N. DYN, *Nonstationary subdivision schemes and multiresolution analysis*, SIAM J. Math. Anal., 27 (1996), pp. 1745–1769.
- [11] A. COHEN AND R. D. RYAN, *Wavelets and Multiscale Signal Processing*, Chapman & Hall, London, 1995.
- [12] I. DAUBECHIES, *Orthonormal bases of compactly supported wavelets*, Comm. Pure Appl. Math., 41 (1988), pp. 909–996.
- [13] I. DAUBECHIES, *Ten Lectures on Wavelets*, CBMS-NSF Regional Conf. Ser. Appl. Math. 61, SIAM, Philadelphia, 1992.
- [14] I. DAUBECHIES, B. HAN, A. RON, AND Z. SHEN, *Framelets: MRA-based constructions of wavelet frames*, Appl. Comput. Harmon. Anal., 14 (2003), pp. 1–46.
- [15] G. DERFEL, N. DYN, AND D. LEVIN, *Generalized functional equations and subdivision processes*, J. Approx. Theory, 80 (1995), pp. 272–297.
- [16] B. DONG AND Z. SHEN, *Pseudosplines, wavelets and framelets*, Appl. Comput. Harmon. Anal., 22 (2007), pp. 78–104.
- [17] N. DYN AND D. LEVIN, *Subdivision schemes in geometric modeling*, Acta Numer., 11 (2002), pp. 73–144.
- [18] S. S. GOH, Z. Y. LIM, AND Z. SHEN, *Symmetric and antisymmetric tight wavelet frames*, Appl. Comput. Harmon. Anal., 20 (2006), pp. 411–421.
- [19] B. HAN, *On dual wavelet tight frames*, Appl. Comput. Harmon. Anal., 4 (1997), pp. 380–413.
- [20] B. HAN, *Vector cascade algorithms and refinable function vectors in Sobolev spaces*, J. Approx. Theory, 124 (2003), pp. 44–88.
- [21] B. HAN, *Compactly supported tight wavelet frames and orthonormal wavelets of exponential decay with a general dilation matrix*, J. Comput. Appl. Math., 155 (2003), pp. 43–67.
- [22] B. HAN, *Refinable functions and cascade algorithms in weighted spaces with Hölder continuous masks*, SIAM J. Math. Anal., 40 (2008), pp. 70–102.
- [23] B. HAN, *Symmetric orthonormal complex wavelets with masks of arbitrarily high linear-phase moments and sum rules*, Adv. Comput. Math., to appear.
- [24] B. HAN AND R.-Q. JIA, *Multivariate refinement equations and convergence of subdivision schemes*, SIAM J. Math. Anal., 29 (1998), pp. 1177–1199.
- [25] B. HAN AND Q. MO, *Tight wavelet frames generated by three symmetric B-spline functions with high vanishing moments*, Proc. Amer. Math. Soc., 132 (2004), pp. 77–86.
- [26] B. HAN AND Q. MO, *Splitting a matrix of Laurent polynomials with symmetry and its application to symmetric framelet filter banks*, SIAM J. Matrix Anal. Appl., 26 (2004), pp. 97–124.
- [27] B. HAN AND Q. MO, *Symmetric MRA tight wavelet frames with three tight generators and high vanishing moments*, Appl. Comput. Harmon. Anal., 18 (2005), pp. 67–93.
- [28] K. JETTER AND D. X. ZHOU, *Order of linear approximation from shift-invariant spaces*, Constr. Approx., 11 (1995), pp. 423–438.
- [29] K. JETTER AND D. X. ZHOU, *Approximation Order of Linear Operators and Finitely Generated Shift-Invariant Spaces*, preprint, 1998.
- [30] W. LAWTON, *Applications of complex valued wavelet transforms to subband decomposition*, IEEE Trans. Signal. Process., 41 (1993), pp. 3566–3568.
- [31] W. LAWTON, S. L. LEE, AND Z. SHEN, *Stability and orthonormality of multivariate refinable functions*, SIAM J. Math. Anal., 28 (1997), pp. 999–1014.
- [32] W. LAWTON, S. L. LEE, AND Z. SHEN, *Convergence of multidimensional cascade algorithm*, Numer. Math., 78 (1998), pp. 427–438.
- [33] A. RON AND Z. SHEN, *Affine systems in $L_2(\mathbb{R}^d)$: The analysis of the analysis operator*, J. Funct. Anal., 148 (1997), pp. 408–447.
- [34] A. RON AND Z. SHEN, *Generalized shift invariant systems*, Constr. Approx., 22 (2005), pp. 1–45.
- [35] V. L. RVACHEV AND V. A. RVACHEV, *A certain finite function*, Dopov. Dokl. Akad. Nauk. Ukraini, 8 (1971), pp. 705–707.
- [36] I. W. SELESNICK, *Smooth wavelet tight frames with zero moments*, Appl. Comp. Harmon. Anal., 10 (2000), pp. 163–181.
- [37] Z. SHEN, *Refinable function vectors*, SIAM J. Math. Anal., 29 (1998), pp. 235–250.

STABILITY AND SYNCHRONISM OF CERTAIN COUPLED DYNAMICAL SYSTEMS*

JOSÉ A. BARRIONUEVO† AND JACQUES A. L. SILVA†

Abstract. We obtain sufficient conditions for the stability of the synchronized solution for certain classes of coupled dynamical systems. These discrete time systems can be used to describe population patches coupled by migration, which are density and/or time dependent. Our results follow from an analytic expression for the transverse Lyapunov exponent obtained through spectral analysis. We then indicate some applications to population dynamics.

Key words. coupled dynamical systems, synchronism

AMS subject classifications. 37A30, 37H15, 37N30

DOI. 10.1137/060658436

1. Introduction. The study of coupled dynamical systems has received considerable attention recently for its interest from the mathematical, physical, and biological points of view; see, for instance, [22], [1], [7], [29], [26], and [3], among others. One concern about such systems is whether or not they will present synchronization phenomena and whether or not such synchronization is stable.

The system below describes the evolution of a system consisting of d identical subsystems where, on every iteration, each subsystem undergoes its common local evolution determined by f , followed by a density-dependent coupling process encoded by C and φ . The system can model a population consisting of d patches, x_j , $j = 1, \dots, d$, where, in the absence of migration, patch j is controlled by a *local dynamics* $x_{t+1}^j = f(x_t^j)$. When migration is present, $\varphi(f(x_t^j))$ individuals leave patch j and are distributed with density c_{ji} on patch i . The *global dynamics* is then

$$(1) \quad x_{t+1}^j = f(x_t^j) - \varphi(f(x_t^j)) + \sum_{i=1}^d c_{ji} \varphi(f(x_t^i)); \quad i, j = 1, \dots, d.$$

Here f is a bounded C^1 map on $[0, \infty)$, $C = [c_{ij}]$ is doubly stochastic, that is, $c_{ij} \geq 0$, and for all i, j , $\sum_{i=1}^d c_{ij} = \sum_{j=1}^d c_{ij} = 1$. Furthermore we will assume φ differentiable a.e. with φ' bounded. The models in population dynamics considered in [1], [7], [8], [27], [13], [26] are all particular cases of (1) for special choices of C and φ .

The condition on C being doubly stochastic reflects that there are no losses during the migration process. It is also necessary for the invariance of the diagonal of the phase space, that is, for $x_t^j = x_t^i = x_t$ to be a solution of (1), where each x_t^j satisfies $x_{t+1}^j = f(x_t^j)$.

In this paper we obtain sufficient conditions for the stability of the aforementioned synchronized solutions. The criteria involve the Lyapunov exponents, to be defined below, of the one-dimensional map f and of the codimension one transverse dynamical system.

*Received by the editors April 28, 2006; accepted for publication (in revised form) May 5, 2008; published electronically September 17, 2008.

<http://www.siam.org/journals/sima/40-3/65843.html>

†Department of Mathematics, Universidade Federal RS, Av. Bento Gonçalves 9500, 91509-900 Porto Alegre, RS, Brazil (josea@mat.ufrgs.br, jaqx@mat.ufrgs.br).

The paper is organized as follows. In the next section we provide a criterion for stability for general systems. In section 3 we improve our result for the case of normal operators. In section 4 we consider a system where the coupling/migration process is time dependent. We formulate and prove the corresponding results of the previous sections in this setting. In the last section we indicate some applications to population dynamics.

Previous results treated only the cases where φ' is a constant or a 2-valued step function, as well as particular examples of matrices C , and these cases are covered by our results. Our treatment only requires φ' to be bounded and C to be doubly stochastic and irreducible. This last condition is easily seen to be necessary; see below. Moreover, we are unaware of any treatment of the systems as in section 4 in the literature. Thus we extend some of the results of [7], [8], [10], [11], [27], [15], [3], [13], and others to these more general situations.

2. Stability: General case. In order to understand the behavior of orbits starting at nearby points of the diagonal of the phase space, we first linearize (1). If $J_t = [\alpha_{ij}]$ denotes the Jacobian matrix of (1) restricted to the synchronized orbit, we have

$$\alpha_{ij} = \begin{cases} f'(x_t)(1 - (1 - c_{ii})\varphi'(f(x_t))) & \text{for } i = j, \\ f'(x_t)\varphi'(f(x_t))c_{ij} & \text{for } i \neq j. \end{cases}$$

We have $J_t = f'(x_t)H_t$, where $H_t = I - \varphi'(f(x_t))B$, and $B = I - C$.

We will assume that C is irreducible. This is almost a necessary condition, for otherwise it permits the existence of uncoupled unsynchronized subsystems that are each synchronized. In this case we can apply the Fröbenius theorem [12] to show that $\lambda = 1$ is the simple dominant eigenvalue of C , associated to the eigenvector $v = (1, \dots, 1)$. This furnishes the decomposition $\mathbb{R}^d = \mathbb{R}v \oplus W$, where W is a C -invariant $(d - 1)$ -dimensional subspace. Under these conditions,

$$(2) \quad B = P^{-1} \begin{bmatrix} 0 & \\ & A \end{bmatrix} P,$$

where P is the matrix of the appropriate change of basis. This decomposition implies that the stability of the synchronized solution of (1) is a consequence of the stability of the trivial solution of the transversal component, w_t , which satisfies

$$(3) \quad w_{t+1} = f'(x_t)(I - \varphi'(f(x_t))A)w_t.$$

We will show that under a certain integrability condition the map above is in fact a contraction, which in turn implies the stability of $w_t \equiv 0$. The analysis of (3) will be based on the Lyapunov exponents (see [17], [18]) of (3). Define

$$K_n(x) = \left\| \prod_{k=0}^{n-1} f'(f^k(x))(I - \varphi'(f^{k+1}(x))A) \right\|^{1/n},$$

where $f^0(x) = x$ and $f^k(x) = f(f^{k-1}(x))$ for $k > 0$. Clearly if $\mathcal{K} = \limsup K_n$ satisfies $\mathcal{K} < 1$, we have that (3) is a contraction. Now observe

$$(4) \quad K_n(x) = L_n(x)\Lambda_n(x) = \left\| \prod_{k=0}^{n-1} f'(f^k(x)) \right\|^{1/n} \left\| \prod_{k=0}^{n-1} I - \varphi'(f^{k+1}(x))A \right\|^{1/n}.$$

$L_n(x)$ depends only on the local dynamics f , while $\Lambda_n(x)$ reflects also the effects of φ and C . Let ρ be an invariant measure of the local system. Define for $x > 0$, $\ln^+(x) = \max(\ln(x), 0)$. By Birkhoff's ergodic theorem, if $\ln^+ |f'| \in L^1(\rho)$, there exists $\lim_n \exp(\frac{1}{n} \sum_{k=0}^{n-1} \ln |f'(f^k(x))|)$ for ρ -a.e. x . For ρ ergodic, this limit, call it L , is independent of x and is given by $\exp(\int_0^\infty \ln |f'(s)| d\rho(s))$. L is the Lyapunov exponent of the local system governed by f .

Similarly the ergodic theorem of Oseledec [18] implies that if

$$(5) \quad \int_0^\infty \ln^+ \|I - \varphi'(s)A\| d\rho(s) < \infty,$$

there exists $\lim_n \Lambda_n(x) =: \Lambda(x)$ for ρ -a.e. x , and this limit is independent of x provided ρ is ergodic. In the rest of the paper we will assume that the f -invariant measure, ρ , is ergodic. Nonergodic ρ can be treated at the cost of additional technicalities. See remark (iii) in section 2.1. Our first result is the following.

THEOREM 1. *Consider the system (1), where f is a C^1 map, φ' is bounded, and C is doubly stochastic and irreducible such that $\ln^+ \|I - \varphi'(s)A\|$ and $\ln^+ |f'(s)|$ are in $L^1(\rho)$, where ρ is an ergodic f -invariant measure. Let $L = \lim_n L_n(x)$ and $\Lambda = \lim_n \Lambda_n(x)$ be as above. If $L \Lambda < 1$, there exists a set E with $\rho(E) = 1$ such that for all $x \in E$, the synchronized solution of (1) is asymptotically stable.*

Proof. By Birkhoff's theorem there exists a set E with $\rho(E) = 1$ such that for all x in E , $\lim_n \Lambda_n(x) = \Lambda$. We claim that

$$(6) \quad \Lambda \leq \exp\left(\int_0^\infty \ln^+ \|I - \varphi'(s)A\| d\rho(s)\right).$$

By the continuity of the norm and the function $\ln^+(\cdot)$, the dominated convergence theorem implies that we need only prove (6) for φ' simple.

Let $\varphi'(x) = \sum_{k=1}^r a_k \chi_{E_k}(x)$, where E_k are measurable and disjoint with $\rho(E_k) > 0$ and such $E \subset \bigcup_k E_k$. For $1 \leq k \leq r$, and x in E , define $\rho_{k,n} = \rho_{k,n}(x)$ by

$$\rho_{k,n} = \frac{\#\{0 \leq j < n : f^j(x) \in E_k\}}{n}.$$

To keep the notation simple we omit the dependence of $\rho_{k,n}$ on x since Birkhoff's ergodic theorem applied to χ_{E_k} shows that for $1 \leq k \leq r$, and all x in E , we have

$$\lim_n \rho_{k,n} = \rho(E_k).$$

This gives

$$\Lambda_n(x) = \left\| \prod_k (I - a_k A)^{n\rho_{k,n}} \right\|^{1/n},$$

which imply

$$\begin{aligned} \Lambda_n(x) &\leq \prod_k \|I - a_k A\|^{\rho_{k,n}} \text{ and} \\ \ln \Lambda_n(x) &\leq \sum_k \rho_{k,n} \ln \|I - a_k A\|. \end{aligned}$$

Therefore, given $\epsilon > 0$, there exists n_0 such that for $n > n_0$, we have $\rho_{k,n} \leq (1 + \epsilon)\rho(E_k)$ and $L_n(x) \leq (1 + \epsilon)L$. Since $\ln(\cdot) \leq \ln^+(\cdot)$,

$$\ln \Lambda_n(x) \leq (1 + \epsilon) \sum_{k=1}^r \ln^+ \|I - a_k A\| \rho(E_k) = (1 + \epsilon) \int_0^\infty \ln^+ \|I - \varphi'(s)A\| d\rho(s).$$

This proves (6). Note that the right-hand side is independent of x . Thus for $n > n_0$,

$$K_n(x) = L_n(x) \Lambda_n(x) \leq (1 + \epsilon)L \Lambda^{1+\epsilon},$$

and since $\epsilon > 0$ is arbitrary we have

$$\lim_n K_n(x) \leq L \Lambda < 1,$$

and thus the transversal map (3) is a contraction. \square

2.1. Remarks.

(i) If $E = [0, \infty)$, $L \Lambda \leq 1$ is necessary for the stability of the synchronized solution.

(ii) In certain cases, which include the case H_t semisimple, we will prove below that the norm $\|I - \varphi'(f^k(x))A\|$ is simply the spectral radius, $\sigma_{-1}(H_t)$, of the restriction of H_t to the subspace W . Therefore, with the same hypothesis on f and φ , we have

$$(7) \quad \Lambda \leq \Lambda_1 = \exp \left(\int_0^\infty \ln \sigma_{-1}(H_{\varphi'(s)}) d\rho(s) \right),$$

and thus $L\Lambda_1 < 1$ is a sufficient condition for the stability of the synchronized solution of (1).

(iii) If ρ fails to be ergodic, a deep result of Choquet (Theorem 31.3 in [6]) gives the representation

$$(8) \quad \rho = \int_{\mathcal{E}(P_1)} \pi d\nu(\pi),$$

where P_1 is the space of f -invariant probability measures and $\mathcal{E}(P_1)$ is the set of extreme points of P_1 , that is, the ergodic f -invariant probability measures, and ν is a probability measure on $\mathcal{E}(P_1)$. Also equality in (8) is, in the weak sense,

$$\int \Psi d\rho = \int_{\mathcal{E}(P_1)} \left(\int \Psi d\pi \right) d\nu(\pi)$$

for all continuous Ψ . Given $\pi \in \mathcal{E}(P_1)$, we apply Theorem 1 to obtain a set E_π of full π measure such that for all $x \in E_\pi$, the limit $\Lambda_\pi = \lim_n \Lambda_n(x)$ exists and

$$\Lambda_\pi \leq \exp \left(\int_0^\infty \ln^+ \|I - \varphi'(s)A\| d\pi(s) \right).$$

Similarly, if $\ln^+ |f'| \in L^1(\pi)$, by Birkhoff's theorem, there exists $L_\pi = \exp(\lim_n \frac{1}{n} \sum_{k=0}^{n-1} \ln |f'(f^k(x))|)$ for $x \in E_\pi$, and it is given by $L_\pi = \exp(\int_0^\infty \ln |f'(s)| d\pi(s))$. From Theorem 1 it follows that if $L_\pi \Lambda_\pi < 1$, then the synchronized solution of (1) is stable for all initial x in $D_\pi = \{(x, x, \dots, x) : x \in E_\pi\}$. Since any two distinct ergodic f -invariant measures are mutually singular, the sets E_π form a partition of the whole space $[0, \infty)$ modulo null sets. Therefore, if $\sup_{\pi \in \mathcal{E}(P_1)} L_\pi \Lambda_\pi < 1$, we have that the synchronized solution of (1) is asymptotically stable. Since the set of ergodic invariant measures can be very complicated, it might be difficult to verify the above condition in general.

3. Stability: Normal operators. In the case of normal operators, that is, $AA^* = A^*A$, one can improve the previous result with the help of the functional calculus of [16] and [23]. Let L_n be as in (4) and define the operator valued cocycle

$$(9) \quad \mathcal{O}_n(x) = \prod_{k=0}^{n-1} (I - \varphi'(f^{k+1}(x))A)$$

so that we have $K_n(x) = L_n(x) \|\mathcal{O}_n\|^{1/n}$. The proof of Oseledec’s theorem in [24] shows that

$$(10) \quad \lim_n (\mathcal{O}_n^*(x)\mathcal{O}_n(x))^{1/2n} = \mathcal{O}(x)$$

exists a.e. in operator norm. The ergodicity of ρ implies that $\mathcal{O}(x)$ is independent of x ; call it \mathcal{O} . Spectral analysis of \mathcal{O} determines the Lyapunov exponents and respective subspaces. Oseledec’s theorem was extended to operator valued cocycles in Hilbert and Banach spaces in [25] and [21] under an additional compactness hypothesis and with the operator norm convergence replaced by strong convergence in (10). The *very special* nature of our cocycles (9) allows us to determine \mathcal{O} above whenever A is a bounded, normal, and not necessarily compact operator on a Hilbert space \mathcal{H} . We present the result in this generality for its own sake and because it involves little extra work. Note, however, that the study of the stability of synchronized solutions of (1) for $d = \infty$ offers additional difficulties; see the remarks in section 3.1.

If A is normal, Theorem 12.23 in [23] gives

$$(11) \quad A = \int_{\sigma(A)} \lambda dP(\lambda),$$

where $dP(\lambda)$ are the spectral projections associated with A ; that is, for all u, v in \mathcal{H} , $d\langle P(\lambda)u, v \rangle$ is a complex measure and

$$(12) \quad \langle Au, v \rangle = \int_{\sigma(A)} \lambda d\langle P(\lambda)u, v \rangle.$$

For a bounded measurable function $F : \sigma(A) \rightarrow \mathbb{C}$, $F(A)$ is defined by

$$(13) \quad F(A) = \int_{\sigma(A)} F(\lambda) dP(\lambda).$$

Observe that

$$(14) \quad \|F(A)\| \leq \|F(\lambda)\|_{L^\infty(\sigma(A))},$$

implying that if $F_n(\lambda) \rightarrow F(\lambda)$ in $L^\infty(\sigma(A))$, $F_n(A) \rightarrow F(A)$ in operator norm.

Setting $\ln 0 = -\infty$, $e^{-\infty} = 0$, and recalling that an *upper semicontinuous* real valued function satisfies that for all $\alpha \in \mathbb{R}$, $\{\lambda : F(\lambda) < \alpha\}$ is open, we now prove the following.

PROPOSITION 2. *The function $F : \mathbb{C} \rightarrow \mathbb{R}$ defined by*

$$(15) \quad F(\lambda) = \exp \left(\int_0^\infty \ln |1 - \lambda \varphi'(s)| d\rho(s) \right)$$

is upper semicontinuous and therefore bounded on compact subsets of \mathbb{C} .

Proof. First note that $F(\lambda) \geq 0$ and

$$(16) \quad F(\lambda) = 0, \text{ if and only if } \int_0^\infty \ln |1 - \lambda\varphi'(s)| d\rho(s) = -\infty.$$

Moreover, f bounded implies that the support of ρ , $\text{supp } \rho$, is compact and Jensen's inequality gives

$$(17) \quad F_M(\lambda) \leq \int_{\text{supp } \rho} |1 - \lambda\varphi'(s)| d\rho(s) \leq 1 + K|\lambda| < \infty,$$

where $K = \sup \{ |\varphi'(s)| : s \in \text{supp } \rho \}$. Thus $F(\lambda)$ is a well-defined nonnegative real number. Defining

$$F_M(\lambda) = \exp \left(\int_0^\infty g_M(\lambda, s) d\rho(s) \right), \text{ where } g_M(\lambda, s) = \max \{ -M, \ln |1 - \lambda\varphi'(s)| \},$$

we have $F_{M+1}(\lambda) \leq F_M(\lambda)$, and monotone convergence gives $F(\lambda) = \lim_M F_M(\lambda)$. Since $g_M(\lambda, s)$ is continuous with respect to λ on $\mathbb{C} \times \text{supp } \rho$ and $\text{supp } \rho$ is compact, F_M is continuous on \mathbb{C} . Thus F is the pointwise limit of a nonincreasing sequence of continuous functions and is therefore upper semicontinuous. \square

The proposition enables us to define $\mathcal{O} = F(A)$,

$$(18) \quad \mathcal{O} = \int_{\sigma(A)} \exp \left(\int_0^\infty \ln |1 - \lambda\varphi'(s)| d\rho(s) \right) dP(\lambda).$$

The spectral mapping theorem then implies

$$(19) \quad \sigma(\mathcal{O}) = F(\sigma(A)) = \{F(\lambda) : \lambda \in \sigma(A)\}.$$

The application of (10) to synchronization is contained in the next theorem.

THEOREM 3. *Let f, φ', C , and L be as in Theorem 1. Assume A is normal and let Λ be the spectral radius of (18). Then, if $L \Lambda < 1$, the synchronized solution of (1) is asymptotically stable.*

Proof. As in Theorem 1 the inequality $L \Lambda < 1$ implies that the transverse component (3) is a contraction, implying the stability of the synchronized solution. Therefore, in order to prove the theorem we need only establish (10) in the strong sense, with \mathcal{O} given by (18). We first consider φ' simple. Let $\varphi' = \sum_{k=1}^N \varphi_k \chi_{E_k}$, where E_k are disjoint measurable. As before, we define

$$\rho_{k,n} = \frac{\#\{0 \leq j < n : f^j(x) \in E_k\}}{n}.$$

Then

$$(20) \quad (\mathcal{O}_n^* \mathcal{O}_n)^{\frac{1}{2n}} = \prod_{k=1}^N |I - \varphi_k A|^{\rho_{k,n}} = \int_{\sigma(A)} \left(\prod_{k=1}^N |1 - \lambda\varphi_k|^{\rho_{k,n}} \right) dP(\lambda).$$

For each n , $F_n(\lambda) = \prod_{k=1}^N |1 - \lambda\varphi_k|^{\rho_{k,n}}$ is continuous. Moreover, for $k = 1, \dots, N$, Birkhoff's theorem gives $\lim_n \rho_{n,k} = \rho(E_k) > 0$. Since $\sigma(A)$ is compact, we have that

$$\lim_n F_n(\lambda) = F(\lambda) = \prod_{k=1}^N |1 - \lambda\varphi_k|^{\rho(E_k)}$$

uniformly for λ in $\sigma(A)$. We can rewrite this as

$$\lim_n F_n(\lambda) = \exp \left(\sum_{k=1}^N \rho(E_k) \ln |1 - \lambda \varphi_k| \right) = \exp \left(\int_0^\infty \ln |1 - \lambda \varphi'(s)| d\rho(s) \right),$$

where the last equality is just the definition of the Lebesgue integral. By (14), we have (10) for φ' simple. For an arbitrary φ' , the right-hand side of (10) is already well defined by (18). We now have

$$(\mathcal{O}_n^* \mathcal{O}_n)^{\frac{1}{2n}} = \int_{\sigma(A)} F_n(\lambda) dP(\lambda),$$

where $F_n(\lambda) = (\prod_{k=1}^n |1 - \lambda \varphi'(x_k)|)^{1/n}$. Since one can approximate φ' uniformly by simple functions, the ergodic theorem and the above result for simple functions imply that for all λ in $\sigma(A)$, $\lim_n F_n(\lambda) = F(\lambda)$ -pointwise convergence.

Thus for all u, v in \mathcal{H} , the Lebesgue dominated convergence theorem implies

$$\lim_n \int_{\sigma(A)} F_n(\lambda) \langle dP(\lambda)u, v \rangle = \int_{\sigma(A)} F(\lambda) \langle dP(\lambda)u, v \rangle.$$

Since $dP(\lambda)$ is a resolution of the identity, this implies that

$$\lim_n \left\| \left((\mathcal{O}_n^* \mathcal{O}_n)^{\frac{1}{2n}} - F(A) \right) u \right\| = 0$$

for all u in \mathcal{H} , which establish strong convergence in (10). This finishes the proof of Theorem 3. \square

3.1. Remarks.

(i) In finite dimension we have $\sigma(A) = \{\lambda_1, \dots, \lambda_M\}$ with corresponding spaces E_j , mutually orthogonal, with $\mathbb{R}^{d-1} = E_1 \oplus \dots \oplus E_M$, and the conclusion of the theorem is that if $u \neq 0$ is in E_j , then $\lim_n \|(\mathcal{O}_n^* \mathcal{O}_n)^{\frac{1}{2n}} u\| = F(\lambda_j)$.

(ii) Since for all $\lambda \in \sigma(A)$, $|1 - \lambda \varphi'(s)| \leq \|I - \varphi'(s)A\|$, we have $\Lambda \leq \exp(\int_0^\infty \ln^+ \|I - \varphi'(s)A\| d\rho(s))$, and therefore Theorem 3 is an improvement of Theorem 1.

(iii) By definition $\Lambda = \sup\{|F(\lambda)| : \lambda \in \sigma(A)\}$, and thus (19) gives (7). All the effect of coupling to synchronization is reflected in Λ and can therefore, in certain cases, be independent of d ; see section 5.

(iv) The case $d = \infty$ is more subtle. First, the choice of a function space for $\{x_j(t)\}_j$ plays an important role. For instance, the synchronized solution belongs to l^p only for $p = \infty$. Second, the Fröbenius theorem is not valid in the same generality. Third, if the function space considered is not a Hilbert space, the functional calculus above is not available.

4. Extensions. The analysis above can be applied to more general systems.

One such instance is the case below, where the coupling/migration process no longer depends on the density but instead obeys a seasonal dynamics; that is, on each cycle a time-dependent fraction μ_t of each patch will mix according to a time-dependent distribution $C_t = [c_{ij}^t]$ as follows:

$$(21) \quad x_{t+1}^j = f(x_t^j) - \mu_t f(x_t^j) + \sum_{i=1}^d c_{ji}^t \mu_t f(x_t^i); \quad i, j = 1, \dots, d.$$

Again f is a C^1 map on $[0, \infty)$. The evolution of the migration parameters μ_t and C_t is governed by maps g and h on $[0, 1]$. The idea is to write $\mu_{t+1} = g(\mu_t)$ and $C_{t+1} = h(C_t)$, but this is not accurate. To make this precise we argue as follows. For g continuous on $[0, 1]$ and μ_0 arbitrary, we assume that $\mu_t = \mu_t(\mu_0)$ is given by $\mu_{t+1} = g(\mu_t)$. Similarly if $\{C_s\}_{s \in [0,1]}$ is a family of doubly stochastic irreducible matrices, h is a continuous map on $[0, 1]$, and s_0 is arbitrary, we define $s_{t+1} = h(s_t)$ and $C_t = C_t(s_0)$ by $C_{t+1} = C_{h(s_t)}$. Assume that $G(\mu, s) = \mu C_s$ is a measurable operator valued map with respect to the product measure $\nu \times \eta$ on $[0, 1] \times [0, 1]$, where ν and η are, respectively, a g -invariant and h -invariant *ergodic* measure on $[0, 1]$. Under these assumptions the diagonal of the phase space, $x_t^i = x_t^j = x_t$, is a synchronized solution and we are interested in its stability. This system falls under the general theory of random dynamical systems as presented in [2] and thus obeys, under the appropriate integrability condition, a multiplicative ergodic theorem. Our approach is direct and avoids the use of this theory. Moreover, due to the particular nature of (21), we are able, as in the previous sections, to obtain much more precise information about the limit operators than is given by the ergodic theorem alone.

The Jacobian matrix of (21), $J_t = [\alpha_{ij}^t]$, is now given by

$$\alpha_{ij} = \begin{cases} f'(x_t) (1 - \mu_t(1 - c_{ii}^t)) & \text{for } i = j, \\ f'(x_t) \mu_t c_{ij}^t & \text{for } i \neq j, \end{cases}$$

yielding $J_t = f'(x_t)(I - \mu_t B_t)$, where $B_t = I - C_t$. Our first task is to decompose the linearized system into diagonal and transversal components and then estimate the respective Lyapunov exponents by the appropriate ergodic theorem. This cannot be done in general for arbitrary families $\{C_t\}$; therefore we will consider special cases of increasing generality.

4.1. Simultaneous diagonalizable matrices. This includes the families of commuting symmetric matrices and of circulant matrices. In this case there exists a matrix P such that for all t , $B_t = P^{-1} M_t P$, where M_t is a diagonal matrix with entries $\{0, \lambda_2(t), \dots, \lambda_d(t)\}$ in its diagonal. Each $\lambda_j(t)$ is measurable and satisfies $0 < |1 - \lambda_j(t)| < 1$ by the Fröbenius theorem (applied to C_t). In the symmetric (commuting) case the $\lambda_j(t)$ are real and lie in $(0, 2)$. We obtain that the transversal component of (21) satisfies

$$(22) \quad w_{t+1} = f'(x_t) (I - \mu_t D_t) w_t,$$

where D_t is the $(d - 1)$ -dimensional $\{\lambda_2(t), \dots, \lambda_d(t)\}$ diagonal matrix. As before, the synchronized solution of (21) is stable if and only if $w_t = 0$ is a stable solution of (22). For each x, μ_0, s_0 in $[0, 1]$, define $L_n(x)$ as in (4), and for $j = 2, \dots, d$, define

$$\Lambda_{j,n}(\mu_0, s_0) = \left(\prod_{k=0}^{n-1} |1 - g^k(\mu_0) \lambda_j(h^k(s_0))| \right)^{1/n}.$$

The analogue of Theorem 3 is as follows.

THEOREM 4. *Consider the system (21) under the above conditions. Assume in addition that*

- (i) $\ln^+ |f'(s)|$ belongs to $L^1(\rho)$, where ρ is an ergodic f -invariant probability measure on $[0, \infty)$.
- (ii) for all $2 \leq j \leq d$, $\ln^+ |1 - \mu \lambda_j(s)|$ belong to $L^1(\nu \times \eta)$.

Then there exist sets $E \subset [0, \infty)$ and $F \subset [0, 1] \times [0, 1]$ with $\rho(E) = (\nu \times \eta)(F) = 1$, such that for all $x \in E$ and $(\mu_0, s_0) \in F$, the limits $L = \lim_n L_n(x)$ and $\Lambda_j = \lim_n \Lambda_{j n}(\mu_0, s_0)$ exist and are independent of x and (μ_0, s_0) . Moreover, if $\Lambda = \max_j \Lambda_j$ satisfies $L\Lambda < 1$, the synchronized solution of (21) is asymptotically stable.

Proof. The existence of $L(x)$ for x in a set of full measure follows, as before, from Birkhoff’s ergodic theorem applied to $\ln L_n(\cdot)$ and (i). In addition, the limit is independent of x since ρ is ergodic. The decomposition (22) allows us to derive the transversal component directly without making use of Oseledec’s theorem. For $2 \leq j \leq d$, $(\mu_0, s_0) \in [0, 1] \times [0, 1]$, we can write

$$\Lambda_{j n}(\mu_0, s_0) = \exp \left(\frac{1}{n} \sum_{k=0}^{n-1} \ln |1 - g^k(\mu_0) \lambda_j(h^k(s_0))| \right).$$

Condition (ii) and a $(d - 1)$ -fold application of Birkhoff’s ergodic theorem imply the existence of a set F , with $\nu \times \eta(F) = 1$ such that for all j and all $(\mu_0, s_0) \in F$, there exists $\Lambda_j(\mu_0, s_0) = \lim_n \Lambda_{j n}(\mu_0, s_0)$. Since, in addition, $\nu \times \eta$ is ergodic, the limit is independent of (μ_0, s_0) and is given by

$$(23) \quad \Lambda_j = \exp \left(\int_{[0,1] \times [0,1]} \ln |1 - \mu \lambda_j(s)| d(\nu \times \eta)(\mu, s) \right).$$

This gives $\text{diag}\{\Lambda_2, \dots, \Lambda_d\}$ as the analogue of (18). Its spectral radius is then $\Lambda = \max_j \Lambda_j$. Thus, if $L\Lambda < 1$, the synchronized solution is asymptotically stable. \square

4.2. Symmetric (noncommuting) matrices. In this case we no longer have a decomposition yielding a diagonal transversal component, like (22), and therefore our conclusions are somewhat weaker than the previous theorem. Since each C_t is symmetric, doubly stochastic, and irreducible, $\lambda = 1$ is the dominant eigenvalue with corresponding eigenvector $v = \{1, \dots, 1\}$. Moreover, the $(d - 1)$ -dimensional subspace $W = (\mathbb{R}v)^\perp$ is invariant for all C_t . This provides a decomposition like (2),

$$(24) \quad B_t = P^{-1} \begin{bmatrix} 0 & \\ & A_t \end{bmatrix} P.$$

Accordingly, the transversal component of (21) now satisfies

$$(25) \quad w_{t+1} = f'(x_t) (I - \mu_t A_t) w_t.$$

The following theorem is the analogue of Theorem 1 in the present situation.

THEOREM 5. *Consider the system (21), where f and μ_t are as in Theorem 4, and assume $\{C_t\}$ as above. For each (μ_0, s_0) , define*

$$\Lambda_n(\mu_0, s_0) = \left\| \prod_{k=0}^{n-1} (I - g^k(\mu_0) A_{h^k(s_0)}) \right\|^{1/n}.$$

Assume that $\ln^+ \|I - \mu A_{h(s)}\|$ belongs to $L^1(\nu \times \eta)$. Then there exists a set $F \subset [0, 1] \times [0, 1]$ of full $\nu \times \eta$ -measure such that for all $(\mu_0, s_0) \in F$, the limit $\Lambda = \lim_n \Lambda_n(\mu_0, s_0)$ exists and is independent of (μ_0, s_0) . Moreover, if Λ satisfies $L\Lambda < 1$, the synchronized solution of (21) is asymptotically stable.

Proof. Consider the map \mathcal{Q} on $[0, 1] \times [0, 1]$, given by $\mathcal{Q}(\mu, s) = (g(\mu), h(s))$. Clearly $\nu \times \eta$ is \mathcal{Q} invariant, and by assumption, $\Lambda_n(\mu, s)$ above defines a $(\nu \times \eta)$ -measurable cocycle. The existence of $\Lambda(\mu_0, s_0)$ for (μ_0, s_0) on a set of full measure

then follows from Oseledec’s ergodic theorem applied to $\Lambda_n(\mu, s)$. We also have that for a.e. (μ_0, s_0) ,

$$(26) \quad \Lambda(\mu_0, s_0) \leq \exp \left(\int_{[0,1] \times [0,1]} \ln^+ \|I - \mu A_{h(s)}\| d(\nu \times \eta)(\mu, s) \right).$$

The proof of (26) follows the same path as the proof of (6), and thus we will omit the details. Note that since each A_t is symmetric, there exists $\lambda(t) \in \sigma(A_t)$ for which $\|I - \mu A_t\| = |1 - \mu \lambda(t)|$. $\lambda(t)$ is the lowest eigenvalue of C_t . Moreover, since $\nu \times \eta$ is ergodic, $\Lambda(\mu_0, s_0)$, is independent of (μ_0, s_0) , and in this case we have

$$\Lambda = \exp \left(\int_{[0,1] \times [0,1]} \ln |1 - \mu \lambda(h(s))| d(\nu \times \eta)(\mu, s) \right)$$

as a consequence of Birkhoff’s theorem. As in the previous theorems, the condition $L \Lambda < 1$ implies the asymptotic stability of the synchronized solution. \square

4.3. Remarks.

(i) The assumptions on the continuity of g and h above can be relaxed. All that is needed is that each map possess an ergodic invariant probability measure, ν and η .

(ii) Similarly to the density-dependent migration, the results above *do not* extend to the $d = \infty$ case even though, once again, the effect of d is encoded in the joint spectrum of $\{C_t\}$.

4.4. Examples. Some special cases of (23) and (26) are worth mentioning.

(i) If s_0 is a periodic point for h , that is, $h^p(s_0) = s_0$, and we take $\eta = \frac{1}{p} \sum_{k=0}^{p-1} \delta_{h^k(s_0)}$, then (23) becomes

$$\Lambda_j = \exp \left(\frac{1}{p} \sum_{k=0}^{p-1} \int_{[0,1]} \ln |1 - \mu \lambda_j(s_k)| d\nu(\mu) \right).$$

A similar expression corresponding to the case of a periodic point for g also holds.

(ii) If $g = h$ and $\eta = \nu$, but $\mu_0 \neq s_0$, we get for (23),

$$\Lambda_j = \exp \left(\int_{[0,1] \times [0,1]} \ln |1 - \mu \lambda_j(s)| d(\nu \times \nu)(\mu, s) \right).$$

(iii) If in (ii) we have, in addition, $\mu_0 = s_0$, that is, the term $\mu_t A_t$ in (25) is of the form $g^t(s_0) A_{g^t(s_0)}$, then (23) becomes

$$\Lambda_j = \exp \left(\int_{[0,1]} \ln |1 - s \lambda_j(s)| d\nu(s) \right).$$

Similar corresponding expressions can be obtained, without difficulty, to represent (26) in the special cases above. We leave the details, as well as the formulation of other special instances of (26), to the interested reader. These formulas are useful in cases where ν and η are known, say Lebesgue, for they permit the exact calculations of the Lyapunov numbers of the transversal dynamical system.

5. Applications. In this section we will restrict ourselves to the density-dependent system (1) and leave the corresponding formulations of the results related to the time-dependent system (21) to the interested reader.

A direct problem consists of determining stability once f , φ , and C are given. In such situations we compute the eigenvalues $\lambda_1, \dots, \lambda_d$ of C . Discarding the eigenvalue 1, we evaluate the integrals defining the spectrum of Λ above by Birkhoff’s ergodic theorem

$$(27) \quad \int_0^\infty \ln |1 - \lambda\varphi'(s)| d\rho(s) = \lim_n \frac{1}{n} \sum_{k=0}^{n-1} \ln |1 - \lambda\varphi'(f^k(x))|,$$

and then the maximum of these values gives the spectral radius Λ and we can then determine stability once the Lyapunov exponent for f is known. We note that in order to use (27) we need the ergodic f -invariant measure, ρ , to be such that for x belonging to a set E of full Lebesgue measure, and all continuous functions Ψ ,

$$\lim_n \frac{1}{n} \sum_{j=0}^{n-1} \Psi(f^j(x)) = \int \Psi(s) d\rho(s).$$

Such measures, called physical or SRB (Sinai–Ruelle–Bowen) measures, are known to exist in certain important cases, such as expansive maps, piecewise monotonic maps of the interval, as well as *axiom A* diffeomorphisms, even though a recent result of Ávila and Bochi [4] shows that they are not typical in the C^1 -topology. Further information can be found in [31], [19], [20]. See also [32] for a survey.

In [5] we consider the following particular instance of (1). The local map f is taken from the Riker family, $f(x) = xe^{r(1-x)}$, $r > 0$, often considered in population dynamics. It is known that there exists an ergodic f -invariant probability SRB measure ρ which describes the asymptotic distribution of almost all orbits of f (see Theorems 25 and 29 in [30]). We let $\varphi(x) = x\mu(x)$, where the dispersal fraction, $\mu(x)$, is a sigmoidal function

$$(28) \quad \mu(x) = \frac{\alpha}{1 + e^{\beta(\gamma-x)}},$$

where the parameter $\alpha \in (0, 1)$ determines the maximal dispersal, and β describes the steepness of dispersal. Moreover, the sign of β determines whether dispersal is positively or negatively density dependent. γ is the inflection point, at which dispersal is half of its maximum. The cases $\beta = \pm\infty$ correspond to threshold population size triggering migration and have been considered in [8]. Here we consider the following two extreme types of coupling:

- (I) A ring configuration of d patches, each patch connected to the $2k$ -nearest neighbors, with uniform dispersion, that is, $c_{ij} = 1/2k$ for $0 < |i - j| \leq k$ and $c_{ij} = 0$ otherwise. In this case the eigenvalues of B are given by (see [27], [28], where the stability of equilibria is analyzed)

$$\lambda_j = 1 - \left(\frac{D_k\left(\frac{2\pi(j-1)}{d}\right) - 1}{2k} \right), \quad j = 1, 2, \dots, d,$$

where $D_k(x) = \sin(k + \frac{1}{2})x / \sin \frac{x}{2}$.

- (II) Global uniform coupling, that is, each of the d patches is connected to the others yielding $c_{ii} = 0$ and $c_{ij} = 1/(d - 1)$ for $i \neq j$. Here the eigenvalues of B are $\lambda_1 = 0$ and $\lambda_2 = \dots = \lambda_d = d/(d - 1)$.

For $i = \text{I, II}$, let $\Lambda_i = \Lambda_i(d, r, \alpha, \beta, \gamma)$ be the value of Λ , obtained using (27), for the respective configuration. In [5] we observe numerically the dependence of Λ_i on the different parameters and its effects on the stability of synchronization. For instance, one can prove (see Corollary 3.1 there) that $\lim_{d \rightarrow \infty} \Lambda_I \geq 1$, independent of the other parameters, which are kept fixed. This easily implies the impossibility of chaotic synchronization for sufficiently large ring configurations such as (I). The situation is different for (II), where, given any d , one can find a local chaotic map f ($L > 1$) and migration parameters, α , β , and γ , yielding a stable synchronized dynamics ($L\Lambda < 1$).

One feature of Theorems 1 and 3 is their continuous dependence on φ' , that is, C^1 -uniform topology. This follows directly from the formulas for Λ above. On the other hand, by analyzing (27) one can see that for certain f, C , it is possible to make Λ arbitrarily large by taking $\|\varphi'\|_{L^\infty(\rho)}$ large. In our case, where φ is given using (28), this implies that there exist situations where threshold triggered migration, $\beta = \pm\infty$, that present chaotic synchronization, while all systems with β greater than a certain β_0 have unstable synchronized orbits. The same reasoning proves that, in general, there is no continuity of Λ in φ in the C^0 -topology.

We can also make use of the results from previous sections to study an inverse problem. In this case f and φ are given and one is interested in finding a double stochastic matrix C , if possible, yielding the desired stability behavior for the synchronized solution of (1). In this regard, Theorem 7 below gives a partial result for symmetric matrices. We will make use of the following result from [14].

PROPOSITION 6. *The following are Corollaries 7 and 8 from [14].*

(i) *If $\{\lambda_2, \dots, \lambda_d\} \in [-1/(d-1), 1]$, then there exists a symmetric doubly stochastic matrix C with $\sigma(C) = \{1, \lambda_2, \dots, \lambda_d\}$.*

(ii) *If $\lambda \in (-1, 1]$, there exists a positive symmetric doubly stochastic matrix C such that $\lambda \in \sigma(C)$.*

The proof of Proposition 6 in [14] provides algorithms to find the matrices C in (i) and (ii) above.

THEOREM 7. *Let f and φ be as in Theorems 1 and 3. Let L denote the Lyapunov exponent of the one-dimensional map f . For $\lambda \in [0, 2]$, let $F(\lambda)$ be as in (15), and define $m = \inf F(\lambda)$ and $M = \sup F(\lambda)$. Then*

(i) *if $Lm > 1$, the synchronized solution of (1) is unstable for all symmetric configurations C .*

(ii) *if $LM < 1$, the synchronized solution of (1) is stable for all symmetric configurations C .*

(iii) *if $L \in (\frac{1}{M}, \frac{1}{m})$, it is possible to find a symmetric doubly stochastic matrix C such that the synchronized solution of (1) has a prescribed stability behavior.*

Proof. Let us first note that any nonnegative symmetric doubly stochastic matrix has its spectrum contained in $(-1, 1]$. This follows from Gershgorin's theorem [9] which states that $\sigma(C)$ is contained in the set $\{\lambda \in \mathbb{C} : \text{for all } i, |\lambda - c_{ii}| \leq \sum_{j \neq i} |c_{ij}|\}$. The hypotheses on C then imply that $\sigma(C) \subset [-1, 1]$. If C is positive, it is not difficult to show that -1 cannot be an eigenvalue of C . Since $B = I - C$, we have $\sigma(B) \subset [0, 2]$. Since $F(\lambda)$ gives the spectrum of $F(A)$, Theorem 7 now follows easily from Theorem 3 and Proposition 6. \square

Acknowledgments. The authors would like thank the anonymous referees for carefully reading an earlier version of this work and for providing a long list of improvements. We also thank A. Lopes for some helpful conversations.

REFERENCES

- [1] J. C. ALLEN, W. M. SCHAFFER, AND D. ROSKO, *Chaos reduces species extinction by amplifying local population noise*, *Nature*, 364 (1993), pp. 229–232.
- [2] L. ARNOLD, *Random Dynamical Systems*, Springer-Verlag, New York, 1998.
- [3] F. ATEY, T. BIYIKOĞLU, AND J. JOST, *Synchronization of Networks with Prescribed Degree Distributions*, preprint available online at <http://xxx.lanl.gov/abs/nlin/0407024> (2005).
- [4] A. ÁVILA AND J. BOCHI, *A generic C^1 map has no absolutely continuous invariant probability measure*, *Nonlinearity*, 19 (2006), pp. 2717–2725.
- [5] J. A. BARRIONUEVO, F. GIORDANI, AND J. A. L. SILVA, *Synchronism in population network with non linear coupling models*, submitted. Preprint available upon request.
- [6] G. CHOQUET, *Lectures on Analysis*, Vol. II: *Representation Theory*, W. A. Benjamin, Inc., New York, 1969.
- [7] D. J. D. EARN, S. LEVIN, AND P. ROHANI, *Coherence and conservation*, *Science*, 290 (2000), pp. 1360–1364.
- [8] J. A. L. SILVA AND F. T. GIORDANI, *Density-dependent migration and synchronism in metapopulations*, *Bull. Math. Biol.*, 68 (2006), pp. 451–465.
- [9] G. HÄMMERLIN AND K. HOFFMAN, *Numerical Mathematics*, Springer-Verlag, New York, 1991.
- [10] M. HASLER AND Y. MAISTRENKO, *An introduction to the synchronization of chaotic systems: coupled skew tent maps*, *IEEE Trans. Circuits Systems*, 10 (1997), pp. 856–866.
- [11] A. HASTINGS, *Complex interactions between dispersal and dynamics: Lessons from coupled logistic equations*, *Ecology*, 74 (1993), pp. 1362–1372.
- [12] A. HOUSEHOLDER, *The Theory of Matrices in Numerical Analysis*, Dover, New York, 1964.
- [13] Y. HUANG AND O. DIEKMANN, *Interspecific influence on mobility and Turing instability*, *Bull. Math. Biol.*, 65 (2003), pp. 143–156.
- [14] S. HWANG AND S. PYO, *The inverse eigenvalue problem for symmetric doubly stochastic matrices*, *Linear Algebra Appl.*, 379 (2004), pp. 77–83.
- [15] S. R. J. JANG AND A. K. MITRA, *Equilibrium stability of single species metapopulations*, *Bull. Math. Biol.*, 62 (2000), pp. 155–161.
- [16] T. KATO, *Perturbation Theory for Linear Operators*, Springer-Verlag, New York, 1995.
- [17] A. KATOK AND B. HASSELBLATT, *Introduction to the Modern Theory of Dynamical Systems*, Cambridge University Press, Cambridge, UK, 1995.
- [18] U. KRENGEL, *Ergodic Theorems*, Walter de Gruyter & Co., Berlin, 1985.
- [19] A. LASOTA AND J. YORKE, *On the existence of invariant measures for piecewise monotonic transformations*, *Trans. Amer. Math. Soc.*, 186 (1973), pp. 481–488.
- [20] T. Y. LI AND J. YORKE, *Ergodic transformations from an interval to itself*, *Trans. Amer. Math. Soc.*, 235 (1978), pp. 183–192.
- [21] R. MAÑÉ, *Lyapunov exponents and stable manifolds for compact transformations*, in *Geometric Dynamics*, Lecture Notes in Math. 1007, Springer, 1983, pp. 522–577.
- [22] A. PIKOVSKY, M. ROSENBLUM, AND J. KURTHS, *Synchronization—A Universal Concept in Nonlinear Sciences*, Cambridge University Press, Cambridge, UK, 2001.
- [23] W. RUDIN, *Functional Analysis*, 2nd ed., McGraw-Hill, New York, 1991.
- [24] D. RUELE, *Ergodic theory of differentiable dynamical systems*, *Publ. Math. Inst. Hautes Études Sci.*, 50 (1979), pp. 27–58.
- [25] D. RUELE, *Characteristic exponents and invariant manifolds in Hilbert space*, *Ann. of Math.*, 115 (1982), pp. 243–290.
- [26] G. D. RUXTON, *Density-dependent migration and stability in a system of linked populations*, *Bull. Math. Biol.*, 58 (1996), pp. 643–660.
- [27] J. A. L. SILVA, M. CASTRO, AND D. JUSTO, *Stability in a metapopulation model with density-dependent dispersal*, *Bull. Math. Biol.*, 63 (2001), pp. 485–506.
- [28] J. A. L. SILVA, M. CASTRO, AND D. JUSTO, *Synchronism in a metapopulation model*, *Bull. Math. Biol.*, 62 (2000), pp. 337–349.
- [29] D. TILMAN AND P. KAREIVA, *Spatial Ecology: The Role of Space in Population Dynamics and Interspecific Interactions*, Princeton University Press, Princeton, NJ, 1997.
- [30] H. THUNBERG, *Periodicity versus chaos in one-dimensional dynamics*, *SIAM Rev.*, 43 (2001), pp. 3–30.
- [31] M. VIANA, *Lecture Notes on Attractors and Physical Measures*, A paper from the 12th Escuela Latinoamericana de Matemáticas, Monografías del IMCA 8, IMCA, Lima, 1999.
- [32] L. S. YOUNG, *What are SRB measures and which dynamical systems have them?*, *J. Statist. Phys.*, 108 (2002), pp. 733–754.

STABILITY OF SINGLE-WAVE-FORM SOLUTIONS IN THE UNDERDAMPED FRENKEL–KONTOROVA MODEL*

WEN-XIN QIN[†], CHUN-LAN XU[†], AND XIN MA[‡]

Abstract. Via the monotonicity approach, we prove that the single-wave-form solution for the underdamped Frenkel–Kontorova model with dc-driving and periodic boundary conditions is globally stable, provided the driving force is large enough.

Key words. single-wave-form solution, Frenkel–Kontorova model, monotonicity

AMS subject classifications. 34C12, 34C15, 34C60, 34D45, 37C65

DOI. 10.1137/070699950

1. Introduction. The dynamics of many physical systems, including charge density waves and Josephson junctions, can be described by the standard Frenkel–Kontorova (F-K) model [5, 6], which can also be seen as a one-dimensional lattice of identical pendula oscillating in parallel planes, coupled to the nearest neighbors. The equation of the F-K model with dc-driving is

$$(1.1) \quad \ddot{u}_j + \Gamma \dot{u}_j + \sin u_j = K(u_{j-1} - 2u_j + u_{j+1}) + F,$$

where $\Gamma > 0$ denotes the damping effect, $K > 0$ measures the coupling strength, and $F \geq 0$ is the constant driving force.

In this paper we are concerned with the single-wave-form solutions of the F-K model (1.1) with periodic boundary conditions

$$(1.2) \quad u_{j+N}(t) = u_j(t) + 2\pi M,$$

where M and N are positive integers. A solution $\{u_j(t)\}$ of system (1.1)–(1.2) is called a single-wave-form solution if there exist a constant $T > 0$ and a waveform function $f: \mathbb{R} \rightarrow \mathbb{R}$ satisfying

$$(1.3) \quad f(t + T) = f(t) + 2\pi, \quad t \in \mathbb{R},$$

such that

$$u_j(t) = f(t + jTM/N).$$

The positive number T is said to be the period of the waveform function f .

A running periodic solution is one for which there is a minimal $T > 0$ such that $u_j(t + T) = u_j(t) + 2\pi$ for $t \in \mathbb{R}$ and $j = 1, \dots, N$. Note that a single-wave-form solution is a special form of running periodic solution with the same waveform and

*Received by the editors August 11, 2007; accepted for publication (in revised form) May 7, 2008; published electronically September 17, 2008. This work was supported by the National Natural Science Foundation of China (grants 10771155, 10571131) and the Natural Science Foundation of Jiangsu Province (grant BK 2006046).

<http://www.siam.org/journals/sima/40-3/69995.html>

[†]Department of Mathematics, Suzhou University, Suzhou 215006, People's Republic of China (wxqin@public1.sz.js.cn, xuchunlan50@163.com).

[‡]Golden Audit College, Nanjing Audit University, Nanjing 210029, People's Republic of China (amethystxin@163.com).

equal phase lags. Single-wave-form solutions are also called ponies on a merry-go-round [1], discrete rotating waves [2], or splay-phase solutions [11] in the literature. Although single-wave-form solutions in the coupled oscillators systems play an important role, it is not easy, as remarked in [16], to give mathematically rigorous analysis of their existence and stability. For the small coupling case, Levi proved by the Brouwer fixed point theorem the existence result in [9]. By reformulating the existence problem as a fixed point problem in a Banach space and applying the Schauder fixed point theorem, Katriel [7] presented a rigorous and detailed analysis on the existence of single-wave-form solutions of system (1.1)–(1.2). The method used in [7] is close in spirit to that in [12], which studied the existence of single-wave-form solutions for the Josephson junction systems. We remark that the stability question was not touched upon in [7] and [12]; see also the open problems in [16]. For the small coupling case, the local stability of the single-wave-form solutions, as remarked by Levi in [9], can be proved by applying the approach used in [8]. Meanwhile, by computing the Floquet multipliers and applying perturbation analysis, Watanabe and Swift investigated the local stability of single-wave-form solutions in series arrays of Josephson junctions [15]. Recently, Baesens and MacKay found the monotonicity for the F-K model under the overdamped condition $\Gamma > 2\sqrt{2K+1}$. It can be inferred from the results of [3, 4] that the single-wave-form solution of the overdamped F-K model (1.1)–(1.2) is globally stable if there is no equilibrium. Note that from [3] we know that there exists $F_d \geq 0$, which depends on Γ and K , such that the F-K model has equilibria for $0 \leq F \leq F_d$ and has no equilibria for $F > F_d$.

The goal of this paper is to study the global stability of the single-wave-form solutions for the underdamped F-K model, the existence of which was demonstrated by Katriel [7]. By global stability we mean that given any initial value $\{u_j(0)\}$ satisfying the periodic boundary condition $u_{j+N}(0) = u_j(0) + 2\pi M$, the solution $\{u_j(t)\}$ approaches the single-wave-form solution $\{\hat{u}_j(t)\}$ as $t \rightarrow +\infty$. Moreover, we will show that $\{\hat{u}_j(t)\}$ attracts each orbit with phase; i.e., for a solution $\{u_j(t)\}$, there exists a $\tau \in \mathbb{R}$, such that $\lim_{t \rightarrow +\infty} |u_j(t) - \hat{u}_j(t + \tau)| = 0$.

THEOREM 1.1. *Assume that $\Gamma > 2\sqrt{2K}$. Then there exists $F_0 > 1$, such that the F-K model (1.1)–(1.2) with $F \geq F_0$ admits a globally stable single-wave-form solution, which attracts each orbit with phase.*

The main idea of the proof is as follows. Inspired by the work of [3, 4, 13], we know that strong monotonicity generally implies stability of running periodic solutions. Therefore, as the first step, we show that the Poincaré map P^{mT} of system (1.1)–(1.2), where T is the period of the waveform function and m is some positive integer, is strongly monotone if $\Gamma > 2\sqrt{2K}$ and the driving force F is large enough. Then we prove via the strong monotonicity that the single-wave-form solution is globally stable.

We should remark that under the overdamped condition, i.e., $\Gamma > 2\sqrt{2K+1}$, the strong monotonicity found by Baesens and MacKay together with the approach presented for finite chains in [3] leads immediately to the stability of single-wave-form solutions. Here we extend the parameter range to the underdamped case, where the damping coefficient Γ could be small if the coupling strength K is small.

2. Trapping region. Let $u = (u_1, \dots, u_N)^T \in \mathbb{R}^N$. Then system (1.1)–(1.2) can be written as

$$(2.1) \quad \ddot{u} + \Gamma \dot{u} + S(u) = -KAu + E,$$

where

$$A = \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 & -1 \\ -1 & 2 & -1 & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ -1 & 0 & \cdots & 0 & -1 & 2 \end{pmatrix}$$

is an $N \times N$ matrix, $S(u) = (\sin u_1, \dots, \sin u_N)^T$, $E = (-2\pi KM + F, F, F, \dots, F, F, 2\pi KM + F)^T$. The eigenvalues of matrix A are

$$\lambda_j = 4 \sin^2 \frac{j\pi}{N}, \quad j = 0, 1, \dots, N - 1.$$

Let

$$\hat{e}_j = \left(1, \cos \frac{2j\pi}{N}, \dots, \cos \frac{(N-1)2j\pi}{N} \right)^T, \quad \tilde{e}_j = \left(0, \sin \frac{2j\pi}{N}, \dots, \sin \frac{(N-1)2j\pi}{N} \right)^T,$$

$j = 1, \dots, [\frac{N}{2}]$. Then \hat{e}_j and \tilde{e}_j are eigenvectors corresponding to $\lambda_j = \lambda_{N-j}$. Note that if N is even, then $\lambda_{N/2}$ is a simple eigenvalue with eigenvector $\hat{e}_{N/2}$. Let $e_j = \hat{e}_j/|\hat{e}_j|$ and $e_{N-j} = \tilde{e}_j/|\tilde{e}_j|$ for $j = 1, \dots, [\frac{N}{2}]$. Then e_j is the normalized eigenvector corresponding to λ_j , $j = 1, \dots, N - 1$. The normalized eigenvector corresponding to $\lambda_0 = 0$ is $e_0 = (1/\sqrt{N}, \dots, 1/\sqrt{N})^T$. Denoting $\mathbf{u} = (u, \dot{u})^T$, we have

$$(2.2) \quad \dot{\mathbf{u}} = C\mathbf{u} + \mathbf{S}(\mathbf{u}) + \mathbf{E},$$

where

$$C = \begin{pmatrix} 0 & I \\ -KA & -\Gamma I \end{pmatrix}, \quad \mathbf{S}(\mathbf{u}) = \begin{pmatrix} 0 \\ -S(u) \end{pmatrix}, \quad \mathbf{E} = \begin{pmatrix} 0 \\ E \end{pmatrix}.$$

The eigenvalues of matrix C are

$$\mu_j^\pm = \frac{-\Gamma \pm \sqrt{\Gamma^2 - 4K\lambda_j}}{2},$$

with the corresponding eigenvector being $(e_j, \mu_j^\pm e_j)^T$, $j = 0, \dots, N - 1$. In particular, $\mu_0^+ = 0$, $\mu_0^- = -\Gamma$.

The conclusions of the following lemma have been proved by Katriel; see Theorem 2 in [7].

LEMMA 2.1. *System (1.1)–(1.2) admits a single-wave-form solution, provided $F > 1$. Moreover, the period T of the waveform function satisfies*

$$2\pi\Gamma/(F + 1) \leq T \leq 2\pi\Gamma/(F - 1).$$

The method of proof involves reformulating the problem as a fixed point problem in a Banach space and applying the Schauder fixed point theorem. The uniqueness of the single-wave-form solution needs some further assumptions; see the last section in [7].

A closed subset Ω of the phase space \mathbb{R}^{2N} is called a trapping region of (2.2) if it is positively invariant for the flow generated by (2.2), and for each point $\mathbf{u}_0 \in \mathbb{R}^{2N}$,

there exists t_0 such that the orbit $\{\mathbf{u}(t)\}$ with $\mathbf{u}(0) = \mathbf{u}_0$ enters Ω and remains there for $t > t_0$.

LEMMA 2.2. *System (2.2) has a trapping region Ω . Moreover, there exist constants $b > 0$, independent of F , v_- , and v^+ , depending upon the driving force F , such that*

$$|u_{j-1} - 2u_j + u_{j+1}| \leq b, \quad v_- \leq \dot{u}_j \leq v^+, \quad j = 1, \dots, N, \quad u_0 = u_N, \quad u_1 = u_{N+1},$$

for $\mathbf{u} = (u_1, \dots, u_n, \dot{u}_1, \dots, \dot{u}_N) \in \Omega$.

Proof. There exists an orthonormal matrix D , the j th column being e_{j-1} , such that $D^{-1}AD = \text{diag}(\lambda_0, \lambda_1, \dots, \lambda_{N-1}) \triangleq \Lambda$. Note that $D^{-1} = D^T$, the transpose of D since D is an orthonormal matrix. Making a transformation $x = D^{-1}u$, where $x = (x_1, x_2, \dots, x_N)^T$, we convert system (2.1) into

$$(2.3) \quad \ddot{x} + \Gamma \dot{x} + D^{-1}S(Dx) = -K\Lambda x + D^{-1}E.$$

Noting that the j th element of $D^{-1}S(Dx)$ is $\langle e_{j-1}, S(Dx) \rangle$, denoted by $g_j(x, t)$, we have the fact that $|g_j(x, t)| \leq \sqrt{N}$. Meanwhile, it follows that

$$D^{-1}E = D^{-1} \begin{pmatrix} F \\ F \\ \vdots \\ F \\ F \end{pmatrix} + D^{-1} \begin{pmatrix} -2\pi KM \\ 0 \\ \vdots \\ 0 \\ 2\pi KM \end{pmatrix} \triangleq \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_{N-1} \\ f_N \end{pmatrix},$$

and hence

$$f_1 = \sqrt{N}F, \quad |f_j| \leq 2\sqrt{2}\pi KM, \quad j = 2, \dots, N.$$

Consequently, we derive that in an x -coordinate system,

$$(2.4) \quad \ddot{x}_j + \Gamma \dot{x}_j + K\lambda_{j-1}x_j = -g_j(x, t) + f_j.$$

There are three cases to discuss.

(1) $\Gamma^2 > 4K\lambda_{j-1}$. In this case the matrix

$$C_j = \begin{pmatrix} 0 & 1 \\ -K\lambda_{j-1} & -\Gamma \end{pmatrix}$$

has two real eigenvalues, $\mu_{j-1}^\pm = (-\Gamma \pm \sqrt{\Gamma^2 - 4K\lambda_{j-1}})/2$. It then follows that there exists a real matrix

$$T_j = \begin{pmatrix} 1 & 1 \\ \mu_{j-1}^+ & \mu_{j-1}^- \end{pmatrix}, \quad \text{such that} \quad T_j^{-1}C_jT_j = \begin{pmatrix} \mu_{j-1}^+ & 0 \\ 0 & \mu_{j-1}^- \end{pmatrix}.$$

Let

$$\begin{pmatrix} y_j \\ z_j \end{pmatrix} = T_j^{-1} \begin{pmatrix} x_j \\ \dot{x}_j \end{pmatrix}.$$

Then from (2.4) we have

$$\begin{pmatrix} \dot{y}_j \\ \dot{z}_j \end{pmatrix} = \begin{pmatrix} \mu_{j-1}^+ & 0 \\ 0 & \mu_{j-1}^- \end{pmatrix} \begin{pmatrix} y_j \\ z_j \end{pmatrix} + T_j^{-1} \begin{pmatrix} 0 \\ -g_j(x, t) + f_j \end{pmatrix}.$$

Denoting

$$Y_j = \begin{pmatrix} y_j \\ z_j \end{pmatrix}, \quad \Lambda_j = \begin{pmatrix} \mu_{j-1}^+ & 0 \\ 0 & \mu_{j-1}^- \end{pmatrix}, \quad \text{and} \quad h_j(x, t) = T_j^{-1} \begin{pmatrix} 0 \\ -g_j(x, t) + f_j \end{pmatrix}$$

yields $\dot{Y}_j = \Lambda_j Y_j + h_j(x, t)$, and hence

$$Y_j(t) = e^{\Lambda_j t} Y_j(0) + \int_0^t e^{\Lambda_j(t-\tau)} h_j(x(\tau), \tau) d\tau.$$

The boundedness of $g_j(x, t)$ and f_j implies that there exists $d_j > 0$, such that $\sup_t |h_j(x(t), t)| \leq d_j$, and

$$|Y_j(t)| \leq e^{(\mu_{j-1}^+)^t} |Y_j(0)| + \frac{d_j}{|\mu_{j-1}^+|} \left(1 - e^{(\mu_{j-1}^+)^t} \right), \quad t \geq 0.$$

Note that d_j ($j = 2, \dots, N$) are independent of F . Let $b_j = (d_j + 1)/|\mu_{j-1}^+|$. Then it follows that $|Y_j(t)| \leq b_j$ for all $t \geq 0$ if $|Y_j(0)| \leq b_j$. Moreover, for each solution of (2.2), there exists $t_0 > 0$ such that $|Y_j(t)| \leq b_j$ for $t > t_0$.

In particular, when $j = 1$, we have the following estimates:

$$T_1 = \begin{pmatrix} 1 & 1 \\ 0 & \mu_0^- \end{pmatrix} \quad \text{and} \quad T_1^{-1} = \begin{pmatrix} 1 & -1/\mu_0^- \\ 0 & 1/\mu_0^- \end{pmatrix}.$$

As a consequence, it follows that

$$\dot{z}_1 = \mu_0^- z_1 - \frac{g_1(x, t)}{\mu_0^-} + \frac{f_1}{\mu_0^-} = -\Gamma z_1 + \frac{1}{\Gamma} (g_1(x, t) - f_1).$$

Let $b_1^- = -\Gamma^{-2} \sqrt{N} (F + 1)$ and $b_1^+ = -\Gamma^{-2} \sqrt{N} (F - 1)$. Then we have $z_1(t) \in [b_1^-, b_1^+]$ for all $t \geq 0$ if $z_1(0) \in [b_1^-, b_1^+]$. Moreover if $z_1(0) \notin [b_1^-, b_1^+]$, then there exists $t_0 > 0$ such that $z_1(t) \in [b_1^-, b_1^+]$ for all $t > t_0$.

(2) $\Gamma^2 = 4K\lambda_{j-1}$. In this case, taking

$$C_j = \begin{pmatrix} 0 & 1 \\ -K\lambda_{j-1} & -\Gamma \end{pmatrix} \quad \text{and} \quad T_j = \begin{pmatrix} 1 & 1 \\ -\Gamma/4 & -3\Gamma/4 \end{pmatrix},$$

we have

$$\begin{aligned} T_j^{-1} C_j T_j &= T_j^{-1} \begin{pmatrix} 0 & 1 \\ -\frac{3}{4} K \lambda_{j-1} & -\Gamma \end{pmatrix} T_j + T_j^{-1} \begin{pmatrix} 0 & 0 \\ -\frac{1}{4} K \lambda_{j-1} & 0 \end{pmatrix} T_j \\ &= \begin{pmatrix} -\frac{\Gamma}{4} & 0 \\ 0 & -\frac{3}{4} \Gamma \end{pmatrix} - \frac{\Gamma}{8} \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix}. \end{aligned}$$

For each $z = (z_1, z_2) \in \mathbb{R}^2$, it follows that

$$\langle T_j^{-1}C_jT_jz, z \rangle = -\frac{\Gamma}{4}z_1^2 - \frac{3}{4}\Gamma z_2^2 - \frac{\Gamma}{8}(z_1^2 - z_2^2) \leq -\frac{\Gamma}{4}(z_1^2 + z_2^2) = -\frac{\Gamma}{4}|z|^2.$$

Let

$$\begin{pmatrix} y_j \\ z_j \end{pmatrix} = T_j^{-1} \begin{pmatrix} x_j \\ \dot{x}_j \end{pmatrix}.$$

Then we have $\dot{Y}_j = \Lambda_j Y_j + h_j(x, t)$, where $\Lambda_j = T_j^{-1}C_jT_j$. Since $\langle \Lambda_j z, z \rangle \leq -\frac{\Gamma}{4}|z|^2$, then $\|e^{\Lambda_j t}\| \leq e^{-\frac{\Gamma}{4}t}$ for $t \geq 0$, and hence

$$\begin{aligned} |Y_j(t)| &\leq e^{-\frac{\Gamma}{4}t}|Y_j(0)| + \int_0^t e^{-\frac{\Gamma}{4}(t-\tau)}|h_j(x(\tau), \tau)|d\tau \\ &\leq e^{-\frac{\Gamma}{4}t}|Y_j(0)| + \frac{4d_j}{\Gamma} \left(1 - e^{-\frac{\Gamma}{4}t}\right), \quad t \geq 0. \end{aligned}$$

Let $b_j = (4d_j + 1)/\Gamma$. Then we obtain the same conclusion as in case (1).

(3) $\Gamma^2 < 4K\lambda_{j-1}$. In this case, the eigenvalues of the matrix C_j are a pair of conjugate complex numbers $\mu_{j-1}^\pm = \alpha_j \pm i\beta_j = (-\Gamma \pm \sqrt{\Gamma^2 - 4K\lambda_{j-1}})/2$, where $\alpha_j = -\Gamma/2 < 0$. Taking

$$T_j = \begin{pmatrix} 1 & 0 \\ \alpha_j & \beta_j \end{pmatrix} \text{ yields } T_j^{-1}C_jT_j = \begin{pmatrix} \alpha_j & \beta_j \\ -\beta_j & \alpha_j \end{pmatrix} \triangleq \Lambda_j.$$

Making a transformation

$$\begin{pmatrix} y_j \\ z_j \end{pmatrix} = T_j^{-1} \begin{pmatrix} x_j \\ \dot{x}_j \end{pmatrix},$$

we obtain from (2.4) that $\dot{Y}_j = \Lambda_j Y_j + h_j(x, t)$. From the fact $\|e^{\Lambda_j t}\| \leq e^{\alpha_j t}$ for $t \geq 0$, it follows that

$$|Y_j(t)| \leq e^{\alpha_j t}|Y_j(0)| + \frac{d_j}{|\alpha_j|}(1 - e^{\alpha_j t}).$$

Taking $b_j = (d_j + 1)/|\alpha_j|$ leads to the same conclusion as in case (1).

Consequently, the set

$$\Omega' = \{(y_1, z_1, \dots, y_N, z_N) \mid b_1^- \leq z_1 \leq b_1^+, |(y_j, z_j)| \leq b_j, j = 2, \dots, N\}$$

is a trapping region. Let Ω denote this trapping region in a \mathbf{u} -coordinate system. For each $\mathbf{u} = (u, \dot{u})^T = (u_1, \dots, u_N, \dot{u}_1, \dots, \dot{u}_N)^T \in \Omega$, let $x = (x_1, \dots, x_N)^T = D^{-1}u$ and $\hat{x} = (0, x_2, \dots, x_N)^T$. Since $\mathbf{u} \in \Omega$, then the corresponding $(y_1, z_1, \dots, y_N, z_N)$ satisfies $|(y_j, z_j)| \leq b_j$ for $j = 2, \dots, N$, and

$$(2.5) \quad |(x_j, \dot{x}_j)| \leq \|T_j\| |(y_j, z_j)| \leq b_j \|T_j\| \triangleq \hat{b}_j, \quad j = 2, \dots, N,$$

implying that $|x_j| \leq \hat{b}_j$, and $|\hat{x}| \leq \sum_{j=2}^N \hat{b}_j$. From $ADx = AD\hat{x}$ we deduce that

$$|u_{j-1} - 2u_j + u_{j+1}| = |Au| = |ADx| = |AD\hat{x}| \leq \|A\| \|D\| |\hat{x}| \leq 4|\hat{x}| \leq 4 \sum_{j=2}^N \hat{b}_j \triangleq b.$$

Note that b_j ($j = 2, \dots, N$) are independent of F . So is b . From the facts $\dot{x}_1 = \mu_0^- z_1 = -\Gamma z_1$ and $z_1 \in [b_1^-, b_1^+]$ we derive that

$$\dot{x}_1 \in \left[\Gamma^{-1}\sqrt{N}(F - 1), \Gamma^{-1}\sqrt{N}(F + 1) \right].$$

Denoting $D = (D_{jk})$, we have from $\dot{u} = D\dot{x}$ that

$$\dot{u}_j = \frac{1}{\sqrt{N}}\dot{x}_1 + D_{j2}\dot{x}_2 + \dots + D_{jN}\dot{x}_N,$$

implying by (2.5) that

$$\frac{F - 1}{\Gamma} - \sum_{k=2}^N |D_{jk}| \hat{b}_k \leq \dot{u}_j \leq \frac{F + 1}{\Gamma} + \sum_{k=2}^N |D_{jk}| \hat{b}_k.$$

Let $b^* = \sqrt{\sum_{k=2}^N \hat{b}_k^2}$. Then b^* is independent of F , and $\sum_{k=2}^N |D_{jk}| \hat{b}_k < b^*$ by the Cauchy inequality. Taking

$$(2.6) \quad v_- = \frac{F - 1}{\Gamma} - b^* \quad \text{and} \quad v^+ = \frac{F + 1}{\Gamma} + b^*$$

yields $v_- \leq \dot{u}_j \leq v^+$, $j = 1, \dots, N$. \square

Let $\hat{\mathbf{u}}(t) = (\hat{u}(t), \dot{\hat{u}}(t))^T = (\hat{u}_1(t), \dots, \hat{u}_N(t), \dot{\hat{u}}_1(t), \dots, \dot{\hat{u}}_N(t))^T$ denote the single-wave-form solution obtained in Lemma 2.1.

LEMMA 2.3. *The single-wave-form solution $\hat{\mathbf{u}}(t) \in \Omega$ for $t \in \mathbb{R}$.*

Proof. Let $\hat{\mathbf{x}}(t) = (\hat{x}_1(t), \dots, \hat{x}_N(t), \dot{\hat{x}}_1(t), \dots, \dot{\hat{x}}_N(t))^T = (D^{-1}\hat{u}(t), D^{-1}\dot{\hat{u}}(t))^T$. Then it is easy to check that $\hat{x}_1(t + T) = \hat{x}_1(t) + 2\pi\sqrt{N}$ for $t \in \mathbb{R}$, and the other components are periodic with period $T > 0$, where T is the period of the waveform function corresponding to $\hat{\mathbf{u}}(t)$. From the proof of Lemma 2.2 we can denote this single-wave-form solution by $(\hat{y}_1(t), \hat{z}_1(t), \dots, \hat{y}_N(t), \hat{z}_N(t))$. It follows that all the components, except $\hat{y}_1(t)$, are periodic, implying that the single-wave-form solution belongs to the trapping region Ω' for all $t \in \mathbb{R}$. \square

In the remainder of this paper, we assume that $F \geq F_1 = \Gamma b^* + 1$ such that $v_- > 0$. Let $T_0 = 2\pi/b^*$. Then from Lemma 2.1 it follows that $0 < T \leq T_0$ if $F \geq F_1$. Let

$$(2.7) \quad c = \frac{2}{v_-} + \frac{v^+ T_0 (Kb + F + 1 + \Gamma v^+ + 2\pi KM)}{v_-^3}.$$

Note that c is a continuous function of F and $c \rightarrow 0$ as $F \rightarrow +\infty$.

LEMMA 2.4. *Assume that $\mathbf{u}(t) = (u_1(t), \dots, u_N(t), \dot{u}_1(t), \dots, \dot{u}_N(t))$ is a solution of (2.2) with $\mathbf{u}(t) \in \Omega$ for $t \geq 0$. Then*

$$\left| \int_{t_0}^{t_1} \cos u_j(t) dt \right| \leq c,$$

where $0 \leq t_0 < t_1$ and $t_1 - t_0 \leq T_0$, $j = 1, \dots, N$.

Proof. From Lemma 2.2 it follows that for $t \geq 0$,

$$\begin{aligned} \left| \frac{\ddot{u}_j}{\dot{u}_j} \right| &= \left| \frac{K(u_{j-1} + u_{j+1} - 2u_j) - \sin u_j + F - \Gamma \dot{u}_j}{\dot{u}_j} \right| \\ &\leq \frac{Kb + F + 1 + \Gamma v^+ + 2\pi KM}{v_-}, \end{aligned}$$

and hence

$$\begin{aligned} & \left| \int_{t_0}^{t_1} \cos u_j(t) dt \right| = \left| \int_{t_0}^{t_1} \frac{\cos u_j}{\dot{u}_j} du_j \right| = \left| \int_{t_0}^{t_1} \frac{1}{\dot{u}_j} d(\sin u_j) \right| \\ & \leq \left| \frac{\sin u_j}{\dot{u}_j} \Big|_{t_0}^{t_1} + \int_{t_0}^{t_1} \frac{\sin u_j}{(\dot{u}_j)^2} d\dot{u}_j \right| \leq \frac{2}{v_-} + \left| \int_{t_0}^{t_1} \frac{\sin u_j}{(\dot{u}_j)^2} \ddot{u}_j du_j \right| \\ & \leq \frac{2}{v_-} + \frac{1}{v_-^2} \left(\frac{Kb + F + 1 + \Gamma v^+ + 2\pi KM}{v_-} \right) \left| \int_{t_0}^{t_1} du_j \right| \\ & \leq \frac{2}{v_-} + \frac{v^+ T_0 (Kb + F + 1 + \Gamma v^+ + 2\pi KM)}{v_-^3} = c. \quad \square \end{aligned}$$

3. Strong monotonicity. We know from [3] that there exists $F_d \geq 0$, which depends on Γ and K , such that the F-K model has equilibria for $0 \leq F \leq F_d$ and has no equilibria for $F > F_d$. Meanwhile, strong monotonicity generally implies stability of running periodic solutions if there is no equilibrium; see [3] and [13]. Therefore, we investigate in this section the strong monotonicity of system (1.1)–(1.2). To proceed, we make a transformation of variables:

$$(3.1) \quad \xi_j = u_j + \Gamma^{-1} \dot{u}_j, \quad \eta_j = \Gamma^{-1} \dot{u}_j, \quad j = 1, 2, \dots, N,$$

i.e., $\Xi = I_\Gamma \mathbf{u}$, in which $\Xi = (\xi, \eta)^T$, $\xi = (\xi_1, \dots, \xi_N)^T$, $\eta = (\eta_1, \dots, \eta_N)^T$,

$$I_\Gamma = \begin{pmatrix} I & \Gamma^{-1}I \\ 0 & \Gamma^{-1}I \end{pmatrix},$$

and I is the identity matrix of order N . Now the system (1.1)–(1.2) becomes

$$(3.2) \quad \left\{ \begin{aligned} \dot{\xi}_1 &= -\Gamma^{-1} \sin(\xi_1 - \eta_1) + \Gamma^{-1}K[\xi_2 - \eta_2 + \xi_N - \eta_N - 2(\xi_1 - \eta_1)] \\ &\quad - 2\pi K\Gamma^{-1}M + \Gamma^{-1}F, \\ \dot{\xi}_2 &= -\Gamma^{-1} \sin(\xi_2 - \eta_2) + \Gamma^{-1}K[\xi_1 - \eta_1 + \xi_3 - \eta_3 - 2(\xi_2 - \eta_2)] \\ &\quad + \Gamma^{-1}F, \\ &\quad \vdots \\ \dot{\xi}_{N-1} &= -\Gamma^{-1} \sin(\xi_{N-1} - \eta_{N-1}) \\ &\quad + \Gamma^{-1}K[\xi_{N-2} - \eta_{N-2} + \xi_N - \eta_N - 2(\xi_{N-1} - \eta_{N-1})] + \Gamma^{-1}F, \\ \dot{\xi}_N &= -\Gamma^{-1} \sin(\xi_N - \eta_N) + \Gamma^{-1}K[\xi_1 - \eta_1 + \xi_{N-1} - \eta_{N-1} - 2(\xi_N - \eta_N)] \\ &\quad + 2\pi K\Gamma^{-1}M + \Gamma^{-1}F, \\ \dot{\eta}_1 &= -\Gamma\eta_1 + \dot{\xi}_1, \\ &\quad \vdots \\ \dot{\eta}_N &= -\Gamma\eta_N + \dot{\xi}_N. \end{aligned} \right.$$

Set $\Phi = (\phi_1, \dots, \phi_N)^T$ and $\Psi = (\psi_1, \dots, \psi_N)^T$. The linearized equations along a solution $\Xi(t) \in I_\Gamma\Omega$ ($t \geq 0$) are

$$(3.3) \quad \begin{cases} \dot{\phi}_j = -\Gamma^{-1}(\cos(\xi_j - \eta_j))(\phi_j - \psi_j) + \Gamma^{-1}K[\phi_{j+1} - \psi_{j+1} \\ \quad + \phi_{j-1} - \psi_{j-1} - 2(\phi_j - \psi_j)] \triangleq L_j(\Phi, \Psi, t), \\ \dot{\psi}_j = -\Gamma\psi_j + \dot{\phi}_j \triangleq H_j(\Phi, \Psi, t), \end{cases}$$

in which $j = 1, 2, \dots, N$, $\phi_0 = \phi_N, \phi_{N+1} = \phi_1, \psi_0 = \psi_N, \psi_{N+1} = \psi_1$.

First we study the properties of the auxiliary equations

$$(3.4) \quad \dot{\phi}_j = L_j(\Phi, \Psi, t) + \varepsilon, \quad \dot{\psi}_j = H_j(\Phi, \Psi, t),$$

in which $\varepsilon > 0$. Let

$$(3.5) \quad \Sigma = \{(\Phi, \Psi) \in \mathbb{R}^{2N} \mid -\phi_j \leq \psi_j \leq \phi_j, \phi_j \geq 0, j = 1, \dots, N\}.$$

Let $\text{int } \Sigma$ and $\partial\Sigma$ denote the interior and boundary of Σ , respectively.

Consider a first order equation

$$(3.6) \quad \dot{w} = \frac{1}{\Gamma} [-\Gamma^2 w - 2K(1 - w)^2 - \cos u_j(t)(1 - w)^2],$$

where $u_j(t)$ is a component of $\mathbf{u}(t)$ and $\mathbf{u}(t) \in \Omega$ ($t \geq 0$) is a solution of (2.2) with external force F . The proof of the following lemma is postponed to the appendix.

LEMMA 3.1. *Assume $\Gamma > 2\sqrt{2K}$. Then there exists $F_0 \geq F_1$, such that each solution $w(t)$ of (3.6), in which $u_j(t)$ is a component of a solution $\mathbf{u}(t) \in \Omega$ ($t \geq 0$) to (2.2) with $F \geq F_0$, satisfies that if $w(0) \in [-1, 1]$, then $1 - \Gamma^2/(4K) < w(t) < 1$ for all $t > 0$ and $-1 < w(t) < 1$ for $t \in [T_0/2, T_0]$.*

We remark that F_0 increases as N and M increase, as can be seen from the value F_1 , which depends on the eigenvalues of matrix A and the proof in the appendix.

LEMMA 3.2. *Assume $\Gamma > 2\sqrt{2K}$ and $F \geq F_0$. Then the solution $(\Phi(\varepsilon, t), \Psi(\varepsilon, t))$ of (3.4), in which $\Xi(t) \in I_\Gamma\Omega$ ($t \geq 0$) is a solution to (3.2), satisfies that $\phi_j(\varepsilon, t) > 0$ and $1 - \Gamma^2/(4K) < \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ for all $t > 0$, and $-1 < \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ for $t \in [T_0/2, T_0]$, $j = 1, \dots, N$, if $(\Phi(\varepsilon, 0), \Psi(\varepsilon, 0)) \in \partial\Sigma \setminus \{0\}$.*

Proof. There are three cases to discuss for $(\Phi(\varepsilon, 0), \Psi(\varepsilon, 0)) \in \partial\Sigma \setminus \{0\}$.

(i) $\phi_j(\varepsilon, 0) > 0$ and $\phi_j(\varepsilon, 0) = \psi_j(\varepsilon, 0)$. Then there is a maximum $t_j^0 > 0$, or $t_j^0 = +\infty$, such that $\phi_j(\varepsilon, t) > 0$ for $t \in [0, t_j^0]$. Let $w_j(\varepsilon, t) = \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t)$ for $t \in [0, t_j^0]$. Then it follows that

$$\begin{aligned} \dot{w}_j &= (\dot{\psi}_j\phi_j - \psi_j\dot{\phi}_j)/\phi_j^2 \\ &= \Gamma^{-1} [-\Gamma^2 w_j - 2K(1 - w_j)^2 - \cos(\xi_j - \eta_j)(1 - w_j)^2] - \varepsilon w_j/\phi_j \\ &\quad + \frac{K}{\Gamma\phi_j}(\phi_{j-1} - \psi_{j-1} + \phi_{j+1} - \psi_{j+1})(1 - w_j), \end{aligned}$$

and hence

$$\dot{w}_j|_{w_j=1, t=0} = -\Gamma - \varepsilon/\phi_j(\varepsilon, 0) < 0,$$

implying that $w_j(\varepsilon, t) < 1$ as long as $\phi_j(\varepsilon, t) > 0$, i.e., $\psi_j(\varepsilon, t) < \phi_j(\varepsilon, t)$ for $t \in (0, t_j^0)$.

(ii) $\phi_j(\varepsilon, 0) = \psi_j(\varepsilon, 0) = 0$. In this case we have that $L_j(\Phi, \Psi, t)|_{(\Phi(\varepsilon, 0), \Psi(\varepsilon, 0))} \geq 0$, i.e., $\dot{\phi}_j|_{t=0} \geq \varepsilon > 0$, and

$$0 \leq \frac{H_j(\Phi, \Psi, t)}{L_j(\Phi, \Psi, t) + \varepsilon} \Big|_{(\Phi(\varepsilon, 0), \Psi(\varepsilon, 0))} = \frac{\Gamma^{-1}K(\phi_{j-1} - \psi_{j-1} + \phi_{j+1} - \psi_{j+1})}{\Gamma^{-1}K(\phi_{j-1} - \psi_{j-1} + \phi_{j+1} - \psi_{j+1}) + \varepsilon} < 1.$$

Consequently, there is a maximum $t_j^0 > 0$, or $t_j^0 = +\infty$, such that $\phi_j(\varepsilon, t) > 0$ and $\psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ for $t \in (0, t_j^0)$.

(iii) $\phi_j(\varepsilon, 0) > 0$ and $\psi_j(\varepsilon, 0) = -\phi_j(\varepsilon, 0)$. Then there exists a maximum $t_j^0 > 0$, or $t_j^0 = +\infty$, such that $\phi_j(\varepsilon, t) > 0$ and $\psi_j(\varepsilon, t) < \phi_j(\varepsilon, t)$ for $t \in (0, t_j^0)$. From the discussions of cases (i) and (ii) we know that

$$\frac{K}{\Gamma\phi_j}(\phi_{j-1} - \psi_{j-1} + \phi_{j+1} - \psi_{j+1})(1 - w_j) \geq 0$$

for $t \in (0, t_0)$, where $t_0 = \min_j\{t_j^0\}$. Furthermore, from the fact $\frac{-\varepsilon w_j}{\phi_j}|_{w_j=-1, t=0} > 0$, Lemma 3.1, and the comparison principle we deduce that $\phi_j(\varepsilon, t) > 0$ and $1 - \Gamma^2/(4K) < w_j(\varepsilon, t) = \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ for $t \in (0, t_0)$. We claim that $t_0 = +\infty$. If this is not true, then we have $\phi_j(\varepsilon, t) \rightarrow 0$ or $w_j(\varepsilon, t) \rightarrow 1$ as $t \rightarrow t_0^-$. In the former case, we have $\psi_j(\varepsilon, t) \rightarrow 0$ as $t \rightarrow t_0^-$, leading to a contradiction of the fact $\lim_{t \rightarrow t_0^-} \dot{\phi}_j(\varepsilon, t) \geq \varepsilon > 0$. The latter case is also impossible. Indeed, if $\phi_j(\varepsilon, t) \geq a > 0$, where a is a positive number, then from part (i) we have that $w_j(\varepsilon, t) < 1 - \delta$ for some small positive number δ since $\dot{w}|_{w_j=1, t=0} < 0$. Consequently, we have $t_0 = +\infty$. Another direct consequence from Lemma 3.1 and the comparison principle is that $-1 < \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ for $t \in [T_0/2, T_0]$. Note that $\frac{-\varepsilon w_j}{\phi_j}|_{w_j=-1} > 0$ and $\frac{-\varepsilon w_j}{\phi_j}|_{w_j=1} < 0$. Then for $t \in [T_0/2, T_0]$, we have from the comparison principle that $-1 < \psi_j(\varepsilon, t)/\phi_j(\varepsilon, t) < 1$ holds true uniformly with respect to $\varepsilon > 0$. \square

LEMMA 3.3. Assume that $\Gamma > 2\sqrt{2K}$ and $F \geq F_0$. Then the solution $(\Phi(t), \Psi(t))$ of (3.3), the linearized equations of (3.2) along $\Xi(t) \in I_\Gamma\Omega$ ($t \geq 0$), satisfies that $\phi_j(t) \geq 0$ and $[1 - \Gamma^2/(4K)]\phi_j(t) \leq \psi_j(t) \leq \phi_j(t)$ for all $t > 0$, $j = 1, \dots, N$, if $(\Phi(0), \Psi(0)) \in \partial\Sigma \setminus \{0\}$.

Proof. Assume that $(\Phi(\varepsilon, t), \Psi(\varepsilon, t))$ is a solution of the auxiliary equations (3.4) with the initial values $\Phi(\varepsilon, 0) = \Phi(0)$ and $\Psi(\varepsilon, 0) = \Psi(0)$. Then from Lemma 3.2 it follows that $[1 - \Gamma^2/(4K)]\phi_j(\varepsilon, t) < \psi_j(\varepsilon, t) < \phi_j(\varepsilon, t)$ and $\phi_j(\varepsilon, t) > 0$ for all $t > 0$. Keeping t fixed and taking $\varepsilon \rightarrow 0$, we complete the proof. \square

In fact, we have further conclusions on $\phi_j(t)$ and $\psi_j(t)$ under the same conditions of Lemma 3.3.

LEMMA 3.4. Assume that $\Gamma > 2\sqrt{2K}$ and $F \geq F_0$. Then the solution $(\Phi(t), \Psi(t))$ of (3.3), the linearized equations of (3.2) along $\Xi(t) \in I_\Gamma\Omega$ ($t \geq 0$), satisfies that $\phi_j(t) > 0$ for all $t > 0$, and $-\phi_j(t) < \psi_j(t) < \phi_j(t)$ for $t \in [T_0/2, T_0]$, $j = 1, \dots, N$, if $(\Phi(0), \Psi(0)) \in \partial\Sigma \setminus \{0\}$.

Proof. Since $(\Phi(0), \Psi(0)) \in \partial\Sigma \setminus \{0\}$, then we may assume that there is some k such that $\phi_k(0) > 0$. From Lemma 3.3 we know that $[1 - \Gamma^2/(4K)]\phi_k(t) \leq \psi_k(t) \leq \phi_k(t)$, $\psi_{k+1}(t) \leq \phi_{k+1}(t)$, and $\psi_{k-1}(t) \leq \phi_{k-1}(t)$ for $t \geq 0$, and thus we obtain from (3.3)

$$\dot{\phi}_k(t) \geq -\Gamma^{-1}\Gamma^2/(4K)\phi_k(t) - 2\Gamma^{-1}K\Gamma^2/(4K)\phi_k(t) = -\frac{\Gamma(2K+1)}{4K}\phi_k(t),$$

implying from the Gronwall inequality that

$$\phi_k(t) \geq \phi_k(0) e^{-\frac{\Gamma(2K+1)t}{4K}} > 0 \text{ for } t > 0.$$

Let $w_k(t) = \psi_k(t)/\phi_k(t)$. Then

$$\begin{aligned} \dot{w}_k &= \Gamma^{-1} [-\Gamma^2 w_k - 2K(1 - w_k)^2 - \cos(\xi_k - \eta_k)(1 - w_k)^2] \\ &\quad + \frac{K}{\Gamma\phi_k}(\phi_{k-1} - \psi_{k-1} + \phi_{k+1} - \psi_{k+1})(1 - w_k). \end{aligned}$$

From the fact $\dot{w}_k|_{w_k=1} = -\Gamma < 0$, we obtain $w_k(t) < 1$, i.e., $\phi_k(t) - \psi_k(t) > 0$ for $t > 0$. Therefore from (3.3) we have that

$$\dot{\phi}_{k+1} \geq -\frac{\Gamma(2K+1)}{4K}\phi_{k+1} + \Gamma^{-1}K(\phi_k(t) - \psi_k(t))$$

and, again by the Gronwall inequality, that

$$\phi_{k+1}(t) \geq e^{-\frac{\Gamma(2K+1)t}{4K}}\phi_{k+1}(0) + \int_0^t e^{-\frac{\Gamma(2K+1)(t-\tau)}{4K}}\Gamma^{-1}K(\phi_k(\tau) - \psi_k(\tau))d\tau > 0, \quad t > 0.$$

With the same reasoning we have that $\phi_{k+2}(t) > 0, t > 0, \dots$, and $\phi_{k-1}(t) > 0, t > 0, \dots$. Consequently, we derive that $\phi_j(t) > 0$ for $t > 0, j = 1, \dots, N$.

Keeping $t \in [T_0/2, T_0]$ fixed, we know from Lemma 3.2 that there exists $\delta > 0$ such that the solution $(\Phi(\varepsilon, t), \Psi(\varepsilon, t))$ of the auxiliary equations (3.4) with the initial values $\Phi(\varepsilon, 0) = \Phi(0)$ and $\Psi(\varepsilon, 0) = \Psi(0)$ has the property that

$$-1 + \delta \leq \frac{\psi_j(\varepsilon, t)}{\phi_j(\varepsilon, t)} \leq 1 - \delta.$$

We remark by virtue of the proof of Lemma 3.2 that δ is independent of ε . Taking $\varepsilon \rightarrow 0$, we derive that $-1 + \delta \leq \psi_j(t)/\phi_j(t) \leq 1 - \delta$. \square

Now we define a partial order in the phase space \mathbb{R}^{2N} . Let $\Xi = (\xi_1, \dots, \xi_N, \eta_1, \dots, \eta_N)^T$ and $\Xi' = (\xi'_1, \dots, \xi'_N, \eta'_1, \dots, \eta'_N)^T$.

DEFINITION 3.5. We define $\Xi \leq \Xi'$ if $\xi_j \leq \xi'_j$ and $|\eta_j - \eta'_j| \leq |\xi_j - \xi'_j|$ for $j = 1, \dots, N$. We also define $\Xi < \Xi'$ if $\Xi \leq \Xi'$ and $\Xi \neq \Xi'$, and

$$\Xi \ll \Xi' \text{ if } \xi_j < \xi'_j \text{ and } |\eta_j - \eta'_j| < |\xi_j - \xi'_j| \text{ for } j = 1, \dots, N.$$

We remark that $\Xi \leq \Xi'$ if and only if $\Xi' - \Xi \in \Sigma$, and $\Xi \ll \Xi'$ if and only if $\Xi' - \Xi \in \text{int } \Sigma$. Let P^t denote the Poincaré map of system (3.2), i.e., $P^t\Xi = \Xi(t)$, in which $\Xi(t)$ is a solution of (3.2) with initial condition $\Xi(0) = \Xi$.

DEFINITION 3.6. The Poincaré map P^t ($t > 0$) of system (3.2) is said to be strongly monotone in the set \mathcal{B} with $P^t(\mathcal{B}) \subset \mathcal{B}$ if

$$\Xi_0 < \Xi'_0 \Rightarrow P^t\Xi_0 \ll P^t\Xi'_0 \text{ for } t > 0 \text{ and } \Xi_0, \Xi'_0 \in \mathcal{B}.$$

Now we can prove strong monotonicity by the above discussion and defining a homotopy as in [4, 14].

LEMMA 3.7. Assume that $\Gamma > 2\sqrt{2K}$ and $F \geq F_0$. Then the Poincaré map P^t of system (3.2) is strongly monotone in $I_\Gamma\Omega$ for $t \in [T_0/2, T_0]$.

Proof. Assume that $\Xi_0, \Xi'_0 \in I_\Gamma\Omega$, that $\Xi_0 < \Xi'_0$, and that $\Xi(t)$ and $\Xi'(t)$ are two solutions of (3.2) with initial values $\Xi(0) = \Xi_0$ and $\Xi'(0) = \Xi'_0$, respectively. Let

$$\Xi(\lambda, 0) = (1 - \lambda)\Xi(0) + \lambda\Xi'(0),$$

where $\lambda \in [0, 1]$. Assume that $\Xi(\lambda, t)$ is a solution of (3.2) with initial value $\Xi(\lambda, 0)$. Then $\Xi(\lambda, t) \in I_\Gamma\Omega$ for $t \geq 0$ since $I_\Gamma\Omega$ is convex and positively invariant for system (3.2). Set

$$\Upsilon(\lambda, t) = \frac{\partial \Xi(\lambda, t)}{\partial \lambda} = (\phi_1(\lambda, t), \dots, \phi_N(\lambda, t), \psi_1(\lambda, t), \dots, \psi_N(\lambda, t))^T.$$

Then

$$\Upsilon(\lambda, 0) = \frac{\partial \Xi(\lambda, 0)}{\partial \lambda} = \Xi'(0) - \Xi(0) \quad \text{and} \quad \Xi'(t) - \Xi(t) = \int_0^1 \Upsilon(\lambda, t) d\lambda.$$

One can easily verify that $\Upsilon(\lambda, t)$ is a solution of the linearized equations (3.3) with initial value $\Upsilon(\lambda, 0)$. Since $\Xi(0) < \Xi'(0)$, then

$$\phi_j(\lambda, 0) \geq 0 \quad \text{and} \quad |\psi_j(\lambda, 0)| \leq \phi_j(\lambda, 0), \quad j = 1, \dots, N,$$

i.e., $\Upsilon(\lambda, 0) \in \Sigma \setminus \{0\}$. From Lemma 3.4 it follows that $\Upsilon(\lambda, t) \in \text{int } \Sigma$ for $t \in [T_0/2, T_0]$, i.e.,

$$\phi_j(\lambda, t) > 0 \quad \text{and} \quad |\psi_j(\lambda, t)| < \phi_j(\lambda, t), \quad j = 1, \dots, N, \quad t \in [T_0/2, T_0],$$

which implies that

$$\xi'_j(t) - \xi_j(t) = \int_0^1 \phi_j(\lambda, t) d\lambda > \int_0^1 |\psi_j(\lambda, t)| d\lambda \geq \left| \int_0^1 \psi_j(\lambda, t) d\lambda \right| = |\eta'_j(t) - \eta_j(t)|,$$

i.e., $\Xi(t) \ll \Xi'(t)$, $t \in [T_0/2, T_0]$. This completes the proof. \square

4. Stability analysis. In this section we assume that $\Gamma > 2\sqrt{2K}$ and $F \geq F_0$ so that the conclusion of Lemma 3.7 holds true. We know from Lemma 2.1 that system (1.1)–(1.2) admits a single-wave-form solution $(\hat{u}_1(t), \dots, \hat{u}_N(t))$ since $F \geq F_0 > 1$. Let us denote it in the Ξ -coordinate system by $\hat{\Xi}(t) = (\hat{\xi}_1(t), \dots, \hat{\xi}_N(t), \hat{\eta}_1(t), \dots, \hat{\eta}_N(t))$, i.e.,

$$\hat{\xi}_j(t) = \hat{u}_j(t) + \Gamma^{-1} \hat{u}_j(t), \quad \hat{\eta}_j(t) = \Gamma^{-1} \hat{u}_j(t), \quad j = 1, \dots, N.$$

Consequently, $\hat{\Xi}(t)$ has the following property:

$$(4.1) \quad \hat{\Xi}(t + T) = \hat{\Xi}(t) + 2\pi \mathbf{e},$$

where $\mathbf{e} = (\mathbf{1}, 0)^T$, and $\mathbf{1}$ denotes a vector in \mathbb{R}^N with all components equal to 1. From Lemma 2.1 it follows that there exists an integer m such that $mT \in [T_0/2, T_0]$. Let

$$\ell = \{\hat{\Xi}(t) | t \in \mathbb{R}\} \quad \text{and} \quad Q(\Xi) = P^{mT} \Xi - 2m\pi \mathbf{e}.$$

The following properties will be used in the subsequent discussion.

(I) Q is strongly monotone in $I_\Gamma\Omega$, and each point in ℓ is a fixed point of the map Q .

(II) For any $\Xi \in \mathbb{R}^{2N}$, there exist $\Xi_1, \Xi_2 \in \ell$ such that $\Xi_1 \ll \Xi \ll \Xi_2$. Indeed, we can take $\Xi_1 = \hat{\Xi}(-nT)$ and $\Xi_2 = \hat{\Xi}(nT)$ with n large enough.

LEMMA 4.1. $\hat{\Xi}(t_1) \ll \hat{\Xi}(t_2)$ if $t_1 < t_2$.

Proof. First we show that every two points in ℓ are ordered by the partial order “ \ll .” Assume that there are two points $\hat{\Xi}(t_1)$ and $\hat{\Xi}(t_2)$ ($t_1 < t_2$) which are unordered, i.e., $\hat{\Xi}(t_2) - \hat{\Xi}(t_1) \notin \Sigma$. From the above property (II) we know that there exists a t' such that $\hat{\Xi}(t') \ll \hat{\Xi}(t_2)$, i.e., $\hat{\Xi}(t_2) - \hat{\Xi}(t') \in \text{int } \Sigma$. Increasing t' , we deduce from the continuity that there exists t_0 such that $\hat{\Xi}(t_2) - \hat{\Xi}(t_0) \in \partial \Sigma$, implying by the strong monotonicity of Q that $Q\hat{\Xi}(t_2) - Q\hat{\Xi}(t_0) \in \text{int } \Sigma$. This is a contradiction since $\hat{\Xi}(t_2)$ and $\hat{\Xi}(t_0)$ are fixed points of Q .

Now we show that $\hat{\Xi}(t_1) \ll \hat{\Xi}(t)$ for all $t > t_1$. The above property (II) implies that there exists t_2 such that $\hat{\Xi}(t_1) \ll \hat{\Xi}(t)$ for $t \geq t_2$. Let

$$t_0 = \sup\{t \geq t_1 \mid \hat{\Xi}(t) - \hat{\Xi}(t_1) \in \partial \Sigma\}.$$

It is evident that $t_0 = t_1$. Indeed, if $t_0 > t_1$, then $\hat{\Xi}(t_0) - \hat{\Xi}(t_1) \in \partial \Sigma$, which is impossible by the strong monotonicity of Q . Therefore, $\hat{\Xi}(t) - \hat{\Xi}(t_1) \in \text{int } \Sigma$, i.e., $\hat{\Xi}(t_1) \ll \hat{\Xi}(t)$ for $t > t_1$. \square

Proof of Theorem 1.1. Let

$$[\Xi_1, \Xi_2] = \{\Xi \mid \Xi_1 \leq \Xi \leq \Xi_2\} \quad \text{for } \Xi_1 \leq \Xi_2.$$

For each $\Xi_0 \in I_\Gamma \Omega$, by the above property (II), there always exist Ξ_1 and $\Xi_2 \in \ell$ such that $\Xi_1 \ll \Xi_2$ and $\Xi_0 \in [\Xi_1, \Xi_2]$. The strong monotonicity of Q leads to the conclusion that $\{Q^n \Xi_0\}_{n=1}^{+\infty} \subset [\Xi_1, \Xi_2]$ since $Q(\Xi_1) = \Xi_1$ and $Q(\Xi_2) = \Xi_2$. From the boundedness of the set $[\Xi_1, \Xi_2]$, it follows that the ω -limit set $\omega(\Xi_0)$ of $\{Q^n \Xi_0\}_{n=1}^{+\infty}$ is nonempty. Define for $n \geq 0$

$$\tau_-(n; \Xi_0) = \sup\{\tau \in \mathbb{R} \mid \hat{\Xi}(\tau) \leq Q^n(\Xi_0)\},$$

$$\tau_+(n; \Xi_0) = \inf\{\tau \in \mathbb{R} \mid Q^n(\Xi_0) \leq \hat{\Xi}(\tau)\}.$$

From Lemma 4.1 and the strong monotonicity of Q we know that $\tau_-(n; \Xi_0)$ is strictly increasing, $\tau_+(n; \Xi_0)$ is strictly decreasing with respect to n , and $\tau_-(n; \Xi_0) \leq \tau_+(n; \Xi_0)$. Moreover, $\tau_\pm(n; \Xi_0)$ are continuous with respect to Ξ_0 . If $\tau_-(0; \Xi_0) = \tau_+(0; \Xi_0) = \tau$, then $\Xi_0 = \hat{\Xi}(\tau)$ and $\Xi_0(t) = \hat{\Xi}(t + \tau)$, $t \in \mathbb{R}$, where $\Xi_0(t)$ is a solution of (3.2) with initial value Ξ_0 at $t = 0$. If $\tau_-(0; \Xi_0) < \tau_+(0; \Xi_0)$, let

$$\tau_-^\infty = \sup_{n \geq 0} \{\tau_-(n; \Xi_0)\}, \quad \tau_+^\infty = \inf_{n \geq 0} \{\tau_+(n; \Xi_0)\}.$$

Then $\tau_-^\infty \leq \tau_+^\infty$. If $\tau_-^\infty < \tau_+^\infty$, denote a limit point of $\{Q^n \Xi_0\}$ by $\tilde{\Xi}_0$. Then $\tau_\pm(0; \tilde{\Xi}_0) = \tau_\pm^\infty$, and

$$\tau_-^\infty = \tau_-(0; \tilde{\Xi}_0) < \tau_-(1; \tilde{\Xi}_0) \leq \tau_+(1; \tilde{\Xi}_0) < \tau_+(0; \tilde{\Xi}_0) = \tau_+^\infty.$$

Take k large enough such that $Q^k(\Xi_0)$ is sufficiently close to $\tilde{\Xi}_0$. Then we have

$$\tau_-^\infty < \tau_-(1; Q^k \Xi_0) \leq \tau_+(1; Q^k \Xi_0) < \tau_+^\infty,$$

i.e.,

$$\tau_-^\infty < \tau_-(k + 1; \Xi_0) \leq \tau_+(k + 1; \Xi_0) < \tau_+^\infty,$$

implying a contradiction. Consequently, $\tau_-^\infty = \tau_+^\infty \triangleq \tau$. We deduce that $|\Xi_0(nmT) - \hat{\Xi}(nmT + \tau)| \rightarrow 0$ as $n \rightarrow +\infty$, leading to the conclusion $|\Xi_0(t) - \hat{\Xi}(t + \tau)| \rightarrow 0$ as $t \rightarrow +\infty$. \square

Appendix. Proof of Lemma 3.1. Since $u_j(t)$ in the right-hand side of (3.6) depends upon the external driving force F , so does the solution $w(t)$ of (3.6). If $w(0) = 1$, then since $\dot{w}|_{w=1} = -\Gamma < 0$, it is easy to check that $w(t) < 1$ for $t > 0$. Now we assume that $w(0) = -1$. First we show by contradiction that $1 - \frac{\Gamma^2}{4K} < w(t) < 1$ for $t \in [0, T_0]$ if F is large enough.

Assume that there exists $\hat{t} \in [0, T_0]$ such that $w(\hat{t}) \leq 1 - \frac{\Gamma^2}{4K} < -1$. Taking

$$t_1 = \inf \{t \mid t \in (0, \hat{t}), w(t) = 1 - \Gamma^2/(4K)\} \quad \text{and} \quad t_0 = \sup\{t \mid t \in [0, t_1), w(t) = -1\}$$

yields that $w(t) \in [1 - \frac{\Gamma^2}{4K}, -1]$ for $t \in [t_0, t_1]$. Note that the function $f(x) = -\Gamma^2 x - 2K(1 - x)^2$ is positive on the interval $[1 - \frac{\Gamma^2}{4K}, -1]$. It follows that

$$\begin{aligned} w(t_1) &= w(t_0) + \int_{t_0}^{t_1} \frac{1}{\Gamma} [-\Gamma^2 w(t) - 2K(1 - w(t))^2 - \cos u_j(t)(1 - w(t))^2] dt \\ &> -1 - \frac{1}{\Gamma} \left| \int_{t_0}^{t_1} \cos u_j(t)(1 - w(t))^2 dt \right|. \end{aligned}$$

Note also that

$$\begin{aligned} \int_{t_0}^{t_1} \cos u_j(t)(1 - w(t))^2 dt &= \int_{t_0}^{t_1} (1 - w(t))^2 d \left(\int_{t_0}^t \cos u_j(\tau) d\tau \right) \\ &= (1 - w(t))^2 \int_{t_0}^t \cos u_j(\tau) d\tau \Big|_{t_0}^{t_1} + 2 \int_{t_0}^{t_1} \left[\int_{t_0}^t \cos u_j(\tau) d\tau \right] (1 - w(t)) \dot{w}(t) dt. \end{aligned}$$

Since $[t_0, t_1] \in [0, T_0]$, we deduce from Lemma 2.4 that

$$(A.1) \quad \left| (1 - w(t))^2 \int_{t_0}^t \cos u_j(\tau) d\tau \Big|_{t_0}^{t_1} \right| \leq c(\Gamma^2/4K)^2.$$

From the fact

$$\max_{x \in [1 - \frac{\Gamma^2}{4K}, -1]} f(x) = \Gamma^2 \left(\frac{\Gamma^2}{8K} - 1 \right),$$

it follows that for $t \in [t_0, t_1]$,

$$\begin{aligned} |\dot{w}| &= \frac{1}{\Gamma} |-\Gamma^2 w - 2K(1 - w)^2 - \cos u_j(1 - w)^2| \\ &\leq \frac{1}{\Gamma} |-\Gamma^2 w - 2K(1 - w)^2| + \frac{1}{\Gamma} |\cos u_j(1 - w)^2| \\ &\leq \Gamma \left(\frac{\Gamma^2}{8K} - 1 \right) + \frac{1}{\Gamma} \left(\frac{\Gamma^2}{4K} \right)^2, \end{aligned}$$

and hence

$$\begin{aligned}
 & \left| \int_{t_0}^{t_1} \left[\int_{t_0}^t \cos u_j(\tau) d\tau \right] (1 - w(t)) \dot{w}(t) dt \right| \\
 & \leq \sup_{t \in [t_0, t_1]} \left\{ |1 - w(t)| \cdot |\dot{w}(t)| \cdot \left| \int_{t_0}^t \cos u_j(\tau) d\tau \right| \right\} \cdot (t_1 - t_0) \\
 \text{(A.2)} \quad & \leq \frac{cT_0\Gamma^2}{4K} \left[\Gamma \left(\frac{\Gamma^2}{8K} - 1 \right) + \frac{1}{\Gamma} \left(\frac{\Gamma^2}{4K} \right)^2 \right].
 \end{aligned}$$

Combining (A.1) and (A.2) yields

$$\begin{aligned}
 & \frac{1}{\Gamma} \left| \int_{t_0}^{t_1} \cos u_j(t) (1 - w(t))^2 dt \right| \\
 \text{(A.3)} \quad & < \frac{c}{\Gamma} \left(\frac{\Gamma^2}{4K} \right)^2 + \frac{cT_0\Gamma}{2K} \left[\Gamma \left(\frac{\Gamma^2}{8K} - 1 \right) + \frac{1}{\Gamma} \left(\frac{\Gamma^2}{4K} \right)^2 \right] = c\Delta,
 \end{aligned}$$

where

$$\Delta = \frac{1}{\Gamma} \left(\frac{\Gamma^2}{4K} \right)^2 + \frac{T_0}{2K} \left[\Gamma^2 \left(\frac{\Gamma^2}{8K} - 1 \right) + \left(\frac{\Gamma^2}{4K} \right)^2 \right].$$

From (2.6) and (2.7) it follows that $c \rightarrow 0$ as $F \rightarrow +\infty$. Choosing F large enough, we have $c\Delta < \frac{\Gamma^2}{4K} - 2$ and hence a contradiction with $w(t_1) = 1 - \frac{\Gamma^2}{4K}$. Consequently, there exists $F_2 \geq F_1$ such that $1 - \frac{\Gamma^2}{4K} < w(t) < 1$ for $t \in [0, T_0]$ if $F \geq F_2$.

Now we construct a first order equation

$$\dot{s} = \Gamma^{-1} [-\Gamma^2 s - 2K(1 - s)^2].$$

It is easy to check by the assumption $\Gamma > 2\sqrt{2K}$ that $s(t) \in (-1, 1)$ for $t > 0$ if $s(0) \in [-1, 1]$. Assume that $s(0) = w(0) = -1$. Then for $t \in [0, T_0]$ it follows that

$$\begin{aligned}
 |w(t) - s(t)| & \leq \frac{1}{\Gamma} \left| \int_0^t (2K(w(\tau) + s(\tau) - 2) + \Gamma^2) (w(\tau) - s(\tau)) d\tau \right| \\
 & \quad + \frac{1}{\Gamma} \left| \int_0^t \cos u_j(\tau) (1 - w(\tau))^2 d\tau \right|.
 \end{aligned}$$

Since $s(t) \in (-1, 1)$ and $1 - \frac{\Gamma^2}{4K} < w(t) < 1$ for $t \in (0, T_0]$, then we have

$$\text{(A.4)} \quad \left| \int_0^t (2K(w(\tau) + s(\tau) - 2) + \Gamma^2) (w(\tau) - s(\tau)) d\tau \right| \leq \int_0^t \Gamma^2 |w(\tau) - s(\tau)| d\tau,$$

and hence by (A.3) and (A.4),

$$|w(t) - s(t)| \leq \int_0^t \Gamma |w(\tau) - s(\tau)| d\tau + c\Delta.$$

The Gronwall inequality implies that

$$|w(t) - s(t)| \leq c \Delta e^{\Gamma t} \leq c \Delta e^{\Gamma T_0}, \quad t \in [0, T_0].$$

As a consequence, there exists $F_3 \geq F_1$, such that if $F \geq F_3$, then

$$c \Delta e^{\Gamma T_0} < s(T_0/2) + 1,$$

and hence $w(t) > -1$ for $t \in [T_0/2, T_0]$. Meanwhile, since $w(T_0) > -1$, then it follows that $1 - \Gamma^2/(4K) < w(t) < 1$ for $t \in [T_0, 2T_0]$ if $F \geq F_2$, and thus $1 - \Gamma^2/(4K) < w(t) < 1$ for all $t > 0$. The proof is completed by taking $F_0 = \max\{F_2, F_3\}$. \square

We should mention here that the above proof was highly inspired by Levi's idea on dealing with a single oscillator [10].

REFERENCES

- [1] D. G. ARONSON, M. GOLUBITSKY, AND J. MALLET-PARET, *Ponies on a merry-go-round in large arrays of Josephson junctions*, Nonlinearity, 4 (1991), pp. 903–910.
- [2] D. G. ARONSON AND Y. S. HUANG, *Single wave form solutions for linear arrays of Josephson junctions*, Phys. D, 101 (1997), pp. 157–177.
- [3] C. BAESENS AND R. S. MACKAY, *Gradient dynamics of tilted Frenkel-Kontorova models*, Nonlinearity, 11 (1998), pp. 949–964.
- [4] C. BAESENS AND R. S. MACKAY, *A novel preserved partial order for cooperative networks of units with overdamped second order dynamics, and application to tilted Frenkel-Kontorova chains*, Nonlinearity, 17 (2004), pp. 567–580.
- [5] O. M. BRAUN AND Y. S. KIVSHAR, *The Frenkel-Kontorova Model. Concepts, Methods, and Applications*, Springer-Verlag, Berlin, 2004.
- [6] L. M. FLORÍA AND J. J. MAZO, *Dissipative dynamics of the Frenkel-Kontorova model*, Adv. Phys., 45 (1996), pp. 505–598.
- [7] G. KATRIEL, *Existence of travelling waves in discrete sine-Gordon rings*, SIAM J. Math. Anal., 36 (2005), pp. 1434–1443.
- [8] M. LEVI, *Caterpillar solutions in coupled pendula*, Ergodic Theory Dynam. Systems, 8 (1988), pp. 153–174.
- [9] M. LEVI, *Dynamics of discrete Frenkel-Kontorova models*, in Analysis, Et Cetera, P. Rabinowitz and E. Zehnder, eds., Academic Press, Boston, 1990, pp. 471–494.
- [10] M. LEVI, *private communications*.
- [11] R. E. MIROLLO, *Splay-phase orbits for equivariant flows on tori*, SIAM J. Math. Anal., 25 (1994), pp. 1176–1180.
- [12] R. MIROLLO AND N. ROSEN, *Existence, uniqueness, and nonuniqueness of single-wave-form solutions to Josephson junction systems*, SIAM J. Appl. Math., 60 (2000), pp. 1471–1501.
- [13] M. QIAN, S. ZHU, AND W.-X. QIN, *Dynamics in a chain of overdamped pendula driven by constant torques*, SIAM J. Appl. Math., 57 (1997), pp. 294–305.
- [14] H. L. SMITH, *Monotone Dynamical Systems*, AMS, Providence, RI, 1995.
- [15] S. WATANABE AND J. W. SWIFT, *Stability of periodic solutions in series arrays of Josephson junctions with internal capacitance*, J. Nonlinear Sci., 7 (1997), pp. 503–536.
- [16] S. WATANABE, H. S. J. VAN DER ZANT, S. H. STROGATZ, AND T. P. ORLANDO, *Dynamics of circular arrays of Josephson junctions and the discrete sine-Gordon equation*, Phys. D, 97 (1996), pp. 429–470.

ERGODICITY AND RATE OF CONVERGENCE FOR A NONSECTORIAL FIBER LAY-DOWN PROCESS*

MARTIN GROTHAUS[†] AND AXEL KLAR[†]

Abstract. A stochastic model for the lay-down of fibers on a conveyor belt in the production process of nonwovens is investigated. In particular, convergence of the stochastic process to the stationary solution is proven and estimates on the speed of convergence are given. Numerical results and examples are presented and compared with the analytical estimates on the speed of convergence.

Key words. fiber dynamics, Fokker–Planck equations, Dirichlet forms, ergodicity, rate of convergence

AMS subject classifications. 37A25, 47D07, 82C31

DOI. 10.1137/070697173

1. Introduction. The understanding of the forms generated by the lay-down of flexible fibers onto a moving conveyor belt is of great interest in the production process of nonwovens that find their applications, e.g., in composite materials (filters), textiles, and the hygiene industry. In the melt-spinning process of nonwoven materials, hundreds of individual endless fibers obtained by the continuous extrusion of a melted polymer are stretched and entangled by highly turbulent air flows to finally form a web on the conveyor belt. The quality of this web and the resulting nonwoven material—in terms of homogeneity and load capacity—depends essentially on the dynamics and the deposition of the fibers.

For the description of the interaction between fibers and turbulent flow a stochastic force model is derived and analyzed in [MW06]. Applying this concept, the fiber fabric can in principle be numerically generated and its quality investigated. However, these or similar simulations usually lead to excessively large computation times, when all physical details of the production process are considered. Thus, simplified models for the lay-down process are needed. In particular, this is true for optimization and control procedures where many different simulations are needed. In [GKM+07] a new simplified stochastic model for the fiber lay-down process, i.e., for the generation of a fiber web on a conveyor belt, has been presented. The process, which takes into account the fiber motion under the influence of turbulence, is described by a stochastic differential system; see section 2.

The solution of the associated Fokker–Planck equation gives the density of the stochastic process. An important criterion for the quality of the web and the resulting nonwoven material is how it converges to equilibrium. In particular, the speed of convergence to the stationary solution is important. The faster this convergence is, the more uniform the produced textile will be. This means that process parameters should be adjusted such that the speed of convergence to equilibrium is optimal.

The trend to equilibrium for Fokker–Planck-type equations has been investigated in many papers. In [AMTU01] a unified presentation of entropy methods for

*Received by the editors July 13, 2007; accepted for publication (in revised form) May 27, 2008; published electronically September 19, 2008. This work was supported by the Kaiserslautern Excellence Cluster *Dependable Adaptive Systems and Mathematical Modeling*.

<http://www.siam.org/journals/sima/40-3/69717.html>

[†]Mathematics Department, University of Kaiserslautern, P.O. Box 3049, 67653 Kaiserslautern, Germany (grothaus@mathematik.uni-kl.de, klar@mathematik.uni-kl.de).

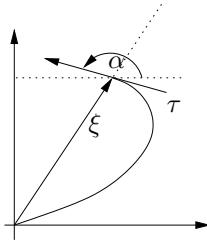


FIG. 1. *Fiber scenario on the conveyor belt.*

nondegenerate linear and nonlinear Fokker–Planck-type equations is given. In [DV01] the linear kinetic Fokker–Planck equation is discussed. In this case the degeneracy of the diffusion operator is similar to the present case. A more general theory for degenerate linear and nonlinear problems has been developed in [Vil06]. A more detailed discussion of these methods with respect to the present problem is given in the next section.

As in the above-mentioned papers, the purpose of the present paper is to investigate analytically the convergence to equilibrium. In particular, we aim at obtaining explicit estimates on the rate of convergence. However, instead of using entropy methods we use here Dirichlet form and operator semigroup techniques.

The paper is organized as follows. In section 2 we present the stochastic model for the fiber dynamics. Furthermore, we relate our convergence result for the solution of the stochastic differential equation with the convergence of the solution of the associated Fokker–Planck equation and outline our strategy for proving ergodicity. In section 3 we prove existence results for the stochastic process as well as several other properties of the process which are important for our analysis. Section 4 investigates several properties of an associated nondegenerate Ornstein–Uhlenbeck-type process. Section 5 proves the ergodicity of the original process and gives explicit bounds on the rates to convergence. Section 6 is devoted to some numerical results comparing the analytical rates of convergence with the numerical ones.

2. The stochastic model for the fiber lay-down process and its associated Fokker–Planck equation. Focusing on a single slender elastic inextensible fiber in a lay-down process, the fiber on the nonmoving conveyor belt can be described by an arc-length parameterized curve $\xi : \mathbb{R}_0^+ \rightarrow \mathbb{R}^2$, as visualized in Figure 1.

Due to its inextensibility, $\|\partial_t \xi\| = 1$ holds. The web-forming is modeled as

$$\begin{aligned} \partial_t \xi &= \tau(\alpha), \\ \partial_t \alpha &= -\nabla \phi(\xi) \cdot \tau^\perp(\alpha), \end{aligned}$$

where $\tau(\alpha) = (\cos \alpha, \sin \alpha)^T$ denotes the normalized tangent on the fiber. Since a curved fiber tends back to its starting point, the change of the angle α is assumed to be proportional to $\nabla \phi(\xi) \cdot \tau^\perp(\alpha)$ with $\tau^\perp(\alpha) = (-\sin \alpha, \cos \alpha)^T$. The potential of this drive is prescribed by a function $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$ with the generic example $\phi(\xi) = \|\xi\|^2/2$. In general, the potential depends on the physical process and needs to be adapted to the experimental parameters.

Considering a turbulent flow in the deposition region of the fiber close to the conveyor belt, the fiber lay-down is additionally affected by a stochastic force that can be modeled by a Wiener process B_t in \mathbb{R} with amplitude σ .

The resulting stochastic differential system reads

$$(2.1) \quad d\xi_t = \tau(\alpha_t) dt,$$

$$(2.2) \quad d\alpha_t = -\nabla\phi(\xi_t) \cdot \tau^\perp(\alpha_t) dt + \sigma dB_t,$$

where the conservation of length $\|\partial_t \xi\| = 1$ is still valid. For details on the mathematical modeling of the process we refer to [GKM+07]. An illustration of the pathwise behavior of the process for different σ and $\phi(\xi) = \|\xi\|^2/2$ is displayed in Figure 2, where the paths of Monte Carlo simulations for different σ are shown. The slow convergence to equilibrium is clearly seen for small and large values of σ .

The density of the above stochastic process can be found by solving the associated Fokker–Planck equation:

$$(2.3) \quad \partial_t u + \cos(\alpha)\partial_{\xi^1} u + \sin(\alpha)\partial_{\xi^2} u - \partial_\alpha \left(\nabla\phi(\xi) \begin{pmatrix} -\sin(\alpha) \\ \cos(\alpha) \end{pmatrix} u \right) = \frac{\sigma^2}{2} \partial_\alpha^2 u.$$

A stationary solution of (2.3) is the function

$$(\xi, \alpha) \mapsto m(\xi) := \frac{1}{N} \exp(-\phi(\xi)), \quad N > 0.$$

For $\phi(\xi) = \|\xi\|^2/2$ we obtain the standard Gaussian as stationary distribution. By μ we denote the measure on $E = \mathbb{R}^2 \times [0, 2\pi]$ having density m with respect to (w.r.t.) the Lebesgue measure. Except for sections 3 and 4 we consider only potentials ϕ such that m is integrable and choose N as a normalizing constant to obtain a probability measure μ . In the following sections we are concerned with the approach to equilibrium of the stochastic process (2.1) or the convergence to the stationary solution of (2.3), respectively.

Remark 2.1. A formal argument for convergence to equilibrium of the solution of (2.3) is given by the following computations; cf. [DV01]. Define the relative entropy

$$H(u/m) = \int_{\mathbb{R}^2} \int_0^{2\pi} u \log \left(\frac{u}{m} \right) d\alpha d\xi.$$

An easy computation gives

$$\partial_t H(u/m) = - \int_{\mathbb{R}^2} \int_0^{2\pi} u \left(\partial_\alpha \log \left(\frac{u}{m} \right) \right)^2 d\alpha d\xi.$$

This vanishes if and only if

$$u = \rho(t, \xi)m(\xi).$$

That means u converges formally to a local equilibrium. Plugging $u = \rho(t, \xi)m(\xi)$ into the evolution equation, one obtains

$$m\partial_t \rho + \cos(\alpha) (\partial_{\xi^1}(\rho m) + \rho m \partial_{\xi^1} \phi) + \sin(\alpha) (\partial_{\xi^2}(\rho m) + \rho m \partial_{\xi^2} \phi) = 0$$

or

$$\partial_t \rho = 0, \quad \partial_{\xi^1} \rho = 0, \quad \partial_{\xi^2} \rho = 0.$$

This means $u = m(\xi)$, i.e., convergence to global equilibrium.

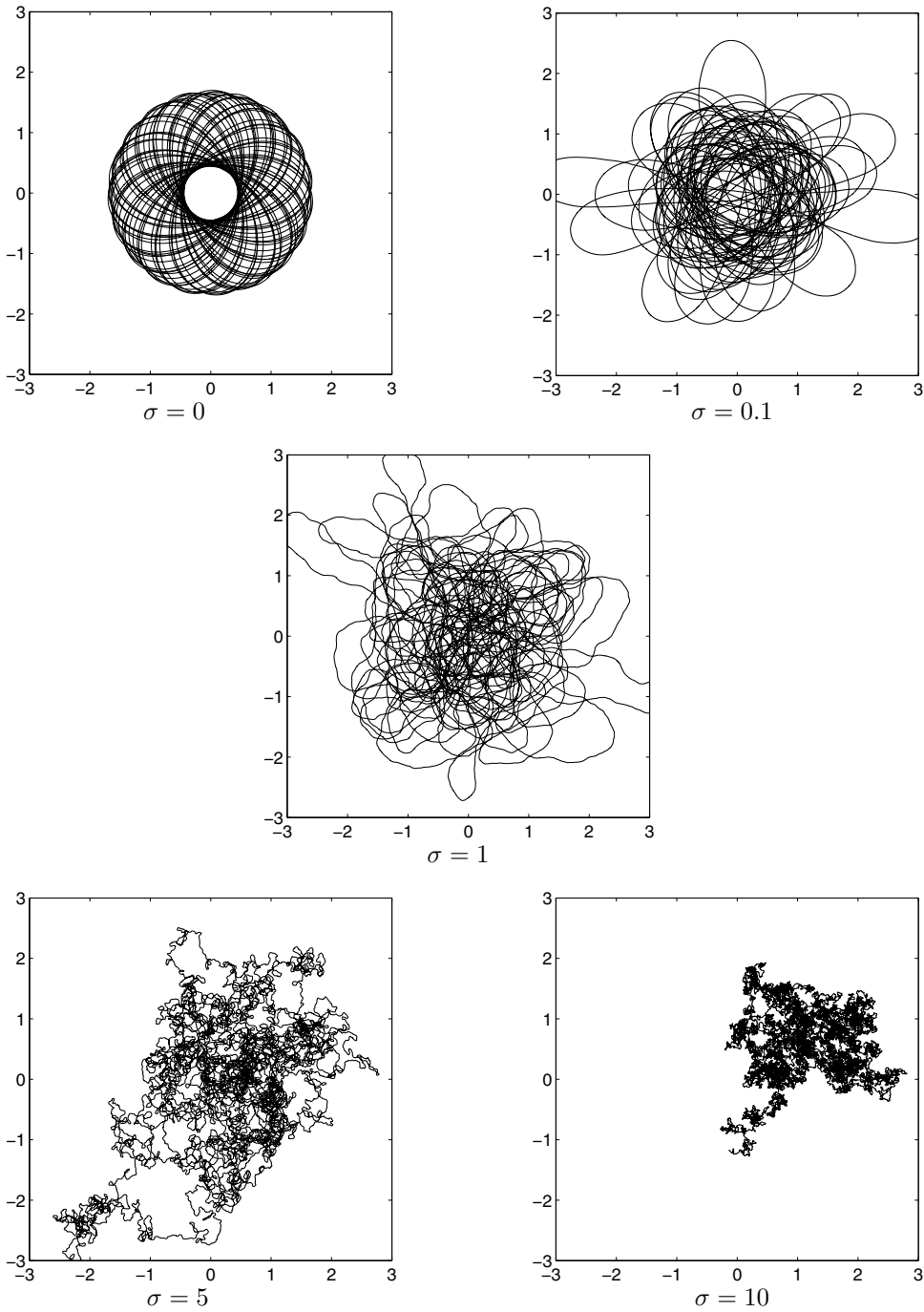


FIG. 2. Representative path behavior for balanced ($\sigma = 1$) as well as deterministic ($\sigma < 1$) and stochastic ($\sigma > 1$) dominated (ξ, α) -systems.

In [Vil06] the arguments described in the above remark were rigorously justified for similar equations. To apply these methods in the context of the present paper among other conditions, the existence of proper Liapunov functions would be necessary. Their existence is another open problem. In the following we choose a different route and consider the stochastic differential equation (2.1) directly.

In section 3 we construct a diffusion process solving (2.1) and having μ as an invariant measure. The law of this process we denote by \mathbf{P}_μ . Then in section 5 we prove that this diffusion is ergodic with rate of convergence

$$(2.4) \quad \left\| \frac{1}{t} \int_0^t f(\mathbf{X}_s) ds - \mathbf{E}_\mu[f] \right\|_{L^2(\mathbf{P}_\mu)} \leq \frac{1}{c^{1/2}} \left(\frac{2}{t} + \frac{1}{t^{1/2}} \left(\frac{A(C)c^{1/2} + B(C)c^{-1/2}}{\sigma} + (1 + 2^{-1/2})\sigma \right) \right) \|f - \mathbf{E}_\mu[f]\|_{L^2(\mu)}, \quad t > 0;$$

see Theorem 5.3 for details. The convergence in (2.4) implies mean ergodicity of the corresponding semigroup $(T_t)_{t \geq 0}$, i.e.,

$$(2.5) \quad \left\| \frac{1}{t} \int_0^t T_s f ds - \mathbf{E}_\mu[f] \right\|_{L^2(\mu)} \leq \frac{1}{c^{1/2}} \left(\frac{2}{t} + \frac{1}{t^{1/2}} \left(\frac{A(C)c^{1/2} + B(C)c^{-1/2}}{\sigma} + (1 + 2^{-1/2})\sigma \right) \right) \|f - \mathbf{E}_\mu[f]\|_{L^2(\mu)}, \quad t > 0.$$

In the case of an analytic semigroup from (2.5) we even could conclude strongly mixing, i.e.,

$$\lim_{t \rightarrow \infty} \|T_t f - \mathbf{E}_\mu[f]\|_{L^2(\mu)} = 0;$$

see, e.g., [Gol85, Exer. 8.24.17]. In our case, however, the corresponding generator

$$(2.6) \quad L = S + A, \quad S = \frac{\sigma^2}{2} \partial_\alpha^2, \quad A = \cos(\alpha) \partial_{\xi^1} + \sin(\alpha) \partial_{\xi^2} - \nabla \phi(\xi) \cdot \tau^{-1}(\alpha) \partial_\alpha$$

is nonsectorial. Hence, $(T_t)_{t \geq 0}$ cannot be the restriction of an analytic semigroup; see, e.g., [Gol85, Thm. 5.9].

However, since the adjoint to L w.r.t. the scalar product in $L^2(\mu)$ is given by $L^* = S - A$, then also for the adjoint process and semigroup we have ergodicity with the same rate of convergence.

Now, strong mixing of $(T_t^*)_{t \geq 0}$ would imply L^1 -convergence of the solution u of the associated Fokker–Planck equation (2.3) with normalized nonnegative initial distribution $u(0) = f$, because

$$\|u(t) - m\|_{L^1(dx)} = \left\| \frac{u(t)}{m} - 1 \right\|_{L^1(\mu)} = \|T_t^* f - 1\|_{L^1(\mu)} \leq \|T_t^* f - 1\|_{L^2(\mu)}.$$

Although $(L, D(L))$ is nonsectorial, we are convinced that for C^∞ -potentials ϕ we have strong mixing by another reasoning not worked out rigorously in the present paper. Since the generator $(L, D(L))$ is hypoelliptic in the sense of Hörmander, it is not too hard to show that $(T_t)_{t \geq 0}$ is strong Feller. Showing that $(T_t)_{t \geq 0}$ additionally is irreducible, strong mixing then follows from Doob’s theorem; see, e.g., [DPZ96, Prop. 4.1.1, Thm. 4.2.1]. But here we would like to stress that this does not give an explicit rate of convergence. Moreover, our approach also applies to C^3 -potentials ϕ .

That $((T_t)_{t \geq 0})$ is weak mixing, i.e.,

$$\lim_{t \rightarrow \infty, t \in I} \|T_t f - \mathbf{E}_\mu[f]\|_{L^2(\mu)} = 0,$$

where $I \subset \mathbb{R}_0^+$ has relative measure 1, follows from weaker spectral properties of $(L, D(L))$ than those giving rise to a corresponding analytic semigroup. Sufficient for $(T_t)_{t \geq 0}$ being weak mixing is the following condition:

Let $f \in L^2(\mu; \mathbb{C})$ and $\lambda \in \mathbb{R}$. If

$$(2.7) \quad T_t f = \exp(i\lambda t) f \quad \text{for all } t \geq 0,$$

then $\lambda = 0$ and f is constant; see, e.g., [DPZ96, Thm. 3.4.1]. Here of course we have to consider the corresponding complexified spaces.

The condition in (2.7) implies that $f \in D(L)$ is an eigenvector to the eigenvalue $i\lambda$, $\lambda \in \mathbb{R}$. If f is now sufficiently smooth, such that we can apply the operators S and A from (2.6) separately, then this condition easily can be verified (S is a symmetric and A an antisymmetric operator in $L^2(\mu)$). Hence it is left to show that eigenvectors to $(L, D(L))$ are sufficiently smooth. This is not clear a priori, since S is degenerated. But $L = S + A$ is of Hörmander form such that we expect to obtain hypoelliptic estimates.

The interplay of S and A is very crucial for our proof of ergodicity. Only the operator S without A would not give ergodic behavior. The idea is to project onto the orthogonal complement of the kernel of $(S, D(S))$. On this subspace $(S, D(S))$ has a bounded inverse. On the kernel of $(S, D(S))$ in turn, L can be associated with a nondegenerated, self-adjoint operator $(G, D(G))$ in $L^2(\gamma)$, where γ is the marginal measure of μ on \mathbb{R}^2 . Assuming that the Dirichlet form corresponding to $(G, D(G))$ fulfills a Poincaré inequality, finally ergodicity can be shown. Such ideas have been used before in, e.g., [OT03] in the context of scaling limit for interacting particles systems.

Crucial for the rate of convergence is Lemma 4.1. There Kato perturbation techniques and a ground state transform are used to show that the functions $h = \nabla_\xi g \cdot \tau$, $g \in D(G)$, are in $D(L)$. These functions are our “key functions” for relating $(L, D(L))$ with $(G, D(G))$; see the proof of Theorem 5.3. The explicit estimates obtained in the proof of Lemma 4.1 finally provide us with the constants for our rate of convergence (2.4).

3. Stochastic differential equation and its solution. We are considering the following stochastic differential equation in $E := \mathbb{R}^2 \times [0, 2\pi]$:

$$(3.1) \quad \begin{aligned} d\xi_t &= \tau(\alpha_t) dt, \\ d\alpha_t &= -\nabla\phi(\xi_t) \cdot \tau^\perp(\alpha_t) dt + \sigma dB_t, \quad \sigma > 0, \quad t \geq 0. \end{aligned}$$

Using Itô’s formula we find the generator of the corresponding Markov process:

$$L = S + A, \quad S = \frac{\sigma^2}{2} \partial_\alpha^2, \quad A = \cos(\alpha) \partial_{\xi^1} + \sin(\alpha) \partial_{\xi^2} - \nabla\phi(\xi) \cdot \tau^\perp(\alpha) \partial_\alpha.$$

Set

$$D := C_{0,\text{pb}}^\infty(E) := \{f|_E \mid f \in C_0^\infty(\mathbb{R}^3), f(\xi, 0) = f(\xi, 2\pi) \text{ for all } \xi \in \mathbb{R}^2$$

and the same holds for all derivatives of $f\}$,

where $C_0^\infty(\mathbb{R}^3)$ denotes the set of compactly supported, infinitely differentiable functions on \mathbb{R}^3 . (S, D) is symmetric and (A, D) is antisymmetric w.r.t. the scalar product in $L^2(\mu)$, where μ is the measure on E having density m w.r.t. the Lebesgue measure.

In the following theorem $H^{1,\infty}(\mathbb{R}^2)$ denotes the Sobolev space of weakly differentiable functions on \mathbb{R}^2 , which are essentially bounded together with their derivative. The adjoint operator $(L^*, D(L^*))$ to (L, D) obviously is a closed extension of $(S - A, D)$.

LEMMA 3.1. *Let $\phi \in H^{1,\infty}(\mathbb{R}^2)$. Then (L, D) and $(L^*|_D, D)$ are essentially m -dissipative Dirichlet operators in $L^2(\mu)$. In particular, $(L^*, D(L^*))$ is the m -dissipative extension of $(L^*|_D, D)$. The m -dissipative extension (closure) of (L, D) we denote by $(L, D(L))$.*

Proof. Since (S, D) is dissipative and (A, D) is antisymmetric, obviously

$$(Lf, f)_{L^2(\mu)} \leq 0 \quad \text{for all } f \in D,$$

i.e., (L, D) is dissipative. Moreover, as in, e.g., [CG08a] one can show that

$$(Lf, (f - 1)^+)_{L^2(\mu)} \leq 0 \quad \text{for all } f \in D,$$

i.e., (L, D) is a Dirichlet operator (here $g^+ := (|g| + g)/2$ for a mapping $g : E \rightarrow \mathbb{R}$). Analogously one obtains that $(L^*|_D, D)$ is also a dissipative Dirichlet operator in $L^2(\mu)$. Hence it is left to show that D is a core for the m -dissipative extensions of (L, D) and $(L^*|_D, D)$.

First we consider the case $\phi = 0$. Then we have

$$Lf = Sf + \cos(\alpha)\partial_{\xi_1} f + \sin(\alpha)\partial_{\xi_2} f, \quad f \in D.$$

In the following we show that in this case L is essentially m -dissipative on

$$F := S(\mathbb{R}^2) \times C_{\text{pb}}^\infty([0, 2\pi]),$$

where $S(\mathbb{R}^2)$ denotes the space of infinitely differentiable, more than polynomial decreasing functions on \mathbb{R}^2 . Since $D \subset F$ is dense in graph norm, this implies that (L, D) is essentially m -dissipative.

Consider the complexified setting, i.e., the corresponding spaces

$$F_{\mathbb{C}} := S(\mathbb{R}^2; \mathbb{C}) \times C_{\text{pb}}^\infty([0, 2\pi]; \mathbb{C}) \subset L^2(\mu; \mathbb{C}),$$

of complex valued functions, and let L act componentwise. Then we apply the Fourier transform \mathcal{F} in the first variable and obtain the corresponding operator

$$\hat{L} := \mathcal{F}L\mathcal{F}^{-1} = Sf + i \cos(\alpha)\xi_1 + i \sin(\alpha)\xi_2.$$

Recall that $\mathcal{F} : F_{\mathbb{C}} \rightarrow F_{\mathbb{C}}$ is a continuous isomorphism and $\mathcal{F} : L^2(\mu; \mathbb{C}) \rightarrow L^2(\mu; \mathbb{C})$ is a unitary isomorphism. Furthermore, by an elementary consideration one shows that $F_{\mathbb{C}}$ is a core for the m -dissipative extension of $(\hat{L}, F_{\mathbb{C}})$. Hence, L is essentially m -dissipative on F .

Finally the case $\phi \in H^{1,\infty}(\mathbb{R}^2)$ follows by a standard Kato perturbation technique. The adjoint operator can be treated analogously. \square

Since $(L, D(L))$ is m -dissipative, it generates a semigroup of contractions $(T_t)_{t \geq 0}$. Let $(T_t^*)_{t \geq 0}$ denote the dual semigroup. It is well known that its generator is $(\bar{L}^*, D(L^*))$. The corresponding stochastic process we can construct as a diffusion process

(i.e., Markov process with continuous sample path) only on the manifold $\mathbb{R}^2 \times S^1$ which is naturally associated with E . Here S^1 denotes the unit circle. Since the measure μ and above function spaces naturally can be lifted on $\mathbb{R}^2 \times S^1$, in what follows we simply identify

$$E = \mathbb{R}^2 \times S^1.$$

By γ we denote the marginal measure of μ on \mathbb{R}^2 .

THEOREM 3.2. *Let m be bounded and continuous, $m > 0$ almost everywhere w.r.t. the Lebesgue measure on \mathbb{R}^2 , and ϕ weakly differentiable on $\{m > 0\}$ with $\nabla\phi \in L^2(\gamma)$. Furthermore, let $\nabla\phi$ be bounded on the sets $\{\phi \leq K\}$ for all $K \in \mathbb{R}$. Then for any probability measure ν which is absolutely continuous and has a bounded density w.r.t. μ there exists a Markov process $(C([0, \infty); E), \mathcal{B}, (\mathbf{F}_t)_{t \geq 0}, (\mathbf{X}_t)_{t \geq 0}, \mathbf{P}_\nu)$ having the following properties:*

- (i) *The initial distribution of the process is given by ν . If μ is a probability measure, then it is an invariant measure for the process.*
- (ii) *Finite-dimensional distributions of \mathbf{P}_ν are given in terms of a contraction semigroup $(T_t)_{t \geq 0}$ having as generator $(L, D(L))$ an m -dissipative extension of (L, D) . This determines the law \mathbf{P}_ν uniquely.*
- (iii) *The process solves (3.1) in the sense of the corresponding martingale problem, i.e., for all $f \in D$,*

$$\mathbf{M}(f)_t = f(\mathbf{X}_t) - f(\mathbf{X}_0) - \int_0^t Lf(\mathbf{X}_s) ds, \quad t \geq 0,$$

is an \mathbf{F}_t -martingale under \mathbf{P}_ν .

Analogous statements hold for the adjoint operators. The associated stochastic process $(C([0, \infty); E), \mathcal{B}, (\mathbf{F}_t)_{t \geq 0}, (\mathbf{X}_t)_{t \geq 0}, \mathbf{P}_\nu^)$ is called the adjoint process.*

Proof. First we consider $\phi \in H^{1,\infty}(\mathbb{R}^2)$. Then from Lemma 3.1 we can conclude that $(L, D(L))$ is an m -dissipative Dirichlet operator with core D . Hence $(L, D(L))$ generates a quasi-regular, generalized Dirichlet form and therefore has an associated Markov process; see [Sta99].

For a general ϕ a corresponding stochastic process can be constructed as in [CG08b, sec. 2.3] by an approximation via Markov process associated to potentials from $H^{1,\infty}(\mathbb{R}^2)$.

Then as in [CG08b] and [CG08a], respectively, one can show that it is a diffusion with the properties (i)–(iii). \square

Remark 3.3. (i) For $\phi \in H^{1,\infty}(\mathbb{R}^2)$ one can obtain much stronger existence results, as shown in [CG08a], where the concepts of generalized Dirichlet forms were applied to a similar problem.

(ii) The process $(\mathbf{X}_t)_{t \geq 0}$ in Theorem 3.2 is the coordinate process, i.e., $X_t(\omega) = \omega(t)$, $t \geq 0$, $\omega \in C([0, \infty); E)$, \mathcal{B} is the Borel σ -algebra, and $(\mathbf{F}_t)_{t \geq 0}$ is the corresponding natural filtration.

COROLLARY 3.4. *Let the assumptions of Theorem 3.2 hold and let μ be a probability measure. Then*

$$(3.2) \quad \mathbf{M}(f)_t = f(\mathbf{X}_t) - f(\mathbf{X}_0) - \int_0^t Lf(\mathbf{X}_s) ds, \quad t \geq 0,$$

is a \mathbf{F}_t -martingale under \mathbf{P}_μ for all f from \overline{D}^L , the closure of D w.r.t. the graph norm to $(L, D(L))$. For $f \in \overline{D}^L$ the quadratic variation process of $\mathbf{M}(f)_t$ is given by

$$(3.3) \quad \langle \mathbf{M}(f) \rangle_t = \sigma^2 \int_0^t |\partial_\alpha f|^2(\mathbf{X}_s) ds, \quad t \geq 0.$$

The corresponding statement holds for the adjoint process. The corresponding martingale we denote by

$$(3.4) \quad \mathbf{M}^*(f)_t := f(\mathbf{X}_t) - f(\mathbf{X}_0) - \int_0^t L^* f(\mathbf{X}_s) ds, \quad t \geq 0,$$

where $f \in \overline{D}^{L^*}$. For $f \in \overline{D}^{L^*}$ again we have

$$(3.5) \quad \langle \mathbf{M}^*(f) \rangle_t = \sigma^2 \int_0^t |\partial_\alpha f|^2(\mathbf{X}_s) ds, \quad t \geq 0.$$

Proof. The enlargement of the class of admissible functions for the martingale problem is obvious. Formulas (3.3) and (3.5) for the corresponding quadratic variation processes can be derived as in [CG08b, sec. 3]. \square

4. Associated nondegenerated generator. In $L^2(\gamma)$ we consider the pre-Dirichlet form

$$\mathcal{E}(f, g) = \frac{1}{2} \int_{\mathbb{R}^2} \nabla_\xi f \cdot (\nabla_\xi g)^T d\gamma, \quad f, g \in C_0^\infty(\mathbb{R}^2).$$

Its closure $(\mathcal{E}, D(\mathcal{E}))$ has a corresponding self-adjoint generator $(G, D(G))$ which on smooth, compactly supported functions is given by

$$Gf = \frac{1}{2} \Delta_\xi f - \frac{1}{2} \nabla_\xi \phi \cdot (\nabla_\xi f)^T, \quad f \in C_0^\infty(\mathbb{R}^2).$$

LEMMA 4.1. *Let $\phi \in C^3(\mathbb{R}^2)$, $\nabla \phi \in L^2(\gamma)$, and assume there exist $0 < C < \infty$ and $K \subset \mathbb{R}^2$ compact such that*

$$\begin{aligned} & (\partial_{\xi_i} \partial_{\xi_j} \partial_{\xi_k} \phi(\xi))^2 + (\partial_{\xi_i} \partial_{\xi_j} \phi(\xi))^2 \\ & \leq C (\partial_{\xi^1} \phi(\xi))^2 + (\partial_{\xi^2} \phi(\xi))^2 \quad \text{for all } i, j, k \in \{1, 2\}, \xi \in \mathbb{R}^2 \setminus K. \end{aligned}$$

Then

$$\nabla_\xi f \cdot \tau \in \overline{D}^L \subset D(L)$$

for all $f \in D(G)$.

Proof. First recall that $D(\mathcal{E}) \subset D(G)$; hence $\nabla_\xi f \cdot \tau \in L^2(\gamma)$ is well-defined. Our assumptions imply that $(G, C_0^\infty(\mathbb{R}^2))$ is essentially self-adjoint; see, e.g., [BKR97, Thm. 7]. Hence there exists a sequence $(f_n)_{n \in \mathbb{N}}$ in $C_0^\infty(\mathbb{R}^2)$ which converges to f in the graph norm. Moreover

$$\nabla_\xi f_n \cdot \tau \in D \subset D(L) \quad \text{for all } n \in \mathbb{N}.$$

Obviously, $(\nabla_\xi f_n \cdot \tau)_{n \in \mathbb{N}}$ converges to $\nabla_\xi f \cdot \tau$ in $L^2(\mu)$. Hence it is sufficient to show that $(L(\nabla_\xi f_n \cdot \tau))_{n \in \mathbb{N}}$ is a Cauchy sequence in $L^2(\mu)$.

We fix $g \in C_0^\infty(\mathbb{R}^2)$. Then

$$\begin{aligned}
 (4.1) \quad \|L(\nabla_\xi g \cdot \tau)\|_{L^2(\mu)}^2 &= \int_{\mathbb{R}^2} \int_{[0,2\pi]} (L(\cos(\alpha)\partial_{\xi^1}g(\xi) + \sin(\alpha)\partial_{\xi^2}g(\xi)))^2 d\alpha d\gamma(\xi) \\
 &= \frac{1}{8} \left(\sigma^4 \mathcal{E}(g, g) + \int_{\mathbb{R}^2} ((2\partial_{\xi^1}\partial_{\xi^2} + \partial_{\xi^1}\phi(\xi)\partial_{\xi^2} + \partial_{\xi^2}\phi(\xi)\partial_{\xi^1})g(\xi))^2 \right. \\
 &\quad + 3((\partial_{\xi^1}^2 - \partial_{\xi^2}\phi(\xi)\partial_{\xi^2})g(\xi))^2 + 3((\partial_{\xi^2}^2 - \partial_{\xi^1}\phi(\xi)\partial_{\xi^1})g(\xi))^2 \\
 &\quad \left. + 2(\partial_{\xi^1}^2 - \partial_{\xi^2}\phi(\xi)\partial_{\xi^2})g(\xi)(\partial_{\xi^2}^2 - \partial_{\xi^1}\phi(\xi)\partial_{\xi^1})g(\xi)d\gamma(\xi) \right).
 \end{aligned}$$

Therefore $(L(\nabla_\xi f_n \cdot \tau))_{n \in \mathbb{N}}$ is a Cauchy sequence if the operators

$$\partial_{\xi^i}\partial_{\xi^j}, \quad \partial_{\xi^i}\phi(\xi)\partial_{\xi^j}, \quad i, j \in \{1, 2\},$$

are Kato bounded by G on $C_0^\infty(\mathbb{R}^2)$.

To show this, it is useful to introduce the operators

$$G_i := \partial_{\xi^i}^2 - \partial_{\xi^i}\phi(\xi)\partial_{\xi^i}, \quad i \in \{1, 2\}.$$

A straightforward integration by parts yields the following for $g \in C_0^\infty(\mathbb{R}^2)$:

$$\begin{aligned}
 (G_1g, G_1g)_{L^2(\gamma)} + (G_2g, G_2g)_{L^2(\gamma)} &\leq C_1(Gg, Gg)_{L^2(\gamma)} + C_2(g, g)_{L^2(\gamma)} \\
 + C_3 \sum_{i=1}^2 ((\partial_{\xi^i}\phi)g, (\partial_{\xi^i}\phi)g)_{L^2(\gamma)} + C_4 \sum_{i,j=1}^2 ((\partial_{\xi^i}\partial_{\xi^j}\phi)g, (\partial_{\xi^i}\partial_{\xi^j}\phi)g)_{L^2(\gamma)} \\
 + C_5 \sum_{i,j,k=1}^2 ((\partial_{\xi^i}\partial_{\xi^j}\partial_{\xi^k}\phi)g, (\partial_{\xi^i}\partial_{\xi^j}\partial_{\xi^k}\phi)g)_{L^2(\gamma)}
 \end{aligned}$$

and

$$\begin{aligned}
 (\partial_{\xi^i}^2g, \partial_{\xi^i}^2g)_{L^2(\gamma)} &\leq C_6(G_ig, G_ig)_{L^2(\gamma)} + C_7(g, g)_{L^2(\gamma)} + C_8((\partial_{\xi^i}^2\phi)g, (\partial_{\xi^i}^2\phi)g)_{L^2(\gamma)} \\
 + C_9((\partial_{\xi^i}^3\phi)g, (\partial_{\xi^i}^3\phi)g)_{L^2(\gamma)}, \quad i \in \{1, 2\},
 \end{aligned}$$

and

$$\begin{aligned}
 (\partial_{\xi^2}\partial_{\xi^1}g, \partial_{\xi^2}\partial_{\xi^1}g)_{L^2(\gamma)} &\leq C_{10}(G_1g, G_1g)_{L^2(\gamma)} + C_{11}(G_2g, G_2g)_{L^2(\gamma)} + C_{12}(g, g)_{L^2(\gamma)} \\
 + C_{13} \sum_{i,j,k=1}^2 ((\partial_{\xi^i}\partial_{\xi^j}\partial_{\xi^k}\phi)g, (\partial_{\xi^i}\partial_{\xi^j}\partial_{\xi^k}\phi)g)_{L^2(\gamma)}
 \end{aligned}$$

and

$$\begin{aligned}
 (\partial_{\xi^i}\phi\partial_{\xi^j}g, \partial_{\xi^i}\phi\partial_{\xi^j}g)_{L^2(\gamma)} &\leq C_{14}(G_jg, G_jg)_{L^2(\gamma)} + C_{15}(g, g)_{L^2(\gamma)} + C_{16}((\partial_{\xi^i}\phi)g, (\partial_{\xi^i}\phi)g)_{L^2(\gamma)} \\
 + C_{17} \sum_{k,l=1}^2 ((\partial_{\xi^k}\partial_{\xi^l}\phi)g, (\partial_{\xi^k}\partial_{\xi^l}\phi)g)_{L^2(\gamma)}, \quad i \neq j, i, j \in \{1, 2\},
 \end{aligned}$$

where $0 < C_1, \dots, C_{17} < \infty$ are constants. Hence it is left to show that the operators given through multiplication by the functions $\partial_{\xi^i} \phi$, $\partial_{\xi^i} \partial_{\xi^j} \phi$ and $\partial_{\xi^i} \partial_{\xi^j} \partial_{\xi^k} \phi$, $i, j, k \in \{1, 2\}$, are Kato bounded by G on $C_0^\infty(\mathbb{R}^2)$.

To get this, we transform the problem to $L^2(dx)$. That is, we use the isometric isomorphism

$$L^2(\gamma) \ni f \mapsto \exp(-\phi/2)f \in L^2(dx)$$

to define

$$\hat{G} := \exp(-\phi/2)G \exp(\phi/2) = -\frac{1}{2}\Delta_\xi + \frac{1}{8} \left((\partial_{\xi^1} \phi)^2 + (\partial_{\xi^2} \phi)^2 \right) - \frac{1}{4} \left((\partial_{\xi^1}^2 \phi) + (\partial_{\xi^2}^2 \phi) \right)$$

on $\hat{D} := \{\exp(-\phi/2)f \mid f \in C_0^\infty(\mathbb{R}^2)\}$. Using our assumptions on ϕ one easily shows that the multiplication operators $\partial_{\xi^i} \phi$, $\partial_{\xi^i} \partial_{\xi^j} \phi$ and $\partial_{\xi^i} \partial_{\xi^j} \partial_{\xi^k} \phi$, $i, j, k \in \{1, 2\}$, are Kato bounded by \hat{G} on \hat{D} . Hence they are also Kato bounded by G on $C_0^\infty(\mathbb{R}^2)$. \square

Remark 4.2. The assumptions in Lemma 4.1 allow a large class of potentials ϕ . For example, the potentials $\phi = \|\cdot\|^p$, $p = 2, 4$ or $p \geq 6$, are admissible. Of course, ϕ may also be nonisotropic.

Remark 4.3. It is interesting to note that the stochastic process associated to the generator G is obtained from the original process (2.1) in the large σ -limit; see [BGK+07].

5. Ergodicity and rate of convergence. From now on we assume that μ is a probability measure on E . We are interested in the convergence to zero of

$$\left\| \frac{1}{t} \int_0^t f(\mathbf{X}_s) ds - \mathbf{E}_\mu[f] \right\|_{L^2(\mathbf{P}_\mu)} = \mathbf{E}_{\mathbf{P}_\mu} \left[\left(\frac{1}{t} \int_0^t f(\mathbf{X}_s) ds - \mathbf{E}_\mu[f] \right)^2 \right]^{1/2} \quad \text{as } t \rightarrow \infty$$

for $f \in L^2(\mu)$. Denote by $\mathbf{E}_\mu[f|p^1]$ the conditional expectation of $f \in L^2(\mu)$ w.r.t. the σ -algebra in $E = \mathbb{R}^2 \times [0, 2\pi]$ generated by the projection on the first variable. Note that

$$\mathbf{E}_\mu[f|p^1](\xi, \alpha) = \frac{1}{2\pi} \int_{[0, 2\pi]} f(\xi, \alpha) d\alpha \quad \text{for } \mu \text{ a.a. } (\xi, \alpha) \in E,$$

i.e., $\mathbf{E}_\mu[f|p^1]$ has a version which is independent of $\alpha \in [0, 2\pi]$. Each $f \in L^2(\mu)$ we can write as

$$(5.1) \quad f = (f - \mathbf{E}_\mu[f|p^1]) + \mathbf{E}_\mu[f|p^1],$$

where

$$(5.2) \quad \mathbf{E}_\mu[(f - \mathbf{E}_\mu[f|p^1])|p^1] = 0.$$

Hence, by the triangle inequality, it is sufficient to consider the following two types of functions: (i) $f \in L^2(\mu)$ with $\mathbf{E}_\mu[f|p^1] = 0$. (ii) $\mathbf{E}_\mu[f|p^1]$ with $f \in L^2(\mu)$.

Denote by $(S, D(S))$ the closure of (S, D) . Of course, $(S, D(S))$ is self-adjoint.

PROPOSITION 5.1. *Let the assumptions of Theorem 3.2 hold and let $f \in D(S)$. Then*

$$\mathbf{E}_{\mathbf{P}_\mu} \left[\left(\frac{1}{t} \int_0^t S f(\mathbf{X}_s) ds \right)^2 \right] \leq \frac{\sigma^2}{t} \|\partial_\alpha f\|_{L^2(\mu)}^2, \quad t > 0.$$

Proof. It is sufficient to consider $f \in D$. We fix $t > 0$ and below we canonically project the laws of the equilibrium processes \mathbf{P}_μ onto $C([0, t], E)$ without expressing this explicitly. We define the time reversal $r_t(\omega) := \omega(t - \cdot)$, $\omega \in C([0, t], E)$.

Using that

$$S = \frac{1}{2}(L + L^*) \quad \text{on } D,$$

we obtain

$$- \int_0^t Sf(\mathbf{X}_s) ds = \frac{1}{2}(\mathbf{M}(f)_t - \mathbf{R}^*(f)_0),$$

where $\mathbf{M}(f)_t$ is the martingale as in (3.2) and

$$\mathbf{R}^*(f)_u := f(\mathbf{X}_t) - f(\mathbf{X}_u) + \int_u^t L^* f(\mathbf{X}_s) ds, \quad 0 \leq u \leq t.$$

Observe that

$$-\mathbf{R}^*(f)_{t-u} = \mathbf{M}^*(f)_u \circ r_t, \quad 0 \leq u \leq t,$$

where $\mathbf{M}^*(f)_t$ is the martingale as in (3.4). Hence

$$\mathbf{E}_{\mathbf{P}_\mu} [(\mathbf{R}^*(f)_0)^2] = \mathbf{E}_{\mathbf{P}_\mu^*} [(\mathbf{M}^*(f)_t)^2],$$

because time reversal gives the adjoint process. Now by the Burkholder–Gundy–Davis inequality and (3.3), (3.5) we obtain

$$\begin{aligned} \mathbf{E}_{\mathbf{P}_\mu} \left[\left(\frac{1}{t} \int_0^t Sf(\mathbf{X}_s) ds \right)^2 \right]^{1/2} &\leq \frac{1}{2t} \mathbf{E}_{\mathbf{P}_\mu} [(\mathbf{M}(f)_t)^2]^{1/2} + \frac{1}{2t} \mathbf{E}_{\mathbf{P}_\mu^*} [(\mathbf{M}^*(f)_t)^2]^{1/2} \\ &\leq \frac{1}{2t} \mathbf{E}_{\mathbf{P}_\mu} [\langle \mathbf{M}(f) \rangle_t]^{1/2} + \frac{1}{2t} \mathbf{E}_{\mathbf{P}_\mu^*} [\langle \mathbf{M}^*(f) \rangle_t]^{1/2} \leq \frac{\sigma}{t^{1/2}} \|\partial_\alpha f\|_{L^2(\mu)}. \quad \square \end{aligned}$$

COROLLARY 5.2. *Let the assumptions of Theorem 3.2 hold and let $f \in L^2(\mu)$ with $\mathbf{E}_\mu[f|p^1] = 0$. Then*

$$\mathbf{E}_{\mathbf{P}_\mu} \left[\left(\frac{1}{t} \int_0^t f(\mathbf{X}_s) ds \right)^2 \right] \leq \frac{4}{\sigma^2 t} \|f\|_{L^2(\mu)}^2, \quad t > 0.$$

Proof. Since $\mathbf{E}_\mu[f|p^1] = 0$, f is in the orthogonal complement of the kernel of $(S, D(S))$. But the orthogonal complement of the kernel of $(S, D(S))$ is the range of $(S, D(S))$, because $(S, D(S))$ is self-adjoint. Hence there exists $S^{-1}f \in D(S)$. Then Proposition 5.1 yields

$$\mathbf{E}_{\mathbf{P}_\mu} \left[\left(\frac{1}{t} \int_0^t f(\mathbf{X}_s) ds \right)^2 \right] \leq \frac{\sigma^2}{t} \|\partial_\alpha S^{-1}f\|_{L^2(\mu)}^2 \leq \frac{4}{\sigma^2 t} \|f\|_{L^2(\mu)}^2,$$

because S^{-1} is bounded by $2\sigma^{-2}$ on the range of $(S, D(S))$. \square

THEOREM 5.3. *Let $\phi \in C^3(\mathbb{R}^2)$, $\nabla\phi \in L^2(\gamma)$, and assume there exist $0 < C < \infty$ and $K \subset \mathbb{R}^2$ compact such that*

$$\begin{aligned} &(\partial_{\xi_i} \partial_{\xi_j} \partial_{\xi_k} \phi(\xi))^2 + (\partial_{\xi_i} \partial_{\xi_j} \phi(\xi))^2 \\ &\leq C(\partial_{\xi_1} \phi(\xi))^2 + (\partial_{\xi_2} \phi(\xi))^2 \quad \text{for all } i, j, k \in \{1, 2\}, \xi \in \mathbb{R}^2 \setminus K. \end{aligned}$$

Furthermore assume that the Dirichlet form $(\mathcal{E}, D(\mathcal{E}))$ corresponding to the nondegenerated generator from section 4 fulfills a Poincaré inequality, i.e., there exists $0 < c < \infty$ such that

$$\mathcal{E}(f - \mathbf{E}_\gamma[f], f - \mathbf{E}_\gamma[f]) \geq c (f - \mathbf{E}_\gamma[f], f - \mathbf{E}_\gamma[f])_{L^2(\gamma)} \quad \text{for all } f \in D(\mathcal{E}).$$

Then

$$(5.3) \quad \left\| \frac{1}{t} \int_0^t f(\mathbf{X}_s) ds - \mathbf{E}_\mu[f] \right\|_{L^2(\mathbf{P}_\mu)} \leq \frac{1}{c^{1/2}} \left(\frac{2}{t} + \frac{1}{t^{1/2}} \left(\frac{A(C)c^{1/2} + B(C)c^{-1/2}}{\sigma} + (1 + 2^{-1/2})\sigma \right) \right) \|f - \mathbf{E}_\mu[f]\|_{L^2(\mu)}$$

for some constants $0 < A(C), B(C) < \infty$ independent of $c, \sigma > 0$, $f \in L^2(\mu)$, and $t > 0$.

Proof. Let $f \in L^2(\mu)$, and without loss of generality we assume $\mathbf{E}_\mu[f] = 0$. Since $(\mathcal{E}, D(\mathcal{E}))$ fulfills a Poincaré inequality, there exists $g \in D(G)$ such that

$$Gg = \mathbf{E}_\mu[f|p^1].$$

Now consider

$$h := \nabla_\xi g \cdot \tau.$$

Then $h \in \overline{D}^L$ by Lemma 4.1 and

$$\begin{aligned} Lh - \mathbf{E}_\mu[f|p^1] &= Lh - Gg = -\frac{\sigma^2}{2} \left(\cos(\alpha)\partial_{\xi^1}g + \sin(\alpha)\partial_{\xi^2}g \right) \\ &\quad + 2\sin(\alpha)\cos(\alpha)\partial_{\xi^2}\partial_{\xi^1}g + \left(\cos(\alpha)^2 - \frac{1}{2} \right) \partial_{\xi^1}^2g + \left(\sin(\alpha)^2 - \frac{1}{2} \right) \partial_{\xi^2}^2g \\ &\quad - \partial_{\xi^1}\phi(\xi) \left(\left(\sin(\alpha)^2 - \frac{1}{2} \right) \partial_{\xi^1}g - \cos(\alpha)\sin(\alpha)\partial_{\xi^2}g \right) \\ &\quad - \partial_{\xi^2}\phi(\xi) \left(\left(\cos(\alpha)^2 - \frac{1}{2} \right) \partial_{\xi^2}g - \cos(\alpha)\sin(\alpha)\partial_{\xi^1}g \right). \end{aligned}$$

It is easy to check that

$$\mathbf{E}_\mu[Lh - \mathbf{E}_\mu[f|p^1]|p^1] = 0.$$

Hence $Lh - \mathbf{E}_\mu[f|p^1]$ can be treated by Corollary 5.2. Thus by (5.1), (5.2) it is only left to consider Lh . Using (3.2), the Burkholder–Gundy–Davis inequality, and (3.3) we obtain

$$(5.4) \quad \begin{aligned} \left\| \frac{1}{t} \int_0^t Lh(\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} &\leq \frac{1}{t} \|h(\mathbf{X}_t) - h(\mathbf{X}_0)\|_{L^2(\mathbf{P}_\mu)} + \frac{1}{t} \|\mathbf{M}_t^h\|_{L^2(\mathbf{P}_\mu)} \leq \frac{2}{t} \|h\|_{L^2(\mu)} \\ &\quad + \frac{1}{t} \mathbf{E}_{\mathbf{P}_\mu} [\langle \mathbf{M}^h \rangle_t]^{1/2} \leq \frac{2}{t} \|\nabla_\xi g\|_{L^2(\gamma)} + \frac{\sigma}{t^{1/2}} \|\partial_\alpha h\|_{L^2(\mu)} \leq \frac{2 + \sigma t^{1/2}}{t} \|\nabla_\xi g\|_{L^2(\gamma)} \\ &= \frac{2 + \sigma t^{1/2}}{t} (G^{-1}\mathbf{E}_\mu[f|p^1], \mathbf{E}_\mu[f|p^1])_{L^2(\gamma)}^{1/2} \leq \frac{(2 + \sigma t^{1/2})}{c^{1/2}t} \|f\|_{L^2(\mu)}, \quad t > 0, \end{aligned}$$

where $0 < c < \infty$ is the spectral gap of G .

Now we obtain by Corollary 5.2 and (5.4) for $t > 0$

$$\begin{aligned} & \left\| \frac{1}{t} \int_0^t f(\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} \\ & \leq \left\| \frac{1}{t} \int_0^t (f - \mathbf{E}_\mu[f|p^1])(\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} + \left\| \frac{1}{t} \int_0^t \mathbf{E}_\mu[f|p^1](\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} \\ & \leq \frac{2}{\sigma t^{1/2}} \|f\|_{L^2(\mu)} + \left\| \frac{1}{t} \int_0^t (\mathbf{E}_\mu[f|p^1] - Lh)(\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} + \left\| \frac{1}{t} \int_0^t Lh(\mathbf{X}_s) ds \right\|_{L^2(\mathbf{P}_\mu)} \\ & \leq \frac{2}{\sigma t^{1/2}} (\|f\|_{L^2(\mu)} + \|Lh\|_{L^2(\mu)}) + \frac{(2 + \sigma t^{1/2})}{c^{1/2}t} \|f\|_{L^2(\mu)}. \end{aligned}$$

To obtain the last inequality we used the fact that

$$\|Lh - \mathbf{E}_\mu[f|p^1]\|_{L^2(\mu)}^2 = (Lh - Gg, Lh - Gg)_{L^2(\mu)} = (Lh, Lh)_{L^2(\mu)} - (Gg, Gg)_{L^2(\mu)},$$

because

$$(Lh, Gg)_{L^2(\mu)} = (Gg, Gg)_{L^2(\mu)}.$$

Finally, using (4.1) and the Kato bound provided in Lemma 4.1, we get the estimate (5.3). \square

Remark 5.4. The assumptions in Theorem 5.3 also allow potentials of the form $\phi = \|\cdot\|^p$, $p = 2, 4$, or $p \geq 6$, as in Remark 4.2, since in these cases also a Poincaré inequality holds. More generally, under our assumptions on the potential a Poincaré inequality holds if ϕ grows as fast as or faster than $\|x\|$ for large $x \in \mathbb{R}^2$; see [RW01].

6. Numerical Results. For a numerical investigation of the rate of convergence to equilibrium for (2.3) we use a semi-Lagrangian method. The method consists of two fractional steps. The first step is the Lagrangian interpretation for the advection part of (2.3) by the modified method of characteristics, while the second step uses Eulerian coordinates for the discretization of the reaction-diffusion part in (2.3). To preserve mass we use the modified method of characteristics with adjusted advection; see [DHP99]. For details on the numerical method we refer to [KRS07]. We consider the cases $\phi(\xi) = \|\xi\|$ and $\phi(\xi) = \|\xi\|^2/2$. We concentrate on the behavior of the relative entropy. One could also consider other observables. However, since convergence of the relative entropy is stronger than convergence of the distribution functions or convergence of time averages, we restrict ourselves to the investigation of the time development of the relative entropy. Figures 3 and 4 show log-plots of the time development of the relative entropy $H(u/m)(t)$ for different values of σ , and Figure 5 shows plots of the inverse decay rate

$$\lim_{T \rightarrow \infty} \int_0^T H(u/m)(t) dt$$

for various values of σ . The quantitative behavior depends on the initial values chosen; however, for other initial values qualitatively similar results are obtained. From Figures 3 and 4 an exponential decay seems possible. Moreover, one observes from Figure 5 that the optimal rate of convergence is found for a finite value of σ . On the other hand, considering the analytical estimate (5.3) one notes that the speed of convergence is dominated by the factor

$$\frac{1}{c^{1/2}t^{1/2}} \left(\frac{A(C)c^{1/2} + B(C)c^{-1/2}}{\sigma} + (1 + 2^{-1/2})\sigma \right).$$

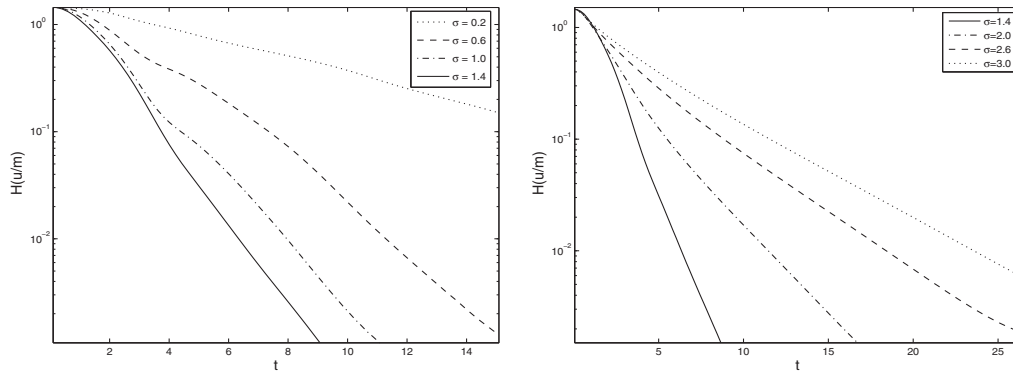


FIG. 3. $H(u/m)(t)$ for $\phi(\xi) = \|\xi\|$.

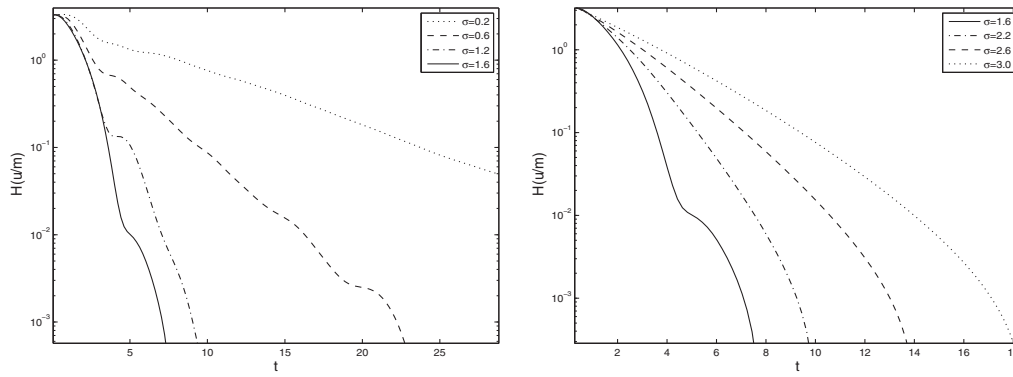


FIG. 4. $H(u/m)(t)$ for $\phi(\xi) = \|\xi\|^2/2$.

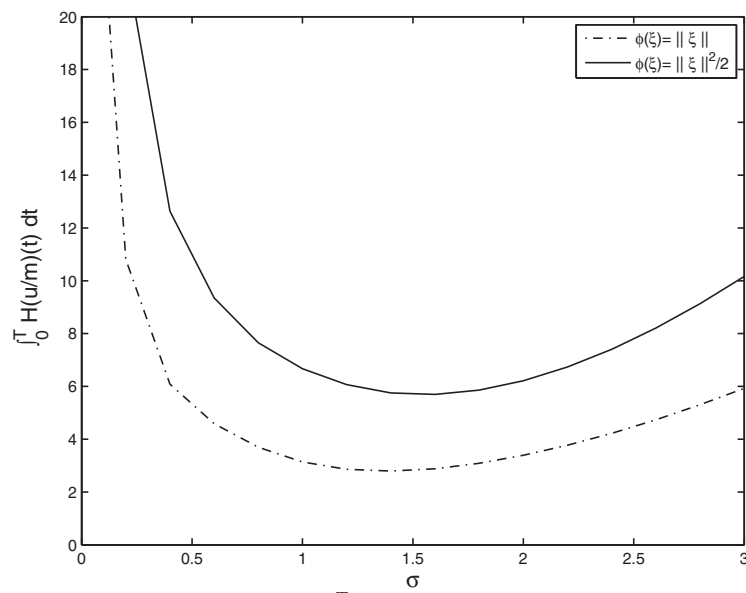


FIG. 5. Rates of convergence $\int_0^T H(u/m)(t)dt$ for different values of σ .

Obviously, a minimal rate of convergence can be determined—as from the numerical results—for a finite value of σ . Moreover, these observations fit the Monte Carlo results shown in Figure 2. From a practical point of view, this means that the process parameters have to be adapted such that σ is in an intermediate range of values to obtain the fastest possible decay to equilibrium and a fiber web which is as uniform as possible.

7. Concluding remarks. A stochastic model for fiber lay-down processes has been investigated analytically. Existence and ergodicity of the process have been shown and estimates on the rates of convergence presented. The results are supported by numerical simulations. According to these results the fastest decay to equilibrium is obtained for process parameters giving a diffusion coefficient σ in an intermediate range of values. We plan to extend this work to problems with moving conveyor belts; see [BGK+07]. In this case an additional difficulty is given by the fact that the stationary solution is not known explicitly.

Acknowledgment. We are grateful to L. Bonilla, F. Conrad, N. Marheineke, M. Seaid, and C. Villani as well as to the unknown referees for helpful discussions.

REFERENCES

- [AMTU01] A. ARNOLD, P. MARKOWICH, G. TOSCANI, AND A. UNTERREITER, *On convex Sobolev inequalities and the rate of convergence to equilibrium for Fokker–Planck type equations*, Comm. Partial Differential Equations, 26 (2001), pp. 43–100.
- [BKR97] V. I. BOGACHEV, N. V. KRYLOV, AND M. RÖCKNER, *Elliptic regularity and essential self-adjointness of Dirichlet operators on \mathbb{R}^n* , Ann. Sc. Norm. Super. Pisa Cl. Sci. (4), 24 (1997), pp. 451–461.
- [BGK+07] L. BONILLA, T. GÖTZ, A. KLAR, N. MARHEINEKE, AND R. WEGENER, *Hydrodynamic limit of the Fokker–Planck equation describing fiber lay-down processes*, SIAM J. Appl. Math., 68 (2007), pp. 648–665.
- [CG08a] F. CONRAD AND M. GROTHAUS, *Construction of N -particle Langevin dynamics for $H^{1,\infty}$ -potentials via generalized Dirichlet forms*, Potential Anal., 28 (2008), pp. 261–282.
- [CG08b] F. CONRAD AND M. GROTHAUS, *Construction, Ergodicity and Rate of Convergence of N -particle Langevin Dynamics with Singular Potentials*, preprint.
- [DPZ96] G. DA PRATO AND J. ZABCZYK, *Ergodicity for Infinite Dimensional Systems*, London Math. Soc. Lecture Notes Ser. 229, Cambridge University Press, Cambridge, UK, 1996.
- [DV01] L. DESVILETTES AND C. VILLANI, *On the trend to global equilibrium for spatially inhomogeneous entropy-dissipating systems: The linear Fokker–Planck equation*, Comm. Pure Appl. Math., 54 (2001), pp. 1–42.
- [DHP99] J. DOUGLAS, C. HUANG, AND F. PEREIRA, *The modified method of characteristics with adjusted advection*, Numer. Math., 83 (1999), pp. 353–369.
- [Gol85] J. A. GOLDSTEIN, *Semigroups of Linear Operators and Applications*, Oxford Math. Monogr., Oxford University Press, New York, 1985.
- [GKM+07] T. GÖTZ, A. KLAR, N. MARHEINEKE, AND R. WEGENER, *A stochastic model and associated Fokker–Planck equation for the fiber lay-down process in nonwoven production processes*, SIAM J. Appl. Math., 67 (2007), pp. 1704–1717.
- [KRS07] A. KLAR, P. REUTERSWÄRD, AND M. SEĀĪD, *A Semi-Lagrangian Method for the Fokker–Planck Equation of Fiber Dynamics*, preprint.
- [MW06] N. MARHEINEKE AND R. WEGENER, *Fiber dynamics in turbulent flows: General modeling framework*, SIAM J. Appl. Math., 66 (2006), pp. 1703–1726.
- [OT03] S. OLLA AND C. TREMOULET, *Equilibrium fluctuations for interacting Ornstein–Uhlenbeck particles*, Comm. Math. Phys., 233 (2003), pp. 463–491.
- [RW01] M. RÖCKNER AND F.-Y. WANG, *Weak Poincaré inequalities and L^2 -convergence rates of Markov semigroups*, J. Funct. Anal., 185 (2001), pp. 564–603.
- [Sta99] W. STANNAT, *The theory of generalized Dirichlet forms and its applications in analysis and stochasticity*, Mem. Amer. Math. Soc., 142 (678) (1999).
- [Vil06] C. VILLANI, *Hypoocoercivity*, preprint.

DETERMINATION OF A LINEAR CRACK IN AN ELASTIC BODY FROM BOUNDARY MEASUREMENTS—LIPSCHITZ STABILITY*

ELENA BERETTA[†], ELISA FRANCONI[‡], AND SERGIO VESSELLA[§]

Abstract. We discuss the stability issue for the problem of determining a thin inclusion in a homogeneous isotropic elastic body from boundary measurements. This problem is severely ill-posed, but in this paper we prove that by restricting ourselves to the class of thin neighborhoods of line segments, a Lipschitz stability estimate holds.

Key words. inverse problems, linear elasticity, cracks

AMS subject classification. 35R30

DOI. 10.1137/070698397

1. Introduction. Let $\Omega \subset \mathbb{R}^2$ be a bounded connected domain, with a sufficiently smooth boundary, representing the region occupied by an elastic, homogeneous, isotropic material.

Let $\sigma \subset \Omega$ be a simple smooth curve and define, for a positive small ϵ , the set

$$\omega_\epsilon = \{x \in \Omega : d(x, \sigma) < \epsilon\},$$

which represents an inclusion of small size made of a different elastic material.

Let \mathbb{C}_0 and \mathbb{C}_1 be the elastic tensor fields in $\Omega \setminus \bar{\omega}_\epsilon$ and ω_ϵ , respectively.

Given a traction field g on $\partial\Omega$, the displacement field u_ϵ , generated by this traction in the body containing the inclusion ω_ϵ , solves the following system of linearized elasticity:

$$(1) \quad \begin{cases} \operatorname{div}(\mathbb{C}_\epsilon \widehat{\nabla} u_\epsilon) = 0 & \text{in } \Omega, \\ (\mathbb{C}_\epsilon \widehat{\nabla} u_\epsilon)\nu = g & \text{on } \partial\Omega, \end{cases}$$

where $\mathbb{C}_\epsilon = \mathbb{C}_0 \chi_{\Omega \setminus \omega_\epsilon} + \mathbb{C}_1 \chi_{\omega_\epsilon}$, $\widehat{\nabla} u_\epsilon = \frac{1}{2}(\nabla u_\epsilon + (\nabla u_\epsilon)^T)$ is the symmetric deformation tensor and ν denotes the outward unit normal to $\partial\Omega$.

Let us also introduce the background displacement field u_0 generated by the same traction g in the absence of the inclusion, namely, the solution of

$$(2) \quad \begin{cases} \operatorname{div}(\mathbb{C}_0 \widehat{\nabla} u_0) = 0 & \text{in } \Omega, \\ (\mathbb{C}_0 \widehat{\nabla} u_0)\nu = g & \text{on } \partial\Omega. \end{cases}$$

*Received by the editors July 25, 2007; accepted for publication (in revised form) May 29, 2008; published electronically September 19, 2008. This work was supported by MIUR, PRIN, grant 2006014115.

<http://www.siam.org/journals/sima/40-3/69839.html>

[†]Dipartimento di Matematica “G. Castelnuovo,” Università di Roma “La Sapienza,” Piazzale Aldo Moro 5, 00185 Roma, Italy (beretta@mat.uniroma1.it).

[‡]Dipartimento di Matematica “U. Dini,” Università degli Studi di Firenze, Viale Morgagni 67A, 50134 Firenze, Italy (francini@math.unifi.it).

[§]Dipartimento di Matematica per le Decisioni, Università degli Studi di Firenze, Via C. Lombroso, 6/17, 50134 Firenze, Italy (sergio.vessella@dmd.unifi.it).

In [BF06] the following asymptotic expansion for $(u_\epsilon - u_0)|_{\partial\Omega}$ as $\epsilon \rightarrow 0$ has been derived:

$$(3) \quad (u_\epsilon - u_0)(y) = 2\epsilon \int_\sigma \mathcal{M}(x) \widehat{\nabla} u_0(x) \cdot \widehat{\nabla} N(x, y) d\sigma(x) + O(\epsilon^{1+\theta}) \quad \text{for } y \in \partial\Omega,$$

where $N(x, y)$ denotes the Neumann function for the operator $\text{div}(\mathbb{C}_0 \widehat{\nabla} \cdot)$ in Ω , \mathcal{M} is a fourth order symmetric tensor which depends on σ and on the elastic properties of the material occupied by the region Ω , and $\theta \in (0, 1)$ is independent of ϵ .

In the following, the symbol “ \cdot ” (as in formula (3)) will denote the usual scalar product between matrices ($A \cdot B = \sum_{ij} a_{ij} b_{ij}$) or vectors ($u \cdot v = \sum_i u_i v_i$).

Clearly, in order to get information on u_ϵ , we need to consider the first order term of the expansion (3). Let us consider

$$(4) \quad u_\sigma(y) = 2 \int_\sigma \mathcal{M}(x) \widehat{\nabla} u_0(x) \cdot \widehat{\nabla} N(x, y) d\sigma(x) \quad \text{for } y \in \overline{\Omega} \setminus \sigma.$$

The inverse problem we are interested in is the following:

Given the trace of the correction term u_σ on some open subset Γ of $\partial\Omega$, determine the curve σ .

Due to the nonlinearity, the problem is severely ill-posed and the reconstruction of an arbitrary smooth curve is expected to fail. On the other hand, restricting the class of admissible curves σ , we expect to regularize the problem. In fact, in the scalar conductivity case, if σ is restricted to be linear, it is possible to show a Lipschitz continuous dependence of the segment σ from the boundary data $u_\sigma|_{\partial\Omega}$ (cf. [ABF04, ABF06]).

In this paper we show that the same type of continuous dependence holds true in the elastic case. More precisely, in Theorem 2.2 we prove a Lipschitz continuous dependence of the segment σ from boundary data of the correction term u_σ . In order to do this we use the qualitative properties of the correction term and quantitative estimates of the unique continuation property for elliptic systems with constant coefficients. In Corollary 2.3 we show that the segment σ depends Lipschitz continuously on the rescaled boundary deviation of the solution u_ϵ .

Our method allows us to detect thin elastic inclusions from knowledge of one boundary measurement. In the book [AK04], the authors deal with a similar problem for diametrically small inclusions. We would like to point out that this analysis is on one side motivated by applications (breast imaging, mine detection), and, on the other side, it is not known if an arbitrary inclusion can be detected with a finite number of boundary measurements, not even in the scalar case.

The plan of the paper is the following: In section 2 we introduce some notation and the main assumptions, and we state the main result. In section 3 we describe some properties and the asymptotic behavior near the endpoints of the crack σ of the function u_σ . Section 4 contains the proof of the main result. The final section is devoted to the proof of some auxiliary lemmas and propositions stated in section 4.

2. The main result. We introduce some notation and assumptions that will be useful in what follows. For any $r > 0$ we denote

$$\Omega_r = \{x \in \Omega \mid d(x, \partial\Omega) > r\},$$

and by $\text{diam}\Omega$ the diameter of Ω . We also denote by $B_r(P)$ the open disc centered in P with radius r .

For any $P, Q \in \mathbb{R}^2$ we denote by $[P, Q]$ the segment with endpoints P and Q ; more precisely,

$$[P, Q] = \{P + t(Q - P) \mid t \in [0, 1]\}.$$

Given positive constants $\gamma, \rho_0, E_0, E_1, L, \gamma_0$ with $L \geq 1$ and $\gamma \in (0, 1]$, which we shall name a priori data,

(H1) we assume that Ω is a bounded simply connected domain of \mathbb{R}^2 such that

$$\text{diam } \Omega \leq E_1.$$

Concerning the regularity of $\partial\Omega$, we assume that

$$\partial\Omega \text{ is of class } C^{2,\gamma} \text{ with constants } \rho_0, E_0.$$

More precisely, for any point $\tilde{P} \in \partial\Omega$, there exists a rigid transformation of coordinates under which we have $\tilde{P} = 0$ and

$$\Omega \cap B_{\rho_0}(0) = \{x \in B_{\rho_0}(0) \mid x_2 > \psi(x_1)\},$$

where ψ is a $C^{2,\gamma}$ function on $(-\rho_0, \rho_0) \subset \mathbb{R}$ satisfying

$$\psi(0) = 0, \quad \psi'(0) = 0, \quad \text{and} \quad \|\psi\|_{C^{2,\gamma}(-\rho_0, \rho_0)} \leq E_0.$$

Moreover we assume that there exists an open subset Γ of $\partial\Omega$ such that for some point $\tilde{Q} \in \Gamma$,

$$\partial\Omega \cap B_{\rho_0}(\tilde{Q}) \subset \Gamma.$$

(H2) We assume that σ is a segment of endpoints P, Q satisfying

$$L^{-1} \leq |P - Q| \leq L \quad \text{and} \quad d(\sigma, \mathbb{R}^2 \setminus \Omega) \geq L^{-1}.$$

(H3) We will assume Ω and ω_ϵ are both homogeneous and isotropic, i.e., the elastic tensor fields \mathbb{C}_0 and \mathbb{C}_1 are of the form

$$(\mathbb{C}_m)_{ijkl} = \lambda_m \delta_{ij} \delta_{kl} + \mu_m (\delta_{ki} \delta_{lj} + \delta_{kj} \delta_{li}) \quad \text{for } i, j, k, l = 1, 2, \quad m = 0, 1,$$

where (λ_0, μ_0) and (λ_1, μ_1) are the Lamé coefficients corresponding to $\Omega \setminus \bar{\omega}_\epsilon$ and ω_ϵ , respectively, and satisfy the monotonicity condition

$$(\lambda_0 - \lambda_1)(\mu_0 - \mu_1) \geq \gamma_0$$

for some positive constant γ_0 .

(H4) There are two positive constants α_0, β_0 such that

$$\min(\mu_0, \mu_1) \geq \alpha_0, \quad \min(\lambda_0 + \mu_0, \lambda_1 + \mu_1) \geq \beta_0/2.$$

We note that hypothesis (H4) ensures that \mathbb{C}_ϵ is strongly convex in Ω ; i.e., if we set $\xi_0 = \min(2\alpha_0, \beta_0)$, then

$$\mathbb{C}_\epsilon A \cdot A \geq \xi_0 |A|^2$$

for any symmetric 2×2 matrix A .

(H5) We shall prescribe the traction field g on $\partial\Omega$ of the form

$$g = (\mathbb{C}_0 W)\nu,$$

where W is a nonzero symmetric 2×2 constant matrix.

Remark 2.1. Under assumptions (H3), (H4) there exist weak solutions u_ϵ and u_0 in $H^1(\Omega, \mathbb{R}^2)$ of (1) and of (2), respectively (see, for example, [F72]). In particular u_0 can be explicitly calculated and corresponds to a pure strain displacement

$$u_0(x) = Wx + c,$$

where c is an arbitrary constant vector. We observe that choosing g as in (H5) is not too restrictive since a pure strain can always be accomplished either by simple extensions in perpendicular directions or by a uniform dilation followed by an isochoric pure strain [G72].

In order to uniquely determine u_ϵ and u_0 , we assume that they satisfy the following normalization conditions:

$$(5) \quad \int_{\partial\Omega} u = 0, \quad \int_{\Omega} \nabla u - (\nabla u)^T = 0.$$

For $y \in \Omega$, we will denote by $N(\cdot, y)$ the Neumann function related to Ω , i.e., the weak solution to the problem

$$\begin{cases} \operatorname{div}(\mathbb{C}_0 \widehat{\nabla} N(\cdot, y)) = -\delta_y \operatorname{Id} & \text{in } \Omega, \\ (\mathbb{C}_0 \widehat{\nabla} N(\cdot, y)) \cdot \nu = -\frac{1}{|\partial\Omega|} \operatorname{Id} & \text{on } \partial\Omega, \end{cases}$$

with the normalization conditions (5). Here Id is the identity matrix in \mathbb{R}^2 . Note that by well-known regularity results for elliptic systems (cf. [C80]) we have

$$(6) \quad N(x, y) = \Gamma(x - y) + w(x, y),$$

where w is a smooth function of x and y and $\Gamma(x, y)$ is the fundamental free space solution of $\operatorname{div}(\mathbb{C}_0 \widehat{\nabla} \cdot)$, given by

$$(7) \quad \Gamma_{ij}(x) = \frac{A}{2\pi} \delta_{ij} \log|x| - \frac{B}{2\pi} \frac{x_i x_j}{|x|^2}, \quad i, j = 1, 2,$$

where $A = \frac{1}{2}(\frac{1}{\mu_0} + \frac{1}{\lambda_0 + 2\mu_0})$ and $B = \frac{1}{2}(\frac{1}{\mu_0} - \frac{1}{\lambda_0 + 2\mu_0})$.

Let us fix an orthonormal system (n, τ) on σ such that n is a unit normal vector field to the segment and τ is a unit tangent vector field.

By Theorem 2.1 of [BF06] we have that, for any $y \in \partial\Omega$, the first order term u_σ of the expansion of $u_\epsilon - u_0$ can be computed explicitly up to a remainder term and has the form

$$(8) \quad u_\sigma(y) = 2 \int_{\sigma} \mathcal{M}W \cdot \widehat{\nabla} N(x, y) d\sigma(x),$$

where

$$\mathcal{M}W = a(\operatorname{tr}W)\operatorname{Id} + bW + c(\tau^T W \tau)\tau \otimes \tau + d(n^T W n)n \otimes n.$$

Here $\text{tr}W$ denotes the trace of the matrix W and the coefficients are given by

$$(9) \quad a = (\lambda_1 - \lambda_0) \frac{\lambda_0 + 2\mu_0}{\lambda_1 + 2\mu_1}, \quad b = 2(\mu_1 - \mu_0) \frac{\mu_0}{\mu_1},$$

$$(10) \quad c = 2(\mu_1 - \mu_0) \left[\left(\frac{2\lambda_1 + 2\mu_1 - \lambda_0}{\lambda_1 + 2\mu_1} - \frac{\mu_0}{\mu_1} \right) \right],$$

and

$$(11) \quad d = 2(\mu_1 - \mu_0) \frac{\mu_1 \lambda_0 - \mu_0 \lambda_1}{\mu_1 (\lambda_1 + 2\mu_1)}.$$

We are now ready to state our main result.

THEOREM 2.2. *Let σ_0 and σ_1 be two segments satisfying assumption (H2). Assume (H1), (H3), (H4), and (H5). Let u_{σ_0} and u_{σ_1} be the functions defined by (8) where σ is replaced by σ_0 and σ_1 , respectively. Then there exists a constant C , depending only on the a priori data, such that*

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq C \|u_{\sigma_0} - u_{\sigma_1}\|_{L^2(\Gamma)}.$$

Here $d_{\mathcal{H}}$ denotes the Hausdorff distance

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) = \min\{\max\{|P_0 - P_1|, |Q_0 - Q_1|\}, \max\{|P_0 - Q_1|, |Q_0 - P_1|\}\},$$

where P_0, Q_0 and P_1, Q_1 are the endpoints of σ_0 and σ_1 , respectively.

The proof of Theorem 2.2 is postponed to section 4.

A straightforward consequence of Theorem 2.2 and of the asymptotic formula (3) is the following.

COROLLARY 2.3. *Let σ_0 and σ_1 be as in Theorem 2.2. For $i = 0, 1$, let u_{ϵ}^i be the solution of (1) corresponding to $\omega_{\epsilon}^i = \{x \in \Omega : d(x, \sigma_i) < \epsilon\}$. Then there exist a positive constant C and $\theta \in (0, 1)$, depending only on the a priori data, such that*

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq C (\epsilon^{-1} \|u_{\epsilon}^0 - u_{\epsilon}^1\|_{L^2(\Gamma)} + \epsilon^{\theta}).$$

Remark 2.4. A similar rescaled stability estimate has been obtained, for example, in [FV89] and in [CMV98] for the case of diametrically small inclusions.

3. Properties of the function u_{σ} . In this section we first establish a representation formula for the function u_{σ} (defined in (4)), and then we deduce some results concerning the regularity of this function and its behavior near the endpoints of σ .

LEMMA 3.1. *Let $\sigma = [P, Q]$, and fix $\tau = \frac{Q-P}{|Q-P|}$ and let $n = \tau^{\perp}$. For $y \in \bar{\Omega} \setminus \sigma$,*

$$(12) \quad u_{\sigma}(y) = \int_{\sigma} \mathbb{C}_0 \widehat{\nabla} N(x, y) n \cdot \varphi_{\sigma} d\sigma(x) + (N(Q, y) \cdot \tau - N(P, y) \cdot \tau) f_{\sigma},$$

where φ_{σ} is the vector whose components in the (n, τ) directions are given by

$$(13) \quad \begin{aligned} \varphi_{\sigma} \cdot n &= \frac{2}{\lambda_1 + 2\mu_1} ((\lambda_1 - \lambda_0) \text{tr}W + 2(\mu_1 - \mu_0) n^T W n), \\ \varphi_{\sigma} \cdot \tau &= \frac{4(\mu_1 - \mu_0)}{\mu_1} n^T W \tau \end{aligned}$$

and f_σ is the constant function

$$(14) \quad f_\sigma = \alpha \operatorname{tr}W - \beta n^T W n,$$

where

$$\alpha = \frac{4}{\lambda_1 + 2\mu_1}((\mu_1 - \mu_0)(\lambda_1 + 2\mu_1) + \mu_1(\lambda_1 - \lambda_0)), \quad \beta = \frac{8(\mu_1 - \mu_0)}{\lambda_1 + 2\mu_1}(\lambda_1 + \mu_1).$$

Proof. By straightforward calculations and by the symmetry of N and W , one can easily check that

$$(15) \quad u_\sigma(y) = \int_\sigma \mathbb{C}_0 \widehat{\nabla} N(x, y) n \cdot \varphi_\sigma \, d\sigma(x) + 2 \int_\sigma \tau^T \widehat{\nabla} N(x, y) \tau f_\sigma \, d\sigma(x).$$

Finally an integration by parts in (15) gives the thesis. \square

Let us consider the double layer potential

$$\mathcal{D}_\sigma \varphi(y) = \int_\sigma \mathbb{C}_0 \widehat{\nabla} \Gamma(x, y) n \cdot \varphi \, d\sigma(x), \quad y \in \mathbb{R}^2 \setminus \sigma,$$

with φ constant vector field. Let us point out the following properties.

LEMMA 3.2. *We have*

$$(16) \quad [\mathcal{D}_\sigma \varphi]_\sigma = \varphi,$$

and there exists a positive constant C such that, for any $y \in \mathbb{R}^2 \setminus \sigma$,

$$|\mathcal{D}_\sigma \varphi(y)| \leq C \left(\left| \log \frac{|Q - y|}{|P - y|} \right| + 1 \right) |\varphi|,$$

$$|\nabla \mathcal{D}_\sigma \varphi(y)| \leq C \left(\frac{1}{|P - y|} + \frac{1}{|Q - y|} + 1 \right) |\varphi|,$$

where $[\mathcal{D}_\sigma \varphi]_\sigma$ denotes the jump of $\mathcal{D}_\sigma \varphi$ across σ .

Proof. Formula (16) follows from standard properties of double layer potentials (see, for example, [AK04]). The behavior of $\mathcal{D}_\sigma \varphi$ at the endpoints of segment σ follows from straightforward computations. \square

We are now ready to state and prove the following.

PROPOSITION 3.3. *Under assumptions (H1)–(H5), u_σ satisfies*

$$(17) \quad \begin{cases} \operatorname{div}(\mathbb{C}_0 \widehat{\nabla} u_\sigma) = 0 & \text{in } \Omega \setminus \sigma, \\ (\mathbb{C}_0 \widehat{\nabla} u_\sigma) \nu = 0 & \text{on } \partial\Omega, \end{cases}$$

and the normalization condition (5).

The function u_σ has a jump on σ given by

$$(18) \quad [u_\sigma]_\sigma = \varphi_\sigma.$$

Moreover there exists a positive constant C , depending only on the a priori data, such that

$$(19) \quad |u_\sigma(y)| \leq C \left(\left| \log \frac{|Q - y|}{|P - y|} \right| + 1 \right).$$

Proof. From the representation formula (12) and by the properties of the function N it is easy to see that u_σ satisfies (17). Recalling that u_σ is the correction term in the expansion of $u_\epsilon - u_0$ and that u_ϵ and u_0 have the same Neumann condition on $\partial\Omega$, and observing that the expansion (3) holds true in a neighborhood of $\partial\Omega$, it follows that u_σ satisfies the homogeneous Neumann datum.

Observe now that u_σ is expressed in terms of a double layer potential and the Neumann function at the endpoints of σ . From (6) and Lemma 3.2, equations (18) and (19) follow. \square

4. Proof of Theorem 2.2. We start by establishing a first rough estimate of $d_{\mathcal{H}}(\sigma_0, \sigma_1)$ in terms of the boundary data error.

LEMMA 4.1. *Let $\sigma_0 = [P_0, Q_0]$ and $\sigma_1 = [P_1, Q_1]$ be two segments satisfying assumption (H2). Assume (H1), (H3), (H4), and (H5). Let $u_{\sigma_0}, u_{\sigma_1}$ be the functions defined by (12) and corresponding to σ_0 and σ_1 , respectively. Let $\epsilon := \|u_{\sigma_0} - u_{\sigma_1}\|_{L^2(\Gamma)}$; then*

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq \omega_1(\epsilon),$$

where $\omega_1(s) = C (\log |\log(\epsilon)|)^{-1/4}$ for some positive constant C depending only on the a priori data.

Proof. We first describe the proof in the case $\epsilon = 0$ (corresponding to the uniqueness).

Let $u := u_{\sigma_0} - u_{\sigma_1}$. Recalling that $(\mathbb{C}_0 \widehat{\nabla} u_{\sigma_0}) \cdot \nu = (\mathbb{C}_0 \widehat{\nabla} u_{\sigma_1}) \cdot \nu = 0$ on $\partial\Omega$, we have that

$$u = (\mathbb{C}_0 \widehat{\nabla} u) \cdot \nu = 0 \quad \text{on} \quad \partial\Omega.$$

Since u is solution of

$$\operatorname{div} \left(\mathbb{C}_0 \widehat{\nabla} u \right) = 0 \quad \text{in} \quad \Omega \setminus (\sigma_0 \cup \sigma_1),$$

by the unique continuation property of solutions to elliptic systems with constant coefficients (see [AM01] and [MR04]) we have that

$$(20) \quad u = 0 \quad \text{in} \quad \Omega \setminus (\sigma_0 \cup \sigma_1).$$

Assume that $\sigma_0 \neq \sigma_1$. Then there exists a portion Σ , for instance, of σ_0 , containing an endpoint P_0 of σ_0 such that Σ is not contained in σ_1 . For the sake of simplicity, let us denote $\varphi_0 := \varphi_{\sigma_0}$ and $f_0 := f_{\sigma_0}$ the functions appearing in (12) corresponding to σ_0 . By (20), $[u]_\Sigma = 0$, and since $[u]_\Sigma = [u_{\sigma_0}]_\Sigma = \varphi_0$ we get

$$(21) \quad \varphi_0 = 0.$$

Furthermore, from the representation formula (12), u_{σ_0} has logarithmic singularities at the endpoints P_0, Q_0 . Hence again by (20) and by (12) we have

$$(22) \quad f_0 = 0.$$

Recalling (13) and (14), conditions (21) and (22) can be written as

$$(23) \quad \begin{cases} n_0^T W \tau_0 = 0, \\ \alpha(\operatorname{tr} W) - \beta n_0^T W n_0 = 0, \\ a(\operatorname{tr} W) + (b + d)n_0^T W n_0 = 0, \end{cases}$$

where n_0 and τ_0 denote the normal and tangent unit vector to σ_0 , the numbers a, b , and d are defined in (9), (10), (11), and α and β are defined in (14). By assumption (H3) it is easy to see that the conditions (23) imply $\text{tr}W = 0$, $n_0^T W n_0 = 0$, from which it follows that $\tau_0^T W \tau_0 = 0$. Finally, by the symmetry of W , since $n_0^T W \tau_0 = \tau_0^T W n_0$ we get that also $\tau_0^T W n_0 = 0$. Hence $W = 0$, which contradicts assumption (H5).

The proof of the case $\varepsilon > 0$ is a quantitative version of the uniqueness part. For the sake of simplicity we postpone this part of the proof to section 5. \square

Denote by σ_t the segment $[P_t, Q_t]$, where

$$P_t = (1 - t)P_0 + tP_1, \quad Q_t = (1 - t)Q_0 + tQ_1,$$

and for $y \in \Omega \setminus \sigma_t$ let

$$u_t(y) := u_{\sigma_t}(y) = \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1 - s)P_t + sQ_t, y) n_t \cdot \varphi_t |Q_t - P_t| ds \\ + (N(Q_t, y) \cdot \tau_t - N(P_t, y) \cdot \tau_t) f_t,$$

where $\tau_t = \frac{Q_t - P_t}{|Q_t - P_t|}$, $n_t = \tau_t^\perp$, $f_t := f_{\sigma_t}$, and $\varphi_t := \varphi_{\sigma_t}$. Clearly, from the above formula, u_t is differentiable for any $t \in [0, 1]$ and its derivative u'_t satisfies the problem

$$(24) \quad \begin{cases} \text{div} \left(\mathbb{C}_0 \widehat{\nabla} u'_t \right) = 0 & \text{in } \Omega \setminus \sigma_t, \\ \left(\mathbb{C}_0 \widehat{\nabla} u'_t \right) \nu = 0 & \text{on } \partial\Omega \end{cases}$$

and the normalization condition (5).

The following statements, whose proof we defer until section 5, allow us to conclude the proof of Theorem 2.2.

LEMMA 4.2. *For any $y \in \partial\Omega$*

$$(25) \quad \begin{aligned} u'_t(y) &:= \frac{d}{dt} u_t(y) \\ &= - \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1 - s)P_t + sQ_t, y) n_t \cdot \varphi_t ((Q_1 - Q_0 - (P_1 - P_0)) \cdot \tau_t) ds \\ &\quad + \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1 - s)P_t + sQ_t, y) n_t \cdot \frac{d}{dt} (\varphi_t |Q_t - P_t|) ds \\ &\quad + \mathbb{C}_0 \widehat{\nabla} N(Q_t, y) ((Q_1 - Q_0) \cdot \tau_t) n_t - (Q_1 - Q_0) \cdot n_t \tau_t \cdot \varphi_t \\ &\quad - \mathbb{C}_0 \widehat{\nabla} N(P_t, y) ((P_1 - P_0) \cdot \tau_t) n_t - (P_1 - P_0) \cdot n_t \tau_t \cdot \varphi_t \\ &\quad + \left(\nabla N(Q_t, y) \frac{dQ_t}{dt} \cdot \tau_t - \nabla N(P_t, y) \frac{dP_t}{dt} \cdot \tau_t \right) f_t \\ &\quad + (N(Q_t, y) - N(P_t, y)) \frac{d(\tau_t f_t)}{dt}. \end{aligned}$$

PROPOSITION 4.3. *There exist constants $C_0, m_0 > 0$, depending only on the a priori data, such that*

$$\|u'_0\|_{L^2(\Gamma)} \geq m_0 d_{\mathcal{H}}(\sigma_0, \sigma_1)$$

and

$$(26) \quad \|u'_t - u'_0\|_{L^2(\Gamma)} \leq C_0(d_{\mathcal{H}}(\sigma_0, \sigma_1))^2.$$

Let us conclude now the proof of Theorem 2.2.

Let $\varepsilon_0 > 0$ such that

$$(27) \quad \omega_1(\varepsilon_0) \leq \frac{m_0}{2C_0},$$

where ω_1 is defined as in Lemma 4.1.

Let us distinguish two cases:

(1) $\varepsilon \in (0, \varepsilon_0]$. By the differentiability of u_t we have that for any $x \in \Gamma$

$$u_{\sigma_1}(x) - u_{\sigma_0}(x) = \int_0^1 u'_t(x) dt = u'_0(x) + \int_0^1 (u'_t(x) - u'_0(x)) dt.$$

By the triangular inequality we get

$$\varepsilon = \|u_{\sigma_1} - u_{\sigma_0}\|_{L^2(\Gamma)} \geq \|u'_0\|_{L^2(\Gamma)} - \|u'_t - u'_0\|_{L^2(\Gamma)}.$$

By Proposition 4.3 we get

$$(28) \quad \varepsilon \geq (m_0 - C_0 d_{\mathcal{H}}(\sigma_0, \sigma_1)) d_{\mathcal{H}}(\sigma_0, \sigma_1).$$

On the other hand, from Lemma 4.1 we have that

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq \omega_1(\varepsilon_0).$$

Hence, by (27), we get

$$m_0 - C_0 d_{\mathcal{H}}(\sigma_0, \sigma_1) \geq m_0 - C_0 \omega_1(\varepsilon_0) \geq \frac{m_0}{2}.$$

By the last inequality and (28) we derive

$$(29) \quad d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq \frac{2}{m_0} \varepsilon.$$

(2) Let us consider the case $\varepsilon \geq \varepsilon_0$. One easily gets

$$(30) \quad d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq \text{diam } \Omega \leq E_1 \frac{\varepsilon}{\varepsilon_0}.$$

Hence, by (29) and (30) we have

$$d_{\mathcal{H}}(\sigma_0, \sigma_1) \leq \max \left\{ \frac{E_1}{\varepsilon_0}, \frac{2}{m_0} \right\} \varepsilon,$$

which concludes the proof of the theorem. \square

5. Proofs of the preliminary results.

5.1. Proof of Lemma 4.1. Let us now consider the case $\varepsilon > 0$. Let $d := d_{\mathcal{H}}(\sigma_0, \sigma_1)$ and let us assume that $d = |P_0 - P_1|$. By classical regularity results in

potential theory and by Lemma 3.2 we have that $u_{\sigma_j} \in C^1(\overline{\Omega} \setminus \sigma_j)$ for $j = 0, 1$ and, for any $r \in (0, \rho_1)$, where $\rho_1 = \min(L^{-1}, \rho_0)$,

$$(31) \quad \begin{aligned} \|u_{\sigma_j}\|_{L^\infty(\Omega \setminus (B_r(P_j) \cup B_r(Q_j)))} &\leq C \left(\log \frac{1}{r} + 1\right), \\ \|\nabla u_{\sigma_j}\|_{L^\infty(\Omega \setminus (\sigma_j \cup B_r(P_j) \cup B_r(Q_j)))} &\leq C \frac{1}{r}, \end{aligned}$$

where C depends only on the a priori data. From stability estimates for the Cauchy problem for solutions of elliptic systems (cf. [MR04]) and by (31) we have

$$(32) \quad \|u_{\sigma_j}\|_{L^\infty(\Omega \cap B_{\frac{\rho_1}{2}}(\overline{P}))} \leq C \varepsilon^\delta,$$

where \overline{P} is a point on Γ and $C > 0, \delta \in (0, 1)$ depend on the a priori data only. Starting from (32) we now establish a propagation error estimate in a neighborhood of $\Omega \setminus (\sigma_0 \cup \sigma_1)$. Let $r \in (0, \frac{\rho_1}{156})$, and define

$$(33) \quad (\sigma_i)_{8r} = \{x \in \Omega : \text{dist}(x, \sigma_i) > 8r\} \quad \text{for } i = 0, 1$$

and $\Omega'_{(r)} = \Omega_{8r} \setminus ((\sigma_0)_{8r}) \cup (\sigma_1)_{8r}$ and let $z_0 = \overline{P} - \frac{\nu}{16} \rho_1$, where ν is the outward unit normal vector to $\partial\Omega$ at the point \overline{P} . Let $x \in \Omega'_{(r)}$ and let $\gamma \subset \Omega'_{(r)}$ be an open arc joining x to z_0 . With the usual procedure we construct a chain of balls with centers along γ where we apply the three-spheres inequality (cf. [AM01] and [MR04]),

$$\int_{B_{3r}} |u|^2 \leq C \left(\int_{B_r} |u|^2 \right)^\tau \left(\int_{B_{4r}} |u|^2 \right)^{1-\tau}$$

($\tau \in (0, 1), C > 0$ absolute constants), in order to give an estimate of the propagation of the error from $B_r(z_0)$ to x . In fact, by (31) and (32) and the above inequality we get

$$\int_{B_{3r}(x)} |u|^2 \leq C r^2 (\varepsilon^\delta)^{\tau N_r} \left(\log \frac{1}{r} \right)^{1-\tau N_r},$$

where $N_r \leq \frac{|\Omega|}{\pi r^2}$. Finally by interpolation we get

$$(34) \quad |u(x)| \leq \omega(\varepsilon, r) := C (\varepsilon^\delta)^{\tau N_r} \left(\log \frac{1}{r} \right)^{1-\tau N_r}$$

for every $x \in \Omega'_{(r)}$.

It is easy to check that there exist a constant $C_0, 0 < C_0 < 1/2$, depending only on the a priori constants such that if $d > 0$ and $\rho \leq C_0 d^2$, then at least one of the following conditions is satisfied:

$$(35) \quad B_\rho(P_0) \cap \sigma_1 = \emptyset \quad \text{or} \quad B_\rho(P_1) \cap \sigma_0 = \emptyset.$$

Without loss of generality we might assume that

$$B_\rho(P_0) \cap \sigma_1 = \emptyset.$$

Up to a rigid transformation we might assume that $P_0 = 0$ and that σ_0 lies along the positive x_1 axis. For $t \in (0, 1], \rho \in (0, C_0 d^2)$, let

$$x_{t\rho}^+ = \bar{x}_\rho + e_2 \frac{\rho t}{2}, \quad x_{t\rho}^- = \bar{x}_\rho - e_2 \frac{\rho t}{2},$$

where $\bar{x}_\rho = (\rho/2, 0)$ and $e_2 = (0, 1)$. Let

$$u_{\sigma_0}^\pm(\bar{x}_\rho) = \lim_{s \rightarrow 0^+} u_{\sigma_0}(\bar{x}_\rho \pm se_2)$$

and

$$u^\pm(\bar{x}_\rho) = \lim_{s \rightarrow 0^+} u(\bar{x}_\rho \pm se_2).$$

We have

$$|[u_{\sigma_0}(\bar{x}_\rho)]_{\sigma_0}| = |[u(\bar{x}_\rho)]_{\sigma_0}| \leq |u^+(\bar{x}_\rho) - u(x_{t\rho}^+)| + |u^-(\bar{x}_\rho) - u(x_{t\rho}^-)| + |u(x_{t\rho}^+)| + |u(x_{t\rho}^-)|.$$

By applying the second of (31) we get

$$|u^+(\bar{x}_\rho) - u(x_{t\rho}^+)| \leq Ct$$

and

$$|u^-(\bar{x}_\rho) - u(x_{t\rho}^-)| \leq Ct,$$

where $C > 0$ depends only on the a priori data. Hence using (34) and the last inequalities we derive

$$|[u_{\sigma_0}(\bar{x}_\rho)]| \leq Ct + 2\omega(\varepsilon; t\rho/2).$$

Therefore, from the properties of double layer potentials, we get

$$(36) \quad |\varphi_0| = |[u_{\sigma_0}(\bar{x}_\rho)]| \leq Ct + 2\omega(\varepsilon; t\rho/2),$$

where C depends only on the a priori constants.

Now we argue by contradiction. Assume

$$d \geq \eta_\varepsilon,$$

where

$$\eta_\varepsilon := \frac{\sqrt{2}}{\sqrt{C_0}} \left(\frac{8E_1^2 |\log \tau|}{\log |\log \varepsilon^\delta|} \right)^{1/8}.$$

Let

$$\rho = \left(\frac{8E_1^2 |\log \tau|}{\log |\log \varepsilon^\delta|} \right)^{1/4} \quad \text{and} \quad t = \left(\frac{8E_1^2 |\log \tau|}{\log |\log \varepsilon^\delta|} \right)^{1/4}.$$

By (36) there exists $\varepsilon_0 > 0$, depending only on the a priori data, such that if $\varepsilon \leq \varepsilon_0$ and

$$d \geq \eta_\varepsilon,$$

then

$$(37) \quad |\varphi_0| \leq \omega_1(\varepsilon) := C (\log |\log(\varepsilon)|)^{-1/4},$$

where C depends only on the a priori constants.

By using (12) for σ_0 and σ_1 and estimating the difference, we can easily see that

$$\begin{aligned} |\Gamma(P_0, y)\tau_0 f_0| &\leq |\Gamma(Q_0, y)\tau_0 f_0| + \left| \int_{\sigma_0} \mathbb{C}_0 \widehat{\nabla} N(x, y) n_0 \cdot \varphi_0 ds(x) \right| \\ &\quad + \left| \int_{\sigma_1} \mathbb{C}_0 \widehat{\nabla} N(x, y) n_1 \cdot \varphi_1 ds(x) \right| + |(N(P_0, y) - \Gamma(P_0, y))\tau_0 f_0| \\ &\quad + |u(y)| + |N(Q_1, y)\tau_1 f_1| + |N(P_1, y)\tau_1 f_1|. \end{aligned}$$

From Lemma 3.2 and (37) we get, for any $y \in \Omega \setminus \sigma_0$

$$\left| \int_{\sigma_0} \mathbb{C}_0 \widehat{\nabla} N(x, y) n_0 \cdot \varphi_0 ds(x) \right| \leq C\omega_1(\varepsilon) \left| \log \frac{|y - P_0|}{|y - Q_0|} \right|$$

and, for any $y \in \Omega \setminus \sigma_1$,

$$\left| \int_{\sigma_1} \mathbb{C}_0 \widehat{\nabla} N(x, y) n_1 \cdot \varphi_1 ds(x) \right| \leq C|\varphi_1| \left| \log \frac{|y - P_1|}{|y - Q_1|} \right|,$$

where C depends only on the a priori data. Furthermore, by regularity estimates, we have

$$|N(P_0, y) - \Gamma(P_0, y)| \leq C \quad \text{for every } y \in \Omega,$$

where C depends only on the a priori constants. It is now straightforward to see that there exists $\varepsilon_1 > 0$, depending only on the a priori data, such that if $\varepsilon \leq \varepsilon_1$ and

$$d \geq (\log(\log |\log \varepsilon^\delta|))^{-1},$$

then

$$(38) \quad |f_0| \leq C\omega_1(\varepsilon).$$

Hence, for $\varepsilon \leq \varepsilon_2 := \min(\varepsilon_0, \varepsilon_1)$ and $d \geq (\log(\log |\log \varepsilon^\delta|))^{-1}$, inequalities (37) and (38) hold true and we get

$$\begin{cases} |n_0^T W \tau_0| \leq C\omega_1(\varepsilon), \\ |\alpha(\text{tr}W) - \beta n_0^T W n_0| \leq C\omega_1(\varepsilon), \\ |a(\text{tr}W) + (b + d)n_0^T W n_0| \leq C\omega_1(\varepsilon), \end{cases}$$

which lead to

$$k_0 = \max_{|x|=1} |x^T W x| \leq C\omega_1(\varepsilon),$$

a contradiction for ε small enough. Hence

$$d \leq (\log(\log |\log \varepsilon^\delta|))^{-1}$$

for every $\varepsilon \leq \varepsilon_2$. If $\varepsilon \geq \varepsilon_2$, then

$$d \leq E_1 \leq E_1 \frac{\varepsilon}{\varepsilon_1}$$

and the proof is complete. \square

5.2. Proof of Lemma 4.2. Taking the derivative of u_t ,

$$\begin{aligned}
 u'_t(y) &= \int_0^1 \left(\frac{d}{dt} \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \right) n_t \cdot \varphi_t |Q_t - P_t| ds \\
 &\quad + \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) n_t \cdot \frac{d}{dt} (\varphi_t |Q_t - P_t|) ds \\
 &\quad + \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \left(\frac{d}{dt} n_t \right) \cdot \varphi_t |Q_t - P_t| ds \\
 (39) \quad &\quad + \left(\nabla N(Q_t, y) \frac{dQ_t}{dt} \cdot \tau_t - \nabla N(P_t, y) \frac{dP_t}{dt} \cdot \tau_t \right) f_t \\
 &\quad + (N(Q_t, y) - N(P_t, y)) \frac{d(\tau_t f_t)}{dt} \\
 &:= I_1 + I_2 + I_3 + \left(\nabla N(Q_t, y) \frac{dQ_t}{dt} \cdot \tau_t - \nabla N(P_t, y) \frac{dP_t}{dt} \cdot \tau_t \right) f_t \\
 &\quad + (N(Q_t, y) - N(P_t, y)) \frac{d(\tau_t f_t)}{dt}.
 \end{aligned}$$

Observe that

$$\begin{aligned}
 &\frac{d}{dt} \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \\
 &= \nabla_x (\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet ((1-s)(P_1 - P_0) + s(Q_1 - Q_0)),
 \end{aligned}$$

where the symbol \bullet should be intended in the following way:

$$\nabla_x M \bullet v := \sum_i \frac{\partial M}{\partial x_i} v_i.$$

On the other hand

$$\begin{aligned}
 \frac{d}{ds} \left(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \right) &= \nabla_x (\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet (Q_t - P_t) \\
 &= \nabla_x (\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet \tau_t |Q_t - P_t|.
 \end{aligned}$$

Set

$$(1-s)(P_1 - P_0) + s(Q_1 - Q_0) = A_t(s)\tau_t + B_t(s)n_t,$$

where

$$A_t(s) = (1-s)(P_1 - P_0) \cdot \tau_t + s(Q_1 - Q_0) \cdot \tau_t$$

and

$$B_t(s) = (1-s)(P_1 - P_0) \cdot n_t + s(Q_1 - Q_0) \cdot n_t.$$

Hence

$$\begin{aligned}
 I_1 &= \int_0^1 \left(\nabla_x(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet (A_t(s)\tau_t + B_t(s)n_t) \right) n_t \cdot \varphi_t ds \\
 (40) \quad &= \int_0^1 A_t(s) \frac{d}{ds} \left(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \right) n_t \cdot \varphi_t ds \\
 &\quad + \int_0^1 B_t(s) \left(\nabla_x(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet n_t \right) n_t \cdot \varphi_t |Q_t - P_t| ds.
 \end{aligned}$$

Integrating by parts the first integral appearing on the right-hand side of (40), we obtain

$$\begin{aligned}
 I_1 &= A_t(1)\mathbb{C}_0 \widehat{\nabla} N(Q_t, y)n_t \cdot \varphi_t - A_t(0)\mathbb{C}_0 \widehat{\nabla} N(P_t, y)n_t \cdot \varphi_t \\
 (41) \quad &- \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \frac{d}{ds} (A_t(s)n_t \cdot \varphi_t) ds \\
 &\quad + \int_0^1 B_t(s) \left(\nabla_x(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet n_t \right) n_t \cdot \varphi_t |Q_t - P_t| ds.
 \end{aligned}$$

Observing now that, for any $y \in \partial\Omega$,

$$\nabla_x \left(\mathbb{C}_0 \widehat{\nabla} N(\cdot, y)n_t \right) \bullet n_t + \nabla_x \left(\mathbb{C}_0 \widehat{\nabla} N(\cdot, y) \cdot \tau_t \right) \bullet \tau_t = 0 \quad \text{in } \Omega,$$

we can rewrite

$$\begin{aligned}
 &\int_0^1 B_t(s) \left(\nabla_x(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet n_t \right) n_t \cdot \varphi_t |Q_t - P_t| ds \\
 &= - \int_0^1 B_t(s) \left(\nabla_x(\mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y)) \bullet \tau_t \right) \tau_t \cdot \varphi_t |Q_t - P_t| ds \\
 &= - \int_0^1 B_t(s) \left(\frac{d}{ds} \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \right) \tau_t \cdot \varphi_t ds \\
 (42) \quad &= -B_t(1)\mathbb{C}_0 \widehat{\nabla} N(Q_t, y)\tau_t \cdot \varphi_t + B_t(0)\mathbb{C}_0 \widehat{\nabla} N(P_t, y)\tau_t \cdot \varphi_t \\
 &\quad + \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \frac{d}{ds} (B_t(s)\tau_t \cdot \varphi_t) ds.
 \end{aligned}$$

Inserting (42) into (41) we derive finally

$$\begin{aligned}
 I_1 &= \mathbb{C}_0 \widehat{\nabla} N(Q_t, y)(A_t(1)n_t - B_t(1)\tau_t) \cdot \varphi_t \\
 (43) \quad &- \mathbb{C}_0 \widehat{\nabla} N(P_t, y)(A_t(0)n_t - B_t(0)\tau_t) \cdot \varphi_t \\
 &- \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \frac{d}{ds} (A_t(s)n_t \cdot \varphi_t - B_t(s)\tau_t \cdot \varphi_t) ds,
 \end{aligned}$$

where

$$\begin{aligned} & \frac{d}{ds} (A_t(s)n_t - B_t(s)\tau_t) \cdot \varphi_t \\ &= \{((Q_1 - Q_0) - (P_1 - P_0)) \cdot \tau_t\} n_t - \{((Q_1 - Q_0) - (P_1 - P_0)) \cdot n_t\} \tau_t \cdot \varphi_t, \\ & \quad A_t(1) = (Q_1 - Q_0) \cdot \tau_t, \quad A_t(0) = (P_1 - P_0) \cdot \tau_t \end{aligned}$$

and

$$B_t(1) = (Q_1 - Q_0) \cdot n_t, \quad B_t(0) = (P_1 - P_0) \cdot n_t.$$

It is easy to calculate

$$\frac{d}{dt} n_t = - \left(\frac{(Q_1 - Q_0) - (P_1 - P_0)}{|Q_t - P_t|} \cdot n_t \right) \tau_t,$$

and, with this, we get

$$(44) \quad I_3 = - \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_t + sQ_t, y) \tau_t \cdot \varphi_t ((Q_1 - Q_0) - (P_1 - P_0)) \cdot n_t \, ds.$$

Inserting (43) and (44) into (39) gives (25), which concludes the proof. \square

5.3. Proof of Proposition 4.3. In Lemma 4.2 we established that

$$\begin{aligned} (45) \quad u'_0(y) &= - \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_0 + sQ_0, y) n_0 \cdot \varphi_0 ((Q_1 - Q_0) - (P_1 - P_0)) \cdot \tau_0 \\ & \quad + \int_0^1 \mathbb{C}_0 \widehat{\nabla} N((1-s)P_0 + sQ_0, y) n_0 \cdot \frac{d}{dt} (\varphi_t |Q_t - P_t|) |_{t=0} \\ & \quad + \mathbb{C}_0 \widehat{\nabla} N(Q_0, y) (((Q_1 - Q_0) \cdot \tau_0) n_0 - ((Q_1 - Q_0) \cdot n_0) \tau_0) \cdot \varphi_0 \\ & \quad - \mathbb{C}_0 \widehat{\nabla} N(P_0, y) (((P_1 - P_0) \cdot \tau_0) n_0 - ((P_1 - P_0) \cdot n_0) \tau_0) \cdot \varphi_0 \\ & \quad + (\nabla N(Q_0, y)(Q_1 - Q_0) - \nabla N(P_0, y)(P_1 - P_0)) \cdot \tau_0 f_0 \\ & \quad + (N(Q_0, y) - N(P_0, y)) \frac{d(\tau_t f_t)}{dt} |_{t=0}. \end{aligned}$$

We divide the proof of the first inequality of Proposition 4.3 into two steps.

(1) We first prove that if

$$(46) \quad \|u'_0\|_{L^2(\partial\Omega)} = 0, \quad \text{then} \quad d_{\mathcal{H}}(\sigma_0, \sigma_1) = 0.$$

Since $d_{\mathcal{H}}(\sigma_0, \sigma_1) = |P_1 - P_0|$ it will be sufficient to show that $P_1 = P_0$.

By (24), (46), and the unique continuation property, we have that

$$u'_0(y) = 0 \quad \text{for} \quad y \in \overline{\Omega} \setminus \sigma_0.$$

In particular u'_0 is zero in a neighborhood of P_0 contained in $\Omega \setminus \sigma_0$. From (45) it is immediate to see that u'_0 has a singularity at P_0 whose leading term is

$$(47) \quad -\mathbb{C}_0 \widehat{\nabla} \Gamma(P_0, y) \left((P_1 - P_0)^\perp \right) \cdot \varphi_0 - \nabla \Gamma(P_0, y) (P_1 - P_0) \cdot \tau_0 f_0.$$

Without loss of generality let us now fix the coordinate system in such a way that $P_0 = 0$, $\tau_0 = (1, 0)$, and $n_0 = (0, 1)$, let $\varphi_0 = (\varphi_0^1, \varphi_0^2)$, and let us set $h_1 = (h_1^1, h_1^2) := (P_1 - P_0)$. Since the leading singularity given by (47) cannot appear and by the definition of \mathbb{C}_0 , we must have, for $i = 1, 2$ and for any $y \in B_r(0) \setminus \sigma_0$ for some small and positive radius r ,

$$(48) \quad \alpha_1 \partial_1 \Gamma_{i1}(y) + \alpha_2 \partial_1 \Gamma_{i2}(y) + \alpha_3 \partial_2 \Gamma_{i1}(y) + \alpha_4 \partial_2 \Gamma_{i2}(y) = 0,$$

where the coefficients are given by

$$(49) \quad \begin{aligned} \alpha_1 &= \lambda_0 (h_1^1 \varphi_0^2 - h_1^2 \varphi_0^1) - 2\mu_0 h_1^2 \varphi_0^1 + h_1^1 f_0, & \alpha_2 &= -\mu_0 h_1^2 \varphi_0^2 + \mu_0 h_1^1 \varphi_0^1, \\ \alpha_3 &= -\mu_0 h_1^2 \varphi_0^2 + \mu_0 h_1^1 \varphi_0^1 + h_1^2 f_0, & \alpha_4 &= \lambda_0 (-h_1^2 \varphi_0^1 + h_1^1 \varphi_0^2) + 2\mu_0 h_1^1 \varphi_0^2. \end{aligned}$$

By calculating explicitly the derivatives of Γ , equations (48) become

$$(50) \quad \begin{aligned} &\frac{1}{2\pi|y|^4} [y_1^3 (A\alpha_1 - B\alpha_4) + y_1^2 y_2 (B\alpha_2 + (A + 2B)\alpha_3) + y_1 y_2^2 ((A - 2B)\alpha_1 + B\alpha_4) \\ &+ y_2^3 (-B\alpha_2 + A\alpha_3)] = 0 \quad \text{for every } y \in B_r(0) \setminus \sigma_0 \end{aligned}$$

for $i = 1$, and

$$(51) \quad \begin{aligned} &\frac{1}{2\pi|y|^4} [y_1^3 (A\alpha_2 - B\alpha_3) + y_1^2 y_2 (B\alpha_1 + (A - 2B)\alpha_4) + y_1 y_2^2 ((A + 2B)\alpha_2 + B\alpha_3) \\ &+ y_2^3 (-B\alpha_1 + A\alpha_4)] = 0 \quad \text{for every } y \in B_r(0) \setminus \sigma_0 \end{aligned}$$

for $i = 2$.

Equations (50) and (51) can be satisfied if and only if each coefficient of the polynomials appearing in the numerators is zero. Since $A + B$ and $A - B$ are different from zero this implies that

$$\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0.$$

Easy computations show that these four equations can be satisfied for some vector $h_1 \neq (0, 0)$ only if

$$f_0 = \varphi_0^1 = \varphi_0^2 = 0,$$

which implies that W is identically zero, leading to a contradiction. Hence $h_1 = P_1 - P_0 = 0$ and the first step is proved.

(2) Since u'_0 depends linearly on $h := (h_1, h_2) := (P_1 - P_0, Q_1 - Q_0)$, in order to prove the estimate from below it is sufficient to show that if $|h|_{\mathbb{R}^4} = (|h_1|^2 + |h_2|^2)^{1/2} = 1$, then

$$\|u'_0\|_{L^2(\partial\Omega)} \geq m_0,$$

where $m_0 > 0$ depends only on the a priori data.

Let $h \in \mathbb{R}^4$ be such that $|h|_{\mathbb{R}^4} = 1$, where $d_{\mathcal{H}}(\sigma_0, \sigma_1) = |h_1|_{\mathbb{R}^2} \geq |h_2|_{\mathbb{R}^2}$. Denote

$$v = u'_0, \quad \text{and} \quad \theta = \|v\|_{L^2(\partial\Omega)}.$$

Recall that v is a solution of (24) and proceed similarly as in Lemma 4.1. If we define $\rho_1 := \min(L^{-1}, \rho_0)$, we have, for any $r \in (0, \rho_1/156)$,

$$|v(x)| \leq \omega(\theta, r) := C(\theta^\delta)^{\tau^{N_r}} \left(\log \frac{1}{r}\right)^{1-\tau^{N_r}}$$

for every $x \in \Omega'_r := \Omega_{8r} \setminus (\sigma_0)_{8r}$, where $(\sigma_0)_{8r}$ is defined as in (33), $\tau, \delta \in (0, 1)$ and $C > 1$ depend on the a priori data only, and $N_r \leq \frac{|\Omega|}{\pi r^2}$. Without loss of generality we might assume that $P_0 = 0$. By (44) we have that

$$(52) \quad v(y) = \frac{H(y)}{|y|^4} + J(y) \quad \text{for } y \in \Omega \setminus \sigma_0,$$

where H is an homogeneous polynomial of degree three whose coefficients are equal to the expressions in (50) and (51), and J satisfies the inequality

$$|J(y)| \leq C |\log |y|| \quad \text{for } y \in \Omega \setminus \sigma_0, |y| \leq \frac{1}{2L},$$

where C depends on the a priori data only. Let $\bar{y} \in \mathbb{R}^2$, $|\bar{y}| = 1$ such that

$$|H(\bar{y})| = \max_{|y|=1} |H(y)|.$$

Denote

$$M = |H(\bar{y})|, \quad y_r = 8r\bar{y}.$$

Since $|H(y_r)| = |H(-y_r)|$ and y_r or $-y_r \in \Omega'_r$, it is not restrictive to assume that $y_r \in \Omega'_r$.

By (52) we have

$$|H(\bar{y})| = \frac{|H(y_r)|}{|y_r|^3} \leq |y_r| |v(y_r)| + |y_r| |J(y_r)|,$$

and hence

$$(53) \quad M \leq 8r\omega(\theta; r) + Cr \log \frac{1}{r}$$

for every $0 \leq r \leq \rho_2$, where $\rho_2 = \frac{1}{256} \min(L^{-1}, \rho_0)$. Denote

$$r_\theta = \left(\frac{E_1^2 |\log \tau|}{|\log |\log \theta^\delta||}\right)^{1/2}.$$

There exists $\theta_0 < e^{-1/\delta}$, θ_0 depending only on the a priori data, such that if $\theta \in (0, \theta_0)$, then $r_\theta \in (0, \rho_2)$. We now choose, for $\theta \in (0, \theta_0)$, $r = r_\theta$ and by (53) we have

$$(54) \quad M \leq C\omega_2(\theta),$$

where

$$\omega_2(\theta) = (\log |\log \theta^\delta|)^{-1/4}$$

and C depends only on the a priori data. This implies that each coefficient of the polynomials appearing in the numerators of (50) and (51) are bounded by $C\omega_2(\theta)$ for some positive constant C depending only on a priori data.

Now, proceeding as in step (1) we derive easily that

$$(55) \quad \text{if } \theta \in (0, \theta_0), \quad \text{then } |h_2| \leq |h_1| \leq C\omega_2(\theta),$$

where C depends only on the a priori data. On the other hand, since $|h_1|^2 + |h_2|^2 = 1$, by (55) we obtain that

$$(56) \quad \text{if } \theta \in (0, \theta_0), \quad \text{then } 1 \leq 2C_0^2(\omega_2(\theta))^2.$$

Denote by $\theta_1 = \min\{\theta_0, \exp(-\frac{1}{\delta} \exp(16C_0^4))\}$. By (55) we get that if $\theta \in (0, \theta_1)$, then $1 \leq 2C_0^2(\omega_2(\theta))^2 \leq 1/2$, a contradiction. Therefore $\theta \geq \theta_1$ and recalling that $v = u'_0$, and $\theta = \|v\|_{L^2(\partial\Omega)}$, we get the thesis for $m_0 = \theta_1$.

Let us finally prove (26). Observe that u'_t is differentiable in t and taking, for instance, the following term appearing in $u'_t - u'_0$ we get

$$\begin{aligned} & |\nabla N(Q_t, y)(Q_1 - Q_0) \cdot \tau_t f_t - \nabla N(Q_0, y)(Q_1 - Q_0) \cdot \tau_0 f_0| \\ & \leq |(\nabla N(Q_t, y) - \nabla N(Q_0, y))|(Q_1 - Q_0) \cdot \tau_t f_t| \\ & \quad + |\nabla N(Q_0, y)(Q_1 - Q_0) \cdot (\tau_t - \tau_0) f_t| \\ & \quad + |\nabla N(Q_0, y)(Q_1 - Q_0) \cdot \tau_0(f_t - f_0)|. \end{aligned}$$

Since

$$\begin{aligned} |\nabla N(Q_t, y) - \nabla N(Q_0, y)| & \leq C|Q_1 - Q_0| \leq Cd_{\mathcal{H}}(\sigma_0, \sigma_1), \\ |\tau_t - \tau_0| & \leq C|P_1 - P_0| = Cd_{\mathcal{H}}(\sigma_0, \sigma_1), \end{aligned}$$

and, also,

$$|f_t - f_0| \leq C|Q_1 - Q_0| \leq Cd_{\mathcal{H}}(\sigma_0, \sigma_1),$$

where C is a constant depending only on the a priori data, hence

$$|\nabla N(Q_t, y)(Q_1 - Q_0) \cdot \tau_t f_t - \nabla N(Q_0, y)(Q_1 - Q_0) \cdot \tau_0 f_0| \leq C(d_{\mathcal{H}}(\sigma_0, \sigma_1))^2.$$

Evaluating the other terms appearing in $u'_t - u'_0$ in a similar way, we get the desired inequality and the proof is concluded. \square

REFERENCES

[ABF04] H. AMMARI, E. BERETTA, AND E. FRANCINI, *Reconstruction of thin conductivity imperfections*, Appl. Anal., 83 (2004), pp. 63–78.
 [ABF06] H. AMMARI, E. BERETTA, AND E. FRANCINI, *Reconstruction of thin conductivity imperfections, II. The case of multiple segments*, Appl. Anal., 85 (2006), pp. 87–105.
 [AK04] H. AMMARI AND H. KANG, *Reconstruction of Small Inhomogeneities from Boundary Measurements*, Lecture Notes in Math. 1846, Springer-Verlag, Berlin, 2004.

- [AKNT02] H. AMMARI, H. KANG, G. NAKAMURA, AND K. TANUMA, *Complete asymptotic expansions of solutions of the system of elastostatics in the presence of an inclusion of small diameter and detection of an inclusion*, J. Elasticity, 67 (2002), pp. 97–129.
- [AM01] G. ALESSANDRINI AND A. MORASSI, *Strong unique continuation for the Lamé system of elasticity*, Comm. Partial Differential Equations, 26 (2001), pp. 1787–1810.
- [BF06] E. BERETTA AND E. FRANCINI, *An asymptotic formula for the displacement field in the presence of thin elastic inhomogeneities*, SIAM J. Math. Anal., 38 (2006), pp. 1249–1261.
- [C80] S. CAMPANATO, *Sistemi ellittici in forma divergenza. Regolarità all'interno*, Quaderni della Scuola Normale Superiore di Pisa, 1980.
- [CMV98] D. J. CEDIO-FENGYA, S. MOSKOW, AND M. VOGELIUS, *Identification of conductivity imperfections of small diameter by boundary measurements. Continuous dependence and computational reconstruction*, Inverse Problems, 14 (1998), pp. 553–595.
- [F72] G. FICHERA, *Existence theorems in elasticity*, in Handbuch der Physik, Vol. VI, Springer-Verlag, Berlin, Heidelberg, New York, 1972, pp. 347–389.
- [FV89] A. FRIEDMAN AND M. VOGELIUS, *Identification of small inhomogeneities of extreme conductivity by boundary measurements: A theorem on continuous dependence*, Arch. Ration. Mech. Anal., 105 (1989), pp. 299–326.
- [G72] M. E. GURTIN, *The linear theory of elasticity*, in Handbuch der Physik, Vol. VI, Springer-Verlag, Berlin, Heidelberg, New York, 1972, pp. 1–295.
- [MR04] A. MORASSI AND E. ROSSET, *Stable determination of cavities in elastic bodies*, Inverse Problems, 20 (2004), pp. 453–480.

NEUTRAL FUNCTIONAL DIFFERENTIAL EQUATIONS WITH APPLICATIONS TO COMPARTMENTAL SYSTEMS*

VÍCTOR MUÑOZ-VILLARRAGUT[†], SYLVIA NOVO[†], AND RAFAEL OBAYA[†]

Abstract. We study the monotone skew-product semiflow generated by a family of neutral functional differential equations with infinite delay and stable D -operator. The stability properties of D allow us to introduce a new order and to take the neutral family to a family of functional differential equations with infinite delay. Next, we establish the 1-covering property of omega-limit sets under the componentwise separating property and uniform stability. Finally, the obtained results are applied to the study of the long-term behavior of the amount of material within the compartments of a neutral compartmental system with infinite delay.

Key words. nonautonomous dynamical systems, monotone skew-product semiflows, neutral functional differential equations, infinite delay, compartmental systems

AMS subject classifications. 37B55, 34K40, 34K14

DOI. 10.1137/070711177

1. Introduction. After the pioneering work of Hale and Meyer [11], the theory of neutral functional differential equations (NFDE) aroused considerable interest and a fast development ensued. At present a wide collection of theoretical and practical results make up the main body of the theory of NFDEs (see Hale [10], Hale and Verduyn Lunel [12], Kolmanovskii and Myshkis [18], and Salamon [22], among many others). In particular, a substantial number of results for delayed functional differential equations (FDEs) have been generalized for NFDEs solving new and challenging problems in these extensions.

In this paper we provide a dynamical theory for nonautonomous monotone NFDEs with infinite delay and autonomous stable D -operator along the lines of the results by Jiang and Zhao [17] and Novo, Obaya, and Sanz [20]. We assume some recurrence properties on the temporal variation of the NFDE. Thus, its solutions induce a skew-product semiflow with minimal flow on the base. In particular, the uniform almost periodic and almost automorphic cases are included in this formulation. The skew-product formalism permits the analysis of the dynamical properties of the trajectories using methods of ergodic theory and topological dynamics.

Novo et al. [20] study the existence of recurrent solutions of nonautonomous FDEs with infinite delay using the phase space $BU \subset C((-\infty, 0], \mathbb{R}^m)$ of bounded and uniformly continuous functions with the supremum norm. Assuming some technical conditions on the vector field, it is shown that every bounded solution is relatively compact for the compact-open topology, and its omega-limit set admits a flow extension. An alternative method for the study of recurrent solutions of almost periodic FDEs with infinite delay makes use of a *fading memory* Banach phase space (see Hino, Murakami, and Naiko [13] for an axiomatic definition and main properties). Since this kind of space contains BU and, under natural assumptions, the restriction of the norm

*Received by the editors December 17, 2007; accepted for publication (in revised form) June 12, 2008; published electronically October 13, 2008. This work was partially supported by Junta de Castilla y León under project VA024A06 and by the M.E.C. under project MTM2005-02144.

<http://www.siam.org/journals/sima/40-3/71117.html>

[†]Departamento de Matemática Aplicada, E.T.S. de Ingenieros Industriales, Universidad de Valladolid, 47011 Valladolid, Spain (vicmun@wmatem.eis.uva.es, sylnov@wmatem.eis.uva.es, rafoba@wmatem.eis.uva.es).

topology to the closure of a bounded solution agrees with the compact-open topology, it seems that the approach considered in Novo et al. [20] becomes natural in many cases of interest.

In this paper we consider NFDEs with linear autonomous operator D defined on BU which is continuous for the norm, continuous for the compact-open topology on bounded sets, and atomic at zero. We obtain an integral representation of D by means of Riesz theorem $Dx = \int_{-\infty}^0 [d\mu]x$, where μ is a real Borel measure with finite total variation. The convolution operator \widehat{D} defined by $\widehat{D}x(s) = \int_{-\infty}^0 [d\mu(\theta)]x(\theta + s)$ maps BU into BU . We prove that if D is stable in the sense of Hale [10], then \widehat{D} is an isomorphism of BU which is continuous for the norm and continuous for the compact-open topology on bounded sets. Moreover, \widehat{D}^{-1} inherits these same properties and is associated to a linear stable operator D^* . In fact, the mentioned behavior of \widehat{D}^{-1} characterizes the stability of the operator D . The proofs are self-contained and require only quantitative estimates associated to the stability of the operator D .

Staffans [24] shows that every NFDE with finite delay and autonomous stable D -operator can be written as an FDE with infinite delay in an appropriate fading memory space. A more systematic study on the inversion of the convolution operator \widehat{D} can be found in the work of Gripenberg, Londen, and Staffans [4]. The papers by Haddock et al. [8, 9], Arino and Bourad [1], among others, make a systematic use of these ideas which have a theoretical and practical interest. We give a version of the above results for infinite delay NFDEs. It is obvious that the inversion of the convolution operator \widehat{D} on BU allows us to transform the original equation into a retarded nonautonomous FDE with infinite delay. In addition, we transfer the dynamical theory of Jiang and Zhao [17] and Novo et al. [20] to nonautonomous monotone NFDEs with infinite delay and autonomous stable D -operator. In an appropriate dynamical framework we assume that the trajectories are bounded, uniformly stable, and satisfy a componentwise separating property, and we show that the omega-limit sets are all copies of the base. It is important to mention that no conditions of strong monotonicity are required, which permits the application of the results under natural physical conditions.

In this paper we provide a detailed description of the long-term behavior of the dynamics in some classes of compartmental systems extensively studied in the literature. Compartmental systems have been used as mathematical models for the study of the dynamical behavior of many processes in biological and physical sciences, which depend on local mass balance conditions (see Jacquez [14], Jacquez and Simon [15, 16], and the references therein). Some initial results for models described by FDEs with finite and infinite delay can be found in Györi [5] and Györi and Eller [6]. The paper by Arino and Haourigui [2] proves the existence of almost periodic solutions for compartmental systems described by almost periodic finite delay FDEs. NFDEs represent systems where the compartments produce or swallow material. Györi and Wu [7] and Wu and Freedman [28] study autonomous NFDEs with finite and infinite delay similar to those considered in this paper. We provide a nonautonomous version, under more general assumptions, of the monotone theory for NFDEs included in Wu and Freedman [28] and Wu [27]. More precise results for the case of scalar NFDEs can be found in Arino and Bourad [1] and Krisztin and Wu [19].

We study the dynamics of monotone compartmental systems in terms of the geometrical structure of the pipes connecting the compartments. Irreducible subsets of the set of indices detect the occurrence of subsystems on the complete system which reduce the dimension of the problem being studied. When the system is closed,

the total mass is an invariant continuous function which implies the stability and boundedness of the solutions. In particular, the omega-limit set of every solution is a minimal set and a copy of the base. In a general compartmental system the existence of a bounded solution assures that every solution is bounded and uniformly stable. We first check that when there is no inflow of material then all the compartments of each irreducible subset with some outflow of material are empty on minimal subsets. On the contrary, when the solutions remain bounded and there is an inflow of material in some compartment of an irreducible subset, then for the indices of this irreducible subset all the minimal sets agree, and all the compartments contain some material. Finally, we describe natural physical conditions which ensure the existence of a unique minimal set asymptotically stable.

This paper is organized as follows. Basic notions in topological dynamics, used throughout the rest of the sections, are stated in section 2. Section 3 is devoted to the study of general and stability properties of linear autonomous operators from BU to \mathbb{R}^m , as well as the behavior of solutions of the corresponding homogeneous and nonhomogeneous equations given by them. In section 4, we study the monotone skew-product semiflow generated by a family of NFDEs with infinite delay and stable D -operator. In particular, we establish the 1-covering property of omega-limit sets under the componentwise separating property and uniform stability. These results are applied in section 5 to show that the solutions of a compartmental system given by a monotone NFDE with infinite delay are asymptotically of the same type as the transport functions. Finally, section 6 deals with the long-term behavior of its solutions in terms of the geometrical structure of the pipes, as explained above.

2. Some preliminaries. Let (Ω, d) be a compact metric space. A real *continuous flow* $(\Omega, \sigma, \mathbb{R})$ is defined by a continuous mapping $\sigma : \mathbb{R} \times \Omega \rightarrow \Omega$, $(t, \omega) \mapsto \sigma(t, \omega)$ satisfying

- (i) $\sigma_0 = \text{Id}$,
- (ii) $\sigma_{t+s} = \sigma_t \circ \sigma_s$ for each $s, t \in \mathbb{R}$,

where $\sigma_t(\omega) = \sigma(t, \omega)$ for all $\omega \in \Omega$ and $t \in \mathbb{R}$. The set $\{\sigma_t(\omega) \mid t \in \mathbb{R}\}$ is called the *orbit* or the *trajectory* of the point ω . We say that a subset $\Omega_1 \subset \Omega$ is *σ -invariant* if $\sigma_t(\Omega_1) = \Omega_1$ for every $t \in \mathbb{R}$. A subset $\Omega_1 \subset \Omega$ is called *minimal* if it is compact, σ -invariant, and its only nonempty compact σ -invariant subset is itself. Every compact and σ -invariant set contains a minimal subset; in particular it is easy to prove that a compact σ -invariant subset is minimal if and only if every trajectory is dense. We say that the continuous flow $(\Omega, \sigma, \mathbb{R})$ is *recurrent* or *minimal* if Ω is minimal.

The flow $(\Omega, \sigma, \mathbb{R})$ is *distal* if for any two distinct points $\omega_1, \omega_2 \in \Omega$ the orbits keep at a positive distance, that is, $\inf_{t \in \mathbb{R}} d(\sigma(t, \omega_1), \sigma(t, \omega_2)) > 0$. The flow $(\Omega, \sigma, \mathbb{R})$ is *almost periodic* when for every $\varepsilon > 0$ there is a $\delta > 0$ such that if $\omega_1, \omega_2 \in \Omega$ with $d(\omega_1, \omega_2) < \delta$, then $d(\sigma(t, \omega_1), \sigma(t, \omega_2)) < \varepsilon$ for every $t \in \mathbb{R}$. If $(\Omega, \sigma, \mathbb{R})$ is almost periodic, it is distal. The converse is not true; even if $(\Omega, \sigma, \mathbb{R})$ is minimal and distal, it does not need to be almost periodic. For the basic properties of almost periodic and distal flows we refer the reader to Ellis [3] and Sacker and Sell [21].

A *flow homomorphism* from another continuous flow (Y, Ψ, \mathbb{R}) to $(\Omega, \sigma, \mathbb{R})$ is a continuous map $\pi : Y \rightarrow \Omega$ such that $\pi(\Psi(t, y)) = \sigma(t, \pi(y))$ for every $y \in Y$ and $t \in \mathbb{R}$. If π is also bijective, it is called a *flow isomorphism*. Let $\pi : Y \rightarrow \Omega$ be a surjective flow homomorphism and suppose (Y, Ψ, \mathbb{R}) is minimal (then so is $(\Omega, \sigma, \mathbb{R})$). (Y, Ψ, \mathbb{R}) is said to be an *almost automorphic extension* of $(\Omega, \sigma, \mathbb{R})$ if there is $\omega \in \Omega$ such that $\text{card}(\pi^{-1}(\omega)) = 1$. Then actually $\text{card}(\pi^{-1}(\omega)) = 1$ for ω in a residual subset $\Omega_0 \subseteq \Omega$; in the nontrivial case $\Omega_0 \subsetneq \Omega$ the dynamics can be very complicated. A minimal flow

(Y, Ψ, \mathbb{R}) is *almost automorphic* if it is an almost automorphic extension of an almost periodic minimal flow $(\Omega, \sigma, \mathbb{R})$. We refer the reader to the work of Shen and Yi [23] for a survey of almost periodic and almost automorphic dynamics.

Let E be a complete metric space and $\mathbb{R}^+ = \{t \in \mathbb{R} \mid t \geq 0\}$. A *semiflow* (E, Φ, \mathbb{R}^+) is determined by a continuous map $\Phi : \mathbb{R}^+ \times E \rightarrow E$, $(t, x) \mapsto \Phi(t, x)$ which satisfies

- (i) $\Phi_0 = \text{Id}$,
- (ii) $\Phi_{t+s} = \Phi_t \circ \Phi_s$ for all $t, s \in \mathbb{R}^+$,

where $\Phi_t(x) = \Phi(t, x)$ for each $x \in E$ and $t \in \mathbb{R}^+$. The set $\{\Phi_t(x) \mid t \geq 0\}$ is the *semiorbit* of the point x . A subset E_1 of E is *positively invariant* (or just Φ -*invariant*) if $\Phi_t(E_1) \subset E_1$ for all $t \geq 0$. A semiflow (E, Φ, \mathbb{R}^+) admits a *flow extension* if there exists a continuous flow $(E, \tilde{\Phi}, \mathbb{R})$ such that $\tilde{\Phi}(t, x) = \Phi(t, x)$ for all $x \in E$ and $t \in \mathbb{R}^+$. A compact and positively invariant subset admits a flow extension if the semiflow restricted to it admits one.

Write $\mathbb{R}^- = \{t \in \mathbb{R} \mid t \leq 0\}$. A *backward orbit* of a point $x \in E$ in the semiflow (E, Φ, \mathbb{R}^+) is a continuous map $\psi : \mathbb{R}^- \rightarrow E$ such that $\psi(0) = x$ and for each $s \leq 0$ it holds that $\Phi(t, \psi(s)) = \psi(s + t)$ whenever $0 \leq t \leq -s$. If for $x \in E$ the semiorbit $\{\Phi(t, x) \mid t \geq 0\}$ is relatively compact, we can consider the *omega-limit set* of x ,

$$\mathcal{O}(x) = \bigcap_{s \geq 0} \text{closure}\{\Phi(t + s, x) \mid t \geq 0\},$$

which is a nonempty compact connected and Φ -invariant set. Namely, it consists of the points $y \in E$ such that $y = \lim_{n \rightarrow \infty} \Phi(t_n, x)$ for some sequence $t_n \uparrow \infty$. It is well known that every $y \in \mathcal{O}(x)$ admits a backward orbit inside this set. Actually, a compact positively invariant set M admits a flow extension if every point in M admits a unique backward orbit which remains inside the set M (see Shen and Yi [23, part II]).

A compact positively invariant set M for the semiflow (E, Φ, \mathbb{R}^+) is *minimal* if it does not contain any nonempty compact positively invariant set other than itself. If E is minimal, we say that the semiflow is minimal.

A semiflow is of *skew-product type* when it is defined on a vector bundle and has a triangular structure; more precisely, a semiflow $(\Omega \times X, \tau, \mathbb{R}^+)$ is a *skew-product semiflow* over the product space $\Omega \times X$, for a compact metric space (Ω, d) and a complete metric space (X, d) , if the continuous map τ is as follows:

$$(2.1) \quad \begin{aligned} \tau : \mathbb{R}^+ \times \Omega \times X &\longrightarrow \Omega \times X \\ (t, \omega, x) &\longmapsto (\omega \cdot t, u(t, \omega, x)), \end{aligned}$$

where $(\Omega, \sigma, \mathbb{R})$ is a real continuous flow $\sigma : \mathbb{R} \times \Omega \rightarrow \Omega$, $(t, \omega) \mapsto \omega \cdot t$, called the *base flow*. The skew-product semiflow (2.1) is *linear* if $u(t, \omega, x)$ is linear in x for each $(t, \omega) \in \mathbb{R}^+ \times \Omega$.

Now, we introduce some definitions concerning the stability of the trajectories. A forward orbit $\{\tau(t, \omega_0, x_0) \mid t \geq 0\}$ of the skew-product semiflow (2.1) is said to be *uniformly stable* if for every $\varepsilon > 0$ there is a $\delta(\varepsilon) > 0$, called the *modulus of uniform stability*, such that if $s \geq 0$ and $d(u(s, \omega_0, x_0), x) \leq \delta(\varepsilon)$ for certain $x \in X$, then for each $t \geq 0$,

$$d(u(t + s, \omega_0, x_0), u(t, \omega_0 \cdot s, x)) = d(u(t, \omega_0 \cdot s, u(s, \omega_0, x_0)), u(t, \omega_0 \cdot s, x)) \leq \varepsilon.$$

A forward orbit $\{\tau(t, \omega_0, x_0) \mid t \geq 0\}$ of the skew-product semiflow (2.1) is said to be *uniformly asymptotically stable* if it is uniformly stable and there is a $\delta_0 > 0$ with

the following property: for each $\varepsilon > 0$ there is a $t_0(\varepsilon) > 0$ such that if $s \geq 0$ and $d(u(s, \omega_0, x_0), x) \leq \delta_0$, then

$$d(u(t + s, \omega_0, x_0), u(t, \omega_0 \cdot s, x)) \leq \varepsilon \quad \text{for each } t \geq t_0(\varepsilon).$$

3. Stable D-operators. We consider the Fréchet space $X = C((-\infty, 0], \mathbb{R}^m)$ endowed with the compact-open topology, i.e., the topology of uniform convergence over compact subsets, which is a metric space for the distance

$$d(x, y) = \sum_{n=1}^{\infty} \frac{1}{2^n} \frac{\|x - y\|_n}{1 + \|x - y\|_n}, \quad x, y \in X,$$

where $\|x\|_n = \sup_{s \in [-n, 0]} \|x(s)\|$ and $\|\cdot\|$ denotes the maximum norm on \mathbb{R}^m .

Let $BU \subset X$ be the Banach space

$$BU = \{x \in X \mid x \text{ is bounded and uniformly continuous}\}$$

with the supremum norm $\|x\|_{\infty} = \sup_{s \in (-\infty, 0]} \|x(s)\|$. Given $r > 0$, we will denote

$$B_r = \{x \in BU \mid \|x\|_{\infty} \leq r\}.$$

As usual, given $I = (-\infty, a] \subset \mathbb{R}$, $t \in I$, and a continuous function $x : I \rightarrow \mathbb{R}^m$, x_t will denote the element of X defined by $x_t(s) = x(t + s)$ for $s \in (-\infty, 0]$.

This section is devoted to the study of general and stability properties of linear autonomous operators $D : BU \rightarrow \mathbb{R}^m$, as well as the behavior of solutions of the corresponding homogeneous equation $Dx_t = 0$, $t \geq 0$, and nonhomogeneous equations $Dx_t = h(t)$, $t \geq 0$, for $h \in C([0, \infty), \mathbb{R}^m)$. We will assume the following:

- (D1) D is linear and continuous for the norm.
- (D2) For each $r > 0$, $D : B_r \rightarrow \mathbb{R}^m$ is continuous when we take the restriction of the compact-open topology to B_r ; i.e., if $x_n \xrightarrow{d} x$ as $n \rightarrow \infty$ with $x_n, x \in B_r$, then $\lim_{n \rightarrow \infty} Dx_n = Dx$.
- (D3) D is atomic at 0 (see the definition in Hale [10] or Hale and Verduyn Lunel [12]).

From (D1) and (D2) we obtain the following representation.

PROPOSITION 3.1. *If $D : BU \rightarrow \mathbb{R}^m$ satisfies (D1) and (D2), then for each $x \in BU$*

$$Dx = \int_{-\infty}^0 [d\mu(s)] x(s),$$

where $\mu = [\mu_{ij}]$ and μ_{ij} is a real regular Borel measure with finite total variation $|\mu_{ij}|(-\infty, 0] < \infty$ for all $i, j \in \{1, \dots, m\}$.

Proof. From Riesz representation theorem we obtain the above relation for each x whose components are of compact support. Moreover, if $x \in BU$, there are an $r > 0$ and a sequence of functions of compact support $\{x_n\}_{n \in \mathbb{N}} \subset B_r$ with $\|x_n\|_{\infty} \leq \|x\|_{\infty}$ such that $x_n \xrightarrow{d} x$ as $n \rightarrow \infty$ and, from hypothesis (D2), $\lim_{n \rightarrow \infty} Dx_n = Dx$. However,

$$Dx_n = \int_{-\infty}^0 [d\mu(s)] x_n(s),$$

and the Lebesgue-dominated convergence theorem yields

$$\lim_{n \rightarrow \infty} Dx_n = \int_{-\infty}^0 [d\mu(s)] x(s),$$

which finishes the proof. □

Since in addition D is atomic at 0, $\det[\mu_{ij}(\{0\})] \neq 0$, and without loss of generality, we may assume that

$$(3.1) \quad Dx = x(0) - \int_{-\infty}^0 [d\nu(s)] x(s),$$

where $\nu = [\nu_{ij}]_{i,j \in \{1, \dots, m\}}$, ν_{ij} is a real regular Borel measure with finite total variation, and $|\nu_{ij}|(\{0\}) = 0$ for all $i, j \in \{1, \dots, m\}$. We will denote by $|\nu|[-r, 0]$ the $m \times m$ matrix $[|\nu_{ij}|[-r, 0]]$ and by $\|\nu\|_{\infty}[-r, 0]$ the corresponding matricial norm.

From now on, we will assume that the operator D satisfying (D1)–(D3) has the form (3.1). First, it is easy to check the following result, whose proof is omitted.

PROPOSITION 3.2. *For all $h \in C([0, \infty), \mathbb{R}^m)$ and $\varphi \in BU$ with $D\varphi = h(0)$, the nonhomogeneous equation*

$$(3.2) \quad \begin{cases} Dx_t = h(t), & t \geq 0, \\ x_0 = \varphi \end{cases}$$

has a solution defined for all $t \geq 0$.

Next we obtain a bound for the solution in a finite interval $[0, T]$, in terms of the initial data and the independent term h , which in particular implies the uniqueness of the solution of (3.2).

LEMMA 3.3. *Given $T > 0$, there are positive constants k_T^1, k_T^2 such that if x is a solution of (3.2), then for each $t \in [0, T]$*

$$(3.3) \quad \|x_t\|_{\infty} \leq k_T^1 \sup_{0 \leq u \leq t} \|h(u)\| + k_T^2 \|\varphi\|_{\infty}.$$

Proof. Since $|\nu_{ij}|[-r, 0] \rightarrow 0$ as $r \rightarrow 0$, for each $i, j \in \{1, \dots, m\}$, there is an $r > 0$ such that $\|\nu\|_{\infty}[-r, 0] < 1/2$. Let x be a solution of (3.2). From (3.1),

$$x(t) = h(t) + \int_{-t}^0 [d\nu(s)] x(t+s) + \int_{-\infty}^{-t} [d\nu(s)] \varphi(t+s)$$

for each $t \geq 0$. Consequently, if $t \in [0, r]$,

$$\|x(t)\| \leq \|h(t)\| + \frac{1}{2} \sup_{0 \leq u \leq t} \|x(u)\| + \|\varphi\|_{\infty} \|\nu\|_{\infty}(-\infty, 0],$$

from which we deduce that if $t \in [0, r]$,

$$(3.4) \quad \sup_{0 \leq u \leq t} \|x(u)\| \leq 2 \sup_{0 \leq u \leq t} \|h(u)\| + 2a \|\varphi\|_{\infty},$$

where $a = \|\nu\|_{\infty}(-\infty, 0]$. Next, let $y(t) = x(t+r)$, which is a solution of

$$\begin{cases} Dy_t = h(t+r), & t \geq 0, \\ y_0 = x_r. \end{cases}$$

As above, we conclude that if $t \in [0, r]$,

$$\sup_{0 \leq u \leq t} \|y(u)\| \leq 2 \sup_{0 \leq u \leq t} \|h(u+r)\| + 2a \|x_r\|_{\infty},$$

which together with $\|x_r\|_\infty \leq \|\varphi\|_\infty + \sup_{0 \leq u \leq r} \|x(u)\|$ and (3.4) yields

$$\sup_{0 \leq u \leq t} \|x(u)\| \leq b \sup_{0 \leq u \leq t} \|h(u)\| + c \|\varphi\|_\infty$$

for $t \in [r, 2r]$ and some positive constants b and c independent of h and φ . This way, the result is obtained in a finite number of steps. \square

Following Hale [10], we introduce the concept of stability for the operator D . Although the initial definition is given for the homogeneous equation, it is easy to deduce quantitative estimates in terms of the initial data for the solution of a nonhomogeneous equation.

DEFINITION 3.4. *The linear operator D is said to be stable if there is a continuous function $c \in C([0, \infty), \mathbb{R}^+)$ with $\lim_{t \rightarrow \infty} c(t) = 0$ such that, for each $\varphi \in BU$ with $D\varphi = 0$, the solution of the homogeneous problem*

$$\begin{cases} Dx_t = 0, & t \geq 0, \\ x_0 = \varphi \end{cases}$$

satisfies $\|x(t)\| \leq c(t) \|\varphi\|_\infty$ for each $t \geq 0$.

PROPOSITION 3.5. *Let us assume that D is stable. Then there is a positive constant $d > 0$ such that, for each $h \in C([0, \infty), \mathbb{R}^m)$ with $h(0) = 0$, the solution of*

$$\begin{cases} Dx_t = h(t), & t \geq 0, \\ x_0 = 0 \end{cases}$$

satisfies $\|x(t)\| \leq d \sup_{0 \leq u \leq t} \|h(u)\|$ for each $t \geq 0$.

Proof. Let $\{e_1, \dots, e_m\}$ be the canonical basis of \mathbb{R}^m . A proof similar to that of Lemma 3.2 of Hale [10, sec. 12] shows that there are m functions $\phi_1, \dots, \phi_m \in BU$ such that $D\phi_j = e_j$ for each $j \in \{1, \dots, m\}$. We will denote by Φ the $m \times m$ matrix function $\Phi = [\phi_1, \dots, \phi_m]$ and by $\|\Phi\|_\infty$ the matricial norm corresponding to the norm $\|\cdot\|_\infty$ on BU .

Let $c \in C([0, \infty), \mathbb{R}^m)$ be the function given in Definition 3.4. Assume that c is decreasing and take $T > 0$ such that $c(T) < 1$. From Lemma 3.3, $\|x(t)\| \leq k_T^1 \sup_{0 \leq u \leq t} \|h(u)\|$, provided that $t \in [0, T]$.

If $t \geq T$, there is a $j \in \mathbb{N}$ such that $t \in [jT, (j+1)T]$ and it is easy to check that $x(t) = x^1(t - (j-1)T) + x^2(t - (j-1)T)$, where x^1 and x^2 are the solutions of

$$\begin{cases} Dx_t^1 = 0, & t \geq 0, \\ x_0^1 = x_{(j-1)T} - \Phi h((j-1)T), \end{cases} \quad \begin{cases} Dx_t^2 = h(t + (j-1)T), & t \geq 0, \\ x_0^2 = \Phi h((j-1)T), \end{cases}$$

respectively. From the stability of D and Lemma 3.3, we deduce that

$$\|x(t)\| \leq c(t - (j-1)T) \|x_{(j-1)T} - \Phi h((j-1)T)\|_\infty + k_{2T} \sup_{(j-1)T \leq u \leq t} \|h(u)\|.$$

In addition, since $t - (j-1)T \geq T$ and c is decreasing we conclude that

$$(3.5) \quad \|x(t)\| \leq c(T) c_j + (c(T) \|\Phi\|_\infty + k_{2T}) \sup_{0 \leq u \leq t} \|h(u)\|, \quad t \in [jT, (j+1)T],$$

where $c_j = \|x_{jT}\|_\infty = \sup_{0 \leq u \leq jT} \|x(u)\|$.

Let $a_T = \max\{k_T^1, c(T) \|\Phi\|_\infty + k_{2T}\}$. We have $c_1 \leq a_T \sup_{0 \leq u \leq T} \|h(u)\|$ and from (3.5), if $j \geq 2$,

$$c_j \leq \max \left\{ c_{j-1}, c(T) c_{j-1} + a_T \sup_{0 \leq u \leq jT} \|h(u)\| \right\}.$$

Hence, we check that for each $j \geq 2$

$$c_j \leq a_T (1 + c(T) + \dots + c(T)^{j-1}) \sup_{0 \leq u \leq jT} \|h(u)\|,$$

and again from (3.5) we finally deduce that for $t \geq 0$ (and hence $t \in [jT, (j+1)T]$ for some $j \geq 0$)

$$\|x(t)\| \leq a_T \sum_{k=0}^j c(T)^k \sup_{0 \leq u \leq t} \|h(u)\| \leq \frac{a_T}{1 - c(T)} \sup_{0 \leq u \leq t} \|h(u)\|,$$

which finishes the proof. \square

THEOREM 3.6. *Let us assume that D is stable. Then there is a continuous function $c \in C([0, \infty), \mathbb{R}^+)$ with $\lim_{t \rightarrow \infty} c(t) = 0$ and a positive constant $k > 0$ such that the solution of (3.2) satisfies*

$$\|x(t)\| \leq c(t) \|\varphi\|_\infty + k \sup_{0 \leq u \leq t} \|h(u)\|$$

for each $t \geq 0$.

Proof. It is not hard to check that $x(t) = x^1(t) + x^2(t)$, where x^1 and x^2 are the solutions of

$$\begin{cases} Dx_t^1 = \psi(t) h(t), & t \geq 0, \\ x_0^1 = \varphi, \end{cases} \quad \begin{cases} Dx_t^2 = (1 - \psi(t)) h(t), & t \geq 0, \\ x_0^2 = 0, \end{cases}$$

respectively, and

$$\begin{aligned} \psi: [0, \infty) &\longrightarrow \mathbb{R} \\ t &\longmapsto \psi(t) = \begin{cases} 1 - t, & 0 \leq t \leq 1, \\ 0, & 1 \leq t. \end{cases} \end{aligned}$$

Moreover, since $y(t) = x^1(t + 1)$ satisfies $Dy_t = 0$, $t \geq 0$, with $y_0 = x_1^1$, the result follows from the application of Definition 3.4, Lemma 3.3, and Proposition 3.5 to y , x^1 on $[0, 1]$, and x^2 , respectively. \square

The conclusions of Theorem 3.6 are essential in what follows. In particular, it allows us to estimate the norm of a function x in terms of the norm of the function $(-\infty, 0] \rightarrow \mathbb{R}^m, s \mapsto Dx_s$.

PROPOSITION 3.7. *Let us assume that D is stable. Then there is a positive constant $k > 0$ such that $\|x^h\|_\infty \leq k \|h\|_\infty$ for all $h \in BU$ and $x^h \in BU$ satisfying $Dx_s^h = h(s)$ for $s \leq 0$.*

Proof. Let $x(t)$ be the solution of

$$\begin{cases} Dx_t = h(0), & t \geq 0, \\ x_0 = x^h, \end{cases}$$

$$\tilde{h}(t) = \begin{cases} h(t), & t \leq 0, \\ h(0), & t \geq 0, \end{cases}$$

and for $s \leq 0$ we define

$$y^s(t) = \begin{cases} x(t+s), & t+s \geq 0, \\ x^h(t+s), & t+s \leq 0. \end{cases}$$

Then

$$\begin{cases} Dy_t^s = \tilde{h}(t+s), & t \geq 0, \\ y_0^s = x_s^h, \end{cases}$$

and Theorem 3.6 yields

$$\|y^s(t)\| \leq c(t) \|x_s^h\|_\infty + k \sup_{0 \leq u \leq t} \|\tilde{h}(u+s)\|_\infty \leq c(t) \|x^h\|_\infty + k \|h\|_\infty$$

for all $t \geq 0$ and $s \leq 0$. Hence, $\|x^h(s)\| = \|y^{s-t}(t)\| \leq c(t) \|x^h\|_\infty + k \|h\|_\infty$, and as $t \rightarrow \infty$ we prove the result. \square

Let D be stable and given by (3.1). We define the linear operator

$$(3.6) \quad \begin{array}{ll} \widehat{D}: BU & \longrightarrow BU \\ x & \mapsto \widehat{D}x: (-\infty, 0] \rightarrow \mathbb{R}^m \\ & s \mapsto Dx_s, \end{array}$$

that is, $\widehat{D}x(s) = x(s) - \int_{-\infty}^0 [d\nu(\theta)] x(\theta + s)$ for each $s \in (-\infty, 0]$, which is well defined, i.e., $h = \widehat{D}x \in BU$, provided that $x \in BU$ because D is bounded and $h(s + \tau) - h(s) = D(x_{s+\tau} - x_s)$ for all $\tau, s \leq 0$. Moreover, it is easy to check that \widehat{D} is bounded for the norm and uniformly continuous when we take the restriction of the compact-open topology to B_r ; i.e., given $\varepsilon > 0$ there is a $\delta(r) > 0$ such that $d(\widehat{D}x_1, \widehat{D}x_2) < \varepsilon$ for all $x_1, x_2 \in B_r$ with $d(x_1, x_2) < \delta(r)$. The next result shows, after proving that \widehat{D} is invertible, that the same happens for \widehat{D}^{-1} .

THEOREM 3.8. *Let us assume that D is stable. Then \widehat{D} is invertible, and \widehat{D}^{-1} is bounded for the norm and uniformly continuous when we take the restriction of the compact-open topology to B_r ; i.e., given $\varepsilon > 0$ there is a $\delta(r) > 0$ such that $d(\widehat{D}^{-1}h_1, \widehat{D}^{-1}h_2) < \varepsilon$ for all $h_1, h_2 \in B_r$ with $d(h_1, h_2) < \delta(r)$.*

Proof. \widehat{D} is injective because from Proposition 3.7 the only solution of $Dx_s = 0$ for $s \leq 0$ is $x = 0$. To show that \widehat{D} is onto, let $h \in BU$ and $\{h_n\}_{n \in \mathbb{N}} \subset B_r$, for some $r > 0$, be a sequence of continuous functions whose components are of compact support such that $h_n \xrightarrow{d} h$ as $n \uparrow \infty$. Moreover, it is easy to choose them with the same modulus of uniform continuity as h . It is not hard to check that for each $n \in \mathbb{N}$ there is an $x^n \in BU$ such that $\widehat{D}x^n = h_n$, that is, $Dx_s^n = h_n(s)$ for $s \leq 0$ and $n \in \mathbb{N}$. From Proposition 3.7, $x^n \in B_{kr}$ because $\|x^n\|_\infty \leq k \|h_n\|_\infty \leq kr$ and $\|x^n - x_\tau^n\|_\infty \leq k \|h_n - (h_n)_\tau\|_\infty$ for each $\tau \leq 0$ and $n \in \mathbb{N}$, which implies that $\{x_n\}_{n \in \mathbb{N}}$ is equicontinuous, and hence relatively compact for the compact-open topology. Hence, there is a convergent subsequence, let us assume the whole sequence; i.e., there is a continuous function x such that $x^n \xrightarrow{d} x$ as $n \uparrow \infty$. From this, $x_s^n \xrightarrow{d} x_s$ for each $s \leq 0$ and (3.1) yields $Dx_s^n = h_n(s) \rightarrow Dx_s$, i.e., $Dx_s = h(s)$ for $s \leq 0$ and $\widehat{D}x = h$. It is immediate to check that $x \in BU$ and then \widehat{D} is onto, as claimed.

Since \widehat{D} is linear, bounded for the norm, and bijective, the continuity of \widehat{D}^{-1} for the norm is immediate. However, it also follows from Proposition 3.7 which reads as $\|\widehat{D}^{-1}h\|_\infty \leq k\|h\|_\infty$. Finally, since \widehat{D}^{-1} is linear, to check the uniform continuity for the metric on each B_r , it is enough to prove the continuity at 0, i.e., $\widehat{D}^{-1}x_n \xrightarrow{d} 0$ as $n \uparrow \infty$, whenever $x_n \xrightarrow{d} 0$ as $n \uparrow \infty$ and $\{x_n\}_{n \in \mathbb{N}} \subset B_r$. Let $y^n = \widehat{D}^{-1}x_n \in B_{kr}$. We extend the definition of y^n to $t \geq 0$ as the solution of

$$\begin{cases} Dy_t^n = x_n(0), & t \geq 0, \\ y_0^n = y^n. \end{cases}$$

The stability of D provides

$$(3.7) \quad \|y^n(t+s)\| \leq c(t)\|y^n\|_\infty + k \sup_{s \leq u \leq 0} \|x_n(u)\|$$

for all $t \geq 0$ and $s \leq 0$. Now we check that $\{y^n\}_{n \in \mathbb{N}}$ converges uniformly to 0 on each compact set $K = [-a, 0]$. Given $\varepsilon > 0$, there is a $t_0 > 0$ such that $c(t_0)\|y^n\|_\infty < \varepsilon/2$ for each $n \in \mathbb{N}$. Moreover, $x_n \rightarrow 0$ in $\widetilde{K} = [-a - t_0, 0]$, and hence there is an n_0 such that for each $n \geq n_0$ we have $k\|x_n\|_{\widetilde{K}} < \varepsilon/2$. Therefore, from (3.7) we deduce that for all $u \in K = [-a, 0]$ and $n \geq n_0$,

$$\|y^n(u)\| = \|y^n(t_0 + u - t_0)\| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon,$$

that is, $\|y^n\|_K < \varepsilon$, and $\widehat{D}^{-1}x_n = y^n \xrightarrow{d} 0$ as $n \uparrow \infty$, which finishes the proof. \square

A more systematic study of the properties of the linear operator \widehat{D} defined in (3.6) can be found in Staffans [25]. The next result provides a necessary and sufficient condition for a continuous operator D to be stable. In particular, if \widehat{D} is invertible and \widehat{D}^{-1} is continuous for the restriction of the compact-open topology to B_r , then D is stable.

THEOREM 3.9. *Let $D: BU \rightarrow \mathbb{R}^m$ be given by (3.1) and let \widehat{D} be the linear operator in BU defined in (3.6). The following statements are equivalent:*

- (i) D is stable.
- (ii) For each $r > 0$ and each sequence $\{x_n\}_{n \in \mathbb{N}}$ in BU such that $\|\widehat{D}x_n\|_\infty \leq r$ and $\widehat{D}x_n \xrightarrow{d} 0$ as $n \uparrow \infty$, $x_n(0) \rightarrow 0$ as $n \uparrow \infty$.

Proof. (i) \Rightarrow (ii) is a consequence of Theorem 3.8.

(ii) \Rightarrow (i) For each $T > 0$ we define $\mathcal{L}_T: \{\varphi \in BU \mid D\varphi = 0\} \rightarrow \mathbb{R}^m$, $\varphi \mapsto x(T)$, where x is the solution of

$$\begin{cases} Dx_t = 0, & t \geq 0, \\ x_0 = \varphi. \end{cases}$$

It is easy to check that \mathcal{L}_T is well defined and linear. In addition, from (3.3) we deduce that $\|\mathcal{L}_T(\varphi)\| \leq \|x_T\|_\infty \leq k_T^2\|\varphi\|_\infty$, and hence it is bounded.

Next we check that $\|\mathcal{L}_T\|_\infty \rightarrow 0$ as $T \rightarrow \infty$, which shows the stability of D because $\|x(T)\| \leq c(T)\|\varphi\|_\infty$ for $c(T) = \|\mathcal{L}_T\|_\infty$. Let us assume, on the contrary, that there exist $\delta > 0$, a sequence $T_n \uparrow \infty$, and a sequence $\{\varphi_n\}_{n \in \mathbb{N}}$ with $\|\varphi_n\|_\infty \leq 1$ and $D\varphi_n = 0$ such that $\|\mathcal{L}_{T_n}(\varphi_n)\| \geq \delta$ for each $n \in \mathbb{N}$. That is, $\|x^n(T_n)\| \geq \delta$, where x^n is the solution of

$$\begin{cases} Dx_t^n = 0, & t \geq 0, \\ x_0^n = \varphi_n. \end{cases}$$

Therefore,

$$\begin{cases} D((x_{T_n}^n)_s) = D(x_{T_n+s}^n) = 0 & \text{if } s \in [-T_n, 0], \\ D((x_{T_n}^n)_s) = D((\varphi_n)_{T_n+s}) & \text{if } s \leq -T_n, \end{cases}$$

and taking $r = \|D\|_\infty$, the sequence $\{x_{T_n}^n\}_{n \in \mathbb{N}} \subset BU$ satisfies $\|\widehat{D}x_n\|_\infty \leq r$ and $\widehat{D}x_{T_n}^n \xrightarrow{d} 0$ as $n \uparrow \infty$. Consequently, $x_{T_n}^n(0) = x^n(T_n) \rightarrow 0$ as $n \uparrow \infty$, which contradicts the fact that $\|x^n(T_n)\| \geq \delta$ and finishes the proof. \square

PROPOSITION 3.10. *Let $D: BU \rightarrow \mathbb{R}^m$ be a stable operator given by (3.1) and let \widehat{D} be the linear operator in BU defined in (3.6). Then*

$$\begin{aligned} D^*: BU &\longrightarrow \mathbb{R}^m \\ x &\longmapsto \widehat{D}^{-1}x(0) \end{aligned}$$

is also stable and satisfies (D1)–(D3).

Proof. From Theorem 3.8, we deduce that D^* satisfies (D1)–(D2). Hence as in Proposition 3.1, there is a real regular Borel measure μ^* with finite total variation such that $D^*x = \int_{-\infty}^0 [d\mu^*(s)]x(s)$ for each $x \in BU$. We can write $\mu^* = A\delta - \nu^*$ with $A = [\mu_{ij}^*(\{0\})]$. We claim that D^* is atomic at 0, i.e., $\det A \neq 0$. Assume, on the contrary, that $\det A = 0$ and let $v \in \mathbb{R}^m$ be a unitary vector with $Av = 0$. For each $\varepsilon > 0$ we take $\varphi_\varepsilon: (-\infty, 0] \rightarrow \mathbb{R}$ with $\|\varphi_\varepsilon\|_\infty = \varphi_\varepsilon(0) = 1$ and $\varphi_\varepsilon(s) = 0$ for each $s \in (-\infty, -\varepsilon]$. Let $x^\varepsilon \in BU$ be defined by $x^\varepsilon(s) = \varphi_\varepsilon(s)v$. The continuity of \widehat{D} yields

$$(3.8) \quad 1 = \|x^\varepsilon\|_\infty \leq c \|\widehat{D}^{-1}x^\varepsilon\|_\infty.$$

However, for each $s \in (-\infty, 0]$

$$\widehat{D}^{-1}x^\varepsilon(s) = D^*x_s^\varepsilon = \varphi_\varepsilon(s)Av - \int_{-\infty}^0 [d\nu^*(\theta)]\varphi_\varepsilon(\theta + s)v$$

and, consequently, $\|\widehat{D}^{-1}x^\varepsilon\|_\infty \leq \|\nu^*\|_\infty(-\varepsilon, 0]$, which tends to 0 as $\varepsilon \rightarrow 0$, contradicts (3.8) and shows that D^* is atomic at 0. Finally, D^* is stable as a consequence of Theorem 3.9. Notice that μ^* is the inverse of the measure μ for the convolution defining the operator \widehat{D} . \square

4. Monotone NFDEs. Throughout this section, we will study the monotone skew-product semiflow generated by a family of NFDEs with infinite delay and stable D -operator. In particular, we establish the 1-covering property of omega-limit sets under the componentwise separating property and uniform stability, as in Jiang and Zhao [17] for FDEs with finite delay, and Novo, Obaya, and Sanz [20] for infinite delay. The main tool in the proof of the result is the transformation of the initial family of NFDEs into a family of FDEs with infinite delay in whose study the results of Novo et al. [20] turn out to be useful.

Let $(\Omega, \sigma, \mathbb{R})$ be a minimal flow over a compact metric space (Ω, d) and denote $\sigma(t, \omega) = \omega \cdot t$ for all $\omega \in \Omega$ and $t \in \mathbb{R}$. In \mathbb{R}^m , we take the maximum norm $\|v\| = \max_{j=1, \dots, m} |v_j|$ and the usual partial order relation

$$\begin{aligned} v \leq w &\iff v_j \leq w_j \quad \text{for } j = 1, \dots, m, \\ v < w &\iff v \leq w \quad \text{and } v_j < w_j \quad \text{for some } j \in \{1, \dots, m\}. \end{aligned}$$

As in section 3, we consider the Fréchet space $X = C((-\infty, 0], \mathbb{R}^m)$ endowed with the compact-open topology, i.e., the topology of uniform convergence over compact subsets, and $BU \subset X$ the Banach space of bounded and uniformly continuous functions with the supremum norm $\|x\|_\infty = \sup_{s \in (-\infty, 0]} \|x(s)\|$.

Let $D: BU \rightarrow \mathbb{R}^m$ be an autonomous and stable linear operator satisfying hypotheses (D1)–(D3) and given by relation (3.1). The subset

$$BU_D^+ = \{x \in BU \mid Dx_s \geq 0 \text{ for each } s \in (-\infty, 0]\}$$

is a positive cone in BU , because it is a nonempty closed subset $BU_D^+ \subset BU$ satisfying $BU_D^+ + BU_D^+ \subset BU_D^+$, $\mathbb{R}^+ BU_D^+ \subset BU_D^+$, and $BU_D^+ \cap (-BU_D^+) = \{0\}$. As usual, a partial order relation on BU is induced, given by

$$\begin{aligned} x \leq_D y &\iff Dx_s \leq Dy_s \text{ for each } s \in (-\infty, 0], \\ x <_D y &\iff x \leq_D y \text{ and } x \neq y. \end{aligned}$$

Remark 4.1. Notice that if we denote the usual partial order of BU

$$x \leq y \iff x(s) \leq y(s) \text{ for each } s \in (-\infty, 0],$$

we have that $x \leq_D y$ if and only if $\widehat{D}x \leq \widehat{D}y$, where \widehat{D} is defined by relation (3.6). Although in some cases they may coincide, this new order is different from the one given by Wu and Freedman in [28].

We consider the family of nonautonomous NFDEs with infinite delay and stable D -operator

$$(4.1)_\omega \quad \frac{d}{dt} Dz_t = F(\omega \cdot t, z_t), \quad t \geq 0, \omega \in \Omega,$$

defined by a function $F: \Omega \times BU \rightarrow \mathbb{R}^m$, $(\omega, x) \mapsto F(\omega, x)$ satisfying the following conditions:

- (F1) F is continuous on $\Omega \times BU$ and locally Lipschitz in x for the norm $\|\cdot\|_\infty$.
- (F2) For each $r > 0$, $F(\Omega \times B_r)$ is a bounded subset of \mathbb{R}^m .
- (F3) For each $r > 0$, $F: \Omega \times B_r \rightarrow \mathbb{R}^m$ is continuous when we take the restriction of the compact-open topology to B_r ; i.e., if $\omega_n \rightarrow \omega$ and $x_n \xrightarrow{d} x$ as $n \rightarrow \infty$ with $x \in B_r$, then $\lim_{n \rightarrow \infty} F(\omega_n, x_n) = F(\omega, x)$.
- (F4) If $x, y \in BU$ with $x \leq_D y$ and $D_j x = D_j y$ holds for some $j \in \{1, \dots, m\}$, then $F_j(\omega, x) \leq F_j(\omega, y)$ for each $\omega \in \Omega$.

From hypothesis (F1), the standard theory of NFDEs with infinite delay (see Wang and Wu [26] and Wu [27]) assures that for each $x \in BU$ and each $\omega \in \Omega$ the system $(4.1)_\omega$ locally admits a unique solution $z(t, \omega, x)$ with initial value x , i.e., $z(s, \omega, x) = x(s)$ for each $s \in (-\infty, 0]$. Therefore, the family $(4.1)_\omega$ induces a local skew-product semiflow

$$(4.2) \quad \begin{aligned} \tau : \mathbb{R}^+ \times \Omega \times BU &\longrightarrow \Omega \times BU \\ (t, \omega, x) &\longmapsto (\omega \cdot t, u(t, \omega, x)), \end{aligned}$$

where $u(t, \omega, x) \in BU$ and $u(t, \omega, x)(s) = z(t + s, \omega, x)$ for $s \in (-\infty, 0]$.

As proved in Theorem 3.8, the operator \widehat{D} defined by relation (3.6) is an isomorphism of BU . Hence, the change of variable $y = \widehat{D}z$ takes $(4.1)_\omega$ to

$$(4.3)_\omega \quad y'(t) = G(\omega \cdot t, y_t), \quad t \geq 0, \omega \in \Omega,$$

with $G: \Omega \times BU \rightarrow \mathbb{R}^m$, $(\omega, x) \mapsto G(\omega, x) = F(\omega, \widehat{D}^{-1}x)$ satisfying the following conditions:

- (H1) G is continuous on $\Omega \times BU$ and locally Lipschitz in x for the norm $\|\cdot\|_\infty$.
- (H2) For each $r > 0$, $G(\Omega \times B_r)$ is a bounded subset of \mathbb{R}^m .
- (H3) For each $r > 0$, $G: \Omega \times B_r \rightarrow \mathbb{R}^m$ is continuous when we take the restriction of the compact-open topology to B_r , i.e., if $\omega_n \rightarrow \omega$ and $x_n \xrightarrow{d} x$ as $n \rightarrow \infty$ with $x \in B_r$, then $\lim_{n \rightarrow \infty} G(\omega_n, x_n) = G(\omega, x)$.
- (H4) If $x, y \in BU$ with $x \leq y$ and $x_j(0) = y_j(0)$ holds for some $j \in \{1, \dots, m\}$, then $G_j(\omega, x) \leq G_j(\omega, y)$ for each $\omega \in \Omega$.

From hypothesis (H1), the standard theory of infinite delay FDEs (see Hino, Murakami, and Naiko [13]) assures that for each $x \in BU$ and each $\omega \in \Omega$ the system $(4.3)_\omega$ locally admits a unique solution $y(t, \omega, x)$ with initial value x , i.e., $y(s, \omega, x) = x(s)$ for each $s \in (-\infty, 0]$. Therefore, the new family $(4.3)_\omega$ induces a local skew-product semiflow

$$(4.4) \quad \begin{aligned} \widehat{\tau} : \mathbb{R}^+ \times \Omega \times BU &\longrightarrow \Omega \times BU \\ (t, \omega, x) &\mapsto (\omega \cdot t, \widehat{u}(t, \omega, x)), \end{aligned}$$

where $\widehat{u}(t, \omega, x) \in BU$ and $\widehat{u}(t, \omega, x)(s) = y(t + s, \omega, x)$ for $s \in (-\infty, 0]$, and it is related to the previous one, (4.2), by

$$(4.5) \quad \widehat{u}(t, \omega, x) = \widehat{D} u(t, \omega, \widehat{D}^{-1}x).$$

As a consequence, most of the results obtained in Novo et al. [20] for the skew-product semiflow (4.4) can now be translated to (4.2).

From hypotheses (F1) and (F2), each bounded solution $z(t, \omega_0, x_0)$ provides a relatively compact trajectory, as deduced from Proposition 4.1 of Novo et al. [20].

PROPOSITION 4.2. *Let $z(t, \omega_0, x_0)$ be a bounded solution of $(4.1)_{\omega_0}$, that is, $r = \sup_{t \in \mathbb{R}} \|z(t, \omega_0, x_0)\| < \infty$. Then $\text{closure}_X \{u(t, \omega_0, x_0) \mid t \geq 0\}$ is a compact subset of BU for the compact-open topology.*

From hypotheses (F1), (F2), and (F3) and Proposition 4.2 and Corollary 4.3 of Novo et al. [20] for the skew-product semiflow (4.4), we can deduce the continuity of the semiflow (4.2) restricted to some compact subsets $K \subset \Omega \times BU$ when the compact-open topology is considered on BU .

PROPOSITION 4.3. *Let $K \subset \Omega \times BU$ be a compact set for the product metric topology and assume that there is an $r > 0$ such that $\tau_t(K) \subset \Omega \times B_r$ for all $t \geq 0$. Then the map*

$$\tau : \mathbb{R}^+ \times K \longrightarrow \Omega \times BU \\ (t, \omega, x) \mapsto (\omega \cdot t, u(t, \omega, x))$$

is continuous when the product metric topology is considered.

From Proposition 4.2, when $z(t, \omega_0, x_0)$ is bounded we can define the omega-limit set of the trajectory of the point (ω_0, x_0) as

$$\mathcal{O}(\omega_0, x_0) = \{(\omega, x) \in \Omega \times BU \mid \exists t_n \uparrow \infty \text{ with } \omega_0 \cdot t_n \rightarrow \omega, u(t_n, \omega_0, x_0) \xrightarrow{d} x\}.$$

Notice that the omega-limit set of a pair $(\omega_0, x_0) \in \Omega \times BU$ makes sense whenever $\text{closure}_X \{u(t, \omega_0, x_0) \mid t \geq 0\}$ is a compact set, because then $\{u(t, \omega_0, x_0)(0) = z(t, \omega_0, x_0) \mid t \geq 0\}$ is a bounded set. Proposition 4.3 implies that the restriction

of the semiflow (4.2) to $\mathcal{O}(\omega_0, x_0)$ is continuous for the compact-open topology. The following result is a consequence of Proposition 4.4 of Novo et al. [20].

PROPOSITION 4.4. *Let $(\omega_0, x_0) \in \Omega \times BU$ be such that $\sup_{t \geq 0} \|z(t, \omega_0, x_0)\| < \infty$. Then $K = \mathcal{O}(\omega_0, x_0)$ is a positively invariant compact subset admitting a flow extension.*

From hypothesis (F4), the monotone character of the semiflow (4.2) is deduced.

PROPOSITION 4.5. *For all $\omega \in \Omega$ and $x, y \in BU$ such that $x \leq_D y$ it holds that*

$$u(t, \omega, x) \leq_D u(t, \omega, y)$$

whenever they are defined.

Proof. From $x \leq_D y$ we know that $\widehat{D}x \leq \widehat{D}y$, and since (F4) \Rightarrow (H4), from Proposition 4.5 of Novo et al. [20] we deduce that $\widehat{u}(t, \omega, \widehat{D}x) \leq \widehat{u}(t, \omega, \widehat{D}y)$ whenever they are defined, that is,

$$u(t, \omega, x) = \widehat{D}^{-1}\widehat{u}(t, \omega, \widehat{D}x) \leq_D \widehat{D}^{-1}\widehat{u}(t, \omega, \widehat{D}y) = u(t, \omega, y),$$

as stated. \square

We establish the 1-covering property of omega-limit sets when in addition to hypotheses (F1)–(F4) the componentwise separating property and uniform stability are assumed:

- (F5) If $x, y \in BU$ with $x \leq_D y$ and $D_i x < D_i y$ holds for some $i \in \{1, \dots, m\}$, then $D_i z_t(\omega, x) < D_i z_t(\omega, y)$ for all $t \geq 0$ and $\omega \in \Omega$.
- (F6) There is an $r > 0$ such that all the trajectories with initial data in $\widehat{D}^{-1}B_r$ are uniformly stable in $\widehat{D}^{-1}B_{r'}$ for each $r' > r$, and relatively compact for the product metric topology.

From relation (4.5) we deduce that the transformed skew-product semiflow (4.4) satisfies the following:

- (H5) If $x, z \in BU$ with $x \leq z$ and $x_i(0) < z_i(0)$ holds for some $i \in \{1, \dots, m\}$, then $y_i(t, \omega, x) < y_i(t, \omega, z)$ for all $t \geq 0$ and $\omega \in \Omega$.
- (H6) There is an $r > 0$ such that all the trajectories with initial data in B_r are uniformly stable in $B_{r'}$ for each $r' > r$, and relatively compact for the product metric topology.

Finally, from Theorem 5.3 of Novo et al. [20] applied to the skew-product semiflow (4.4) satisfying hypotheses (H1)–(H6), we obtain the next result for NFDEs with infinite delay.

THEOREM 4.6. *Assume that hypotheses (F1)–(F6) hold and let $(\omega_0, x_0) \in \Omega \times \widehat{D}^{-1}B_r$ be such that $K = \mathcal{O}(\omega_0, x_0) \subset \Omega \times \widehat{D}^{-1}B_r$. Then $K = \{(\omega, c(\omega)) \mid \omega \in \Omega\}$ is a copy of the base and*

$$\lim_{t \rightarrow \infty} d(u(t, \omega_0, x_0), c(\omega_0 \cdot t)) = 0,$$

where $c : \Omega \rightarrow BU$ is a continuous equilibrium, i.e., $c(\omega \cdot t) = u(t, \omega, c(\omega))$ for any $\omega \in \Omega, t \geq 0$, and it is continuous for the compact-open topology on BU .

Remark 4.7. It is easy to check that it is enough to ask for property (F5) (and for (H5) in the case of FDEs with infinite delay) for initial data in BU whose trajectories are globally defined on \mathbb{R} .

5. Compartmental systems. We consider compartmental models for the mathematical description of processes in which the transport of material between compartments takes a nonnegligible length of time, and each compartment produces or swallows material. We provide a nonautonomous version, without strong monotonicity assumptions, of previous autonomous results by Wu and Freedman [28] and Wu [26].

First, we introduce the model with which we are going to deal as well as some notation. Let us suppose that we have a system formed by m compartments C_1, \dots, C_m . Denote by C_0 the environment surrounding the system, and by $z_i(t)$ the amount of material within compartment C_i at time t for each $i \in \{1, \dots, m\}$. Material flows from compartment C_j into compartment C_i through a pipe P_{ij} having a transit time distribution given by a positive regular Borel measure μ_{ij} with finite total variation $\mu_{ij}(-\infty, 0] = 1$ for each $i, j \in \{1, \dots, m\}$. Let $\tilde{g}_{ij} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be the so-called *transport function* determining the volume of material flowing from C_j to C_i given in terms of the time t and the value of $z_j(t)$ for $i \in \{0, \dots, m\}, j \in \{1, \dots, m\}$. For each $i \in \{1, \dots, m\}$, we will assume that there exists an incoming flow of material \tilde{I}_i from the environment into compartment C_i which depends only on time. For each $i \in \{1, \dots, m\}$, at time $t \geq 0$, the compartment C_i produces material itself at a rate $\sum_{j=1}^m \int_{-\infty}^0 z'_j(t+s) d\nu_{ij}(s)$, where ν_{ij} is a positive regular Borel measure with finite total variation $\nu_{ij}(-\infty, 0] < \infty$ and $\nu_{ij}(\{0\}) = 0$ for all $i, j \in \{1, \dots, m\}$.

Once the destruction and creation of material is taken into account, the change of the amount of material of any compartment $C_i, 1 \leq i \leq m$, equals the difference between the amount of total influx into and total outflux out of C_i , and we obtain a model governed by the following system of infinite delay NFDEs:

$$(5.1) \quad \frac{d}{dt} \left[z_i(t) - \sum_{j=1}^m \int_{-\infty}^0 z_j(t+s) d\nu_{ij}(s) \right] = -\tilde{g}_{0i}(t, z_i(t)) - \sum_{j=1}^m \tilde{g}_{ji}(t, z_i(t)) + \sum_{j=1}^m \int_{-\infty}^0 \tilde{g}_{ij}(t+s, z_j(t+s)) d\mu_{ij}(s) + \tilde{I}_i(t),$$

$i = 1, \dots, m$. For simplicity, we denote $\tilde{g}_{i0} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+, (t, v) \mapsto \tilde{I}_i(t)$ for $i \in \{1, \dots, m\}$ and let $\tilde{g} = (\tilde{g}_{ij})_{i,j} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^{m(m+2)}$. We will assume that

- (C1) \tilde{g} is C^1 -admissible, i.e., \tilde{g} is C^1 in its second variable and $\tilde{g}, \frac{\partial}{\partial v} \tilde{g}$ are uniformly continuous and bounded on $\mathbb{R} \times \{v_0\}$ for all $v_0 \in \mathbb{R}$; all its components are monotone in the second variable, and $\tilde{g}_{ij}(t, 0) = 0$ for each $t \in \mathbb{R}$;
- (C2) \tilde{g} is a recurrent function, i.e., its hull is minimal;
- (C3) $\mu_{ij}(-\infty, 0] = 1$ and $\int_{-\infty}^0 |s| d\mu_{ij}(s) < \infty$;
- (C4) $\nu_{ij}(\{0\}) = 0$ and $\sum_{j=1}^m \nu_{ij}(-\infty, 0] < 1$, which implies that the operator $D : BU \rightarrow \mathbb{R}^m$, with $D_i x = x_i(0) - \sum_{j=1}^m \int_{-\infty}^0 x_j(s) d\nu_{ij}(s), i = 1, \dots, m$, is stable and satisfies (D1)–(D3);
- (C5) the measures $d\eta_{ij} = c_{ij} d\mu_{ij} - \sum_{k=0}^m d_{ki} d\nu_{ij}$ are positive, where

$$c_{ij} = \inf_{(t,v) \in \mathbb{R}^2} \frac{\partial \tilde{g}_{ij}}{\partial v}(t, v) \text{ and } d_{ij} = \sup_{(t,v) \in \mathbb{R}^2} \frac{\partial \tilde{g}_{ij}}{\partial v}(t, v).$$

In practical cases, in which the solutions with physical interest belong to the positive cone and the functions g_{ij} are only defined on $\mathbb{R} \times \mathbb{R}^+$, we can extend them to $\mathbb{R} \times \mathbb{R}$ by $g_{ij}(t, -v) = -g_{ij}(t, v)$ for all $v \in \mathbb{R}^+$. Note that (C5) is a condition for controlling the material produced in the compartments in terms of the material transported through the pipes.

The above formulation includes some particularly interesting cases. When the measures ν_{ij} and μ_{ij} are concentrated on a compact set, then (5.1) is an NFDE with finite delay. When the measures $\nu_{ij} \equiv 0$, then (5.1) is a family of FDEs with finite or infinite delay.

As usual, we include the nonautonomous system (5.1) in a family of nonautonomous NFDEs with infinite delay and stable D -operator of the form (4.1) $_{\omega}$ as follows.

Let Ω be the *hull* of \tilde{g} , namely, the closure of the set of mappings $\{\tilde{g}_t \mid t \in \mathbb{R}\}$, with $\tilde{g}_t(s, v) = \tilde{g}(t+s, v)$, $(s, v) \in \mathbb{R}^2$, with the topology of uniform convergence on compact sets, which from (C1) is a compact metric space (more precisely from the admissibility of \tilde{g} ; see Hino et al. [13]). Let $(\Omega, \sigma, \mathbb{R})$ be the continuous flow defined on Ω by translation, $\sigma : \mathbb{R} \times \Omega \rightarrow \Omega$, $(t, \omega) \mapsto \omega \cdot t$, with $\omega \cdot t(s, v) = \omega(t + s, v)$. By hypothesis (C2), the flow $(\Omega, \sigma, \mathbb{R})$ is minimal. In addition, if \tilde{g} is almost periodic (resp., almost automorphic) the flow will be almost periodic (resp., almost automorphic). Notice that these two cases are included in our formulation.

Let $g : \Omega \times \mathbb{R} \rightarrow \mathbb{R}^{m(m+2)}$, $(\omega, v) \mapsto \omega(0, v)$, continuous on $\Omega \times \mathbb{R}$ and denote $g = (g_{ij})_{i,j}$. It is easy to check that, for all $\omega = (\omega_{ij})_{i,j} \in \Omega$ and all $i \in \{1, \dots, m\}$, ω_{i0} is a function dependent only on t ; thus, we can define $I_i = \omega_{i0}$, $i \in \{1, \dots, m\}$. Let $F : \Omega \times BU \rightarrow \mathbb{R}^m$ be the map defined by

$$F_i(\omega, x) = -g_{0i}(\omega, x_i(0)) - \sum_{j=1}^m g_{ji}(\omega, x_j(0)) + \sum_{j=1}^m \int_{-\infty}^0 g_{ij}(\omega \cdot s, x_j(s)) d\mu_{ij}(s) + I_i(\omega)$$

for $(\omega, x) \in \Omega \times BU$ and $i \in \{1, \dots, m\}$. Hence, the family

$$(5.2)_{\omega} \quad \frac{d}{dt} Dz_t = F(\omega \cdot t, z_t), \quad t \geq 0, \omega \in \Omega,$$

where the stable operator D is defined in (C4) and satisfies (D1)–(D3), includes system (5.1) when $\omega = \tilde{g}$.

It is easy to check that this family satisfies hypotheses (F1)–(F3). The following lemma will be useful when proving (F4) and (F5). We omit its proof, which is analogous to the one given in Wu and Freedman [28] for the autonomous case with finite delay.

LEMMA 5.1. *For all $\omega \in \Omega$, $x, y \in BU$ with $x \leq_D y$, and $i = 1, \dots, m$*

$$(5.3) \quad F_i(\omega, y) - F_i(\omega, x) \geq - \sum_{j=0}^m d_{ji} [D_i y - D_i x] + \sum_{j=1}^m \int_{-\infty}^0 (y_j(s) - x_j(s)) d\eta_{ij}(s),$$

where the measures η_{ij} are defined in (C5).

Condition (C5) is essential to proving the monotone character of the semiflow. It can be improved in some cases (see Arino and Bourad [1] for the scalar case).

PROPOSITION 5.2. *Under assumptions (C1)–(C5), the family (5.2) $_{\omega}$ satisfies hypotheses (F4), (F5) and $\Omega \times BU_D^+$ is positively invariant.*

Proof. Let $x, y \in BU$ with $x \leq_D y$ and $D_i x = D_i y$ for some $i \in \{1, \dots, m\}$. From (C4), apart from the stability of the operator D , it is easy to prove that the inverse operator of \widehat{D} defined by (3.6) is positive. Hence, from $x \leq_D y$, that is, $\widehat{D}x \leq \widehat{D}y$, we also deduce that $x \leq y$, which, together with $D_i x = D_i y$, relation (5.3), and hypothesis (C5), yields $F_i(\omega, y) \geq F_i(\omega, x)$, that is, hypothesis (F4) holds.

Next, we check hypothesis (F5). Let $x, y \in BU$ with $x \leq_D y$ and $D_i x < D_i y$ for some $i \in \{1, \dots, m\}$. Since (F4) holds, from Proposition 4.5 $u(t, \omega, x) \leq_D u(t, \omega, y)$ and, as before, we deduce in this case that $u(t, \omega, x) \leq u(t, \omega, y)$, i.e., $z_t(\omega, x) \leq z_t(\omega, y)$ for all $t \geq 0$ and $\omega \in \Omega$. Let $h(t) = D_i z_t(\omega, y) - D_i z_t(\omega, x)$. From (5.2) $_{\omega}$ and

Lemma 5.1

$$\begin{aligned}
 h'(t) &= F_i(\omega \cdot t, z_t(\omega, y)) - F_i(\omega \cdot t, z_t(\omega, x)) \\
 &\geq - \sum_{j=0}^m d_{ji} h(t) + \sum_{j=1}^m \int_{-\infty}^0 (z_j(t+s, \omega, y) - z_j(t+s, \omega, x)) d\eta_{ij}(s),
 \end{aligned}$$

and again from hypothesis (C5) we deduce that $h'(t) \geq -dh(t)$ for some $d \geq 0$, which together with $h(0) > 0$ yields $h(t) = D_i z_t(\omega, y) - D_i z_t(\omega, x) > 0$ for each $t \geq 0$ and (F5) holds. Finally, since $I_i(\omega) \geq 0$ for each $\omega \in \Omega$ and $i \in \{1, \dots, m\}$, and the semiflow is monotone, a comparison argument shows that $\Omega \times BU_D^+$ is positively invariant, as stated. \square

Next we will study some cases in which hypothesis (F6) is satisfied. In order to do this, we define $M: \Omega \times BU \rightarrow \mathbb{R}$, the total mass of the system (5.2) $_\omega$, as

$$(5.4) \quad M(\omega, x) = \sum_{i=1}^m D_i x + \sum_{i=1}^m \sum_{j=1}^m \int_{-\infty}^0 \left(\int_s^0 g_{ji}(\omega \cdot \tau, x_i(\tau)) d\tau \right) d\mu_{ji}(s)$$

for all $\omega \in \Omega$ and $x \in BU$, which is well defined from condition (C3). The next result shows the continuity properties of M and its variation along the flow.

PROPOSITION 5.3. *The total mass M is a continuous function on all the sets of the form $\Omega \times B_r$ with $r > 0$ for the product metric topology. Moreover, for each $t \geq 0$*

$$(5.5) \quad \frac{d}{dt} M(\tau_t(\omega, x)) = \sum_{i=1}^m [I_i(\omega \cdot t) - g_{0i}(\omega \cdot t, z_i(t, \omega, x))] .$$

Proof. The continuity follows from (D2), (C1), and (C3). A straightforward computation similar to that given in Wu and Freedman [28] shows that

$$(5.6) \quad M(\omega \cdot t, z_t(\omega, x)) = M(\omega, x) + \sum_{i=1}^m \int_0^t [I_i(\omega \cdot s) - g_{0i}(\omega \cdot s, z_i(s, \omega, x))] ds ,$$

from which (5.5) is deduced. \square

The following lemma is essential in the proof of the stability of solutions.

LEMMA 5.4. *Let $x, y \in BU$ with $x \leq_D y$. Then*

$$0 \leq D_i z_t(\omega, y) - D_i z_t(\omega, x) \leq M(\omega, y) - M(\omega, x)$$

for each $i = 1, \dots, m$ and whenever $z(t, \omega, x)$ and $z(t, \omega, y)$ are defined.

Proof. From Propositions 5.2 and 4.5 the skew-product semiflow induced by (5.2) $_\omega$ is monotone. Hence, if $x \leq_D y$, then $u(t, \omega, x) \leq_D u(t, \omega, y)$ whenever they are defined. From this, as before, since \widehat{D}^{-1} is positive we also deduce that $x \leq y$ and $u(t, \omega, x) \leq u(t, \omega, y)$. Therefore, $D_i z_t(\omega, x) \leq D_i z_t(\omega, y)$ and $z_i(t, \omega, x) \leq z_j(t, \omega, y)$ for each $i = 1, \dots, m$. In addition, the monotonicity of transport functions yields $g_{ij}(\omega, z_j(t, \omega, x)) \leq g_{ij}(\omega, z_j(t, \omega, y))$ for each $\omega \in \Omega$. From all these inequalities and (5.4) and (5.6) we deduce that

$$\begin{aligned}
 0 \leq D_i z_t(\omega, y) - D_i z_t(\omega, x) &\leq \sum_{i=1}^m [D_i z_t(\omega, y) - D_i z_t(\omega, x)] \\
 &\leq M(\omega \cdot t, z_t(\omega, y)) - M(\omega \cdot t, z_t(\omega, x)) \leq M(\omega, y) - M(\omega, x),
 \end{aligned}$$

as stated. \square

PROPOSITION 5.5. *Fix $r > 0$. Then given $\varepsilon > 0$ there exists $\delta > 0$ such that if $x, y \in B_r$ with $d(x, y) < \delta$, then $\|z(t, \omega, x) - z(t, \omega, y)\| \leq \varepsilon$ whenever they are defined.*

Proof. Let $c = \max_i \sum_{j=1}^m \nu_{ij}(-\infty, 0] < 1$. From the continuity of M , given $\varepsilon_0 = \varepsilon(1 - c) > 0$ there exists $0 < \delta < \varepsilon_0$, such that if $x, y \in B_r$ with $d(x, y) < \delta$, then $|M(\omega, y) - M(\omega, x)| < \varepsilon_0$. Therefore, if $x, y \in B_r$ and $x \leq_D y$, from Lemma 5.4 we deduce that $0 \leq D_i z_t(\omega, y) - D_i z_t(\omega, x) < \varepsilon_0$ whenever $d(x, y) < \delta$. The definition of D_i yields

$$\begin{aligned} 0 \leq z_i(t, \omega, y) - z_i(t, \omega, x) &< \varepsilon_0 + \sum_{j=1}^m \int_{-\infty}^0 [z_j(t + s, \omega, y) - z_j(t + s, \omega, x)] dv_{ij}(s) \\ &\leq \varepsilon_0 + \|z_t(\omega, y) - z_t(\omega, x)\|_\infty \sum_{j=1}^m \nu_{ij}(-\infty, 0], \end{aligned}$$

from which we deduce that $\|z_t(\omega, y) - z_t(\omega, x)\|_\infty(1 - c) < \varepsilon_0 = \varepsilon(1 - c)$, that is, $\|z(t, \omega, x) - z(t, \omega, y)\| \leq \varepsilon$ whenever they are defined. The case in which x and y are not ordered follows easily from this one. \square

As a consequence, from the existence of a bounded solution for one of the systems of the family, the boundedness of all solutions is inferred, and this is the case in which hypothesis (F6) holds.

THEOREM 5.6. *Under assumptions (C1)–(C5), if there exists $\omega_0 \in \Omega$ such that (5.2) $_{\omega_0}$ has a bounded solution, then all solutions of (5.2) $_\omega$ are bounded as well, hypothesis (F6) holds, and all omega-limit sets are copies of the base.*

Proof. The boundedness of all solutions is an easy consequence of the previous proposition and the continuity of the semiflow. Let $(\omega, x) \in \Omega \times BU$ and $r' > 0$ such that $z_t(\omega, x) \in B_{r'}$ for all $t \geq 0$. Then also from Proposition 5.5, we deduce that given $\varepsilon > 0$ there exists a $\delta > 0$ such that

$$\|z(t + s, \omega, x) - z(t, \omega \cdot s, y)\| = \|z(t, \omega \cdot s, z_s(\omega, x)) - z(t, \omega \cdot s, y)\| < \varepsilon$$

for all $t \geq 0$ whenever $y \in B_{r'}$ and $d(z_s(\omega, x), y) < \delta$, which shows the uniform stability of the trajectories in $B_{r'}$ for each $r' > 0$. Moreover, for each $r > 0$ there is an $r' > 0$ such that $\hat{D}^{-1}B_r \subset B_{r'}$. Hence, hypothesis (F6) holds for all $r > 0$ and Theorem 4.6 applies for all initial data, which finishes the proof. \square

Concerning the solutions of the original compartmental system, we obtain the following result providing a nontrivial generalization of the autonomous case, in which the asymptotical constancy of the solutions was shown (see Wu and Freedman [28]). Although the theorem is stated in the almost periodic case, similar conclusions are obtained changing almost periodic to periodic, almost automorphic, or recurrent, that is, all solutions are asymptotically of the same type as the transport functions.

THEOREM 5.7. *Under assumptions (C1)–(C5) and in the almost periodic case, if there is a bounded solution of (5.1), then there is at least an almost periodic solution and all the solutions are asymptotically almost periodic. For closed systems, i.e., $\tilde{I}_i \equiv 0$ and $\tilde{g}_{0i} \equiv 0$ for each $i = 1, \dots, m$, there are infinitely many almost periodic solutions and the rest of them are asymptotically almost periodic.*

Proof. The first statement is an easy consequence of the previous theorem. Let $\omega_0 = \tilde{g}$. The omega-limit of each solution $z(t, \omega_0, x_0)$ is a copy of the base $\mathcal{O}(\omega_0, x_0) = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$, and hence $z(t, \omega_0, x(\omega_0)) = x(\omega_0 \cdot t)(0)$ is an almost periodic solution of (5.1) and

$$\lim_{t \rightarrow \infty} \|z(t, \omega_0, x_0) - z(t, \omega_0, x(\omega_0))\| = 0.$$

The statement for closed systems follows in addition from (5.6), which implies that the mass is constant along the trajectories. Hence, there are infinitely many minimal subsets because from the definition of the mass and (C4), given $c > 0$ there is an $(\omega_0, x_0) \in \Omega \times BU_D^+$ such that $M(\omega_0, x_0) = c$ and hence $M(\omega, x) = c$ for each $(\omega, x) \in \mathcal{O}(\omega_0, x_0)$. \square

6. Long-term behavior of compartmental systems. This section deals with the long-term behavior of the amount of material within the compartments of the compartmental system (5.1) satisfying hypotheses (C1)–(C5). As in the previous section, the study of the minimal sets for the corresponding skew-product semiflow (4.2) induced by the family $(5.2)_\omega$ will be essential. In addition to hypotheses (C1)–(C5) we will assume the following hypothesis:

- (C6) Given $i \in \{0, \dots, m\}$ and $j \in \{1, \dots, m\}$ either $\tilde{g}_{ij} \equiv 0$ on $\mathbb{R} \times \mathbb{R}^+$ (and hence $g_{ij} \equiv 0$ on $\Omega \times \mathbb{R}^+$), i.e., *there is not a pipe from compartment C_j to compartment C_i* , or for each $v > 0$ there is a $\delta_v > 0$ such that $\tilde{g}_{ij}(t, v) \geq \delta_v$ for all $t \in \mathbb{R}$ (and hence $g_{ij}(\omega, v) > 0$ for all $\omega \in \Omega$ and $v > 0$). In this case we will say that the pipe P_{ij} carries material (or that there is a pipe from compartment C_j to compartment C_i).

Let $I = \{1, \dots, m\}$. $\mathcal{P}(I)$ denotes, as usual, the set of all subsets of I .

DEFINITION 6.1. Let $\zeta : \mathcal{P}(I) \rightarrow \mathcal{P}(I)$, $J \mapsto \cup_{j \in J} \{i \in I \mid P_{ij} \text{ carries material}\}$. A subset J of I is said to be irreducible if $\zeta(J) \subset J$ and no proper subset of J has that property. System (5.1) is irreducible if the whole set I is irreducible.

Note that $\zeta(I) \subset I$, so there is always some irreducible subset of I . Irreducible sets detect the occurrence of dynamically independent subsystems. Our next result gives a useful property of irreducible sets with more than one element.

PROPOSITION 6.2. If a subset J of I is irreducible, then, for all $i, j \in J$ with $i \neq j$, there exist $p \in \mathbb{N}$ and $i_1, \dots, i_p \in J$ such that $P_{i_1 i}$, $P_{i_2 i_1}, \dots, P_{i_p i_{p-1}}$, and $P_{j i_p}$ carry material.

Proof. Let us assume, on the contrary, that $j \notin \cup_{n=1}^\infty \zeta^n(\{i\}) = \tilde{J}_i$. Then $\tilde{J}_i \subsetneq J$ and, obviously, $\zeta(\tilde{J}_i) \subset \tilde{J}_i$, which contradicts the fact that J is irreducible. \square

Let J_1, \dots, J_k be all the irreducible subsets of I and let $J_0 = I \setminus \cup_{l=1}^k J_l$. These sets reflect the geometry of the compartmental system in a good enough way as to describe the long-term behavior of the solutions, as we will see below.

Let K be any minimal subset of $\Omega \times BU$ for the skew-product semiflow induced by $(5.2)_\omega$. From Theorem 5.6, K is of the form $K = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$, where x is a continuous map from Ω into BU . All of the subsequent results give qualitative information about the long-term behavior of the solutions. Let us see that, provided that we are working on a minimal set K , if there is no inflow from the environment, then the total mass is constant on K , all compartments out of an irreducible subset are empty, and, in an irreducible subset, either all compartments are empty or all are never empty. In particular, in any irreducible subset with some outflow of material, all compartments are empty.

THEOREM 6.3. Assume that $\tilde{I}_i \equiv 0$ for each $i \in I$ and let $K = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ be a minimal subset of $\Omega \times BU$ with $K \subset \Omega \times BU_D^+$. Then the following hold:

- (i) There exists $c \geq 0$ such that $M|_K \equiv c$.
- (ii) $x_i \equiv 0$ for each $i \in J_0$.
- (iii) If, for some $l \in \{1, \dots, k\}$, there exists $j_l \in J_l$ such that $x_{j_l} \equiv 0$, then $x_i \equiv 0$ for each $i \in J_l$. In particular, this happens if there is a $j_l \in J_l$ such that there is an outflow of material from C_{j_l} .

Proof. We first suppose that the system is closed, i.e., $\tilde{g}_{0i} \equiv 0$, $\tilde{I}_i \equiv 0$ for all $i \in I$, from which we deduce $g_{0i} \equiv 0$ and $I_i \equiv 0$ for all $i \in I$.

(i) From (5.6) the total mass M is constant along the trajectories, and hence $M(\omega \cdot t, x(\omega \cdot t)) = M(\omega, x(\omega))$ for all $t \geq 0$ and $\omega \in \Omega$, which together with the fact that Ω is minimal and M continuous shows the statement.

(ii) Let $i \in J_0$. The set $\tilde{J}_i = \cup_{n=1}^\infty \zeta^n(\{i\})$ satisfies $\zeta(\tilde{J}_i) \subset \tilde{J}_i$ and hence contains an irreducible set J_l for some $l \in \{1, \dots, k\}$. Consequently, there are $i_1, \dots, i_p \in J_0$ and $j_l \in J_l$ such that $P_{j_l i_p}$ carry material.

It is easy to prove that there is an $r > 0$ such that $\|x(\omega)\|_\infty \leq r$ for each $\omega \in \Omega$. We define $M_l: \Omega \times BU \rightarrow \mathbb{R}$, the *mass restricted to J_l* , as

$$(6.1) \quad M_l(\omega, y) = \sum_{i \in J_l} D_i y + \sum_{i, j \in J_l} \int_{-\infty}^0 \left(\int_s^0 g_{ji}(\omega \cdot \tau, y_i(\tau)) d\tau \right) d\mu_{ji}(s),$$

which is continuous on $\Omega \times B_r$. From $x(\omega) \geq_D 0$, which also implies $x(\omega) \geq 0$, and (C1), we have $0 \leq M_l(\omega, x(\omega)) \leq M(\omega, x(\omega)) = c$ for each $\omega \in \Omega$.

Since J_l is irreducible, for all $i \in J_l$ and $\omega \in \Omega$

$$\frac{d}{dt} D_i x(\omega \cdot t) = - \sum_{j \in J_l} g_{ji}(\omega \cdot t, x_i(\omega \cdot t)(0)) + \sum_{j \in J_l \cup J_0} \int_{-\infty}^0 g_{ij}(\omega \cdot (s+t), x_j(\omega \cdot t)(s)) d\mu_{ij}(s)$$

because the rest of the terms vanish. Consequently,

$$(6.2) \quad \frac{d}{dt} M_l(\omega \cdot t, x(\omega \cdot t)) = \sum_{i \in J_l} \sum_{j \in J_0} \int_{-\infty}^0 g_{ij}(\omega \cdot (s+t), x_j(\omega \cdot t)(s)) d\mu_{ij}(s) \geq 0$$

for each $\omega \in \Omega$. We claim that $M_l(\omega, x(\omega))$ is constant for each $\omega \in \Omega$. Assume, on the contrary, that there are $\omega_1, \omega_2 \in \Omega$ such that $M_l(\omega_1, x(\omega_1)) < M_l(\omega_2, x(\omega_2))$, and let $t_n \uparrow \infty$ such that $\lim_{n \rightarrow \infty} \omega_2 \cdot t_n = \omega_1$. From (6.2) we deduce that $M_l(\omega_2, x(\omega_2)) \leq M_l(\omega_2 \cdot t_n, x(\omega_2 \cdot t_n))$ for each $n \in \mathbb{N}$, and taking limits as $t \rightarrow \infty$ we conclude that $M_l(\omega_2, x(\omega_2)) \leq M_l(\omega_1, x(\omega_1))$, a contradiction. Hence $M_l(\omega, x(\omega))$ is constant and from (6.2)

$$(6.3) \quad \sum_{i \in J_l} \sum_{j \in J_0} \int_{-\infty}^0 g_{ij}(\omega \cdot (s+t), x_j(\omega \cdot t)(s)) d\mu_{ij}(s) = 0.$$

Next we check that $x_{i_p} \equiv 0$. From (6.3) we deduce that for each $\omega \in \Omega$

$$(6.4) \quad \int_{-\infty}^0 g_{j_l i_p}(\omega \cdot s, x_{i_p}(\omega)(s)) d\mu_{j_l i_p}(s) = 0.$$

Assume that there is an $\omega_0 \in \Omega$ such that $x_{i_p}(\omega_0)(0) > 0$. Hence there is an $\varepsilon > 0$ with $x_{i_p}(\omega_0)(s) > 0$ for each $s \in (-\varepsilon, 0]$, and since $P_{j_l i_p}$ carries material $g_{j_l i_p}(\omega_0 \cdot s, x_{i_p}(\omega_0)(s)) > 0$ for $s \in (-\varepsilon, 0]$. In addition, from $\mu_{j_l i_p}(-\infty, 0] = 1$ there is a $b \leq 0$ such that $\mu_{j_l i_p}(b - \varepsilon, b] > 0$. Hence, denoting $\omega_0 \cdot (-b) = \omega_1$ we deduce that

$$\int_{b-\varepsilon}^b g_{j_l i_p}(\omega_1 \cdot s, x_{i_p}(\omega_1)(s)) d\mu_{j_l i_p}(s) > 0,$$

which contradicts (6.4) and shows that $x_{i_p} \equiv 0$, as claimed. Since $x(\omega) \geq_D 0$, we have $D_{i_p} x(\omega) \geq 0$ and from the definition of D_{i_p} we deduce that $D_{i_p} x(\omega) = 0$ for each

$\omega \in \Omega$. Therefore,

$$0 = \frac{d}{dt} D_{i_p} x(\omega \cdot t) = \sum_{j=1}^m \int_{-\infty}^0 g_{i_p j}(\omega \cdot (t+s), x_j(\omega \cdot t)(s)) d\mu_{i_p j}(s),$$

from which $\int_{-\infty}^0 g_{i_p i_{p-1}}(\omega \cdot s, x_{i_{p-1}}(\omega)(s)) d\mu_{i_p i_{p-1}}(s) = 0$, and as before $x_{i_{p-1}} \equiv 0$. In a finite number of steps we check that $x_i \equiv 0$, as stated.

(iii) From Proposition 6.2, given $i, j_l \in J_l$ there exist $p \in \mathbb{N}$ and $i_1, \dots, i_p \in J_l$ such that $P_{i_1 i}, P_{i_2 i_1}, \dots, P_{i_p i_{p-1}}$, and $P_{j_l i_p}$ carry material. If $x_{j_l} \equiv 0$, the same argument given in the last part of (ii) shows that $x_i \equiv 0$, which finishes the proof for closed systems.

Next we deal with the case when $\tilde{I}_i \equiv 0$ for each $i \in I$ but the system is not necessarily closed. We also have $I_i \equiv 0$ and from (5.5) we deduce that the total mass M is decreasing along the trajectories. In particular,

$$(6.5) \quad \frac{d}{dt} M(\omega \cdot t, x(\omega \cdot t)) = - \sum_{i=1}^n g_{0i}(\omega \cdot t, x_i(\omega \cdot t)(0)) \leq 0.$$

Assume that there are $\omega_1, \omega_2 \in \Omega$ such that $M(\omega_1, x(\omega_1)) < M(\omega_2, x(\omega_2))$, and let $t_n \uparrow \infty$ such that $\lim_{n \rightarrow \infty} \omega_1 \cdot t_n = \omega_2$. From relation (6.5) we deduce that $M(\omega_1 \cdot t_n, x(\omega_1 \cdot t_n)) \leq M(\omega_1, x(\omega_1))$ for each $n \in \mathbb{N}$, and taking limits as $n \uparrow \infty$ we conclude that $M(\omega_2, x(\omega_2)) \leq M(\omega_1, x(\omega_1))$, a contradiction, which shows that M is constant on K , as stated in (i). Consequently, the derivative in (6.5) vanishes and $g_{0i}(\omega \cdot t, x_i(\omega \cdot t)(0)) = 0$ for all $i \in I, \omega \in \Omega$, and $t \geq 0$. This means that $z(t, \omega, x(\omega)) = x(\omega \cdot t)(0)$ is a solution of a closed system, and (ii) and the first part of (iii) follow from the previous case.

Finally, let $j_l \in J_l$ be such that there is an outflow of material from C_{j_l} , that is, $g_{0j_l}(\omega, v) > 0$ for all $\omega \in \Omega$ and $v > 0$. Moreover, as before, $g_{0j_l}(\omega, x_{j_l}(\omega)(0)) = 0$ for each $\omega \in \Omega$, which implies that $x_{j_l} \equiv 0$ and completes the proof. \square

Remark 6.4. Notice that, concerning the solutions of the family of systems (5.2) $_{\omega}$ and hence the solutions of the original system (5.1) when $\omega = \tilde{g}$, we deduce that in the case of no inflow from the environment, $\lim_{t \rightarrow \infty} z_i(t, \omega, x_0) = 0$ for all $i \in J_0, i \in J_l$ for compartments J_l with some outflow, and each $x_0 \geq_D 0$.

Remark 6.5. If there is no inflow from the environment and for all $l \in \{1, \dots, k\}$ there is a $j_l \in J_l$ such that there is outflow of material from C_{j_l} , then the only minimal set in $\Omega \times BU_D^+$ is $K = \{(\omega, 0) \mid \omega \in \Omega\}$ and all the solutions $z(t, \omega, x_0)$ with initial data $x_0 \geq_D 0$ tend to 0 as $t \rightarrow \infty$.

In a nonclosed system, that is, a system which may have any inflow and any outflow of material, if there exists a bounded solution, i.e., all solutions are bounded as shown above, and an irreducible set which has *some inflow*, then, working on a minimal set, all compartments of that irreducible set are nonempty and there must be some outflow from the irreducible set.

THEOREM 6.6. *Assume that there exists a bounded solution of (5.1) and let $K = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ be a minimal subset of $\Omega \times BU_D^+$. If, for some $l \in \{1, \dots, k\}$, there is a $j_l \in J_l$ such that $\tilde{I}_{j_l} \neq 0$, i.e., there is some inflow into C_{j_l} , then*

- (i) $x_i \neq 0$ for each $i \in J_l$, and
- (ii) there is a $j \in J_l$ such that there is outflow of material from C_j .

Proof. (i) Let us assume, on the contrary, that there is an $i \in J_l$ such that $x_i \equiv 0$. Then since $x(\omega) \geq_D 0$ we have that $0 \leq D_i x(\omega)$, and from the definition of D_i given

in (C4) we deduce that $D_i x(\omega) = 0$ for each $\omega \in \Omega$. Therefore,

$$(6.6) \quad 0 = \frac{d}{dt} D_i x(\omega \cdot t) = \sum_{j=1}^m \int_{-\infty}^0 g_{ij}(\omega \cdot (t+s), x_j(\omega \cdot t)(s)) d\mu_{ij}(s) + I_i(\omega \cdot t)$$

for all $\omega \in \Omega$, $t \geq 0$, and, as in (ii) of Theorem 6.3, we check that $x_{j_l} \equiv 0$. However, since $\tilde{I}_{j_l} \not\equiv 0$, there is an $\omega_0 \in \Omega$ such that $I_{j_l}(\omega_0) > 0$, which contradicts (6.6) for $\omega = \omega_0$, $i = j_l$ at $t = 0$.

(ii) Assume, on the contrary, that $g_{0j} \equiv 0$ for each $j \in J_l$. Then if we consider (6.1) the restriction of the mass to J_l , we check that

$$\frac{d}{dt} M_l(\omega \cdot t, x(\omega \cdot t)) = \sum_{i \in J_l} \left[I_i(\omega \cdot t) + \sum_{j \in J_0} \int_{-\infty}^0 g_{ij}(\omega \cdot (s+t), x_j(\omega \cdot t)(s)) d\mu_{ij}(s) \right] \geq 0$$

for all $\omega \in \Omega$ and $t \geq 0$. A similar argument to the one given in (ii) of Theorem 6.3 shows that $M_l(\omega, x(\omega))$ is constant for each $\omega \in \Omega$, which contradicts the fact that the above derivative is strictly positive for $\omega = \omega_0$ at $t = 0$ and proves the statement. \square

Finally, we will change hypothesis (C6) to the following, slightly stronger one.

(C6)* Given $i \in \{0, \dots, m\}$ and $j \in \{1, \dots, m\}$ either $\tilde{g}_{ij} \equiv 0$ on $\mathbb{R} \times \mathbb{R}^+$ (and hence $g_{ij} \equiv 0$ on $\Omega \times \mathbb{R}^+$), i.e., there is not a pipe from compartment C_j to compartment C_i , or for each $v \geq 0$ there is a $\delta_v > 0$ such that $\frac{\partial}{\partial v} \tilde{g}_{ij}(t, v) \geq \delta_v$ for each $t \in \mathbb{R}$ (and hence $\frac{\partial}{\partial v} g_{ij}(\omega, v) > 0$ for all $\omega \in \Omega$ and $v \geq 0$). In this case we will say that the pipe P_{ij} carries material (or that there is a pipe from compartment C_j to compartment C_i).

In this case, we are able to prove that if there exists a bounded solution, then all the minimal sets coincide both on irreducible sets having some outflow and out of irreducible sets. Concerning the solutions of the initial compartmental system (5.1),

$$\lim_{t \rightarrow \infty} |z_i(t, x_0) - z_i(t, y_0)| = 0$$

for all $i \in J_0$, $i \in J_l$ for compartments J_l with some outflow, and all $x_0, y_0 \geq_D 0$.

THEOREM 6.7. *Let us assume that hypotheses (C1)–(C5) and (C6)* hold and that there exists a bounded solution of system (5.1). Let $K_1 = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ and $K_2 = \{(\omega, y(\omega)) \mid \omega \in \Omega\}$ be two minimal subsets of $\Omega \times BU_D^+$. Then*

- (i) $x_i \equiv y_i$ for each $i \in J_0$;
- (ii) if, for some $l \in \{1, \dots, k\}$, there is a $j_l \in J_l$ such that there is outflow of material from C_{j_l} , then $x_i \equiv y_i$ for each $i \in J_l$.

Proof. For each $i \in \{0, \dots, m\}$ and each $j \in \{1, \dots, m\}$ we define $h_{ij} : \Omega \rightarrow \mathbb{R}^+$ as

$$h_{ij}(\omega) = \int_0^1 \frac{\partial g_{ij}}{\partial v}(\omega, s x_j(\omega)(0) + (1-s) y_j(\omega)(0)) ds \geq 0,$$

and we consider the family of monotone linear compartmental systems

$$(6.7)_\omega \quad \begin{aligned} \frac{d}{dt} D_i \hat{z}_i &= -h_{0i}(\omega \cdot t) \hat{z}_i(t) - \sum_{j=1}^m h_{ji}(\omega \cdot t) \hat{z}_j(t) \\ &+ \sum_{j=1}^m \int_{-\infty}^0 h_{ij}(\omega \cdot (s+t)) \hat{z}_j(t+s) d\mu_{ij}(s), \quad \omega \in \Omega, \end{aligned}$$

satisfying the corresponding hypotheses (C1)–(C4) and (C6). Moreover, (C5) for each of the systems $(6.7)_\omega$, follows from

$$\inf_{\omega \in \Omega} h_{ij}(\omega) \geq \inf_{v \geq 0, \omega \in \Omega} \frac{\partial g_{ij}}{\partial v}(\omega, v), \quad \sup_{\omega \in \Omega} h_{ij}(\omega) \leq \sup_{v \geq 0, \omega \in \Omega} \frac{\partial g_{ij}}{\partial v}(\omega, v),$$

and (C5) for (5.1). From the definition of h_{ij} and (C6)* we deduce that the irreducible sets for the families $(6.7)_\omega$ and $(5.2)_\omega$ coincide. Consequently, Theorem 6.3 (see Remark 6.4) applies to this case, and we deduce that if $z_0 \geq_D 0$ and J_l is a compartment with some outflow of material, then

$$\lim_{t \rightarrow \infty} \hat{z}_i(t, \omega, z_0) = 0 \quad \text{for each } i \in J_0 \cup J_l.$$

The same happens for $z_0 \leq_D 0$ because the systems are linear.

Let $z(\omega) = x(\omega) - y(\omega)$ for each $\omega \in \Omega$. It is easy to check $\hat{z}(t, \omega, z(\omega)) = z(\omega \cdot t)(0)$ for all $\omega \in \Omega$ and $t \geq 0$. Moreover, we can find $z_1 \leq_D 0$ and $z_0 \geq_D 0$ such that $z_1 \leq_D z(\omega) \leq_D z_0$ for each $\omega \in \Omega$. Hence, the monotonicity of the induced skew-product semiflow and the positivity of \hat{D}^{-1} yields

$$\hat{z}(t, \omega, z_1) \leq z(\omega \cdot t)(0) \leq \hat{z}(t, \omega, z_0) \quad \text{for all } \omega \in \Omega \ t \geq 0,$$

from which we deduce that $z_i \equiv 0$ for all $i \in J_0, i \in J_l$ and (i) and (ii) follow. \square

As a consequence, under the same assumptions of the previous theorem, when for all $l \in \{1, \dots, k\}$ there is an outflow of material from one of the compartments in J_l , there is a unique minimal set $K = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ in $\Omega \times BU_D^+$ attracting all the solutions with initial data in BU_D^+ ; i.e.,

$$\lim_{t \rightarrow \infty} \|z(t, \omega, x_0) - x(\omega \cdot t)(0)\| = 0, \quad \text{whenever } x_0 \geq_D 0.$$

Moreover, $x \not\equiv 0$ if and only if there is some $j \in \{1, \dots, m\}$ such that $\tilde{I}_j \neq 0$; i.e., there is some inflow into one of the compartments C_j .

For the next result, in addition to hypotheses (C1)–(C5) and (C6)* we will assume the following hypothesis:

(C7) If $K_1 = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ and $K_2 = \{(\omega, y(\omega)) \mid \omega \in \Omega\}$ are two minimal subsets of $\Omega \times BU_D^+$ such that $x(\omega) \leq_D y(\omega)$ and $D_i x(\omega_0) = D_i y(\omega_0)$ for some $\omega_0 \in \Omega$ and $i \in \{1, \dots, m\}$, then $x(\omega) = y(\omega)$ for each $\omega \in \Omega$, i.e., $K_1 = K_2$.

Note that if $D_i x(\omega_0) = D_i y(\omega_0)$ holds for some $\omega_0 \in \Omega$ and $i \in \{1, \dots, m\}$, then from hypothesis (F5) we deduce that it holds for each $\omega \in \Omega$.

Hypothesis (C7) is relevant when it applies to closed systems, and it holds in many cases studied in the literature. A closed system satisfying (C7) is irreducible. Systems with a unique compartment, studied by Arino and Bourad [1] and Krisztin and Wu [19], satisfy (C7). It follows from Theorem 6.3 that irreducible closed systems described by FDEs (see Arino and Haourigui [2]) satisfy (C7). Closed systems given by Wu [27], and Wu and Freedman [28] in the strongly ordered case, also satisfy (C7).

DEFINITION 6.8. Let $K_1 = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ and $K_2 = \{(\omega, y(\omega)) \mid \omega \in \Omega\}$ be two minimal subsets. It is said that $K_1 <_D K_2$ if $x(\omega) <_D y(\omega)$ for each $\omega \in \Omega$.

Hypothesis (C7) allows us to classify the minimal subsets in terms of the value of their total mass, as shown in the next result.

THEOREM 6.9. Assume that system (5.1) is closed (i.e., $\tilde{I}_i \equiv 0$ and $\tilde{g}_{0i} \equiv 0$ for each $i \in \{1, \dots, m\}$), and hypotheses (C1)–(C5), (C6)*, and (C7) hold. Then for

each $c > 0$ there is a unique minimal subset K_c such that $M|_{K_c} = c$. Moreover, $K_c \subset \Omega \times BU_D^+$ and $K_{c_1} <_D K_{c_2}$ whenever $c_1 < c_2$.

Proof. Since the minimal subsets are copies of the base, and the total mass (5.4) is constant along the trajectories and increasing for the D -order because \widehat{D}^{-1} is positive, it is easy to check that given $c > 0$ there is a minimal subset $K_c \subset \Omega \times BU_D^+$ such that $M|_{K_c} = c$.

Let \widehat{D} be the isomorphism of BU defined by the relation (3.6). For each $x \in BU$ we define $x^+ = \widehat{D}^{-1} \sup(0, \widehat{D}x)$. Hence $0 \leq_D x^+$, $x \leq_D x^+$, and if $y \in BU$ with $x \leq_D y$ and $0 \leq_D y$, then $x^+ \leq_D y$.

Since the semiflow is monotone, from $x \leq_D x^+$ we deduce that $u(t, \omega, x) \leq_D u(t, \omega, x^+)$. Since the system is closed, $u(t, \omega, 0) = 0$, and from $0 \leq_D x^+$ we check that $0 \leq_D u(t, \omega, x^+)$. Consequently $u(t, \omega, x)^+ \leq_D u(t, \omega, x^+)$ for each $t \geq 0$.

Next we check that if $K = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ is minimal, the same happens for $K^+ = \{(\omega, x(\omega)^+) \mid \omega \in \Omega\}$. Since $x(\omega \cdot t) = u(t, \omega, x)$ for each $t \geq 0$, we deduce that $x(\omega \cdot t)^+ = u(t, \omega, x(\omega))^+ \leq_D u(t, \omega, x(\omega)^+)$, and the fact that \widehat{D}^{-1} is positive yields $x(\omega \cdot t)^+ \leq u(t, \omega, x(\omega)^+)$ for each $t \geq 0$. In addition, since the total mass (5.4) is constant along the trajectories and increasing for the D -order, we deduce that $M(\omega, x(\omega)^+) = M(\omega \cdot t, u(t, \omega, x(\omega)^+)) \geq M(\omega \cdot t, u(t, \omega, x(\omega))^+) = M(\omega \cdot t, x(\omega \cdot t)^+)$ for each $t \geq 0$. Moreover, since $x(\omega)^+$ is a continuous function in ω and Ω is minimal, a similar argument to the one given in (ii) of Theorem 6.3 shows that $M(\omega, x(\omega)^+)$ is constant on Ω and, consequently, $M(\omega \cdot t, u(t, \omega, x(\omega)^+)) = M(\omega \cdot t, x(\omega \cdot t)^+)$ for each $\omega \in \Omega$ and $t \geq 0$. Hence, from (5.4) we conclude that

$$0 = \sum_{i=1}^m D_i(u(t, \omega, x(\omega)^+) - x(\omega \cdot t)^+),$$

that is, $D(u(t, \omega, x(\omega)^+)) = D(x(\omega \cdot t)^+)$ for each $\omega \in \Omega$ and $t \geq 0$. In addition, it is easy to check that $(\varphi_s)^+ = (\varphi^+)_s$ whenever $\varphi \in BU$ and $s \leq 0$, from which we deduce that $D((u(t, \omega, x(\omega)^+))_s) = D((x(\omega \cdot t)^+)_s)$ for each $s \leq 0$, $t \geq 0$, and $\omega \in \Omega$. That is, $\widehat{D}(u(t, \omega, x(\omega)^+)) = \widehat{D}(x(\omega \cdot t)^+)$ for each $t \geq 0$ and $\omega \in \Omega$, and since \widehat{D} is an isomorphism $u(t, \omega, x(\omega)^+) = x(\omega \cdot t)^+$ for each $t \geq 0$ and $\omega \in \Omega$, which shows that K^+ is a minimal subset, as stated. Let $K_1 = \{(\omega, x(\omega)) \mid \omega \in \Omega\}$ and $K_2 = \{(\omega, y(\omega)) \mid \omega \in \Omega\}$ be two minimal subsets such that $M|_{K_i} = c$ for $i = 1, 2$. We fix $\omega \in \Omega$. The change of variable $\widehat{z}(t) = z(t) - y(\omega \cdot t)$ takes (5.2) $_\omega$ to

$$\frac{d}{dt} D\widehat{z}_t = G(\omega \cdot t, \widehat{z}_t), \quad t \geq 0, \omega \in \Omega,$$

where $G(\omega \cdot t, \widehat{z}_t) = F(\omega \cdot t, \widehat{z}_t + y(\omega \cdot t)) - F(\omega \cdot t, y(\omega \cdot t))$. It is not hard to check that this is a new family of compartmental systems satisfying the corresponding hypotheses (C1)–(C5) and (C6)*, and

$$\widehat{K} = \{(\omega, x(\omega) - y(\omega)) \mid \omega \in \Omega\}$$

is one of its minimal subsets. As before

$$\widehat{K}^+ = \{(\omega, (x(\omega) - y(\omega))^+) \mid \omega \in \Omega\}$$

is also a minimal subset, and hence

$$K^+ = \{(\omega, y(\omega) + (x(\omega) - y(\omega))^+) \mid \omega \in \Omega\} = \{(\omega, z(\omega)) \mid \omega \in \Omega\}$$

is a minimal set for the initial family. For each $\omega \in \Omega$ we have $z(\omega) \geq_D y(\omega)$.

Let us assume that $Dz(\omega) \gg Dy(\omega)$ for each $\omega \in \Omega$, which implies that $D((x(\omega) - y(\omega))^+) \gg 0$ for each $\omega \in \Omega$. Consequently, $D((x(\omega) - y(\omega))_s^+) = D(((x(\omega) - y(\omega))_s^+)_s^+) = D((x(\omega \cdot s) - y(\omega \cdot s))^+) \gg 0$ for each $s \leq 0$, and we deduce that $\widehat{D}x(\omega) > \widehat{D}y(\omega)$, i.e., $x(\omega) >_D y(\omega)$ for each $\omega \in \Omega$, and $M|_{K_1} > M|_{K_2}$, a contradiction. Hence, there are an $\omega_0 \in \Omega$ and an $i \in \{1, \dots, m\}$ such that $D_i z(\omega_0) = D_i y(\omega_0)$, and hypothesis (C7) provides that $z(\omega) = y(\omega)$ for each $\omega \in \Omega$. That is, $(x(\omega) - y(\omega))^+ \equiv 0$ for each $\omega \in \Omega$, or equivalently $x(\omega) - y(\omega) \leq_D 0$ for each $\omega \in \Omega$. Finally, as before, from $M|_{K_1} = M|_{K_2}$ we conclude by contradiction that $x(\omega) = y(\omega)$ for each $\omega \in \Omega$, and the minimal set K_c is unique, as stated. The same argument shows that $K_{c_1} <_D K_{c_2}$ whenever $c_1 < c_2$ and finishes the proof. \square

REFERENCES

- [1] O. ARINO AND F. BOURAD, *On the asymptotic behavior of the solutions of a class of scalar neutral equations generating a monotone semiflow*, J. Differential Equations, 87 (1990), pp. 84–95.
- [2] O. ARINO AND E. HAOURIGUI, *On the asymptotic behavior of solutions of some delay differential systems which have a first integral*, J. Math. Anal. Appl., 122 (1987), pp. 36–46.
- [3] R. ELLIS, *Lectures on Topological Dynamics*, W. A. Benjamin, New York, 1969.
- [4] G. GRIPENBERG, S.-O. LONDEN, AND O. STAFFANS, *Volterra Integral and Functional Equations*, Encyclopedia Math. Appl., Cambridge University Press, Cambridge, New York, 1990.
- [5] I. GYÖRI, *Connections between compartmental systems with pipes and integro-differential equations*, Math. Modelling, 7 (1986), pp. 1215–1238.
- [6] I. GYÖRI AND J. ELLER, *Compartmental systems with pipes*, Math. Biosci., 53 (1981), pp. 223–247.
- [7] I. GYÖRI AND J. WU, *A neutral equation arising from compartmental systems with pipes*, J. Dynam. Differential Equations, 3 (1991), pp. 289–311.
- [8] J. R. HADDOCK, T. KRISZTIN, AND J. WU, *Asymptotic equivalence of neutral and infinite retarded differential equations*, Nonlinear Anal., 14 (1990), pp. 369–377.
- [9] J. R. HADDOCK, T. KRISZTIN, J. TERJÉKI, AND J. WU, *An invariance principle of Lyapunov-Razumikhin type for neutral functional differential equations*, J. Differential Equations, 107 (1994), pp. 395–417.
- [10] J. K. HALE, *Theory of Functional Differential Equations*, Appl. Math. Sci. 3, Springer-Verlag, Berlin, Heidelberg, New York, 1977.
- [11] J. K. HALE AND K. R. MEYER, *A Class of Functional Equations of Neutral Type*, Mem. Amer. Math. Soc. 76, AMS, Providence, RI, 1967.
- [12] J. K. HALE AND S. M. VERDUYN LUNEL, *Introduction to Functional Differential Equations*, Appl. Math. Sci. 99, Springer-Verlag, Berlin, Heidelberg, New York, 1993.
- [13] Y. HINO, S. MURAKAMI, AND T. NAIKO, *Functional-Differential Equations with Infinite Delay*, Lecture Notes in Math. 1473, Springer-Verlag, Berlin, Heidelberg, 1991.
- [14] J. A. JACQUEZ, *Compartmental Analysis in Biology and Medicine*, 3rd ed., Thomson-Shore Inc., Ann Arbor, MI, 1996.
- [15] J. A. JACQUEZ AND C. P. SIMON, *Qualitative theory of compartmental systems*, SIAM Rev., 35 (1993), pp. 43–79.
- [16] J. A. JACQUEZ AND C. P. SIMON, *Qualitative theory of compartmental systems with lags*, Math. Biosci., 180 (2002), pp. 329–362.
- [17] J. JIANG AND X.-Q. ZHAO, *Convergence in monotone and uniformly stable skew-product semiflows with applications*, J. Reine Angew. Math, 589 (2005), pp. 21–55.
- [18] V. KOLMANOVSKII AND A. MYSHKIS, *Introduction to the Theory and Applications of Functional Differential Equations*, Math. Appl., Kluwer Academic, Dordrecht, The Netherlands, 1999.
- [19] T. KRISZTIN AND J. WU, *Asymptotic periodicity, monotonicity, and oscillation of solutions of scalar neutral functional-differential equations*, J. Math. Anal. Appl., 199 (1996), pp. 502–525.
- [20] S. NOVO, R. OBAYA, AND A. M. SANZ, *Stability and extensibility results for abstract skew-product semiflows*, J. Differential Equations, 235 (2007), pp. 623–646.
- [21] R. J. SACKER AND G. R. SELL, *Lifting Properties in Skew-Products Flows with Applications to Differential Equations*, Mem. Amer. Math. Soc. 190, AMS, Providence, RI, 1977.

- [22] D. SALAMON, *Control and Observation of Neutral Systems*, Pitman, London, 1984.
- [23] W. SHEN AND Y. YI, *Almost Automorphic and Almost Periodic Dynamics in Skew-Product Semiflows*, Mem. Amer. Math. Soc. 647, AMS, Providence, RI, 1998.
- [24] O. J. STAFFANS, *A neutral FDE with stable D-operator is retarded*, J. Differential Equations, 49 (1983), pp. 208–217.
- [25] O. J. STAFFANS, *On a neutral functional differential equation in a fading memory space*, J. Differential Equations, 50 (1983), pp. 183–217.
- [26] Z. WANG AND J. WU, *Neutral functional differential equations with infinite delay*, Funkcial. Ekvac., 28 (1985), pp. 157–170.
- [27] J. WU, *Unified treatment of local theory of NFDEs with infinite delay*, Tamkang J. Math., 22 (1991), pp. 51–72.
- [28] J. WU AND H. I. FREEDMAN, *Monotone semiflows generated by neutral functional differential equations with application to compartmental systems*, Canad. J. Math., 43 (1991), pp. 1098–1120.

COARSENING IN NONLOCAL INTERFACIAL SYSTEMS*

DEJAN SLEPČEV†

Abstract. We consider coarsening in interfacial systems driven by nonlocal energies. Of particular interest are the nonlocal Cahn–Hilliard equation and models of biological aggregation. The energies considered cause the system to separate into phases. The pattern of interfaces evolves under nonlocal surface-tension-type effects. The typical length scales grow and the pattern coarsens. We prove a rigorous upper bound on the coarsening rate. The proof uses the energy-based approach to estimates on rate of coarsening introduced by Kohn and Otto [*Comm. Math. Phys.*, 229 (2002), pp. 375–395]. To show the required estimates on the flatness of the energy landscape we develop a geometric approach which is applicable to a wider class of problems, which includes ones based on local, gradient-type energies.

Key words. coarsening, interpolation inequality, nonlocal equation, phase segregation, biological aggregation, Wasserstein distance

AMS subject classifications. 45K05, 35B99, 82C24, 92D25

DOI. 10.1137/080713598

1. Introduction. Our main focus is systems driven by nonlocal energies. Of particular interest are nonlocal Cahn–Hilliard-type equations and equations modeling biological aggregation. The nonlocal Cahn–Hilliard equations that we consider were derived by Giacomini and Lebowitz [8] as limits of the lattice-gas dynamics modeling phase segregation in binary alloys. In this setting, they represent a refinement in modeling over the standard, fourth order Cahn–Hilliard equations. The equations modeling biological aggregation were derived by Topaz, Bertozzi, and Lewis [18].

The mathematical descriptions of these systems are rather similar. Both equations are gradient flows of the same general energy:

$$(1) \quad E(u) := \iint (u(x) - u(y))^2 K(x - y) dx dy + \int W(u(x)) dx.$$

Here $K \geq 0$ is the interaction kernel, and W is a double-well potential whose minima are at 0 and 1. We leave the domain of integration vague at the moment. Heuristically it is convenient to consider the domains to be $\mathbb{R}^N \times \mathbb{R}^N$ and \mathbb{R}^N , respectively. However, for technical reasons, when stating and proving rigorous results we consider the problem on a finite domain.

The equations we study are gradient flows of the energy, in the appropriate metrics:

$$(2) \quad u_t - \nabla \cdot \left(\mu(u) \nabla \left(\frac{\delta E}{\delta u} \right) \right) = 0.$$

That is,

$$(3) \quad u_t - \nabla \cdot \left(\mu(u) \nabla \left(4 \int K(y) dy u - 4K * u + W'(u) \right) \right) = 0.$$

*Received by the editors January 17, 2008; accepted for publication (in revised form) June 10, 2008; published electronically October 13, 2008. The research of this work was partially supported by NSF grant DMS-0638481 and by the Center for Nonlinear Analysis through NSF grants DMS-0405343 and DMS-0635983.

<http://www.siam.org/journals/sima/40-3/71359.html>

†Department of Mathematical Sciences, Carnegie Mellon University, Pittsburgh, PA 15213 (slepcev@math.cmu.edu).

For both equations the mobility μ is a nonnegative function. More precisely, for the aggregation equation $\mu(u) = u$, while for the nonlocal Cahn–Hilliard equation $\mu > 0$ on $[0, 1]$.

The second term of the energy causes the system to separate into phases, while the first term penalizes the existence of interfaces. The energy E is a nonlocal counterpart of the energy

$$(4) \quad E_{loc}(u) := \int \frac{1}{2} |\nabla u(x)|^2 + W(u(x)) dx.$$

Roughly speaking, both of the energies measure interfacial area. The longer the length scale in the system, the better the approximation to the interfacial area. More precisely, for both energies, the Γ -limit of the appropriately rescaled energy is the functional measuring perimeter of the set occupied by one of the phases

$$E \xrightarrow{\Gamma} \text{const. } E_{per}.$$

E_{per} is defined for BV functions with the range $\{0, 1\}$. For E_{loc} this is the result of Modica and Mortola [15] (see also [14]), while for E it was proven by Alberti et al. [2] (see also [1]). Moreover matched asymptotics arguments (by Giacomini and Lebowitz [9] and by Bertozzi and Slepčev [4]) show that the sharp-interface limits of the dynamics described by (2) are the Mullins–Sekerka (MS) equation for the nonlocal Cahn–Hilliard equation and the Hele–Shaw (HS) equation for the aggregation model.

After the interfaces have formed, the system slowly evolves, reducing the interfacial area. During this process the length scales that characterize the coarseness of a configuration grow. We are interested in the rate at which these length scales grow—the rate of coarsening. The fact that the sharp interface limits, (MS) and (HS), are both invariant under the scaling $x \rightarrow \lambda x$, $t \rightarrow \lambda^3 t$ suggests that the typical length scale grows as $t^{1/3}$. We prove a weak formulation of this statement, following the technique of Kohn and Otto [10], who proved the result for gradient flows of local energies (4). We use the energy as the measure of the coarseness of the system. In particular let \bar{E} be the energy density, that is, the energy per unit volume. Note that \bar{E} has units of $1/\text{length}$. We show a weak version of the statement

$$\bar{E} \gtrsim t^{-1/3}.$$

This provides an upper bound on rate of coarsening as it shows that the interfacial area cannot decay faster than the given rate.

Outline. In the remainder of the introduction we discuss the gradient-flow structure of the equations, and the framework for obtaining rigorous result on coarsening rates introduced by Kohn and Otto. We also introduce the two applications we have in mind in more detail. In section 2 we list the assumptions needed and give the precise formulation of the main results on the rate of coarsening. In particular we present a relaxed formulation of the abstract Kohn–Otto framework that is applicable to weak solutions of gradient-flow equations. We also show that the “geodesic distance” associated to the metric of the configuration space can be utilized even when the mobility is nonlinear. In section 3 we present the proof of the bound on the rate of coarsening. The main technical ingredient, the interpolation inequality, is proved in section 4. The approach we take in proving the interpolation inequality is general; essentially the same proof covers both types of mobilities and both nonlocal and local energies. We illustrate the application to local energies in subsection 4.1.

1.1. Gradient flow structure. We now introduce the geometric structure of (2). It is based on the formal Riemannian viewpoint developed by Otto [16]. Equation (2) can be understood as a gradient flow of the energy (1) on the manifold of configurations. Since the equation is in divergence form it preserves the integral of u over the space. Thus the solution of the equation is a path on the manifold of functions with the same integral.

At each point the tangent space is the set of possible perturbations, all of which have mean zero. The local metric is defined as follows: Let s_1, s_2 be two tangent vectors at u . Then

$$(5) \quad g_u(s_1, s_2) = \int \mu(u) \nabla p_1 \cdot \nabla p_2,$$

where

$$-\nabla \cdot (\mu(u) \nabla p_i) = s_i \quad \text{for } i = 1, 2.$$

Equation (2) is the gradient flow of energy (1) with respect to the metric (5), that is, for every tangent vector s

$$g(u_t, s) = -\frac{\delta E}{\delta u}[s].$$

Considering the configuration space as a manifold enables us to measure the steepness of the energy landscape in a certain sense. This information provides bounds on the speed of the dynamics.

In particular the local metric gives rise to a global metric on the manifold. Given a regular enough path, $v(s)$ for, say, $s \in [0, 1]$, on the manifold, we can measure its length

$$\text{length}(v) = \int_0^1 \sqrt{g_{v(s)}(v', v')} ds.$$

We can then define the global metric on the manifold: Let the distance of u_1 and u_2 be

$$d(u_1, u_2) = \inf\{\text{length}(v) : v \text{ is a path connecting } u_1 \text{ and } u_2\}.$$

It turns out that when $\mu = \text{const.}$, then $d(u_1, u_2)$ is a multiple of the H^{-1} norm, while for $\mu(u) = u$ the distance becomes the Wasserstein metric.

1.2. Kohn–Otto framework. Kohn and Otto [10] introduced an approach for obtaining information on the flatness of the energy landscape, and consequently on the rate of coarsening. The approach is robust and has been applied to studies of coarsening in epitaxial growth, mean-field models, thin-liquid films, and other systems [6, 7, 5, 11, 12, 17]. See also [13] for a related result.

We first present it in the abstract setting used in [17], which applies to gradient flows. Consider an energy E on a Riemannian manifold (\mathcal{M}, g) . The metric g introduces a global distance on \mathcal{M} , which we denote by d .

PROPOSITION 1. *Let $h^* \in \mathcal{M}$. Let $h : \mathbb{R}_+ \rightarrow \mathcal{M}$ be a solution of*

$$(6) \quad h_t = -\text{grad}E(h),$$

and let $h(0) = h_0$.

Assume that for some $\alpha \geq 0$ the interpolation inequality

$$(7) \quad E(h) \operatorname{dist}(h, h^*)^\alpha \geq 1 \quad \text{for all } h \in \mathcal{M} \text{ with } E(h) \leq \varepsilon$$

holds. Then for $\sigma \in (1, 1 + \frac{2}{\alpha})$

$$(8) \quad \int_0^T E(h(t))^\sigma dt \gtrsim \int_0^T (t^{-\frac{\alpha}{\alpha+2}})^\sigma dt,$$

provided $T \gg \operatorname{dist}(h_0, h^*)^{\alpha+2}$ and $E(h(0)) \leq \varepsilon$.

Remark 1. The precise meaning of \gtrsim and \gg is the following: For all $\sigma \in (1, 1 + \frac{2}{\alpha})$ there exists a constant $C = C(\alpha, \sigma)$ such that for all $\delta > 0$ there exists $C_\delta = C(\alpha, \sigma, \delta)$:

$$(9) \quad \int_0^T E(h(t))^\sigma dt \geq (1 - \delta)C \int_0^T (t^{-\frac{\alpha}{\alpha+2}})^\sigma dt,$$

provided $T \geq C_\delta \operatorname{dist}(h_0, h^*)^{\alpha+2}$.

Proof of the proposition is based on the ODE arguments of [10] and can be found in [17].

We adapt Proposition 1 to our setting in section 2 (Theorem 2). In particular since the Riemannian structure is only formal, the theorem states precisely which elements of the geometric structure are needed. Furthermore it applies to weak solutions of the gradient flow equations.

The main technical difficulty in applying the proposition is showing the interpolation inequality. Section 4 is devoted to proving interpolation inequalities relevant to nonlocal energies.

1.3. Nonlocal Cahn–Hilliard equation. Equation (2) is a rescaled version of the model by Giacomini and Lebowitz [8, 9]. We introduce it in original variables below and, for completeness, present the rescaling needed.

The free energy is given by

$$\mathcal{E} = \frac{1}{4} \int_\Omega \int_{\mathbb{R}^N} K(x - y) (\rho(x) - \rho(y))^2 dx dy + \int_\Omega f_c(\rho(x)) dx.$$

The associated gradient flow is

$$\rho_t - \nabla \cdot \left(\sigma(\rho) \nabla \left(\frac{\delta \mathcal{E}}{\delta \rho} \right) \right) = 0.$$

Here $K \geq 0$ is a smooth kernel with symmetry $K(x) = K(-x)$. Giacomini and Lebowitz assume that K is compactly supported, with support contained in Ω . This assumption is physically quite reasonable, but it is not necessary from the mathematical point of view. The function f_c is a double-well potential, symmetric about $\frac{1}{2}$ with minima at $\frac{1}{2} \pm m$. The mobility function σ is assumed to be smooth, symmetric about $\frac{1}{2}$, and positive on $(0, 1)$, with

$$(GL1) \quad \sigma(0) = 0 \quad \text{and} \quad \sigma(1) = 0.$$

Let

$$f(\rho) = f_c(\rho) + \frac{\int_{\mathbb{R}^N} K(x) dx}{2} \left(\rho - \frac{1}{2} \right)^2.$$

Giacomin and Lebowitz assume that for some $c > 0$ and for all $\rho \in (0, 1)$

$$(GL2) \quad \frac{1}{c} \leq \sigma(\rho) f''(\rho) \leq c.$$

This assumption is needed for the existence/uniqueness theory they use (it makes the equation uniformly parabolic), but is not directly required for the coarsening estimates.

Under the assumptions above, Giacomin and Lebowitz show that for initial datum $0 < \rho_0 < 1$, there exists a unique weak solution $\rho \in L^2([0, T], H^1(\Omega))$ with $\rho_t \in L^2([0, T], H^{-1}(\Omega))$ for any $T > 0$. Furthermore $0 < \rho < 1$.

Rescaling. We rescale the dependent variable, the potential, and the mobility so that the wells of the new potential are 0 and 1. Let

$$\begin{aligned} u &= \frac{1}{2m} \left(\rho - \frac{1}{2} \right) + \frac{1}{2}, \\ W_m(u) &= \frac{1}{4m^2} f_c(\rho) = \frac{1}{4m^2} f_c \left(2m \left(u - \frac{1}{2} \right) + \frac{1}{2} \right), \\ \mu_m(u) &= \sigma(\rho) = \sigma \left(2m \left(u - \frac{1}{2} \right) + \frac{1}{2} \right). \end{aligned}$$

Under this rescaling u solves (2) and is hence a gradient flow of (1).

1.4. Biological aggregation. Topaz, Bertozzi, and Lewis [18] introduced a model of biological aggregation that emerges due to “social forces” between individuals. That is, the individuals are attracted to other individuals of their species, but avoid overcrowding. The population is modeled by its density u . The velocity of individuals is modeled as

$$v = v_a + v_r = \nabla(K * u) - \nabla g(u),$$

where $v_a = \nabla(K * u)$ is the term modeling attraction to other individuals which are being sensed through the kernel K . The term modeling repulsion, v_r , is given by a local operator $v_r = -\nabla g(u)$, where g is an increasing function. The continuity equation then reads

$$u_t + \nabla \cdot (u v) = u_t + \nabla \cdot (u \nabla(K * u - g(u))) = 0.$$

From a biological perspective it is reasonable to assume that $g'(0) = 0$ and g is strictly convex. However, it is sufficient to assume that

$$(BA) \quad \begin{aligned} &\text{the function } g'(z) - \int_{\mathbb{R}^N} K(x) dx \text{ has exactly one zero on } \mathbb{R}^+, \\ &g'(0) < \int_{\mathbb{R}^N} K(x) dx \text{ and } \liminf_{z \rightarrow \infty} g'(z) - \int_{\mathbb{R}^N} K(x) dx > 0. \end{aligned}$$

Under this assumption u solves (3) for some double-well potential W . More precisely, let $G(z) := \int_0^z g(s) ds$ and $\tilde{W}(z) := G(z) - \frac{1}{2} \int_{\mathbb{R}^N} K(x) dx z^2$. The condition (BA) implies that \tilde{W} is concave at 0, and has exactly one inflection point. Thus we can define

$$W(z) := 4 \left(\tilde{W}(z) - \left(\min_{s>0} \frac{\tilde{W}(s)}{s} \right) z \right).$$

Then W is a double-well potential on $[0, \infty)$ with one well at 0. It follows that

$$u_t - \nabla \cdot \left(u \nabla \left(\int_{\mathbb{R}^N} K(y) dy u - K * u + \frac{1}{4} W'(u) \right) \right) = 0,$$

which after scaling the time by factor 4 is (3).

On the level of the model, the equation provides information on why herds (or other animal groups) form, why they have an almost constant density, why they have sharp boundaries, and how they evolve. Numerical simulations conducted in one dimension in [18] also observe the coarsening phenomenon. The primary driving force for coarsening in one dimension is the nonlocal interaction via kernel K , as there are no surface-tension-like effect. The rate of coarsening depends on the decay of K at ∞ . Nevertheless the rigorous bounds we prove still apply and are in fact optimal for certain kernels.

2. Statement of the result. When thinking about coarsening we have in mind an infinite domain on which coarsening persists for all time. However, building the theory for such solutions poses major challenges. We instead consider domains of finite size and prove results that are independent of the domain size. In particular we consider the domain $\Omega = [0, \Lambda]^N$. We investigate the dynamics of periodic configurations on \mathbb{R}^N with period cell Ω . Thus we consider Ω with the topology of the torus $\mathbb{R}^N / (\Lambda\mathbb{Z})^N$. In particular the distances on Ω are measured on the torus, and thus may be different from the ones measured in \mathbb{R}^N .

Throughout the paper we use the following, somewhat nonstandard notation. For $U \subseteq \Omega$ and a function u ,

$$\int_U f(x) dx := \frac{1}{|\Omega|} \int_U f(x) dx \quad \text{and} \quad \|U\| := \frac{|U|}{|\Omega|}.$$

Let P be the maximal interval containing 1 on which $\mu > 0$:

$$P = \{z : z \leq 1, \mu|_{[z,1]} > 0\} \cup \{z : z \geq 1, \mu|_{[1,z]} > 0\}.$$

The configuration space is

$$\mathcal{M} := L^1(\Omega, \bar{P}).$$

While the configurations are functions defined on Ω , when convenient we also consider them as periodic functions of \mathbb{R}^N .

To a configuration, u , we associate the energy density

$$(10) \quad \bar{E}(u) := \int_{\Omega} \int_{\mathbb{R}^N} (u(x) - u(y))^2 K(x - y) dx dy + \int_{\Omega} W(u(x)) dx.$$

If the expression is not defined, we say that the energy density is infinite. Conditions on the interaction kernel, K , and the double-well potential, W , are described below.

We make the following assumptions on the interaction kernel K :

- (K1) K is nonnegative and $K \in L^1(\mathbb{R}^N) \cap C^2(\mathbb{R}^N)$.
- (K2) $K(x) = K(-x)$ for all $x \in \mathbb{R}^N$. (This condition ensures the symmetry of the interaction term in (10) with respect to x and y .)
- (K3) $K(0) > 0$.
- (K4) $K \in W^{2,1}(\mathbb{R}^N)$ and $\|K\|_{C^2(\mathbb{R}^N)} < \infty$.

The last condition is needed only for the existence theory [4]. Condition (K3) is not essential either, but significantly simplifies parts of the presentation. In particular it enables us to associate a length scale to a kernel in the following way: For $r > 0$ let

$$(11) \quad \kappa(x) := \frac{1}{|B(0,1)|} \chi_{B(0,1)}(x) \quad \text{and} \quad \kappa_r(x) := \frac{1}{r^N} \kappa\left(\frac{x}{r}\right),$$

where χ_U is the characteristic function of the set U . Given $r > 0$ let $h_K(r) := \sup\{c : K \geq c\kappa_r\}$. Note that by assumption (K3), $h_K(r) > 0$ for r small. It is not hard to prove that $h_K(r) \rightarrow 0$ as $r \rightarrow 0$ and also as $r \rightarrow \infty$. Consider the location of the maximum of $h_K(r)$. If there is more than one maximum, we pick the first one. More precisely let

$$(12) \quad r_K := \min\{r_{max} \mid h_K(r_{max}) = \max_{r>0} h_K(r)\}.$$

We let

$$H_K = h_K(r_K).$$

To state the conditions on the potential W we define

$$a := \int_{\Omega} u(x, 0) dx.$$

To simplify the presentation, from here on we restrict our attention only to

$$(13) \quad 0 < a \leq \frac{1}{2}.$$

We assume the following:

(W1) W is a nonnegative continuous function.

(W2) $W(0) = W(1) = 0$ and $W > 0$ on $\overline{P} \setminus \{0, 1\}$.

(W3) At least linear growth at $\pm\infty$: There exists a constant h_W such that $W(z) \geq 8h_W(z-1)$ for all $z \in P \cap (\frac{9}{8}, \infty)$ and $W(z) \geq h_W|z|$ for all $z \in P \cap (-\infty, -\frac{a}{4})$.

We can furthermore require

$$h_W \leq W(z) \quad \text{for all } z \in \left[\frac{a}{4}, \frac{7}{8}\right].$$

The condition on growth of W is needed only if P is infinite.

We now state the main results. Theorem 2 is an adaptation of Proposition 1 that applies to the configuration spaces we are investigating. In particular the paths are given as weak solutions of the continuity equation. Having the gradient-flow structure is reduced to requiring a dissipation inequality. The theorem relies on interpolation inequalities we establish in section 4. In Corollary 3 we apply the theorem to the main equations of our interest. The classes of solutions of (3) studied in [9] and [4] satisfy the conditions of the theorem, and thus the bounds on coarsening hold.

Below, by C^{weak} we mean continuous with respect to weak topology of the target space.

THEOREM 2. *Let $\Omega = [0, \Lambda]^N$. Assume that conditions (13), (K1)–(K3), and (W1)–(W3) hold and that mobility $\mu(z) \equiv z$ or that $\mu \in C(P, (0, c_\mu])$. Suppose that $u \in C^{weak}([0, \infty), L^1(\Omega, \overline{P}))$ is a weak solution of*

$$u_t + \nabla \cdot J = 0$$

for some flux $J \in L^1(\Omega, \mathbb{R}^N)$. Assume that the energy dissipation inequality holds: For almost all $0 \leq t_1 < t_2$

$$(14) \quad \int_{t_1}^{t_2} \int_{\Omega} \frac{1}{\mu(u)} |J|^2 dx dt \leq -(E(u(t_2)) - E(u(t_1))).$$

In the case $\mu(z) \not\equiv z$ we also need the regularity assumption that $u(t) \in L^2(\Omega)$ for all $t \geq 0$ and that the range of $u(\cdot, 0)$ is contained in P . Assume further that $\limsup_{t \rightarrow 0^+} \bar{E}(u(\cdot, t)) \ll 1$.

Then for all $\sigma \in (1, 2)$ there exists a constant $C = C(\sigma, a, H_K, r_K, c_\mu)$ such that for all $T \gg 1$

$$(15) \quad \int_0^T \bar{E}(u(t))^\sigma dt \geq C \int_0^T \left(t^{-\frac{1}{3}}\right)^\sigma dt.$$

The precise meaning of condition $\bar{E} \ll 1$ is that \bar{E} has to be small enough for the interpolation inequality to hold. Precise values can be found in the statement of Theorem 6.

By weak solution we mean that for all $\phi \in C_c^\infty(\Omega \times [0, \infty))$

$$\iint_{[0, \infty) \times \Omega} u \phi_t + J \cdot \nabla \phi dx dt + \int_{\Omega} u(x, 0) \phi(x, 0) dx = 0.$$

Recall that we consider Ω with the topology of a torus, which means that the test functions used in the definition of a weak solution are also defined on the torus; in other words they are periodic.

The form of the equation ensures (via testing against test functions with only time dependence) that

$$\int_{\Omega} u(x, t) dx = \text{const.}$$

Note that the dissipation inequality (14) is an equality if u is a classical solution of the gradient flow (2), when $J = -\mu(u) \nabla \frac{\delta E}{\delta u}$.

COROLLARY 3. *Let $\Omega = [0, \Lambda]^N$. Assume that conditions (K1)–(K4) and (W1)–(W3) hold. Assume only one of the following holds:*

- (i) $\mu(z) \equiv z$, $W \in C^2([0, \infty))$ with $W'' > -4 \int_{\mathbb{R}^N} K(y) dy$ on \mathbb{R}^+ , and u is the solution of (3) in the sense of [4] with $u(\cdot, 0) \in L^\infty(\Omega)$.
- (ii) $\mu(z) \leq c_\mu$ for all z , conditions (GL1) and (GL2) hold, and u is a solution of (3) in the sense of [9]. Furthermore assume that the range of $u(\cdot, 0)$ is contained in the interior of P .
- (iii) $0 < \mu(z) \leq c_\mu$ for all z , (GL2) holds on \mathbb{R} , and u is a solution of (3) in the sense of [9]. Furthermore assume that $u(\cdot, 0) \in L^\infty(\Omega)$.

Assume that $\limsup_{t \rightarrow 0^+} \bar{E}(u(\cdot, t)) \ll 1$ and (13) holds. Then for all $\sigma \in (1, 2)$ there exists a constant $C = C(\sigma, a, H_K, r_K, c_\mu)$ such that for all $T \gg 1$

$$(16) \quad \int_0^T \bar{E}(u(t))^\sigma dt \geq C \int_0^T \left(t^{-\frac{1}{3}}\right)^\sigma dt.$$

Proof. In the first case the conditions of W imply that associated $g(z) := \int_{\mathbb{R}^N} K(x) dx + \frac{1}{4} W'(z)$ is an increasing function. This in turn implies that the conditions under which Bertozzi and Slepčev [4] proved existence of solutions of (3) hold.

Properties of solutions ensuring that the assumptions of Theorem 2 hold were also established. This implies the claim of the corollary.

Case (ii) the existence theory needed for Theorem 2 to apply was established by Giacomini and Lebowitz [9].

The case (iii) is in principle simpler than case (ii) and includes the constant mobility case. The only technical issue is that the L^∞ bounds used in [9] follow from condition (GL1). In our case appropriate bounds can be established, for example, as in [4]. \square

3. Proof of Theorem 2. We seek to apply the framework of Proposition 1. However, the configuration space, \mathcal{M} , is not a true Riemannian manifold and the only remnant of the gradient flow structure is the energy dissipation inequality (14). Nevertheless arguments of the proof of the proposition can be adapted to include this setting. We define the “geodesic distance” on \mathcal{M} as follows: Given $u_0, u_1 \in \mathcal{M}$ let us first define a representation of admissible paths between u_0 and u_1 :

$$\mathcal{A}(u_0, u_1) := \left\{ (u, J) : u : [0, 1] \rightarrow \mathcal{M}, J \in L^1(\Omega \times [0, 1], \mathbb{R}^N) \text{ such that} \right. \\ \left. \begin{aligned} &u_t + \nabla \cdot J = 0 \quad \text{on } \Omega \times [0, 1] \text{ weakly,} \\ &u \in C^{weak}([0, 1], L^1(\Omega)) \text{ and } u(0) = u_0, u(1) = u_1, \text{ and} \\ &\int_0^1 \int_\Omega \frac{1}{\mu(u(x, t))} |J(x, t)|^2 dx dt < \infty \end{aligned} \right\}.$$

We define

$$(17) \quad d^2(u_0, u_1) := \inf_{(u, J) \in \mathcal{A}} \int_0^1 \int_\Omega \frac{1}{\mu(u(x, t))} |J(x, t)|^2 dx dt.$$

Here $\frac{0}{0} = 0$. We note that d may, in general, be infinite. It follows from the definition that d satisfies the triangle inequality. This can be shown by concatenating the appropriate test flows (with optimally rescaled times).

Let u be as in the statement of the theorem. We define

$$(18) \quad L(t) := d(u(t), a), \quad \bar{L}(t) := \frac{1}{\sqrt{|\Omega|}} d(u(t), a).$$

In the case $\mu(u) = u$ it follows from the characterization of d given below in (19) that $L(t)$ is finite for all t . In the other case, from the assumption on the range of u_0 it follows that $L(0)$ is finite. To see this it is enough to consider the test pair (\tilde{u}, \tilde{J}) with $\tilde{u}(s) = u_0 + s(a - u_0)$ for $s \in [0, 1]$ and $\tilde{J} = \nabla p$, where p solves $-\Delta p = a - u_0$. The fact that $L(t)$ is finite for all t then follows from the argument for continuity of L given in Lemma 4.

Let

$$\bar{E}(t) := \bar{E}(u(t)).$$

From (14) we have that \bar{E} is nonincreasing almost everywhere. We now modify \bar{E} on a set of measure 0 to ensure that it is nonincreasing:

$$\bar{E}_{new}(t) = \min\{\liminf_{s \rightarrow t^-} \bar{E}(s), \bar{E}(t)\}.$$

Inspecting the proof of Proposition 1 from [17] shows that in addition to the interpolation inequality, one only needs the inequality

$$\left(\frac{d\bar{L}}{d\bar{E}}\right)^2 \leq -\frac{dt}{d\bar{E}}$$

(since \bar{E} is nonincreasing, \bar{L} can be considered as a function of \bar{E}), which follows from the more familiar form of the dissipation inequality,

$$\left(\frac{d\bar{L}}{dt}\right)^2 \leq -\frac{d\bar{E}}{dt},$$

where both inequalities are to be understood as a comparison of measures (with given densities). The latter inequality in turn follows from the assumption (14). We prove these claims in Lemmas 4 and 5.

To be able to prove the interpolation inequalities we need a more workable form of d . In the case $\mu(u) = u$ the distance d is nothing else than the Wasserstein distance. This was shown by Benamou and Brenier [3] (see also section 8.1 in Villani’s book [19]):

(19)

$$d_W(u_0, u_1)^2 = \inf \left\{ \iint_{\Omega \times \Omega} |x - y|^2 d\pi(x, y) \mid \int_{\Omega} d\pi(\cdot, y) = u_0, \int_{\Omega} d\pi(x, \cdot) = u_1 \right\}.$$

The distance above, $|x - y|$, is taken on torus Ω .

In the case $\mu(u) \leq c_\mu$ first note that from the definition of the distance (17), it follows that the distance corresponding to $\mu(u)$ is greater than the distance corresponding to constant mobility c_μ . Thus

$$\bar{L}_\mu \geq \bar{L}_{c_\mu} = \frac{1}{c_\mu} \bar{L}_1.$$

Thus it is enough to establish the interpolation inequality for $\bar{L} = \bar{L}_1$. But for mobility equal to one, the distance is the H^{-1} norm. More precisely, for $u_0, u_1 \in \mathcal{M} \cap L^2(\Omega)$, $\int_{\Omega} u_1 - u_0 dx = 0$, and hence we can consider the following representation of the H^{-1} norm:

$$d(u_0, u_1)^2 = \|u_0 - u_1\|_{H^{-1}}^2 = \int_{\Omega} |\nabla p|^2 dx,$$

where $p \in H^2(\Omega)$ (p is periodic by topology of Ω) is a solution of

$$-\Delta p = u_1 - u_0.$$

Proof of this claim is straightforward; it relies on convexity in J of the functional on the right-hand side of (17) and the observation that $J = \nabla p$, along with $u(t) = u_0 + t(u_1 - u_0)$, minimizes the action functional. One can also show that for $u \in \mathcal{M} \cap L^2(\Omega)$

$$(20) \quad \bar{L}(u) = \max_{\xi \in H^1(\Omega), \xi \neq \text{const.}} \frac{\int_{\Omega} (u - a)\xi dx}{\sqrt{\int_{\Omega} |\nabla \xi|^2 dx}}$$

by using Cauchy–Schwarz inequality to show that $\xi = p$ (with $u_1 = u$ and $u_2 = a$) is the maximizing function. Given that $u(t) \in \mathcal{M} \cap L^2(\Omega)$ for all t we can use this characterization.

Thus to complete the proof of the theorem one only needs the interpolation inequalities established in section 4.

We now prove the two lemmas that were used in the arguments above.

LEMMA 4. *Assume that u satisfies the conditions of Theorem 2 and \bar{L} is defined in (18). Then \bar{L} is a continuous function and for almost all $t \geq 0$ and $h > 0$*

$$\left(\frac{\bar{L}(u(t+h)) - \bar{L}(u(t))}{h}\right)^2 \leq -\frac{\bar{E}(u(t+h)) - \bar{E}(u(t))}{h}.$$

Proof. By the assumptions of the theorem u is a distributional solution of

$$u_t + \nabla \cdot J = 0 \quad \text{on } \Omega \times [t, t+h].$$

Moreover it follows from assumption (14) for all $t \geq 0$ and all $h > 0$ that

$$\int_t^{t+h} \int_{\Omega} \frac{1}{\mu(u)} |J|^2 dx dt \leq \limsup_{s \rightarrow 0^+} E(u(s)) < \infty.$$

Note that it was also assumed that $u \in C^{weak}([t, t+h], L^1(\Omega))$. Thus (u, J) , after appropriate rescaling in time, belongs to $\mathcal{A}(u(t), u(t+h))$. By the triangle inequality,

$$\begin{aligned} (\bar{L}(u(t+h)) - \bar{L}(u(t)))^2 &\leq \inf_{(\tilde{u}, \tilde{J}) \in \mathcal{A}(u(t), u(t+h))} \int_0^1 \int_{\Omega} \frac{1}{\mu(\tilde{u}(x, s))} |\tilde{J}(x, s)|^2 dx ds \\ &\leq h \int_0^h \int_{\Omega} \frac{1}{\mu(u(x, t+s))} |J|^2 dx ds. \end{aligned}$$

Thus \bar{L} is a continuous function. Dividing the above by h^2 and using (14) gives that for almost all $t \geq 0$ and $h > 0$

$$\begin{aligned} \left(\frac{\bar{L}(u(t+h)) - \bar{L}(u(t))}{h}\right)^2 &\leq \frac{1}{h} \int_0^h \int_{\Omega} \frac{1}{\mu(u(x, t+s))} |J|^2 dx ds \\ &\leq -\frac{\bar{E}(u(t+h)) - \bar{E}(u(t))}{h}. \quad \square \end{aligned}$$

For a function e on \mathbb{R} let us define $e(t+) := \lim_{s \rightarrow t^+} e(s)$ and $e(t-) := \lim_{s \rightarrow t^-} e(s)$.

LEMMA 5. *Let e be a nonnegative, nonincreasing function on $[0, \infty)$. Let l be a continuous function on $[0, \infty)$, such that*

$$(21) \quad \left(\frac{l(t_2) - l(t_1)}{t_2 - t_1}\right)^2 \leq -\frac{e(t_2) - e(t_1)}{t_2 - t_1} \quad \text{for almost all } t_2 > t_1 \geq 0.$$

Then $l(t)$ is an absolutely continuous function on $[0, \infty)$ and for all $\tau_2 > \tau_1 \geq 0$

$$\int_{\tau_1}^{\tau_2} \left(\frac{dl}{dt}\right)^2 dt \leq e(\tau_1+) - e(\tau_2-).$$

Furthermore, consider $t(e) := \sup\{t : e(t) \geq e\}$, the “inverse” of the function e and $l(e) := l(t(e))$. Then

$$\left(\frac{l(e_2) - l(e_1)}{e_2 - e_1}\right)^2 \leq -\frac{t(e_2) - t(e_1)}{e_2 - e_1} \quad \text{for all } e(0) \geq e_1 > e_2 \geq 0.$$

Consequently l is an absolutely continuous function of e , and for all $e(0) \geq e_1 > e_2 \geq 0$

$$\int_{e_2}^{e_1} \left(\frac{dl}{de} \right)^2 de \leq t(e_2+) - t(e_1-).$$

Proof. If $e(0) = 0$, the proof is trivial. So assume $e(0) > 0$. Continuity of l implies that

$$\left(\frac{l(t_2) - l(t_1)}{t_2 - t_1} \right)^2 \leq -\frac{e(t_2-) - e(t_1+)}{t_2 - t_1} \quad \text{for all } t_2 > t_1 \geq 0.$$

Let $\varepsilon > 0$ and $\delta := \varepsilon^2/e(0)$. Let $[x_i, y_i]$ for $i = 1, \dots, m$ be a family of disjoint intervals on $[0, \infty)$ of total length less than δ :

$$\sum_{i=1}^m y_i - x_i < \delta.$$

Then

$$\begin{aligned} \sum_{i=1}^m |l(y_i) - l(x_i)| &\leq \sum_{i=1}^m \sqrt{e(x_i) - e(y_i)} \sqrt{y_i - x_i} \\ &\leq \sqrt{\sum_{i=1}^m e(x_i) - e(y_i)} \sqrt{\sum_{i=1}^m y_i - x_i} \\ &\leq \sqrt{e(0)} \sqrt{\delta} = \varepsilon. \end{aligned}$$

So l is absolutely continuous.

To prove the second claim note that for any $h > 0$

$$\begin{aligned} \int_{\tau_1}^{\tau_2} \left(\frac{l(t+h) - l(t)}{h} \right)^2 dt &\leq \int_{\tau_1}^{\tau_2} \frac{e(t+) - e(t+h-)}{h} dt \\ &\leq \frac{1}{h} \left(\int_{\tau_1}^{\tau_1+h} e(t+) dt - \int_{\tau_2}^{\tau_2+h} e(t-) dt \right). \end{aligned}$$

By taking the $\liminf_{h \rightarrow 0}$ and using Fatou's lemma we obtain

$$\int_{\tau_1}^{\tau_2} \left(\frac{dl}{dt} \right)^2 dt \leq e(\tau_1+) - e(\tau_2+).$$

Now use this claim on interval $(\tau_1, \tau_2 - \varepsilon)$ and take the limit as $\varepsilon \rightarrow 0$.

To prove the remaining claims, note that

$$(l(t(e_2)) - l(t(e_1)))^2 \leq (e(t(e_2)-) - e(t(e_1)+))(t(e_2) - t(e_1)).$$

Observing that $e(t(e_2)-) \geq e_2$ and $e(t(e_1)+) \leq e_1$ yields the desired inequality. The claims then follow from the arguments presented above. \square

4. Interpolation inequalities. In this section we prove the interpolation inequalities needed. As shown in section 3 we only need to consider the \bar{L} corresponding Wasserstein distance (19) and to H^{-1} norm (20).

The proof we present is general and extends to local energies, which we discuss in subsection 4.1. It also captures the improved constants established in [5]; see Remark 2. It is based on simple geometric heuristic. Consider function \tilde{u} with range $\{0, 1\}$ and $\kappa_r * \tilde{u}$, its average over ball of radius r . Then for r small

$$\int |\tilde{u} - \kappa_r * \tilde{u}| dx$$

contains information about the interfacial area, while for r large it carries information on the distance $d(\tilde{u}, a)$. This allows us to interpolate between the energy and the distance. We divided the proof into steps and present the motivation at their beginning.

THEOREM 6 (interpolation inequality). *Let $0 < a \leq \frac{1}{2}$. Assume K satisfies conditions (K1)–(K3) and W satisfies (W1)–(W3). There exists a constant $C = C(a, h_W, r_K, H_K) > 0$ such that for all $\Lambda > 0$ and all configurations $u \in \mathcal{M}$ for which $\bar{E}(u) < \frac{a}{64} (\frac{h_W}{1+h_W}) H_K$, and in the H^{-1} case also $\bar{E}(u) < \frac{1}{2^{N+2}} \frac{ah_W}{20} H_K$, the following holds:*

$$(22) \quad \bar{E}(u) \bar{L}(u) \geq C.$$

The constant $C = c(N)r_K (\frac{h_W}{1+h_W}) a^{3/2} H_K$.

Proof. Step 1: Reduction. Let κ_{r_K} be as defined by (11) and (12). We use the notation $\kappa_r := \kappa_{r_K}$. To make the distinction between energies, let \bar{E}_K and \bar{E}_{κ_r} be the energy densities corresponding to kernels K and κ_r , respectively. Note that $\bar{E}_K \geq H_K \bar{E}_{\kappa_r}$. So it is enough to show the above claim for κ_r , with $r_K = r$ and $H_K = 1$. Therefore from here on we consider only $K = \kappa_r$.

Step 2: u is separated into phases. We show that any low-energy-density configuration, u , has a significant portion of the mass on the set where values of u are close to 1. More precisely:

Claim. Let $A := \{x : u(x) \geq \frac{7}{8}\}$ and $\underline{A} := \chi_A$. For future reference let $\underline{A} := \{x : u(x) < \frac{a}{4}\}$ and let I be the interfacial region, $I := \Omega \setminus (A \cup \underline{A})$. Assume $\bar{E}(u) < \frac{3}{32} ah_W$. Then

$$(23) \quad \|A\| = \frac{|A|}{|\Omega|} = \int_{\Omega} \tilde{u}(x) dx > \frac{1}{2} \int_{\Omega} u(x) dx = \frac{a}{2} \quad \text{and} \quad \|A\| \leq 2a.$$

Proof. Due to assumption (W3)

$$\bar{E}(u) \geq h_W \left(\left\| \left\{ \frac{a}{4} < u \leq \frac{7}{8} \right\} \right\| + \left\| \left\{ \frac{9}{8} < u \right\} \right\| \right).$$

Consequently

$$\begin{aligned} a &= \int_{\Omega} u dx = \int_{\{u \leq \frac{a}{4}\}} u dx + \int_{\{\frac{a}{4} < u \leq \frac{7}{8}\}} u dx + \int_{\{\frac{7}{8} < u \leq \frac{9}{8}\}} u dx + \int_{\{\frac{9}{8} < u\}} u dx \\ &\leq \frac{a}{4} + \frac{7}{8} \left\| \left\{ \frac{a}{4} < u \leq \frac{7}{8} \right\} \right\| + \frac{9}{8} \int_{\Omega} \tilde{u} dx + \left\| \left\{ \frac{9}{8} < u \right\} \right\| + \frac{1}{h_W} \int_{\{\frac{9}{8} < u\}} W(u) dx \\ &\leq \frac{a}{4} + \frac{\bar{E}}{h_W} + \frac{9}{8} \int_{\Omega} \tilde{u} dx + \frac{\bar{E}}{h_W} \\ &< \frac{a}{4} + \frac{3a}{16} + \frac{9}{8} \int_{\Omega} \tilde{u} dx. \end{aligned}$$

Therefore $\int_{\Omega} \tilde{u} dx > \frac{a}{2}$.

To prove the second claim, note that

$$\begin{aligned} a &= \int_{\Omega} u(x) dx \geq \frac{7}{8} \|A\| + \int_{\{-\frac{a}{4} \leq a < 0\}} u(x) dx + \int_{\{u < -\frac{a}{4}\}} u(x) dx \\ &\geq \frac{7}{8} \|A\| - \frac{a}{4} - \frac{1}{h_W} \int_{\Omega} W(u(x)) dx \\ &\geq \frac{7}{8} \|A\| - \frac{a}{4} - \frac{\bar{E}}{h_W} \\ &\geq \frac{7}{8} \|A\| - \frac{a}{2}. \quad \square \end{aligned}$$

Step 3: Energy bounds a measure of interfacial area.

Claim.

$$\int_{\Omega} |\tilde{u} - \kappa_r * \tilde{u}| dx \leq \left(\frac{16}{9} + \frac{2}{h_W} \right) \bar{E}.$$

Heuristically, when r is small the expression on the left-hand side measures r times the area of the boundary of $\{\tilde{u} = 1\}$. More precisely the area is measured in such a way that features of size less than r are smoothed out, and thus neglected, to some extent.

Proof. We use the following notation for the sum of sets $X + Y := \{x + y \mid x \in X, y \in Y\}$. Using the fact that \tilde{u} takes only values 0 and 1, we obtain

$$\begin{aligned} &\int_{\Omega} \left| \tilde{u}(x) - \int_{\mathbb{R}^N} \kappa_r(x - y) \tilde{u}(y) dy \right| dx \\ &= \int_{\Omega} \int_{\mathbb{R}^N} |\tilde{u}(x) - \tilde{u}(y)| \kappa_r(x - y) dy dx \\ &\leq \frac{1}{|\Omega|} \left[\int_A \int_{\underline{A} + \Lambda \mathbb{Z}^N} + \int_{\underline{A}} \int_{A + \Lambda \mathbb{Z}^N} |\tilde{u}(x) - \tilde{u}(y)|^2 \kappa_r(x - y) dy dx \right. \\ &\quad \left. + \int_I \int_{\mathbb{R}^N} + \int_{\Omega} \int_{I + \Lambda \mathbb{Z}^N} \kappa_r(x - y) dy dx \right] \\ &\leq \left(\frac{3}{4} \right)^{-2} \int_{\Omega} \int_{\mathbb{R}^N} |u(x) - u(y)|^2 \kappa_r(x - y) dy dx + 2 \left\| \left\{ \frac{a}{4} \leq u \leq \frac{7}{8} \right\} \right\| \\ &\leq \left(\frac{16}{9} + \frac{2}{h_W} \right) \bar{E}. \quad \square \end{aligned}$$

Step 4. Claim. $\phi : (0, \infty) \rightarrow [0, \infty)$ defined by

$$\phi(s) := \int_{\Omega} |\tilde{u} - \tilde{u} * \kappa_s| dx$$

is subadditive.

One should note that some other possible measures of surface area (for example, the volume of appropriate tubular neighborhood) do not have this property in general and may have a superlinear growth (for appropriate range of r).

Proof. Let $s = p + q$ for some $p, q > 0$. As in Step 3 we have

$$\phi(s) = \int_{\Omega} \int_{\mathbb{R}^N} |\tilde{u}(x) - \tilde{u}(y)| \kappa_s(y) dy dx.$$

Now let $z = x - \frac{p}{s}y$. Using periodicity and the scaling properties of kernel κ_s one finds

$$\begin{aligned} \phi(s) &\leq \int_{\Omega} \int_{\mathbb{R}^N} (|\tilde{u}(x) - \tilde{u}(z)| + |\tilde{u}(z) - \tilde{u}(y)|) \kappa_s(y) dy dx \\ &= \int_{\Omega} \int_{\mathbb{R}^N} \left| \tilde{u}(x) - \tilde{u}\left(x - \frac{p}{s}y\right) \right| \kappa_s(y) dy dx + \int_{\Omega} \int_{\mathbb{R}^N} \left| \tilde{u}(z) - \tilde{u}\left(z - \frac{q}{s}y\right) \right| \kappa_s(y) dy dz; \end{aligned}$$

substitute $\tilde{y} = \frac{p}{s}y$ in the first integral and $\tilde{y} = \frac{q}{s}y$ and $x = z$ in the second to obtain

$$\begin{aligned} &= \int_{\Omega} \int_{\mathbb{R}^N} |\tilde{u}(x) - \tilde{u}(\tilde{y})| \kappa_p(\tilde{y}) d\tilde{y} dx + \int_{\Omega} \int_{\mathbb{R}^N} |\tilde{u}(x) - \tilde{u}(\tilde{y})| \kappa_q(\tilde{y}) d\tilde{y} dx \\ &= \phi(p) + \phi(q). \quad \square \end{aligned}$$

Step 5. $\kappa_l * \tilde{u} \sim \tilde{u}$ for some l of size $\frac{1}{\bar{E}}$. More precisely:

Claim. Let $\mu > 2$ be a constant, which we specify later. If $\bar{E} < \frac{a}{\mu} \left(\frac{h_W}{1+h_W}\right)$, then for

$$(24) \quad l = \left\lfloor \frac{a}{\mu} \left(\frac{h_W}{1+h_W}\right) \frac{1}{\bar{E}} \right\rfloor r =: i r$$

the following holds:

$$(25) \quad \phi(l) = \int_{\Omega} |\tilde{u} - \kappa_l * \tilde{u}| dx < \frac{2}{\mu} a.$$

Proof. The assumption on \bar{E} implies that $l > \frac{1}{2} \frac{a}{\mu} \left(\frac{h_W}{1+h_W}\right) \frac{1}{\bar{E}} r > 0$. Subadditivity of ϕ established in Step 4 and the bound of Step 3 imply

$$\int_{\Omega} |\tilde{u} - \kappa_l * \tilde{u}| dx \leq i \int_{\Omega} |\tilde{u} - \kappa_r * \tilde{u}| dx \leq \left\lfloor \frac{a}{\mu} \left(\frac{h_W}{1+h_W}\right) \frac{1}{\bar{E}} \right\rfloor 2 \left(1 + \frac{1}{h_w}\right) \bar{E} \leq \frac{2}{\mu} a. \quad \square$$

Step 6. If for some $l > 1$,

$$\kappa_l * \tilde{u} \sim \tilde{u}, \quad \text{then} \quad \bar{L} \gtrsim l.$$

More precisely:

Claim. Set $\mu = 64$. There exists a constant c , depending only on dimension N and on a , such that for any $l > 1$,

$$(26) \quad \text{if} \quad \int_{\Omega} |\tilde{u} - \kappa_l * \tilde{u}| dx < \frac{2}{\mu} a, \quad \text{then} \quad \bar{L} > cl.$$

We split the proof of this claim into four parts. First we establish two auxiliary claims. Then we prove claim (26) for the Wasserstein metric case and for the H^{-1} metric case separately.

Step 6a. Let $A_l := \{x \in \Omega : \tilde{u} * \kappa_l > \frac{7}{8}\}$. If

$$\int_{\Omega} |\tilde{u} - \kappa_l * \tilde{u}| dx < \frac{2}{\mu} a,$$

then

$$(27) \quad \|A_l\| > \frac{\mu - 32}{2\mu} a.$$

By the assumption

$$\frac{2}{\mu} a > \int_{\Omega} |\tilde{u} - \kappa_l * \tilde{u}| dx \geq \int_{A \setminus A_l} \frac{1}{8} dx = \frac{1}{8} \|A \setminus A_l\|.$$

From (23) we have $\|A\| > \frac{a}{2}$. Combining the two inequalities gives

$$\|A_l\| \geq \|A\| - \|A \setminus A_l\| > \left(\frac{1}{2} - \frac{16}{\mu}\right) a.$$

Remark. From this point on the proof does not require the closeness of \tilde{u} and $\kappa_l * \tilde{u}$ explicitly, but rather it uses only the fact that A_l is large, as described by (27). That is, we require only that after \tilde{u} is averaged over radius l it still have well-developed interfaces.

Step 6b. A significant subset of A_l can be well approximated by balls of radius l . More precisely:

Claim. Set $\mu = 64$. There exists a finite subset, J , of A_l such that for $A_{ball} = \cup_{x \in J} B(x, l)$,

$$(28) \quad \frac{8}{7} \|A_{ball} \cap A\| \geq \|A_{ball}\| > \frac{1}{2^{N+2}} a.$$

This claim has its roots in [5]; see also [17]. Let J be a maximal family of points in A_l such that balls in $\{B(x, l)\}_{x \in J}$ are disjoint. Then $A_l \subset \cup_{x \in J} B(x, 2l)$, by definition. Therefore, using (27)

$$\|A_{ball}\| \geq \frac{1}{2^N} \|A_l\| > \frac{1}{2^{N+2}} a.$$

Since $J \subset A_l$, for all $x \in J$ we have that

$$\frac{7}{8} \leq \kappa_l * \tilde{u}(x) \leq \frac{|B(x, l) \cap A|}{|B(x, l)|}.$$

Summing over $x \in J$ gives the first inequality.

Step 6c (Wasserstein). Let $\gamma = \left(\frac{9}{8}\right)^{1/N} - 1$. For set U and $\lambda \geq 0$, let

$$U^\lambda := \{x \in \Omega : \text{dist}(x, U) \leq \lambda\}.$$

Let $\lambda = \gamma l$. Using Lemma 8 and the fact that $u \geq 0$,

$$\begin{aligned} \bar{L}^2 &= \frac{d_{Wass}(u, a)^2}{|\Omega|} \geq \lambda^2 \left(\int_{A_{ball}} u(x) dx - a \|A_{ball}^\lambda\| \right) \\ &\geq \lambda^2 \left(\frac{7}{8} \|A \cap A_{ball}\| - a \left(1 + \frac{\lambda}{l}\right)^N \|A_{ball}\| \right) \\ &\geq \gamma^2 l^2 \left(\frac{49}{64} - \frac{1}{2} (1 + \gamma)^N \right) \|A_{ball}\| \\ &\geq \gamma^2 \frac{1}{5} \frac{1}{2^{N+2}} a l^2. \end{aligned}$$

Combining the conclusions of Steps 5 and 6 now proves the interpolation inequality (22). In particular

$$\bar{L} \geq c(N)r \left(\frac{h_W}{1 + h_W} \right) a^{3/2} \frac{1}{\bar{E}}.$$

Step 6d (H^{-1}). Assume $\bar{E} < \frac{1}{20} \frac{1}{2^{N+2}} a h_w$. To obtain a lower bound on \bar{L} , given by (20), we first build a local test function. For $\gamma > 1$ to be determined, let $\eta : [0, \infty) \rightarrow [0, 1]$ be defined by

$$\eta(z) := \begin{cases} 1 & \text{if } z \in [0, l], \\ l - \frac{z-l}{\gamma l - l} & \text{if } l < z < \gamma l, \\ 0 & \text{if } \gamma l \leq z. \end{cases}$$

Let $\bar{\xi}(x) := \eta(|x|)$. This is the local test function. Assume, for the moment, that $0 \in A_l$. Let $\hat{u}(x) := \max\{u(x), -\frac{a}{4}\}$. Then

$$(29) \quad \int_{B(0, \gamma l)} |\nabla \bar{\xi}|^2 dx = (\gamma^N - 1) |B(0, l)| \frac{1}{(\gamma - 1)^2 l^2}.$$

Also

$$\begin{aligned} &\int_{B(0, \gamma l)} (\hat{u} - a) \bar{\xi} dx \\ &\geq \frac{7}{8} |A \cap B(0, l)| - \frac{a}{4} (|B(0, l) \setminus A| + (\gamma^N - 1) |B(0, l)|) - a \gamma^N |B(0, l)| \\ &\geq \left(\left(\frac{7}{8}\right)^2 - \frac{a}{32} - \frac{a(\gamma^N - 1)}{4} - a \gamma^N \right) |B(0, l)|. \end{aligned}$$

Let us now set $\gamma = \left(\frac{8}{7}\right)^{1/N}$. Then

$$(30) \quad \int_{B(0, \gamma l)} (\hat{u} - a) \bar{\xi} dx \geq \frac{3}{20} |B(0, l)|.$$

Now let us construct the (global) test function on Ω . Let

$$\xi(x) = \sup_{y \in J} \bar{\xi}(x - y).$$

Using (29), (23), and (28) we obtain

$$(31) \quad \int_{\Omega} |\nabla \xi|^2 dx \leq \frac{\gamma^N - 1}{(\gamma - 1)^2 l^2} \|A_{ball}\| \leq \frac{a}{(\gamma - 1)^2 l^2}.$$

Using the fact that balls of radius l centered at points in J are disjoint we obtain

$$\begin{aligned} \int_{\Omega} (u - a)\xi dx &\geq \int_{\Omega} (\hat{u} - a)\xi dx + \int_{\{u < -\frac{a}{4}\}} u dx \\ &\geq \frac{3}{20} \|A_{ball}\| - \frac{1}{h_W} \int_{\Omega} W(u) dx \\ &\geq \frac{3}{20} \frac{1}{2^{N+2}} a - \frac{\bar{E}}{h_w} \\ &\geq \frac{1}{10} \frac{1}{2^{N+2}} a. \end{aligned}$$

Here we used the assumption $\bar{E} < \frac{1}{20} \frac{1}{2^{N+2}} a h_w$. Therefore

$$\bar{L} \geq \frac{\int_{\Omega} (u - a)\xi dx}{\sqrt{\int_{\Omega} |\nabla \xi|^2 dx}} \geq \tilde{c}(N) \sqrt{al}.$$

The definition of l now implies

$$\bar{L} \geq c(N)r \left(\frac{h_W}{1 + h_W} \right) a^{3/2} \frac{1}{\bar{E}},$$

which proves the interpolation inequality. \square

Remark 2. In the case $N = 2$, one can obtain a sharper result with respect to scaling in a (as $a \rightarrow 0^+$) by considering more carefully constructed test functions. This was done for the Mullins–Sekerka evolution by Conti, Niethammer, and Otto in [5]. In particular if γ is taken of size $a^{-1/2}$, and on $[l, cl]$, we replace the linear η by the optimal one, $\eta(z) = (\ln \gamma l - \ln z) / \ln \gamma$. By using such a test function one obtains that

$$\bar{L} \geq c(N)r \left(\frac{h_W}{1 + h_W} \right) a^{3/2} |\ln a|^{1/2} \frac{1}{\bar{E}}.$$

4.1. Interpolation inequalities for the local energy. Let us now consider the case of the local energy density:

$$(32) \quad \bar{E}(u) := \int_{\Omega} \frac{1}{2} |\nabla u(x)|^2 + W(u(x)) dx.$$

The method developed in the proof of Theorem 6 applies to the local energy with minor modifications.

COROLLARY 7 (interpolation inequality). *Let $0 < a \leq \frac{1}{2}$. Assume W satisfies (W1)–(W3). There exists a constant $C = C(a, h_W) > 0$ such that for all $\Lambda > 0$ and all configurations $u \in \mathcal{M}$ for which $\bar{E}(u) < \frac{a}{64} \left(\frac{h_W}{1 + h_W} \right)$, and in the H^{-1} case also $\bar{E}(u) < \frac{1}{2^{N+2}} \frac{a h_W}{20}$, the following holds:*

$$(33) \quad \bar{E}(u) \bar{L}(u) \geq C.$$

The constant $C = c(N) \left(\frac{h_W}{1 + h_W} \right) a^{3/2}$.

Proof. Note that Step 1 is not needed, while the estimate of Step 2 used only the W term, which is the same for both the local and the nonlocal energy. The main fact we need to check is the statement of Step 3, which we prove below for $r^2 = \frac{1}{2}$. Steps 4, 5, and 6 do not require any modifications.

To prove that

$$\int_{\Omega} |\tilde{u} - \kappa_r * \tilde{u}| dx \leq \left(\frac{16}{9} + \frac{2}{h_W} \right) \bar{E}$$

with $r^2 = \frac{1}{2}$ we begin as in Step 3 of the proof of Theorem 6:

$$\begin{aligned} & \int_{\Omega} \left| \tilde{u}(x) - \int_{\mathbb{R}^N} \kappa_r(x-y) \tilde{u}(y) dy \right| dx \\ & \leq \left(\frac{3}{4} \right)^{-2} \int_{\Omega} \int_{\mathbb{R}^N} |u(x) - u(y)|^2 \kappa_r(x-y) dy dx + 2 \left\| \left\{ \frac{a}{4} \leq u \leq \frac{7}{8} \right\} \right\| \\ & \leq \frac{16}{9} \int_{\Omega} \frac{1}{|B(x,r)|} \int_{B(x,r)} |u(x) - u(y)|^2 dy dx + \frac{2}{h_W} \bar{E}. \end{aligned}$$

It remains to further estimate the first term:

$$\begin{aligned} & \int_{\Omega} \frac{1}{|B(x,r)|} \int_{B(x,r)} |u(x) - u(y)|^2 dy dx \\ & \leq \int_{\Omega} \frac{1}{|B(0,r)|} \int_{B(x,r)} |x-y|^2 \left| \int_0^1 \nabla u(x - s(x-y)) ds \right|^2 dy dx \\ & \leq \frac{1}{2} \int_{\Omega} \frac{1}{|B(0,r)|} \int_{B(0,r)} \int_0^1 |\nabla u(x - sz)|^2 ds dz dx \\ & \leq \frac{1}{2} \frac{1}{|B(0,r)|} \int_{B(0,r)} \int_0^1 2\bar{E} ds dz = \bar{E}. \quad \square \end{aligned}$$

5. Appendix.

5.1. A property of Wasserstein distance. The following lemma is analogous to Lemma 5 in [17]. We state it in large generality, the reason being that we want to consider Ω with metric from the torus $R^N / (\Lambda\mathbb{Z})^N$. In applications σ is the Lebesgue measure.

LEMMA 8. *Let (Ω, d, σ) be a metric space endowed with finite measure σ . Let $u \in L^1(\Omega)$ be a nonnegative function with average $a := \int_{\Omega} u(x) d\sigma(x)$. Let $A \subset \Omega$ measurable, and let $A^l := \{x \in \Omega : d(x, A) \leq l\}$. Then*

$$d_{W_{ass}}^2(u, a) \geq l^2 \left(\int_A u(x) d\sigma(x) - a\sigma(A^l) \right).$$

Proof. We use the definition of Wasserstein distance. Let π be an admissible transportation plan, that is, a measure on $\Omega \times \Omega$ with marginals $u(x)d\sigma(x)$ and $a\sigma(y)$.

Then

$$\begin{aligned}
 \int_{\Omega \times \Omega} |x - y|^2 d\pi(x, y) &= \int_{A \times (\Omega \setminus A^l)} |x - y|^2 d\pi(x, y) \\
 &\geq l^2 \pi(A \times (\Omega \setminus A^l)) \\
 &\geq l^2 (\pi(A \times \Omega) - \pi(\Omega \times A^l)) \\
 &= l^2 \left(\int_A u(x) d\sigma(x) - \int_{A^l} a d\sigma(y) \right) \\
 &= l^2 \left(\int_A u(x) d\sigma(x) - a\sigma(A^l) \right). \quad \square
 \end{aligned}$$

Acknowledgments. The author would like to thank Andrea Bertozzi and Chad Topaz for helpful discussions, as well as the referees for carefully reading the manuscript and providing valuable suggestions. The author would also like to thank the Center for Nonlinear Analysis for its support during the preparation of this paper.

REFERENCES

- [1] G. ALBERTI AND G. BELLETTINI, *A non-local anisotropic model for phase transitions: Asymptotic behaviour of rescaled energies*, European J. Appl. Math., 9 (1998), pp. 261–284.
- [2] G. ALBERTI, G. BELLETTINI, M. CASSANDRO, AND E. PRESUTTI, *Surface tension in Ising systems with Kac potentials*, J. Statist. Phys., 82 (1996), pp. 743–796.
- [3] J.-D. BENAMOU AND Y. BRENIER, *A numerical method for the optimal time-continuous mass transport problem and related problems*, in Monge Ampère Equation: Applications to Geometry and Optimization (Deerfield Beach, FL, 1997), Contemp. Math. 226, AMS, Providence, RI, 1999, pp. 1–11.
- [4] A. BERTOZZI AND D. SLEPČEV, *On nonlocal equations modeling biological aggregation*, in preparation (2008).
- [5] S. CONTI, B. NIETHAMMER, AND F. OTTO, *Coarsening rates in off-critical mixtures*, SIAM J. Math. Anal., 37 (2006), pp. 1732–1741.
- [6] S. DAI AND R. L. PEGO, *Universal bounds on coarsening rates for mean-field models of phase transitions*, SIAM J. Math. Anal., 37 (2005), pp. 347–371.
- [7] S. DAI AND R. L. PEGO, *An upper bound on the coarsening rate for mushy zones in a phase-field model*, Interfaces Free Bound., 7 (2005), pp. 187–197.
- [8] G. GIACOMIN AND J. L. LEBOWITZ, *Phase segregation dynamics in particle systems with long range interactions. I. Macroscopic limits*, J. Statist. Phys., 87 (1997), pp. 37–61.
- [9] G. GIACOMIN AND J. L. LEBOWITZ, *Phase segregation dynamics in particle systems with long range interactions. II. Interface motion*, SIAM J. Appl. Math., 58 (1998), pp. 1707–1729.
- [10] R. V. KOHN AND F. OTTO, *Upper bounds on coarsening rates*, Comm. Math. Phys., 229 (2002), pp. 375–395.
- [11] R. V. KOHN AND X. YAN, *Upper bound on the coarsening rate for an epitaxial growth model*, Comm. Pure Appl. Math., 56 (2003), pp. 1549–1564.
- [12] R. V. KOHN AND X. YAN, *Coarsening rates for models of multicomponent phase separation*, Interfaces Free Bound., 6 (2004), pp. 135–149.
- [13] B. LI AND J.-G. LIU, *Epitaxial growth without slope selection: Energetics, coarsening, and dynamic scaling*, J. Nonlinear Sci., 14 (2004), pp. 429–451.
- [14] L. MODICA, *The gradient theory of phase transitions and the minimal interface criterion*, Arch. Rational Mech. Anal., 98 (1987), pp. 123–142.
- [15] L. MODICA AND S. MORTOLA, *Un esempio di Γ^- -convergenza*, Boll. Un. Mat. Ital. B (5), 14 (1977), pp. 285–299.
- [16] F. OTTO, *The geometry of dissipative evolution equations: The porous medium equation*, Comm. Partial Differential Equations, 26 (2001), pp. 101–174.
- [17] F. OTTO, T. RUMP, AND D. SLEPČEV, *Coarsening rates for a droplet model: Rigorous upper bounds*, SIAM J. Math. Anal., 38 (2006), pp. 503–529.
- [18] C. M. TOPAZ, A. L. BERTOZZI, AND M. A. LEWIS, *A nonlocal continuum model for biological aggregation*, Bull. Math. Biol., 68 (2006), pp. 1601–1623.
- [19] C. VILLANI, *Topics in Optimal Transportation*, Grad. Stud. Math. 58, AMS, Providence, RI, 2003.

UNIQUENESS OF POSITIVE BOUND STATES TO SCHRÖDINGER SYSTEMS WITH CRITICAL EXPONENTS*

CONGMING LI[†] AND LI MA[‡]

Abstract. We prove the uniqueness of the positive solutions of the following elliptic system: (1) $-\Delta(u(x)) = u(x)^\alpha v(x)^\beta$, (2) $-\Delta(v(x)) = u(x)^\beta v(x)^\alpha$. Here $x \in R^n$, $n \geq 3$, and $1 \leq \alpha < \beta \leq \frac{n+2}{n-2}$ with $\alpha + \beta = \frac{n+2}{n-2}$. In the special case when $n = 3$ and $\alpha = 2, \beta = 3$, the system is closely related to the ones from the stationary Schrödinger system with critical exponents for the Bose–Einstein condensate. As the first step, we prove the radial symmetry of the positive solutions to the elliptic system above with critical exponents. We then prove that $u = v$, which is a key point for our uniqueness result.

Key words. moving plane, positive solutions, radial symmetric, uniqueness

AMS subject classifications. 35J45, 35J60, 45G05, 45G15

DOI. 10.1137/080712301

1. Introduction. In this paper, we consider the uniqueness of positive solutions to the following stationary Schrödinger system:

$$(1) \quad \begin{cases} -\Delta(u(x)) = u(x)^\alpha v(x)^\beta, \\ -\Delta(v(x)) = u(x)^\beta v(x)^\alpha. \end{cases}$$

Here, in the special case when $n = 3$ and $\alpha = 2, \beta = 3$, and $u = v$, system (1) is reduced to the quintic Schrödinger equation considered by Bourgain [1]. We prove in this paper that there is a uniqueness result for the system above. In general, our system is related to the ones from the stationary Schrödinger system with critical exponents for the Bose–Einstein condensate (see [15], [19], [20], and [23]). In the earlier works [19], [20], and [23], more attention was paid to the elliptic system (1) with subcritical exponents. Very interestingly, Chen and Li proved that the best constant in the weighted Hardy–Littlewood–Sobolev inequality can be achieved by explicit radial symmetric functions (see [5] and [18]). As a consequence of their work, the uniqueness of positive solutions to the corresponding elliptic system (it is (1) in the case when $\alpha = 0$ and $\beta = \frac{n+2}{n-2}$) has been settled. However, when $0 < \alpha, \beta \leq \frac{n+2}{n-2}$, the uniqueness of smooth positive solutions to the stationary Schrödinger system (1) is an open question. Generally speaking, there are very few results even for the uniqueness of positive solutions to ordinary differential systems. The aim of this paper is to prove the radial symmetry and uniqueness of positive solutions to (1) with critical exponents and $1 \leq \alpha < \beta \leq \frac{n+2}{n-2}$.

*Received by the editors January 2, 2008; accepted for publication (in revised form) June 17, 2008; published electronically October 13, 2008.

<http://www.siam.org/journals/sima/40-3/71230.html>

[†]Department of Applied Mathematics, University of Colorado at Boulder, Campus Box 526, Boulder, CO 80309 (cli@colorado.edu). The work of this author was partially supported by NSF grant DMS-0401174.

[‡]Department of Mathematical Sciences, Tsinghua University, Beijing 100084, China (lma@math.tsinghua.edu.cn). The work of this author was partially supported by the National Natural Science Foundation of China through grant 10631020, a key grant KZ200710025012 of the Education Committee of Beijing, and SRFDP grant 20060003002.

As one can expect, just like in the work of Weinstein [28] in the scalar case with subcritical exponent, there is a close relationship between the stationary Schrödinger system with critical exponent and the Hardy–Littlewood–Sobolev inequality. As we show below, this is true.

Since we shall use Hardy–Littlewood–Sobolev inequality to prove radial symmetry of our solutions, let’s first explore the recent progress of Lieb’s conjecture. We begin by recalling the well-known Hardy–Littlewood–Sobolev inequality. Let $0 < \lambda < n$, $1 < s, r < \infty$, and $\|f\|_p$ be the $L^p(\mathbb{R}^n)$ norm of the function f . We shall write by $\|f\|_{p,\Omega}$ the L^p norm of the function f on the domain Ω . Then the classical Hardy–Littlewood–Sobolev inequality states that

$$(2) \quad \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} \frac{f(x)g(y)}{|x - y|^\lambda} dx dy \leq C_{s,\lambda,n} \|f\|_r \|g\|_s$$

for any $f \in L^r(\mathbb{R}^n)$, $g \in L^s(\mathbb{R}^n)$, and for $\frac{1}{r} + \frac{1}{s} + \frac{\lambda}{n} = 2$. Hardy and Littlewood also introduced the double weighted inequality, which was later generalized by Stein and Weiss in [26] in the following form:

$$(3) \quad \left| \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} \frac{f(x)g(y)}{|x|^{\alpha_0} |x - y|^\lambda |y|^{\beta_0}} dx dy \right| \leq C_{\alpha_0,\beta_0,s,\lambda,n} \|f\|_r \|g\|_s,$$

where $\alpha_0 + \beta_0 \geq 0$,

$$(4) \quad 1 - \frac{1}{r} - \frac{\lambda}{n} < \frac{\alpha_0}{n} < 1 - \frac{1}{r} \quad \text{and} \quad \frac{1}{r} + \frac{1}{s} + \frac{\lambda + \alpha_0 + \beta_0}{n} = 2.$$

The best constant in the weighted inequality (3) can be obtained by maximizing the functional

$$(5) \quad J(f, g) = \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} \frac{f(x)g(y)}{|x|^{\alpha_0} |x - y|^\lambda |y|^{\beta_0}} dx dy$$

under the constraints $\|f\|_r = \|g\|_s = 1$. Then the corresponding Euler–Lagrange equations are the system of integral equations

$$(6) \quad \begin{cases} \lambda_1 r f(x)^{r-1} = \frac{1}{|x|^{\alpha_0}} \int_{\mathbb{R}^n} \frac{g(y)}{|y|^{\beta_0} |x - y|^\lambda} dy, \\ \lambda_2 s g(x)^{s-1} = \frac{1}{|x|^{\beta_0}} \int_{\mathbb{R}^n} \frac{f(y)}{|y|^{\alpha_0} |x - y|^\lambda} dy, \end{cases}$$

where $f, g \geq 0$, $x \in \mathbb{R}^n$, and $\lambda_1 r = \lambda_2 s = J(f, g)$.

Let $u = c_1 f^{r-1}$, $v = c_2 g^{s-1}$, $p = \frac{1}{r-1}$, $q = \frac{1}{s-1}$, $pq \neq 1$, and for a proper choice of constants c_1 and c_2 , system (6) becomes

$$(7) \quad \begin{cases} u(x) = \frac{1}{|x|^{\alpha_0}} \int_{\mathbb{R}^n} \frac{v(y)^q}{|y|^{\beta_0} |x - y|^\lambda} dy, \\ v(x) = \frac{1}{|x|^{\beta_0}} \int_{\mathbb{R}^n} \frac{u(y)^p}{|y|^{\alpha_0} |x - y|^\lambda} dy, \end{cases}$$

where $u, v \geq 0$, $0 < p, q < \infty$, $0 < \lambda < n$, $\frac{\alpha_0}{n} < \frac{1}{p+1} < \frac{\lambda + \alpha_0}{n}$, and $\frac{1}{p+1} + \frac{1}{q+1} = \frac{\lambda + \alpha_0 + \beta_0}{n}$.

Note that in the special case where $\alpha_0 = 0$ and $\beta_0 = 0$, system (7) reduces to the following system:

$$(8) \quad \begin{cases} u(x) = \int_{R^n} \frac{v^q(y)}{|x-y|^\lambda} dy, \\ v(x) = \int_{R^n} \frac{u^p(y)}{|x-y|^\lambda} dy \end{cases}$$

with

$$(9) \quad \frac{1}{q+1} + \frac{1}{p+1} = \frac{\lambda}{n}.$$

It is well known that this integral system is closely related to the system of partial differential equations

$$(10) \quad \begin{cases} (-\Delta)^{\gamma/2} u = v^q, & u > 0, \text{ in } R^n, \\ (-\Delta)^{\gamma/2} v = u^p, & v > 0, \text{ in } R^n, \end{cases}$$

where $\gamma = n - \lambda$.

When $p = q = \frac{n+\gamma}{n-\gamma}$ and $u(x) = v(x)$, system (8) becomes the single equation

$$(11) \quad u(x) = \int_{R^n} \frac{u(y)^{\frac{n+\gamma}{n-\gamma}}}{|x-y|^{n-\gamma}} dy, \quad u > 0, \text{ in } R^n.$$

The corresponding PDE is the well-known family of semilinear equations

$$(12) \quad (-\Delta)^{\gamma/2} u = u^{(n+\gamma)/(n-\gamma)}, \quad u > 0, \text{ in } R^n.$$

In particular, when $n \geq 3$, and $\gamma = 2$, (12) becomes

$$(13) \quad -\Delta u = u^{(n+2)/(n-2)}, \quad u > 0, \text{ in } R^n.$$

The classification of the solutions of (13) has provided an important ingredient in the study of the well-known Yamabe problem and the prescribing scalar curvature problem. Equation (13) was studied by Gidas, Ni, and Nirenberg [11], Caffarelli, Gidas, and Spruck [2], Chen and Li [3], and Li [16]. They classified all the positive solutions. Recently, Wei and Xu [27] generalized this result to the solutions of the more general equation (12) with γ being any even number between 0 and n . Some results of Chen, Li, and Ou have been improved by Hang in [12]. One may see related results in [14] and [17].

Although the systems for other real values of α, β between 0 and n are of interest to some, we shall concentrate in this paper only on the system (1) with critical exponents when $1 \leq \alpha, \beta \leq \frac{n+2}{n-2}$ and $\alpha + \beta = \frac{n+2}{n-2}$.

Our main results are the following two theorems.

THEOREM 1. *Assume that $n \geq 3$, $1 \leq \alpha, \beta \leq \frac{n+2}{n-2}$, and $\alpha + \beta = \frac{n+2}{n-2}$. Then any $L^{\frac{2n}{n-2}}(\mathbf{R}^n) \times L^{\frac{2n}{n-2}}(\mathbf{R}^n)$ positive solution pair (u, v) to system (1) is radial symmetric.*

THEOREM 2. *Assume that $n \geq 3$, $1 \leq \alpha \leq \beta \leq \frac{n+2}{n-2}$, and $\alpha + \beta = \frac{n+2}{n-2}$. Then any $L^{\frac{2n}{n-2}}(\mathbf{R}^n) \times L^{\frac{2n}{n-2}}(\mathbf{R}^n)$ radial symmetric solution pair (u, v) to system (1) is unique and $u = v$.*

In the proof of Theorem 1, we shall use the equivalent integral form for the non-negative solutions

$$(u, v) \in L^{\frac{2n}{n-2}}(\mathbf{R}^n) \times L^{\frac{2n}{n-2}}(\mathbf{R}^n)$$

to system (1):

$$u(x) = \int_{\mathbf{R}^n} \frac{u^\alpha v^\beta(y)}{|x-y|^{n-2}} dy$$

and

$$v(x) = \int_{\mathbf{R}^n} \frac{v^\alpha u^\beta(y)}{|x-y|^{n-2}} dy.$$

This equivalent form can be proved in the same way as in section 4 in the work of Chen, Li, and Ou [8].

We point out that when $u = v$, the elliptic system (1) reduces to the elliptic equation (13) with critical exponent. Then $u = v$ is in a special family of functions:

$$(14) \quad \phi_{x_o, t}(x) = c \left(\frac{t}{t^2 + |x - x_o|^2} \right)^{(n-2)/2},$$

where $t > 0$, $x_o \in \mathbf{R}^n$, with some positive constant c such that each $\phi_{x_o, t}(x)$ solves (13). This special family of functions is important in the study of (1).

Our results are motivated by the previous work [6], where Chen, Li, and Ou considered the more general system (8) and established the symmetry and monotonicity of the solutions. In [4], Chen and Li also obtained a regularity result of the solutions to (8). To establish the symmetry of the solutions to (8), Chen, Li, and Ou [8], [6], [7] introduced a new idea: an integral form of the method of moving planes. It is entirely different from the traditional method used for PDEs. Instead of relying on maximum principles, certain integral norms were estimated. The new method is a very powerful tool in studying qualitative properties of other integral equations and systems. In fact, following Chen, Li, and Ou's work, Jin and Li [13] studied the symmetry of the solutions to the more general system (7).

Using the maximum principle trick, de Figueiredo and Felmer [10] found an interesting Liouville-type result for certain type of elliptic systems. Recently Ma and Chen [21] discussed the Liouville-type theorem for the positive solutions to the elliptic system (10). They also discussed the radial symmetry property of nonnegative solutions to a higher-order differential equation [22]. Naito and Usami [24] studied the existence of nonoscillatory solutions to second-order elliptic systems of Emden–Fowler type. Some nonexistence results for the Emden–Fowler system can be found in [25]. We refer to [13], [9], and [24] for more references about elliptic systems.

We first prove the radial symmetry of the solutions to (1) with critical exponents. It is obvious that the radial symmetry of the solutions reduces (1) to a system of ODEs, which has the special solution pair $(\phi_{o, t}(x), \phi_{o, t}(x))$. To prove the uniqueness, we prove that $u(0) = v(0)$. Then by the uniqueness of the initial value problem for ODEs, we conclude that $u = v = \phi_{o, t}$. This is the key observation in establishing the uniqueness of positive solutions for (1) with critical exponents.

Theorems 1 and 2 will be proved in the next two sections.

2. Proof of the radial symmetry. We use the moving plane method introduced by Chen, Li, and Ou in [8] to study system (1). We recall the Hardy–Littlewood–Sobolev inequality:

$$(15) \quad |Tf|_p \leq C(n, p) |f|_{\frac{np}{n+2p}},$$

where $C(n, p)$ is a uniform positive constant and

$$Tf(x) = \int_{\mathbf{R}^n} |x - y|^{2-n} f(y) dy.$$

Proof of Theorem 1. For each $\lambda \in \mathbf{R}$, we denote

$$H_\lambda = \{x \in \mathbf{R}^n; x_1 < \lambda\}.$$

For each $x = (x_1, x') \in \mathbf{R}^n$, we let

$$x_\lambda = (2\lambda - x_1, x')$$

be the reflection point of x with respect to the hyperplane ∂H_λ . We let $e_1 = (1, 0, \dots, 0)$.

We define

$$u_\lambda(x) = u(x_\lambda), \quad B_\lambda^u = \{x \in H_\lambda; u_\lambda(x) > u(x)\}$$

and

$$v_\lambda(x) = v(x_\lambda), \quad B_\lambda^v = \{x \in H_\lambda; v_\lambda(x) > v(x)\}.$$

To do the moving plane method, we need the following formula, which is obtained by a change of variables.

$$u(x) = \int_{H_\lambda} \frac{u^\alpha v^\beta(y)}{|x - y|^{n-2}} dy + \int_{H_\lambda} \frac{u_\lambda^\alpha v_\lambda^\beta(y)}{|x_\lambda - y|^{n-2}} dy$$

and

$$u_\lambda(x) = \int_{H_\lambda} \frac{u^\alpha v^\beta(y)}{|x_\lambda - y|^{n-2}} dy + \int_{H_\lambda} \frac{u_\lambda^\alpha v_\lambda^\beta(y)}{|x - y|^{n-2}} dy.$$

Then we have

$$(16) \quad u_\lambda(x) - u(x) = \int_{H_\lambda} (u_\lambda^\alpha v_\lambda^\beta - u^\alpha v^\beta)(y) \left(\frac{1}{|x - y|^{n-2}} - \frac{1}{|x_\lambda - y|^{n-2}} \right) dy.$$

Note that for $x \in H_\lambda$, we have

$$\frac{1}{|x - y|^{n-2}} > \frac{1}{|x_\lambda - y|^{n-2}}.$$

Then for $x \in B_\lambda^u$, using the mean value theorem to

$$u_\lambda^\alpha v_\lambda^\beta - u^\alpha v^\beta = (u_\lambda^\alpha - u^\alpha) v_\lambda^\beta + u^\alpha (v_\lambda^\beta - v^\beta)$$

and after dropping the term of the integration on the region $v_\lambda - v \leq 0$, we have

$$(17) \quad \begin{cases} 0 \leq u_\lambda(x) - u(x) \\ \leq \alpha \int_{B_\lambda^u} \frac{u_\lambda^{\alpha-1} v_\lambda^\beta (u_\lambda - u)}{|x-y|^{n-2}} dy + \beta \int_{B_\lambda^v} \frac{u_\lambda^\alpha v_\lambda^{\beta-1} (v_\lambda - v)}{|x-y|^{n-2}} dy \\ := I + II. \end{cases}$$

Let $p = \frac{2n}{n-2}$. Using the Hardy–Littlewood–Sobolev inequality (15) we can bound the first term I in (17) by

$$(18) \quad \begin{cases} |I|_p \leq C(n, p) |u_\lambda^{\alpha-1} v_\lambda^\beta (u_\lambda - u)|_{\frac{2n}{n+2}} \\ \leq C(n, p) |u_\lambda|_p^{\alpha-1} |v_\lambda|_p^\beta |u_\lambda - u|_p. \end{cases}$$

Here the integrations are over the set B_λ^u .

Using again the Hardy–Littlewood–Sobolev inequality (15) we can bound the first term II in (17) by

$$(19) \quad \begin{cases} |II|_p \leq C(n, p) |u_\lambda^\alpha v_\lambda^{\beta-1} (v_\lambda - v)|_{\frac{2n}{n+2}} \\ \leq C(n, p) |u_\lambda|_p^\alpha |v_\lambda|_p^{\beta-1} |v_\lambda - v|_p. \end{cases}$$

Here the integrations are over the domain B_λ^v . Hence, we have

$$(20) \quad \begin{aligned} & |u_\lambda - u|_{p, B_\lambda^u} \\ & \leq C(n, p) (|u_\lambda|_{p, B_\lambda^u}^{\alpha-1} |v_\lambda|_{p, B_\lambda^u}^\beta |u_\lambda - u|_{p, B_\lambda^u} + |u_\lambda|_{p, B_\lambda^v}^\alpha |v_\lambda|_{p, B_\lambda^v}^{\beta-1} |v_\lambda - v|_{p, B_\lambda^v}). \end{aligned}$$

Similarly, we have the following formulae for v and v_λ :

$$v(x) = \int_{H_\lambda} \frac{v^\alpha u^\beta(y)}{|x-y|^{n-2}} dy + \int_{H_\lambda} \frac{v_\lambda^\alpha u_\lambda^\beta(y)}{|x_\lambda - y|^{n-2}} dy$$

and

$$v_\lambda(x) = \int_{H_\lambda} \frac{v^\alpha u^\beta(y)}{|x_\lambda - y|^{n-2}} dy + \int_{H_\lambda} \frac{v_\lambda^\alpha u_\lambda^\beta(y)}{|x - y|^{n-2}} dy.$$

Then we have the following estimate:

$$(21) \quad \begin{aligned} & |v_\lambda - v|_{p, B_\lambda^v} \\ & \leq C(n, p) (|v_\lambda|_{p, B_\lambda^v}^{\alpha-1} |u_\lambda|_{p, B_\lambda^v}^\beta |v_\lambda - v|_{p, B_\lambda^v} + |v_\lambda|_{p, B_\lambda^v}^\alpha |u_\lambda|_{p, B_\lambda^u}^{\beta-1} |u_\lambda - u|_{p, B_\lambda^u}). \end{aligned}$$

After these preparations, we can use the moving plane method as developed in [8] to prove the radial symmetry of the solutions.

At first, let's start the plane from infinity. Indeed, for $\lambda \gg 1$ large enough, we know that the quantities

$$|v_\lambda|_{p, B_\lambda^u}, |u_\lambda|_{p, B_\lambda^u}, |v_\lambda|_{p, B_\lambda^v}$$

and

$$|u_\lambda|_{p, B_\lambda^v}$$

all are small, which give us that

$$|u_\lambda - u|_{p, B_\lambda^u} \leq \frac{1}{4}(|v_\lambda - v|_{p, B_\lambda^v} + |u_\lambda - u|_{p, B_\lambda^u})$$

and

$$|v_\lambda - v|_{p, B_\lambda^v} \leq \frac{1}{4}(|v_\lambda - v|_{p, B_\lambda^v} + |u_\lambda - u|_{p, B_\lambda^u}).$$

These imply that $|u_\lambda - u|_{p, B_\lambda^u} = 0$ and $|v_\lambda - v|_{p, B_\lambda^v} = 0$. In other words, $B_\lambda^u = \phi$ and $B_\lambda^v = \phi$.

Next we define

$$\lambda_0 = \inf\{\lambda \in \mathbf{R}; B_{\lambda'}^u = B_{\lambda'}^v = \phi \forall \lambda' \geq \lambda\}.$$

Then it comes from the fact that $u(x) \rightarrow 0$ as $|x| \rightarrow \infty$ and $u(x) > 0$ in \mathbf{R}^n that $\lambda_0 < +\infty$. By the definition of λ_0 , we have $u_{\lambda_0}(x) \leq u(x)$ for $x \in H_{\lambda_0}$. If u or v is not symmetric in the x_1 direction at λ_0 , then using the expression (16), we see that $u_{\lambda_0}(x) < u(x)$ and $v_{\lambda_0}(x) < v(x)$ for $x \in H_{\lambda_0}$. This implies that both of $B_{\lambda_0}^u$ and $B_{\lambda_0}^v$ are empty. Then by the absolute continuity of the integral with respect to the domain we see that $|v_\lambda|_{p, B_\lambda^v}$, $|u_\lambda|_{p, B_\lambda^u}$, $|v_\lambda - v|_{p, B_\lambda^v}$, and $|u_\lambda - u|_{p, B_\lambda^u}$ are small for λ close to λ_0 and enable us to repeat the above argument showing that both of B_λ^u and B_λ^v are empty for λ close to λ_0 . This contradicts the definition of λ_0 .

3. Proof of the uniqueness. In some sense, the proof of Theorem 2 is just at hand by using the integral expression of the solution pair (u, v) .

Proof of Theorem 2. Let $(u, v) \in L^{\frac{2n}{n-2}}(\mathbf{R}^n) \times L^{\frac{2n}{n-2}}(\mathbf{R}^n)$ be a pair of solutions to system (1). By Theorem 1, we know that u and v are radial symmetric about some point x_0 , say, $x_0 = 0$. We will consider only the case when $1 \leq \alpha < \beta < \frac{n+2}{n-2}$, and the case when $\alpha = \beta$ can be done similarly.

Since $u \in L^{\frac{2n}{n-2}}(\mathbf{R}^n)$ and $v \in L^{\frac{2n}{n-2}}(\mathbf{R}^n)$, using the same method in [8], we have that $u \in C^2(\mathbf{R}^n)$ and $v \in C^2(\mathbf{R}^n)$ with

$$u(x) \rightarrow 0, \quad v(x) \rightarrow 0,$$

as $|x| \rightarrow \infty$.

Since our solution u is radial symmetric, we can write, in polar coordinates, the first equation in (1) as

$$(r^{n-1}u'(r))' = -r^{n-1}u(r)^\alpha v(r)^\beta,$$

where $r = |x|$.

Integrating both sides of the above equation from 0 to r yields

$$r^{n-1}u'(r) = - \int_0^r s^{n-1}u^\alpha v^\beta(s) ds.$$

It follows by another integration that

$$(22) \quad u(r) = u(0) - \int_0^r \frac{1}{\tau^{n-1}} \int_0^\tau s^{n-1}u^\alpha v^\beta ds d\tau.$$

Similarly, for $v(r)$, we have

$$(23) \quad v(r) = v(0) - \int_0^r \frac{1}{\tau^{n-1}} \int_0^\tau s^{n-1}v^\alpha u^\beta ds d\tau.$$

As we mentioned in the introduction, we need to show only that $u(0) = v(0)$. Otherwise, suppose that

$$(24) \quad u(0) < v(0).$$

Then by continuity, for all small $r > 0$,

$$(25) \quad u(r) < v(r).$$

In other words, there exists an $R > 0$, such that

$$(26) \quad u(r) < v(r) \quad \forall r \in (0, R).$$

Let R_o be the supreme value of R such that (26) holds. Then $R_o \leq \infty$ and $u(R_o) = v(R_o)$, where we have used the fact that $u(+\infty) = v(+\infty) = 0$. By the definition of R_o and $\alpha < \beta$, we have that

$$(27) \quad u(r)^\alpha v(r)^\beta > v(r)^\alpha u(r)^\beta \quad \forall r \in (0, R_o).$$

Then we have from (22) and (23) that

$$0 > u(0) - v(0) = \int_0^{R_o} \frac{1}{\tau^{n-1}} \int_0^\tau s^{n-1} (u^\alpha v^\beta - u^\beta v^\alpha)(s) ds d\tau > 0.$$

This is impossible.

Similarly, one can show that $u(0) > v(0)$ is impossible too. Therefore, we must have

$$u(0) = v(0).$$

Finally, by the standard ODE theory, we arrive at

$$u(r) \equiv v(r).$$

Hence, our elliptic system (1) has been reduced to the elliptic equation (13) with the critical exponent. By now, it is standard that our solution pair u and v is of the form (14). This completes the proof of Theorem 2.

Acknowledgments. The second author expresses his deep appreciation to Prof. W. X. Chen for helpful discussions on the method of moving planes. This work was done while the first author was visiting the Department of Mathematical Sciences in Tsinghua University, Beijing, which he would like to thank for their hospitality. Both authors thank the referees for detailed comments and suggestions on the original manuscript.

REFERENCES

- [1] J. BOURGAIN, *Global Solutions of Nonlinear Schrödinger Equations*, Amer. Math. Soc. Colloq. Publ. 46, AMS, Providence, RI, 1999.
- [2] L. CAFFARELLI, B. GIDAS, AND J. SPRUCK, *Asymptotic symmetry and local behavior of semi-linear elliptic equations with critical Sobolev growth*, Comm. Pure Appl. Math., 42 (1989), pp. 271–297.
- [3] W. CHEN AND C. LI, *Classification of solutions of some nonlinear elliptic equations*, Duke Math. J., 63 (1991), pp. 615–622.

- [4] W. CHEN AND C. LI, *Regularity of solutions for a system of integral equations*, Comm. Pure Appl. Anal., 4 (2005), pp. 1–8.
- [5] W. CHEN AND C. LI, *The best constant in weighted Hardy-Littlewood-Sobolev inequality*, Proc. Amer. Math. Soc., 136 (2008), pp. 955–962.
- [6] W. CHEN, C. LI, AND B. OU, *Classification of solutions for a system of integral equations*, Comm. Partial Differential Equations, 30 (2005), pp. 59–65.
- [7] W. CHEN, C. LI, AND B. OU, *Qualitative properties of solutions for an integral equation*, Discrete Contin. Dyn. Syst., 12 (2005), pp. 347–354.
- [8] W. CHEN, C. LI, AND B. OU, *Classification of solutions for an integral equation*, Comm. Pure Appl. Math., 59 (2006), pp. 330–343.
- [9] D. G. DE FIGUEIREDO AND P. L. FELMER, *On superquadratic elliptic systems*, Trans. Amer. Math. Soc., 343 (1994), pp. 99–116.
- [10] D. G. DE FIGUEIREDO AND P. L. FELMER, *A Liouville-type theorem for elliptic systems*, Ann. Sc. Norm. Super. Pisa Cl. Sci. (4), 21 (1994), pp. 387–397.
- [11] B. GIDAS, W. M. NI, AND L. NIRENBERG, *Symmetry of positive solutions of nonlinear elliptic equations in \mathbf{R}^n* , in Mathematical Analysis and Applications, Part A, Adv. Math. Suppl. Stud. 7A, Academic Press, New York, 1981, pp. 369–402.
- [12] F. B. HANG, *On the integral systems related to Hardy-Littlewood-Sobolev inequality*, Math. Res. Lett., 14 (2007), pp. 373–383.
- [13] C. JIN AND C. LI, *Symmetry of solutions to some integral equations*, Proc. Amer. Math. Soc., 134 (2006), pp. 1661–1670.
- [14] C. JIN AND C. LI, *Quantitative analysis of some system of integral equations*, Calc. Var. Partial Differential Equations, 26 (2006), pp. 447–457.
- [15] T. KANNA AND M. LAKSHMANAN, *Exact soliton solutions, shape changing collisions, and partially coherent solitons in coupled nonlinear Schrödinger equations*, Phys. Rev. Lett., 86 (2001), pp. 5043–5046.
- [16] C. LI, *Local asymptotic symmetry of singular solutions to nonlinear elliptic equations*, Invent. Math., 123 (1996), pp. 221–231.
- [17] C. LI AND J. LIM, *The singularity analysis of solutions to some integral equations*, Commun. Pure Appl. Anal., 6 (2007), pp. 453–464.
- [18] E. LIEB, *Sharp constants in the Hardy-Littlewood-Sobolev and related inequalities*, Ann. of Math., 118 (1983), pp. 349–374.
- [19] T. C. LIN AND J. WEI, *Ground state of N coupled nonlinear Schrödinger equations in \mathbf{R}^n , $n \leq 3$* , Commun. Math. Phys., 255 (2005), pp. 629–653.
- [20] T. C. LIN AND J. WEI, *Spikes in two coupled nonlinear Schrödinger equations*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 22 (2005), pp. 403–439.
- [21] L. MA AND D. Z. CHEN, *A Liouville type theorem for an integral system*, Comm. Pure Appl. Anal., 5 (2006), pp. 855–859.
- [22] L. MA AND D. Z. CHEN, *Radial symmetry and monotonicity for an integral equation*, J. Math. Anal. Appl., 342 (2008), pp. 943–949.
- [23] L. MA AND L. ZHAO, *Sharp thresholds of blow up and global existence for the coupled nonlinear Schrödinger system*, J. Math. Phys., 49 (2008), article 062103.
- [24] M. NAITO AND H. USAMI, *Existence of nonoscillatory solutions to second-order elliptic systems of Emden-Fowler type*, Indiana Univ. Math. J., 55 (2006), pp. 317–339.
- [25] J. SERRIN AND H. ZOU, *Non-existence of positive solutions of Lane-Emden systems*, Differential Integral Equations, 9 (1996), pp. 635–653.
- [26] E. M. STEIN AND G. WEISS, *Fractional integrals in n -dimensional Euclidean space*, J. Math. Mech., 7 (1958), pp. 503–514.
- [27] J. WEI AND X. XU, *Classification of solutions of higher order conformally invariant equations*, Math. Ann., 313 (1999), pp. 207–228.
- [28] M. I. WEINSTEIN, *Nonlinear Schrödinger equations and sharp interpolate estimates*, Comm. Math. Phys., 87 (1983), pp. 567–576.

STABILITY OF TRAVELING WAVES IN QUASI-LINEAR HYPERBOLIC SYSTEMS WITH RELAXATION AND DIFFUSION*

TONG LI†

Abstract. We establish the existence and the stability of traveling wave solutions of a quasi-linear hyperbolic system with both relaxation and diffusion. The traveling wave solutions are shown to be asymptotically stable under small disturbances and under the subcharacteristic condition using a weighted energy method. The delicate balance between the relaxation and the diffusion that leads to the stability of the traveling waves is identified; namely, the diffusion coefficient is bounded by a constant multiple of the relaxation time. Such a result provides an important first step toward the understanding of the transition from stability to instability as parameters vary in physical problems involving both relaxation and diffusion.

Key words. stability, traveling wave, relaxation, diffusion, weighted energy method, traffic flow

AMS subject classifications. 35B30, 35B40, 35L65, 76L05, 90B20

DOI. 10.1137/070690638

1. Introduction. The phenomenon of relaxation arises in many physical problems such as kinetic relaxation to fluid dynamics, gases not in local thermodynamic equilibrium, elasticity with memory, phase transitions, shallow water waves, and traffic flows. A remarkable development of the stability theory for various relaxation systems has appeared in past decades; see, e.g., [1, 2, 8, 10, 11, 13, 15, 16, 18, 23]. The real physical problems usually involve both relaxation and diffusion. It was found by asymptotic analysis and numerical simulations that the fine interplay between the relaxation and the diffusion may enhance physically interesting behavior such as soliton waves and oscillatory solutions [4, 3, 5, 6, 7, 9, 14, 17, 20, 22, 24]. However, the rigorous stability theory for such systems has not been well studied. In the current paper, we establish rigorously the existence and the stability of traveling wave solutions of a quasi-linear hyperbolic system with both relaxation and diffusion. The delicate balance between the relaxation and the diffusion that leads to the nonlinear stability of the traveling wave fronts is identified to occur when the diffusion coefficient is bounded by a constant multiple of the relaxation time. Such a result provides an important first step toward understanding the transition from stability to instability as the diffusion coefficient and the relaxation time vary in the physical problems.

Consider the following quasi-linear hyperbolic system with relaxation and diffusion:

$$(1) \quad \begin{cases} v_t - u_x = 0, \\ u_t + p(v)_x = \frac{1}{\tau}(u_e(v) - u) + \mu u_{xx} \end{cases}$$

subject to the initial data

$$(2) \quad (v, u)(x, 0) = (v_0, u_0)(x) \rightarrow (v_{\pm}, u_{\pm}) \text{ as } x \rightarrow \pm\infty, \quad u_{\pm} = u_e(v_{\pm}),$$

where, in the context of traffic flows, v is specific volume and u velocity. The first equation in (1) is a conservation law, while the second describes drivers' acceleration

*Received by the editors May 6, 2007; accepted for publication (in revised form) June 20, 2008; published electronically October 17, 2008.

<http://www.siam.org/journals/sima/40-3/69063.html>

†Department of Mathematics, University of Iowa, Iowa City, IA 52242 (tli@math.uiowa.edu).

behavior. The first term on the right-hand side of the second equation in (1) expresses the tendency of traffic at a given specific volume v to relax to some equilibrium speed u_e satisfying

$$u'_e(v) > 0.$$

The parameter $\tau > 0$ corresponds to drivers' response time to the traffic. The second term on the left-hand side of the second equation in (1) is an anticipation factor: drivers slow down at the sight of an increase in traffic density ahead. The function p is the so-called traffic pressure satisfying

$$(3) \quad p'(v) < 0.$$

The last term on the right-hand side of the second equation in (1) models viscosity with coefficient $\mu > 0$, a presumed tendency to adjust one's speed to that of the surrounding traffic.

A strict subcharacteristic condition,

$$(4) \quad -\sqrt{-p'(v)} < u'_e(v) < \sqrt{-p'(v)},$$

is imposed for all v under consideration. Subcharacteristic condition (4) is a necessary condition for linear stability (Whitham [23]) and for nonlinear stability (Li and Liu [10] and Liu [15]) when $\mu = 0$.

The purpose of this paper is to show the existence and stability of the traveling wave solutions $(V, U)(x - st)$ of (1) under appropriate conditions on the relaxation time τ and the diffusion coefficient μ .

A traveling wave solution is a solution of the form

$$(v, u)(x, t) = (V, U)(x - st) \equiv (V, U)(z),$$

where $z = x - st$, satisfying

$$(5) \quad (V, U)(z) \rightarrow (v_{\pm}, u_{\pm}) \text{ as } z \rightarrow \pm\infty, \quad u_{\pm} = u_e(v_{\pm}),$$

where v_+ , v_- , and s satisfy the Rankine–Hugoniot condition

$$(6) \quad -s(v_+ - v_-) - u_e(v_+) + u_e(v_-) = 0$$

and the entropy condition

$$(7) \quad (v_+ - v_-)(u_e(v) - u_e(v_{\pm}) + s(v - v_{\pm})) < 0$$

for all v in between v_+ and v_- . Indeed, the corresponding jump (v_-, v_+) is an admissible shock of the equilibrium equation

$$(8) \quad v_t - u_e(v)_x = 0.$$

Furthermore, the speed s of the traveling wave must be subcharacteristic; i.e., for all $z \in R$ it holds

$$(9) \quad -\sqrt{-p'(V(z))} < s < \sqrt{-p'(V(z))}.$$

For a weight function $w \geq 0$, L^2_w denotes the space of measurable functions f satisfying $\sqrt{w}f \in L^2$ with norm

$$\|f\|_{L^2_w} = \left(\int w(x)|f(x)|^2 dx \right)^{1/2}.$$

$H_w^j, j \geq 0$, denotes the weighted Sobolev space with norm

$$\|f\|_{H_w^j} = \left(\sum_{k=0}^j \|\partial_x^k f\|_w^2 \right)^{1/2}.$$

We now state our main results.

THEOREM 1.1. *Suppose that subcharacteristic conditions (4) and (9), the Rankine–Hugoniot condition (6), and the entropy condition (7) hold, and suppose that the diffusion coefficient is appropriately small,*

$$(10) \quad 0 < \mu \leq m\tau,$$

for some $m > 0$ critically depending on subcharacteristic conditions (4) and (9). Then there exists a traveling wave solution $(V, U)(x - st)$ of (1) and (5), which is unique up to a shift.

Moreover, there exists a constant $\varepsilon_0 > 0$ such that if

$$|v_- - v_+| + \|v_0(\cdot) - V(\cdot + x_0)\|_2 + \|u_0(\cdot) - U(\cdot + x_0)\|_2 + \|(\phi_0, \phi_{0,x}, \psi_0)\|_{L_w^2} \leq \varepsilon_0,$$

where x_0 is determined by

$$(11) \quad \int_{-\infty}^{+\infty} (v_0 - V)(x)dx = x_0(v_+ - v_-)$$

and

$$(\phi_0, \psi_0)(x) = \left(\int_{-\infty}^x (v_0(y) - V(y + x_0))dy, u_0(x) - U(x + x_0) \right)$$

and where the weight function w is defined as

$$(12) \quad w(V(x + x_0)) = \frac{(V(x + x_0) - v_+)(V(x + x_0) - v_-)}{Q(V(x + x_0))},$$

then the Cauchy problem (1), (2) has a unique global solution $(v, u)(x, t)$ satisfying

$$v(x, t) - V(x + x_0 - st), u(x, t) - U(x + x_0 - st) \in C^0(0, \infty; H^2 \cap L_w^2) \cap L^2(0, \infty; H^2 \cap L_w^2)$$

and

$$(13) \quad \sup_{x \in R} |(v, u)(x, t) - (V, U)(x + x_0 - st)| \rightarrow 0 \quad \text{as } t \rightarrow +\infty.$$

Previously, Jin and Liu [4] derived that the weakly nonlinear limit of the relaxation system of Jin and Xin [2] with a viscous term is governed by the Kortweg–de Vries (KdV) equation and dispersive waves are enhanced when the diffusion coefficient dominates in the sense that

$$\tau \ll O(\tau^{\frac{1}{4}}) \ll \mu.$$

When the diffusion coefficient is not too small, Kerner and Konhäuser [6], Kurtze and Hong [7], Lee, Lee, and Kim [9], and Ou et al. [22] were able to derive a modified KdV equation and obtain soliton-like solutions and oscillatory solutions for traffic flow models similar to (1) by asymptotic analysis and numerical simulations. The

phenomena occur in other applications. Keener [5] studied the diffusion induced insulin oscillatory secretion. He found that oscillatory secretion of insulin results from an important interplay between flow rate of the reactor and insulin diffusion.

In the absence of the diffusion term,

$$\mu = 0,$$

Li and Liu [10] established the nonlinear stability of traveling waves of (1) under the subcharacteristic conditions (4) and (9). The stability result is a consequence of the nonlinearity, the relaxation, and the subcharacteristic conditions (4) and (9).

In the current paper, we prove that, in addition to assuming the subcharacteristic conditions (4) and (9), if the diffusion coefficient is appropriately small (10), then traveling wave solutions of (1) are stable. In the context of traffic flows, our results indicate that the stability of the traveling fronts is guaranteed if the diffusion mechanism, the tendency to adjust one’s speed to that of the surrounding traffic, is dominated by the relaxation mechanism, the tendency of traffic to relax to some equilibrium speed.

The plan of this paper is the following. In section 2, we prove the existence of traveling wave solutions of the relaxation system with viscosity (1) by the phase plane analysis provided that the diffusion coefficient is appropriately small (10). In section 3, the stability problem is reformulated in terms of perturbations to the underlying traveling wave. Section 4 is devoted to establishing the desired a priori estimates for the nonlinear stability of the traveling wave solutions. When the equilibrium function $u_e(v)$ is nonconvex, there are traveling wave solutions that decay at an algebraic decay rate at infinity. A weighted energy method was developed to establish the stability of such degenerate traveling wave profiles [12, 10, 16, 19]. We adapt the weighted energy method developed in [12, 10, 16, 19] to a system of conservation laws with both relaxation and diffusion (1). It is under condition (10) that we prove the desired weighted energy estimates that lead to the stability of the traveling wave solutions.

2. Existence of traveling wave profiles. In this section, we prove the existence of traveling wave solutions of the relaxation system with viscosity (1). We show that, under the subcharacteristic conditions (4) and (9), and when the diffusion coefficient is appropriately small (10), there is a traveling wave profile of (1) connecting two states satisfying the Rankine–Hugoniot condition (6) and the entropy condition (7).

LEMMA 2.1. *Assume that v_+ , v_- , and s satisfy the Rankine–Hugoniot condition (6), the entropy condition (7), and the subcharacteristic conditions (4), (9), and that the diffusion coefficient is appropriately small (10). There exists a traveling wave solution $(V, U)(x - st)$ of (1) with boundary condition (5), which is unique up to a shift.*

Moreover, the profile $V = V(z)$ is monotone and

$$(14) \quad (v_+ - v_-)(u_e(V(z)) - u_e(v_{\pm}) + s(V(z) - v_{\pm})) < 0$$

for all $z \in R$.

Proof. Substituting

$$(v, u)(x, t) = (V, U)(z), \quad z = x - st,$$

into (1), we have

$$(15) \quad \begin{cases} -sV_z - U_z = 0, \\ -sU_z + p(V)_z = \frac{u_e(V)-U}{\tau} + \mu U_{zz}. \end{cases}$$

Integrating the first equation of (15) over $(\pm\infty, z)$ and using boundary condition (5) yield

$$(16) \quad -sV - U = -sv_{\pm} - u_{\pm} = -sv_{\pm} - u_e(v_{\pm}).$$

Therefore the Rankine–Hugoniot condition (6) for the equilibrium equation (8) is satisfied by the end states of the traveling wave of the relaxation system (1).

We derive from (15), (16), and (5) that V satisfies a second order nonlinear differential equation

$$(17) \quad V_{zz} + \frac{p'(V) + s^2}{s\mu} V_z - \frac{u_e(V) - u_e(v_{\pm}) + s(V - v_{\pm})}{\tau s\mu} = 0$$

and

$$(18) \quad V(z) \rightarrow v_{\pm} \text{ as } z \rightarrow \pm\infty.$$

We will establish the existence of solutions for (17) and (18) by phase plane analysis.

Let $W = V'$. Equation (17) is rewritten as a system of first order differential equations

$$(19) \quad \begin{cases} V' = W, \\ W' = -\frac{p'(V)+s^2}{s\mu}W + \frac{u_e(V)-u_e(v_{\pm})+s(V-v_{\pm})}{\tau s\mu}. \end{cases}$$

$(v_{\pm}, 0)$ are two equilibrium points of (19). Without loss of generality, we assume that $v_- > v_+$. Other cases can be treated similarly.

The entropy condition (7) implies that

$$(20) \quad Q(V) \equiv -\frac{1}{\tau}(u_e(V) - u_e(v_{\pm}) + s(V - v_{\pm})) < 0, \quad v_+ < V < v_-,$$

and consequently

$$-u'_e(v_+) \leq s < -u'_e(v_-) \leq 0.$$

Consider

$$(21) \quad -u'_e(v_+) < s < -u'_e(v_-) < 0.$$

The degenerate cases $-u'_e(v_+) = s$ and $s = -u'_e(v_-)$ can be treated similarly.

Linearize the system of nonlinear equations (19) at equilibrium points $(v_{\pm}, 0)$ to obtain

$$(22) \quad \begin{cases} (V - v_{\pm})' = W, \\ W' = -\frac{p'(v_{\pm})+s^2}{s\mu}W + \frac{u'_e(v_{\pm})+s}{\tau s\mu}(V - v_{\pm}). \end{cases}$$

The eigenvalues of the Jacobian of (22) satisfy

$$\lambda^2 + \lambda \frac{p'(v_{\pm}) + s^2}{s\mu} - \frac{u'_e(v_{\pm}) + s}{\tau s\mu} = 0.$$

The entropy condition (21), the subcharacteristic condition (9), and $s < 0$ imply that

$$(23) \quad -\frac{u'_e(v_-) + s}{\tau s\mu} < 0, \quad \frac{p'(v_-) + s^2}{s\mu} > 0.$$

Thus $(v_-, 0)$ is a saddle with eigenvalues

$$\lambda_{1-} > 0 > \lambda_{2-}$$

and the corresponding right eigenvectors

$$r_{1-} = (1, \lambda_{1-})^T, \quad r_{2-} = (1, \lambda_{2-})^T.$$

Therefore the unstable manifold of (19) at $(v_-, 0)$ is tangent to r_{1-} .

Similarly, the entropy condition (21), the subcharacteristic condition (9), and $s < 0$ imply that

$$(24) \quad -\frac{u'_e(v_+) + s}{\tau s \mu} > 0, \quad \frac{p'(v_+) + s^2}{s \mu} > 0.$$

Thus $(v_+, 0)$ is a stable node provided that

$$\left(\frac{p'(v_+) + s^2}{s \mu}\right)^2 + 4\frac{u'_e(v_+) + s}{\tau s \mu} > 0$$

or

$$0 < \frac{\mu}{\tau} < -\frac{(p'(v_+) + s^2)^2}{4s(u'_e(v_+) + s)}.$$

The above inequality holds provided that the diffusion coefficient is appropriately small (10).

The eigenvalues of the Jacobian of (22) satisfy

$$0 > \lambda_{1+} > \lambda_{2+},$$

and the corresponding right eigenvectors are

$$r_{1+} = (1, \lambda_{1+})^T, \quad r_{2+} = (1, \lambda_{2+})^T.$$

The nullclines of (19) consist of

$$(25) \quad W = 0$$

and

$$(26) \quad W = \frac{u_e(V) - u_e(v_{\pm}) + s(V - v_{\pm})}{\tau(p'(V) + s^2)} \equiv g(V).$$

The nullclines intersect at the equilibrium points $(v_{\pm}, 0)$, where

$$g(v_{\pm}) = 0.$$

Moreover, by using the entropy condition (21) and the subcharacteristic condition (9), we have

$$(27) \quad g'(v_-) = \frac{u'_e(v_-) + s}{(p'(v_-) + s^2)\tau} > 0, \quad g'(v_+) = \frac{u'_e(v_+) + s}{(p'(v_+) + s^2)\tau} < 0,$$

and

$$g(V) < 0, \quad v_+ < V < v_-.$$

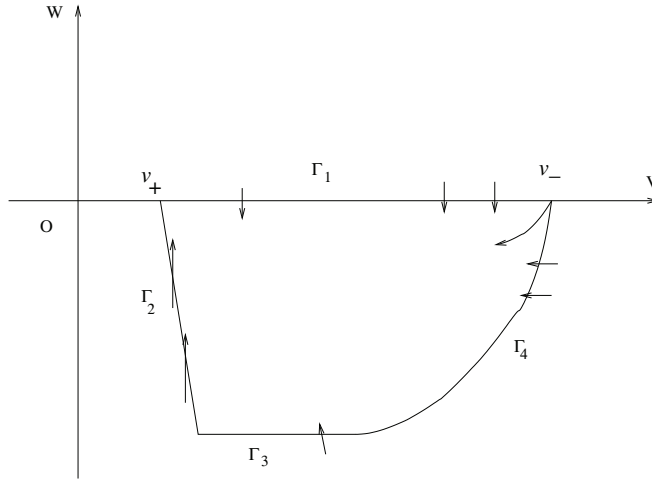


FIG. 1. The invariant domain D .

Let $v_0 \in (v_+, v_-)$ be the minimum point of g on $[v_+, v_-]$ and

$$W_0 = g(v_0) = \min_{v_+ < V < v_-} g(V).$$

Then

$$W_0 < 0, \quad g'(v_0) = 0.$$

Let D be a region bounded by the following curves on the phase plane (see Figure 1):

$$\Gamma_1 : W = 0, \quad v_+ \leq V \leq v_-,$$

$$\Gamma_2 : W = -a(V - v_+), \quad v_+ \leq V \leq v_1,$$

$$\Gamma_3 : W = W_0 = g(v_0), \quad v_1 \leq V \leq v_0,$$

$$\Gamma_4 : W = g(V), \quad v_0 \leq V \leq v_-,$$

where $a > 0$ is to be determined and (v_1, W_0) is the intersection point of Γ_2 and Γ_3 satisfying

$$W_0 = -a(v_1 - v_+), \quad v_+ \leq v_1 \leq v_0.$$

We will show that, under the condition that the diffusion coefficient is appropriately small (10), a can be chosen such that D is an invariant region of (19).

(i) From (19) and (21) we derive that on the top boundary of D , Γ_1 , $W' < 0$, the flow points downward toward the interior of D .

(ii) On the bottom boundary of D , Γ_3 , the flow points in the upper left direction, which points toward the interior of D . This is because Γ_3 lies below the nullcline $W = g(V)$ and thus $W' > 0$.

(iii) On the right boundary of D , Γ_4 , which is part of the nullcline $W = g(V)$, the flow points leftward toward the interior of D . This is because $V' = W \leq 0$ and $W' = 0$.

(iv) It remains to show that the flow points toward the interior of D on the left boundary of D , Γ_2 .

Statement (iv) is true if we can choose $a > 0$ such that

$$(28) \quad \left. \frac{dW}{dV} \right|_{\Gamma_2} < -a < 0, \quad v_+ \leq V \leq v_1.$$

We show that such $a > 0$ exists if the diffusion coefficient μ is appropriately small (10). Noting (24) and (27), we have

$$\lambda_{1+} = \frac{2 \frac{u'_e(v_+) + s}{\tau s \mu}}{\sqrt{\left(\frac{p'(v_+) + s^2}{s \mu}\right)^2 + 4 \frac{u'_e(v_+) + s}{\tau s \mu} + \frac{p'(v_+) + s^2}{s \mu}}} < g'(v_+) < 0.$$

Thus Γ_2 is below $W = g(V)$ for V near v_+ , $V \in (v_+, v_+ + \epsilon)$ for some $\epsilon > 0$.

Indeed, if $\lambda_{1+} < \min_{[v_+, v_0]} g'(V)$, we will choose a satisfying

$$\lambda_{1+} < -a < \min_{[v_+, v_0]} g'(V) < 0.$$

Thus Γ_2 is below $W = g(V)$ for $v_+ \leq V \leq v_1$.

We now show that (28) is satisfied. Since

$$\lambda_{2+} < \lambda_{1+} < -a < g'(v_+) < 0,$$

there is an $\epsilon > 0$, such that

$$(29) \quad \left. \frac{dW}{dV} \right|_{\Gamma_2} < -a < 0, \quad v_+ \leq V \leq v_+ + \epsilon.$$

If V is away from v_+ , then there is a $\delta > 0$ such that

$$0 < \frac{g(V)}{W} \leq 1 - \delta, \quad v_+ + \epsilon \leq V \leq v_1.$$

From (19) we have

$$\left. \frac{dW}{dV} \right|_{\Gamma_2} = \frac{W'}{V'} = -\frac{p'(V) + s^2}{s \mu} \left(1 - \frac{g(V)}{W}\right) \leq -\delta \frac{p'(V) + s^2}{s \mu} < 0$$

for $v_+ + \epsilon \leq V \leq v_1$. Noting the definition of g and (26) and assuming that the diffusion coefficient μ is appropriately small (10), we have

$$\left. \frac{dW}{dV} \right|_{\Gamma_2} \leq -\delta \frac{p'(V) + s^2}{s \mu} < g'(V) < 0, \quad v_+ + \epsilon \leq V \leq v_1.$$

Thus

$$(30) \quad \left. \frac{dW}{dV} \right|_{\Gamma_2} < -a < 0, \quad v_+ + \epsilon \leq V \leq v_1.$$

Combining (29) and (30), we prove that there is an a satisfying (28). If $\min_{[v_+, v_0]} g'(V) < \lambda_{1+}$, we consider a piecewise linear curve Γ_2 instead. The above argument can be modified to prove (28).

Therefore D is an invariant region of (19).

Now we look for a trajectory connecting the two equilibrium points $(v_{\pm}, 0)$. Noting (23) and (27), we have

$$0 < \lambda_{1-} = \frac{2 \frac{u'_e(v_-)+s}{\tau s \mu}}{\sqrt{\left(\frac{p'(v_-)+s^2}{s \mu}\right)^2 + 4 \frac{u'_e(v_-)+s}{\tau s \mu} + \frac{p'(v_-)+s^2}{s \mu}}} < g'(v_-).$$

Since the unstable manifold of (19) at $(v_-, 0)$ is tangent to $r_{1-} = (1, \lambda_{1-})^T$ with $0 < \lambda_{1-} < g'(v_-)$, there is a trajectory of (19) originating from $(v_-, 0)$ that enters the domain D . Since D is an invariant region, such a trajectory is trapped inside D . Therefore the trajectory flows into the stable node $(v_+, 0)$ as $t \rightarrow +\infty$.

We thus find a solution V of (17) and (18) which is monotone decreasing. The entropy condition (7) follows from (20) and the above monotonicity. The existence of the traveling wave profile (V, U) of (15) and (5) is proved by using (16).

Lemma 2.1 is proved. \square

3. Reformulation of the stability problem. In this section, we reformulate the problem that the solution $(v, u)(x, t)$ of (1), (2) exists globally and approaches a shifted traveling wave solution $(V, U)(x+x_0-st)$ under the subcharacteristic condition (4) and under condition (10) as $t \rightarrow \infty$.

We will look for a solution of the following form:

$$(31) \quad (v, u)(x, t) = (V, U)(z + x_0) + (\phi_z, \psi)(z, t),$$

where $z = x - st$.

For initial data satisfying (11), the conservation law in (1) implies that

$$\begin{aligned} \phi(\pm\infty, t) &= \int_{-\infty}^{+\infty} (v(x, t) - V(x + x_0 - st)) dx = \int_{-\infty}^{+\infty} (v_0(x) - V(x + x_0)) dx \\ &= x_0(v_+ - v_-) - \int_{-\infty}^{+\infty} (V(x + x_0) - V(x)) dx = 0. \end{aligned}$$

For simplicity of notation, we assume the shift $x_0 = 0$ in the rest of the paper.

We substitute (31) into (1), by virtue of (15), and integrate the first equation once with respect to z , to obtain that the perturbation (ϕ, ψ) satisfies

$$(32) \quad \begin{cases} \phi_t - s\phi_z - \psi = 0, \\ \psi_t - s\psi_z + (p(V + \phi_z) - p(V))_z \\ = \frac{1}{\tau}(u_e(V + \phi_z) - u_e(V) - \psi) + \mu\psi_{zz}. \end{cases}$$

The first equation of (32) gives

$$(33) \quad \psi = \phi_t - s\phi_z.$$

Substituting (33) into the second equation of (32), we get a closed equation for ϕ ,

$$\begin{aligned} L(\phi) &\equiv (\phi_t - s\phi_z)_t - s(\phi_t - s\phi_z)_z + (p'(V)\phi_z)_z + \frac{1}{\tau}\phi_t + \lambda\phi_z - \mu(\phi_t - s\phi_z)_{zz} \\ (34) \quad &= -F(V, \phi_z), \end{aligned}$$

where $Q < 0$ is defined in (20) and

$$\lambda = Q'(V) = -\frac{1}{\tau}(u'_e(V) + s),$$

and

$$F(V, \phi_z) = F_1 + F_2$$

with

$$(35) \quad F_1 := -\frac{1}{\tau}(u_e(V + \phi_z) - u_e(V) - u'_e(V)\phi_z),$$

$$(36) \quad F_2 := (p(V + \phi_z) - p(V) - p'(V)\phi_z)_z = (G(V, \phi_z)\phi_z^2)_z,$$

and

$$G(V, \phi_z) := \int_0^1 \int_0^1 p''(V + \theta\eta\phi_z)\theta d\theta d\eta$$

which is the error term due to nonlinearity of the function p .

The corresponding initial data for (34) become

$$(37) \quad \phi(z, 0) = \phi_0(z), \quad \phi_t(z, 0) = s\phi'_0(z) - \psi_0 = \phi_1(z).$$

The asymptotic stability of the profile (V, U) means that the perturbation (ϕ_z, ψ) decays to zero as $t \rightarrow +\infty$.

First, by noting (14), we have that the weight function defined in (12) satisfies $w(V) > 0$.

Now we introduce the solution space of the problem (34), (37) as follows:

$$X(0, T) = \{\phi(z, t) : \phi \in C^0([0, T]; H^3 \cap L^2_w) \cap C^1(0, T; H^2 \cap L^2_w), \phi_z, \phi_t \in L^2(0, T; H^2 \cap L^2_w)\}$$

with $0 < T \leq +\infty$.

By virtue of (33), we have

$$\psi \in C^0([0, T]; H^2 \cap L^2_w) \cap L^2(0, T; H^2 \cap L^2_w).$$

By the Sobolev embedding theorem, if we let

$$(38) \quad N(t) = \sup_{0 \leq s \leq t} \{ \|\phi(s)\|_3 + \|\phi_t(s)\|_2 + \|\phi(s)\|_{H^1_w} + \|\phi_t(s)\|_{L^2_w} \},$$

then

$$(39) \quad \sup_{z \in R} \{ |\phi|, |\phi_z|, |\phi_{zz}|, |\phi_t|, |\phi_{tz}| \} \leq CN(t).$$

Thus Theorem 1.1 is a consequence of the following theorem.

THEOREM 3.1. *Under the conditions of Theorem 1.1, there exist positive constants δ_1, C , and C_1 such that if $N(0) \leq \delta_1$, then the problem (34), (37) has a unique global solution $\phi \in X(0, +\infty)$ satisfying*

$$(40) \quad \begin{aligned} & \|\phi(t)\|_3^2 + \|\phi_t\|_2^2 + \|\phi\|_{H^1_w}^2 + \|\phi_t\|_{L^2_w}^2 + \int_0^t \|(\phi_t, \phi_z)(s)\|_2^2 ds \\ & + C_1\mu \left(\int_0^t \|(\phi_t - s\phi_z)_z\|_{L^2_w}^2 ds + \int_0^t \|(\phi_t - s\phi_z)_{zz}\|_1^2 ds \right) \leq CN^2(0) \end{aligned}$$

for $t \in [0, +\infty)$. Furthermore,

$$(41) \quad \sup_{z \in R} |(\phi_z, \phi_t)(z, t)| \rightarrow 0 \text{ as } t \rightarrow \infty.$$

If ϕ is the global solution in the above theorem, then (ϕ, ψ) , defined in (33), becomes a global solution of the problem (32) with $(\phi, \psi)(z, 0) = (\phi_0, \psi_0)(z)$, and consequently we have the desired solution of the problem (1), (2) through the relation (31). On the other hand, the solution of (1) is unique in the space $C^0(0, T; H^2 \cap L_w^2)$; therefore Theorem 1.1 follows from Theorem 3.1. The estimate (40) gives

$$\begin{aligned} \phi_t^2 + \phi_z^2 &= \int_{-\infty}^z (2\phi_t\phi_{tz} + 2\phi_z\phi_{zz})(y, t)dy \\ &\leq \left(\int_{-\infty}^{+\infty} (\phi_t^2 + \phi_z^2)dy \right)^{1/2} \left(\int_{-\infty}^{+\infty} (\phi_{tz}^2 + \phi_{zz}^2)dy \right)^{1/2} \rightarrow 0 \text{ as } t \rightarrow \infty \end{aligned}$$

as claimed in Theorem 1.1.

Global existence for ϕ will be derived from the following local existence theorem for ϕ combined with a priori estimates.

PROPOSITION 3.2 (local existence). *For any $\delta_0 > 0$, there exists a positive constant T_0 depending on δ_0 , such that if $\phi_0 \in H^3 \cap H_w^1$ and $\phi_1 \in H^2 \cap L_w^2$, with $N(0) < \delta_0/2$, then the problem (34), (37) has a unique solution $\phi \in X(0, T_0)$ satisfying*

$$(42) \quad N(t) < 2N(0)$$

for any $0 \leq t \leq T_0$.

PROPOSITION 3.3 (a priori estimates). *Let $\phi \in X(0, T)$ be a solution for a positive constant T ; then there exists a positive constant δ_2 independent of T such that if*

$$N(t) < \delta_2, \quad t \in [0, T],$$

then ϕ satisfies (40) for any $0 \leq t \leq T$.

The local existence is classical, so we omit the proof; cf. [21]. Proving Proposition 3.3 is our main task in the following section.

4. Energy estimates. In this section, we will complete the proof of our stability theorem by establishing the desired a priori estimates. The stability result is a consequence of the compressibility of the traveling wave profile (14), the subcharacteristic conditions (4) and (9), the fact that the diffusion coefficient is bounded by a constant multiple of the relaxation time (10), and the weighted energy estimates.

LEMMA 4.1. *Under the conditions of Theorem 1.1, there are positive constants C and C_1 such that if subcharacteristic condition (4) is satisfied for $v \in [v_+, v_-]$, then any solution $\phi \in X(0, T)$ of problem (34), (37) satisfies*

$$\begin{aligned} &\|\phi(t)\|_{H_w^1}^2 + \|\phi_t(t)\|_{L_w^2}^2 + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds \\ &+ \int_0^t \int_R |U_z|\phi^2 dz ds + C_1\mu \int_0^t \int_R w(\phi_t - s\phi_z)_z^2 dz ds \end{aligned}$$

$$(43) \quad \leq C \left\{ \|\phi_0\|_{H^1_w}^2 + \|\phi_1\|_{L^2_w}^2 + \int_0^t \int_R w|F|(|\phi| + |(\phi_t, \phi_z)|) dz ds \right\}$$

for $t \in [0, T]$.

Proof. Let $w = w(V) > 0$ be the weight function defined in (12).

Multiplying (34) by $2w(V)\phi$, we obtain

$$(44) \quad 2w(V)\phi L(\phi) = -2Fw(V)\phi.$$

The left-hand side of (44) is reduced to

$$\begin{aligned} & 2[(\phi_t - s\phi_z)_t - s(\phi_t - s\phi_z)_z + (p'(V)\phi_z)_z]w\phi + 2\left(\frac{1}{\tau}\phi_t + \lambda\phi_z\right)w\phi \\ & - 2\mu w\phi(\phi_t - s\phi_z)_{zz} \\ = & \left[\frac{1}{\tau}w\phi^2 + 2w\phi(\phi_t - s\phi_z)\right]_t - 2w(\phi_t - s\phi_z)^2 - 2wp'(V)\phi_z^2 + (p'(V)w_z)_z\phi^2 \\ & - (\lambda w)_z\phi^2 + sw_z(\phi^2)_t - s^2\{w_z(\phi^2)\}_z + s^2w_{zz}\phi^2 \\ & + \{-2sw\phi(\phi_t - s\phi_z) + 2p'(V)w\phi\phi_z - p'(V)w_z\phi^2 + \lambda w\phi^2 - 2\mu w\phi(\phi_t - s\phi_z)_z\}_z \\ & + 2\mu w_z\phi(\phi_t - s\phi_z)_z + 2\mu w\phi_z(\phi_t - s\phi_z)_z \\ = & \left[\frac{1}{\tau}w\phi^2 + 2w\phi(\phi_t - s\phi_z) + sw_z\phi^2\right]_t - 2w(\phi_t - s\phi_z)^2 - 2p'(V)w\phi_z^2 + A(V)\phi^2 \\ (45) \quad & + 2\mu w_z\phi(\phi_t - s\phi_z)_z + 2\mu w\phi_z(\phi_t - s\phi_z)_z + 2\mu sw'(V)V_{zz}\phi\phi_z + \{\dots\}_z, \end{aligned}$$

where $\{\dots\}_z$ denotes the terms which will disappear after integration with respect to $z \in R$ and

$$(46) \quad \begin{aligned} A(V) &= -\{(-p'(V) - s^2)w'(V)V_z + \lambda w - \mu sw'(V)V_{zz}\}_z \\ &= -\{w'(V)Q(V) + Q'(V)w\}_z \\ &= -\{wQ\}''(V)V_z, \end{aligned}$$

where the second equality is due to (17).

Since V is monotone decreasing and

$$(wQ)''(V) = 2,$$

therefore

$$(47) \quad A(V) = -2V_z > 0.$$

Then we calculate

$$(48) \quad 2(\phi_t - s\phi_z)wL(\phi) = -2F(\phi_t - s\phi_z)w.$$

The left-hand side of (48) is

$$\begin{aligned}
 & 2[(\phi_t - s\phi_z)_t - s(\phi_t - s\phi_z)_z + (p'(V)\phi_z)_z]w(\phi_t - s\phi_z) \\
 & + \frac{2}{\tau}w(\phi_t - s\phi_z)(\phi_t - s\phi_z - uV'_e(V)\phi_z) - 2\mu w(\phi_t - s\phi_z)(\phi_t - s\phi_z)_{zz} \\
 = & [w(\phi_t - s\phi_z)^2]_t + \left(\frac{2}{\tau}w + sw_z\right)(\phi_t - s\phi_z)^2 - 2p'(V)w_z\phi_z(\phi_t - s\phi_z) \\
 & - \frac{2}{\tau}wu'_e(V)\phi_z(\phi_t - s\phi_z) - [p'(V)w\phi_z^2]_t + [swp'(V)\phi_z^2]_z - s(wp'(V))_z\phi_z^2 \\
 & - [sw(\phi_t - s\phi_z)^2 - 2p'(V)w\phi_z(\phi_t - s\phi_z)]_z \\
 & - (2\mu w(\phi_t - s\phi_z)(\phi_t - s\phi_z)_z + 2\mu w_z(\phi_t - s\phi_z)(\phi_t - s\phi_z)_z + 2\mu w(\phi_t - s\phi_z)_z^2) \\
 = & [-wp'(V)\phi_z^2 + w(\phi_t - s\phi_z)^2]_t + \left(\frac{2}{\tau}w + sw_z\right)(\phi_t - s\phi_z)^2 \\
 & - s(wp'(V))_z\phi_z^2 - \frac{2}{\tau}u'_e(V)w\phi_z(\phi_t - s\phi_z) - 2p'(V)w_z\phi_z(\phi_t - s\phi_z) \\
 & + 2\mu w_z(\phi_t - s\phi_z)(\phi_t - s\phi_z)_z + 2\mu w(\phi_t - s\phi_z)_z^2 + \{\dots\}_z.
 \end{aligned}$$

Hence, the combination (44) + (48) × 2τ yields

$$\begin{aligned}
 & \{E_1(\phi, (\phi_t - s\phi_z)) + E_3(\phi_z)\}_t + E_2(\phi_z, (\phi_t - s\phi_z)) + E_4(\phi) + 4\mu\tau w(\phi_t - s\phi_z)_z^2 \\
 & \quad + 2\mu w_z\phi(\phi_t - s\phi_z)_z + 2\mu w\phi_z(\phi_t - s\phi_z)_z + 2\mu sw'(V)V_{zz}\phi\phi_z \\
 (49) \quad & + 4\mu\tau w_z(\phi_t - s\phi_z)(\phi_t - s\phi_z)_z + \{\dots\}_z = -2Fw\{\phi + 2\tau(\phi_t - s\phi_z)\},
 \end{aligned}$$

where

$$\begin{aligned}
 E_1(\phi, (\phi_t - s\phi_z)) &= 2\tau w(\phi_t - s\phi_z)^2 + 2w\phi(\phi_t - s\phi_z) + \left(\frac{1}{\tau}w + sw_z\right)\phi^2, \\
 E_3(\phi_z) &= -2\tau wp'(V)\phi_z^2, \\
 E_2(\phi_z, (\phi_t - s\phi_z)) &= 2(w + \tau sw_z)(\phi_t - s\phi_z)^2 \\
 & \quad + 4(-u'_e(V)w - \tau p'(V)w_z)\phi_z(\phi_t - s\phi_z) \\
 & \quad + 2(-wp'(V) - \tau s(wp'(V))_z)\phi_z^2, \\
 (50) \quad E_4(\phi) &= A(V)\phi^2,
 \end{aligned}$$

where $A(V)$ is defined in (46).

The discriminants of the quadratics $E_j (j = 1, 2)$ are, respectively,

$$D_1 = -4w[w + 2\tau sw_z],$$

$$D_2 = 16\{(-u'_e(V)w - \tau p'(V)w_z)^2 + (w + \tau sw_z)(wp'(V) + \tau s(wp'(V))_z)\}.$$

We claim that

$$(51) \quad D_1 < 0, \quad D_2 < 0,$$

provided that $|v_+ - v_-|$ is suitably small and subcharacteristic condition (4) is satisfied.

Indeed, the two inequalities in (51) are equivalent to

$$(52) \quad 1 + 2s\tau \frac{w_z}{w} > 0,$$

$$(53) \quad \left(-u'_e(V) - \tau p'(V) \frac{w_z}{w}\right)^2 < \left(1 + s\tau \frac{w_z}{w}\right) \left(-p'(V) - s\tau \frac{(wp'(V))_z}{w}\right).$$

Noticing that $|\tau(w(V))_z|$ and $|\tau(w(V)p'(V))_z|$ are small provided that $|v_+ - v_-|$ is suitably small and subcharacteristic condition (4) is satisfied, we derive inequalities (52) and (53). Thus condition (51) is satisfied.

Therefore there exist positive constants M_0 and M such that

$$(54) \quad \begin{cases} M_0 w \{\phi^2 + (\phi_t - s\phi_z)^2\} \leq E_1 \leq M w \{\phi^2 + (\phi_t - s\phi_z)^2\}, \\ M_0 w \{\phi_z^2 + (\phi_t - s\phi_z)^2\} \leq E_2. \end{cases}$$

Furthermore, (3) and (47) yield

$$(55) \quad \begin{cases} E_3 = -2\tau w p'(V) \phi_z^2 \geq 0, \\ E_4 \geq 2|V_z| \phi^2 \geq 0. \end{cases}$$

Now we estimate terms containing μ in (49).

First, note that $4\mu\tau w(\phi_t - s\phi_z)_z^2$ is a good term, namely, nonnegative. Other terms containing μ are estimated by using the Cauchy–Schwarz inequality, estimate (54), the fact that $|w(V)_z|$ and $|V_{zz}|$ are small, and (10). Indeed, let us estimate the second term

$$(56) \quad \begin{aligned} & |2\mu w \phi_z (\phi_t - s\phi_z)_z| \\ & \leq \mu\tau w (\phi_t - s\phi_z)_z^2 + \frac{\mu}{4\tau} w \phi_z^2 \leq \mu\tau w (\phi_t - s\phi_z)_z^2 + \frac{1}{4} M_0 w \phi_z^2 \end{aligned}$$

provided that $\frac{\mu}{\tau} \leq m \leq M_0$ (10).

Substituting estimates (54), (55), and (56) into (49) and integrating the result with respect to t and z , we arrive at the desired estimate (43).

This completes the proof of Lemma 4.1. \square

Next we estimate the higher order derivatives of ϕ .

Let $\phi_z = \Phi$; then

$$(57) \quad \begin{aligned} \partial_z L(\phi) &= (\phi_{zt} - s\phi_{zz})_t - s(\phi_{zt} - s\phi_{zz})_z + (p'(V)\phi_z)_{zz} \\ &\quad + \phi_{zt} + \lambda\phi_{zz} + \lambda_z\phi_z - \mu(\phi_{zt} - s\phi_{zz})_{zz} \\ &= L(\phi_z) + \lambda_z\phi_z + (p'(V)_z\phi_z)_z = L(\Phi) + \lambda_z\Phi + (p'(V)_z\Phi)_z. \end{aligned}$$

Multiplying the derivative of (34) with respect to z by $2\phi_z$ and $2(\phi_t - s\phi_z)_z$, respectively, we have

$$\begin{aligned} 2\partial_z L(\phi)\phi_z &= -2F_z\phi_z, \\ 2\partial_z L(\phi)(\phi_t - s\phi_z)_z &= -2F_z(\phi_t - s\phi_z)_z. \end{aligned}$$

By a similar argument used in obtaining estimates for the first order derivatives with $w = 1$, we have

$$(58) \quad \begin{aligned} & \left[\frac{1}{\tau} \Phi^2 + 2\Phi(\Phi_t - s\Phi_z) \right]_t - 2p'(V)\Phi_z^2 - 2(\Phi_t - s\Phi_z)^2 + \lambda_z\Phi^2 \\ & + 2\mu\Phi_z(\Phi_t - s\Phi_z)_z - 2p'(V)_z\Phi\Phi_z + \{\dots\}_z = -2F_z\Phi \end{aligned}$$

and

$$\begin{aligned}
 & [(\Phi_t - s\Phi_z)^2 - p'(V)\Phi_z^2]_t + \frac{2}{\tau}(\Phi_t - s\Phi_z)^2 - \frac{2}{\tau}u'_e(V)\Phi_z(\Phi_t - s\Phi_z) \\
 & - sp'(V)_z\Phi_z^2 + 2\lambda_z\Phi(\Phi_t - s\Phi_z) + 2(p'(V)\Phi)_z(\Phi_t - s\Phi_z) + \{\dots\}_z \\
 (59) \quad & + 2\mu(\Phi_t - s\Phi_z)_z^2 = -2F_z(\Phi_t - s\Phi_z).
 \end{aligned}$$

The combination (58) + (59) $\times 2\tau$ yields

$$\begin{aligned}
 & \{E_1(\Phi, (\Phi_t - s\Phi_z)) + E_3(\Phi_z)\}_t + E_2(\Phi_z, (\Phi_t - s\Phi_z)) + H(\Phi) + \{\dots\}_z \\
 (60) \quad & + 4\mu\tau(\Phi_t - s\Phi_z)_z^2 + 2\mu\Phi_z(\Phi_t - s\Phi_z)_z = -F_z\{2\Phi + 4\tau(\Phi_t - s\Phi_z)\},
 \end{aligned}$$

where

$$\begin{aligned}
 (61) \quad & H(\Phi) = \lambda_z\Phi^2 + 4\tau\lambda_z\Phi(\Phi_t - s\Phi_z) - 2p'(V)_z\Phi\Phi_z + 4\tau(p'(V)\Phi)_z(\Phi_t - s\Phi_z), \\
 & E_1(\Phi, (\Phi_t - s\Phi_z)) = 2\tau(\Phi_t - s\Phi_z)^2 + 2\Phi(\Phi_t - s\Phi_z) + \frac{1}{\tau}\Phi^2, \\
 & E_3(\Phi_z) = -2p'(V)\Phi_z^2, \\
 & E_2(\Phi_z, (\Phi_t - s\Phi_z)) = 2(\Phi_t - s\Phi_z)^2 - 4u'_e(V)\Phi_z(\Phi_t - s\Phi_z) \\
 & \quad - 2(p'(V) + s\tau p'(V)_z)\Phi_z^2.
 \end{aligned}$$

Integrating (60) with respect to t and z and taking $\frac{\mu}{\tau}$ and $|v_+ - v_-|$ suitably small, we have the following estimate:

$$\begin{aligned}
 (62) \quad & \|\Phi(t)\|_1^2 + \|\Phi_t(t)\|^2 + \int_0^t \|(\Phi_t, \Phi_z)(s)\|^2 ds + C_1\mu \int_0^t \|(\Phi_t - s\Phi_z)_z(s)\|^2 ds \\
 & \leq C \left\{ \|\Phi_0\|_1^2 + \|\Phi_1\|^2 + \int_0^t \int |H(\Phi)| dz ds + \int_0^t \int_R |F_z|(|\Phi| + |(\Phi_t, \Phi_z)|) dz ds \right\},
 \end{aligned}$$

where $\Phi_0 = \phi'_0$ and $\Phi_1 = \phi'_1$, and C and C_1 are positive constants.

Using the estimate (43), we obtain

$$\begin{aligned}
 & \int_0^t \int_R |H(\Phi)| dz ds \leq \frac{1}{2} \int_0^t \|(\Phi_t, \Phi_z)(s)\|^2 ds + C \int_0^t \int_R \Phi^2 dz ds \\
 & \leq \frac{1}{2} \int_0^t \|(\Phi_t, \Phi_z)(s)\|^2 ds \\
 (63) \quad & + C \left\{ \|\phi_0\|_{H^1_w}^2 + \|\phi_1\|_{L^2_w}^2 + \int_0^t \int w|F|(|\phi| + |(\phi_t, \phi_z)|) dz ds \right\}.
 \end{aligned}$$

Substituting (63) into (62) and replacing Φ by $\partial_z\phi$, we have the following lemma.

LEMMA 4.2. *Under the conditions of Theorem 1.1, there are positive constants C and C_1 such that if subcharacteristic condition (4) is satisfied for all $v \in [v_+, v_-]$, then any solution $\phi \in X(0, T)$ of problem (34), (37) satisfying*

$$\begin{aligned}
 & \|\partial_z\phi\|_1^2 + \|\partial_z\phi_t\|^2 + \frac{1}{2} \int_0^t \|(\partial_z\phi_t, \partial_z\phi_z)(s)\|^2 ds + C_1\mu \int_0^t \|(\phi_t - s\phi_z)_{zz}(s)\|^2 ds \\
 & \leq C \left\{ \|\phi_0\|_2^2 + \|\phi_1\|_1^2 + \int_0^t \int |F_z|(|\partial_z\phi| + |(\partial_z\phi_t, \partial_z\phi_z)|) dz ds \right. \\
 (64) \quad & \left. + \|\phi_0\|_{H^1_w}^2 + \|\phi_1\|_{L^2_w}^2 + \int_0^t \int w|F|(|\phi| + |(\phi_t, \phi_z)|) dz ds \right\}
 \end{aligned}$$

holds for $t \in [0, T]$.

Next we calculate the equality

$$\partial_z^2 \phi \partial_z^2 L(\phi) + 2\partial_z^2(\phi_t - s\phi_z) \partial_z^2 L(\phi) = -F_{zz}(\partial_z^2 \phi + 2\partial_z^2(\phi_t - s\phi_z))$$

in the same way as in the proof of Lemma 4.2. Setting $\Psi := \partial_z^2 \phi$, we have

$$\partial_z^2 L(\phi) = L(\Psi) + (p'(V)\Phi)_{zzz} - (p'(V)\Psi_z)_z + (\lambda\Phi)_{zz} - \lambda\Psi_z.$$

A straightforward calculation gives

$$\begin{aligned} & \left[2\tau(\Psi_t - s\Psi_z)^2 - 2\tau p'(V)\Psi_z^2 + 2\tau\Psi(\Psi_t - s\Psi_z) + \frac{1}{\tau}\Psi^2 \right]_t + 2(\Psi_t - s\Psi_z)^2 \\ & - 4u'_e(V)\Psi_z(\Psi_t - s\Psi_z) + 2(-p'(V) - s\tau p'(V)_z)\Psi_z^2 + 8\tau\lambda_z\Psi(\Psi_t - s\Psi_z) \\ & + 3\tau\lambda_z\Psi^2 + 2\tau\lambda_{zz}\Psi\phi_z + 4\tau\lambda_{zz}\phi_z(\Psi_t - s\Psi_z) + 4\mu\tau(\Psi_t - s\Psi_z)_z^2 \\ (65) \quad & + 2\mu\Psi_z(\Psi_t - s\Psi_z)_z + \{\dots\}_z = -2F_{zz}[\Psi + 2\tau(\Psi_t - s\Psi_z)] + J, \end{aligned}$$

where J satisfies

$$\begin{aligned} |J| &= |[-(p'(V)\Phi)_{zzz} + (p'(V)\Psi_z)_z][\Psi + 2\tau(\Psi_t - s\Psi_z)]| \\ &\leq \frac{1}{3}|(\Psi_t, \Psi_z)|^2 + C|(\Phi, \Psi)|^2. \end{aligned}$$

Thus, noting $\Psi = \phi_{zz}$ and (10), we derive from (65) and $|v_+ - v_-| \ll 1$ that

$$\begin{aligned} & \|\partial_z^2 \phi(t)\|_1^2 + \|\partial_z^2 \phi_t\|_1^2 + \frac{2}{3} \int_0^t \|(\partial_z^2 \phi_t, \partial_z^2 \phi_z)(s)\|_1^2 ds - C \int_0^t (\|\partial_z^2 \phi\|_1^2 + \|\phi_z\|_1^2) ds \\ & + C_1\mu \int_0^t \|(\partial_z^2 \phi_t - s\partial_z^2 \phi_z)_{zz}(s)\|_1^2 ds \\ (66) \quad & \leq C \left\{ \|\phi_0\|_3^2 + \|\phi_1\|_2^2 + \left| \int_0^t \int F_{zz}(\partial_z^2 \phi + 2\partial_z^2(\phi_t - s\phi_z)) dz ds \right| \right\}, \end{aligned}$$

where C and C_1 are positive constants.

Combining successively estimates (43), (64), and (66), we have

$$\begin{aligned} & \|\phi(t)\|_3^2 + \|\phi_t(t)\|_2^2 + \|\phi(t)\|_{H_w^1}^2 + \|\phi_t(t)\|_{L_w^2}^2 + \int_0^t \int |\lambda_z|\phi^2 dz ds + \int_0^t \|(\phi_t, \phi_z)\|_2^2 ds \\ & + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds + C_1\mu \left(\int_0^t \|(\phi_t - s\phi_z)_z\|_{L_w^2}^2 ds + \int_0^t \|(\phi_t - s\phi_z)_{zz}\|_1^2 ds \right) \\ & \leq C \left\{ \|\phi_0\|_3^2 + \|\phi_1\|_2^2 + \|\phi_0\|_{H_w^1}^2 + \|\phi_1\|_{L_w^2}^2 \right. \\ & + \int_0^t \int (w|F|(|\phi| + |(\phi_t, \phi_z)|) + |F_z|(|\partial_z \phi| + |(\partial_z \phi_t, \partial_z \phi_z)|)) dz ds \\ (67) \quad & \left. + \left| \int_0^t \int F_{zz}(\partial_z^2 \phi + 2\partial_z^2(\phi_t - s\phi_z)) dz ds \right| \right\}, \end{aligned}$$

where

$$F = F_1 + F_2 = O(1)(\phi_z^2 + \phi_{zz}^2)$$

as defined in (35), (36). There are terms containing fourth order derivatives of ϕ in $F_{2,zz}$. The energy estimates for the derivatives of ϕ up to fourth order can be established; see [10].

Thus by virtue of (39), the integrals on the right-hand side of (67) are majored by

$$CN(t) \left(\int_0^t \|(\phi_t, \phi_z)\|_2^2 ds + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds \right);$$

then we have

$$\begin{aligned} & N^2(t) + \int_0^t \|(\phi_t, \phi_z)\|_2^2 ds + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds \\ & \leq CN^2(0) + CN(t) \left(\int_0^t \|(\phi_t, \phi_z)\|_2^2 ds + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds \right). \end{aligned}$$

Therefore, by assuming $N(T) \leq \frac{1}{2C} = \delta_2$, we obtain the desired estimate

$$N^2(t) + \int_0^t \|(\phi_t, \phi_z)\|_2^2 ds + \int_0^t \|(\phi_t, \phi_z)(s)\|_{L_w^2}^2 ds \leq CN^2(0) \quad \text{for } t \in [0, T].$$

By choosing $N(0) \leq \delta_1$ small, we arrive at $N(T) \leq \delta_2$ for any $T > 0$.

Thus the proof of Proposition 3.3 is completed. \square

In summary, we established the weighted energy estimates by assuming the smallness of $|v_+ - v_-|$, the compressibility of the traveling wave profile (14), the subcharacteristic conditions (4) and (9), and the fact that the diffusion coefficient is bounded by a constant multiple of the relaxation time (10).

Acknowledgment. The author would like to acknowledge Professor Doochul Kim for insightful conversations which motivated the current work.

REFERENCES

- [1] G.-Q. CHEN, C. D. LEVERMORE, AND T. P. LIU, *Hyperbolic conservation laws with stiff relaxation terms and entropy*, Comm. Pure Appl. Math., 47 (1994), pp. 787–830.
- [2] S. JIN AND Z. XIN, *The relaxing schemes for systems of conservation laws in arbitrary space dimensions*, Comm. Pure Appl. Math., 48 (1995), pp. 235–276.
- [3] W. L. JIN AND H. M. ZHANG, *The formation and structure of vehicle clusters in the Payne–Whitham traffic flow model*, Transportation Res. Part B, 37 (2003), pp. 207–223.
- [4] S. JIN AND J.-G. LIU, *Relaxation and diffusion enhanced dispersive waves*, Proc. R. Soc. Lond. Ser. A, 446 (1994), pp. 555–563.
- [5] J. P. KEENER, *Diffusion induced oscillatory insulin secretion*, Bull. Math. Biol., 63 (2001), pp. 625–641.
- [6] B. S. KERNER AND P. KONHÄUSER, *Structure and parameters of clusters in traffic flow*, Phys. Rev. E, 50 (1994), pp. 54–83.
- [7] D. A. KURTZE AND D. C. HONG, *Traffic jams, granular flow, and soliton selection*, Phys. Rev. E, 52 (1995), pp. 218–221.
- [8] C. LATTANZIO AND P. MARCATI, *The zero relaxation limit for the hydrodynamic Whitham traffic flow model*, J. Differential Equations, 141 (1997), pp. 150–178.
- [9] H. Y. LEE, H.-W. LEE, AND D. KIM, *Steady-state solutions of hydrodynamic traffic models*, Phys. Rev. E, 69 (2004), 016118.
- [10] T. LI AND H. LIU, *Stability of a traffic flow model with nonconvex relaxation*, Commun. Math. Sci., 3 (2005), pp. 101–118.
- [11] T. LI AND H. LIU, *Critical thresholds in a relaxation model for traffic flows*, Indiana Univ. Math. J., 57 (2008), pp. 1409–1431.

- [12] T. LI, *Rigorous asymptotic stability of a CJ detonation wave in the limit of small resolved heat release*, *Combust. Theory Model.*, 1 (1997), pp. 259–270.
- [13] T. LI, *Global solutions of nonconcave hyperbolic conservation laws with relaxation arising from traffic flow*, *J. Differential Equations*, 190 (2003), pp. 131–149.
- [14] T. LI, *Nonlinear dynamics of traffic jams*, *Phys. D*, 207 (2005), pp. 41–51.
- [15] T. P. LIU, *Hyperbolic conservation laws with relaxation*, *Comm. Math. Phys.*, 108 (1987), pp. 153–175.
- [16] H. LIU, J. WANG, AND T. YANG, *Stability of a relaxation model with a nonconvex flux*, *SIAM J. Math. Anal.*, 29 (1998), pp. 18–29.
- [17] M. MURAMATSU AND T. NAGATANI, *Soliton and kink jams in traffic flow with open boundaries*, *Phys. Rev. E*, 60 (1999), pp. 180–187.
- [18] C. MASCIA AND K. ZUMBRUN, *Stability of large-amplitude shock profiles of general relaxation systems*, *SIAM J. Math. Anal.*, 37 (2005), pp. 889–913.
- [19] A. MATSUMURA AND K. NISHIHARA, *Asymptotic stability of traveling waves for scalar viscous conservation laws with non-convex nonlinearity*, *Comm. Math. Phys.*, 165 (1994), pp. 83–96.
- [20] T. NAGATANI, *The physics of traffic jams*, *Rep. Progr. Phys.*, 65 (2002), pp. 1331–1386.
- [21] T. NISHIDA, *Nonlinear Hyperbolic Equations and Related Topics in Fluid Dynamics*, *Publ. Math. d’Orsay 78-02*, Département de Mathématique, Université de Paris-Sud, Orsay, France, 1978.
- [22] Z.-H. OU, S.-Q. DAI, P. ZHANG, AND L.-Y. DONG, *Nonlinear analysis in the Aw-Rascle anticipation model of traffic flow*, *SIAM J. Appl. Math.*, 67 (2007), pp. 605–618.
- [23] G. B. WHITHAM, *Linear and Nonlinear Waves*, Wiley, New York, 1974.
- [24] H. M. ZHANG, *Driver memory, traffic viscosity and a viscous vehicular traffic flow model*, *Transportation Res. Part B*, 37 (2003), pp. 27–41.

NONEXISTENCE OF SUPERSONIC TRAVELING WAVES FOR NONLINEAR SCHRÖDINGER EQUATIONS WITH NONZERO CONDITIONS AT INFINITY*

MIHAI MARIȘ†

Abstract. We prove that the nonexistence of supersonic finite-energy traveling waves for nonlinear Schrödinger equations with nonzero conditions at infinity is a general phenomenon which holds for a large class of equations. The same is true for sonic traveling waves in two dimensions. In higher dimensions we prove that sonic traveling waves, if they exist, must approach their limit at infinity in a very rigid way. In particular, we infer that there are no sonic traveling waves with finite energy and finite momentum.

Key words. nonlinear Schrödinger equation, nonzero conditions at infinity, traveling wave, integral identities, Gross–Pitaevskii equations and systems, cubic–quintic nonlinear Schrödinger equation

AMS subject classifications. 35Q51, 35Q55, 35Q40, 35B65, 35J15, 35J20, 35J50, 37K40, 37K05

DOI. 10.1137/070711189

1. Introduction. The aim of this paper is to study traveling wave solutions for nonlinear Schrödinger equations

$$(1.1) \quad i \frac{\partial \Phi}{\partial t} + \Delta \Phi + F(x, |\Phi|^2) \Phi = 0 \quad \text{in } \mathbf{R}^N,$$

where F is a real-valued function defined on $\mathbf{R}^N \times \mathbf{R}_+$, Φ is a complex-valued function on \mathbf{R}^N satisfying the “boundary condition” $|\Phi| \rightarrow r_0$ as $|x| \rightarrow \infty$, and r_0 is a positive constant verifying $\lim_{|x| \rightarrow \infty, s \rightarrow r_0^2} F(x, s) = 0$.

The above equation with the considered nonzero conditions at infinity arises in a large variety of physical problems, such as superconductivity, superfluidity in helium II, phase transitions, and Bose–Einstein condensates. Two important particular cases of (1.1) have been extensively studied both by physicists and by mathematicians: the Gross–Pitaevskii equation (where $F(x, s) = 1 - s$) and the so-called cubic–quintic Schrödinger equation (where $F(x, s) = -\alpha_1 + \alpha_3 s - \alpha_5 s^2$, $\alpha_1, \alpha_3, \alpha_5$ are positive, and $\frac{3}{16} < \frac{\alpha_1 \alpha_5}{\alpha_3^2} < \frac{1}{4}$).

Equation (1.1) has a Hamiltonian structure: denoting $V(x, s) = \int_s^{r_0^2} F(x, \tau) d\tau$, it is easy to see that, at least formally, the “energy”

$$(1.2) \quad E(\Phi) = \int_{\mathbf{R}^N} |\nabla \Phi|^2 dx + \int_{\mathbf{R}^N} V(x, |\Phi|^2) dx$$

is a conserved quantity. There is another important (vector) quantity associated to (1.1), namely, the momentum. It is given by

$$(1.3) \quad P(\Phi) = (P_1(\Phi), \dots, P_N(\Phi)), \text{ where } P_k(\Phi) = \int_{\mathbf{R}^N} \left(i \frac{\partial \Phi}{\partial x_k}, \Phi \right) dx = \int_{\mathbf{R}^N} \operatorname{Re} \left(i \frac{\partial \Phi}{\partial x_k} \bar{\Phi} \right) dx.$$

*Received by the editors December 17, 2007; accepted for publication (in revised form) June 17, 2008; published electronically October 17, 2008.

<http://www.siam.org/journals/sima/40-3/71118.html>

†Département de Mathématiques, UMR 6623, Université de Franche-Comté, 16, Route de Gray, 25030 Besançon Cedex, France (mihai.maris@univ-fcomte.fr).

Note that, in general, the momentum is not well-defined for any function Φ of finite energy. In the case where F does not depend on the variable x_k , the momentum with respect to the x_k -direction, P_k , is conserved by those solutions of (1.1) for which it can be well-defined.

It is worth noting that (1.1) can be put into a hydrodynamical form by using Madelung’s transformation $\Phi(x, t) = \sqrt{\rho(x, t)}e^{i\theta(x, t)}$ (which is singular when $\Phi = 0$). A straightforward computation shows that in the region where $\Phi \neq 0$, the functions $\rho = |\Phi|^2$ and θ satisfy the system

$$(1.4) \quad \rho_t + 2\operatorname{div}(\rho\nabla\theta) = 0,$$

$$(1.5) \quad \theta_t + |\nabla\theta|^2 - \frac{\Delta\rho}{2\rho} + \frac{|\nabla\rho|^2}{4\rho} - F(x, \rho) = 0.$$

Equation (1.4) and the derivatives with respect to x_1, \dots, x_N of (1.5) are, respectively, the equation of conservation of mass and Euler’s equations for a compressible inviscid fluid of density ρ and velocity $2\nabla\theta$.

Let us assume that F admits a partial derivative with respect to the last variable (in what follows, this derivative will be denoted by $\partial_{N+1}F$ or by $\frac{\partial F}{\partial s}$) and that $\lim_{|x| \rightarrow \infty, \rho \rightarrow r_0^2} \partial_{N+1}F(x, \rho) = -L$, where L is a positive constant. Taking the derivative with respect to t of (1.5) and substituting ρ_t from (1.4) we obtain

$$(1.6) \quad \theta_{tt} + 2\partial_{N+1}F(x, \rho)(\rho\Delta\theta + \nabla\rho \cdot \nabla\theta) + \frac{\partial}{\partial t} \left(|\nabla\theta|^2 - \frac{\Delta\rho}{2\rho} + \frac{|\nabla\rho|^2}{4\rho} \right) = 0.$$

For a small oscillatory motion (i.e., a sound wave), all nonlinear terms in (1.6), except $2\rho\Delta\theta$, may be neglected. In view of the behavior of ρ and $\partial_{N+1}F(x, \rho)$ for large $|x|$, we find that in a neighborhood of infinity, the velocity potential θ essentially obeys the wave equation $\theta_{tt} - 2r_0^2L\Delta\theta = 0$. It is well known that the solutions of the wave equation propagate with a finite speed; in the present situation, we infer that the velocity of sound waves at infinity is $r_0\sqrt{2L}$. In what follows we will always assume that $\partial_{N+1}F(x, \rho) \rightarrow -L$ as $|x| \rightarrow \infty$ and $\rho \rightarrow r_0^2$ (the convergence being in a sense to be defined) and we will denote by $v_s = r_0\sqrt{2L}$ the sound velocity at infinity.

For a fixed $y \in S^{N-1}$, a traveling wave for (1.1) moving with velocity c in direction y is a solution of the form $\Phi(x, t) = \psi(x - cty)$. Without loss of generality we will assume that $y = (1, 0, \dots, 0)$, i.e., traveling waves move in the x_1 -direction. The traveling wave profile satisfies the equation

$$(1.7) \quad -ic\frac{\partial\psi}{\partial x_1} + \Delta\psi + F(x, |\psi|^2)\psi = 0 \quad \text{in } \mathbf{R}^N.$$

In a series of papers, Grant, Jones, Putterman, and Roberts, among others, formally and numerically studied traveling waves for the Gross–Pitaevskii equation and related systems (see, e.g., [16], [19], [21], [22], [5], and the references therein). In particular, they conjectured that such solutions exist if and only if their speed c belongs to the interval $(-v_s, v_s)$. For the cubic-quintic nonlinear Schrödinger equation, the existence of subsonic traveling waves in one dimension has been proved in [1] and their stability has been studied in [2]. The nonexistence of such solutions for sonic and supersonic speeds has also been conjectured in any space dimension. In the case of the Gross–Pitaevskii equation, it has been shown in [17] that any traveling wave of finite energy and speed $c > v_s$ must be constant. It has also been proved in [18] that the same result is true if $N = 2$ and $c^2 = v_s^2$. The proofs in

[17], [18] strongly depend on the special algebraic structure of the nonlinearity in the Gross–Pitaevskii equation. In the present paper we show that the nonexistence of finite-energy traveling waves moving faster than the sound velocity is a general phenomenon which holds for a large class of equations and systems of the form (1.1). We also prove that there are no finite-energy sonic traveling waves in space dimension two. In higher dimensions we show that any finite-energy sonic traveling wave ψ must satisfy $|\psi|^2 - r_0^2 \in L^p(\mathbf{R}^N)$ for any $p > \frac{2N-1}{2N-3}$. On the other hand, if a sonic traveling wave satisfies $|\psi|^2 - r_0^2 \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$, then it must be constant.

This article is organized as follows: In the next section we prove that traveling waves, whenever they exist, are smooth functions. If their speed is supersonic (or sonic, provided they converge sufficiently fast at infinity), then they must satisfy a special integral identity. This will be proved in section 3. In section 4 we show how this identity implies, under general assumptions, the nonexistence of traveling waves with finite energy. We apply our results to the Gross–Pitaevskii equation, to the cubic–quintic Schrödinger equation, and to a Gross–Pitaevskii–Schrödinger system, which describes the motion of an uncharged impurity in a Bose condensate. In the last section we describe all supersonic and sonic traveling waves (with finite or infinite energy) for one-dimensional equations with nonlinearities independent on the space variable.

2. Basic properties of traveling waves. We keep the previous notation and consider the following set of assumptions:

- **(H1)** $F : \mathbf{R}^N \times [0, \infty) \rightarrow \mathbf{R}$ is a measurable function which has the following properties:
 - (a) for any $s \in [0, \infty)$, $F(\cdot, s)$ is measurable;
 - (b) for any $x \in \mathbf{R}^N$, $F(x, \cdot)$ is continuous;
 - (c) F is bounded on bounded subsets of $\mathbf{R}^N \times [0, \infty)$.
- **(H2)** There exist $\alpha > 0$, $C > 0$, and $r_* > 0$ such that for any $x \in \mathbf{R}^N$ and for any $s \geq r_*$ we have $F(x, s) \leq -Cs^\alpha$.
- **(H3)** $\lim_{|x| \rightarrow \infty} F(x, r_0^2) = 0$ and $F(\cdot, r_0^2) \in L^1(\mathbf{R}^N)$.
- **(H4)** F admits a partial derivative with respect to the last variable and $\partial_{N+1}F$ is bounded on bounded subsets of $\mathbf{R}^N \times [0, \infty)$. Moreover, $\lim_{|x| \rightarrow \infty} \partial_{N+1}F(x, r_0^2) = -L$, where $L > 0$ and $\partial_{N+1}F(\cdot, r_0^2) + L \in L^{p_0}(\mathbf{R}^N)$ for some $p_0 \in [1, 2]$.
- **(H5)** There are some positive constants R_0, η, M such that ∂_{N+1}^2F exists on $(\mathbf{R}^N \setminus \overline{B}(0, R_0)) \times (r_0^2 - \eta, r_0^2 + \eta)$ and

$$|\partial_{N+1}^2F(x, s)| \leq M \quad \text{for all } (x, s) \in (\mathbf{R}^N \setminus \overline{B}(0, R_0)) \times (r_0^2 - \eta, r_0^2 + \eta).$$

DEFINITION 2.1. A traveling wave (of speed c) for (1.1) is a function $\psi \in L^1_{loc}(\mathbf{R}^N)$ that satisfies (1.7) in $\mathcal{D}'(\mathbf{R}^N)$ together with the “boundary condition” $|\psi| \rightarrow r_0$ as $|x| \rightarrow \infty$.

In view of (1.2), we say that a traveling wave ψ has finite energy if $\nabla\psi \in L^2(\mathbf{R}^N)$ and $V(\cdot, |\psi|^2) \in L^1(\mathbf{R}^N)$.

We have the following result concerning the regularity of traveling waves.

PROPOSITION 2.2. Let ψ be a finite-energy traveling wave for (1.1).

- (i) Assume that $F : \mathbf{R}^N \times \mathbf{R}_+ \rightarrow \mathbf{R}$ is measurable and satisfies (H1a), (H1b), (H2), the function $x \mapsto \int_{r_0^2}^{r_*} F(x, \tau) d\tau$ belongs to $L^1_{loc}(\mathbf{R}^N)$ (where r_* is given by (H2)), and $F(\cdot, |\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. Then $\psi \in L^\infty(\mathbf{R}^N)$.

If, in addition, F satisfies (H1c), then $\psi \in W_{loc}^{2,p}(\mathbf{R}^N)$ for any $p \in [1, \infty)$. In particular, $\psi \in C^{1,\alpha}(\mathbf{R}^N)$ for any $\alpha \in [0, 1)$.

(ii) Suppose that $F \in C^k(\mathbf{R}^N \times [0, \infty))$ for some $k \in \mathbf{N}^*$, (H2) holds, and $F(\cdot, |\psi|^2)\psi \in L_{loc}^1(\mathbf{R}^N)$. Then $\psi \in W_{loc}^{k+2,p}(\mathbf{R}^N)$ for any $p \in [1, \infty)$. In particular, if F is C^∞ , then $\psi \in C^\infty(\mathbf{R}^N)$.

Proof. (i) The proof relies upon the ideas and methods developed by Farina in [13], [14]. By (H2) we have

$$\begin{aligned} V(x, s) &= - \int_{r_0^2}^s F(x, \tau) d\tau \geq - \int_{r_0^2}^{r_*^2} F(x, \tau) d\tau + \int_{r_*^2}^s C\tau^\alpha d\tau \\ &= - \int_{r_0^2}^{r_*^2} F(x, \tau) d\tau + \frac{C}{\alpha + 1} (s^{\alpha+1} - r_*^{\alpha+1}). \end{aligned}$$

Consequently, for any $s \geq r_*$ we get $s^{\alpha+1} \leq r_*^{\alpha+1} + \frac{\alpha+1}{C}(V(x, s) + \int_{r_0^2}^{r_*^2} F(x, \tau) d\tau)$, so that

$$|\psi|^{2\alpha+2}(x) \leq \max \left(r_*^{\alpha+1}, r_*^{\alpha+1} + \frac{\alpha + 1}{C} \left(V(x, |\psi|^2(x)) + \int_{r_0^2}^{r_*^2} F(x, \tau) d\tau \right) \right).$$

Since $V(\cdot, |\psi|^2)$ and $\int_{r_0^2}^{r_*^2} F(\cdot, \tau) d\tau$ belong to $L_{loc}^1(\mathbf{R}^N)$, we infer that $\psi \in L_{loc}^{2\alpha+2}(\mathbf{R}^N)$.

We will use a well-known inequality of Kato (see Lemma A, p. 138, in [23]):

If $u \in L_{loc}^1(\mathbf{R}^N)$ is a real-valued function and $\Delta u \in L_{loc}^1(\mathbf{R}^N)$, then

$$(2.1) \quad \Delta(u^+) \geq \text{sgn}^+(u)\Delta u \quad \text{in } \mathcal{D}'(\mathbf{R}^N).$$

Let $\varphi(x) = e^{-\frac{icx_1}{2}}\psi(x)$. Then $\varphi \in L_{loc}^{2\alpha+2}(\mathbf{R}^N) \subset L_{loc}^1(\mathbf{R}^N)$ and an easy computation shows that φ satisfies

$$(2.2) \quad \Delta\varphi + \left(F(x, |\varphi|^2) + \frac{c^2}{4} \right) \varphi = 0 \quad \text{in } \mathcal{D}'(\mathbf{R}^N).$$

It is clear that $F(\cdot, |\varphi|^2)\varphi \in L_{loc}^1(\mathbf{R}^N)$ (because $F(x, |\psi|^2)\psi \in L_{loc}^1(\mathbf{R}^N)$ by hypothesis) and it follows from (2.2) that $\Delta\varphi \in L_{loc}^1(\mathbf{R}^N)$. Choose $\tilde{r} \geq r_*$ and $C_1 > 0$ such that $Cs^{2\alpha} - \frac{c^2}{4} \geq C_1(s - \tilde{r})^{2\alpha}$ for any $s \geq \tilde{r}$. Denoting $\varphi_1 = \text{Re}(\varphi)$, $\varphi_2 = \text{Im}(\varphi)$ and using Kato's inequality for $\varphi_i - \tilde{r}$, $i = 1, 2$, then using (2.2) and (H2) we get

$$\begin{aligned} (2.3) \quad \Delta(\varphi_i - \tilde{r})^+ &\geq \text{sgn}^+(\varphi_i - \tilde{r})\Delta(\varphi_i - \tilde{r}) = \text{sgn}^+(\varphi_i - \tilde{r}) \left[- \left(F(x, |\varphi|^2) + \frac{c^2}{4} \right) \varphi_i \right] \\ &\geq \text{sgn}^+(\varphi_i - \tilde{r}) \left[C|\varphi|^{2\alpha} - \frac{c^2}{4} \right] \varphi_i \geq \text{sgn}^+(\varphi_i - \tilde{r}) \left[C|\varphi_i|^{2\alpha} - \frac{c^2}{4} \right] \varphi_i \\ &\geq C_1 \text{sgn}^+(\varphi_i - \tilde{r})(\varphi_i - \tilde{r})^{2\alpha+1} = C_1 [(\varphi_i - \tilde{r})^+]^{2\alpha+1}. \end{aligned}$$

Next we use the following result of Brézis (Lemma 2, p. 273, in [9]).

LEMMA 2.3 (see [9]). Let $p \in (1, \infty)$. Assume that $u \in L_{loc}^p(\mathbf{R}^N)$ satisfies

$$-\Delta u + |u|^{p-1}u \leq 0 \quad \text{in } \mathcal{D}'(\mathbf{R}^N).$$

Then $u \leq 0$ a.e. on \mathbf{R}^N .

It follows from (2.3) that the function $u_i = (C_1)^{\frac{1}{2\alpha}}(\varphi_i - \tilde{r})^+$ satisfies $-\Delta u_i + |u_i|^{2\alpha}u_i \leq 0$ in $\mathcal{D}'(\mathbf{R}^N)$. Since $u_i \in L_{loc}^{2\alpha+1}(\mathbf{R}^N)$, we may use Lemma 2.3 and thus get $u_i \leq 0$ a.e. in \mathbf{R}^N , that is, $\varphi_i \leq \tilde{r}$ a.e. in \mathbf{R}^N .

It is obvious that both φ and $-\varphi$ satisfy (2.2). Repeating the above argument for $-\varphi$, we infer that $-\varphi_i \leq \tilde{r}$ a.e. on \mathbf{R}^N . Therefore we have $|\varphi_i| \leq \tilde{r}$ a.e. on \mathbf{R}^N , $i = 1, 2$, which implies that $\varphi \in L^\infty(\mathbf{R}^N)$. Since $|\varphi| = |\psi|$, we have proved that $\psi \in L^\infty(\mathbf{R}^N)$.

Using (H1c) and (2.2) we infer that $\Delta\varphi \in L^\infty(B(x, 2R)) \subset L^p(B(x, 2R))$ for any $x \in \mathbf{R}^N$, $R > 0$, and $p \geq 1$. By standard elliptic estimates we obtain $\varphi \in W^{2,p}(B(x, R))$ for any $x \in \mathbf{R}^N$, $R > 0$, and $p \in (1, \infty)$. Thus $\psi = e^{\frac{icx_1}{2}}\varphi \in W_{loc}^{2,p}(\mathbf{R}^N)$ for any $p \in (1, \infty)$, and consequently ψ belongs to $C_{loc}^{1,\alpha}(\mathbf{R}^N)$ for any $\alpha \in [0, 1)$ by the Sobolev embedding theorem.

(ii) Assume $F \in C^1(\mathbf{R}^N \times [0, \infty))$. Differentiating (1.7) with respect to x_k we get

$$(2.4) \quad -ic\psi_{x_1x_k} + \Delta\psi_{x_k} + \frac{\partial F}{\partial x_k}(x, |\psi|^2)\psi + 2\partial_{N+1}F(x, |\psi|^2) \left(\psi \cdot \frac{\partial \psi}{\partial x_k} \right) \psi + F(x, |\psi|^2) \frac{\partial \psi}{\partial x_k} = 0$$

in $\mathcal{D}'(\mathbf{R}^N)$. Hence $\Delta\psi_{x_k} \in L_{loc}^p(\mathbf{R}^N)$ for $1 \leq p < \infty$. By standard elliptic regularity theory we get $\psi_{x_k} \in W_{loc}^{2,p}(\mathbf{R}^N)$ for $1 < p < \infty$, $1 \leq k \leq N$; therefore $\psi \in W_{loc}^{3,p}(\mathbf{R}^N)$ for $1 \leq p < \infty$. If $F \in C^k(\mathbf{R}^N \times [0, \infty))$, we may differentiate (2.4) further and repeat the above arguments. After an easy induction, we get $\psi \in W_{loc}^{k+2,p}(\mathbf{R}^N)$ for any $p \in (1, \infty)$. \square

LEMMA 2.4. Assume that (H1), (H3), (H4), (H5) hold and $u \in L_{loc}^4(\mathbf{R}^N, \mathbf{C})$ satisfies $|u(x)| \rightarrow r_0$ as $|x| \rightarrow \infty$ and $V(\cdot, |u|^2) \in L^1(\mathbf{R}^N)$.

Then $|u|^2 - r_0^2 \in L^2(\mathbf{R}^N)$.

Proof. Let R_0, η, M be as in (H5). From (H4) and the fact that $|u(x)| \rightarrow r_0$ as $|x| \rightarrow \infty$ it follows that there exists $R_1 > R_0$ such that

$$\partial_{N+1}F(x, r_0^2) < -\frac{L}{2} \quad \text{and} \quad |u(x)|^2 \in (r_0^2 - \eta, r_0^2 + \eta) \quad \text{for any } x \text{ satisfying } |x| \geq R_1.$$

For $(x, s) \in (\mathbf{R}^N \setminus B(0, R_1)) \times (r_0^2 - \eta, r_0^2 + \eta)$ we get, by Taylor's formula with respect to the $(N + 1)$ th variable,

$$V(x, s) = -(s - r_0^2)F(x, r_0^2) - \frac{1}{2}(s - r_0^2)^2\partial_{N+1}F(x, r_0^2) - \frac{1}{2} \int_{r_0^2}^s (s - \tau)^2 \partial_{N+1}^2 F(x, \tau) d\tau.$$

In particular, for $s = |u(x)|^2$ we obtain

$$(2.5) \quad -\frac{1}{2}(|u(x)|^2 - r_0^2)^2\partial_{N+1}F(x, r_0^2) = V(x, |u(x)|^2) + (|u(x)|^2 - r_0^2)F(x, r_0^2) + \frac{1}{2} \int_{r_0^2}^{|u(x)|^2} (|u(x)|^2 - \tau)^2 \partial_{N+1}^2 F(x, \tau) d\tau.$$

For $x \in \mathbf{R}^N \setminus B(0, R_1)$ we get by (H5)

$$\left| \int_{r_0^2}^{|u(x)|^2} (|u(x)|^2 - \tau)^2 \partial_{N+1}^2 F(x, \tau) d\tau \right| \leq M \left| \int_{r_0^2}^{|u(x)|^2} (|u(x)|^2 - \tau)^2 d\tau \right| = \frac{M}{3} \left| (|u(x)|^2 - r_0^2) \right|^3.$$

It is clear that there exists $R_2 \geq R_1$ such that $\frac{M}{3}||u(x)|^2 - r_0^2| \leq \frac{L}{4}$ on $\mathbf{R}^N \setminus B(0, R_2)$. Using (H4) and (2.5) we infer that

$$\begin{aligned} \frac{L}{4}(|u(x)|^2 - r_0^2)^2 &\leq -\frac{1}{2}(|u(x)|^2 - r_0^2)^2 \partial_{N+1} F(x, r_0^2) \\ &\leq V(x, |u(x)|^2) + (|u(x)|^2 - r_0^2)F(x, r_0^2) + \frac{1}{2} \cdot \frac{M}{3} \left| |u(x)|^2 - r_0^2 \right|^3 \\ &\leq V(x, |u(x)|^2) + (|u(x)|^2 - r_0^2)F(x, r_0^2) + \frac{L}{8} \left| |u(x)|^2 - r_0^2 \right|^2 \quad \text{on } \mathbf{R}^N \setminus B(0, R_2). \end{aligned}$$

Consequently

$$(2.6) \quad \frac{L}{8}(|u(x)|^2 - r_0^2)^2 \leq V(x, |u(x)|^2) + (|u(x)|^2 - r_0^2)F(x, r_0^2) \text{ on } \mathbf{R}^N \setminus B(0, R_2).$$

Since $F(\cdot, r_0^2) \in L^1(\mathbf{R}^N)$ by (H3), $V(\cdot, |u|^2) \in L^1(\mathbf{R}^N)$, and $||u(x)|^2 - r_0^2| \leq \frac{3L}{4M}$ on $\mathbf{R}^N \setminus B(0, R_2)$, using (2.6) we get $(|u|^2 - r_0^2)^2 \in L^1(\mathbf{R}^N \setminus B(0, R_2))$. It is obvious that $(|u|^2 - r_0^2)^2 \in L^1(B(0, R_2))$ because $u \in L^4_{loc}(\mathbf{R}^N)$. Hence $(|u|^2 - r_0^2)^2 \in L^1(\mathbf{R}^N)$ and Lemma 2.4 is proved. \square

PROPOSITION 2.5. Assume that (H1)–(H5) hold and let ψ be a finite-energy traveling wave for (1.1) (in the sense of Definition 2.1) such that $F(\cdot, |\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. Then the following hold:

(i) $\nabla\psi \in W^{1,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$.

(ii) Let $R_* \geq 0$ be such that $|\psi(x)| \geq \frac{r_0}{2}$ for $|x| \geq R_*$. There exists a real-valued function θ such that $\theta \in W^{2,p}_{loc}(\mathbf{R}^N \setminus \bar{B}(0, R_*))$ for any $p < \infty$, $\nabla\theta \in W^{1,p}(\mathbf{R}^N \setminus \bar{B}(0, R_*))$ for any $p \in [2, \infty)$, and

$$\psi(x) = |\psi(x)|e^{i\theta(x)} \quad \text{on } \mathbf{R}^N \setminus B(0, R_*).$$

Proof. (i) We already know by Proposition 2.2(i) and Lemma 2.4 that ψ is bounded, $\psi \in W^{2,p}_{loc}(\mathbf{R}^N)$ for any $p \in [1, \infty)$, and $|\psi|^2 - r_0^2 \in L^2(\mathbf{R}^N)$.

Let R_0, η, M be as in (H5). Choose $R_1 > R_0$ such that $|\psi|^2(x) \in (r_0^2 - \eta, r_0^2 + \eta)$ for $x \in \mathbf{R}^N \setminus B(0, R_1)$.

By using Taylor’s formula with respect to the last variable for the function F we get

$$(2.7) \quad F(x, s) = F(x, r_0^2) + (s - r_0^2)\partial_{N+1}F(x, r_0^2) + \int_{r_0^2}^s (s - \tau)\partial_{N+1}^2F(x, \tau) d\tau$$

if $(x, s) \in (\mathbf{R}^N \setminus \bar{B}(0, R_0)) \times (r_0^2 - \eta, r_0^2 + \eta)$; hence

(2.8)

$$\begin{aligned} F(x, |\psi|^2(x))\psi(x) &= F(x, r_0^2)\psi(x) + (|\psi|^2(x) - r_0^2)\partial_{N+1}F(x, r_0^2)\psi(x) \\ &\quad + \psi(x) \int_{r_0^2}^{|\psi|^2(x)} (|\psi|^2(x) - \tau)\partial_{N+1}^2F(x, \tau) d\tau \quad \text{for any } |x| \geq R_1. \end{aligned}$$

We analyze the three terms on the right-hand side of (2.8). Assumptions (H1) and (H3) imply $F(\cdot, r_0^2) \in L^1 \cap L^\infty(\mathbf{R}^N)$. Since $\psi \in L^\infty(\mathbf{R}^N)$, it follows that $F(\cdot, r_0^2)\psi \in L^1 \cap L^\infty(\mathbf{R}^N)$.

We may write $(|\psi|^2 - r_0^2)\partial_{N+1}F(\cdot, r_0^2)\psi = -L(|\psi|^2 - r_0^2)\psi + (|\psi|^2 - r_0^2)(L + \partial_{N+1}F(\cdot, r_0^2))\psi$. We know that $\psi \in L^\infty(\mathbf{R}^N)$, $|\psi|^2 - r_0^2 \in L^2 \cap L^\infty(\mathbf{R}^N)$ and by (H4)

we have $L + \partial_{N+1}F(\cdot, r_0^2) \in L^{p_0} \cap L^\infty(\mathbf{R}^N)$ for some $p_0 \in [1, 2]$, so we infer that $(|\psi|^2 - r_0^2)\partial_{N+1}F(\cdot, r_0^2)\psi \in L^2 \cap L^\infty(\mathbf{R}^N)$.

As in the proof of Lemma 2.4, for $x \in \mathbf{R}^N \setminus B(0, R_1)$ we have

$$(2.9) \quad \left| \int_{r_0^2}^{|\psi|^2(x)} (|\psi|^2(x) - \tau)\partial_{N+1}^2F(x, \tau) d\tau \right| \leq M \left| \int_{r_0^2}^{|\psi|^2(x)} |\psi|^2(x) - \tau d\tau \right| = \frac{M}{2} (|\psi|^2(x) - r_0^2)^2.$$

Consequently the function $x \mapsto \int_{r_0^2}^{|\psi|^2(x)} (|\psi|^2(x) - \tau)\partial_{N+1}^2F(x, \tau) d\tau$ belongs to $L^1 \cap L^\infty(\mathbf{R}^N \setminus B(0, R_1))$.

Summing up, we have proved that $F(\cdot, |\psi|^2)\psi \in L^2 \cap L^\infty(\mathbf{R}^N \setminus B(0, R_1))$. From (H1) and the fact that ψ is bounded on \mathbf{R}^N it follows that $F(\cdot, |\psi|^2)\psi$ is bounded on $B(0, R_1)$, and hence $F(\cdot, |\psi|^2)\psi \in L^2 \cap L^\infty(\mathbf{R}^N)$.

We have $\frac{\partial\psi}{\partial x_k} \in L^2(\mathbf{R}^N)$ because ψ has finite energy. Coming back to (1.7), we get

$$\Delta\psi = ic \frac{\partial\psi}{\partial x_1} - F(\cdot, |\psi|^2)\psi \in L^2(\mathbf{R}^N).$$

It is well known that $\Delta\psi \in L^p(\mathbf{R}^N)$ with $1 < p < \infty$ implies $\frac{\partial^2\psi}{\partial x_j \partial x_k} \in L^p(\mathbf{R}^N)$ for any $j, k \in \{1, \dots, N\}$ (this follows, e.g., from the fact that $\frac{\xi_j \xi_k}{|\xi|^2}$ is a Fourier multiplier on $L^p(\mathbf{R}^N)$ if $1 < p < \infty$; see Theorem 3, p. 96, in [27]). Therefore all second derivatives of ψ are in $L^2(\mathbf{R}^N)$, so that $\frac{\partial\psi}{\partial x_k} \in H^1(\mathbf{R}^N) = W^{1,2}(\mathbf{R}^N)$ for $k = 1, \dots, N$.

The rest of the proof is an easy bootstrap argument. Assume that $\nabla\psi \in W^{1,p}(\mathbf{R}^N)$ for some $p \geq 2$. In case $p < N$, it follows from the Sobolev embedding theorem that $\nabla\psi \in L^{p^*}(\mathbf{R}^N)$, where $\frac{1}{p^*} = \frac{1}{p} - \frac{1}{N}$. From (1.7) we have $\Delta\psi = ic \frac{\partial\psi}{\partial x_1} - F(\cdot, |\psi|^2)\psi \in L^{p^*}(\mathbf{R}^N)$ and infer as previously that $\nabla\psi \in W^{1,p^*}(\mathbf{R}^N)$. Repeating this argument if necessary, after a finite number of steps we get $\nabla\psi \in W^{1,q}(\mathbf{R}^N)$ for some $q \geq N$. Then by Sobolev embedding we get $\nabla\psi \in L^r(\mathbf{R}^N)$ for any $r \in [q, \infty)$. From (1.7) we obtain $\Delta\psi \in L^p(\mathbf{R}^N)$ for $p \in [2, \infty)$ and infer that $\nabla\psi \in W^{1,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$.

(ii) Take $R_* > 0$ such that $|\psi(x)| \geq \frac{r_0}{2}$ on $\mathbf{R}^N \setminus B(0, R_*)$ and denote $\tilde{\psi}(x) = \frac{\psi(x)}{|\psi(x)|}$. It is then standard to prove that $\tilde{\psi} \in W_{loc}^{2,p}(\mathbf{R}^N \setminus \bar{B}(0, R_*))$ for $p \in [1, \infty)$ and $\nabla\tilde{\psi} \in W^{1,p}(\mathbf{R}^N \setminus \bar{B}(0, R_*))$ for any $p \in [2, \infty)$ (see, e.g., Lemma C1, p. 66, in [11]).

Let us consider first the case $N \geq 3$. For $R_* \leq R_1 < R_2$, the domain $\Omega_{R_1, R_2} = B(0, R_2) \setminus \bar{B}(0, R_1)$ is simply connected in \mathbf{R}^N . It follows from Theorem 3, p. 38, in [11] that there exists a real-valued function $\theta_{R_1, R_2} \in W^{2,p}(\Omega_{R_1, R_2})$ ($1 < p < \infty$) such that $\tilde{\psi} = e^{i\theta_{R_1, R_2}}$ on Ω_{R_1, R_2} . If $R_* \leq R_1 < R_2, R_* \leq R_3 < R_4$, and $(R_1, R_2) \cap (R_3, R_4) \neq \emptyset$, then $\tilde{\psi} = e^{i\theta_{R_1, R_2}} = e^{i\theta_{R_3, R_4}}$ on $\Omega_{R_1, R_2} \cap \Omega_{R_3, R_4}$, and thus $\theta_{R_3, R_4} - \theta_{R_1, R_2} \in 2\pi\mathbf{Z}$ on $\Omega_{R_1, R_2} \cap \Omega_{R_3, R_4}$. Since functions in $W^{s,p}(\Omega_{R_1, R_2} \cap \Omega_{R_3, R_4})$ with values in \mathbf{Z} are constant when $sp \geq 1$ (see Theorem B1, p. 65, in [11]), there exists $k \in \mathbf{Z}$ such that $\theta_{R_3, R_4} - \theta_{R_1, R_2} = 2\pi k$ on $\Omega_{R_1, R_2} \cap \Omega_{R_3, R_4}$. Let $(R_n)_{n \geq 1}$ be an increasing sequence such that $R_* < R_1$ and $R_n \rightarrow \infty$. Let $k_n \in \mathbf{Z}$ be such that $\theta_{R_*, R_n} = \theta_{R_*, R_1} + 2\pi k_n$ on Ω_{R_*, R_1} . Define $\theta(x) = \theta_{R_*, R_n}(x) - 2\pi k_n$ for $x \in \Omega_{R_*, R_n}$. It is clear that θ is well-defined on $\mathbf{R}^N \setminus \bar{B}(0, R_*)$, $\tilde{\psi} = e^{i\theta}$, and $\theta \in W_{loc}^{2,p}(\mathbf{R}^N \setminus \bar{B}(0, R_*))$ for any $p \in [1, \infty)$.

Next we consider the case $N = 2$. Since ψ is C^1 and $|\psi| \geq \frac{r_0}{2}$ on $\mathbf{R}^2 \setminus \bar{B}(0, R_*)$, the topological degree $deg(\psi, \partial B(0, R))$ is well-defined for any $R \geq R_*$ and does not depend on R . It is well known that ψ admits a C^1 lifting θ (i.e., $\psi = |\psi|e^{i\theta}$)

on $\mathbf{R}^2 \setminus \overline{B}(0, R_*)$ if and only if $\text{deg}(\psi, \partial B(0, R)) = 0$ for $R \geq R_*$. Denoting by $\tau = (-\sin \zeta, \cos \zeta)$ the unit tangent vector at $\partial B(0, R)$ at a point $Re^{i\zeta}$, we get

$$(2.10) \quad \begin{aligned} |\text{deg}(\psi, \partial B(0, R))| &= \left| \frac{1}{2i\pi} \int_0^{2\pi} \frac{\frac{\partial}{\partial \zeta}(\psi(Re^{i\zeta}))}{\psi(Re^{i\zeta})} d\zeta \right| = \left| \frac{R}{2i\pi} \int_0^{2\pi} \frac{\frac{\partial \psi}{\partial \tau}(Re^{i\zeta})}{\psi(Re^{i\zeta})} d\zeta \right| \\ &\leq \frac{R}{2\pi} \int_0^{2\pi} \frac{2}{r_0} |\nabla \psi(Re^{i\zeta})| d\zeta \leq \frac{R}{\pi r_0} \sqrt{2\pi} \left(\int_0^{2\pi} |\nabla \psi(Re^{i\zeta})|^2 d\zeta \right)^{\frac{1}{2}}. \end{aligned}$$

On the other hand,

$$\int_{\mathbf{R}^2 \setminus \overline{B}(0, R_*)} |\nabla \psi(x)|^2 dx = \int_{R_*}^{\infty} R \int_0^{2\pi} |\nabla \psi(Re^{i\zeta})|^2 d\zeta dR.$$

We have $\int_{\mathbf{R}^2 \setminus \overline{B}(0, R_*)} |\nabla \psi(x)|^2 dx < \infty$ (because ψ has finite energy) and infer that there exists $R_1 > R_*$ such that $R_1 \int_0^{2\pi} |\nabla \psi(R_1 e^{i\zeta})|^2 d\zeta < \frac{\pi r_0^2}{8} \frac{1}{R_1}$. From (2.10) we get

$$|\text{deg}(\psi, \partial B(0, R_1))| < \frac{R_1}{\pi r_0} \sqrt{2\pi} \left(\frac{\pi r_0^2}{8} \frac{1}{R_1} \right)^{\frac{1}{2}} = \frac{1}{2}.$$

Since the topological degree is an integer, we have necessarily $\text{deg}(\psi, \partial B(0, R_1)) = 0$. Consequently $\text{deg}(\psi, \partial B(0, R)) = 0$ for any $R \geq R_*$ and ψ admits a C^1 lifting θ . In fact, $\theta \in W_{loc}^{2,p}(\mathbf{R}^2 \setminus \overline{B}(0, R_*))$ because $\psi \in W_{loc}^{2,p}(\mathbf{R}^2 \setminus \overline{B}(0, R_*))$ (see Theorem 3, p. 38, in [11]).

If $N = 1$, the existence of a lifting $\psi = |\psi|e^{i\theta}$ follows immediately from Theorem 1, p. 27, in [11].

Finally, it is easy to see that $|\frac{\partial \bar{\psi}}{\partial x_j}| = |\frac{\partial \theta}{\partial x_j}|$ and $|\frac{\partial^2 \bar{\psi}}{\partial x_j \partial x_k}|^2 = |\frac{\partial^2 \theta}{\partial x_j \partial x_k}|^2 + |\frac{\partial \theta}{\partial x_j}|^2 |\frac{\partial \theta}{\partial x_k}|^2 \geq |\frac{\partial^2 \theta}{\partial x_j \partial x_k}|^2$, and (i) implies $\nabla \theta \in W^{1,p}(\mathbf{R}^N \setminus \overline{B}(0, R_*))$ for any $p \in [2, \infty)$. \square

3. An integral identity. The main result of this section is given by the next theorem.

THEOREM 3.1. *Assume that (H1)–(H5) hold. Let $\psi = \psi_1 + i\psi_2$ be a finite-energy traveling wave for (1.1) such that $F(\cdot, |\psi|^2) \in L^1_{loc}(\mathbf{R}^N)$. Let R_* be sufficiently big, so that $|\psi| \geq \frac{r_0}{2}$ on $\mathbf{R}^N \setminus B(0, R_*)$, and let θ be the lifting given by Proposition 2.5 (ii). Let $\chi \in C^\infty(\mathbf{R}^N)$ be a cut-off function such that $\chi = 0$ on $B(0, 2R_*)$ and $\chi = 1$ on $\mathbf{R}^N \setminus B(0, 3R_*)$. Then the following hold:*

(i) *The functions $F(\cdot, |\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2)$ and $G_j = \psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} - r_0^2 \frac{\partial}{\partial x_j}(\chi\theta)$, $j = 1, \dots, N$, belong to $L^1 \cap L^\infty(\mathbf{R}^N)$. (We always extend $\chi\theta$ by zero on $\overline{B}(0, R_*)$.)*

(ii) *If $N \geq 2$ and $c^2 > v_s^2$, we have the identity*

$$(3.1) \quad \begin{aligned} &\int_{\mathbf{R}^N} |\nabla \psi|^2 - F(x, |\psi|^2)|\psi|^2 - \frac{v_s^2}{2}(|\psi|^2 - r_0^2) dx \\ &= c \left(1 - \frac{v_s^2}{c^2} \right) \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial}{\partial x_1}(\chi\theta) dx. \end{aligned}$$

(iii) *Identity (3.1) holds if $c^2 = v_s^2$ and either*

- $N = 2$ or
- $N \geq 3$, and we assume in addition that $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$.

Proof. (i) Let R_0, η, M be as in (H5) and take $R_1 > R_0$ such that $|\psi|^2(x) \in (r_0^2 - \eta, r_0^2 + \eta)$ for $x \in \mathbf{R}^N \setminus B(0, R_1)$. Using (2.7) and the fact that $v_s^2 = 2Lr_0^2$ we get

$$\begin{aligned}
 & F(x, |\psi|^2(x))|\psi|^2(x) + \frac{v_s^2}{2}(|\psi|^2(x) - r_0^2) = F(x, r_0^2)|\psi|^2(x) \\
 (3.2) \quad & + (|\psi|^2(x) - r_0^2)[\partial_{N+1}F(x, r_0^2) + L]|\psi|^2(x) - L(|\psi|^2(x) - r_0^2)^2 \\
 & + |\psi|^2(x) \int_{r_0^2}^{|\psi|^2(x)} (|\psi|^2(x) - \tau) \partial_{N+1}^2 F(x, \tau) d\tau \quad \text{for any } |x| \geq R_1.
 \end{aligned}$$

Since $\psi \in L^\infty(\mathbf{R}^N)$ by Proposition 2.2 (i) and $F(\cdot, r_0^2) \in L^1 \cap L^\infty(\mathbf{R}^N)$ by (H1) and (H3), we infer that $F(\cdot, r_0^2)|\psi|^2 \in L^1 \cap L^\infty(\mathbf{R}^N)$.

We have $\psi \in L^\infty(\mathbf{R}^N)$, $\partial_{N+1}F(\cdot, r_0^2) + L \in L^{p_0} \cap L^\infty(\mathbf{R}^N)$ by (H4) and $|\psi|^2 - r_0^2 \in L^2 \cap L^\infty(\mathbf{R}^N)$ by Lemma 2.4; hence $(|\psi|^2 - r_0^2)[\partial_{N+1}F(\cdot, r_0^2) + L]|\psi|^2 \in L^1 \cap L^\infty(\mathbf{R}^N)$.

From Proposition 2.2 (i), Lemma 2.4, and (2.9) it follows that the last two terms on the right-hand side of (3.2) are in $L^1 \cap L^\infty(\mathbf{R}^N \setminus \overline{B}(0, R_1))$. Hence $F(\cdot, |\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2) \in L^1 \cap L^\infty(\mathbf{R}^N \setminus \overline{B}(0, R_1))$. Clearly, the function $F(\cdot, |\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2)$ is bounded on $\overline{B}(0, R_1)$; therefore this function belongs to $L^1 \cap L^\infty(\mathbf{R}^N)$.

Since $\psi_1 = |\psi| \cos \theta$ and $\psi_2 = |\psi| \sin \theta$, a straightforward computation gives

$$(3.3) \quad \psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} = (\psi_1^2 + \psi_2^2) \frac{\partial \theta}{\partial x_j} \quad \text{on } \mathbf{R}^N \setminus \overline{B}(0, R_*).$$

Therefore

$$(3.4) \quad \psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} - r_0^2 \frac{\partial}{\partial x_j}(\chi \theta) = (|\psi|^2 - r_0^2) \frac{\partial \theta}{\partial x_j} \quad \text{on } \mathbf{R}^N \setminus \overline{B}(0, 3R_*).$$

From Lemma 2.4, Proposition 2.5 (ii), and the Sobolev embedding theorem we have $|\psi|^2 - r_0^2 \in L^2 \cap L^\infty(\mathbf{R}^N)$ and $\frac{\partial \theta}{\partial x_j} \in L^2 \cap L^\infty(\mathbf{R}^N \setminus \overline{B}(0, R_*))$, respectively. Identity (3.4) implies $G_j \in L^1 \cap L^\infty(\mathbf{R}^N \setminus \overline{B}(0, 3R_*))$. Since G_j is continuous on \mathbf{R}^N , we conclude that $G_j \in L^1 \cap L^\infty(\mathbf{R}^N)$.

(ii) Equation (1.7) is equivalent to the system

$$(3.5) \quad c \frac{\partial \psi_2}{\partial x_1} + \Delta \psi_1 + F(x, |\psi|^2) \psi_1 = 0 \quad \text{in } \mathcal{D}'(\mathbf{R}^N),$$

$$(3.6) \quad -c \frac{\partial \psi_1}{\partial x_1} + \Delta \psi_2 + F(x, |\psi|^2) \psi_2 = 0 \quad \text{in } \mathcal{D}'(\mathbf{R}^N).$$

In view of Proposition 2.2 (i), equalities (3.5) and (3.6) hold in $L^p_{loc}(\mathbf{R}^N)$ for $1 \leq p < \infty$. Multiplying (3.5) by ψ_2 and (3.6) by ψ_1 and then subtracting the resulting equalities gives us

$$(3.7) \quad \frac{c}{2} \frac{\partial}{\partial x_1} (|\psi|^2 - r_0^2) = \operatorname{div}(\psi_1 \nabla \psi_2 - \psi_2 \nabla \psi_1).$$

We multiply (3.5) by ψ_1 and (3.6) by ψ_2 , then add the corresponding equalities to obtain

$$(3.8) \quad |\nabla \psi_1|^2 + |\nabla \psi_2|^2 - F(x, |\psi|^2)|\psi|^2 - c \left(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \right) = \frac{1}{2} \Delta (|\psi|^2 - r_0^2).$$

From (3.7) and (3.8) we get

$$(3.9) \quad \frac{c}{2} \frac{\partial}{\partial x_1} (|\psi|^2 - r_0^2) = \operatorname{div}(\psi_1 \nabla \psi_2 - \psi_2 \nabla \psi_1 - r_0^2 \nabla(\chi\theta)) + r_0^2 \Delta(\chi\theta),$$

respectively,

$$(3.10) \quad \begin{aligned} \frac{1}{2} \Delta (|\psi|^2 - r_0^2) - \frac{v_s^2}{2} (|\psi|^2 - r_0^2) &= |\nabla \psi_1|^2 + |\nabla \psi_2|^2 - F(x, |\psi|^2) |\psi|^2 - \frac{v_s^2}{2} (|\psi|^2 - r_0^2) \\ &\quad - c \left(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial}{\partial x_1} (\chi\theta) \right) - cr_0^2 \frac{\partial}{\partial x_1} (\chi\theta). \end{aligned}$$

Since $\psi \in W_{loc}^{2,p}(\mathbf{R}^N)$, equalities (3.7)–(3.10) hold in $L_{loc}^p(\mathbf{R}^N)$ for $1 \leq p < \infty$. We denote

$$H = |\nabla \psi_1|^2 + |\nabla \psi_2|^2 - F(x, |\psi|^2) |\psi|^2 - \frac{v_s^2}{2} (|\psi|^2 - r_0^2) - c \left(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial}{\partial x_1} (\chi\theta) \right).$$

We take the derivative of (3.9) with respect to x_1 (in $\mathcal{D}'(\mathbf{R}^N)$), multiply it by c , and then take the Laplacian of (3.10) (in $\mathcal{D}'(\mathbf{R}^N)$). Summing up the resulting equalities, we obtain

$$(3.11) \quad \frac{1}{2} \left(\Delta^2 - v_s^2 \Delta + c^2 \frac{\partial^2}{\partial x_1^2} \right) (|\psi|^2 - r_0^2) = \Delta H + c \frac{\partial}{\partial x_1} (\operatorname{div}(G)) \quad \text{in } \mathcal{D}'(\mathbf{R}^N).$$

From (i) we have $H, G_1, \dots, G_N \in L^1 \cap L^\infty(\mathbf{R}^N)$, and we know from Lemma 2.4 that $|\psi|^2 - r_0^2 \in L^2 \cap L^\infty(\mathbf{R}^N)$. Therefore $H, G_1, \dots, G_N, |\psi|^2 - r_0^2 \in \mathcal{S}'(\mathbf{R}^N)$, and we infer that, in fact, equality (3.11) holds in $\mathcal{S}'(\mathbf{R}^N)$. Taking the Fourier transform of (3.11) we get

$$(3.12) \quad \frac{1}{2} (|\xi|^4 + v_s^2 |\xi|^2 - c^2 \xi_1^2) \mathcal{F}(|\psi|^2 - r_0^2) = -|\xi|^2 \widehat{H} - c \sum_{k=1}^N \xi_1 \xi_k \widehat{G}_k \quad \text{in } \mathcal{S}'(\mathbf{R}^N).$$

We have $\widehat{H}, \widehat{G}_k \in L^\infty \cap C^0(\mathbf{R}^N)$ because $H, G_k \in L^1(\mathbf{R}^N)$. Thus the right-hand side of (3.12) is a continuous function on \mathbf{R}^N . Since $|\psi|^2 - r_0^2 \in L^2(\mathbf{R}^N)$, we have $\mathcal{F}(|\psi|^2 - r_0^2) \in L^2(\mathbf{R}^N)$ and infer that the left-hand side of (3.12) belongs to $L_{loc}^2(\mathbf{R}^N)$ and (3.12) holds a.e. on \mathbf{R}^N .

We denote

$$\Gamma = \{ \xi \in \mathbf{R}^N \mid |\xi|^4 + v_s^2 |\xi|^2 - c^2 \xi_1^2 = 0 \}.$$

If $c^2 \leq v_s^2$, we have $\Gamma = \{0\}$. If $c^2 > v_s^2$, it is easy to see that Γ is a nontrivial submanifold of \mathbf{R}^N . In the latter case, we claim that

$$(3.13) \quad |\xi|^2 \widehat{H}(\xi) + c \sum_{k=1}^N \xi_1 \xi_k \widehat{G}_k(\xi) = 0 \quad \text{for any } \xi \in \Gamma.$$

To prove this claim, we argue by contradiction and suppose that there exists $\xi^0 \in \Gamma$ such that $|\xi^0|^2 \widehat{H}(\xi^0) + c \sum_{k=1}^N \xi_1^0 \xi_k^0 \widehat{G}_k(\xi^0) \neq 0$. By continuity, there exist

$m > 0$ and a neighborhood U of ξ_0 such that $|\xi|^2 \widehat{H} + c \sum_{k=1}^N \xi_1 \xi_k \widehat{G}_k| \geq m$ on U . From (3.12) we infer that

$$|\mathcal{F}(|\psi|^2 - r_0^2)(\xi)| \geq \frac{2m}{\left| |\xi|^4 + v_s^2 |\xi|^2 - c^2 \xi_1^2 \right|} \quad \text{a.e. on } U \setminus \Gamma.$$

Since 0 and $(\sqrt{c^2 - v_s^2}, 0, \dots, 0)$ are not isolated points of Γ , we may assume that $\xi^0 \neq 0$ and $\xi^0 \neq (\sqrt{c^2 - v_s^2}, 0, \dots, 0)$. A straightforward computation (details can be found in [17, p. 98] in the case $v_s^2 = 2$; the general case is similar) shows that

$$\int_{U \setminus \Gamma} \frac{1}{\left| |\xi|^4 + v_s^2 |\xi|^2 - c^2 \xi_1^2 \right|^2} d\xi = \infty;$$

consequently $\int_{U \setminus \Gamma} |\mathcal{F}(|\psi|^2 - r_0^2)(\xi)|^2 d\xi = \infty$. But this is in contradiction with $\mathcal{F}(|\psi|^2 - r_0^2) \in L^2(\mathbf{R}^N)$ and the claim is proved.

It is not difficult to see that $\Gamma = \{(\xi_1, \xi') \in \mathbf{R} \times \mathbf{R}^{N-1} \mid |\xi'|^2 = \frac{1}{2}(-v_s^2 - 2\xi_1^2 + \sqrt{v_s^4 + 4c^2 \xi_1^2})\}$. Let $f(t) = \sqrt{\frac{1}{2}(-v_s^2 - 2t^2 + \sqrt{v_s^4 + 4c^2 t^2})}$. The function f is well-defined for $t \in [-\sqrt{c^2 - v_s^2}, \sqrt{c^2 - v_s^2}]$, $f(0) = 0$, and $\lim_{t \rightarrow 0} \frac{f^2(t)}{t^2} = -1 + \frac{c^2}{v_s^2}$. Fix $j \in \{2, \dots, N\}$. For $t \in (0, \sqrt{c^2 - v_s^2}]$, let $\xi(t) = (t, 0, \dots, 0, f(t), 0, \dots, 0)$ and $\tilde{\xi}(t) = (t, 0, \dots, 0, -f(t), 0, \dots, 0)$, where $f(t)$, respectively, $-f(t)$, stand at the j th place. It is obvious that $\xi(t), \tilde{\xi}(t) \in \Gamma$. From (3.13) we obtain, respectively,

$$(3.14) \quad (t^2 + f^2(t))\widehat{H}(\xi(t)) + ct^2 \widehat{G}_1(\xi(t)) + ct f(t) \widehat{G}_j(\xi(t)) = 0,$$

$$(3.15) \quad (t^2 + f^2(t))\widehat{H}(\tilde{\xi}(t)) + ct^2 \widehat{G}_1(\tilde{\xi}(t)) - ct f(t) \widehat{G}_j(\tilde{\xi}(t)) = 0.$$

We multiply (3.14) and (3.15) by $\frac{1}{t^2}$, then pass to the limit as $t \downarrow 0$ to obtain, respectively,

$$(3.16) \quad \frac{c^2}{v_s^2} \widehat{H}(0) + c \widehat{G}_1(0) + c \sqrt{-1 + \frac{c^2}{v_s^2}} \widehat{G}_j(0) = 0,$$

$$(3.17) \quad \frac{c^2}{v_s^2} \widehat{H}(0) + c \widehat{G}_1(0) - c \sqrt{-1 + \frac{c^2}{v_s^2}} \widehat{G}_j(0) = 0.$$

From (3.16) and (3.17) we infer that $\frac{c^2}{v_s^2} \widehat{H}(0) + c \widehat{G}_1(0) = 0$ and $\widehat{G}_j(0) = 0$, that is, $\int_{\mathbf{R}^N} H(x) + \frac{v_s^2}{c} G_1(x) dx = 0$ and $\int_{\mathbf{R}^N} G_j(x) dx = 0$. The first of these integral identities is exactly (3.1) and the latter can be written as

$$(3.18) \quad \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} - r_0^2 \frac{\partial}{\partial x_j} (\chi \theta) dx = 0 \quad \text{for } j = 2, \dots, N.$$

(iii) Assume that $c^2 = v_s^2$. Then (3.1) is equivalent to $\widehat{H}(0) + c \widehat{G}_1(0) = 0$. Denoting $\xi = (\xi_1, \xi')$, where $\xi' = (\xi_2, \dots, \xi_N)$, identity (3.12) implies

$$(3.19) \quad \begin{aligned} \mathcal{F}(|\psi|^2 - r_0^2)(\xi) &= -2 \frac{\xi_1^2}{|\xi|^4 + c^2 |\xi'|^2} (\widehat{H}(\xi) + c \widehat{G}_1(\xi)) \\ &- 2c \sum_{k=2}^N \frac{\xi_1 \xi_k}{|\xi|^4 + c^2 |\xi'|^2} \widehat{G}_k(\xi) - 2 \frac{|\xi'|^2}{|\xi|^4 + c^2 |\xi'|^2} \widehat{H}(\xi) \quad \text{a.e. } \xi \in \mathbf{R}^N. \end{aligned}$$

For $\varepsilon \in (0, 1]$, we denote $\Omega_\varepsilon = \{(\xi_1, \xi') \in \mathbf{R} \times \mathbf{R}^{N-1} \mid \xi_1 \in [0, \varepsilon], 0 \leq |\xi'| \leq \xi_1\}$. We will use the following lemma.

LEMMA 3.2. *Let $N \geq 2$ and $k \in \{2, \dots, N\}$.*

- (i) *The function $\xi \mapsto \frac{\xi_1^2}{\xi_1^4 + c^2|\xi'|^2}$ belongs to $L^p(\Omega_\varepsilon)$ if and only if $p < N - \frac{1}{2}$.*
- (ii) *The function $\xi \mapsto \frac{\xi_1 \xi_k}{\xi_1^4 + c^2|\xi'|^2}$ belongs to $L^p(\Omega_\varepsilon)$ for any $p \in [1, 2N - 1)$.*

Proof of Lemma 3.2. (i) Using Fubini's theorem for positive functions, then passing to spherical coordinates in \mathbf{R}^{N-1} and making the change of variables $r = \xi_1^2 t$ we get

$$\begin{aligned}
 (3.20) \quad & \int_{\Omega_\varepsilon} \left(\frac{\xi_1^2}{\xi_1^4 + c^2|\xi'|^2} \right)^p d\xi = \int_0^\varepsilon \xi_1^{2p} \int_{\{|\xi'| \leq \xi_1\}} \frac{1}{(\xi_1^4 + c^2|\xi'|^2)^p} d\xi' d\xi_1 \\
 & = \int_0^\varepsilon \xi_1^{2p} |S^{N-2}| \int_0^{\xi_1} \frac{r^{N-2}}{(\xi_1^4 + c^2 r^2)^p} dr d\xi_1 \\
 & = |S^{N-2}| \int_0^\varepsilon \xi_1^{2p} \int_0^{\frac{1}{\xi_1}} \frac{(\xi_1^2 t)^{N-2}}{(\xi_1^4 + c^2 \xi_1^4 t^2)^p} \xi_1^2 dt d\xi_1 \quad (\text{change of variables } r = \xi_1^2 t) \\
 & = |S^{N-2}| \int_0^\varepsilon \xi_1^{2(N-1-p)} \int_0^{\frac{1}{\xi_1}} \frac{t^{N-2}}{(1 + c^2 t^2)^p} dt d\xi_1.
 \end{aligned}$$

Assume that $p < N - \frac{1}{2}$. Obviously $\frac{t^{N-2}}{(1+c^2t^2)^p} \leq 1$ for $t \in [0, 1]$ and $\frac{t^2}{1+c^2t^2} \leq \frac{1}{c^2}$, and thus we have

$$\int_0^{\frac{1}{\xi_1}} \frac{t^{N-2}}{(1 + c^2 t^2)^p} dt \leq 1 + \frac{1}{c^{2p}} \int_1^{\frac{1}{\xi_1}} t^{N-2p-2} dt = \begin{cases} C_1 + \frac{C_2}{\xi_1^{N-2p-1}} & \text{if } p \neq \frac{N-1}{2}, \\ C_3 + C_4 \ln \xi_1 & \text{if } p = \frac{N-1}{2}, \end{cases}$$

where C_j are some positive constants. This estimate implies that the right-hand side of (3.20) is finite if $p < N - \frac{1}{2}$.

If $p \geq N - \frac{1}{2}$, denote $c_p = \int_0^1 \frac{t^{N-2}}{(1+c^2t^2)^p} dt > 0$. Since $\frac{1}{\xi_1} > 1$ for $\xi_1 \in (0, \varepsilon)$, the right-hand side of (3.20) is greater than $|S^{N-2}| c_p \int_0^\varepsilon \xi_1^{2(N-1-p)} d\xi_1 = \infty$.

(ii) Proceeding as above, we have

$$\begin{aligned}
 (3.21) \quad & \int_{\Omega_\varepsilon} \left| \frac{\xi_1 \xi_k}{\xi_1^4 + c^2|\xi'|^2} \right|^p d\xi \leq \int_{\Omega_\varepsilon} \frac{\xi_1^p |\xi'|^p}{(\xi_1^4 + c^2|\xi'|^2)^p} d\xi = \int_0^\varepsilon \xi_1^p |S^{N-2}| \int_0^{\xi_1} \frac{r^{p+N-2}}{(\xi_1^4 + c^2 r^2)^p} dr d\xi_1 \\
 & = |S^{N-2}| \int_0^\varepsilon \xi_1^p \int_0^{\frac{1}{\xi_1}} \frac{(\xi_1^2 t)^{p+N-2}}{(\xi_1^4 + c^2 \xi_1^4 t^2)^p} \xi_1^2 dt d\xi_1 \quad (\text{change of variables } r = \xi_1^2 t) \\
 & = |S^{N-2}| \int_0^\varepsilon \xi_1^{2N-p-2} \int_0^{\frac{1}{\xi_1}} \frac{t^{p+N-2}}{(1 + c^2 t^2)^p} dt d\xi_1.
 \end{aligned}$$

As previously,

$$\int_0^{\frac{1}{\xi_1}} \frac{t^{p+N-2}}{(1 + c^2 t^2)^p} dt < \frac{1}{c^{2p}} \int_0^{\frac{1}{\xi_1}} t^{N-p-2} dt = \frac{1}{c^{2p(N-p-1)}} \frac{1}{\xi_1^{N-p-1}} \quad \text{if } N-p-1 > 0.$$

Therefore in the case $p < N - 1$, the right-hand side of (3.21) is less than $C \int_0^\varepsilon \xi_1^{N-1} d\xi_1 < \infty$. If $p > N - 1$, the integral $\int_0^\infty \frac{t^{p+N-2}}{(1+c^2t^2)^p} dt$ converges. Let a_p be its value. If

$N - 1 < p < 2N - 1$, we get $\int_{\Omega_\varepsilon} \left| \frac{\xi_1 \xi_k}{\xi_1^4 + c^2 |\xi'|^2} \right|^p d\xi \leq |S^{N-2}| a_p \int_0^\varepsilon \xi_1^{2N-2-p} d\xi_1 < \infty$ by (3.21). \square

Remark. It can be proved that the function $\xi \mapsto \frac{\xi_1 \xi_k}{\xi_1^4 + c^2 |\xi'|^2}$ does not belong to $L^p(\Omega_\varepsilon)$ if $p \geq 2N - 1$, but we will not make use of this fact here.

Now we come back to the proof of Theorem 3.1. All we have to do is to show that $\widehat{H}(0) + c\widehat{G}_1(0) = 0$. We argue by contradiction and assume that $\widehat{H}(0) + c\widehat{G}_1(0) \neq 0$. Since the functions \widehat{H} and \widehat{G}_j are continuous, there exists $\varepsilon \in (0, 1)$ such that $|\widehat{H}(\xi) + c\widehat{G}_1(\xi)| \geq \frac{1}{2} |\widehat{H}(0) + c\widehat{G}_1(0)|$ for any $\xi \in \Omega_\varepsilon$. Taking a smaller ε if necessary, we may also assume that $|\xi|^4 + c^2 |\xi'|^2 \leq 2(\xi_1^4 + c^2 |\xi'|^2)$ for any $\xi \in \Omega_\varepsilon$. By (3.19) we have

$$(3.22) \quad \begin{aligned} & \frac{1}{2} \frac{\xi_1^2}{\xi_1^4 + c^2 |\xi'|^2} |\widehat{H}(0) + c\widehat{G}_1(0)| \leq 2 \frac{\xi_1^2}{|\xi|^4 + c^2 |\xi'|^2} |\widehat{H}(\xi) + c\widehat{G}_1(\xi)| \\ & \leq |\mathcal{F}(|\psi|^2 - r_0^2)(\xi)| + 2|c| \sum_{k=2}^N \frac{|\xi_1 \xi_k|}{\xi_1^4 + c^2 |\xi'|^2} |\widehat{G}_k(\xi)| + 2 \frac{|\xi'|^2}{|\xi|^4 + c^2 |\xi'|^2} |\widehat{H}(\xi)| \quad \text{a.e. on } \Omega_\varepsilon. \end{aligned}$$

Consider first the case $N = 2$. We know that $\mathcal{F}(|\psi|^2 - r_0^2) \in L^2(\mathbf{R}^2)$, and consequently $\mathcal{F}(|\psi|^2 - r_0^2) \in L^p(\Omega_\varepsilon)$ for any $p \in [1, 2]$. Since \widehat{G}_k are continuous and bounded, by Lemma 3.2 (ii) we infer that the functions $\xi \mapsto \frac{\xi_1 \xi_k}{\xi_1^4 + c^2 |\xi'|^2} \widehat{G}_k(\xi)$ belong to $L^p(\Omega_\varepsilon)$ for any $p \in [1, 3]$. It is obvious that $\frac{|\xi'|^2}{|\xi|^4 + c^2 |\xi'|^2} |\widehat{H}(\xi)| \leq \frac{1}{c^2} |\widehat{H}(\xi)|$ and \widehat{H} is continuous and bounded on \mathbf{R}^N . We conclude that the right-hand side of (3.22) belongs to $L^p(\Omega_\varepsilon)$ for any $p \in [1, 2]$. Then (3.22) implies that $\xi \mapsto \frac{\xi_1^2}{\xi_1^4 + c^2 |\xi'|^2}$ belongs to $L^2(\Omega_\varepsilon)$, which contradicts Lemma 3.2 (i). This contradiction proves that $\widehat{H}(0) + c\widehat{G}_1(0) = 0$.

Next we assume that $N \geq 3$ and $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$. Equation (3.8) can be written as

$$(3.23) \quad \begin{aligned} & -\frac{1}{2} \Delta(|\psi|^2 - r_0^2) + \frac{v_s^2}{2} (|\psi|^2 - r_0^2) \\ & = -|\nabla \psi_1|^2 - |\nabla \psi_2|^2 + F(x, |\psi|^2) |\psi|^2 + \frac{v_s^2}{2} (|\psi|^2 - r_0^2) + c \left(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \right). \end{aligned}$$

We have already proved that $F(\cdot, |\psi|^2) |\psi|^2 + \frac{v_s^2}{2} (|\psi|^2 - r_0^2) \in L^1 \cap L^\infty(\mathbf{R}^N)$. From Proposition 2.5 (i) we have $|\nabla \psi|^2 \in L^p(\mathbf{R}^N)$ for any $p \in [1, \infty]$. Using the assumption $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$, we infer that the right-hand side of (3.23) belongs to $L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$. By the Hausdorff-Young inequality, for any function $f \in L^p(\mathbf{R}^N)$ with $1 \leq p \leq 2$ we have $\mathcal{F}(f) \in L^{p'}(\mathbf{R}^N)$, where $\frac{1}{p} + \frac{1}{p'} = 1$ (see, e.g., Theorem 1.2.1, p. 6, in [4]). Passing to Fourier transforms in (3.23) we get

$$(3.24) \quad \begin{aligned} \mathcal{F}(|\psi|^2 - r_0^2)(\xi) &= \frac{2}{|\xi|^2 + v_s^2} \mathcal{F} \left[-|\nabla \psi|^2 + \left(F(\cdot, |\psi|^2) |\psi|^2 + \frac{v_s^2}{2} (|\psi|^2 - r_0^2) \right) \right. \\ & \quad \left. + c \left(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \right) \right] (\xi) \quad \text{a.e. } \xi \in \mathbf{R}^N. \end{aligned}$$

We obtain from (3.24) that $\mathcal{F}(|\psi|^2 - r_0^2) \in L^{N-\frac{1}{2}}(\mathbf{R}^N)$. Combined with Lemma 3.2 (ii) and the fact that \widehat{H} , \widehat{G}_j , and $\xi \mapsto \frac{|\xi'|^2}{|\xi|^4 + c^2 |\xi'|^2}$ are bounded, this implies that the last

expression in (3.22) is in $L^{N-\frac{1}{2}}(\Omega_\varepsilon)$. We infer that the function $\xi \mapsto \frac{\xi_1^2}{\xi_1^4+c^2|\xi'|^2}|\widehat{H}(0)+c\widehat{G}_1(0)|$ must be in $L^{N-\frac{1}{2}}(\Omega_\varepsilon)$ for any sufficiently small ε . If $\widehat{H}(0)+c\widehat{G}_1(0) \neq 0$, this contradicts Lemma 3.2 (i). Thus necessarily $\widehat{H}(0)+c\widehat{G}_1(0) = 0$ and the proof of Theorem 3.1 is complete. \square

It is an open problem whether any finite-energy traveling wave ψ of (1.1) moving with speed $c = \pm v_s$ satisfies $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$. Even for very particular cases of (1.1), such as the Gross–Pitaevskii equation, the answer to this question is not known. However, we have the following.

PROPOSITION 3.3. *Assume that (H1)–(H5) hold and let $\psi = \psi_1 + i\psi_2$ be a finite-energy traveling wave for (1.1) such that $F(\cdot, |\psi|^2) \in L^1_{loc}(\mathbf{R}^N)$. Let R_* be sufficiently big so that $|\psi| \geq \frac{r_0}{2}$ on $\mathbf{R}^N \setminus B(0, R_*)$, let θ be the lifting given by Proposition 2.5 (ii), and let $\chi \in C^\infty(\mathbf{R}^N)$ be a cut-off function as in Theorem 3.1. Then the following hold:*

- (i) *Let $p \in (1, \infty)$. The following assertions are equivalent:*
 - (a) $\nabla(\chi\theta) \in L^p(\mathbf{R}^N)$;
 - (b) $\psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} \in L^p(\mathbf{R}^N)$ for any $j \in \{1, \dots, N\}$;
 - (c) $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^p(\mathbf{R}^N)$;
 - (d) $|\psi|^2 - r_0^2 \in W^{2,p}(\mathbf{R}^N)$;
 - (e) $|\psi|^2 - r_0^2 \in L^p(\mathbf{R}^N)$.
- (ii) *If $N \geq 3$, there exists $\theta_0 \in \mathbf{R}$ such that $\chi\theta - \theta_0 \in W^{2,q}(\mathbf{R}^N)$ for any $q \in [\frac{2N}{N-2}, \infty)$.*

Moreover, if $c^2 = v_s^2$, we have the following:

- (iii) $|\psi|^2 - r_0^2 \in L^p(\mathbf{R}^N)$ and $\psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} \in L^p(\mathbf{R}^N)$ for any $p > \frac{2N-1}{2N-3}$ and $j \in \{1, \dots, N\}$.
- (iv) $\nabla(|\psi|^2 - r_0^2) \in L^p(\mathbf{R}^N)$ for any $p > \frac{2N-1}{2N-2}$.
- (v) $\partial_{j,k}^2(|\psi|^2 - r_0^2) \in L^p(\mathbf{R}^N)$ for any $p \in (1, \infty)$.

Proof. (i) Since $\psi \in L^\infty(\mathbf{R}^N)$ and (3.3) holds, the equivalence (a) \Leftrightarrow (b) is clear. It is also obvious that (b) \Rightarrow (c).

From the classical Marcinkiewicz theorem (see Theorem 3, p. 96, in [27]) it follows that the functions $\frac{1}{|\xi|^2+v_s^2}$, $\frac{\xi_j}{|\xi|^2+v_s^2}$, and $\frac{\xi_j \xi_k}{|\xi|^2+v_s^2}$ are L^p -multipliers for $1 < p < \infty$. Assume that $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^p(\mathbf{R}^N)$. Since $|\nabla\psi|^2 \in L^1 \cap L^\infty(\mathbf{R}^N)$ and $F(\cdot, |\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2) \in L^1 \cap L^\infty(\mathbf{R}^N)$ by Theorem 3.1 (i), we have $-\nabla\psi|^2 + (F(\cdot, |\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2)) + c(\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1}) \in L^p(\mathbf{R}^N)$ and infer from (3.24) that $|\psi|^2 - r_0^2 \in W^{2,p}(\mathbf{R}^N)$. Hence (c) \Rightarrow (d). It is obvious that (d) \Rightarrow (e).

It follows from Proposition 2.5 (ii) that $\partial_k(\chi\theta) \in \mathcal{S}'(\mathbf{R}^N)$. It is then clear that all terms appearing in (3.9) belong to $\mathcal{S}'(\mathbf{R}^N)$. We take the derivative of (3.9) with respect to x_k (in $\mathcal{S}'(\mathbf{R}^N)$), then take the Fourier transform of the resulting equality to obtain

$$r_0^2 \mathcal{F} \left(\frac{\partial}{\partial x_k}(\chi\theta) \right) = - \sum_{j=1}^N \frac{\xi_j \xi_k}{|\xi|^2} \widehat{G}_j + \frac{c}{2} \frac{\xi_1 \xi_k}{|\xi|^2} \mathcal{F}(|\psi|^2 - r_0^2)$$

or, equivalently,

$$(3.25) \quad r_0^2 \frac{\partial}{\partial x_k}(\chi\theta) = \sum_{j=1}^N R_j R_k(G_j) - \frac{c}{2} R_1 R_k(|\psi|^2 - r_0^2),$$

where R_j is the Riesz transform, $R_j \phi = \mathcal{F}^{-1}(i \frac{\xi_j}{|\xi|} \widehat{\phi})$. It is well known that the Riesz transform maps continuously $L^p(\mathbf{R}^N)$ into $L^p(\mathbf{R}^N)$ for $1 < p < \infty$ (see, e.g., Theorem 3, p. 96, and Example (iii), p. 95, in [27]). From Theorem 3.1 (i) we have $G_j \in L^1 \cap L^\infty(\mathbf{R}^N)$; therefore $R_j R_k(G_j) \in L^q(\mathbf{R}^N)$ for any $q \in (1, \infty)$. Assume that $|\psi|^2 - r_0^2 \in L^p(\mathbf{R}^N)$ for some $p \in (1, \infty)$. Then $R_1 R_k(|\psi|^2 - r_0^2) \in L^p(\mathbf{R}^N)$ and from (3.25) we infer that $\frac{\partial}{\partial x_k}(\chi\theta) \in L^p(\mathbf{R}^N)$ for any $k \in \{1, \dots, N\}$. Thus (e) \Rightarrow (a) and (i) is proved.

(ii) It is well known that for any function ϕ satisfying $\nabla\phi \in L^p(\mathbf{R}^N)$ with $p < N$, there exists a constant λ such that $\phi - \lambda \in L^{p^*}(\mathbf{R}^N)$, where $\frac{1}{p^*} = \frac{1}{p} - \frac{1}{N}$ (see Theorem 4.5.9 in [20] or Lemma 7 and Remark 4.2 in [15, pp. 774–775] for a different proof). From Proposition 2.5 (ii) we have $\nabla(\chi\theta) \in W^{1,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$. If $N \geq 3$, we infer that there exists $\theta_0 \in \mathbf{R}$ such that $\chi\theta - \theta_0 \in L^q(\mathbf{R}^N)$ for $q \in [\frac{2N}{N-2}, \infty)$. Therefore $\chi\theta - \theta_0 \in W^{2,q}(\mathbf{R}^N)$ for any $q \in [\frac{2N}{N-2}, \infty)$ and, in particular, $\chi\theta - \theta_0 \rightarrow 0$ as $|x| \rightarrow \infty$.

(iii) We will use the following result due to Lizorkin (see Theorem 8, p. 288, in [24]).

THEOREM 3.4 (see [24]). *Let $\beta \in [0, 1)$ and let $K \in L^\infty(\mathbf{R}^N) \cap C^N(\mathbf{R}^N \setminus \{0\})$. Assume that*

$$\left(\prod_{j=1}^N |\xi_j|^{k_j+\beta} \right) \partial_1^{k_1} \dots \partial_N^{k_N} K \in L^\infty(\mathbf{R}^N) \quad \text{for any } k_1, \dots, k_N \in \{0, 1\}.$$

Then K is a Fourier multiplier from $L^p(\mathbf{R}^N)$ to $L^{\frac{p}{1-\beta p}}(\mathbf{R}^N)$ for any $p \in (1, \frac{1}{\beta})$.

Let $K(\xi) = \frac{|\xi|^2}{|\xi|^4 + c^2 |\xi'|^2}$, where $\xi' = (\xi_2, \dots, \xi_N)$. A straightforward but tedious computation shows that K satisfies the assumptions of Lizorkin’s theorem for $\beta = \frac{1}{2N-1}$. From (3.19) we obtain

$$(3.26) \quad |\psi|^2 - r_0^2 = 2R_1^2 \left(\mathcal{F}^{-1} \left(K(\widehat{H} + c\widehat{G}_1) \right) \right) + 2c \sum_{j=2}^N R_1 R_j \left(\mathcal{F}^{-1}(K\widehat{G}_j) \right) + 2 \sum_{j=2}^N R_j^2 \left(\mathcal{F}^{-1}(K\widehat{H}) \right),$$

where R_j ’s denote Riesz transforms. Since $H, G_1, \dots, G_N \in L^1 \cap L^\infty(\mathbf{R}^N)$, by (3.26) and Lizorkin’s theorem we infer that $|\psi|^2 - r_0^2 \in L^p(\mathbf{R}^N)$ for any $p \in (\frac{2N-1}{2N-3}, \infty)$. The rest of (iii) follows from part (i), (b) \Leftrightarrow (e).

(iv) and (v) From (iii) and (i), (d) \Leftrightarrow (e) it follows immediately that $|\psi|^2 - r_0^2 \in W^{2,p}(\mathbf{R}^N)$ for any $p \in (\frac{2N-1}{2N-3}, \infty)$. Using (3.19) we obtain

$$(3.27) \quad \begin{aligned} \partial_{k\ell}^2 (|\psi|^2 - r_0^2) &= 2R_k R_\ell R_1^2 \left(\mathcal{F}^{-1} \left(|\xi|^2 K(\widehat{H} + c\widehat{G}_1) \right) \right) \\ &+ 2c \sum_{j=2}^N R_k R_\ell R_1 R_j \left(\mathcal{F}^{-1}(|\xi|^2 K\widehat{G}_j) \right) \\ &+ 2 \sum_{j=2}^N R_k R_\ell R_j^2 \left(\mathcal{F}^{-1}(|\xi|^2 K\widehat{H}) \right) \quad \text{in } \mathcal{S}'(\mathbf{R}^N). \end{aligned}$$

It can be proved by direct computation that the function $|\xi|^2 K$ satisfies the assumptions of Lizorkin’s theorem for $\beta = 0$. Consequently $|\xi|^2 K$ is an L^p -multiplier for

$1 < p < \infty$. Since $H, G_j \in L^1 \cap L^\infty(\mathbf{R}^N)$, it follows from (3.27) that $\partial_{k\ell}^2 (|\psi|^2 - r_0^2) \in L^p(\mathbf{R}^N)$ for $1 < p < \infty$.

By using the Gagliardo–Nirenberg inequality

$$\|\nabla\phi\|_{L^p}^2 \leq C\|\phi\|_{L^q}\|\nabla^2\phi\|_{L^r} \quad \text{if} \quad \frac{1}{p} = \frac{1}{2}\left(\frac{1}{q} + \frac{1}{r}\right),$$

we infer that $\nabla(|\psi|^2 - r_0^2) \in L^p(\mathbf{R}^N)$ for any $p > \frac{2N-1}{2N-2}$. \square

COROLLARY 3.5. *Under the assumptions of Theorem 3.1, assume that $N \geq 3$, $c^2 = v_s^2$, and the momentum of ψ with respect to the x_1 -direction is well-defined, that is, $\psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} \in L^1(\mathbf{R}^N)$. Then ψ satisfies (3.1).*

Proof. From Proposition 3.3 (iii) and (i) we have $\psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} \in L^p(\mathbf{R}^N)$ for $p \in (\frac{2N-1}{2N-3}, \infty)$. Then the assumption $\psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} \in L^1(\mathbf{R}^N)$ implies $\psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} \in L^p(\mathbf{R}^N)$ for any $p \in [1, \infty)$. Now the conclusion follows from Theorem 3.1 (iii). \square

4. Nonexistence results. In this section we show how Theorem 3.1 may be used to prove nonexistence of supersonic and sonic traveling waves with finite energy for some equations of type (1.1).

4.1. Equations invariant by translations. We consider the equation

$$(4.1) \quad i\frac{\partial\Phi}{\partial t} + \Delta\Phi + G(|\Phi|^2)\Phi = 0 \quad \text{in } \mathbf{R}^N.$$

We assume that the function $G : [0, \infty) \rightarrow \mathbf{R}$ satisfies the following assumptions:

- **(A1)** $G \in C^2([0, \infty), \mathbf{R})$ and there exists $r_0 > 0$ such that $G(r_0^2) = 0$ and $G'(r_0^2) < 0$.
- **(A2)** There exists $\alpha > 0$ such that $\limsup_{s \rightarrow \infty} \frac{G(s)}{s^\alpha} < 0$.

Obviously, (4.1) is of the form (1.1). As previously, we associate to (4.1) the “boundary condition” $|\Phi| \rightarrow r_0^2$ as $|x| \rightarrow \infty$. In this context, the sound velocity at infinity is $v_s = r_0\sqrt{-2G'(r_0^2)}$. The energy corresponding to (4.1) is $E(\Phi) = \int_{\mathbf{R}^N} |\nabla\Phi|^2 dx + \int_{\mathbf{R}^N} V(|\Phi|^2) dx$, where $V(s) = \int_s^{r_0^2} G(\tau) d\tau$. Let ψ be a finite-energy traveling wave for (4.1) (in the sense of Definition 2.1) moving with speed c . Then ψ satisfies the equation

$$(4.2) \quad -ic\frac{\partial\psi}{\partial x_1} + \Delta\psi + G(|\psi|^2)\psi = 0 \quad \text{in } \mathcal{D}'(\mathbf{R}^N), \quad |\psi| \rightarrow r_0 \quad \text{as } |x| \rightarrow \infty.$$

If G satisfies (A1)–(A2), it is easy to see that $F(x, s) := G(s)$ satisfies the assumptions (H1)–(H5) in section 2 (with $L = -G'(r_0^2)$). It is then clear that the conclusions of Propositions 2.2 and 2.5 and Theorem 3.1 (i) are valid for ψ . Moreover, we have the following.

PROPOSITION 4.1 (Pohozaev identities). *Let ψ be as above. Choose $R_* > 0$ such that $|\psi| \geq \frac{r_0}{2}$ on $\mathbf{R}^N \setminus B(0, R_*)$. Let θ be the lifting of $\frac{\psi}{|\psi|}$ on $\mathbf{R}^N \setminus B(0, R_*)$ (as given by Proposition 2.5 (ii)) and let χ be a cut-off function as in Theorem 3.1. The following identities hold:*

$$(4.3) \quad - \int_{\mathbf{R}^N} \left| \frac{\partial\psi}{\partial x_1} \right|^2 dx + \int_{\mathbf{R}^N} \sum_{j=2}^N \left| \frac{\partial\psi}{\partial x_j} \right|^2 dx + \int_{\mathbf{R}^N} V(|\psi|^2) dx = 0,$$

$$(4.4) \quad - \int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx + \int_{\mathbf{R}^N} \sum_{j=1, j \neq k}^N \left| \frac{\partial \psi}{\partial x_j} \right|^2 dx + \int_{\mathbf{R}^N} V(|\psi|^2) dx - c \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial}{\partial x_1} (\chi \theta) dx = 0 \quad \text{for } k = 2, \dots, N.$$

It is worth noting that Proposition 4.1 is valid for any speed $c \in \mathbf{R}$.

Proof. Since the arguments are rather classical, we only sketch the proof.

Formally, traveling waves are critical points of the functional $E_c = E + cP_1$, where E is the energy and P_1 is the momentum with respect to the x_1 -direction (see (1.3)). Identities (4.3) and (4.4) are simple consequences of the behavior of E_c with respect to dilations in \mathbf{R}^N . To be more precise, define $\psi_{k,t}(x) = \psi(x_1, \dots, x_{k-1}, tx_k, x_{k+1}, \dots, x_N)$ and $g_k(t) = E_c(\psi_{k,t})$. If ψ is a critical point of E_c , one would expect that $g'_k(1) = \frac{d}{dt}(E_c(\psi_{k,t}))|_{t=1} = 0$, and this is precisely (4.3) if $k = 1$, respectively, (4.4) if $k \geq 2$. However, this argument is not rigorous for at least two reasons. First, it is not clear what function space one should consider to define E_c (and this could not be a vector space because of the boundary conditions at infinity). Second, even if an appropriate function space is found, we do not know whether $\frac{d}{dt}(\psi_{k,t})|_{t=1} = x_k \frac{\partial \psi}{\partial x_k}$ belong to the tangent space at ψ of the considered function space.

The most convenient way to prove Pohozaev identities is to use a truncation argument. Fix a function $\eta \in C_c^\infty(\mathbf{R}^N)$ such that $\eta = 1$ on $B(0, 1)$ and $\eta = 0$ on $\mathbf{R}^N \setminus B(0, 2)$. For $n \geq 1$, define $\eta_n(x) = \eta(\frac{x}{n})$. We take the scalar product of (4.2) by $x_k \eta_n(x) \frac{\partial \psi}{\partial x_k}$ and integrate by parts the resulting equality. It is standard (see, e.g., Proposition 1, p. 320, in [3] or Lemma 2.4, p. 104, in [10]) to prove that

$$(4.5) \quad \lim_{n \rightarrow \infty} \int_{\mathbf{R}^N} \left(\Delta \psi, x_k \eta_n(x) \frac{\partial \psi}{\partial x_k} \right) dx = - \int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx + \frac{1}{2} \int_{\mathbf{R}^N} |\nabla \psi|^2 dx$$

and

$$(4.6) \quad \lim_{n \rightarrow \infty} \int_{\mathbf{R}^N} \left(G(|\psi|^2) \psi, x_k \eta_n(x) \frac{\partial \psi}{\partial x_k} \right) dx = \frac{1}{2} \int_{\mathbf{R}^N} V(|\psi|^2) dx.$$

It is obvious that $(ic \frac{\partial \psi}{\partial x_1}, \eta_n(x) x_1 \frac{\partial \psi}{\partial x_1}) = c \eta_n(x) x_1 (i \frac{\partial \psi}{\partial x_1}, \frac{\partial \psi}{\partial x_1}) = 0$. Thus taking the scalar product of (4.2) by $x_1 \eta_n(x) \frac{\partial \psi}{\partial x_1}$, integrating, and using (4.5) and (4.6) we get (4.3).

By (3.3) we have $(-i \frac{\partial \psi}{\partial x_j}, \psi) = \psi_1 \frac{\partial \psi_2}{\partial x_j} - \psi_2 \frac{\partial \psi_1}{\partial x_j} = |\psi|^2 \frac{\partial \theta}{\partial x_j}$ on $\mathbf{R}^N \setminus \bar{B}(0, R_*)$. Using the convention $\partial^\alpha (\chi \theta) = 0$, $(\partial^\alpha \chi) \theta = 0$ on $B(0, 2R_*)$, we have

$$(4.7) \quad \begin{aligned} \left(-i \frac{\partial \psi}{\partial x_j}, \psi \right) &= (1 - \chi) \left(-i \frac{\partial \psi}{\partial x_j}, \psi \right) + \chi |\psi|^2 \frac{\partial \theta}{\partial x_j} \\ &= (1 - \chi) \left(-i \frac{\partial \psi}{\partial x_j}, \psi \right) + |\psi|^2 \frac{\partial (\chi \theta)}{\partial x_j} - |\psi|^2 \theta \frac{\partial \chi}{\partial x_j} \quad \text{on } \mathbf{R}^N. \end{aligned}$$

Therefore we get for $k = 2, \dots, N$

$$\begin{aligned}
 & \int_{\mathbf{R}^N} \left(-ic \frac{\partial \psi}{\partial x_1}, x_k \eta_n(x) \frac{\partial \psi}{\partial x_k} \right) dx \\
 &= \frac{c}{2} \int_{\mathbf{R}^N} x_k \eta_n(x) \left[\frac{\partial}{\partial x_1} \left(-i\psi, \frac{\partial \psi}{\partial x_k} \right) + \frac{\partial}{\partial x_k} \left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) \right] dx \\
 &= -\frac{c}{2} \int_{\mathbf{R}^N} x_k \frac{\partial \eta_n}{\partial x_1}(x) \left(-i\psi, \frac{\partial \psi}{\partial x_k} \right) + \left(\eta_n(x) + x_k \frac{\partial \eta_n}{\partial x_k}(x) \right) \left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) dx \\
 (4.8) \quad &= \frac{c}{2} \int_{\mathbf{R}^N} x_k \frac{\partial \eta_n}{\partial x_1}(x) \left[(1 - \chi) \left(-i \frac{\partial \psi}{\partial x_k}, \psi \right) + |\psi|^2 \frac{\partial(\chi\theta)}{\partial x_k} - |\psi|^2 \theta \frac{\partial \chi}{\partial x_k} \right] dx \\
 &\quad - \frac{c}{2} \int_{\mathbf{R}^N} \eta_n(x) \left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) dx \\
 &\quad - \frac{c}{2} \int_{\mathbf{R}^N} x_k \frac{\partial \eta_n}{\partial x_k}(x) \left[(1 - \chi) \left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) + |\psi|^2 \frac{\partial(\chi\theta)}{\partial x_1} - |\psi|^2 \theta \frac{\partial \chi}{\partial x_1} \right] dx \\
 &= \frac{c}{2} \int_{\mathbf{R}^N} x_k |\psi|^2 \left(\frac{\partial \eta_n}{\partial x_1} \frac{\partial(\chi\theta)}{\partial x_k} - \frac{\partial \eta_n}{\partial x_k} \frac{\partial(\chi\theta)}{\partial x_1} \right) - \eta_n(x) \left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) dx \quad \text{if } n > 3R_*
 \end{aligned}$$

because $\text{supp}(1 - \chi) \subset \bar{B}(0, 3R_*)$ and $\text{supp} \nabla \eta_n \subset \bar{B}(0, 2n) \setminus B(0, n)$, and consequently $(1 - \chi) \frac{\partial \eta_n}{\partial x_j} = 0$ and $\frac{\partial \chi}{\partial x_\ell} \frac{\partial \eta_n}{\partial x_j} = 0$ on \mathbf{R}^N for $n > 3R_*$.

It is obvious that

$$\begin{aligned}
 (4.9) \quad & \int_{\mathbf{R}^N} x_k \left(\frac{\partial \eta_n}{\partial x_1} \frac{\partial(\chi\theta)}{\partial x_k} - \frac{\partial \eta_n}{\partial x_k} \frac{\partial(\chi\theta)}{\partial x_1} \right) dx \\
 &= \int_{\mathbf{R}^N} x_k \left[\frac{\partial}{\partial x_1} \left(\eta_n \frac{\partial(\chi\theta)}{\partial x_k} \right) - \frac{\partial}{\partial x_k} \left(\eta_n \frac{\partial(\chi\theta)}{\partial x_1} \right) \right] dx = \int_{\mathbf{R}^N} \eta_n \frac{\partial(\chi\theta)}{\partial x_1} dx.
 \end{aligned}$$

Since $|\psi|^2 - r_0^2$ and $\nabla(\chi\theta)$ belong to $L^2(\mathbf{R}^N)$, using the dominated convergence theorem we obtain

$$\begin{aligned}
 & \left| \int_{\mathbf{R}^N} x_k (|\psi|^2 - r_0^2) \left(\frac{\partial \eta_n}{\partial x_1} \frac{\partial(\chi\theta)}{\partial x_k} - \frac{\partial \eta_n}{\partial x_k} \frac{\partial(\chi\theta)}{\partial x_1} \right) dx \right| \\
 & \leq 4 \|\nabla \eta\|_{L^\infty(\mathbf{R}^N)} \int_{B(0, 2n) \setminus B(0, n)} (|\psi|^2 - r_0^2) \cdot |\nabla(\chi\theta)| dx \longrightarrow 0 \quad \text{as } n \longrightarrow \infty.
 \end{aligned}$$

(4.10)

Recall that $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} \in L^1(\mathbf{R}^N)$ by Theorem 3.1 (i), and by dominated convergence we get

$$\begin{aligned}
 & \int_{\mathbf{R}^N} \eta_n \left[\left(-i \frac{\partial \psi}{\partial x_1}, \psi \right) - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} \right] dx = \int_{\mathbf{R}^N} \eta_n \left[\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} \right] dx \\
 & \longrightarrow \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} dx \quad \text{as } n \longrightarrow \infty.
 \end{aligned}$$

(4.11)

Combining (4.8)–(4.11) we find

$$(4.12) \quad \lim_{n \rightarrow \infty} \int_{\mathbf{R}^N} \left(-ic \frac{\partial \psi}{\partial x_1}, x_k \eta_n(x) \frac{\partial \psi}{\partial x_k} \right) dx \\ = -\frac{c}{2} \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} dx.$$

Taking the scalar product of (4.2) by $\eta_n(x)x_k \frac{\partial \psi}{\partial x_k}$, integrating over \mathbf{R}^N , and using (4.5), (4.6), and (4.12) we obtain (4.4). \square

THEOREM 4.2. *Assume that $N \geq 2$, (A1), (A2) hold and let ψ be a finite-energy traveling wave for (3.1) such that $G(|\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. Suppose that either*

- $c^2 > v_s^2$, where $v_s = r_0 \sqrt{-2G'(r_0^2)}$ is the sound velocity at infinity, or
- $N = 2$ and $c^2 = v_s^2$, or
- $N \geq 3$, $c^2 = v_s^2$, and $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$.

Moreover, assume that G satisfies

- **(A3)** there exists $\alpha \in [-1 + \frac{N-3}{N-1}(1 - \frac{v_s^2}{c^2}), \frac{v_s^2}{c^2}]$ such that

$$sG(s) + \frac{v_s^2}{2}(s - r_0^2) + \left(1 - \alpha - \frac{v_s^2}{c^2}\right)V(s) \leq 0 \quad \text{for any } s \geq 0.$$

Then ψ is constant.

Proof. It follows from Propositions 2.2 and 2.5 that ψ is smooth, and Proposition 4.1 implies that ψ satisfies (4.3) and (4.4). Summing up the identities (4.4) for $k = 2, \dots, N$ we get

$$(4.13) \quad \int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_1} \right|^2 + \frac{N-3}{N-1} \sum_{k=2}^N \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx + \int_{\mathbf{R}^N} V(|\psi|^2) dx \\ - c \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial(\chi\theta)}{\partial x_1} dx = 0.$$

On the other hand, from Theorem 3.1 we have

$$(4.14) \quad \int_{\mathbf{R}^N} |\nabla \psi|^2 - G(|\psi|^2)|\psi|^2 - \frac{v_s^2}{2}(|\psi|^2 - r_0^2) dx \\ - c \left(1 - \frac{v_s^2}{c^2}\right) \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - r_0^2 \frac{\partial}{\partial x_1}(\chi\theta) dx = 0.$$

We multiply (4.13) by $-1 + \frac{v_s^2}{c^2}$ and add the resulting equality to (4.14) to get

$$(4.15) \quad \int_{\mathbf{R}^N} \frac{v_s^2}{c^2} \left| \frac{\partial \psi}{\partial x_1} \right|^2 + \left(1 - \left(1 - \frac{v_s^2}{c^2}\right) \frac{N-3}{N-1}\right) \sum_{k=2}^N \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx \\ - \int_{\mathbf{R}^N} G(|\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2) + \left(1 - \frac{v_s^2}{c^2}\right)V(|\psi|^2) dx = 0.$$

Let α satisfy (A3). Multiplying (4.3) by α and adding it to (4.15) we obtain

$$(4.16) \quad \int_{\mathbf{R}^N} \left(\frac{v_s^2}{c^2} - \alpha\right) \left| \frac{\partial \psi}{\partial x_1} \right|^2 + \left(\alpha + 1 - \left(1 - \frac{v_s^2}{c^2}\right) \frac{N-3}{N-1}\right) \sum_{k=2}^N \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx \\ = \int_{\mathbf{R}^N} G(|\psi|^2)|\psi|^2 + \frac{v_s^2}{2}(|\psi|^2 - r_0^2) + \left(1 - \alpha - \frac{v_s^2}{c^2}\right)V(|\psi|^2) dx.$$

By (A3), the right-hand side of (4.16) is less than or equal to zero. If $\alpha \in (-1 + (1 - \frac{v_s^2}{c^2})\frac{N-3}{N-1}, \frac{v_s^2}{c^2})$, it follows from (4.16) that $\int_{\mathbf{R}^N} |\frac{\partial \psi}{\partial x_k}|^2 dx = 0$ for $k = 1, \dots, N$, which implies $\nabla \psi = 0$ on \mathbf{R}^N , i.e., ψ is constant. If $\alpha = -1 + (1 - \frac{v_s^2}{c^2})\frac{N-3}{N-1}$, we infer from (4.16) that $\int_{\mathbf{R}^N} |\frac{\partial \psi}{\partial x_1}|^2 dx = 0$; consequently $\frac{\partial \psi}{\partial x_1} = 0$ on \mathbf{R}^N , which implies that ψ does not depend on x_1 . Since $\int_{\mathbf{R}^N} |\nabla \psi|^2 dx$ is finite, we have necessarily $\nabla \psi = 0$ on \mathbf{R}^N , which means that ψ is constant. A similar argument shows that ψ is constant in the case $\alpha = \frac{v_s^2}{c^2}$. \square

Remark. Let α, C_1 , and \tilde{r} be positive constants satisfying $G(s^2) + \frac{c^2}{4} \leq -C_1(s - \tilde{r})^{2\alpha}$ for any $s \geq \tilde{r}$ (such constants exist by assumption (A2)). Let ψ be as in Theorem 4.2. It follows from the proof of Proposition 2.2 (i) that $|\psi(x)| \leq \tilde{r}\sqrt{2}$ for any x . Therefore the proof of Theorem 4.2 is still valid if the inequality in (A3) holds only for all $s \in [0, 2\tilde{r}^2]$.

If $c^2 = v_s^2, N \geq 3$, and ψ is as above, we already know from Proposition 3.3 (iii) that $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^p(\mathbf{R}^N)$ for any $p \in (\frac{2N-1}{2N-3}, \infty)$. Therefore we have the following.

COROLLARY 4.3. *Assume that (A1), (A2), (A3) hold, $N \geq 3$, and $c^2 = v_s^2$. Let ψ be a traveling wave for (4.1) having finite energy, finite momentum with respect to the x_1 -direction (i.e., $\psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} \in L^1(\mathbf{R}^N)$), and such that $G(|\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. Then ψ is constant.*

Example 4.4. The Gross–Pitaevskii equation is of type (4.1) with $G(s) = 1 - s$. In this case we have $r_0 = 1, V(s) = \frac{1}{2}(s - 1)^2$, and $v_s = \sqrt{2}$. For any finite-energy function ψ we have $\int_{\mathbf{R}^N} (|\psi|^2 - 1)^2 dx < \infty$; hence $\psi \in L^4_{loc}(\mathbf{R}^N)$ and consequently $G(|\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. Assumptions (A1) and (A2) are clearly satisfied. We find $sG(s) + \frac{v_s^2}{2}(s - r_0^2) + (1 - \alpha - \frac{v_s^2}{c^2})V(s) = -(\frac{1}{2} + \alpha + \frac{v_s^2}{c^2})(1 - s)^2$. The last expression is nonpositive for any s if $\alpha \geq -\frac{1}{2} - \frac{v_s^2}{c^2}$, and thus assumption (A3) is also satisfied. Hence the conclusion of Theorem 4.2 holds for the Gross–Pitaevskii equation. In particular, we recover the nonexistence results in [17], [18].

Example 4.5. The cubic–quintic Schrödinger equation is of the form (4.1) with $G(s) = -\alpha_1 + \alpha_3 s - \alpha_5 s^2$, where $\alpha_1, \alpha_3, \alpha_5$ are positive and $\frac{3}{16} < \frac{\alpha_1 \alpha_5}{\alpha_3^2} < \frac{1}{4}$. The nonlinearity G can be written as $G(s) = -\alpha_5(s - r_1^2)(s - r_0^2)$, where $0 < r_1 < r_0$. In this case we have $v_s^2 = -2r_0^2 G'(r_0) = 2\alpha_5 r_0^2 (r_0^2 - r_1^2)$ and $V(s) = \frac{\alpha_5}{3}(s - r_0^2)^2 (s + \frac{1}{2}r_0^2 - \frac{3}{2}r_1^2)$. For any function ψ with finite energy we have $V(|\psi|^2) \in L^1(\mathbf{R}^N)$, which implies $\psi \in L^6_{loc}(\mathbf{R}^N)$ and consequently $G(|\psi|^2)\psi \in L^1_{loc}(\mathbf{R}^N)$. It is obvious that G satisfies (A1) and (A2). If $c^2 \geq v_s^2$, we have $-\frac{v_s^2}{c^2} \in [-1 + \frac{N-3}{N-1}(1 - \frac{v_s^2}{c^2}), \frac{v_s^2}{c^2}]$ and an easy computation shows that $sG(s) + \frac{v_s^2}{2}(s - r_0^2) + V(s) = -\frac{\alpha_5}{6}(4s + 5r_0^2 - 3r_1^2) \leq 0$ for any $s \geq 0$. Hence assumption (A3) holds for $\alpha = -\frac{v_s^2}{c^2}$; therefore the conclusion of Theorem 4.2 is valid for the cubic–quintic Schrödinger equation.

Remark. The proof of nonexistence of supersonic and sonic traveling waves for equations of type (1.1) relies on identity (3.1), combined with Pohozaev identities. We have proved (3.1) in an “indirect” way, starting from (3.11), using the Fourier transform and analyzing the behavior near the origin of the symbols of the differential operators involved. A natural question is whether (3.1) could be proved “directly” by multiplying the equations by appropriate functions and integrating by parts (and it is very tempting to try to do so because of the form of (3.7) and (3.8)!). We suspect that it is not possible to find such a proof, a heuristical reason being the following: If a “direct” proof of (3.1) could be found, it should be valid for any value of c .

Since Pohozaev identities are also valid for any c , one could infer that, for any c , equation (4.1) and system (4.17)–(4.18) below do not admit nontrivial finite-energy traveling waves. However, in the case of the Gross–Pitaevskii equation the existence of nontrivial, finite-energy traveling waves moving with sufficiently small speed c has been proved in [5] in dimension $N = 2$, respectively, in [7] and [12] in dimension $N = 3$. In a recent work [6], existence of traveling waves has been proved in space dimensions $N = 2$ and $N = 3$ for a wider range of speeds, including speeds c close to (and less than) v_s if $N = 2$. For Schrödinger equations of cubic–quintic type, the existence of small velocity traveling waves has been proved in [25] in any space dimension $N \geq 4$. Even for these particular cases, the question of whether such solutions exist for any speed $c \in (-v_s, v_s)$ is, to our knowledge, still open.

4.2. A Gross–Pitaevskii–Schrödinger system. Our second application concerns the system

$$(4.17) \quad i \frac{\partial \Psi}{\partial t} + \Delta \Psi - \frac{1}{\varepsilon^2} \left(|\Psi|^2 + \frac{1}{\varepsilon^2} |\Phi|^2 - 1 \right) \Psi = 0 \quad \text{in } \mathbf{R}^N,$$

$$(4.18) \quad i \delta \frac{\partial \Phi}{\partial t} + \Delta \Phi - \frac{1}{\varepsilon^2} (q^2 |\Psi|^2 - \varepsilon^2 k^2) \Phi = 0 \quad \text{in } \mathbf{R}^N,$$

which describes the motion of an uncharged impurity in a Bose condensate (see [16]). Here Ψ and Φ are the wavefunctions for bosons, respectively, for the impurity, and ε , δ , q , k are dimensionless physical constants. Assuming that the condensate is at rest at infinity, the functions Ψ and Φ must satisfy the “boundary conditions” $|\Psi| \rightarrow 1$ and $|\Phi| \rightarrow 0$ as $|x| \rightarrow \infty$.

System (4.17)–(4.18) has a Hamiltonian structure, the associated energy being

$$(4.19) \quad E(\Psi, \Phi) = \int_{\mathbf{R}^N} |\nabla \Psi|^2 + \frac{1}{\varepsilon^2 q^2} |\nabla \Phi|^2 + \frac{1}{2\varepsilon^2} (|\Psi|^2 - 1)^2 + \frac{1}{\varepsilon^4} |\Psi|^2 |\Phi|^2 - \frac{k^2}{\varepsilon^2 q^2} |\Phi|^2 dx.$$

We are interested in traveling wave solutions for (4.17)–(4.18), i.e., solutions of the form $\Psi(x, t) = \psi(x_1 - ct, x_2, \dots, x_N)$, $\Phi(x, t) = \varphi(x_1 - ct, x_2, \dots, x_N)$. Such solutions must satisfy the equations

$$(4.20) \quad -ic \frac{\partial \psi}{\partial x_1} + \Delta \psi - \frac{1}{\varepsilon^2} \left(|\psi|^2 + \frac{1}{\varepsilon^2} |\varphi|^2 - 1 \right) \psi = 0,$$

$$(4.21) \quad -ic\delta \frac{\partial \varphi}{\partial x_1} + \Delta \varphi - \frac{1}{\varepsilon^2} (q^2 |\psi|^2 - \varepsilon^2 k^2) \varphi = 0,$$

together with the boundary conditions $|\psi| \rightarrow 1$ and $|\varphi| \rightarrow 0$ as $|x| \rightarrow \infty$.

Equation (4.17) is of type (1.1). In view of the analysis in the introduction, the associated sound velocity at infinity is $\frac{\sqrt{2}}{\varepsilon}$.

In one space dimension, system (4.20)–(4.21) with the considered boundary conditions has been studied in [26]. It was proved that it admits nontrivial solutions if c is less than the sound velocity at infinity; in this case the structure of the set of traveling waves has been investigated, and it was proved that it contains global subcontinua in appropriate (weighted) Sobolev spaces.

Here we study the finite-energy traveling waves for (4.17)–(4.18) in dimension $N \geq 2$. In view of (4.19), by *finite-energy traveling wave* we mean a couple of functions

$(\psi, \varphi) \in L^1_{loc}(\mathbf{R}^N) \times L^1_{loc}(\mathbf{R}^N)$ which satisfy (4.20)–(4.21) in $\mathcal{D}'(\mathbf{R}^N)$, the boundary conditions $|\psi| \rightarrow 1, \varphi \rightarrow 0$ as $|x| \rightarrow \infty$, and such that $\nabla\psi, \nabla\varphi, \varphi \in L^2(\mathbf{R}^N)$, $(|\psi|^2 - 1)^2 + \frac{2}{\varepsilon^2}|\psi|^2|\varphi|^2 \in L^1(\mathbf{R}^N)$. As before, we denote $\psi_1 = \text{Re}(\psi), \psi_2 = \text{Im}(\psi), \varphi_1 = \text{Re}(\varphi), \varphi_2 = \text{Im}(\varphi)$. We have the following.

PROPOSITION 4.6. *Let $c \in \mathbf{R}$ and let (ψ, φ) be a finite-energy traveling wave for (4.17)–(4.18). Then the following hold.*

(i) *The function ψ is bounded and C^∞ and $\varphi, \nabla\psi \in W^{k,p}(\mathbf{R}^N)$ for any $k \in \mathbf{N}$ and $p \geq 2$.*

(ii) *There exist $R_* \geq 0$ and a real-valued function θ such that $\psi = |\psi|e^{i\theta}$ on $\mathbf{R}^N \setminus B(0, R_*)$ and $\nabla\theta \in W^{k,p}(\mathbf{R}^N \setminus B(0, R_*))$ for any $k \in \mathbf{N}$ and $p \geq 2$.*

(iii) *Let $\chi \in C^\infty(\mathbf{R}^N)$ be a cut-off function such that $\chi = 0$ on $B(0, 2R_*)$ and $\chi = 1$ on $\mathbf{R}^N \setminus B(0, 3R_*)$. We have $\psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} - \frac{\partial}{\partial x_1}(\chi\theta) \in L^1(\mathbf{R}^N)$ and the following Pohozaev-type identities hold:*

$$(4.22) \quad \int_{\mathbf{R}^N} -\left|\frac{\partial\psi}{\partial x_1}\right|^2 - \frac{1}{\varepsilon^2 q^2} \left|\frac{\partial\varphi}{\partial x_1}\right|^2 + \sum_{j=2}^N \left(\left|\frac{\partial\psi}{\partial x_j}\right|^2 + \frac{1}{\varepsilon^2 q^2} \left|\frac{\partial\varphi}{\partial x_j}\right|^2 \right) dx + \int_{\mathbf{R}^N} \frac{1}{2\varepsilon^2} (|\psi|^2 - 1)^2 + \frac{1}{\varepsilon^4} |\psi|^2 |\varphi|^2 - \frac{k^2}{\varepsilon^2 q^2} |\varphi|^2 dx = 0,$$

and for any $k \in \{2, \dots, N\}$,

$$(4.23) \quad \int_{\mathbf{R}^N} -\left|\frac{\partial\psi}{\partial x_k}\right|^2 - \frac{1}{\varepsilon^2 q^2} \left|\frac{\partial\varphi}{\partial x_k}\right|^2 + \sum_{j=1, j \neq k}^N \left(\left|\frac{\partial\psi}{\partial x_j}\right|^2 + \frac{1}{\varepsilon^2 q^2} \left|\frac{\partial\varphi}{\partial x_j}\right|^2 \right) dx + \int_{\mathbf{R}^N} \frac{1}{2\varepsilon^2} (|\psi|^2 - 1)^2 + \frac{1}{\varepsilon^4} |\psi|^2 |\varphi|^2 - \frac{k^2}{\varepsilon^2 q^2} |\varphi|^2 dx - c \int_{\mathbf{R}^N} \psi_1 \frac{\partial\psi_2}{\partial x_1} - \psi_2 \frac{\partial\psi_1}{\partial x_1} - \frac{\partial}{\partial x_1}(\chi\theta) dx - \frac{2c\delta}{\varepsilon^2 q^2} \int_{\mathbf{R}^N} \varphi_1 \frac{\partial\varphi_2}{\partial x_1} dx = 0.$$

Proof. Putting $F(x, s) = -\frac{1}{\varepsilon^2}(s + \frac{1}{\varepsilon^2}|\varphi(x)|^2 - 1)$, equation (4.20) is a particular case of (1.7). Clearly, in this case we have $r_0 = 1$.

It is obvious that F satisfies the assumptions (H1a) and (H1b) in section 2. Clearly, $F(x, s) \leq -\frac{1}{\varepsilon^2}(s - 1) \leq -\frac{1}{2\varepsilon^2}s$ for any $s \geq 2$ and $x \in \mathbf{R}^N$, and hence F satisfies (H2) for $r_* = 2$. Moreover, $\int_{r_0}^{r_*} F(x, \tau) d\tau = -\frac{1}{\varepsilon^2}(\frac{1}{2} + \frac{1}{\varepsilon^2}|\varphi(x)|^2)$ is a locally integrable function of x . We have $|\psi|^4 \leq 2(|\psi|^2 - 1)^2 + 2$ and $(|\psi|^2 - 1)^2 \in L^1(\mathbf{R})$ because (ψ, φ) has finite energy, and hence $\psi \in L^4_{loc}(\mathbf{R}^N)$. We also have $|\varphi|^2\psi \leq \frac{1}{2}(|\varphi|^2 + |\varphi|^2|\psi|^2)$ and $|\varphi|^2, |\varphi|^2|\psi|^2 \in L^1(\mathbf{R})$. It is then clear that $F(\cdot, |\psi|^2)\psi = -\frac{1}{\varepsilon^2}|\psi|^2\psi - \frac{1}{\varepsilon^4}|\varphi|^2\psi + \frac{1}{\varepsilon^2}\psi$ belongs to $L^1_{loc}(\mathbf{R}^N)$. Hence we may use Proposition 2.2 (i) and infer that $\psi \in L^\infty(\mathbf{R}^N)$.

By hypothesis we have $\varphi \in L^2(\mathbf{R}^N)$ and $\nabla\varphi \in L^2(\mathbf{R}^N)$, that is, $\varphi \in W^{1,2}(\mathbf{R}^N)$. Assume that $\varphi \in W^{1,p}(\mathbf{R}^N)$ for some $p \in (1, \infty)$. Since ψ is bounded, by (4.21) we find $\Delta\varphi \in L^p(\mathbf{R}^N)$, and we infer that $\varphi \in W^{2,p}(\mathbf{R}^N)$. If $p < N$, by the Sobolev embedding we have $\varphi \in L^{p^*}(\mathbf{R}^N)$ and $\nabla\varphi \in L^{p^*}(\mathbf{R}^N)$ (where $\frac{1}{p^*} = \frac{1}{p} - \frac{1}{N}$), and hence $\varphi \in W^{1,p^*}(\mathbf{R}^N)$. Repeating the above argument if necessary, after a finite number of steps we find $\varphi \in W^{2,q}(\mathbf{R}^N)$ for some $q \geq N$ and the Sobolev embedding implies $\varphi \in L^r(\mathbf{R}^N)$ and $\nabla\varphi \in L^r(\mathbf{R}^N)$ for any $r \in [q, \infty)$. Using (4.21) again, we conclude that $\Delta\varphi \in L^r(\mathbf{R}^N)$, and hence $\varphi \in W^{2,r}(\mathbf{R}^N)$ for any $r \in [2, \infty)$.

It follows that $\varphi \in C^1(\mathbf{R}^N)$, which implies $F \in C^1(\mathbf{R}^N)$ (and consequently F satisfies (H1c)). By Proposition 2.2 (ii) we get $\psi \in W_{loc}^{3,p}(\mathbf{R}^N)$ for any $p \in [1, \infty)$. In particular, $\psi \in C^2(\mathbf{R}^N)$.

We have $F(x, 1) = -\frac{1}{\varepsilon^4}|\varphi(x)|^2$, and F clearly satisfies assumption (H3). It is obvious that $\partial_{N+1}F(x, s) = -\frac{1}{\varepsilon^2}$ and $\partial_{N+1}^2F(x, s) = 0$ on $\mathbf{R}^N \times \mathbf{R}_+$, and therefore F satisfies (H4) and (H5). Thus we may use Proposition 2.5 (i) and infer that $\nabla\psi \in W^{1,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$.

The rest of the proof is a very easy induction. For $k \in \mathbf{N}^*$, assume that $\nabla\psi \in W^{k,p}(\mathbf{R}^N)$ and $\varphi \in W^{k+1,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$. Consider $\alpha \in \mathbf{N}^N$ such that $|\alpha| = k$. Differentiating (4.20) and (4.21) we obtain, respectively,

$$\begin{aligned} \Delta(\partial^\alpha\psi) &= ic\partial^\alpha\frac{\partial\psi}{\partial x_1} + \frac{1}{\varepsilon^2}\partial^\alpha\left(\left(|\psi|^2 + \frac{1}{\varepsilon^2}|\varphi|^2 - 1\right)\psi\right), \\ \Delta(\partial^\alpha\varphi) &= ic\delta\partial^\alpha\frac{\partial\varphi}{\partial x_1} + \frac{1}{\varepsilon^2}\partial^\alpha\left((q^2|\psi|^2 - \varepsilon^2k^2)\varphi\right). \end{aligned}$$

We infer that $\Delta(\partial^\alpha\psi), \Delta(\partial^\alpha\varphi) \in L^p(\mathbf{R}^N)$ for any $p \in [2, \infty)$. By hypothesis we have $\partial^\alpha\psi, \partial^\alpha\varphi \in L^p(\mathbf{R}^N)$, and therefore $\partial^\alpha\psi, \partial^\alpha\varphi \in W^{2,p}(\mathbf{R}^N)$ for any $p \in [2, \infty)$. Since this is true for any α with $|\alpha| = k$, we have $\nabla\psi \in W^{k+1,p}(\mathbf{R}^N)$ and $\varphi \in W^{k+2,p}(\mathbf{R}^N)$. We conclude that $\nabla\psi$ and φ belong to $W^{k,p}(\mathbf{R}^N)$ for any $k \in \mathbf{N}$ and $p \in [2, \infty)$.

(ii) is an immediate corollary of Proposition 2.5 (ii).

(iii) It follows directly from Theorem 3.1 (i) that $\psi_1\frac{\partial\psi_2}{\partial x_1} - \psi_2\frac{\partial\psi_1}{\partial x_1} - \frac{\partial}{\partial x_1}(\chi\theta) \in L^1(\mathbf{R}^N)$. The proof of (4.22) and (4.23) is similar to that of (4.3) and (4.4) (multiply (4.20) by $x_j\eta_n\frac{\partial\psi}{\partial x_j}$ and (4.21) by $\frac{1}{\varepsilon^2q^2}x_j\eta_n\frac{\partial\varphi}{\partial x_j}$, where $\eta_n(x) = \eta(\frac{x}{n})$ is a cut-off function, add the resulting equalities, integrate by parts, and pass to the limit as $n \rightarrow \infty$). We omit the details. \square

We have the following result concerning the nonexistence of supersonic traveling waves for (4.17)–(4.18).

THEOREM 4.7. *Let $N \geq 2$ and let (ψ, φ) be a finite-energy traveling wave for the system (4.17)–(4.18), moving with velocity c . Assume that either*

- $c^2 > \frac{2}{\varepsilon^2}$, or
- $N = 2$ and $c^2 = \frac{2}{\varepsilon^2}$, or
- $N \geq 3$, $c^2 = \frac{2}{\varepsilon^2}$, and $\psi_1\frac{\partial\psi_2}{\partial x_1} - \psi_2\frac{\partial\psi_1}{\partial x_1} \in L^{\frac{2N-1}{2N-3}}(\mathbf{R}^N)$.

Then $\varphi = 0$ and ψ is constant on \mathbf{R}^N .

Proof. Let θ, χ be as in Proposition 4.6 and let $F(x, s) = -\frac{1}{\varepsilon^2}(s + \frac{1}{\varepsilon^2}|\varphi(x)|^2 - 1)$. We have already seen that F satisfies assumptions (H1)–(H5), and it follows that identity (3.1) holds. Taking into account the particular form of F , this identity can be written as

$$\begin{aligned} (4.24) \quad & \int_{\mathbf{R}^N} |\nabla\psi|^2 + \frac{1}{\varepsilon^2}(|\psi|^2 - 1)^2 + \frac{1}{\varepsilon^4}|\varphi|^2|\psi|^2 dx \\ & = c\left(1 - \frac{2}{\varepsilon^2c^2}\right) \int_{\mathbf{R}^N} \psi_1\frac{\partial\psi_2}{\partial x_1} - \psi_2\frac{\partial\psi_1}{\partial x_1} - \frac{\partial}{\partial x_1}(\chi\theta) dx. \end{aligned}$$

We take the scalar product of (4.21) by φ , then integrate the resulting equality to get

$$(4.25) \quad \int_{\mathbf{R}^N} |\nabla\varphi|^2 dx + \frac{q^2}{\varepsilon^2} \int_{\mathbf{R}^N} |\varphi|^2|\psi|^2 dx - k^2 \int_{\mathbf{R}^N} |\varphi|^2 dx - 2c\delta \int_{\mathbf{R}^N} \varphi_1\frac{\partial\varphi_2}{\partial x_1} dx = 0.$$

Summing up the identities (4.23) for $k = 2, 3, \dots, N$, we find

$$\begin{aligned}
 & \int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_1} \right|^2 + \frac{1}{\varepsilon^2 q^2} \left| \frac{\partial \varphi}{\partial x_1} \right|^2 + \frac{N-3}{N-1} \sum_{j=2}^N \left(\left| \frac{\partial \psi}{\partial x_j} \right|^2 + \frac{1}{\varepsilon^2 q^2} \left| \frac{\partial \varphi}{\partial x_j} \right|^2 \right) dx \\
 (4.26) \quad & + \int_{\mathbf{R}^N} \frac{1}{2\varepsilon^2} (|\psi|^2 - 1)^2 + \frac{1}{\varepsilon^4} |\psi|^2 |\varphi|^2 - \frac{k^2}{\varepsilon^2 q^2} |\varphi|^2 dx \\
 & - c \int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - \frac{\partial}{\partial x_1} (\chi \theta) dx - \frac{2c\delta}{\varepsilon^2 q^2} \int_{\mathbf{R}^N} \varphi_1 \frac{\partial \varphi_2}{\partial x_1} dx = 0.
 \end{aligned}$$

Next we combine equalities (4.24)–(4.26) in order to eliminate the terms $\int_{\mathbf{R}^N} \varphi_1 \frac{\partial \varphi_2}{\partial x_1} dx$ and $\int_{\mathbf{R}^N} \psi_1 \frac{\partial \psi_2}{\partial x_1} - \psi_2 \frac{\partial \psi_1}{\partial x_1} - \frac{\partial}{\partial x_1} (\chi \theta) dx$. We find

$$\begin{aligned}
 & \frac{2}{\varepsilon^2 c^2} \int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_1} \right|^2 dx + \left(1 - \left(1 - \frac{2}{\varepsilon^2 c^2} \right) \frac{N-3}{N-1} \right) \int_{\mathbf{R}^N} \sum_{j=2}^N \left| \frac{\partial \psi}{\partial x_j} \right|^2 dx \\
 (4.27) \quad & + \frac{2}{(N-1)\varepsilon^2 q^2} \left(1 - \frac{2}{\varepsilon^2 c^2} \right) \int_{\mathbf{R}^N} \sum_{j=2}^N \left| \frac{\partial \varphi}{\partial x_j} \right|^2 dx \\
 & + \frac{1}{2\varepsilon^2} \left(1 + \frac{2}{\varepsilon^2 c^2} \right) \int_{\mathbf{R}^N} (|\psi|^2 - 1)^2 dx + \frac{1}{\varepsilon^4} \int_{\mathbf{R}^N} |\varphi|^2 |\psi|^2 dx = 0.
 \end{aligned}$$

Obviously, all integrals in (4.27) are nonnegative. If $c^2 \geq \frac{2}{\varepsilon^2}$, all coefficients are also nonnegative, and therefore each term in (4.27) must be zero. In particular, $\int_{\mathbf{R}^N} \left| \frac{\partial \psi}{\partial x_k} \right|^2 dx = 0$ for any $k \in \{1, \dots, N\}$, which implies $\nabla \psi = 0$ on \mathbf{R}^N , i.e., ψ is constant. Since $\int_{\mathbf{R}^N} (|\psi|^2 - 1)^2 dx = 0$, necessarily $|\psi| = 1$. We also have $0 = \int_{\mathbf{R}^N} |\varphi|^2 |\psi|^2 dx = \int_{\mathbf{R}^N} |\varphi|^2 dx$, and hence $\varphi = 0$ on \mathbf{R}^N . \square

5. The one-dimensional case. Since most of the proofs in the preceding section are not valid in space dimension $N = 1$ (in particular, we do not have identities analogous to (4.4) and (4.23)), we treat separately the one-dimensional case. It turns out that some integrations can be performed explicitly and some of the results are stronger than in higher dimensions.

Let $G : [0, \infty) \rightarrow \mathbf{R}$ be a function satisfying the following assumption:

- **(A)** $G \in C([0, \infty))$ and there exists $r_0 > 0$ such that $G(r_0^2) = 0$. Moreover, $G \in C^1([r_0^2 - \eta, r_0^2 + \eta])$ for some $\eta > 0$ and $G'(r_0^2) = -L < 0$.

We consider the Schrödinger equation

$$(5.1) \quad i \frac{\partial \Psi}{\partial t} + \Psi_{xx} + G(|\Psi|^2) \Psi = 0 \quad \text{in } \mathbf{R},$$

together with the “boundary condition” $|\Psi| \rightarrow r_0$ as $x \rightarrow \pm\infty$. We have seen in the introduction that the sound velocity at infinity associated to (5.1) and to the considered boundary condition is $v_s = r_0 \sqrt{2L}$. As usually, a traveling wave moving with velocity c is a solution of the form $\Psi(x, t) = \psi(x - ct)$. It must satisfy

$$(5.2) \quad -ic\psi' + \psi'' + G(|\psi|^2)\psi = 0 \quad \text{in } \mathbf{R}, \quad |\psi(x)| \rightarrow r_0 \quad \text{as } x \rightarrow \pm\infty.$$

We have the following result concerning supersonic and sonic traveling waves.

THEOREM 5.1. *Let $\psi \in L^1_{loc}(\mathbf{R})$ be a solution of (5.2) in $\mathcal{D}'(\mathbf{R})$ such that $G(|\psi|^2)\psi \in L^1_{loc}(\mathbf{R})$. Assume that G satisfies (A) and either*

- (i) $c^2 > v_s^2$, or
- (ii) $c^2 = v_s^2$ and, denoting $V(s) = \int_s^{r_0^2} G(\tau) d\tau$ and $W(s) = v_s^2 s^2 - 4(s + r_0^2)V(s + r_0^2)$, there exists $\varepsilon > 0$ such that one of the following conditions is verified:
 - (a) $W(s) > 0$ on $(-\varepsilon, 0) \cup (0, \varepsilon)$;
 - (b) $W(s) > 0$ on $(-\varepsilon, 0)$ and $W(s) < 0$ on $(0, \infty)$;
 - (c) $W(s) > 0$ on $(0, \varepsilon)$ and $W(s) < 0$ on $[-r_0^2, 0)$.

Then either ψ is constant or $\psi(x) = r_0 e^{i(cx + \theta_0)}$, where θ_0 is a real constant.

Remark. Theorem 5.1 gives all supersonic and sonic traveling waves for (5.1), no matter whether their energy is finite or not (and we see that finite-energy traveling waves must be constant).

It is easy to see that W is C^2 near 0 and $W(0) = W'(0) = W''(0) = 0$. Condition (ii)(a) is satisfied, for instance, if G is C^3 near r_0^2 (this clearly implies that W is C^4 near 0) and $W'''(0) = 0$, $W^{(iv)}(0) > 0$, or equivalently $r_0^2 G'''(r_0^2) = 3L$ and $4G''(r_0^2) + r_0^2 G'''(r_0^2) > 0$. The condition $W(s) > 0$ on $(-\varepsilon, 0)$ in (ii)(b), respectively, $W(s) > 0$ on $(0, \varepsilon)$ in (ii)(c), is satisfied if G is C^3 near r_0^2 and $W'''(0) < 0$ (respectively, $W'''(0) > 0$); however, in these cases only the information on the behavior of G in a neighborhood of r_0^2 is not sufficient to get the conclusion of Theorem 5.1.

Proof of Theorem 5.1. Let $\varphi(x) = e^{-\frac{icx}{2}} \psi(x)$. Then $\varphi \in L^1_{loc}(\mathbf{R})$ and it is easy to see that

$$(5.3) \quad \varphi'' + \left(G(|\varphi|^2) + \frac{c^2}{4} \right) \varphi = 0 \quad \text{in } \mathcal{D}'(\mathbf{R}).$$

From (5.3) we get $\varphi'' \in L^1_{loc}(\mathbf{R})$. This implies that φ' is a continuous function on \mathbf{R} (see, e.g., Lemma VIII.2, p. 123, in [8]). Thus $\varphi \in C^1(\mathbf{R})$. Since $|\varphi| \rightarrow r_0$ as $x \rightarrow \pm\infty$, we infer that φ is bounded on \mathbf{R} . Coming back to (5.3) we see that φ'' is continuous and bounded on \mathbf{R} . In particular $\varphi \in C^2(\mathbf{R})$, and this implies $\psi \in C^2(\mathbf{R})$.

Denoting $\psi_1 = \text{Re}(\psi)$, $\psi_2 = \text{Im}(\psi)$, equation (5.2) is equivalent to the system

$$(5.4) \quad c\psi'_2 + \psi''_1 + G(|\psi|^2)\psi_1 = 0,$$

$$(5.5) \quad -c\psi'_1 + \psi''_2 + G(|\psi|^2)\psi_2 = 0 \quad \text{in } \mathbf{R}.$$

We multiply (5.4) by $2\psi'_1$ and (5.5) by $2\psi'_2$, then add the resulting equalities to get $[(\psi'_1)^2 + (\psi'_2)^2]' - (V(|\psi|^2))' = 0$. Hence there exists $k_1 \in \mathbf{R}$ such that

$$(5.6) \quad |\psi'|^2(x) - V(|\psi|^2)(x) = k_1 \quad \text{for any } x \in \mathbf{R}.$$

Multiplying (5.4) by ψ_2 and (5.5) by $-\psi_1$, then summing up the corresponding equations we obtain $\frac{c}{2}(|\psi|^2 - r_0^2)' - (\psi_1\psi'_2 - \psi_2\psi'_1)' = 0$. Consequently there is some $k_2 \in \mathbf{R}$ such that

$$(5.7) \quad \frac{c}{2}(|\psi|^2 - r_0^2) - (\psi_1\psi'_2 - \psi_2\psi'_1) = k_2 \quad \text{in } \mathbf{R}.$$

Next we multiply (5.4) by $2\psi_1$ and (5.5) by $2\psi_2$, then add the resulting equalities to find

$$(5.8) \quad 2c(\psi_1\psi'_2 - \psi_2\psi'_1) + (|\psi|^2 - r_0^2)'' - 2|\psi'|^2 + 2G(|\psi|^2)|\psi|^2 = 0.$$

Taking into account (5.6) and (5.7), equation (5.8) can be written as

$$(5.9) \quad (|\psi|^2 - r_0^2)'' + c^2(|\psi|^2 - r_0^2) - 2V(|\psi|^2) + 2G(|\psi|^2)|\psi|^2 = 2k_1 + 2ck_2.$$

Denote $v(x) = |\psi|^2(x) - r_0^2$. Then v is real-valued, C^2 , and tends to zero as $x \rightarrow \pm\infty$, and hence there exists a sequence $x_n \rightarrow \infty$ such that $v''(x_n) \rightarrow 0$. Writing (5.9) for x_n , then passing to the limit as $n \rightarrow \infty$ we see that necessarily $k_1 + ck_2 = 0$ and v satisfies the equation

$$(5.10) \quad v'' + c^2v - 2V(v + r_0^2) + 2(v + r_0^2)G(v + r_0^2) = 0 \quad \text{in } \mathbf{R}.$$

Next we multiply (5.10) by $2v'$, then integrate the resulting equation and obtain $(v')^2 + c^2v^2 - 4(v + r_0^2)V(v + r_0^2) = k_3$ in \mathbf{R} , where k_3 is a constant. It is clear that there exists a sequence $y_n \rightarrow \infty$ such that $v'(y_n) \rightarrow 0$; consequently $k_3 = 0$ and we have

$$(5.11) \quad (v')^2(x) + c^2v^2(x) - 4(v + r_0^2)V(v + r_0^2)(x) = 0 \quad \text{for any } x \in \mathbf{R}.$$

Our aim is to prove that, under the assumptions of Theorem 5.1, we have $v = 0$ on \mathbf{R} .

Suppose first that $c^2 > v_s^2 = 2Lr_0^2$. Since G satisfies (A), it follows that $V \in C^2([r_0^2 - \eta, r_0^2 + \eta])$ and we have by Taylor's formula

$$V(r_0^2 + s) = V(r_0^2) + sV'(r_0^2) + \frac{1}{2}s^2V''(r_0^2) + s^2h(s) = \frac{1}{2}Ls^2 + s^2h(s) \quad \text{for } s \in [-\eta, \eta],$$

where $h(s) \rightarrow 0$ as $s \rightarrow 0$. Take $\varepsilon_1 \in (0, \eta]$ such that $c^2 - v_s^2 - 2Ls - 4(s + r_0^2)h(s) > 0$ for any $s \in [-\varepsilon_1, \varepsilon_1]$. Suppose that $v(x_0) \in [-\varepsilon_1, 0) \cup (0, \varepsilon_1]$ for some $x_0 \in \mathbf{R}$. By (5.11) we obtain

$$0 = (v')^2(x_0) + v^2(x_0)[c^2 - v_s^2 - 2Lv(x_0) - 4(v(x_0) + r_0^2)h(v(x_0))] > 0,$$

a contradiction. Consequently we cannot have $v(x) \in [-\varepsilon_1, 0) \cup (0, \varepsilon_1]$. Since v is continuous and $v(x) \rightarrow 0$ as $x \rightarrow \pm\infty$, we infer that necessarily $v(x) = 0$ for any $x \in \mathbf{R}$.

Next assume that $c^2 = v_s^2$. Equation (5.11) can be written as

$$(5.12) \quad (v')^2(x) + W(v(x)) = 0 \quad \text{on } \mathbf{R}.$$

If assumption (iia) is verified, we cannot have $v(x) \in (-\varepsilon, 0) \cup (0, \varepsilon)$ and we infer, as above, that $v = 0$ on \mathbf{R} . In case (iib), we cannot have $v(x) \in (-\varepsilon, 0)$ and we infer that $v(x) \geq 0$ for any $x \in \mathbf{R}$. Since $v(x) \rightarrow 0$ as $x \rightarrow \infty$, there is some x_0 such that v achieves a nonnegative maximum at x_0 . Then $v'(x_0) = 0$ and from (5.12) we get $W(v(x_0)) = 0$. But $W(s) < 0$ for $s > 0$ by (iib); hence $v(x_0) = 0$ and consequently $v = 0$ on \mathbf{R} . Similarly we have $v = 0$ in the case (iic) (note that $v = |\psi|^2 - r_0^2 \geq -r_0^2$ and it suffices to know that $W < 0$ on $[-r_0^2, 0)$).

Thus we always have $v = 0$, that is, $|\psi|^2 = r_0^2$ on \mathbf{R} . Consequently there exists a lifting $\theta \in C^2(\mathbf{R}, \mathbf{R})$ such that $\psi(x) = r_0e^{i\theta(x)}$ for any $x \in \mathbf{R}$. It is clear that $\psi_1\psi_2' - \psi_2\psi_1' = |\psi|^2\theta' = r_0^2\theta'$ (see (3.3)). On the other hand we have $\psi_1\psi_2' - \psi_2\psi_1' = -k_2$ by (5.7), and hence $\theta' = -\frac{k_2}{r_0^2}$ is constant; therefore $\theta(x) = -\frac{k_2}{r_0^2}x + \theta_0$, where θ_0 is a real constant. Since $\psi = r_0e^{i(-\frac{k_2}{r_0^2}x + \theta_0)}$ satisfies (5.2), we find $-c\frac{k_2}{r_0^2} - (\frac{k_2}{r_0^2})^2 = 0$, and thus either $\frac{k_2}{r_0^2} = 0$ or $\frac{k_2}{r_0^2} = -c$. Finally we have either $\psi(x) = e^{i\theta_0}$ or $\psi(x) = e^{i(cx + \theta_0)}$ and the proof is complete. \square

Example 5.2. In the case of the Gross–Pitaevskii equation we have $G(s) = 1 - s$ and obtain $W(s) = -2s^3$ (see Example 4.4). In the case of the cubic–quintic nonlinearity we have $G(s) = -\alpha_5(s - r_1^2)(s - r_0^2)$, where $\alpha_5 > 0$, $0 < r_1 < r_0$ (see Example 4.5)

and a simple computation gives $W(s) = -2\alpha_5 s^3 (\frac{4}{3}r_0^2 - r_1^2 + \frac{1}{3}s)$. Therefore both the Gross–Pitaevskii and the cubic–quintic nonlinearities satisfy assumption (iib) and Theorem 5.1 gives all sonic and supersonic traveling waves for these equations.

Remark. The proof of Theorem 5.1 provides a method to find subsonic traveling waves for (5.1). With the above notation, it follows from (5.11) that on any interval where $v' \neq 0$ we have $v'(x) = \pm \sqrt{4(v+r_0^2)V(v+r_0^2)(x) - c^2v^2(x)}$. In many interesting applications this equation can be integrated, and we obtain explicitly $v = |\psi|^2 - r_0^2$. Then it is not hard to find (up to a constant) the corresponding phase θ .

Remark. Assume that $N = 1$ and let (ψ, φ) be a finite-energy traveling wave for system (4.17)–(4.18). It follows from the proof of Proposition 4.6 that ψ and φ are C^∞ functions and $\psi', \varphi \in W^{k,p}(\mathbf{R})$ for any $k \in \mathbf{N}$ and $p \geq 2$. If $c^2 \geq \frac{2}{\varepsilon^2}$ (recall that $\frac{\sqrt{2}}{\varepsilon}$ is the sound velocity at infinity associated to (3.21)–(3.22)) and if there is a lifting $\psi(x) = v(x)e^{i\alpha(x)}$, $\varphi(x) = u(x)e^{i\beta(x)}$, where v, u, α, β are real-valued functions of class C^2 , [26, Proposition 3.1, p. 1545] implies that $v = 1$, α is constant, and $\varphi = 0$ on \mathbf{R} .

REFERENCES

- [1] I. V. BARASHENKOV AND V. G. MAKHANKOV, *Soliton-like “bubbles” in a system of interacting bosons*, Phys. Lett. A, 128 (1988), pp. 52–56.
- [2] I. V. BARASHENKOV, A. D. GOICHEVA, V. G. MAKHANKOV, AND I. V. PUZYININ, *Stability of soliton-like bubbles*, Phys. D, 34 (1989), pp. 240–254.
- [3] H. BERESTYCKI AND P.-L. LIONS, *Nonlinear scalar field equations, I. Existence of a ground state*, Arch. Rational Mech. Anal., 82 (1983), pp. 313–345.
- [4] J. BERGH AND J. LÖFSTRÖM, *Interpolation Spaces. An Introduction*, Springer-Verlag, Berlin, Heidelberg, New York, 1976.
- [5] F. BÉTHUEL AND J.-C. SAUT, *Travelling-waves for the Gross-Pitaevskii equation I*, Ann. Inst. H. Poincaré Phys. Théor., 70 (1999), pp. 147–238.
- [6] F. BÉTHUEL, P. GRAVEJAT, AND J.-C. SAUT, *Travelling-Waves for the Gross-Pitaevskii Equation II*, preprint; available online from <http://arxiv.org/abs/0711.2408v1>.
- [7] F. BÉTHUEL, G. ORLANDI, AND D. SMETS, *Vortex rings for the Gross-Pitaevskii equation*, J. Eur. Math. Soc. (JEMS), 6 (2004), pp. 17–94.
- [8] H. BRÉZIS, *Analyse Fonctionnelle*, Masson, Paris, 1983.
- [9] H. BRÉZIS, *Semilinear equations in \mathbf{R}^N without condition at infinity*, Appl. Math. Optim., 12 (1984), pp. 271–282.
- [10] H. BRÉZIS AND E. H. LIEB, *Minimum action solutions for some vector field equations*, Comm. Math. Phys., 96 (1984), pp. 97–113.
- [11] H. BRÉZIS, J. BOURGAIN, AND P. MIRONESCU, *Lifting in Sobolev spaces*, J. Anal. Math., 80 (2000), pp. 37–86.
- [12] D. CHIRON, *Travelling-waves for the Gross-Pitaevskii equation in dimension larger than two*, Nonlinear Anal., 58 (2004), pp. 175–204.
- [13] A. FARINA, *Finite-energy solutions, quantization effects and Liouville-type results for a variant of the Ginzburg-Landau systems in \mathbf{R}^k* , Differential Integral Equations, 11 (1998), pp. 875–893.
- [14] A. FARINA, *From Ginzburg-Landau to Gross-Pitaevskii*, Monatsh. Math., 139 (2003), pp. 265–269.
- [15] P. GÉRARD, *The Cauchy problem for the Gross-Pitaevskii equation*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 23 (2006), pp. 765–779.
- [16] J. GRANT AND P. H. ROBERTS, *Motions in a Bose condensate III. The structure and effective masses of charged and uncharged impurities*, J. Phys. A, 7 (1974), pp. 260–279.
- [17] P. GRAVEJAT, *A non-existence result for supersonic travelling-waves in the Gross-Pitaevskii equation*, Comm. Math. Phys., 243 (2003), pp. 93–103.
- [18] P. GRAVEJAT, *Limit at infinity and non-existence results for sonic travelling-waves in the Gross-Pitaevskii equation*, Differential Integral Equations, 17 (2004), pp. 1213–1232.
- [19] E. P. GROSS, *Hydrodynamics of a superfluid condensate*, J. Math. Phys., 4 (1963), pp. 195–207.
- [20] L. HÖRMANDER, *The Analysis of Linear Partial Differential Operators, Vol. 1*, Springer-Verlag, Berlin, 1983.

- [21] C. A. JONES AND P. H. ROBERTS, *Motions in a Bose condensate IV. Axisymmetric solitary waves*, J. Phys. A, 15 (1982), pp. 2599–2619.
- [22] C. A. JONES, S. J. PUTTERMAN, AND P. H. ROBERTS, *Motions in a Bose condensate V. Stability of wave solutions of nonlinear Schrödinger equations in two and three dimensions*, J. Phys. A, 19 (1986), pp. 2991–3011.
- [23] T. KATO, *Schrödinger operators with singular potentials*, Israel J. Math., 13 (1972), pp. 135–148.
- [24] P. I. LIZORKIN, *On multipliers of Fourier integrals in the spaces $L_{p,\theta}$* , Proc. Steklov Inst. Math., 89 (1967), pp. 269–290.
- [25] M. MARIŞ, *Existence of nonstationary bubbles in higher dimensions*, J. Math. Pures Appl., 81 (2002), pp. 1207–1239.
- [26] M. MARIŞ, *Global branches of travelling-waves to a Gross–Pitaevskii–Schrödinger system in one dimension*, SIAM J. Math. Anal., 37 (2006), pp. 1535–1559.
- [27] E. STEIN, *Singular Integrals and Differentiability Properties of Functions*, Princeton University Press, Princeton, NJ, 1970.

EULERIAN CALCULUS FOR THE DISPLACEMENT CONVEXITY IN THE WASSERSTEIN DISTANCE*

SARA DANERI[†] AND GIUSEPPE SAVARÉ[‡]

Abstract. In this paper we give a new proof of the (strong) displacement convexity of a class of integral functionals defined on a compact Riemannian manifold satisfying a lower Ricci curvature bound. Our approach does not rely on existence and regularity results for optimal transport maps on Riemannian manifolds, but it is based on the Eulerian point of view recently introduced by Otto and Westdickenberg [*SIAM J. Math. Anal.*, 37 (2005), pp. 1227–1255] and on the metric characterization of the gradient flows generated by the functionals in the Wasserstein space.

Key words. gradient flows, displacement convexity, heat and porous medium equation, nonlinear diffusion, optimal transport, Kantorovich-Rubinstein-Wasserstein distance, Riemannian manifolds with a lower Ricci curvature bound

AMS subject classifications. 58J35, 58C21, 49J40

DOI. 10.1137/08071346X

1. Introduction. In this paper we give a new proof, based on a gradient flow approach and on the Eulerian point of view introduced by [19], of the so-called “displacement convexity” for integral functionals as

$$(1.1) \quad \mathcal{E}(\mu) := \int_{\mathbb{M}} e(\rho) \, dV + e'(\infty) \mu^\perp(\mathbb{M}), \quad \rho = \frac{d\mu}{dV},$$

where μ is a Borel probability measure on a compact, connected Riemannian manifold without boundary (\mathbb{M}, \mathbf{g}) , V is the volume measure on \mathbb{M} induced by the metric tensor \mathbf{g} , μ^\perp is the singular part of μ with respect to V , $e : [0, +\infty) \rightarrow \mathbb{R}$ is a smooth convex function satisfying the so-called McCann conditions (see (1.7) below), and $e'(\infty) = \lim_{r \rightarrow +\infty} \frac{e(r)}{r}$. When e has a superlinear growth, $e'(\infty) = +\infty$ so that μ should be absolutely continuous with respect to V when $\mathcal{E}(\mu)$ is finite.

Displacement convexity for integral functionals. The notion of *displacement convexity* has been introduced by [15] to study the behavior of integral functionals like (1.1) along optimal transportation paths, i.e., geodesics in the space of Borel probability measures $\mathcal{P}(\mathbb{M})$ endowed with the L^2 -Kantorovich-Rubinstein-Wasserstein distance.

Recall that (the square of) this distance can be defined by the following optimal transport problem

$$(1.2) \quad \left. \begin{aligned} W_2^2(\mu^0, \mu^1) &:= \min \left\{ \int_{\mathbb{M} \times \mathbb{M}} d^2(x, y) \, d\sigma(x, y) : \sigma \in \mathcal{P}(\mathbb{M} \times \mathbb{M}), \right. \\ &\left. \sigma(B \times \mathbb{M}) = \mu^0(B), \sigma(\mathbb{M} \times B) = \mu^1(B) \quad \forall B \text{ Borel set in } \mathbb{M} \right\} \end{aligned}$$

*Received by the editors January 15, 2008; accepted for publication (in revised form) May 1, 2008; published electronically October 17, 2008.

<http://www.siam.org/journals/sima/40-3/71346.html>

[†]S.I.S.S.A., Via Beirut 2-4, 34014, Trieste, Italy (daneri@sissa.it).

[‡]Università di Pavia, Department of Mathematics, Via Ferrata 1, 27100, Pavia, Italy (giuseppe.savare@unipv.it). Partially supported by grants of M.I.U.R., PRIN '06.

for the cost function induced by the Riemannian distance d on the manifold \mathbb{M} . We keep the usual notation to denote by $\mathcal{P}_2(\mathbb{M})$ the metric space $(\mathcal{P}(\mathbb{M}), W_2)$, which is called Wasserstein space; being \mathbb{M} compact, W_2 induces the topology of the weak convergence of probability measures (i.e., the weak* topology associated with the duality of $\mathcal{P}(\mathbb{M})$ with $C^0(\mathbb{M})$).

As in any metric space (minimal, constant speed) *geodesics* can be defined as curves $\mu : s \in [0, 1] \mapsto \mu^s \in \mathcal{P}_2(\mathbb{M})$ between μ^0 and μ^1 , satisfying

$$(1.3) \quad W_2(\mu^r, \mu^s) = |s - r| W_2(\mu^0, \mu^1) \quad \forall 0 \leq r \leq s \leq 1.$$

A functional $\mathcal{E} : \mathcal{P}(\mathbb{M}) \rightarrow (-\infty, +\infty]$ is then (*strongly displacement convex* (or, more generally, (*strongly displacement λ -convex* for some $\lambda \in \mathbb{R}$) if, for all Wasserstein geodesics $\{\mu^s\}_{0 \leq s \leq 1} \subset \mathcal{P}_2(\mathbb{M})$, we have

$$(1.4) \quad \mathcal{E}(\mu^s) \leq (1 - s)\mathcal{E}(\mu^0) + s\mathcal{E}(\mu^1) - \frac{\lambda}{2}s(1 - s)W_2^2(\mu^0, \mu^1), \quad \forall s \in [0, 1].$$

A weaker notion is also often considered: one can ask that there exists *at least one* geodesic connecting μ^0 to μ^1 along which (1.4) holds. Figalli and Villani [12] has recently proved that for the class of integral functionals (1.1) considered in the present paper these two notions are equivalent.

The term “displacement convexity” arises from the strictly related concept of “displacement interpolation” introduced by [15] in the Euclidean case $\mathbb{M} = \mathbb{R}^d$; in a general metric setting, property (1.4) is simply called, as in the Riemannian case, “ λ -geodesic convexity” (or “geodesic convexity” if $\lambda = 0$).

It is possible to show [4] that, in the Euclidean case, the measures μ^s can also be defined through the formula

$$(1.5) \quad \mu^s(B) := \sigma(\{(x, y) \in \mathbb{R}^d \times \mathbb{R}^d : (1 - s)x + sy \in B\}),$$

where σ is a minimizer of (1.2).

A similar construction can also be performed in a Riemannian manifold [14, 20, 13]: the segments $s \mapsto (1 - s)x + sy$ should be substituted by a Borel map $\gamma : \mathbb{M} \times \mathbb{M} \rightarrow C^0([0, 1]; \mathbb{M})$ that at each couple $(x, y) \in \mathbb{M} \times \mathbb{M}$ associate a (minimal, constant speed) geodesic $s \mapsto \gamma^s(x, y)$ in \mathbb{M} connecting x to y . We have the representation formula

$$(1.6) \quad \mu^s(B) := \sigma(\{(x, y) \in \mathbb{M} \times \mathbb{M} : \gamma^s(x, y) \in B\}), \quad \text{where } \sigma \text{ is a minimizer of (1.2).}$$

After the pioneering paper [15], the notion of displacement convexity for integral functionals found applications in many different fields, as functional inequalities [18, 2, 9], generation, contraction, and asymptotic properties of diffusion equations and gradient flows [17, 1, 19, 4, 8, 5], Riemannian geometry, and synthetic study of metric-measure spaces [20, 21, 14].

In the context of Riemannian manifolds it turns out that displacement λ -convexity of certain classes of entropy functionals is equivalent to a lower bound for the Ricci curvature of the manifold. The connection between displacement convexity and Ricci curvature, introduced by [18], was then further deeply studied by [9, 10, 20]; the equivalence has been proved by Sturm and Von Renesse in [24], who considered the case in which the domain of the functional consists only of measures that are absolutely continuous with respect to the volume measure, and then completed by Lott and Villani [14] (with the remarks made in [12], where convexity in the strong form has

been proved), who extended the previous results to the functionals defined by (1.1) on all $\mathcal{P}(\mathbb{M})$. We refer to the forthcoming monograph [23] for further references, details, and discussions.

The strategy followed by the authors of [9] (and by all the following contributions) in order to characterize the displacement convexity of entropy functionals relies on a characterization of optimal transportation and Wasserstein geodesics [16] and on a careful study of the Jacobian properties of the exponential function which are crucial to estimate the integral functionals along this class of curves. The lack of regularity of Wasserstein geodesics and the lack of global smoothness of the squared distance function d^2 on the manifold \mathbb{M} (due to the existence of the cut-locus) require a careful use of nonsmooth analysis arguments and nontrivial approximation processes to extend the results to geodesics between arbitrary measures (see [14, 12]).

The main result is the following.

THEOREM 1.1. (I) *If $e \in C^\infty(0, +\infty)$ satisfies the McCann conditions*

$$(1.7) \quad U(\rho) := \rho e'(\rho) - (e(\rho) - e(0_+)) \geq 0, \quad \rho U'(\rho) - \left(1 - \frac{1}{n}\right)U(\rho) \geq 0, \quad n := \dim(\mathbb{M}) > 1$$

and \mathbb{M} has nonnegative Ricci curvature, then the functional \mathcal{E} defined by (1.1) is (strongly) displacement convex.

(II) *If \mathcal{E} is the relative entropy functional, corresponding to $e(\rho) = \rho \log \rho$ (which satisfies (1.7) in any dimension) in (1.1), and there exists $\lambda \in \mathbb{R}$ such that*

$$(1.8) \quad \text{Ric}_{\mathbf{g}_x}(\xi, \xi) \geq \lambda \langle \xi, \xi \rangle_{\mathbf{g}_x} \quad \forall x \in \mathbb{M}, \quad \forall \xi \in T_x \mathbb{M},$$

then the functional \mathcal{E} defined by (1.1) is (strongly) displacement λ -convex.

Remark 1.2. Besides the logarithmic entropy corresponding to $e(\rho) = \rho \log \rho$ (and $U(\rho) = \rho$), typical examples of functionals that satisfy properties (1.7) are

$$(1.9) \quad e(\rho) = \frac{1}{m-1} \rho^m, \quad U(\rho) = \rho^m, \quad m \geq 1 - \frac{1}{n}.$$

We recall that assumptions (1.7) imply the convexity of the function $\rho \mapsto e(\rho)$ (since the dimension n is greater than 1, they are in fact more restrictive).

Aim of the paper: an Eulerian approach to displacement convexity. In this paper we present an alternative proof of Theorem 1.1, which does not rely on the existence and smoothness of optimal transport maps and geodesics for the Wasserstein distance.

Our strategy can be described in three steps:

1. Following the approach suggested by Otto and Westdickenberg in [19], we work in the subspace $\mathcal{P}_2^{ar}(\mathbb{M})$ of measures with smooth and positive densities, and we use the ‘‘Riemannian’’ formula for the Wasserstein distance, originally introduced in the Euclidean framework by Benamou and Brenier [6]: if $\mu^i = \rho^i \mathbf{V} \in \mathcal{P}_2^{ar}(\mathbb{M})$, $i = 0, 1$, then [19, Proposition 4.3]

$$(1.10) \quad W_2^2(\mu^0, \mu^1) = \inf_{\mathcal{C}(\mu^0, \mu^1)} \left\{ \int_0^1 \int_{\mathbb{M}} |\nabla \phi^s|^2 \rho^s \, dV \, ds \right\} \quad \forall \mu^0, \mu^1 \in \mathcal{P}_2^{ar}(\mathbb{M})$$

where

$$(1.11) \quad \mathcal{C}(\mu^0, \mu^1) = \left\{ (\rho, \phi) : \rho \in C^\infty([0, 1] \times \mathbb{M}; \mathbb{R}_+), \quad \phi \in C^\infty([0, 1] \times \mathbb{M}) \right. \\ \left. \partial_s \rho^s + \nabla \cdot (\rho^s \nabla \phi^s) = 0 \text{ in } (0, 1) \times \mathbb{M}, \quad \mu^i = \rho^i \mathbf{V} \right\}.$$

Even though the Wasserstein space can't be endowed with a smooth Riemannian structure, (1.11) still shows a “Riemannian” characterization of the Wasserstein distance on $\mathcal{P}_2^{ar}(\mathbb{M})$.

2. The second important fact, originally showed by the so-called “Otto calculus” in [17], is that the nonlinear diffusion equation

$$(1.12) \quad \partial_t \rho_t - \Delta_g U(\rho_t) = 0 \quad \text{in } [0, +\infty) \times \mathbb{M}, \quad \rho|_{t=0} = \rho_0,$$

where $U : \mathbb{R}^+ \rightarrow \mathbb{R}$ is the function defined in (1.7) and Δ_g is the Laplace–Beltrami operator on \mathbb{M} , is the gradient flow of the functional (1.1) in $\mathcal{P}_2(\mathbb{M})$. Indeed, (1.12) corresponds to the heat equation if U is the logarithmic entropy and to the porous medium equation if U is defined by (1.9).

Starting directly from (1.10) and owing to the fact that the flow generated by (1.12) preserves smooth and positive densities, when $\text{Ric}_g(\mathbb{M}) \geq 0$ we shall show that the measures $\mu_t = \rho_t \mathbb{V} \in \mathcal{P}_2^{ar}(\mathbb{M})$ associated to the solutions of (1.12) also solve the Evolution Variational Inequality (E.V.I.)

$$(1.13) \quad \frac{1}{2} \frac{d^+}{dt} W_2^2(\nu, \mu_t) \leq \mathcal{E}(\nu) - \mathcal{E}(\mu_t) \quad \forall t \geq 0, \nu \in \mathcal{P}_2^{ar}(\mathbb{M}),$$

which has been introduced in [4] as a purely metric characterization of the gradient flows of geodesically convex functionals in metric spaces (and, in particular, in $\mathcal{P}_2(\mathbb{R}^d)$); here

$$(1.14) \quad \frac{d^+}{dt} \zeta(t) = \limsup_{h \downarrow 0} \frac{\zeta(t+h) - \zeta(t)}{h}$$

for every real function $\zeta : [0, +\infty) \rightarrow \mathbb{R}$.

When $\text{Ric}(\mathbb{M}) \geq \lambda$ (a shorthand for (1.8)), we also show that the solutions of the heat equation satisfy the modified inequality

$$(1.15) \quad \frac{1}{2} \frac{d^+}{dt} W_2^2(\nu, \mu_t) + \frac{\lambda}{2} W_2^2(\nu, \mu_t) \leq \mathcal{E}(\nu) - \mathcal{E}(\mu_t) \quad \forall t \geq 0, \nu \in \mathcal{P}_2^{ar}(\mathbb{M}),$$

where \mathcal{E} is the relative entropy functional whose integrand function is $e(\rho) = \rho \log \rho$. Note that (1.15) reduces to (1.13) when $\lambda = 0$. In order to prove (1.13) and (1.15), we propose an “Eulerian” strategy which could be adapted to more general situations.

3. The third crucial fact is the following: whenever a functional \mathcal{E} satisfies (1.13) (or, more generally, (1.15)) for a given semigroup $\mathcal{S}_t : \mu_0 = \rho_0 \mathbb{V} \mapsto \mu_t = \rho_t \mathbb{V}$ in $\mathcal{P}_2^{ar}(\mathbb{M})$, \mathcal{E} is displacement convex (resp., displacement λ -convex). Thus, the question of the behavior of \mathcal{E} along geodesics can be reduced to a differential estimate of \mathcal{E} along the smooth and positive solutions of its gradient flow.

Plan of the paper. In section 2 we present the main ideas of our approach in the simplified (finite-dimensional and smooth) setting of geodesically convex functions on Riemannian manifolds. We think that these ideas are sufficiently general to be useful in other circumstances, at least for distances which admits a Riemannian characterization as (1.10), see, e.g., [11, 7].

After a brief review of the definition of (gradient) λ -flows in arbitrary metric spaces (basically following the ideas of [4]), we present in section 3 our first result, showing that the existence of a flow satisfying the E.V.I. (1.15) (even on a dense subset of initial data, such as $\mathcal{P}_2^{ar}(\mathbb{M})$) entails the (strong) displacement λ -convexity of the functional \mathcal{E} .

Following the strategy explained in the second section, in the last two sections we prove the differential estimates showing that (1.12) satisfies (1.13) (in section 4) or, in the case of the Heat equation, (1.15) (in section 5).

2. Gradient flows and geodesic convexity in a smooth setting.

Contraction semigroups and action integrals. In order to explain the main point of our strategy, let us first consider the simple setting of a smooth function $F : X \rightarrow \mathbb{R}$ on a complete Riemannian manifold X with metric $\langle \cdot, \cdot \rangle_g$, (squared) norm $|\xi|_g^2 = \langle \xi, \xi \rangle_g$, and the endowed Riemannian distance

$$(2.1) \quad d^2(u, v) := \min \left\{ \int_0^1 |\dot{\gamma}^s|_g^2 ds, \quad \gamma : [0, 1] \rightarrow X, \gamma^0 = v, \gamma^1 = u \right\}.$$

In a smooth setting, the geodesic λ -convexity of F can be expressed through the differential condition

$$(2.2) \quad \frac{d^2}{ds^2} F(\gamma^s) \geq \lambda |\dot{\gamma}^s|_g^2$$

along any geodesic curve γ minimizing (2.1). As we discussed in the introduction, the direct computation of (2.2) could be difficult in a nonsmooth, infinite dimensional setting; it is therefore important to find equivalent conditions which avoid twofold differentiation along geodesics. One possibility, suggested in [19], is to find equivalent conditions to geodesic λ -convexity in terms of the gradient flow generated by F .

Let us recall that the gradient flow of F is a continuous semigroup of (time-dependent) maps $S_t : X \rightarrow X$, $t \in [0, +\infty)$, which at every initial datum u associate the curve $u_t := S_t(u)$ solution of the differential equation

$$(2.3) \quad \dot{u}_t = -\nabla F(u_t) \quad \forall t \geq 0, \quad u_0 = u.$$

It is well known that, when F is geodesically λ -convex, S_t is λ -contracting; i.e.,

$$(2.4) \quad d^2(S_t(u), S_t(v)) \leq e^{-2\lambda t} d^2(u, v), \quad \forall t \geq 0, \quad \forall u, v \in X.$$

By the semigroup property, (2.4) is also equivalent to the differential inequality (see (1.14))

$$(2.5) \quad \left. \frac{d^+}{dt} d^2(S_t(u), S_t(v)) \right|_{t=0} \leq -2\lambda d^2(u, v) \quad \forall u, v \in X.$$

Otto and Westdickenberg [19] revert this argument and observe that it could be easier to directly prove (2.5) by a differential estimate involving only the action of the semigroup along smooth curves; as a byproduct, one should obtain the convexity of F . To this aim, they consider a smooth curve γ^s , $s \in [0, 1]$, connecting v to u , and the action integral \mathcal{A}_t associated with its smooth perturbation

$$(2.6) \quad \gamma_t^s := S_t(\gamma^s), \quad A_t^s := |\partial_s \gamma_t^s|_g^2, \quad \mathcal{A}_t := \int_0^1 A_t^s ds,$$

where $\partial_s \gamma$ denotes the tangent vector in $T_\gamma X$ obtained by differentiating w.r.t. s . Since, by the very definition of d ,

$$(2.7) \quad d^2(S_t(v), S_t(u)) \leq \mathcal{A}_t$$

and for every $\varepsilon > 0$ one can always find a curve γ^s so that $\mathcal{A}_0 \leq d^2(u, v) + \varepsilon$ (in a smooth setting one can take $\varepsilon = 0$), (2.5) surely holds if one can prove that

$$(2.8) \quad \frac{d^+}{dt} \mathcal{A}_t \Big|_{t=0} \leq -2\lambda \mathcal{A}_0, \quad \text{or its pointwise version} \quad \frac{\partial^+}{\partial t} \Big|_{t=0} A_t^s \leq -2\lambda A_0^s.$$

Having obtained the contraction property from (2.8), it still remains open how to deduce that F is geodesically convex. Notice that along an arbitrary curve η^s

$$(2.9) \quad \frac{\partial}{\partial s} F(\eta^s) = \langle \nabla F(\eta^s), \partial_s \eta^s \rangle_g = -\langle \partial_r \mathbf{S}_r(\eta^s) \Big|_{r=0}, \partial_s \eta^s \rangle_g;$$

applied to $\eta^s := \gamma_t^s$, (2.9) and the semigroup property $\mathbf{S}_r(\gamma_t^s) = \gamma_{t+r}^s$ yield

$$(2.10) \quad \frac{\partial}{\partial s} F(\gamma_t^s) = -\langle \partial_t \gamma_t^s, \partial_s \gamma_t^s \rangle_g.$$

In a smooth setting we can assume that γ^s is a minimal geodesic; operating a further differentiation with respect to s , we obtain

$$(2.11) \quad \begin{aligned} \frac{\partial^2}{\partial s^2} F(\gamma^s) &\stackrel{(2.10)}{=} -\frac{\partial}{\partial s} \langle \partial_t \gamma_t^s, \partial_s \gamma_t^s \rangle_g \Big|_{t=0} = -\langle D_{\partial_s} \partial_t \gamma_t^s, \partial_s \gamma_t^s \rangle_g \Big|_{t=0} \\ &\quad - \langle \partial_t \gamma_t^s, D_{\partial_s} \partial_s \gamma_t^s \rangle_g \Big|_{t=0} \\ &= -\langle D_{\partial_s} \partial_t \gamma_t^s, \partial_s \gamma_t^s \rangle_g \Big|_{t=0} = -\langle D_{\partial_t} \partial_s \gamma_t^s, \partial_s \gamma_t^s \rangle_g \Big|_{t=0} \\ &= -\frac{1}{2} \frac{\partial}{\partial t} \langle \partial_s \gamma_t^s, \partial_s \gamma_t^s \rangle_g \Big|_{t=0} \\ (2.12) \quad &\stackrel{(2.6)}{=} -\frac{1}{2} \frac{\partial}{\partial t} \Big|_{t=0} A_t^s \stackrel{(2.8)}{\geq} \lambda |\partial_s \gamma^s|_g^2, \end{aligned}$$

where we used the standard properties of the covariant differentiations $D_{\partial_s}, D_{\partial_t}$ and, in (2.11), the fact that at $t = 0$ $D_{\partial_s} \partial_s \gamma_t^s = 0$, being $\gamma_t^s = \gamma^s$ a geodesic.

A metric derivation of convexity. Even if the previous differential argument shows that (2.8) implies geodesic λ -convexity, it still requires nice smooth properties on geodesics and covariant differentiation, which could be hard to extend to a nonsmooth setting.

This is not at all surprising, since the contraction property (2.5) and its action-differential characterization (2.8) do not carry all the information linking the semigroup \mathbf{S} to F : in order to conclude the argument in (2.11), we had therefore to insert the information coming from (2.9).

To overcome these difficulties, we shall deal with a more precise metric characterization of \mathbf{S} than (2.4). As it has been proposed and studied in [4], gradient flows of geodesically λ -convex functionals in “almost” Euclidean settings should satisfy a purely metric formulation in terms of the E.V.I.

$$(2.13) \quad \frac{1}{2} \frac{d^+}{dt} d^2(\mathbf{S}_t(u), v) + \frac{\lambda}{2} d^2(\mathbf{S}_t(u), v) + F(\mathbf{S}_t(u)) \leq F(v), \quad \forall v \in X, t > 0.$$

It can be proved (see [5]) that (2.13) characterizes \mathbf{S} and implies the contractivity property (2.4).

As we discussed before, here we invert the usual procedure (starting from a convex functional, construct its gradient flow) and we suppose that there exists a smooth flow

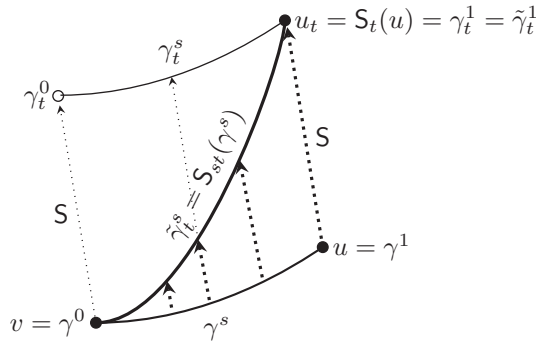


FIG. 1. Variation of the curve γ^s under the action of the semigroup S .

S_t satisfying (2.13). The following result, whose proof will be postponed (in a more general form) to Theorem 3.2 in the next section, shows that F is geodesically λ -convex.

THEOREM 2.1. *Suppose that there exists a continuous semigroup of maps $S_t \in C^0(X; X)$, $t \geq 0$, satisfying (2.13). Then for every (minimal, constant speed) geodesic $\gamma : [0, 1] \rightarrow X$*

$$(2.14) \quad F(\gamma^s) \leq (1 - s)F(\gamma^0) + sF(\gamma^1) - \frac{\lambda}{2}s(1 - s)d^2(\gamma^0, \gamma^1), \quad \forall s \in [0, 1];$$

i.e., F is (strongly) geodesically λ -convex.

E.V.I. through action-differential estimates. Thanks to Theorem 2.1, it is possible to prove the geodesic λ -convexity of F by exhibiting a flow S satisfying the E.V.I. (2.13). According to the general strategy suggested by [19], we want to reduce (2.13) to a suitable family of differential inequalities satisfied by the action A_t^s of (2.6).

The idea here is to consider a different family of perturbations of a given smooth curve $\gamma : [0, 1] \rightarrow X$, still induced by the semigroup S . In fact, differently from the contraction estimate (2.5) where we are flowing both the points u, v through S_t , in (2.13) we want to keep the point $v := \gamma^0$ fixed and to vary only $u := \gamma^1$. If γ^s is a smooth curve connecting them, it is then natural to consider the new families (see Figure 1)

$$(2.15) \quad \tilde{\gamma}_t^s := S_{st}(\gamma^s) = \gamma_{st}^s, \quad \tilde{F}_t^s := F(\tilde{\gamma}_t^s) \quad s \in [0, 1], \quad t \geq 0.$$

Notice that $\tilde{\gamma}_0^s = \gamma^s$, $\tilde{\gamma}_t^0 = \gamma^0 = v$, and $\tilde{\gamma}_t^1 = S_t(\gamma^1) = S_t(u)$. As before, we introduce the quantities

$$(2.16) \quad \tilde{A}_t^s := |\partial_s \tilde{\gamma}_t^s|_g^2, \quad \tilde{\mathcal{A}}_t := \int_0^1 \tilde{A}_t^s ds.$$

THEOREM 2.2 (A differential inequality linking action and flow). *Suppose that for every smooth curve $\gamma : [0, 1] \rightarrow X$, the quantities $\tilde{A}_t^s, \tilde{F}_t^s$ induced by the flow S through (2.15), (2.16) satisfy*

$$(2.17) \quad \frac{1}{2} \frac{\partial}{\partial t} \tilde{A}_t^s + \frac{\partial}{\partial s} \tilde{F}_t^s \leq -\lambda s \tilde{A}_t^s, \quad \forall t \geq 0.$$

Then S satisfies (2.13), it is the gradient flow of F , and F is geodesically λ -convex. Moreover, it is sufficient to check (2.17) at $t = 0$.

Proof. Let us first observe that (2.17) yields, after an integration with respect to s in $[0, 1]$,

$$(2.18) \quad \frac{1}{2} \frac{d}{dt} \tilde{A}_t + \tilde{F}_t^1 - \tilde{F}_t^0 \leq -\lambda \int_0^1 s \tilde{A}_t^s ds.$$

By the semigroup property, it is sufficient to prove (2.13) at $t = 0$. We choose a geodesic γ^s connecting v to u , and we consider the curves given by (2.15). Since

$$(2.19) \quad d^2(S_t(u), v) \leq \int_0^1 \tilde{A}_t^s ds = \tilde{A}_t,$$

$$d^2(v, u) = \int_0^1 \tilde{A}_0^s ds = \tilde{A}_0, \quad \tilde{F}_t^1 = F(S_t(u)), \quad \tilde{F}_t^0 = F(v),$$

by (2.18) at $t = 0$ we obtain

$$(2.20) \quad \frac{1}{2} \frac{d^+}{dt} d^2(S_t(u), v) \Big|_{t=0} + F(u) - F(v) \leq -\lambda \int_0^1 s \tilde{A}_0^s ds = -\frac{\lambda}{2} d^2(u, v),$$

where in the last identity we used the fact that γ^s is a geodesic, and therefore $\tilde{A}_0^s = |\partial_s \gamma^s|_g^2$ is constant in $[0, 1]$ and takes the value $d^2(\gamma^0, \gamma^1) = d^2(v, u)$.

Since $\tilde{\gamma}_{t_0+t}^s = S_{st} \tilde{\gamma}_{t_0}^s$ by the semigroup property, if S satisfies (2.17) at the initial time $t = 0$ for an arbitrary smooth curve γ , then it also satisfies (2.17) for $t > 0$. \square

Our last result provides a simple criterion to check (2.17):

THEOREM 2.3. *Suppose that $S : [0, +\infty) \times X \rightarrow X$ is the “differential” gradient flow of F satisfying (2.9) for any smooth curve γ^s , let $\gamma_t^s, \tilde{\gamma}_t^s, A_t^s, \tilde{A}_t^s, \tilde{F}_t^s$ be defined as in (2.6), (2.15), and (2.16), and let us set*

$$(2.21) \quad \tilde{D}_r^s := \frac{1}{2} \lim_{h \downarrow 0} h^{-1} \left(|\partial_s \gamma_{sr+h}^s|_g^2 - |\partial_s \gamma_{sr}^s|_g^2 \right).$$

Then

$$(2.22) \quad \frac{1}{2} \frac{\partial}{\partial t} \tilde{A}_t^s + \frac{\partial}{\partial s} \tilde{F}_t^s = s \tilde{D}_t^s.$$

Furthermore, if (2.8) holds, then

$$(2.23) \quad \tilde{D}_t^s \leq -\lambda \tilde{A}_t^s,$$

and (2.17) holds, too, so that F is geodesically λ -convex, and S is also its “metric” gradient flow, characterized by the E.V.I. (2.13).

Proof. Let us set

$$(2.24) \quad \tilde{\gamma}_{t,\tau}^s := S_\tau \tilde{\gamma}_t^s = \gamma_{st+\tau}^s, \quad \tilde{A}_{t,\tau}^s := |\partial_s \tilde{\gamma}_{t,\tau}^s|_g^2,$$

so that

$$(2.25) \quad \tilde{\gamma}_{t+h}^s = \tilde{\gamma}_{t,sh}^s, \quad \partial_s \tilde{\gamma}_{t+h}^s = \partial_s \tilde{\gamma}_{t,\tau}^s + h \partial_\tau \tilde{\gamma}_{t,\tau}^s \Big|_{\tau=sh}, \quad \tilde{D}_t^s = \frac{1}{2} \frac{\partial}{\partial \tau} \tilde{A}_{t,\tau}^s \Big|_{\tau=0}.$$

Observe that the identity

$$(2.26) \quad |x + y|_g^2 = 2\langle x + y, y \rangle_g + |x|_g^2 - |y|_g^2, \quad \forall x, y \in T_\gamma X$$

yields

$$\begin{aligned} \tilde{A}_{t+h}^s &= \left| \partial_s \tilde{\gamma}_{t+h}^s \right|_g^2 \stackrel{(2.25)}{=} \left| \partial_s \tilde{\gamma}_{t,\tau}^s + h \partial_\tau \tilde{\gamma}_{t,\tau}^s \right|_g^2 \Big|_{\tau=sh} \\ &\stackrel{(2.26)}{=} \left[2h \langle \partial_s \tilde{\gamma}_{t,\tau}^s + h \partial_\tau \tilde{\gamma}_{t,\tau}^s, \partial_\tau \tilde{\gamma}_{t,\tau}^s \rangle + \left| \partial_s \tilde{\gamma}_{t,\tau}^s \right|_g^2 - h^2 \left| \partial_\tau \tilde{\gamma}_{t,\tau}^s \right|_g^2 \right] \Big|_{\tau=sh} \\ &= 2h \langle \partial_s \tilde{\gamma}_{t+h}^s, \partial_\theta \mathbf{S}_\theta(\tilde{\gamma}_{t+h}^s) \rangle \Big|_{\theta=0} + \tilde{A}_{t,sh}^s - o(h) \stackrel{(2.9)}{=} -2h \frac{\partial}{\partial s} F(\tilde{\gamma}_{t+h}^s) + \tilde{A}_{t,sh}^s - o(h). \end{aligned}$$

We thus get

$$(2.27) \quad \frac{1}{2h} (\tilde{A}_{t+h}^s - \tilde{A}_t^s) + \frac{\partial}{\partial s} F(\tilde{\gamma}_{t+h}^s) = \frac{1}{2h} (\tilde{A}_{t,sh}^s - \tilde{A}_t^s) - o(1),$$

so that, passing to the limit as $h \downarrow 0$ we get (2.22). \square

Remark 2.4. Notice that the remainder term $o(1)$ in (2.27) is nonnegative, so it can be simply neglected, if one is just interested in the inequality (2.17).

3. Gradient flows and geodesic convexity in a metric setting. In this section we will briefly recall some basic definitions and properties of gradient flows in a metric setting, and we will prove Theorem 2.1 in a slightly more general framework.

Let (X, d) be a metric space (not necessarily complete) and let $F : X \rightarrow (-\infty, +\infty]$ be a lower semicontinuous functional, whose proper domain $D(F) := \{w \in X : F(w) < +\infty\}$ is dense in X (otherwise we can always restrict all the next statements to the closure of $D(F)$ in X). We also assume that F is bounded from below, i.e., $F_{\inf} := \inf_{u \in X} F(u) > -\infty$.

A C^0 -semigroup S in $C^0(X; X)$ is a family $S_t, t \geq 0$, of continuous maps in X such that

$$(3.1) \quad S_{t+h}(u) = S_h(S_t(u)), \quad \lim_{t \downarrow 0} S_t(u) = S_0(u) = u \quad \forall u \in X, t, h \geq 0.$$

Given a real number $\lambda \in \mathbb{R}$, we say that S is the λ -(gradient) flow of F if it satisfies

$$(3.2a) \quad S_t(X) \subset D(F) \text{ for every } t > 0;$$

$$(3.2b) \quad \text{the map } t \mapsto F(S_t(u)) \text{ is not increasing in } (0, +\infty);$$

$$(3.2c)$$

$$\frac{1}{2} \frac{d^+}{dt} d^2(S_t(u), v) + \frac{\lambda}{2} d^2(S_t(u), v) + F(S_t(u)) \leq F(v), \quad \forall u \in X, v \in D(F), t \geq 0.$$

Clearly, if S is a λ -flow for F , then it is also a λ' -flow for every $\lambda' \leq \lambda$. The next proposition collects some useful properties of λ -flows.

PROPOSITION 3.1 (Integral characterization of flows and contraction). *A C^0 -semigroup S satisfies (3.2a,b,c) if and only if it satisfies the following integrated form*

$$(3.3) \quad \frac{e^{\lambda(t_1-t_0)}}{2} d^2(S_{t_1}(u), v) - \frac{1}{2} d^2(S_{t_0}(u), v) \leq E_\lambda(t_1-t_0) (F(v) - F(S_{t_1}(u))) \quad \forall 0 \leq t_0 < t_1,$$

for every $u \in X, v \in D(F)$, where

$$E_\lambda(t) := \int_0^t e^{\lambda r} dr = \begin{cases} \frac{e^{\lambda t} - 1}{\lambda} & \text{if } \lambda \neq 0, \\ t & \text{if } \lambda = 0. \end{cases}$$

In particular, S satisfies the uniform regularization bound

$$(3.4) \quad F(S_t(u)) \leq F(v) + \frac{1}{2E_\lambda(t)} d^2(u, v) \quad \forall u \in X, v \in D(F), t > 0,$$

the uniform continuity estimate

$$(3.5) \quad d^2(S_{t_1}(u), S_{t_0}(u)) \leq 2E_{-\lambda}(t_1 - t_0) (F(S_{t_0}u) - F_{\text{inf}}) \quad \forall u \in D(F), 0 \leq t_0 \leq t_1,$$

and the λ -contraction property, i.e.,

$$(3.6) \quad d(S_t(u), S_t(v)) \leq e^{-\lambda t} d(u, v) \quad \forall u, v \in X, t \geq 0.$$

Proof. Clearly (3.3) yields (3.2a), being $D(F) \neq \emptyset$; (3.2b) and (3.5) follow by taking $v := S_{t_0}(u)$, and (3.2c) can be proved by dividing both sides of (3.3) by $t_1 - t_0$ and passing to the limit as $t_1 \downarrow t_0$. In order to prove the converse implication, let us first observe that for a continuous real function $\zeta : [0, +\infty) \rightarrow \mathbb{R}$

$$(3.7) \quad \liminf_{h \downarrow 0} \frac{\zeta(t+h) - \zeta(t)}{h} \leq 0 \quad \forall t > 0 \quad \implies \quad \zeta \text{ is not increasing.}$$

In fact, if $0 \leq t_0 < t_0 + \tau$ existed with $\delta := \tau^{-1}(\zeta(t_0 + \tau) - \zeta(t_0)) > 0$, then a minimum point $\bar{t} \in [t_0, t_0 + \tau)$ of $t \mapsto \zeta(t) - \zeta(t_0) - \delta(t - t_0)$ would satisfy

$$\liminf_{h \downarrow 0} \frac{\zeta(\bar{t} + h) - \zeta(\bar{t})}{h} - \delta \geq 0, \quad \text{which contradicts (3.7).}$$

(3.3) then follows by (3.2c), after a multiplication by $e^{\lambda t}$ and choosing

$$\zeta(t) := \frac{e^{\lambda t}}{2} d^2(S_t(u), v) + \int_{\bar{t}}^t e^{\lambda r} (F(S_r(u)) - F(v)) dr, \quad \bar{t} > 0,$$

and recalling the monotonicity property (3.2b). A similar argument shows that

$$(3.8) \quad \frac{1}{2} d^2(S_{t_1}(u), v) - \frac{1}{2} d^2(S_{t_0}(u), v) + \frac{\lambda}{2} \int_{t_0}^{t_1} d^2(S_r(u), v) dr \leq (t_1 - t_0) (F(v) - F(S_{t_1}(u))),$$

for every $0 \leq t_0 < t_1$, $u \in X$, and $v \in D(F)$. In order to prove the λ -contracting property, we apply (3.8) obtaining

$$\begin{aligned} d^2(S_h(u), S_h(v)) - d^2(u, v) &= d^2(S_h(u), S_h(v)) - d^2(S_h(u), v) + d^2(S_h(u), v) - d^2(u, v) \\ &\leq -\lambda \int_0^h (d^2(S_h(u), S_r(v)) + d^2(S_r(u), v)) dr + 2h (F(v) - F(S_h(v))). \end{aligned}$$

We divide this inequality by h , and we pass to the limit as $h \downarrow 0$; the continuity of S_t , the lower semicontinuity of F , and the semigroup property of S yield

$$(3.9) \quad \frac{d^+}{dt} d^2(S_t(u), S_t(v)) \leq -2\lambda d^2(u, v) \quad \forall u, v \in X, t > 0,$$

which yields (3.6) thanks to (3.7). \square

We can now prove the main result of this section: if a functional F admits a λ -flow, then F is geodesically λ -convex.

THEOREM 3.2 (Geodesic convexity via E.V.I.). *Let us suppose that S is a λ -flow for the functional F , according to (3.2a,b,c), and let $\gamma : [0, 1] \rightarrow X$ be a Lipschitz curve satisfying*

$$(3.10) \quad d(\gamma^r, \gamma^s) \leq L|r - s|, \quad L^2 \leq d^2(\gamma^0, \gamma^1) + \varepsilon^2 \quad \forall r, s \in [0, 1],$$

for some constant $\varepsilon \geq 0$. Then for every $t > 0$ and $s \in [0, 1]$

$$(3.11) \quad F(S_t(\gamma^s)) \leq (1 - s)F(\gamma^0) + sF(\gamma^1) - \frac{\lambda}{2}s(1 - s)d^2(\gamma^0, \gamma^1) + \frac{\varepsilon^2}{2E_\lambda(t)}s(1 - s).$$

In particular, when γ is a geodesic (i.e., γ satisfies (3.10) with $L = d(\gamma^0, \gamma^1)$, $\varepsilon = 0$), we have

$$(3.12) \quad F(\gamma^s) \leq (1 - s)F(\gamma^0) + sF(\gamma^1) - \frac{\lambda}{2}s(1 - s)d^2(\gamma^0, \gamma^1);$$

i.e., F is (strongly) geodesically λ -convex.

Proof. Let γ be satisfying (3.10) and let us set $\gamma_t^s := S_t(\gamma^s)$. Choosing $t_0 = 0$, $t_1 = t$, $u := \gamma^s$, and taking a convex combination of (3.3) written for $v := \gamma^0$, and $v := \gamma^1$, we get

$$(3.13) \quad \begin{aligned} & \frac{e^{\lambda t}}{2} ((1 - s)d^2(\gamma_t^s, \gamma^0) + s d^2(\gamma_t^s, \gamma^1)) - \frac{1}{2} ((1 - s)d^2(\gamma^s, \gamma^0) + s d^2(\gamma^s, \gamma^1)) \\ & \leq E_\lambda(t) ((1 - s)F(\gamma^0) + sF(\gamma^1) - F(\gamma_t^s)). \end{aligned}$$

We now observe that the elementary inequality

$$(3.14) \quad (1 - s)a^2 + sb^2 \geq s(1 - s)(a + b)^2 \quad \forall a, b \in \mathbb{R}, \quad s \in [0, 1],$$

and the triangular inequality yield

$$(3.15) \quad \begin{aligned} (1 - s)d^2(\gamma_t^s, \gamma^0) + sd^2(\gamma_t^s, \gamma^1) & \stackrel{(3.14)}{\geq} s(1 - s)(d(\gamma_t^s, \gamma^0) + d(\gamma_t^s, \gamma^1))^2 \\ & \geq s(1 - s)d^2(\gamma^0, \gamma^1). \end{aligned}$$

On the other hand, (3.10) yields

$$(3.16) \quad (1 - s)d^2(\gamma^s, \gamma^0) + sd^2(\gamma^s, \gamma^1) \leq L^2s(1 - s).$$

Inserting (3.16) and (3.15) in (3.13) we obtain

$$(3.17) \quad \frac{e^{\lambda t} - 1}{2}s(1 - s)d^2(\gamma^0, \gamma^1) - \frac{\varepsilon^2}{2}s(1 - s) \leq E_\lambda(t) ((1 - s)F(\gamma^0) + sF(\gamma^1) - F(\gamma_t^s)).$$

Dividing then both sides of (3.17) by $E_\lambda(t)$ we get (3.11); when $\varepsilon = 0$ we can pass to the limit as $t \downarrow 0$ obtaining (3.12). \square

We conclude this section by considering the case when the flow S is defined only on a *dense* subset X_0 of $D(F)$ (which is dense in X). In order to prove the geodesic convexity of F in X by Theorem 3.2 we first have to extend S to the whole space X . This can be achieved by a density argument, if X is complete and the lower semicontinuous functional F satisfies the following approximation property:

$$(3.18) \quad \forall u \in D(F) \quad \exists u_n \in X_0: \quad \lim_{n \rightarrow \infty} d(u_n, u) = 0, \quad \lim_{n \rightarrow \infty} F(u_n) = F(u).$$

We state the precise extension result in the next theorem.

THEOREM 3.3. *Suppose that the functional F and the subset $X_0 \subset D(F)$ satisfy (3.18) and let S be a λ -flow for F in X_0 . If X is complete, S can be extended to a unique λ -flow \bar{S} in X , and therefore F is (strongly) geodesically λ -convex in X .*

Proof. Given $u \in X$ and a sequence $u_n \in X_0$ converging to u (X_0 is also dense in X), we can define

$$(3.19) \quad \bar{S}_t(u) := \lim_{n \rightarrow \infty} S_t(u_n) \quad \forall t > 0,$$

where it is clear that the limit in (3.19) exists (being X complete and S_t Lipschitz by (3.6)) and does not depend on the particular sequence u_n we used to approximate u . Moreover, \bar{S}_t is a semigroup and satisfies the estimate (3.5) and the λ -contracting property (3.6); being X_0 dense in X , it is not difficult to combine (3.5), (3.6), and (3.18) to show that $\lim_{t \downarrow 0} \bar{S}_t(u) = u$ for every $u \in X$.

In order to prove that \bar{S} is still a λ -flow for F in X we have to check (3.3) in X : we fix $v \in D(F)$ and a sequence $v_n \in X_0$ converging to v with $F(v_n) \rightarrow F(v)$, and we pass to the limit as $n \rightarrow \infty$ in the inequalities

$$(3.20) \quad \frac{e^{\lambda(t_1-t_0)}}{2} d^2(S_{t_1}(u_n), v_n) - \frac{1}{2} d^2(S_{t_0}(u_n), v_n) \leq E_\lambda(t_1 - t_0)(F(v_n) - F(S_{t_1}(u_n))),$$

using the lower semicontinuity of F . \square

4. Nonlinear diffusion equations as gradient flows of entropy functionals in $\mathcal{P}_2(\mathbb{M})$. We apply the strategy described in the section 2 to prove the geodesic convexity of the integral functional (1.1) in the case of a Riemannian manifold of nonnegative Ricci curvature. We therefore exhibit a smooth flow (induced by the nonlinear diffusion equation (1.12) on the dense subset $\mathcal{P}_2^{gr}(\mathbb{M})$) which satisfies the E.V.I. (1.13).

Before stating the main theorem of this section, let us recall a fundamental result on this kind of evolution equations, which can be found in [22, 19]:

THEOREM 4.1 (Classical solutions of nonlinear diffusion equations). *Let $e \in C^\infty(\mathbb{R}^+)$ and U be functions that satisfy the assumptions (1.7) of Theorem 1.1. For every $\rho_0 \in C^\infty(\mathbb{M})$ with $\rho_0 > 0$, there exists a unique smooth positive solution $\rho \in C^\infty([0, +\infty) \times X)$ to the Cauchy problem*

$$(4.1) \quad \partial_t \rho_t = \Delta_g U(\rho_t), \quad \rho|_{t=0} = \lim_{t \downarrow 0} \rho_t = \rho_0.$$

Moreover, given a one-parameter family of positive initial data $s \mapsto \rho_0^s \in C^\infty([0, 1] \times \mathbb{M})$, the corresponding solutions ρ_t^s of the equation (4.1) depend smoothly on s, t .

For every $\mu_0 = \rho_0 V \in \mathcal{P}_2^{gr}(\mathbb{M})$ we denote by $\mathcal{S}_t(\mu_0) \in \mathcal{P}_2^{gr}(\mathbb{M})$ the measure $\mu_t = \rho_t V$. The main result that we show in this section is the following.

THEOREM 4.2. *Let $e \in C^\infty(\mathbb{R}^+)$ and U be functions that satisfy the assumptions (1.7) of Theorem 1.1 and let us suppose that*

$$(4.2) \quad \text{Ric}_g(x) \geq 0 \quad \forall x \in \mathbb{M}.$$

The semigroup \mathcal{S} induced by (4.1) in $\mathcal{P}_2^{gr}(\mathbb{M})$ is a 0-flow in $\mathcal{P}_2^{gr}(\mathbb{M})$ for the functional

$$(4.3) \quad \mathcal{E}(\mu) = \int_{\mathbb{M}} e(\rho) \, dV, \quad \forall \mu = \rho V \in \mathcal{P}_2^{gr}(\mathbb{M}).$$

In particular, for every $\mu_0 = \rho_0 V, \nu \in \mathcal{P}_2^{ar}(\mathbb{M})$, the measures $\mu_t = \mathcal{S}_t(\mu_0) = \rho_t V \in \mathcal{P}_2^{ar}(\mathbb{M})$ solving (4.1) satisfy the E.V.I.

$$(4.4) \quad \frac{1}{2} \frac{d^+}{dt} W_2^2(\nu, \mu_t) \leq \mathcal{E}(\nu) - \mathcal{E}(\mu_t) \quad \forall t \in [0, +\infty).$$

In order to prove Theorem 4.2, thanks to the ‘‘Riemannian-like’’ characterization of the Wasserstein distance provided by (1.10), we can follow the strategy presented in section 2; in particular we want to prove the differential inequality of Theorem 2.2. Following Otto’s formalism [17], we collect in the next table the formal correspondences between the various objects:

X, Riemannian manifold, with distance d	$\mathcal{P}_2^{ar}(\mathbb{M})$ with distance W_2
a smooth curve γ^s in X	a smooth family $\mu^s = \rho^s V \in \mathcal{P}_2^{ar}(\mathbb{M})$
the tangent vector $\partial_s \gamma^s$ in $T_{\gamma^s} X$	the vector field $\nabla \phi^s$ where $-\nabla \cdot (\rho^s \nabla \phi^s) = \frac{\partial}{\partial s} \rho^s$
$ \partial_s \gamma^s _g^2$	$\int_{\mathbb{M}} \nabla \phi^s(x) _g^2 \rho^s(x) dV(x)$
$\gamma_t^s := \mathcal{S}_t(\gamma^s), \tilde{\gamma}_t^s := \gamma_{st}^s = \mathcal{S}_{st}(\gamma^s)$	$\mu_t^s = \rho_t^s V := \mathcal{S}_t(\mu^s), \tilde{\mu}_t^s = \tilde{\rho}_t^s V := \mu_{st}^s = \mathcal{S}_{st}(\mu^s)$
$\tilde{A}_t^s = \partial_s \tilde{\gamma}_t^s _g^2$	$\int_{\mathbb{M}} \nabla \tilde{\phi}_t^s(x) _g^2 \tilde{\rho}_t^s(x) dV(x)$
$F(\gamma^s)$	$\mathcal{E}(\mu^s) = \int_{\mathbb{M}} e(\rho^s) dV$
$(\partial_\theta \mathcal{S}_\theta \gamma^s) _{\theta=0} = -\nabla F(\gamma^s)$	$-\nabla U(\rho^s)/\rho^s = -\nabla e'(\rho^s).$

The core of the proof of Theorem 4.2 lies in the following lemma:

LEMMA 4.3. Let $\mu^s = \rho^s V, s \in [0, 1]$, be a smooth family of measures in $\mathcal{P}_2^{ar}(\mathbb{M})$ and let $\tilde{\mu}_t^s = \tilde{\rho}_t^s V = \mathcal{S}_{st}(\mu^s)$ be obtained by flowing ρ^s along the flow (4.1); i.e., $\tilde{\rho}_t^s = \rho_{st}^s$ where ρ_t^s satisfies

$$(4.5) \quad \frac{\partial}{\partial t} \rho_t^s - \Delta_g U(\rho_t^s) = 0 \text{ in } \mathbb{M}, \quad \forall s \in [0, 1], t > 0; \quad \rho_{t=0}^s = \rho^s.$$

Let $\tilde{\phi}_t^s \in C^\infty([0, 1] \times [0, +\infty) \times \mathbb{M})$ be the functions defined by the equation

$$(4.6) \quad -\nabla \cdot (\tilde{\rho}_t^s \nabla \tilde{\phi}_t^s) = \partial_s \tilde{\rho}_t^s \quad \text{in } \mathbb{M}, \quad \int_{\mathbb{M}} \tilde{\phi}_t^s(x) dV(x) = 0 \quad \forall s \in [0, 1], t \in [0, +\infty),$$

and let us set

$$(4.7) \quad \begin{aligned} \tilde{A}_t^s &:= \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s(x)|_g^2 \tilde{\rho}_t^s(x) dV(x), \\ \tilde{D}_t^s &:= - \int_{\mathbb{M}} \left[(|\text{Hess } \tilde{\phi}_t^s|_g^2 + \text{Ric}_g(\nabla \tilde{\phi}_t^s, \nabla \tilde{\phi}_t^s)) U(\tilde{\rho}_t^s) + (\Delta_g \tilde{\phi}_t^s)^2 (\tilde{\rho}_t^s U'(\tilde{\rho}_t^s) - U(\tilde{\rho}_t^s)) \right] dV. \end{aligned}$$

Then, we have the formula

$$(4.8) \quad \frac{\partial}{\partial t} \frac{1}{2} \tilde{A}_t^s + \frac{\partial}{\partial s} \mathcal{E}(\tilde{\rho}_t^s V) = s \tilde{D}_t^s, \quad \forall t \in [0, +\infty), \forall s \in [0, 1].$$

In particular, if \mathbb{M} has nonnegative Ricci curvature, then $\tilde{D}_t^s \leq 0$ and therefore

$$(4.9) \quad \frac{\partial}{\partial t} \frac{1}{2} \tilde{A}_t^s + \frac{\partial}{\partial s} \mathcal{E}(\tilde{\rho}_t^s V) \leq 0.$$

Proof. Being $\tilde{\rho}_t^s := \rho_\tau^\sigma|_{\sigma=s, \tau=st}$ we get

$$(4.10) \quad \frac{\partial}{\partial s} \tilde{\rho}_t^s = \left(\frac{\partial}{\partial \sigma} \rho_\tau^\sigma + t \frac{\partial}{\partial \tau} \rho_\tau^\sigma \right)_{\sigma=s, \tau=st}, \quad \frac{\partial}{\partial t} \tilde{\rho}_t^s = s \partial_\tau \rho_\tau^s|_{\tau=st} = s \Delta_{\mathbf{g}} U(\tilde{\rho}_t^s),$$

$$(4.11) \quad \frac{\partial^2}{\partial t \partial s} \tilde{\rho}_t^s \stackrel{(4.6)}{=} -\nabla \cdot \left(\frac{\partial}{\partial t} \tilde{\rho}_t^s \nabla \tilde{\phi}_t^s \right) - \nabla \cdot \left(\tilde{\rho}_t^s \frac{\partial}{\partial t} \nabla \tilde{\phi}_t^s \right),$$

$$(4.12) \quad \frac{\partial^2}{\partial s \partial t} \tilde{\rho}_t^s \stackrel{(4.10)}{=} s \Delta_{\mathbf{g}} \left(U'(\tilde{\rho}_t^s) \frac{\partial}{\partial s} \tilde{\rho}_t^s \right) + \Delta_{\mathbf{g}} U(\tilde{\rho}_t^s) \\ \stackrel{(4.6)}{=} -s \Delta_{\mathbf{g}} \left(U'(\tilde{\rho}_t^s) \nabla \cdot \left(\tilde{\rho}_t^s \nabla \tilde{\phi}_t^s \right) \right) + \Delta_{\mathbf{g}} U(\tilde{\rho}_t^s).$$

Differentiation and integration by parts yield

$$\frac{\partial}{\partial t} \int_{\mathbb{M}} \frac{1}{2} |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \tilde{\rho}_t^s \, dV = \int_{\mathbb{M}} \left\langle \frac{\partial}{\partial t} \nabla \tilde{\phi}_t^s, \nabla \tilde{\phi}_t^s \right\rangle_{\mathbf{g}} \tilde{\rho}_t^s \, dV + \frac{1}{2} \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \frac{\partial}{\partial t} \tilde{\rho}_t^s \, dV \\ = - \int_{\mathbb{M}} \nabla \cdot \left(\tilde{\rho}_t^s \frac{\partial}{\partial t} \nabla \tilde{\phi}_t^s \right) \tilde{\phi}_t^s \, dV \stackrel{(4.10)}{=} \frac{1}{2} s \int_{\mathbb{M}} \Delta_{\mathbf{g}} \left(|\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \right) U(\tilde{\rho}_t^s) \, dV \\ \stackrel{(4.11)}{=} \int_{\mathbb{M}} \frac{\partial^2}{\partial t \partial s} \tilde{\rho}_t^s \tilde{\phi}_t^s \, dV + \int_{\mathbb{M}} \left(\nabla \cdot \left(\frac{\partial}{\partial t} \tilde{\rho}_t^s \nabla \tilde{\phi}_t^s \right) \right) \tilde{\phi}_t^s \, dV + \frac{1}{2} s \int_{\mathbb{M}} \Delta_{\mathbf{g}} \left(|\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \right) U(\tilde{\rho}_t^s) \, dV$$

$$(4.13) \quad \stackrel{(4.12)}{=} \int_{\mathbb{M}} \left(\Delta_{\mathbf{g}} U(\tilde{\rho}_t^s) - s \Delta_{\mathbf{g}} \left(U'(\tilde{\rho}_t^s) \nabla \cdot \left(\tilde{\rho}_t^s \nabla \tilde{\phi}_t^s \right) \right) \right) \tilde{\phi}_t^s \, dV \\ - s \int_{\mathbb{M}} \Delta_{\mathbf{g}} U(\tilde{\rho}_t^s) |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \, dV + \frac{s}{2} \int_{\mathbb{M}} \Delta_{\mathbf{g}} \left(|\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \right) U(\tilde{\rho}_t^s) \, dV \\ = \int_{\mathbb{M}} U(\tilde{\rho}_t^s) \Delta_{\mathbf{g}} \tilde{\phi}_t^s \, dV - s \int_{\mathbb{M}} \left(\left\langle \nabla U(\tilde{\rho}_t^s), \nabla \tilde{\phi}_t^s \right\rangle_{\mathbf{g}} \Delta_{\mathbf{g}} \tilde{\phi}_t^s + \tilde{\rho}_t^s U'(\tilde{\rho}_t^s) \left(\Delta_{\mathbf{g}} \tilde{\phi}_t^s \right)^2 \right) \, dV \\ - \frac{s}{2} \int_{\mathbb{M}} \Delta_{\mathbf{g}} \left(|\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \right) U(\tilde{\rho}_t^s) \, dV \\ = - \int_{\mathbb{M}} \left\langle \nabla U(\tilde{\rho}_t^s), \nabla \tilde{\phi}_t^s \right\rangle_{\mathbf{g}} \, dV + s \int_{\mathbb{M}} \left[-\frac{1}{2} \Delta_{\mathbf{g}} \left(|\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \right) + \left\langle \nabla \tilde{\phi}_t^s, \nabla \Delta_{\mathbf{g}} \tilde{\phi}_t^s \right\rangle_{\mathbf{g}} \right] U(\tilde{\rho}_t^s) \, dV \\ (4.14) \quad + s \int_{\mathbb{M}} \left(\Delta_{\mathbf{g}} \tilde{\phi}_t^s \right)^2 \left(U(\tilde{\rho}_t^s) - \tilde{\rho}_t^s U'(\tilde{\rho}_t^s) \right) \, dV$$

Applying Bochner formula:

$$(4.15) \quad \left\langle \nabla \phi, \nabla \Delta_{\mathbf{g}} \phi \right\rangle_{\mathbf{g}} - \frac{1}{2} \Delta_{\mathbf{g}} \left(|\nabla \phi|_{\mathbf{g}}^2 \right) = -|\text{Hess } \phi|_{\mathbf{g}}^2 - \text{Ric}_{\mathbf{g}}(\nabla \phi, \nabla \phi),$$

we get

$$(4.16) \quad \frac{\partial}{\partial t} \frac{1}{2} \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \tilde{\rho}_t^s \, dV + \int_{\mathbb{M}} \left\langle \nabla U(\tilde{\rho}_t^s), \nabla \tilde{\phi}_t^s \right\rangle_{\mathbf{g}} \, dV = s \tilde{D}_t^s.$$

Now we observe that the second term in the right-hand side of (4.16) is the derivative of the functional (4.3) along the curve $s \mapsto \tilde{\rho}_t^s V \in \mathcal{P}_2^{gr}(\mathbb{M})$:

$$(4.17) \quad \frac{\partial}{\partial s} \mathcal{E}(\tilde{\mu}_t^s) = \int_{\mathbb{M}} e'(\tilde{\rho}_t^s) \frac{\partial}{\partial s} \tilde{\rho}_t^s \, dV = - \int_{\mathbb{M}} e'(\tilde{\rho}_t^s) \nabla \cdot \left(\tilde{\rho}_t^s \nabla \tilde{\phi}_t^s \right) \, dV = \int_{\mathbb{M}} \nabla U(\tilde{\rho}_t^s) \cdot \nabla \tilde{\phi}_t^s \, dV,$$

and we eventually obtain (4.8).

Finally, when $\text{Ric}_g(\mathbb{M}) \geq 0$, using the inequality $(\Delta_g \phi)^2 \leq n|\text{Hess } \phi|_g^2$ and (1.7) we easily get $\tilde{D}_t^s \leq 0$ and (4.9). \square

Proof of Theorem 4.2. We argue as in the proof of Theorem 2.2: we fix $\varepsilon > 0$ and we choose a smooth curve $(\rho, \phi) \in \mathcal{C}(\nu, \mu)$ such that

$$(4.18) \quad \int_0^1 \tilde{A}_0^s \, ds = \int_0^1 \int_{\mathbb{M}} |\nabla \phi^s|_g^2 \rho^s \, dV \, ds \leq W_2^2(\nu, \mu) + \varepsilon.$$

Let $(\tilde{\rho}, \tilde{\phi})$ be a smooth variation defined as in Lemma 4.3; since $\tilde{\rho}_t^0 V = \rho^0 V = \nu$ and $\tilde{\rho}_t^1 V = \mu_t$, for every $t > 0$ we have $(\tilde{\rho}_t^s, \tilde{\phi}_t^s) \in \mathcal{C}(\nu, \mu_t)$, and therefore

$$(4.19) \quad W_2^2(\nu, \mu_t) \leq \int_0^1 \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s|_g^2 \tilde{\rho}_t^s \, dV \, ds = \int_0^1 \tilde{A}_t^s \, ds.$$

Integrating (4.9) for $s \in [0, 1]$ and $t \in [0, \tau]$ and recalling that $t \mapsto \mathcal{E}(\mu_t)$ is not increasing, we get

$$(4.20) \quad \frac{1}{2} \int_0^1 \tilde{A}_\tau^s \, ds - \frac{1}{2} \int_0^1 \tilde{A}_0^s \, ds \leq \tau (\mathcal{E}(\nu) - \mathcal{E}(\mu_\tau)).$$

Combining (4.20) with (4.19) and (4.18) we get

$$(4.21) \quad \frac{1}{2} W_2^2(\nu, \mu_\tau) - \frac{1}{2} W_2^2(\nu, \mu) \leq \tau (\mathcal{E}(\nu) - \mathcal{E}(\mu_\tau)) + \varepsilon,$$

and, as ε is arbitrary,

$$(4.22) \quad \frac{1}{2} W_2^2(\nu, \mu_\tau) - \frac{1}{2} W_2^2(\nu, \mu) \leq \tau (\mathcal{E}(\nu) - \mathcal{E}(\mu_\tau)).$$

Since the semigroup associated with (4.1) is translation invariant, (4.22) is the integral formulation (3.3) of (4.4). \square

Remark 4.4. Taking into account Theorem 2.3, (4.8) perfectly fits with the calculation performed by [19, Lemma 4.4], which provides the same expression for \tilde{D}_t^s . Applying now Theorem 3.3, with the choices $X := \mathcal{P}_2(\mathbb{M})$, $X_0 := \mathcal{P}_2^{ar}(\mathbb{M})$, and $F := \mathcal{E}$ (which satisfies the approximation condition (3.18); see [3]), we can prove the first part of Theorem 1.1.

COROLLARY 4.5. *Let $\mathcal{E} : \mathcal{P}_2(\mathbb{M}) \rightarrow (-\infty, +\infty]$ be the functional defined in (1.1). If e satisfies McCann conditions (1.7) and $\text{Ric}_g(\mathbb{M}) \geq 0$, then \mathcal{E} is (strongly) displacement convex along every geodesic $\mu : s \in [0, 1] \mapsto \mu^s \in \mathcal{P}_2(\mathbb{M})$, i.e.,*

$$(4.23) \quad \mathcal{E}(\mu^s) \leq (1 - s)\mathcal{E}(\mu^0) + s\mathcal{E}(\mu^1) \quad \forall s \in [0, 1].$$

5. The heat equation and the displacement λ -convexity of the logarithmic entropy. In this last section we prove the second part of Theorem 1.1: we thus assume that the Riemannian manifold \mathbb{M} satisfies the lower Ricci curvature bound

$$(5.1) \quad \text{Ric}_g(\mathbb{M}) \geq \lambda \quad \text{i.e.,} \quad \text{Ric}_{g_x}(\xi, \xi) \geq \lambda |\xi|_g^2 \quad \forall \xi \in T_x \mathbb{M},$$

and we consider the logarithmic entropy functional

$$(5.2) \quad \mathcal{E}(\mu) = \int_{\mathbb{M}} \rho \log \rho \, dV, \quad \rho = \frac{d\mu}{dV},$$

corresponding to $e(\rho) := \rho \log \rho$. Since $U(\rho) = \rho$, the Wasserstein gradient flow associated to \mathcal{E} is the Heat equation

$$(5.3) \quad \frac{\partial}{\partial t} \rho_t - \Delta_{\mathbf{g}} \rho_t = 0 \quad \text{in } \mathbb{M}, \quad \rho|_{t=0} = \rho_0.$$

The main result of this section is the following.

THEOREM 5.1. *The semigroup $\mathcal{S}_t : \mu_0 = \rho_0 \mathbf{V} \mapsto \mu_t = \rho_t \mathbf{V}$, generated by the solution of the Heat equation (5.3), is a λ -flow in $\mathcal{P}_2^{ar}(\mathbb{M})$ for the logarithmic entropy functional; i.e., μ_t satisfies the inequality*

$$(5.4) \quad \frac{1}{2} \frac{d^+}{dt} W_2^2(\nu, \mu_t) + \frac{\lambda}{2} W_2^2(\nu, \mu_t) \leq \mathcal{E}(\nu) - \mathcal{E}(\mu_t) \quad \forall t \in [0, +\infty), \nu \in \mathcal{P}_2^{ar}(\mathbb{M}).$$

In particular, the logarithmic entropy functional (5.2) is (strongly) displacement λ -convex; i.e., for every geodesic $\mu^s : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{M})$ between μ^0 and μ^1 , we have

$$(5.5) \quad \mathcal{E}(\mu^s) \leq (1-s)\mathcal{E}(\mu^0) + s\mathcal{E}(\mu^1) - \frac{\lambda}{2}s(1-s)W_2^2(\mu^0, \mu^1), \quad \forall s \in [0, 1].$$

Proof. By Theorem 3.3, if \mathcal{S} is a λ -flow for the functional (5.2) in $\mathcal{P}_2^{ar}(\mathbb{M})$, then \mathcal{E} is (strongly) displacement λ -convex. In order to prove that \mathcal{S} is a λ -flow, since (3.2a,b) are immediate, we check that \mathcal{S} satisfies the E.V.I. (3.2c), and we argue as in the proof of Theorem 4.2 and Theorem 2.2. We thus fix $\varepsilon > 0$, and we choose a smooth curve $(\rho, \phi) \in \mathcal{C}(\nu, \mu)$

$$(5.6) \quad \int_0^1 \tilde{A}_0^s ds = \int_0^1 \int_{\mathbb{M}} |\nabla \phi^s|_{\mathbf{g}}^2 \rho^s dV ds \leq W_2^2(\nu, \mu) + \varepsilon^2.$$

By a standard reparametrization technique (see next Lemma 5.1), we can also assume that

$$(5.7) \quad W_2(\mu^{s_0}, \mu^{s_1}) \leq L|s_0 - s_1|, \quad L^2 := W_2^2(\nu, \mu) + \varepsilon^2 \quad \forall s_0, s_1 \in [0, 1]; \quad \mu^s := \rho^s \mathbf{V}.$$

We keep the same notation of Theorem 4.2 and Lemma 4.3; i.e.,

$$(5.8) \quad \tilde{\mu}_t^s = \tilde{\rho}_t^s \mathbf{V} := \mathcal{S}_{st}(\mu^s), \quad \tilde{A}_t^s := \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \tilde{\rho}_t^s dV, \quad \tilde{F}_t^s = \mathcal{E}(\tilde{\mu}_t^s),$$

where $\tilde{\phi}_t^s$ is a family of potentials associated with $\tilde{\rho}_t^s$ as in (4.6). Since $U(\rho) = \rho$, the term $\rho U'(\rho) - U(\rho)$ in the definition of \tilde{D}_t^s vanishes, so that in the present case

$$(5.9) \quad \tilde{D}_t^s = - \int_{\mathbb{M}} \left(|\text{Hess } \tilde{\phi}_t^s|_{\mathbf{g}}^2 + \text{Ric}_{\mathbf{g}} \left(\nabla \tilde{\phi}_t^s, \nabla \tilde{\phi}_t^s \right) \right) \tilde{\rho}_t^s dV \stackrel{(5.1)}{\leq} -\lambda \int_{\mathbb{M}} |\nabla \tilde{\phi}_t^s|_{\mathbf{g}}^2 \tilde{\rho}_t^s dV = -\lambda \tilde{A}_t^s,$$

and (4.8) yields the differential inequality

$$(5.10) \quad \frac{1}{2} \frac{\partial}{\partial t} \tilde{A}_t^s + \lambda s \tilde{A}_t^s + \frac{\partial}{\partial s} \tilde{F}_t^s \leq 0 \quad \forall s \in [0, 1], \quad \forall t > 0.$$

Multiplying inequality (5.10) by $e^{2\lambda st} > 0$ we obtain

$$(5.11) \quad \frac{1}{2} \frac{\partial}{\partial t} \left(e^{2\lambda st} \tilde{A}_t^s \right) + \frac{\partial}{\partial s} \left(e^{2\lambda st} \tilde{F}_t^s \right) \leq 2\lambda t e^{2\lambda st} \tilde{F}_t^s.$$

Integrating with respect to s from 0 to 1 we get

$$(5.12) \quad \frac{d}{dt} \left(\frac{1}{2} \int_0^1 e^{2\lambda st} \tilde{A}_t^s ds \right) + e^{2\lambda t} \tilde{F}_t^1 - \tilde{F}_t^0 \leq \int_0^1 2\lambda t e^{2\lambda st} \tilde{F}_t^s ds,$$

and a further integration with respect to t yields

$$(5.13) \quad \frac{1}{2} \int_0^1 e^{2\lambda st} \tilde{A}_t^s ds - \frac{1}{2} \int_0^1 A_0^s ds + E_{2\lambda}(t)\mathcal{E}(\mu_t) - t\mathcal{E}(\nu) \leq \int_0^t \int_0^1 2\lambda r e^{2\lambda sr} \tilde{F}_r^s ds dr.$$

Applying the next Lemma 5.1, since for $\lambda \neq 0$ $\int_0^1 \frac{1}{e^{2\lambda st}} ds = \frac{1-e^{-2\lambda t}}{2\lambda t} = \frac{1}{e^{\lambda t} \mathfrak{s}(\lambda t)}$, $\mathfrak{s}(t) := \frac{t}{\sinh(t)}$, we get

$$(5.14) \quad \begin{aligned} & \frac{e^{\lambda t} \mathfrak{s}(\lambda t)}{2} W_2^2(\mu_t, \nu) - \frac{1}{2} W_2^2(\mu, \nu) + E_{2\lambda}(t)\mathcal{E}(\mu_t) - t\mathcal{E}(\nu) \\ & \leq \int_0^t \int_0^1 2\lambda r e^{2\lambda sr} \tilde{F}_r^s ds dr + \varepsilon^2/2. \end{aligned}$$

Let us first consider the case $\lambda \leq 0$: being \mathcal{E} nonnegative, the right-hand side in (5.14) is less or equal than ε ; since $\varepsilon > 0$ is arbitrary, we obtain the same inequality with 0 in the right-hand side. Since $t^{-1}E_{2\lambda}(t) \rightarrow 1$ as $t \downarrow 0$ and $\mathfrak{s}(0) = 1$, we thus obtain

$$(5.15) \quad \frac{1}{2} \frac{d^+}{dt} \left(e^{\lambda t} \mathfrak{s}(\lambda t) W_2^2(\mu_t, \nu) \right) \Big|_{t=0} + \mathcal{E}(\mu) \leq \mathcal{E}(\nu).$$

Being $\mathfrak{s}'(0) = 0$ it is then easy to check that

$$\frac{d^+}{dt} \left(e^{\lambda t} \mathfrak{s}(\lambda t) W_2^2(\mu_t, \nu) \right) \Big|_{t=0} = \frac{d^+}{dt} \left(W_2^2(\mu_t, \nu) \right) \Big|_{t=0} + \lambda W_2^2(\mu, \nu),$$

which yields (5.4).

Let us now consider the case $\lambda > 0$. Notice that we already know that \mathcal{S} is a 0-flow; by (5.7) we can apply the estimate (3.11) with $\lambda = 0$ obtaining

$$\begin{aligned} r \tilde{F}_r^s &= r \mathcal{E}(\mathcal{S}_{rs}(\mu^s)) \stackrel{(3.11)}{\leq} r \left((1-s)\mathcal{E}(\mu^0) + s\mathcal{E}(\mu^1) + \frac{\varepsilon^2}{2rs} s(1-s) \right) \\ &\leq r (\mathcal{E}(\mu^0) + \mathcal{E}(\mu^1)) + \varepsilon^2/2, \end{aligned}$$

since $s \in [0, 1]$. We thus get

$$(5.16) \quad \int_0^t \int_0^1 2\lambda r e^{2\lambda sr} \tilde{F}_r^s ds dr \leq \lambda t e^{2\lambda t} \left(t(\mathcal{E}(\mu_0) + \mathcal{E}(\mu_1)) + \varepsilon^2 \right);$$

inserting this bound in (5.14) and passing to the limit as $\varepsilon \downarrow 0$ we find

$$(5.17) \quad \frac{e^{\lambda t} \mathfrak{s}(\lambda t)}{2} W_2^2(\mu_t, \nu) - \frac{1}{2} W_2^2(\mu, \nu) + E_{2\lambda}(t)\mathcal{E}(\mu_t) - t\mathcal{E}(\nu) \leq \lambda t^2 e^{2\lambda t} (\mathcal{E}(\mu_0) + \mathcal{E}(\mu_1)).$$

Dividing by t and letting t tend to 0, the second term vanishes, so we obtain the E.V.I. also in the case in which $\lambda > 0$. \square

LEMMA 5.1. *Let $\nu, \mu \in \mathcal{P}_2^{ar}(\mathbb{M})$ and let $(\rho, \phi) \in \mathcal{C}(\nu, \mu)$ be a smooth solution of the continuity equation*

$$\frac{\partial}{\partial s} \rho^s + \nabla \cdot (\rho^s \nabla \phi^s) = 0 \quad \text{in } [0, 1] \times \mathbb{M} \quad \text{with}$$

$$\rho^0 \mathbf{V} = \nu, \quad \rho^1 \mathbf{V} = \mu, \quad \text{and} \quad A^s := \int_{\mathbb{M}} |\nabla \phi^s|_{\mathbb{g}}^2 \rho^s \, dV.$$

For every positive function $f \in C^\infty[0, 1]$

$$(5.18) \quad W_2^2(\nu, \mu) \leq L_f \int_0^1 f(s) A^s \, ds, \quad \text{where} \quad L_f := \int_0^1 \frac{1}{f(s)} \, ds.$$

Moreover, for every $\varepsilon > 0$ there exists a smooth rescaling $\mathbf{s}_\varepsilon : [0, 1] \rightarrow [0, 1]$ so that the reparametrized families

$$(5.19) \quad \bar{\rho}^r := \rho^{\mathbf{s}_\varepsilon(r)}, \quad \bar{\phi}^r := \mathbf{s}'_\varepsilon(r) \phi^{\mathbf{s}_\varepsilon(r)}, \quad \bar{\mu}^r := \bar{\rho}^r \mathbf{V}$$

satisfy

$$(5.20) \quad (\bar{\rho}, \bar{\phi}) \in \mathcal{C}(\nu, \mu), \quad W_2(\bar{\mu}^{r_0}, \bar{\mu}^{r_1}) \leq L|r_0 - r_1|, \quad L^2 \leq \int_0^1 A^s \, ds + \varepsilon^2.$$

Proof. Let us consider the smooth increasing map $r : [0, 1] \rightarrow [0, 1]$

$$r(s) := L_f^{-1} \int_0^s \frac{1}{f(s)} \, ds \quad \text{and its inverse } \mathbf{s} := r^{-1} \quad \text{with} \quad \mathbf{s}'(r(s)) = L_f f(s).$$

It is immediate to check that the smooth (reparametrized) curve

$$(5.21) \quad \bar{\rho}^r(x) := \rho^{\mathbf{s}(r)}(x), \quad \bar{\phi}^r(x) := \mathbf{s}'(r) \phi^{\mathbf{s}(r)}(x)$$

belongs to $\mathcal{C}(\nu, \mu)$. It follows that

$$W_2^2(\nu, \mu) \leq \int_0^1 \bar{A}^r \, dr, \quad \text{where} \quad \bar{A}^r := \int_{\mathbb{M}} |\nabla \bar{\phi}^r|_{\mathbb{g}}^2 \bar{\rho}^r \, dV \stackrel{(5.21)}{=} (\mathbf{s}'(r))^2 A^{\mathbf{s}(r)},$$

so that

$$\int_0^1 \bar{A}^r \, dr = \int_0^1 A^{\mathbf{s}(r)} (\mathbf{s}'(r))^2 \, dr = \int_0^1 A^s \mathbf{s}'(r(s)) \, ds = L_f \int_0^1 f(s) A^s \, ds.$$

Choosing now the reparametrization \mathbf{s}_ε corresponding to the choice

$$(5.22) \quad f_\varepsilon(s) := \frac{1}{\sqrt{\varepsilon^2 + A^s}}, \quad L_{f_\varepsilon} := \int_0^1 \sqrt{\varepsilon^2 + A^s} \, ds, \quad L_{f_\varepsilon}^2 \leq \varepsilon^2 + \int_0^1 A^s \, ds,$$

we get

$$\begin{aligned} W^2(\bar{\mu}^{r_0}, \bar{\mu}^{r_1}) &\leq |r_1 - r_0| \int_{r_0}^{r_1} \bar{A}^r \, dr = |r_1 - r_0| L_{f_\varepsilon}^2 \int_{r_0}^{r_1} A^{\mathbf{s}(r)} f_\varepsilon^2(\mathbf{s}(r)) \, dr \\ &\leq (r_1 - r_0)^2 L_{f_\varepsilon}^2, \end{aligned}$$

which yields (5.20). \square

REFERENCES

- [1] M. AGUEH, *Existence of solutions to degenerate parabolic equations via the Monge-Kantorovich theory*, Adv. Differential Equations, 10 (2005), pp. 309–360.
- [2] M. AGUEH, N. GHOUSSOUB, AND X. KANG, *Geometric inequalities via a general comparison principle for interacting gases*, Geom. Funct. Anal., 14 (2004), pp. 215–244.
- [3] L. AMBROSIO AND G. BUTTAZZO, *Weak lower semicontinuous envelope of functionals defined on a space of measures*, Ann. Mat. Pura Appl., 150 (1988), pp. 311–339.
- [4] L. AMBROSIO, N. GIGLI, AND G. SAVARÉ, *Gradient flows in metric spaces and in the space of probability measures*, Lectures Math. ETH Zürich, Birkhäuser Verlag, Basel, 2005.
- [5] L. AMBROSIO AND G. SAVARÉ, *Gradient flows of probability measures*, in Handbook of Evolution Equations (III), Elsevier, New York, 2006.
- [6] J.-D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem*, Numer. Math., 84 (2000), pp. 375–393.
- [7] J. A. CARRILLO, S. LISINI, AND G. SAVARÉ, *The porous medium flow and generalized displacement convexity*, Technical report, in preparation (2008).
- [8] J. A. CARRILLO, R. J. MCCANN, AND C. VILLANI, *Contractions in the 2-Wasserstein length space and thermalization of granular media*, Arch. Ration. Mech. Anal., 179 (2006), pp. 217–263.
- [9] D. CORDERO-ERAUSQUIN, R. J. MCCANN, AND M. SCHMUCKENSCHLÄGER, *A Riemannian interpolation inequality à la Borell, Brascamp and Lieb*, Invent. Math., 146 (2001), pp. 219–257.
- [10] D. CORDERO-ERAUSQUIN, R. J. MCCANN, AND M. SCHMUCKENSCHLÄGER, *Prékopa-Leindler type inequalities on Riemannian manifolds, Jacobi fields, and optimal transport*, Ann. Fac. Sci. Toulouse Math. (6), 15 (2006), pp. 613–635.
- [11] J. DOLBEAULT, B. NAZARET, AND G. SAVARÉ, *A new class of “dynamic” transport distances between measures*, Calc. Var. Partial Differential Equations, to appear.
- [12] A. FIGALLI AND C. VILLANI, *Strong displacement convexity on Riemannian manifolds*, Math. Z., 257 (2007), pp. 251–259.
- [13] S. LISINI, *Characterization of absolutely continuous curves in Wasserstein spaces*, Calc. Var. Partial Differential Equations, 28 (2007), pp. 85–120.
- [14] J. LOTT AND C. VILLANI, *Ricci curvature for metric-measure spaces via optimal transport*, Ann. of Math. (2), to appear.
- [15] R. J. MCCANN, *A convexity principle for interacting gases*, Adv. Math., 128 (1997), pp. 153–179.
- [16] R. J. MCCANN, *Polar factorization of maps on Riemannian manifolds*, Geom. Funct. Anal., 11 (2001), pp. 589–608.
- [17] F. OTTO, *The geometry of dissipative evolution equations: The porous medium equation*, Comm. Partial Differential Equations, 26 (2001), pp. 101–174.
- [18] F. OTTO AND C. VILLANI, *Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality*, J. Funct. Anal., 173 (2000), pp. 361–400.
- [19] F. OTTO AND M. WESTDICKENBERG, *Eulerian calculus for the contraction in the Wasserstein distance*, SIAM J. Math. Anal., 37 (2005), pp. 1227–1255.
- [20] K.-T. STURM, *On the geometry of metric measure spaces. I*, Acta Math., 196 (2006), pp. 65–131.
- [21] K.-T. STURM, *On the geometry of metric measure spaces. II*, Acta Math., 196 (2006), pp. 133–177.
- [22] J. L. VÁZQUEZ, *The Porous Medium Equation*, Oxford Mathematical Monographs, The Clarendon Press Oxford University Press, Oxford, 2007. Mathematical theory.
- [23] C. VILLANI, *Optimal Transport, Old and New*, Springer-Verlag, 2009, to appear.
- [24] M.-K. VON RENESSE AND K.-T. STURM, *Transport inequalities, gradient estimates, entropy, and Ricci curvature*, Comm. Pure Appl. Math., 58 (2005), pp. 923–940.

POSITIVITY PROPERTIES OF THE FOURIER TRANSFORM AND THE STABILITY OF PERIODIC TRAVELLING-WAVE SOLUTIONS*

JAIME ANGULO PAVA[†] AND FÁBIO M. A. NATALI[‡]

Abstract. In this paper we establish a method to obtain the stability of periodic travelling-wave solutions for equations of Korteweg–de Vries-type $u_t + u^p u_x - M u_x = 0$, with M being a general pseudodifferential operator and where $p \geq 1$ is an integer. Our approach uses the theory of totally positive operators, the Poisson summation theorem, and the theory of Jacobi elliptic functions. In particular we obtain the stability of a family of periodic travelling waves solutions for the Benjamin–Ono equation. The present technique gives a new way to obtain the existence and stability of cnoidal and dnoidal waves solutions associated with the Korteweg–de Vries and modified Korteweg–de Vries equations, respectively. The theory has prospects for the study of periodic travelling-wave solutions of other partial differential equations.

Key words. dispersive equations, Korteweg–de Vries-type equations, periodic travelling waves, Jacobi elliptic functions, nonlinear stability

AMS subject classifications. 76B25, 35Q51, 35Q53

DOI. 10.1137/080718450

1. Introduction. One of the main properties of dispersive nonlinear evolution equations is that usually they sustain steadily translating waves called travelling waves. These solutions imply a balance between the effects of nonlinearity and of frequency dispersion. By depending on specific boundary conditions on the wave's shape, for instance, in the case of water waves, these special states of motion can arise either solitary or periodic waves. The study of these special steady waveforms is essential to the explanation of many wave phenomena observed in the practice, for instance, in surface water waves propagating in a canal, in propagation of internal waves, or in shallow-water ocean surface waves (see Benjamin [12], [13], [14] and Osborne et al. [44]).

In the water wave context, Constantin in [21] and Constantin and Escher in [22] analyzed a free boundary problem for harmonic functions and showed that periodic or solitary travelling waves possess stability properties within the shallow-water regime (see also Toland [47] and Constantin and Strauss [23] and the citations therein). Moreover, various nonlinear dispersive model equations are an accurate approximation to the governing equations for water waves (see [5]). From these considerations, questions about the stability of travelling waves and their existence as exact solutions of the dynamical equations are very important.

The solitary waves are in general single crested, symmetric, localized travelling waves, whose hyperbolic sech profiles are well known (see Ono [43] and Benjamin [16] for the existence of solitary waves of algebraic type or with a finite number of oscillations). The study of the nonlinear stability or instability in the form of solitary waves has had a terrific development and refinement in recent years. The proofs have

*Received by the editors March 13, 2008; accepted for publication (in revised form) June 17, 2008; published electronically October 22, 2008.

<http://www.siam.org/journals/sima/40-3/71845.html>

[†]Department of Mathematics, IME-USP Rua do Matão 1010, CEP 05508-090, São Paulo - SP - Brazil (angulo@ime.usp.br). This author's research was partially supported by CNPq/Brazil.

[‡]Department of Mathematics, Universidade Estadual de Maringá Avenida Colombo, 5790, CEP 87020-900, Maringá - PR - Brazil (fmanatali@uem.br). This author's research was supported by FAPESP/Brazil under grant 06/61310-3.

been simplified and sufficient conditions were obtained to ensure the stability to small localized perturbations in the waveform. Those conditions have shown to be effective in a variety of circumstances; see, for example, [2], [3], [4], [14], [17], [27], [28], [49], and [48].

The situation regarding periodic travelling waves is very different. The stability and the existence of explicit formulas of these progressive wavetrains have received comparatively little attention. A first study of these waveforms was determined by Benjamin in [16] with regard to the periodic steady solutions called *cnoidal waves*, which were found initially by Korteweg and de Vries in [34] for the equation currently called the Korteweg–de Vries equation (KdV henceforth):

$$(1.1) \quad u_t + uu_x + u_{xxx} = 0,$$

where $u = u(x, t)$ is a real-valued function of the two variables $x, t \in \mathbb{R}$. Benjamin put forward an approach to the stability of cnoidal waves in the form

$$\varphi(\xi) = \beta_2 + (\beta_3 - \beta_2)\text{cn}^2\left(\sqrt{\frac{\beta_3 - \beta_1}{12}}\xi; k\right),$$

but did not provide a detailed justification of his assertions, and several aspects seem problematic. The first result of stability for periodic solutions of the KdV was obtained by McKean in [37], who considered the orbital stability of all periodic finite-genus solutions with respect to perturbations of the same period. McKean's approach was based on the integrable structure of the KdV. More recently Angulo, Bona, and Scialom in [10] returned to Benjamin's original question and gave a complete theory of stability of cnoidal waves for (1.1) with respect to perturbations of the same period (see also [7]). The approach for obtaining this result was based on the ideas of Bona, Weinstein, and Grillakis, Shatah, and Strauss (see [17], [27], [48]) but adapted to the periodic context. So new theories of stability for other dispersive equations such as the focusing nonlinear Schrödinger equation and the modified KdV has been obtained (see Angulo [9], [8]). It is remarkable to see that in all these works the use of an elaborated spectral theory for the periodic eigenvalue problem was necessary (see [29], [36]),

$$(1.2) \quad \begin{cases} \frac{d^2}{dx^2}\Psi + [\rho - n(n+1)k^2\text{sn}^2(x; k)]\Psi = 0, \\ \Psi(0) = \Psi(2K(k)), \quad \Psi'(0) = \Psi'(2K(k)), \end{cases}$$

with specific values of $n \in \mathbb{N}$. In (1.2), $\text{sn}(\cdot; k)$ represents the Jacobi elliptic function of type snoidal with *modulus* k , $k \in (0, 1)$, and K represents the complete elliptic integral of the first kind defined by

$$K(k) = \int_0^1 \frac{dt}{\sqrt{(1-t^2)(1-k^2t^2)}}.$$

We recall that the second order differential equation in (1.2) is known as the *Jacobi form of the Lamé equation*.

We note also that Gardner in [26] provided a theory for determining that the large wavelength periodic waves are linearly unstable whenever the limiting homoclinic wave (solitary wave) is unstable. He applied it for diverse types of nonlinear evolution equations in one space variable, in the case of the generalized KdV equations

$$u_t + u^p u_x + u_{xxx} = 0,$$

$p \in \mathbb{N}$, and assuming that this equation admits a family of large wavelength periodic waves U^α such that the period T_α tends to infinity as α tends to zero, then they are unstable whenever $p > 4$ and $\alpha > 0$ is sufficiently small.

Then, the main focus in this paper will be the study of the existence and stability of periodic travelling-wave solutions for equations of the form

$$(1.3) \quad u_t + u^p u_x - M u_x = 0,$$

where $p \geq 1$ is an integer and M is a differential or pseudodifferential operator in the framework of periodic functions. M is defined as a Fourier multiplier operator by

$$(1.4) \quad \widehat{Mg}(k) = \alpha(k)\widehat{g}(k), \quad k \in \mathbb{Z},$$

where the symbol α of M is assumed to be a measurable, locally bounded, even function on \mathbb{R} , satisfying the conditions

$$(1.5) \quad A_1|k|^{m_1} \leq \alpha(k) \leq A_2(1 + |k|)^{m_2}$$

for $m_1 \leq m_2$, $|k| \geq k_0$, $\alpha(k) > b$ for all $k \in \mathbb{Z}$ and $A_i > 0$. The travelling-wave solutions for (1.3) will have the form

$$u(x, t) = \varphi_c(x - ct),$$

where the profile φ_c is a smooth periodic function with an a priori fundamental period $2L$, $L > 0$. Hence substituting this form of u into (1.3) and integrating once (with the integration constant being considered zero throughout our theory), one obtains that $\varphi = \varphi_c$ is the solution of the equation

$$(1.6) \quad (M + c)\varphi - \frac{1}{p+1}\varphi^{p+1} = 0.$$

Associated with (1.6) we consider the linear, closed, unbounded, self-adjoint operator $\mathcal{L} : D(\mathcal{L}) \rightarrow L^2_{per}([-L, L])$ defined on a dense subspace of $L^2_{per}([-L, L])$ by

$$(1.7) \quad \mathcal{L}u = (M + c)u - \varphi^p u.$$

From the theory of compact symmetric operators applied to the periodic eigenvalue problem

$$(1.8) \quad \begin{cases} \mathcal{L}\psi &= \lambda\psi, \\ \psi(-L) &= \psi(L), \quad \psi'(-L) = \psi'(L), \end{cases}$$

it is possible to see that the spectrum of \mathcal{L} is a countable infinite set of eigenvalues, $\{\lambda_n\}$, with

$$(1.9) \quad \lambda_0 \leq \lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \dots,$$

where $\lambda_n \rightarrow \infty$ as $n \rightarrow \infty$ (see Proposition 3.1 below for a proof of this assertion). In particular, from (1.6) we obtain that \mathcal{L} has zero as an eigenvalue with eigenfunction

$d\varphi/dx$. As is well known this property of \mathcal{L} is deduced from the invariance of the solutions of (1.3) by translations.

A set of sharp conditions is available in the literature to imply the stability of the orbit generated by φ_c , namely, $\Omega_{\varphi_c} = \{\varphi_c(\cdot + y) : y \in \mathbb{R}\}$. So, we say that Ω_{φ_c} (or φ_c) is stable in $H_{per}^{m_2}([-L, L])$ by the periodic flow generated by (1.3) if, for any $\epsilon > 0$, there exists a $\delta > 0$ such that, for $u_0 \in H_{per}^{m_2}([-L, L])$ with $d(u_0, \Omega_{\varphi_c}) \equiv \inf_{y \in \mathbb{R}} \|u_0 - \varphi_c(\cdot + y)\|_{H_{per}^{m_2}} < \delta$, the solution u of (1.3) with $u(x, 0) = u_0$ is global in time and satisfies $d(u(t), \Omega_{\varphi_c}) < \epsilon$ for all $t \in \mathbb{R}$. Thus, from [14], [17], [49], [27] the conditions that imply stability are the following:

(1.10)

(P_0) there is a nontrivial smooth curve of periodic solutions for (1.6) of the form

$$c \in I \subseteq \mathbb{R} \rightarrow \varphi_c \in H_{per}^{m_2}([-L, L]);$$

(P_1) \mathcal{L} has an unique negative eigenvalue λ , and it is simple;

(P_2) the eigenvalue 0 is simple;

$$(P_3) \frac{d}{dc} \int_{-L}^L \varphi_c^2(x) dx > 0.$$

The problem about the existence of a nontrivial smooth curve of periodic solutions in the form required by (P_0) above presents new and delicate aspects that need to be handled. The possibility of finding explicit solutions for (1.6) will depend naturally on the form of M . If it is a differential operator of the form $M = -\partial_x^2$, then we use the quadrature method (it means writing (1.6) in the form $[\varphi'_c]^2 = F(\varphi_c)$), and the theory of elliptic functions has shown to be a main tool. So, the solutions will depend on the Jacobi elliptic functions of *snoidal*, *cnoidal*, and *dnoidal* types (see [9], [8], [10], [15]). Now, since the period of these functions depends on the complete elliptic integral $K(k)$, we have that the elliptic modulus k will depend on the velocity c , and therefore we have that a priori the period of φ_c will depend on c . Hence, by using the implicit function theorem, the wanted smooth branch of periodic solutions with a fixed minimal period has been obtained in many cases. We note that the procedure of the quadrature method in general does not work if M is a pseudodifferential operator such as the nonlocal operator $\mathcal{H}\partial_x$, with \mathcal{H} being the Hilbert transform. In this paper we will make a different approach to obtain explicit solutions of (1.6) for a specific form of M and values of p . This approach will be based on the classical *Poisson summation theorem* (see [38], [45], [46]). At least two important advantages of this approach can be obtained: The first one is that it can be used for obtaining solutions when M is a pseudodifferential operator, for example, in the case of $\mathcal{H}\partial_x$. The other one is that related with computing the integral in (1.10). In general obtaining property (P_3) can be very difficult in the periodic case, as the results that have appeared in the literature have shown, since the use of nontrivial identities for the complete elliptic integrals of the first and the second kinds sometimes come on the scene as a fundamental piece in the analysis. As we will see our method to obtain property (P_3) can be very easy as a combination of the *Poisson summation theorem* and the *Parseval theorem*.

With regard to conditions (P_1) and (P_2), the problem is very delicate. One of the most remarkable results in the theory of stability of solitary wave solutions was given by Albert [2] and Albert and Bona [3], where sufficient conditions to obtain properties (P_1) and (P_2) were given. The advantage of that approach is that it does not require an explicit computation of the spectrum of the linear operator (1.7), since (P_1) and

(P_2) are obtained exclusively from positivity properties of the Fourier transform of the solitary wave in question. The present paper establishes an extension of the theory in [2] and [3] in the case of periodic travelling-wave solutions. The periodic problem has new points not encountered when considering issues related to the solitary waves. Our analysis also relies upon the theory of totally positive operators, and so the class $PF(2)$ defined by Karlin in [32] (see also [2]) is basic in our study.

Our theory leads to a significant simplification of some recent proofs of stability of periodic travelling-wave solutions of KdV-type equations (see [9], [7], [10]) such as in the case of the KdV and the modified KdV equations, since in those cases the verification of properties (P_1) and (P_2) required the determination of the instability intervals associated with the Lamé equation in (1.2) and of an explicit formula of at least the first three eigenvalues ρ (see [36]). *Our analysis does not need this information.*

Such as will be shown in section 4, this method establishes the first result about the stability of periodic travelling-wave solutions found by Benjamin in [15] for the Benjamin–Ono equation

$$(1.11) \quad u_t + uu_x - \mathcal{H}u_{xx} = 0,$$

where \mathcal{H} denotes the periodic Hilbert transform defined by

$$(1.12) \quad \mathcal{H}f(x) = \frac{1}{2L} \text{p.v.} \int_{-L}^L \cot g \left[\frac{\pi(x-y)}{2L} \right] f(y) dy.$$

The associated periodic waves for (1.11) with a minimal period $2L$ are given for $c > \frac{\pi}{L}$ as

$$(1.13) \quad \varphi_c(x) = \frac{2\pi}{L} \left(\frac{\sinh(\gamma)}{\cosh(\gamma) - \cos\left(\frac{\pi x}{L}\right)} \right),$$

such that $\gamma > 0$ satisfies $\tanh(\gamma) = \frac{\pi}{cL}$.

It is important to note that the stability results now presented are obtained by periodic initial disturbances having *the same minimal period* of our periodic solutions. This method cannot be extended for obtaining stability results with more general periodic perturbations, for instance, by periodic disturbances of two times the minimal period of our periodic solutions. In section 6 we give an explanation of this fact.

The plan of this paper is as follows. The next section is devoted to describing briefly the notation that will be used and making a few preliminary remarks regarding periodic Sobolev spaces, the Poisson summation theorem, and some results of global well-posedness in the periodic case of the KdV, modified KdV, and the Benjamin–Ono equations. Sections 3 and 4 contain our full theory which relates positivity properties of periodic travelling-wave solutions to the stability theory of [27]. Applications of section 4 to specific periodic travelling waves are presented in section 5. In section 6 some comments about the $PF(2)$ property in the periodic case are established. Finally, the appendix contains some basic properties of the elliptic integrals which are relevant to the theory of section 5, and the proof of an inequality of Poincaré–Wintinger type for nonlocal operators is established.

2. Notation and preliminaries.

2.1. Function classes. Let Ω be an open set of the real line \mathbb{R} and $1 \leq p \leq \infty$; then $L^p(\Omega)$ is the usual Banach space of (equivalence classes of) real-

complex-valued Lebesgue measurable functions defined on Ω provided with the norm

$$\|f\|_p = \left(\int_{\Omega} |f(x)|^p dx \right)^{\frac{1}{p}}$$

if $1 \leq p < \infty$. When $p = \infty$ we have $\|f\|_{\infty} = \sup_{x \in \Omega} |f(x)|$. When $p = 2$ the Banach space $L^p(\Omega)$ is a Hilbert space with inner product defined by $(f, g)_2 = \int_{\Omega} f(x)\overline{g(x)}dx$, where $f, g \in L^2(\Omega)$. The L^2 -based Sobolev spaces of periodic functions are defined as follows [30]. If $\mathcal{P} = C_{per}^{\infty}$ denotes the collection of all of the functions $f : \mathbb{R} \rightarrow \mathbb{C}$ which are C^{∞} and periodic with period $2l > 0$, the collection \mathcal{P}' of all continuous linear functionals from \mathcal{P} into \mathbb{C} is the set of periodic distributions. If $\Psi \in \mathcal{P}'$, we denote the evaluation of Ψ at φ by $\Psi(\varphi) = \langle \Psi, \varphi \rangle$ for $\varphi \in \mathcal{P}$. For $k \in \mathbb{Z}$, let $\Theta_k(x) = e^{\frac{ik\pi x}{l}}$, $x \in \mathbb{R}$. The Fourier transform of $\Psi \in \mathcal{P}'$ is a function $\widehat{\Psi} : \mathbb{Z} \rightarrow \mathbb{C}$ defined by $\widehat{\Psi}(k) = \frac{1}{2l} \langle \Psi, \Theta_k \rangle$, $k \in \mathbb{Z}$. $\widehat{\Psi}(k)$ are called the Fourier coefficients of Ψ . As usual, a function $\psi \in L^p([-l, l])$, $p \geq 1$, is an element of \mathcal{P}' by defining

$$\langle \psi, \varphi \rangle = \frac{1}{2l} \int_{-l}^l \psi(x)\varphi(x)dx, \quad \varphi \in \mathcal{P}.$$

If $\psi \in L^p([-l, l])$ for some $p \geq 1$, then, for $k \in \mathbb{Z}$,

$$\widehat{\psi}(k) = \frac{1}{2l} \int_{-l}^l \psi(x)e^{-\frac{ik\pi x}{l}} dx.$$

The Fourier inverse transform of a sequence $\alpha = (\alpha_k)_{k \in \mathbb{Z}} \in \mathcal{S}(\mathbb{Z})$, where $\mathcal{S}(\mathbb{Z})$ denotes the space of the rapidly decreasing sequences, is the function $\check{\alpha} \in \mathcal{P}$ defined by $\check{\alpha}(x) = \sum_{k \in \mathbb{Z}} \alpha_k \Theta_k(x)$. We consider the space \mathcal{P}' provided with the usual weak-star topology, but it will not be needed here. We denote by C_{2l} the space of the continuous and $2l$ -periodic functions. Let $\alpha = (\alpha_k)_{k \in \mathbb{Z}}$ be a sequence of complex value. The Hilbert space $\ell^2 := \ell^2(\mathbb{Z})$ is defined by

$$\ell^2 = \left\{ \alpha; \|\alpha\|_{\ell^2} := \left(\sum_{k=-\infty}^{+\infty} |\alpha_k|^2 \right)^{\frac{1}{2}} < \infty \right\}.$$

For $s \in \mathbb{R}$, the Sobolev space $H_{per}^s([-l, l]) := H_{2l}^s$ is the set of all $f \in \mathcal{P}'$ such that

$$\|f\|_{H_{2l}^s}^2 \equiv 2l \sum_{k=-\infty}^{+\infty} (1 + |k|^2)^s |\widehat{f}(k)|^2 < \infty.$$

The collection H_{2l}^s is a Hilbert space with inner product

$$(f, g)_{H_{2l}^s} = 2l \sum_{k=-\infty}^{+\infty} (1 + |k|^2)^s \widehat{f}(k)\overline{\widehat{g}(k)}.$$

When $s = 0$, H_{2l}^s is a Hilbert space that is isometrically isomorphic to a subspace of $L^2([-l, l])$ and $(f, g)_{H_{2l}^0} = (f, g) = \int_{-l}^l f(x)\overline{g(x)}dx$. The space H_{2l}^0 will be denoted by L_{2l}^2 , and its norm will be $\|\cdot\|_{L_{2l}^2}$. Of course, $H_{2l}^s \subseteq L_{2l}^2$ for all $s \geq 0$, and we have

for $s > 1/2$ the Sobolev embedding $H_{2l}^s \hookrightarrow C_{2l}$ (see [30]). The space $\ell_{s,2l}^2 := \ell_{s,2l}^2(\mathbb{Z})$, $s \in \mathbb{R}$, is defined by

$$\ell_{s,2l}^2(\mathbb{Z}) := \left\{ \alpha = (\alpha_k)_{k \in \mathbb{Z}}; \|\alpha\|_{\ell_s^2} := \left(2l \sum_{k=-\infty}^{+\infty} (1 + |k|^2)^s |\alpha_k|^2 \right)^{\frac{1}{2}} < +\infty \right\}.$$

$\ell_{s,2l}^2$ is a Hilbert space with inner product

$$(\alpha, \beta)_{\ell_{s,2l}^2} = 2l \sum_{k=-\infty}^{+\infty} (1 + |k|^2)^s \alpha_k \overline{\beta_k},$$

where $\alpha = (\alpha_k)_{k \in \mathbb{Z}}$ and $\beta = (\beta_k)_{k \in \mathbb{Z}}$. Then we have that $f \in H_{2l}^s$ if and only if $(\widehat{f}(k))_{k \in \mathbb{Z}} \in \ell_{s,2l}^2$, and so from the Parseval theorem (see [30]), $\|\widehat{f}\|_{\ell^2}^2 = \frac{1}{2l} \|f\|_{L^2}^2$, it follows that $\|f\|_{H_{2l}^s} = \|\widehat{f}\|_{\ell_{s,2l}^2}$. The convolution of two sequences α and β is the sequence $\alpha * \beta$ defined, for all $k \in \mathbb{Z}$, by $(\alpha * \beta)_k = \sum_{j=-\infty}^{+\infty} \alpha_{k-j} \beta_j$, whenever the right-hand side of the identity above makes sense. Next, we present some results that we will need throughout this work. We start with the Young inequality (see [38]).

PROPOSITION 2.1. *Let $\alpha \in \ell^1(\mathbb{Z})$ and $\beta \in \ell^2(\mathbb{Z})$. Then $\alpha * \beta \in \ell^2(\mathbb{Z})$. Moreover,*

$$\|\alpha * \beta\|_{\ell^2} \leq \|\alpha\|_{\ell^1} \|\beta\|_{\ell^2}.$$

*In particular, for every $\alpha \in \ell^1$ fixed, the linear operator $\beta \in \ell^2 \mapsto \alpha * \beta \in \ell^2$ is continuous.*

Now, we present the Poisson summation theorem. It will be used to find the explicit form of the periodic travelling-wave solutions for some equations.

THEOREM 2.1. *Let $\widehat{f}(x) = \int_{-\infty}^{+\infty} f(y)e^{-ixy} dy$ and $f(y) = \int_{-\infty}^{\infty} \widehat{f}(x)e^{ixy} dx$ satisfying*

$$|f(y)| \leq \frac{A}{(1 + |y|)^{1+\delta}} \quad \text{and} \quad |\widehat{f}(x)| \leq \frac{A}{(1 + |x|)^{1+\delta}},$$

where $\delta > 0$ and $A > 0$ (then \widehat{f} and f can be assumed to be continuous functions). Thus, for $L > 0$,

$$\sum_{n=-\infty}^{+\infty} f(x + 2Ln) = \frac{1}{2L} \sum_{n=-\infty}^{+\infty} \widehat{f}\left(\frac{n}{2L}\right) e^{\frac{\pi i n x}{L}}.$$

The two series above converge absolutely.

Proof. See [38], [45], and [46]. □

2.2. Results of global well-posedness. Some results about local and global well-posedness associated with (1.3) in the periodic case were initially established in [1]. Here we establish two results that we will use in our theory.

THEOREM 2.2. *Let $s \geq 1$ be given. For each $u_0 \in H_{2l}^s$ there is a unique solution of (1.3), for the cases $p = 1, 2$ and $M = -\partial_x^2$, that for each $T > 0$ lies in $C(0, T; H_{2l}^s)$. Moreover, the correspondence $u_0 \mapsto u$ is an analytic function of the relevance function spaces.*

Proof. See [20]. □

THEOREM 2.3. *Let $s \geq \frac{1}{2}$ be given. For each $u_0 \in H_{2l}^s$ there is a unique solution of (1.3), for the cases $p = 1$ and $M = \mathcal{H}\partial_x$, that for each $T > 0$ lies in $C(0, T; H_{2l}^s)$.*

Moreover, the correspondence $u_0 \mapsto u$ is a continuous function of the relevance function spaces.

Proof. See [39], [41], and [40]. \square

3. Basic functional spaces. In this section we establish the main spaces in our study of the stability of periodic wave solutions associated with (1.6) for the case of the Benjamin–Ono (BO), Korteweg–de Vries (KdV), and modified Korteweg–de Vries (mKdV) equations. For the KdV and mKdV cases we have in (1.6), $M = -\partial_x^2$ with symbol $(\frac{\pi}{L})^2 n^2$, $p = 1$, and $p = 2$, respectively. Next, the BO equation is obtained with $M = \mathcal{H}\partial_x$, with the symbol being $\frac{\pi}{L}|n|$ and $p = 1$. Here, \mathcal{H} denotes the periodic Hilbert transform defined in (1.12) and such that via the Fourier transform satisfies $\widehat{\mathcal{H}f}(n) = -i\operatorname{sgn}(n)\widehat{f}(n)$ for all $n \in \mathbb{Z}$. Then, our three periodic travelling-wave equations are

$$\begin{aligned}
 (3.1) \quad & \mathcal{H}\varphi'_c + c\varphi_c - \frac{1}{2}\varphi_c^2 = 0 && \text{(BO),} \\
 & \varphi''_c + \frac{1}{2}\varphi_c^2 - c\varphi_c = 0 && \text{(KdV),} \\
 & \varphi''_c + \frac{1}{3}\varphi_c^3 - c\varphi_c = 0 && \text{(mKdV).}
 \end{aligned}$$

Remark 3.1. (a) Here, we shall consider a more convenient form for the mKdV travelling-wave equation, namely,

$$(3.2) \quad \varphi''_c + \varphi_c^3 - c\varphi_c = 0.$$

(b) The periodic travelling-wave solutions for the KdV and mKdV will be considered of period L , but the periodic travelling-wave solution for the BO equation will be considered of period $2L$ only by convenience.

(c) After a “bootstrap” argument, we can conclude that every φ_c belongs to H_{2L}^s for all $s \in \mathbb{R}$. Thus, φ is infinitely differentiable, with all derivatives in L_{2L}^2 .

We will suppose that $c > -b$, where b satisfies $\alpha(k) > b$ for all $k \in \mathbb{Z}$. With this condition $M + c$ represents a positive operator. Then, by using the spectral theorem for compact and self-adjoint operators we have the following characterization of the spectrum of \mathcal{L} defined in (1.7).

PROPOSITION 3.1. *The operator \mathcal{L} in (1.7) is a closed, unbounded, self-adjoint operator on L_{2L}^2 whose spectrum consists of an enumerable (infinite) set of eigenvalues $\{\lambda_k\}_{k=0}^\infty$ satisfying $\lambda_0 \leq \lambda_1 \leq \lambda_2 \leq \dots$, and $\lambda_k \rightarrow \infty$ as $k \rightarrow \infty$. In particular, \mathcal{L} has 0 as an eigenvalue with eigenfunction $\frac{d}{dx}\varphi_c$.*

Proof. We suppose that our periodic functions have period L . Clearly \mathcal{L} defined on $H_L^{m_2}$ is a closed, unbounded, self-adjoint operator on $L_{per}^2([0, L])$. Let us proof that the spectrum of $\mathcal{T} := M + c$ is a countable infinite set of eigenvalues, $\{\gamma_n\}$, with

$$(3.3) \quad \gamma_0 \leq \gamma_1 \leq \gamma_2 \leq \gamma_3 \leq \dots,$$

where $\gamma_n \rightarrow \infty$ as $n \rightarrow \infty$. In fact, let $R_c = (M + c)^{-1}$, whose symbol is $\frac{1}{c+\alpha(k)}$ for $k \in \mathbb{Z}$. Since $\frac{1}{c+\alpha(k)} \in \ell^2(\mathbb{Z})$ we have that there is a unique $G_c \in L_{per}^2([0, L])$ such that $\widehat{G}_c(k) = \frac{1}{c+\alpha(k)}$, and, because of this, we have the action

$$R_c f(x) = \frac{1}{L} \int_0^L G_c(x - y)f(y)dy,$$

defined for $f \in L_{per}^2([0, L])$. Since $[0, L]$ is a bounded set we have that the kernel $\widetilde{G}_c(x, y) := G_c(x - y) \in L^2([0, L] \times [0, L])$. So, R_c is a Hilbert–Schmidt operator on

$L^2_{per}([0, L])$ (see [33]), and therefore R_c is a compact operator on $L^2_{per}([0, L])$ for all $c > 0$ (here we supposed without loss of generality that $b = 0$), and so we obtain (3.3).

Next, we will show that there is a μ_1 (large enough) such that $\mathcal{M} = (\mathcal{L} + \mu_1)^{-1}$ exists and is a bounded, positive, compact, and self-adjoint operator. In fact, first of all, it is easy to see that \mathcal{L} is limited below; that is, if $f \in D(\mathcal{L})$, we have $\langle \mathcal{L}f, f \rangle \geq -\beta \langle f, f \rangle$, where $\beta = \|\varphi_c\|_{L^\infty_{per}} + c$. Then, we can choose a μ_1 such that $\mathcal{L} + \mu_1 > 0$; that is, \mathcal{M} is positive. We denote $\mu_1 := \mu$ only by convenience. Let ν be a positive number such that $\nu + \varphi_c - c > 0$ and $\nu + \mu > 0$. Thus, for $\mu > 0$ we have $f = (\mathcal{L} + \mu)g \Leftrightarrow (I - M)g = \Upsilon f$, where $Mg = R_{\mu+\nu}[(\nu + \varphi_c - c)g]$, $\Upsilon = R_{\nu+\mu}$, and we denote $h = \nu + \varphi_c - c$. Next, from the Parseval theorem, it follows that

$$\|Mg\|_{L^2_{per}} \leq \sup_{k \in \mathbb{Z}} \left\{ \frac{1}{\alpha(k) + \nu + \mu} \right\} \|h\|_{L^\infty_{per}} \|g\|_{L^2_{per}}.$$

Thus, we can choose μ such that $\|M\|_{B(L^2_{per})} < 1$ and $\mathcal{L} + \mu > 0$. Then, $I - M$ is invertible, and we have $g = (I - M)^{-1}\Upsilon f$ and write $\mathcal{M} = (\mathcal{L} + \mu)^{-1} = (I - M)^{-1}\Upsilon$. Υ being a compact operator and $(I - M)^{-1} \in B(L^2_{per})$ we have that \mathcal{M} is a compact operator. Then, there is an orthonormal basis $\{\varphi_k\}_{k=0}^\infty$ of L^2_{per} consisting of eigenfunctions of \mathcal{M} with nonzero eigenvalues $\{\mu_k\}_{k=0}^\infty$ satisfying $\mu_1 \geq \mu_2 \geq \mu_3 \geq \dots > 0$, and $\mu_k \rightarrow 0$ as $k \rightarrow \infty$. Since $\mathcal{M}\varphi_k = \mu_k\varphi_k \in D(\mathcal{L} + \mu)$ we have that

$$\mathcal{L}\varphi_k = \left(\frac{1}{\mu_k} - \mu \right) \varphi_k := \lambda_k \varphi_k.$$

Thus, there is a sequence of eigenvalues of \mathcal{L} , $\{\lambda_k\}_{k=0}^\infty$, satisfying $\lambda_0 \leq \lambda_1 \leq \lambda_2 \leq \dots$, and $\lambda_k \rightarrow \infty$ as $k \rightarrow \infty$. This argument shows what is desired. \square

The next step will be centralized in the study of some specific spectral properties of the operator \mathcal{L} . For this, let us define two families of linear operators. They will be related with \mathcal{L} but with the advantage that both of them are compacts. The results that will be present below are extensions of the results about stability of solitary waves solutions in Albert [2] and Albert and Bona [3] to the periodic case.

For every $\theta \geq 0$ define the operator $S_\theta : \ell^2(\mathbb{Z}) \rightarrow \ell^2(\mathbb{Z})$ by considering

$$S_\theta \alpha(n) = \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\alpha_j = \frac{1}{\omega_\theta(n)} (K * \alpha)_n,$$

where $\omega_\theta(n) = \alpha(n) + \theta + c$, $K(n) = \widehat{\varphi_c^p}(n)$, $n \in \mathbb{Z}$. Since $\omega_\theta(n) > 0$ for all $n \in \mathbb{Z}$, it follows that the space X defined by

$$X = \left\{ \alpha \in \ell^2(\mathbb{Z}); \|\alpha\|_{X,\theta} := \left(\sum_{n=-\infty}^\infty |\alpha_n|^2 \omega_\theta(n) \right)^{\frac{1}{2}} < \infty \right\}$$

is a Hilbert space with norm $\|\alpha\|_{X,\theta}$ and corresponding inner product $\langle \alpha, \beta \rangle_{X,\theta} = \sum_{n=-\infty}^\infty \alpha_n \overline{\beta_n} \omega_\theta(n)$.

PROPOSITION 3.2. (a) *If $\alpha \in \ell^2$ is an eigensequence of S_θ for a nonzero eigenvalue, then $\alpha \in X$.*

(b) *The restriction of S_θ to X is a compact, self-adjoint operator on X with respect to the norm $\|\cdot\|_{X,\theta}$.*

Proof. In fact, we denote the norm in X simply by $\|\cdot\|_X$, and the operator S_θ by S , and $\mu := \mu_\theta = \frac{1}{\omega_\theta}$. It will be shown that $S^2g = S(Sg) \in X$, since $g = \frac{1}{\gamma}Sg = \frac{1}{\gamma^2}S^2g$;

this will prove the proposition. By the Minkowski, Hölder, and Young inequalities we obtain

$$\begin{aligned}
 (3.4) \quad \|S^2\alpha\|_X &= \left\| \sum_j K(\cdot - j)\mu(\cdot)S\alpha(j) \right\|_X = \left\| \sum_j \sum_m K(\cdot - j)\mu(\cdot)K(j - m)\mu(j)\alpha(m) \right\|_X \\
 &\leq \sum_j \sum_m \left[\sum_n K(n - j)^2\mu(n) \right]^{\frac{1}{2}} K(j - m)\mu(j)|\alpha(m)| \\
 &\leq \|K^2\|_{\ell^1} \|\mu\|_{\ell^2}^{\frac{1}{2}} \|\mu^{\frac{3}{2}}\|_{\ell^1}^{\frac{1}{4}} \|\mu^{\frac{5}{4}}\|_{\ell^1}^{\frac{1}{2}} \|\alpha\|_{\ell^2}.
 \end{aligned}$$

Since $|\mu(x)| \leq B(1 + |x|)^{-1}$ for every $x \in \mathbb{R}$ and some $B > 0$, all of the quantities in the last expression on the right-hand side in (3.4) are finite. The proof of (a) is complete.

Next, we prove that $S|_X$ determines a compact and self-adjoint operator. Let

$$\tilde{K}(n, j) = \frac{K(n - j)}{\omega_\theta(n)\omega_\theta(j)};$$

then, by the Cauchy–Schwarz inequality we have

$$\begin{aligned}
 &\sum_n \sum_j \tilde{K}^2(n, j)\omega_\theta(n)\omega_\theta(j) = \sum_n \sum_j \frac{K^2(n - j)}{\omega_\theta(n)\omega_\theta(j)} \\
 &= \sum_n (K^2 * \mu)(n)\mu(n) \leq \|K^2 * \mu\|_{\ell^2} \|\mu\|_{\ell^2} \leq \|K^2\|_{\ell^1} \|\mu\|_{\ell^2}^2 < \infty.
 \end{aligned}$$

That is, $S_\theta\alpha(n) = \sum_{j=-\infty}^\infty \tilde{K}(n, j)\alpha_n\omega_\theta(j)$ is a Hilbert–Schmidt operator, and thus $S|_X$ is compact. To prove that S_θ is self-adjoint it is necessary to observe that ω_θ and $\widehat{\varphi}_c$ being even, the real kernel K is symmetric. \square

The next two results are immediate consequences of Proposition 3.2 and the spectral theorem for compact, self-adjoint operators on a Hilbert space.

COROLLARY 3.1. *Suppose $\theta \geq 0$. Then, 1 is an eigenvalue of S_θ (as an operator of X) if and only if $-\theta$ is an eigenvalue of \mathcal{L} (as an operator of L^2_{2L}). Furthermore, both eigenvalues have the same multiplicity.* \square

COROLLARY 3.2. *For every $\theta \geq 0$, S_θ has a family of eigensequences $\{\psi_{i,\theta}\}_{i=0}^\infty$ forming an orthonormal basis of X with respect to the norm $\|\cdot\|_{X,\theta}$. The eigensequences correspond to real eigenvalues $\{\lambda_i(\theta)\}_{i=0}^\infty$ whose only possible accumulation point is zero.* \square

In this way, the eigenvalues can be enumerated in order of decreasing absolute value, that is, $|\lambda_0(\theta)| \geq |\lambda_1(\theta)| \geq |\lambda_2(\theta)| \geq \dots$.

4. Positivity properties for periodic travelling-wave solutions. In this section we give sufficient conditions to obtain properties (P_1) and (P_2) in (1.10).

DEFINITION 4.1. *We say that a sequence $\alpha = (\alpha_n)_{n \in \mathbb{Z}} \subseteq \mathbb{R}$ is in the class $PF(2)$ discrete if*

- (i) $\alpha_n > 0$ for all $n \in \mathbb{Z}$,
- (ii) $\alpha_{n_1 - m_1}\alpha_{n_2 - m_2} - \alpha_{n_1 - m_2}\alpha_{n_2 - m_1} > 0$ for $n_1 < n_2$ and $m_1 < m_2$.

The definition above is a particular case of the continuous ones which appear in [2] and [32]; namely, we say that a function $g : \mathbb{R} \rightarrow \mathbb{R}$ is in the class $PF(2)$ if

- (i) $g(x) > 0$ for all $x \in \mathbb{R}$,
- (ii) $g(x_1 - y_1)g(x_2 - y_2) - g(x_1 - y_2)g(x_2 - y_1) > 0$ for $x_1 < x_2$ and $y_1 < y_2$.

As an example, consider $g(x) = \operatorname{sech}^2(x)$.

The next result gives us a sufficient condition for a function g belonging to the $PF(2)$ continuous class. This result is very useful for our purpose (see [3]).

LEMMA 4.1. *Suppose g is a positive, twice-differentiable function on \mathbb{R} satisfying*

$$\frac{d^2}{dx^2}(\log g(x)) < 0 \quad \text{for } x \neq 0.$$

Then $g \in PF(2)$. □

The main result of this paper is now presented.

THEOREM 4.1. *Suppose that φ_c is a positive even solution of (1.6) such that $\widehat{\varphi}_c > 0$ and $K = \widehat{\varphi}_c^p \in PF(2)$ discrete. Then (P_1) and (P_2) in (1.10) hold for the operator \mathcal{L} in (1.7).*

Proof. First, we noticed that, S_θ being a compact operator on X , we get a set of eigenvalues $\{\lambda_i(\theta)\}_{i=0}^\infty$ and the corresponding set of eigenfunctions $\{\psi_{i,\theta}\}_{i=0}^\infty$, which form an orthonormal basis for X . Moreover, we have $|\lambda_0(\theta)| \geq |\lambda_1(\theta)| \geq |\lambda_2(\theta)| \geq \dots$. It will show that the eigenvalues $\lambda_0(\theta)$ and $\lambda_1(\theta)$ are positive, distinct, and simple. In fact, since $S_\theta|_X$ is a compact, self-adjoint operator it follows that

$$(4.1) \quad \lambda_0(\theta) = \pm \sup_{\|\alpha\|_X=1} |\langle S_\theta \alpha, \alpha \rangle_X|.$$

Let $\psi(\theta) := \psi$ be an eigensequence of S_θ corresponding to $\lambda_0(\theta) := \lambda_0$. We will show that ψ is one-signed; that is, either $\psi(n) \leq 0$ or $\psi(n) \geq 0$. By contradiction, suppose ψ takes both negative and positive values. Since by hypotheses the kernel K is positive we have

$$\begin{aligned} S_\theta|\psi|(n) &= \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\psi^+(j) + \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\psi^-(j) \\ &> \left| \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\psi^+(j) - \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\psi^-(j) \right|, \end{aligned}$$

where ψ^+ e ψ^- are the positive and negative parts of ψ , respectively. Then,

$$S_\theta|\psi|(n) > \left| \frac{1}{\omega_\theta(n)} \sum_{j=-\infty}^\infty K(n-j)\psi(j) \right| = |S_\theta\psi(n)| = |\lambda_0||\psi(n)|,$$

where “>” holds because ψ , by supposition, takes both positive and negative values. From the last inequality we conclude that

$$\begin{aligned} \langle S_\theta(|\psi|), |\psi| \rangle_{X,\theta} &= \sum_{n=-\infty}^\infty S_\theta|\psi|(n)|\psi(n)|\omega_\theta(n) \\ &> \sum_{n=-\infty}^\infty |\lambda_0||\psi(n)|^2\omega_\theta(n) = |\lambda_0|\|\psi\|_{X,\theta}^2. \end{aligned}$$

Hence, if we assume that $\|\psi\|_X = 1$, we obtain $\langle S_\theta(|\psi|), |\psi| \rangle_X |\lambda_0|$, which contradicts (4.1). Then, there is an eigensequence ψ_0 which is nonnegative. Since K is a positive sequence and $S_\theta(\psi_0) = \lambda_0\psi_0$, we have $\psi_0(n) > 0$ for all $n \in \mathbb{Z}$. Now, such a ψ_0 cannot be orthogonal to any nontrivial one-signed eigensequence in X , and so λ_0 is a simple eigenvalue. Notice that the preceding argument also shows that $-\lambda_0$ cannot be an eigenvalue of S_θ ; therefore it follows that $|\lambda_1| < \lambda_0$.

Next, we will study the eigenvalue $\lambda_1(\theta) = \lambda_1$. But, first, we need some definitions and results. We consider the following set of indices:

$$\Delta = \{(n_1, n_2) \in \mathbb{Z} \times \mathbb{Z}; n_1 < n_2\}.$$

Denoting $\bar{n} = (n_1, n_2)$ and $\bar{m} = (m_1, m_2)$, we define for $\bar{n}, \bar{m} \in \Delta$ the following sequence:

$$K_2(\bar{n}, \bar{m}) := K(n_1 - m_1)K(n_2 - m_2) - K(n_1 - m_2)K(n_2 - m_1).$$

By hypothesis $K \in PF(2)$ discrete it follows that $K_2 > 0$. Next, let $\ell^2(\Delta)$ be defined as

$$\ell^2(\Delta) = \left\{ \alpha = (\alpha_{\bar{n}})_{\bar{n} \in \Delta}; \sum_{\Delta} \sum_{\Delta} |\alpha_{\bar{n}}|^2 := \sum_{n_1 \in \mathbb{N}} \sum_{\substack{n_1 < n_2 \\ n_2 \in \mathbb{Z}}} |\alpha(n_1, n_2)|^2 < +\infty \right\},$$

and define the operator $S_{2,\theta} : \ell^2(\Delta) \rightarrow \ell^2(\Delta)$ by

$$S_{2,\theta}g(\bar{n}) = \sum_{\Delta} \sum_{\Delta} G_{2,\theta}(\bar{n}, \bar{m})g(\bar{m}),$$

where $G_{2,\theta}(\bar{n}, \bar{m}) = \frac{K_2(\bar{n}, \bar{m})}{\omega_\theta(n_1)\omega_\theta(n_2)}$. We also consider the space

$$W = \left\{ \alpha \in \ell^2(\Delta); \|\alpha\|_{W,\theta} := \left(\sum_{\Delta} \sum_{\Delta} |\alpha(\bar{n})|^2 \omega_\theta(n_1)\omega_\theta(n_2) \right)^{\frac{1}{2}} < \infty \right\}.$$

Then W is a Hilbert space with norm $\|\cdot\|_{W,\theta}$ given above and with inner product

$$\langle \alpha, \beta \rangle_{W,\theta} = \sum_{\Delta} \sum_{\Delta} \alpha(\bar{n})\overline{\beta(\bar{n})}\omega_\theta(n_1)\omega_\theta(n_2).$$

Remark 4.1. (1) We can show, in an analogous way to Proposition 3.2, that $S_{2,\theta}|_W$ is a self-adjoint, compact operator. Therefore, the associated eigenvalues can be enumerated in order of decreasing absolute value, that is, $|\mu_0(\theta)| \geq |\mu_1(\theta)| \geq |\mu_2(\theta)| \geq \dots$.

(2) A similar argument can be used to show that $\mu_0(\theta) := \mu_0$ is positive and simple and $|\mu_1| < \mu_0$.

DEFINITION 4.2. Let $\alpha^1, \alpha^2 \in \ell^2(\mathbb{Z})$; we define the wedge product $\alpha^1 \wedge \alpha^2$ in Δ by

$$(\alpha^1 \wedge \alpha^2)(n_1, n_2) = \alpha^1(n_1)\alpha^2(n_2) - \alpha^1(n_2)\alpha^2(n_1).$$

We have the following results from Definition 4.2.

LEMMA 4.2. Let $A = \{\alpha^1 \wedge \alpha^2; \text{ for } \alpha^1, \alpha^2 \in X, \alpha^1 \wedge \alpha^2 \in \ell^2(\Delta)\}$. Then A is dense in W .

Proof. See [31] and [32]. \square

LEMMA 4.3. *Let $\alpha^1, \alpha^2 \in \ell^2(\mathbb{Z})$. Then $S_{2,\theta}(\alpha^1 \wedge \alpha^2) = S_\theta \alpha^1 \wedge S_\theta \alpha^2$.*

Proof. See [31] and [32]. \square

In what follows, we will represent by S_θ^X the restriction of S_θ on the Hilbert space X . We shall use some spectral results in Kato [33]. In fact, we decompose S_θ^X into the form $S_\theta^X = \lambda_0 P_\theta + Q_\theta$, where P_θ is the orthogonal projection on $M_0 = [\psi_0]$, $P_\theta Q_\theta = Q_\theta P_\theta = 0$, and $\text{spectrum}(Q_\theta) = \text{spectrum}(P_\theta) \setminus \{\lambda_0\}$. Moreover, in this case we have that the spectral radius of Q_θ when restricted to the subspace $N = \text{Ker} P_\theta$ is exactly $|\lambda_1|$. Furthermore, we know that the eigenvalue $\lambda_0 = \lambda_0(\theta)$ is differentiable with respect to $\theta \geq 0$. The same argument can be applied to μ_0 . Then, as a consequence ψ_0 and $\tau_0 = \text{eigensequence}$ associated with μ are differentiable eigensequences with respect to θ .

LEMMA 4.4. (a) *In the notation given above we have*

$$\frac{(S_\theta^X)^m}{\lambda_0^m} \longrightarrow P_\theta,$$

where $m \rightarrow +\infty$ on the strong topology of $B(X, X)$.

(b) *The statement in part (a) is valid if S_θ^X is replaced by $S_{2,\theta}$, the eigenvalues λ_0 and λ_1 are replaced by μ_0 and μ_1 , and P_θ, M_0 , and Q_θ are replaced by appropriate operators and subspaces of X .*

Proof. The proof is analogous as viewed in [2], [3] in the case of sequences. \square

LEMMA 4.5. (a) $\mu_0(\theta) = \lambda_0(\theta)\lambda_1(\theta)$. *Then, we can conclude that $\lambda_1 > 0$.*

(b) λ_1 *is simple.*

Proof. (a) In fact, from Lemma 4.3 we have that $\lambda_0\lambda_1$ is an eigenvalue of $S_{2,\theta}$ whose eigensequence is $\psi_0 \wedge \psi_1$, where $\psi_1(\theta) := \psi_1$ is the eigensequence associated with λ_1 . Hence $\mu_0 \geq \lambda_0|\lambda_1|$. Then, since $-\mu_0$ cannot be an eigenvalue of $S_{2,\theta}$, we will show that $\mu_0 \leq \lambda_0|\lambda_1|$. If $|\lambda_1| < \frac{\mu_0}{\lambda_0}$, let P_θ be as in Lemma 4.4 and write $W = M_0 \oplus N$. Let $\alpha^1 = r_1\psi_0 + \omega\gamma^1$ and $\alpha^2 = r_2\psi_0 + \omega\gamma^2$, where $r_1, r_2 \in \mathbb{R}$ and $\gamma^1, \gamma^2 \in N$. For instance, from the induction principle we have

(4.2)

$$\begin{aligned} \left(\frac{S_{2,\theta}}{\mu_0}\right)^m (\alpha^1 \wedge \alpha^2)(n_1, n_2) &= r_1 \left[\psi_0(n_1) \left(\frac{S_\theta}{\beta}\right)^m \gamma^2(n_2) - \psi_0(n_2) \left(\frac{S_\theta}{\beta}\right)^m \gamma^2(n_1) \right] \\ &\quad + r_2 \left[\psi_0(n_2) \left(\frac{S_\theta}{\beta}\right)^m \gamma^1(n_1) - \psi_0(n_1) \left(\frac{S_\theta}{\beta}\right)^m \gamma^1(n_2) \right] \\ &\quad + \left(\frac{S_\theta}{\lambda_0}\right)^m \gamma^1(n_1) \left(\frac{S_\theta}{\beta}\right)^m \gamma^2(n_2) \\ &\quad - \left(\frac{S_\theta}{\lambda_0}\right)^m \gamma^1(n_2) \left(\frac{S_\theta}{\beta}\right)^m \gamma^2(n_1), \end{aligned}$$

where $\beta = \frac{\mu_0}{\lambda_0} > |\lambda_1|$. Since $\left(\frac{S_\theta^X}{\lambda_0}\right)^m \rightarrow P_\theta$ with $P_\theta \equiv 0$ on N and β is strictly greater than the spectral radius of the restriction of S_θ^X to N , each term on the right-hand side of the preceding equality tends to zero as $m \rightarrow \infty$. But the set A defined in Lemma 4.2 is dense in W . Then, we can conclude from (4.2), after a computation, that $\left(\frac{S_{2,\theta}}{\mu_0}\right)^m g \rightarrow 0$ strongly. But this contradicts Lemma 4.4. The proof of item (a) is completed.

(b) The next step is to show that λ_1 is simple. We write $\psi_1 = \psi_1^P + \psi_1^I$, where ψ_1^P and ψ_1^I denote, respectively, the even and odd parts of ψ_1 . We begin by showing that $\psi_1^P \equiv 0$. Since the kernel K of S_θ is symmetric and ω_θ is even, we have S_θ map even sequences into even sequences and odd sequences into odd sequences. Then, we get that $\psi_0 \wedge \psi_1^P$ satisfies $S_{2,\theta}(\psi_0 \wedge \psi_1^P) = \lambda_0 \lambda_1 (\psi_0 \wedge \psi_1^P)$. Since $\mu_0 > 0$ and simple we obtain $\psi_0 \wedge \psi_1^P \in [\tau_0]$. Therefore, either $\psi_0 \wedge \psi_1^P \equiv 0$ or $\psi_0 \wedge \psi_1^P \neq 0$; that is, $\psi_0 \wedge \psi_1^P$ is one-signed. With this fact, if ψ_1^P does not vanish, then it will have at most one zero. Then, ψ_1^P being an even sequence, either $\psi_1^P \equiv 0$ or else ψ_1^P is one-signed, except possibly in $n = 0$. If the second case holds, we have to consider three cases: $\psi_1^P(0) = 0$, $\psi_1^P(0) > 0$, and $\psi_1^P(0) < 0$. If $\psi_1^P(0) = 0$, then we should have, from the definition of wedge product above, that $0 < (\psi_0 \wedge \psi_1^P)(0, n) = \psi_0(0)\psi_1^P(n)$ for $n > 0$ (where we are assuming that $\psi_0 > 0$ and $\psi_0 \wedge \psi_1^P > 0$). Then, $\psi_1^P(n) > 0$ for all $n \in \mathbb{Z}$ because $\psi_1^P(n)$ is an even sequence. Next, if $\psi_1^P(0) > 0$, then for $n > 0$ we get $0 < (\psi_0 \wedge \psi_1^P)(0, n) = \psi_0(0)\psi_1^P(n) - \psi_0(n)\psi_1^P(0) < \psi_0(0)\psi_1^P(n)$ and therefore $\psi_1^P > 0$. The last case is similar to the second one. These considerations make ψ_1^P a one-signed eigensequence of S_θ (except possibly in $n \neq 0$), and so the inner product $\langle \psi_1^P, \psi_0 \rangle_{X,\theta}$ is also one-signed. But this a contradiction because we have two eigensequences of a self-adjoint operator associated with distinct eigenvalues whose inner product is nonzero. Therefore, $\psi_1^P \equiv 0$, and so ψ_1 is odd. A similar argument shows that ψ_1 can have at most one zero. Of course, this one must be in $n = 0$.

The sequence ψ_1 shown previously was an arbitrary sequence, and this is associated with eigenvalue λ_1 . Then, we show that any eigensequence ψ associated with λ_1 must be odd and $\psi(n) = 0 \Leftrightarrow n = 0$. But, two eigensequences of this kind cannot be orthogonal since the product of them is even, and thus λ_1 is simple. This fact completes the proof of the lemma. \square

We turn back to the proof of Theorem 4.1. Let us consider $\|\psi_i(\theta)\|_{X,\theta} = 1$ for $i = 0, 1$. From Lemma 4.5, $\mu_0(\theta) = \lambda_0(\theta)\lambda_1(\theta)$ with μ_0 and λ_0 differentiable with respect to θ , and we have that λ_1 is also differentiable with respect to this parameter. Hence, the associated eigensequence ψ_1 is also differentiable since $S_\theta\psi_1(\theta) = \lambda_1(\theta)\psi_1(\theta)$. Next, we will show that

$$(4.3) \quad \frac{d}{d\theta}\lambda_i(\theta) < 0, \quad i = 0, 1, \quad \theta \geq 0.$$

In fact, writing $\psi_i(\theta)$, $i = 0, 1$, instead of $\psi_{i,\theta}(n)$, we have

$$\frac{d}{d\theta}\lambda_i(\theta) = \frac{d}{d\theta} \sum_{n=-\infty}^{\infty} S_\theta\psi_i(\theta)\psi_i(\theta)\omega_\theta(n).$$

Thus,

$$\begin{aligned} \frac{d}{d\theta}\lambda_i(\theta) &= 2 \sum_{n=-\infty}^{\infty} \frac{d}{d\theta}\psi_i(\theta)S_\theta\psi_i(\theta)\omega_\theta(n) = 2\lambda_i(\theta) \sum_{n=-\infty}^{\infty} \left(\frac{d}{d\theta}\psi_i(\theta)\right)\psi_i(\theta)\omega_\theta(n) \\ &= 2\lambda_i(\theta) \left\{ \frac{d}{d\theta} \frac{1}{2} \left(\sum_{n=-\infty}^{\infty} \psi_i(\theta)^2 \omega_\theta(n) \right) - \frac{1}{2} \sum_{n=-\infty}^{\infty} \psi_i(\theta)^2 \right\} \\ &= -\lambda_i(\theta) \sum_{n=-\infty}^{\infty} \psi_i(\theta)^2 < 0, \end{aligned}$$

which shows the affirmation. Next, for $\theta \geq 0$,

$$\begin{aligned} \lambda_0(\theta) &= r(S_\theta^X) = \|S_\theta^X\|_{B(X,X)} \leq \left(\sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \left\{ \frac{K(n-m)}{\omega_\theta(n)} \right\}^2 \right)^{\frac{1}{2}} \\ &= \left\| K^2 * \frac{1}{\omega_\theta^2} \right\|_{\ell^1(\mathbb{Z})}^{\frac{1}{2}} \leq \|K\|_{\ell^2(\mathbb{Z})} \left\| \frac{1}{\omega_\theta} \right\|_{\ell^2(\mathbb{Z})}. \end{aligned}$$

Since $\frac{1}{\omega_\theta} \rightarrow 0$ as $\theta \rightarrow +\infty$ and $(\frac{1}{\omega_\theta})^2 \in \ell^1$ with $|\frac{1}{\omega_\theta}|^2 \leq (\frac{B}{1+|n|})^2$, for some $B > 0$, $\|\frac{1}{\omega_\theta}\|_{\ell^2} \rightarrow 0$ as $\theta \rightarrow +\infty$. Therefore,

$$(4.4) \quad \lim_{\theta \rightarrow +\infty} \lambda_0(\theta) = 0.$$

The next step is to show that

$$(4.5) \quad \lambda_1(0) = 1.$$

In fact, since $\frac{d}{dx}\varphi_c$ is an eigenfunction of \mathcal{L} with eigenvalue $\theta = 0$, we can conclude from Corollary 3.1 that $\widehat{\frac{d}{dx}\varphi_c}$ is an eigensequence of S_θ with eigenvalue 1. On the other hand, $\widehat{\frac{d}{dx}\varphi_c}(n) = -in\frac{\pi}{L}\widehat{\varphi_c}(n)$ is odd and vanishes only at $n = 0$. Since ψ_1 is also odd and vanishes at $n = 0$ we must have $\langle \psi_1, \widehat{\frac{d}{dx}\varphi_c} \rangle_{X,\theta} \neq 0$. It follows that ψ_1 and $\widehat{\frac{d}{dx}\varphi_c}$ cannot be eigensequences of S_θ for distinct eigenvalues. Then, ψ_1 and $\widehat{\frac{d}{dx}\varphi_c}$ are associated with the same eigenvalue, and therefore $\lambda_1(0) = 1$. With this fact, since $\lambda_0(0) > \lambda_1(0) = 1$, from (4.3) and (4.4) it follows that there is a unique $\theta_0 \in (0, +\infty)$ such that $\lambda_0(\theta_0) = 1$. Then, from Corollary 3.1, if we consider $\kappa = -\theta_0 < 0$, then \mathcal{L} has a negative eigenvalue which is simple. Now, for $i \geq 2$ and $\theta > 0$ we have from (4.5) that

$$\lambda_i(\theta) \leq \lambda_1(\theta) < \lambda_1(0) = 1.$$

It is straightforward to see that 1 cannot be an eigenvalue of S_θ for all $\theta \in (0, +\infty) \setminus \{\theta_0\}$, since 1 is an eigenvalue only for $\theta = 0$ and $\theta = \theta_0$. Then we obtain (P_1) .

Because $\lambda_1(0) = 1$ and λ_1 is a simple eigenvalue it follows that $\theta = 0$ is a simple eigenvalue of \mathcal{L} by the Corollary 3.1. This fact shows (P_2) and as a consequence the theorem. \square

Remark 4.2. The Fourier transform in Theorem 4.1 needs to be evaluated in the minimal period of φ_c (see section 6).

5. Stability of periodic travelling-wave solutions. In this section we are interested in applying the theory in section 4 to obtain the stability of specific periodic travelling waves associated with the KdV, mKdV, and BO equations. Our approach to obtain condition (P_0) in (1.10) will be based on the Poisson summation theorem and the implicit function theorem. This new approach, in the periodic context, will give a simple way to calculate condition (P_3) in (1.10). We start with the definition of stability.

DEFINITION 5.1. *Let φ be a periodic travelling-wave solution with period $2L$ of (1.6), and consider $\tau_r\varphi(x) = \varphi(x+r)$, $x \in \mathbb{R}$, and $r \in \mathbb{R}$. We define the set $\Omega_\varphi \subset H_{2L}^{\frac{m-2}{2}}$, the orbit generated by φ , as*

$$\Omega_\varphi = \{g; g = \tau_r\varphi \text{ for some } r \in \mathbb{R}\}.$$

And, for any $\eta > 0$, define the set $U_\eta \subset H_{2L}^{\frac{m_2}{2}}$ by

$$U_\eta = \left\{ f; \inf_{g \in \Omega_\varphi} \|f - g\|_{H_{2L}^{\frac{m_2}{2}}} < \eta \right\}.$$

With this terminology, we say that φ is (orbitally) stable in $H_{2L}^{\frac{m_2}{2}}$ by the flow generated by (1.3) if the following hold:

(i) There is s_0 such that $H_{2L}^{s_0} \subseteq H_{2L}^{\frac{m_2}{2}}$ and the initial value problem associated with (1.3) is globally well-posed in $H_{2L}^{s_0}$ (see Theorems 2.2 and 2.3).

(ii) For every $\varepsilon > 0$, there is $\delta > 0$ such that, for all $u_0 \in U_\delta \cap H_{2L}^{s_0}$, the solution u of (1.3) with $u(0, x) = u_0(x)$ satisfies $u(t) \in U_\varepsilon$ for all $t > 0$.

The proof of the following general stability theorem can be shown by using the techniques in Grillakis, Shatah, and Strauss [27] (see also Angulo [6]).

THEOREM 5.1. *Let φ_c be a periodic travelling-wave solution of (1.6), and suppose that part (i) of the definition of stability holds. Suppose also that the operator \mathcal{L} defined previously in (1.7) has properties (P_1) and (P_2) in (1.10). Choose $\chi \in L_{2L}^2$ such that $\mathcal{L}\chi = \varphi_c$, and define $I = (\chi, \varphi_c)_{L_{2L}^2}$. If $I < 0$, then φ_c is stable.*

Remark 5.1. In our cases the function χ in Theorem 5.1 satisfies that $\chi = -\frac{d}{dc}\varphi_c$. Then, we need to verify properties (P_0) and (P_3) in (1.10).

5.1. Stability of periodic travelling-wave solutions for the BO equation.

This section is concerned with the stability theory of periodic travelling-wave solutions to the BO equation found initially by Benjamin in [13]. Next, we will present an interesting method for obtaining an explicit solution to the BO equation in (3.1), by using the Poisson summation theorem. In fact, consider the following equation:

$$\mathcal{H}\phi'_\omega + \omega\phi_\omega - \frac{1}{2}\phi_\omega^2 = 0.$$

This equation determines solitary travelling-wave solutions to the BO equation on \mathbb{R} in the form

$$\phi_\omega(x) = \frac{4\omega}{1 + \omega^2 x^2}, \quad \omega > 0.$$

Its Fourier transform is given by

$$\widehat{\phi_\omega}^{\mathbb{R}}(x) = 4\pi e^{-\frac{2\pi}{\omega}|x|}.$$

Then, by the Poisson summation theorem, we obtain that

$$\begin{aligned} \psi_\omega(x) &\equiv \sum_{n=-\infty}^{+\infty} \phi_\omega(x + 2Ln) = \frac{2\pi}{L} \sum_{n=-\infty}^{+\infty} e^{-\frac{\pi|n|}{\omega L}} e^{\frac{\pi i n x}{L}} \\ (5.1) \quad &= \frac{2\pi}{L} \sum_{n=0}^{+\infty} \varepsilon_n e^{-\frac{\pi n}{\omega L}} \cos\left(\frac{n\pi x}{L}\right) = \frac{2\pi}{L} \operatorname{Re} \left[\coth\left(\frac{\pi}{2\omega L} + \frac{i\pi x}{2L}\right) \right] \\ &= \frac{2\pi}{L} \left(\frac{\sinh\left(\frac{\pi}{\omega L}\right)}{\cosh\left(\frac{\pi}{\omega L}\right) - \cos\left(\frac{\pi x}{L}\right)} \right), \end{aligned}$$

where

$$\varepsilon_n = \begin{cases} 1 & \text{if } n = 0, \\ 2 & \text{if } n = 1, 2, 3, \dots \end{cases}$$

Let $\varphi_c, c > 0$, be a smooth periodic solution of the first equation in (3.1). Thus, φ_c can be expressed as a Fourier series

$$(5.2) \quad \varphi_c(x) = \sum_{n=-\infty}^{+\infty} a_n e^{\frac{in\pi x}{L}}.$$

Substituting the expression above into the BO equation in (3.1), we get

$$\left[\frac{\pi|n|}{L} + c \right] a_n = \frac{1}{2} \sum_{m=-\infty}^{+\infty} a_{n-m} a_m.$$

Next, from (5.1) we consider $a_n \equiv \frac{2\pi}{L} e^{-\gamma|n|}, n \in \mathbb{Z}, \gamma \in \mathbb{R}$. Substituting a_n into the last identity we have

$$\sum_{m=-\infty}^{+\infty} a_{n-m} a_m = \frac{4\pi^2}{L^2} e^{-\gamma|n|} \left[|n| + 1 + 2 \sum_{k=1}^{+\infty} e^{-2\gamma k} \right] = \frac{4\pi^2}{L^2} e^{-\gamma|n|} (|n| + \coth\gamma).$$

Then, we conclude that

$$(5.3) \quad c + \frac{\pi|n|}{L} = \frac{2\pi}{L} \cdot \frac{1}{2} (|n| + \coth\gamma).$$

We denote $\gamma = \frac{\pi}{\omega L}$ and consider $c > \frac{\pi}{L}$. Then, if we choose $\omega = \omega(c) > 0$ such that $\tanh(\gamma) = \frac{\pi}{cL}$, we obtain from (5.3) that $\psi_{\omega(c)} = \varphi_c$ (hence, φ_c is given by (5.1)). Therefore, we obtain that φ_c has the form in (1.13) with $\gamma > 0$ satisfying $\tanh(\gamma) = \frac{\pi}{cL}$.

Thus, from (5.1) we have that $\varphi_c > 0$, and $\gamma := \gamma(c) = \tanh^{-1} \left(\frac{\pi}{cL} \right)$ being a differentiable function for $c > \frac{\pi}{L}$, it follows that

$$c \in \left(\frac{\pi}{L}, +\infty \right) \mapsto \varphi_c \in H_{2L}^n$$

is a smooth curve of periodic travelling-wave solutions for the BO equation for all $n \in \mathbb{N}$. Then, by defining χ in Theorem 5.1 as $\chi = -\frac{d}{dc} \varphi_c$, we obtain from the first equation in (3.1) that $\mathcal{L}\chi = \varphi_c$. Then $I = (\chi, \varphi_c)_{L^2_{2L}}$ becomes

$$(5.4) \quad I = -\frac{1}{2} \frac{d}{dc} \|\varphi_c\|_{L^2_{2L}}^2.$$

We will show that $I < 0$. Indeed, from (5.2) and (5.1) we get

$$(5.5) \quad \varphi_c(x) = \frac{2\pi}{L} \sum_{n=-\infty}^{+\infty} e^{-\gamma|n|} e^{\frac{in\pi x}{L}};$$

then, from the Parseval theorem we conclude that

$$(5.6) \quad \begin{aligned} I &= -\frac{1}{2} \frac{d}{dc} \|\varphi_c\|_{L^2_{2L}}^2 = -\frac{1}{2} \frac{d}{dc} \|\widehat{\varphi}_c\|_{\ell^2}^2 \cdot 2L = -\frac{1}{2} \frac{d}{dc} \left(\frac{4\pi^2}{L^2} \sum_{n=-\infty}^{\infty} e^{-2\gamma|n|} \right) \cdot 2L \\ &= -\frac{4\pi^3}{c^2 L^3} \left(\frac{1}{1 - \left(\frac{\pi}{cL}\right)^2} \right) \left(\sum_{n=-\infty}^{\infty} |n| e^{-2\gamma|n|} \right) \cdot 2L < 0. \end{aligned}$$

Finally, we will verify that conditions (P_1) and (P_2) are true for the operator

$$\mathcal{L}_{BO} = \mathcal{H}\partial_x + c - \varphi_c.$$

Let $\widehat{\varphi}_c(n) = \frac{2\pi}{L}e^{-\gamma|n|}$ be the Fourier coefficients of φ_c . In Albert [2] it has already been seen that the function $f(x) = e^{-\gamma|x|}$ belongs to the $PF(2)$ class in the continuous case, and so $\widehat{\varphi}_c$ is in the $PF(2)$ class in the sense of Definition 4.1. Hence, we obtain from Theorem 2.3 the stability of the periodic solutions (1.13) in $H_{2L}^{\frac{1}{2}}$ by the periodic flow of the BO equation.

5.1.1. Stability of the constant solutions. To complete the investigation about periodic travelling-wave solutions for the BO equation, we will study the stability of the constant solutions. Hence, if $\varphi_c(x) \equiv \tau$ is a constant solution to the BO equation, we have $\varphi_c \equiv 2c$ and $\varphi_c \equiv 0$ as solutions. We consider only the first case. The next result will resume our purpose.

PROPOSITION 5.1. *Let $L > 0$ and $c > 0$ be given. Consider $\psi_0 \equiv 2c$ to be a nontrivial constant solution of (3.1). Then, ψ_0 is stable in $H_{2L}^{\frac{1}{2}}([-L, L])$, provided $c < \frac{\pi}{L}$.*

Proof. The proof of this proposition follows from standard ideas (see [10], [15]) and from the following nonlocal Poincaré–Wirtinger-type inequality: for $f \in H_{2L}^{\frac{1}{2}}$ such that $\int_{-L}^L f(x)dx = 0$, we have $\int_{-L}^L [D^{\frac{1}{2}}f(x)]^2 dx \geq \frac{\pi}{L} \int_{-L}^L f^2(x)dx$, where $D = \mathcal{H}\partial_x$ (see the appendix for a proof of this inequality). \square

5.2. Stability of periodic travelling-wave solutions for the mKdV equation. Next, we will establish the existence of a smooth curve of periodic travelling-wave solutions for the mKdV equation

$$(5.7) \quad u_t + 3u^2u_x + u_{xxx} = 0,$$

of the form $u(x, t) = \varphi(x - ct) := \varphi_c(\xi)$, where $\xi = x - ct$, $c \in \mathbb{R}$, and is of period L . The equation which determines the periodic travelling-wave solutions is

$$(5.8) \quad \varphi_c'' + \varphi_c^3 - c\varphi_c = 0.$$

Next, we obtain an explicit solution for (5.8) using the Poisson summation theorem. It considers for $\omega > 0$ the solitary wave solutions for the mKdV equation on \mathbb{R} :

$$\phi_\omega(x) = \sqrt{2\omega} \operatorname{sech}(\sqrt{\omega}x), \quad x \in \mathbb{R}.$$

Its Fourier transform is $\widehat{\phi}_\omega(x) = \sqrt{2\pi} \operatorname{sech}(\frac{\pi x}{2\sqrt{\omega}})$, where $\omega > 0$, and it will be chosen later. From the Poisson summation theorem we obtain the following periodic function of period L :

$$(5.9) \quad \psi_\omega(\xi) = \frac{\sqrt{2\pi}}{L} \sum_{n=0}^{\infty} \varepsilon_n \operatorname{sech}\left(\frac{\pi n}{2\sqrt{\omega}L}\right) \cos\left(\frac{2\pi n\xi}{L}\right),$$

where

$$\varepsilon_n = \begin{cases} 1, & n = 0, \\ 2, & n = 1, 2, 3, \dots \end{cases}$$

On the other hand, it considers the Fourier expansion of the Jacobi elliptic function *dnoidal*, dn , of period L (see [19], [35], [42]),

$$\frac{2K}{L} \operatorname{dn} \left(\frac{2K\xi}{L}; k \right) = \frac{\pi}{L} + \frac{4\pi}{L} \sum_{n=1}^{+\infty} \frac{q^n}{1+q^{2n}} \cos \left(\frac{2n\pi\xi}{L} \right),$$

where $K = K(k)$ is the complete elliptic integral of the first kind and $q = e^{(-\frac{\pi K'}{K})}$, which is called the “nome.” Here, $K'(k) = K(\sqrt{1-k^2})$. We can conclude that

$$\frac{q^n}{1+q^{2n}} = \frac{1}{2} \operatorname{sech} \left(\frac{n\pi K'}{K} \right).$$

Therefore,

$$\frac{2K}{L} \operatorname{dn} \left(\frac{2K\xi}{L}; k \right) = \frac{\pi}{L} + \frac{2\pi}{L} \sum_{n=1}^{+\infty} \operatorname{sech} \left(\frac{n\pi K'}{K} \right) \cos \left(\frac{2n\pi\xi}{L} \right).$$

Because of the shape of the series that determines ψ_ω given above (see [35]), let $\varphi_c(\xi) = \eta \operatorname{dn}(\frac{\eta\xi}{\sqrt{2}}; k)$ be a periodic solution of period L for (5.8), with $\eta > 0$ and $k \in (0, 1)$ fixed. Then, the following identities should be satisfied:

$$(5.10) \quad c = \frac{\eta^2}{2}(1+k'^2) \quad \text{and} \quad \eta = \frac{2\sqrt{2}K(k)}{L}$$

with $k'^2 = 1-k^2$. Thus, for $k \in (0, 1)$ we should have that $\eta \in (\sqrt{c}, \sqrt{2c})$ and from the asymptotic properties of K that $c > \frac{2\pi^2}{L^2}$. Next, for $k \in (0, 1)$ fixed, η is immediately defined from (5.10), and so

$$(5.11) \quad c = \frac{4K^2(k)}{L^2}(2-k^2) > \frac{2\pi^2}{L^2}$$

since $k \rightarrow K^2(k)(2-k^2)$ is a strictly increasing function. Therefore, with $k \in (0, 1)$ and c defined in (5.11) we define $\omega = \omega(c)$ as

$$\omega = \frac{c}{16(2-k^2)K'^2(k)}.$$

Therefore, from (5.9) it follows that $\psi_{\omega(c)} = \varphi_c$ is a solution of (5.8).

Next, it is necessary to build a smooth curve, $c \rightarrow \varphi_c$, of dnoidal wave solutions. Initially, we obtain the following a priori estimate for the fundamental period of φ_c . Namely,

$$(5.12) \quad T_{\varphi_c}(\eta) = \frac{2K(k(\eta))}{\sqrt{c}} \sqrt{2-k^2} > \frac{\sqrt{2}\pi}{\sqrt{c}},$$

where $k^2(\eta) = 2 - \frac{2c}{\eta^2}$. In fact, for $\eta \rightarrow \sqrt{c}$ it follows that $k \rightarrow 0^+$ and thus $T_{\varphi_c}(\eta) \rightarrow \frac{\pi\sqrt{2}}{\sqrt{c}}$. For $\eta \rightarrow \sqrt{2c}$ we have $k \rightarrow 1^-$ and therefore $T_{\varphi_c}(\eta) \rightarrow +\infty$. But, $\eta \mapsto T_{\varphi_c}(\eta)$ being a strictly increasing function (see Theorem 2.1 in Angulo [9]) we get (5.12).

Remark 5.2. The function φ_c obtained previously, that is,

$$(5.13) \quad \varphi_c(\xi) = \eta \operatorname{dn} \left(\frac{\eta}{\sqrt{2}} \xi; k \right),$$

is a positive function and has been built by the periodization of the solitary wave solution associated with (5.7). Thus, it is natural to ask if we can again obtain this solitary wave. Indeed, this fact can be determined by (5.13), since for $\eta \rightarrow \sqrt{2c}$ we have $k \rightarrow 1^-$ and then $\text{dn}(u; 1^-) = \text{sech}(u)$. Hence we have formally that $\varphi_c(\xi) = \sqrt{2c} \text{sech}(\sqrt{c}\xi)$. The other limit case, that is, $\eta \rightarrow \sqrt{c}$, we have $k \rightarrow 0^+$ and so $\text{dn}(u; 0^+) = 1$; then we get $\varphi_c(\xi) = \sqrt{c}$, the nontrivial constant solutions for the mKdV.

Now, we construct a family of dnoidal waves solutions with period L . Let $c > 0$ such that $\sqrt{c} > \frac{\pi\sqrt{2}}{L}$. Since $\eta \in (\sqrt{c}, \sqrt{2c}) \mapsto T_{\varphi_c}(\eta)$ is a strictly increasing mapping, it follows from (5.12) and from Theorem 2.1 in Angulo [9] that there is a unique $\eta \equiv \eta(c) \in (\sqrt{c}, \sqrt{2c})$ such that the fundamental period of the dnoidal wave φ_c will be $T_{\varphi_c}(\eta(c)) = L$. Moreover, we have $c \in (2\pi^2/L^2, +\infty) \rightarrow \varphi_c \in H_{per}^n([0, L])$ is a smooth curve, $c \in (2\pi^2/L^2, +\infty) \rightarrow \eta = \eta(c)$ is a strictly increasing function, and its derivative with respect to the velocity c is given by

$$\frac{d\eta}{dc} = \frac{\eta}{2c} + \frac{k^2 k'^2 \eta^3 (2 - k^2) K}{\sqrt{c^3} ((2 - k^2)E - 2(1 - k^2)K)}.$$

The next result gives us a relation between the velocity of the solitary and periodic waves associated with the mKdV, and it will be useful later.

THEOREM 5.2. *We consider the mapping $\omega : (2\pi^2/L^2, +\infty) \mapsto \mathbb{R}$, given by*

$$(5.14) \quad \omega(c) = c / (16(2 - k^2)K'^2(k));$$

then $\frac{d\omega}{dc} > 0$.

Proof. Indeed, $\frac{d\omega}{dc} = [4(2 - k^2)K'^2 + 8K'c \frac{dk}{dc} (kK' - (2 - k^2) \frac{dK'}{dk})] / [16(2 - k^2)^2 K'^4]$. Since $\frac{dK'}{dk} = -(\frac{E' - k^2 K'}{kk'^2}) < 0$, it suffices to show that $\frac{dk}{dc} > 0$. So, since

$$\frac{dk}{dc} = \frac{2}{\eta^3} \left(2c \frac{d\eta}{dc} - \eta \right),$$

we need to verify the sign of the expression $2c \frac{d\eta}{dc} - \eta$. Next, we calculate the exact value of $\frac{d\eta}{dc}$. In fact,

$$\frac{d\eta}{dc} = \frac{\eta}{2c} + \underbrace{\frac{k^2 k'^2 \eta^3 (2 - k^2) K}{\sqrt{c^3} ((2 - k^2)E - 2(1 - k^2)K)}}_A.$$

Thus, $2c \frac{d\eta}{dc} - \eta = 2cA$ and therefore $\frac{dk}{dc} > 0$. This fact completes the proof of the theorem. \square

Next, we will show that $\widehat{\varphi}_c > 0$ and $K = \widehat{\varphi}_c^2$ belongs to $PF(2)$ discrete. It is easy to see that $\widehat{\varphi}_c > 0$ because of the form of the Fourier coefficients of φ_c given by (5.9). Moreover, $\widehat{\varphi}_c \in PF(2)$ discrete because the function $f(x) = \mu \text{sech}(\nu x)$ belongs to $PF(2)$ continuous (see [2]). So, since the convolution of even sequences in $PF(2)$ discrete is a sequence in $PF(2)$ discrete (see [32]) we can conclude that $K = \widehat{\varphi}_c^2 \in PF(2)$ discrete. Now, by choosing $\chi = -\frac{d}{dc} \varphi_c$ we have that $\mathcal{L}\chi = \varphi_c$. Then, by the Parseval theorem, $I = -\frac{1}{2} \frac{d}{dc} \|\varphi_c\|_{L_{per}^2}^2 = -\frac{L}{2} \frac{d}{dc} \|\widehat{\varphi}_c\|_{\ell^2}^2$. Since $\|\widehat{\varphi}_c\|_{\ell^2}^2 =$

$2\frac{\pi^2}{L^2} \sum_{n=-\infty}^{+\infty} \operatorname{sech}^2\left(\frac{\pi n}{\sqrt{\omega(c)}L}\right)$, we have that

$$\frac{d}{dc} \|\widehat{\varphi}_c\|_{\ell^2}^2 = \frac{C_1(L)}{\sqrt{\omega(c)}^3} \frac{d\omega}{dc} \sum_{n=-\infty}^{+\infty} \operatorname{sech}^2\left(\frac{\pi n}{\sqrt{\omega(c)}L}\right) n \operatorname{tgh}\left(\frac{\pi n}{\sqrt{\omega(c)}L}\right).$$

Since $(n \operatorname{tgh}(\frac{\pi n}{\sqrt{\omega(c)}L}))_{n \in \mathbb{Z}}$ is a positive sequence we have from Theorem 5.2 that $\frac{d}{dc} \|\widehat{\varphi}_c\|_{\ell^2}^2 > 0$. Thus, we obtain that the dnoidal wave φ_c in (5.14) is stable in $H^1_{per}([0, L])$ by the flow of the mKdV.

5.3. Stability of periodic travelling-wave solutions for the KdV equation. Now, we apply the results obtained previously to the proof of the stability of periodic travelling-wave solutions of cnoidal type associated with the KdV equation and satisfying

$$(5.15) \quad \varphi_c'' + \frac{1}{2}\varphi_c^2 - c\varphi_c = 0.$$

We consider the solitary wave solutions $\phi_\omega(x) = 3\omega \operatorname{sech}^2(\frac{\sqrt{\omega}x}{2})$, whose Fourier transform is given by $\widehat{\phi}_\omega(x) = \frac{12\pi x}{\sinh(\frac{\pi x}{\sqrt{\omega}})}$; then from the Poisson summation theorem we consider

$$(5.16) \quad \psi_\omega(\xi) = \frac{12\sqrt{\omega}}{L} + \frac{12\pi}{L^2} \sum_{n \neq 0} n \operatorname{csch}\left(\frac{\pi n}{\sqrt{\omega}L}\right) e^{\frac{2\pi i n \xi}{L}}.$$

Since ω is arbitrary, consider $\omega := \omega(k)$ such that $\sqrt{\omega(k)} = \frac{K(k)}{K(k')L}$, $k \in (0, 1)$, $k'^2 = 1 - k^2$. Then, we obtain

$$(5.17) \quad \psi_{\omega(k)}(\xi) = \frac{12\sqrt{\omega}}{L} + \frac{24\pi}{L^2} \sum_{n=1}^{+\infty} n \operatorname{csch}\left(\frac{\pi n K'}{K}\right) \cos\left(\frac{2\pi n \xi}{L}\right).$$

Now, we invoke the Fourier expansion of dn^2 (see [19], [42]), that is,

$$K^2 \left(\operatorname{dn}^2\left(\frac{2K\xi}{L}; k\right) - \frac{E}{K} \right) = 2\pi \sum_{n=1}^{+\infty} \frac{nq^n}{1 - q^{2n}} \cos\left(\frac{2\pi n \xi}{L}\right),$$

where $q = e^{-\frac{\pi K'}{K}}$. We can conclude that

$$\frac{q^n}{1 - q^{2n}} = \frac{1}{2} \operatorname{csch}\left(\frac{n\pi K'}{K}\right).$$

Then, we get from (5.17)

$$(5.18) \quad \psi_{\omega(k)}(\xi) = \frac{12\sqrt{\omega(k)}}{L} + \frac{24K^2}{L^2} \left(\operatorname{dn}^2\left(\frac{2K\xi}{L}; k\right) - \frac{E}{K} \right)$$

for $k \in (0, 1)$.

Next, because of the equality (5.18), we consider $\varphi_c(\xi) = a + b \left(\operatorname{dn}^2(dx; k) - \frac{E}{K} \right)$ a periodic travelling-wave solution for (5.15) of period L . Then, the following nonlinear system is obtained:

$$(5.19) \quad \begin{cases} \frac{b^2}{2} - 6d^2b = 0, & 4bd^2(1 + k'^2) + ab - b^2 \frac{E}{K} - cb = 0, \\ \frac{a^2}{2} - \frac{abE}{K} + \frac{b^2}{2} \left(\frac{E}{K} \right)^2 - ac - cb \frac{E}{K} - 2bd^2k'^2 = 0. \end{cases}$$

Since φ_c is periodic of period L it follows that $d = \frac{2K(k)}{L}$. Then, from the first equation of the system above we have that $b = \frac{48K^2}{L^2}$. Substituting those values at the second equation we get

$$(5.20) \quad c = \frac{16K}{L^2} [(1 + k'^2)K - 3E] + a.$$

From the third equation in (5.19) and the value of c in (5.20) we have the quadratic equation in terms of a ,

$$(5.21) \quad a^2 + \frac{32K}{L^2} [(1 + k'^2)K - 3E] a - \frac{(1 + k'^2)1536K^3E}{L^4} + \frac{768K^4k'^2}{L^4} + \frac{2304K^2E^2}{L^4} = 0,$$

whose positive solution is

$$a = -\frac{16K}{L^2} [(1 + k'^2)K - 3E] + \frac{16K^2}{L^2} \sqrt{1 - k^2 + k^4}.$$

Thus, the value of c is $c = \frac{16K^2}{L^2} \sqrt{1 - k^2 + k^4}$. Hence, for $k \in (0, 1)$ we have that $c \in (\frac{4\pi^2}{L^2}, +\infty)$. Therefore, writing φ_c in a convenient form, in terms of cn^2 , we obtain

$$\varphi_c(\xi) = \frac{16K^2}{L^2} \left[\sqrt{1 - k^2 + k^4} + 1 - 2k^2 \right] + \frac{48K^2k^2}{L^2} \operatorname{cn}^2 \left(\frac{2K}{L} \xi; k \right).$$

We can see that this formula is the same that was obtained by Angulo in [7] and it can be rewritten as

$$(5.22) \quad \varphi_c(\xi) = \beta_2 + (\beta_3 - \beta_2) \operatorname{cn}^2 \left(\sqrt{\frac{\beta_3 - \beta_1}{12}} \xi; k \right),$$

where

$$\beta_2 = \frac{16K^2}{L^2} \left[\sqrt{1 - k^2 + k^4} + 1 - 2k^2 \right], \quad \beta_3 = \frac{16K^2}{L^2} \left[\sqrt{1 - k^2 + k^4} + 1 + k^2 \right],$$

and β_1 is such that

$$\beta_3 - \beta_1 = \frac{48K^2}{L^2}.$$

By making a similar analysis such as in the case of the mKdV equation, we can obtain a smooth curve of positive cnoidal waves with the form in (5.22), $c \in$

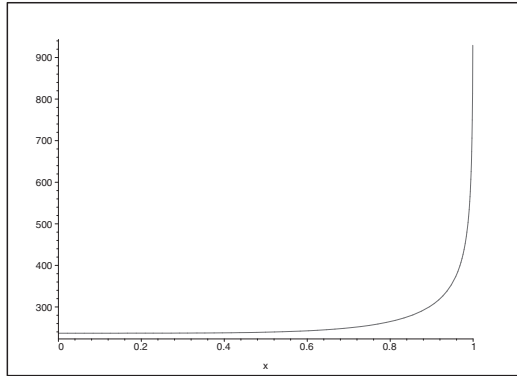


FIG. 1. Graphic of the function $a(k)$ with period $L = 1$.

$(\frac{4\pi^2}{L^2}, +\infty) \mapsto \varphi_c \in H_{per}^n([0, L])$, such that $k := k(c)$ is a strictly increasing smooth function (see [7]) for all $n \in \mathbb{N}$. Moreover, we can determine that for $k \in (0, 1)$ there is a unique $c \in (\frac{4\pi^2}{L^2}, +\infty)$ such that $k(c) = k$. Therefore, the function $\omega(k)$ defined above can be expressed as a function of c , $\omega = \omega(k(c))$, and it is a strictly increasing function (it will be seen later). Then, since $\frac{K(k)}{K(k')}$ $\in (0, +\infty)$ it follows that for $c \in (\frac{4\pi^2}{L^2}, +\infty)$ we obtain $\omega(k(c)) \in (0, +\infty)$. Therefore, the mapping $c \in (\frac{4\pi^2}{L^2}, +\infty) \mapsto \psi_{\omega(k(c))} \in H_{per}^n([0, L])$ is a smooth curve for all $n \in \mathbb{N}$. Next, we note that $2\psi_{\omega(k(c))}(\xi) - \varphi_c(\xi) = \frac{24\sqrt{\omega(k(c))}}{L} - a(k(c))$, where for $k = k(c)$

$$a(k) = \frac{16K^2}{L^2} \left[\sqrt{1 - k^2 + k^4} + 2 - k^2 + 3\frac{E}{K} \right].$$

Thus, for $s(k(c)) \equiv a(k(c)) - \frac{24\sqrt{\omega(k(c))}}{L}$, we can write $\varphi_c(\xi) \equiv s(k(c)) + \psi_{\omega(k(c))}(\xi)$. Therefore, we obtain immediately that the Fourier coefficients of φ_c are

$$\widehat{\varphi}_c(n) = \begin{cases} a(k), & n = 0, \\ \frac{12\pi}{L^2} n \operatorname{csch} \left(\frac{\pi n}{\sqrt{\omega(k)}L} \right), & n \neq 0. \end{cases}$$

Remark 5.3. After some calculations, we can obtain that $s(k)$ is a positive function defined in $(0, 1)$. Making use of Maple, we can also determine that $s(k)$ does not have any root on the extremes of the interval $(0, 1)$. We can also determine that the function $a(k)$ is a positive strictly increasing function (see Figure 1).

Since $s(k) > 0$ and the function $f : \mathbb{R} \rightarrow \mathbb{R}$ defined by $f(x) = \frac{12\pi}{L^2} x \operatorname{csch}(\frac{\pi x}{\sqrt{\omega}L})$, belongs to $PF(2)$ in the continuous case, we can use Lemma 4.1 for obtaining that $\widehat{\varphi}_c$ belongs to the $PF(2)$ discrete case. Indeed, since

$$a(k) > \frac{24\sqrt{\omega(k)}}{L} > \frac{12\sqrt{\omega(k)}}{L} = f(0) > f(x), \quad x \neq 0,$$

we can redefine f by a smooth function $h : \mathbb{R} \rightarrow \mathbb{R}$ such that $h(0) = a(k)$, $h(x) \equiv f(x)$ on $(-\infty, -1] \cup [1, +\infty)$ and on the interval $(-1, 1)$ we “complete” f in a differentiable

way, such that h belongs to $PF(2)$ continuous. Therefore, the sequence to be obtained (if we look only at the set of integers numbers) will be $h(n) = \widehat{\varphi}_c(n)$.

Next, let $\chi = -\frac{d}{dc}\varphi_c$ such that $\mathcal{L}\chi = \varphi_c$. Then by the Parseval theorem, it follows that $I = -\frac{L}{2} \frac{d}{dc} (\|\widehat{\varphi}_c\|_{\ell^2}^2)$. Hence,

$$\begin{aligned} \frac{d}{dc} \|\widehat{\varphi}_c\|_{L^2_{per}}^2 &= C_1 a(k) \frac{da}{dk} \frac{dk}{dc} \\ &+ \frac{C_2}{\sqrt{\omega(k)^3}} \frac{d\omega}{dk} \frac{dk}{dc} \underbrace{\sum_{\substack{n=-\infty \\ n \neq 0}}^{+\infty} n^3 \operatorname{csch}^2\left(\frac{\pi n}{L\sqrt{\omega(k)}}\right) \coth\left(\frac{\pi n}{L\sqrt{\omega(k)}}\right)}_{b_n}, \end{aligned}$$

where $C_1 := C_1(L)$, $C_2 := C_2(L) > 0$. Next, we need to show only that the quantities $\frac{da}{dk}$ and $\frac{d\omega}{dk}$ are positive because $k := k(c)$ is a strictly increasing function and $(b_n)_{n \in \mathbb{Z}}$ is obviously a positive sequence. Hence, we have

$$\frac{d\omega}{dk} = 2 \frac{K \left(\frac{dK}{dk} K' - K \frac{dK'}{dk} \right)}{K'^3}.$$

Since $\frac{dK}{dk} > 0$ and $\frac{dK'}{dk} < 0$ we get that $\frac{d\omega}{dk} > 0$. By making use of a similar argument, we can also show that $\frac{da}{dk} > 0$ because we have that

$$a(k) = \frac{16K(k)^2}{L^2} \left[\sqrt{1 - k^2 + k^4} + 2 - k^2 \right] + 48 \frac{E(k)K(k)}{L^2}.$$

Therefore, $I < 0$ and the positive cnoidal waves φ_c are stable in $H^1_{per}([0, L])$ by the periodic flow of the KdV equation.

6. Comments. In this section we make some basic remarks about the results contained in the body of this paper.

6.1. General perturbations. In contrast to the case of solitary waves for which the natural class of disturbance in the stability problem is that of localized disturbances, for periodic waves there are several classes of disturbance for which stability needs to be addressed. Here we consider periodic perturbations with the same fundamental period of the periodic travelling waves. For disturbance, for example, with a double period, stability results remain open in the context of KdV-type equations. If we consider the KdV equation, our stability result in $H^1_{per}([0, L])$ was based on the spectral structure of the operator $\mathcal{L}_{cn} = -\frac{d^2}{dx^2} + c - \varphi_c$, with φ_c defined in (5.22). Now, if we consider this operator with domain $H^2_{per}([0, 2L])$, then the number of negative eigenvalues will be exactly 3 (see [10]). Moreover, since the function $\frac{d}{dc} \int_{-L}^L \varphi_c^2(x) dx$ is even positive, the abstract setting of Grillakis, Shatah, and Strauss in [27] and [28] cannot be applied. We note that in the case of the focusing Schrödinger equation

$$iu_t + u_{xx} + |u|^2 u = 0,$$

Angulo in [9] showed an instability result for dnoidal waves solutions when the class of periodic disturbance is two times the minimal period of the periodic travelling wave in question (see also Gally and Hărăguş [24], [25] for recent new results of stability for periodic travelling waves of Schrödinger equations).

6.2. Property $PF(2)$ discrete. We consider the BO equation (but we can do an analogous analysis with the other two equations); then we have seen throughout this paper that the spectral properties of the operator $\mathcal{L} = \mathcal{H}\partial_x + c - \varphi_c$ were obtained from the fact that the Fourier coefficients of the periodic travelling-wave solution given by (1.13) are in $PF(2)$. The Poisson summation theorem was used to find such a wave with a minimal period $2L$. Moreover, the Fourier coefficients were calculated with this period. If we double the period, that is, if we consider \mathcal{L} with domain $D(\mathcal{L}) = H_{4L}^1$, the property $PF(2)$ is not satisfied to $\tilde{K} = \widehat{\varphi}_c^{(4L)}$, where $\widehat{\varphi}_c^{(4L)}$ denotes the periodic Fourier transform of φ_c but with period $4L$. Indeed, we consider the Fourier expansion in the form $\varphi_c(x) = \sum_{n=0}^{+\infty} \widehat{\varphi}_c^{(4L)}(n) \cos\left(\frac{n\pi x}{2L}\right)$. Since $\varphi_c(0) = \varphi_c(2L)$, $\sum_{n=0}^{+\infty} \widehat{\varphi}_c^{(4L)}(n)(\cos(n\pi) - 1) = 0$. Thus, $-2 \sum_{\substack{n=2k+1 \\ k \in \mathbb{N}}} \widehat{\varphi}_c^{(4L)}(n) = 0$. Therefore, $\tilde{K} = \widehat{\varphi}_c^{(4L)}$ cannot belong to $PF(2)$. Then, we cannot affirm anything about the stability of the wave φ_c when this case is considered. So, Theorem 4.1 cannot be applied here.

In Theorem 4.1, the Fourier transform needs to be evaluated in the minimal period for the solution φ_c . In fact, let L be this minimal period, and we evaluate the Fourier transform of φ_c as being of period $2L$; then

$$\begin{aligned} \widehat{\varphi}_c(k) &= \frac{1}{2L} \int_{-L}^L \varphi_c(x) e^{-\frac{ik\pi x}{L}} dx = \frac{1}{2L} \int_{-L}^L \varphi_c(x+L) e^{-\frac{ik\pi x}{L}} dx \\ &= \frac{1}{2L} \int_0^{2L} \varphi_c(y) e^{-\frac{ik\pi y}{L}} e^{i\pi k} dy = (-1)^k \widehat{\varphi}_c(k). \end{aligned}$$

Then, for k being odd we have $\widehat{\varphi}_c(k) = 0$. Thus, we cannot apply our theory if the minimal period is not fixed.

6.3. Positivity of the periodic travelling waves. Another fact that deserves a special mention is about the condition that the solution φ_c in Theorem 4.1 needs to be a positive solution. Indeed, the classical Fourier theorem (see [30]) told us that this is necessary. Suppose without loss of generality that $\varphi_c(0) = 0$. Then, φ_c being smooth it follows that $\varphi_c^p(0) = \sum_{n=-\infty}^{+\infty} \widehat{\varphi}_c^p(n) = 0$. In other words some Fourier coefficient of φ_c^p must be negative.

6.4. Stability and instability of periodic travelling waves for the critical KdV and NLS. In a forthcoming paper [11] we apply the theory in section 4 to obtain the stability/instability of a special family of periodic travelling-wave solutions for the critical KdV equation

$$u_t + 5u^4 u_x + u_{xxx} = 0$$

and for the critical nonlinear Schrödinger (NLS) equation

$$iu_t + u_{xx} + |u|^4 u = 0.$$

7. Appendix.

7.1. Jacobi elliptic functions. We establish some basic properties of Jacobian elliptic integrals (see [18] and [19]). The normal elliptic integral of the first kind is

$$\int_0^y \frac{dt}{\sqrt{(1-t^2)(1-k^2 t^2)}} = \int_0^\varphi \frac{d\theta}{\sqrt{1-k^2 \sin^2 \theta}} \equiv F(\varphi, k),$$

where $y = \sin \varphi$, whereas the normal elliptic integral of the second kind is

$$\int_0^y \sqrt{\frac{1 - k^2 t^2}{1 - t^2}} dt = \int_0^\varphi \sqrt{1 - k^2 \sin^2 \theta} d\theta \equiv E(\varphi, k).$$

The number k is called the modulus and belongs to the interval $(0, 1)$. The number $k' = \sqrt{1 - k^2}$ is called the complementary modulus. The parameter φ is called the argument of the normal elliptic integrals. It is usually understood that $0 \leq y \leq 1$ or $0 \leq \varphi \leq \frac{\pi}{2}$. For $y = 1$, the integrals above are said to be complete. In this case, one writes

$$\int_0^1 \frac{dt}{\sqrt{(1 - t^2)(1 - k^2 t^2)}} = \int_0^{\frac{\pi}{2}} \frac{d\theta}{\sqrt{1 - k^2 \sin^2 \theta}} \equiv F\left(\frac{\pi}{2}, k\right) \equiv K(k) \equiv K$$

and

$$\int_0^1 \sqrt{\frac{1 - k^2 t^2}{1 - t^2}} dt = \int_0^{\frac{\pi}{2}} \sqrt{1 - k^2 \sin^2 \theta} d\theta \equiv E\left(\frac{\pi}{2}, k\right) \equiv E(k) \equiv E.$$

Clearly, we have $K(0) = E(0) = \frac{\pi}{2}$, while $E(1) = 1$ and $K(1) = +\infty$. For $k \in (0, 1)$, $\frac{dK}{dk} > 0$, $\frac{d^2K}{dk^2} > 0$, $\frac{dE}{dk} < 0$, $\frac{d^2E}{dk^2} < 0$, and $E(k) < K(k)$. Moreover, $E(k) + K(k)$ and $E(k)K(k)$ are strictly increasing functions for every $k \in (0, 1)$. Next, we have some derivatives of the complete elliptical integrals K and E used in this work,

$$\frac{dK}{dk} = \frac{E - k'^2 K}{kk'^2}, \quad \frac{dE}{dk} = \frac{E - K}{k}.$$

The Jacobian elliptic functions are usually defined as follows. It considers the elliptic integral

$$u(y_1; k) \equiv u = \int_0^{y_1} \frac{dt}{\sqrt{(1 - t^2)(1 - k^2 t^2)}} = \int_0^\varphi \frac{d\theta}{\sqrt{1 - k^2 \sin^2 \theta}} \equiv F(\varphi, k),$$

which is a strictly increasing function of the variable y_1 . Its inverse function is written $y_1 = \sin \varphi \equiv \text{sn}(u; k)$, or briefly $y_1 = \text{snu}$ when it is not necessary to emphasize the modulus k . The other two basic elliptic functions, the cnoidal and dnoidal functions, are defined in terms of sn by

$$\text{cn}(u; k) = \sqrt{1 - y_1^2} = \sqrt{1 - \text{sn}^2(u; k)}, \quad \text{dn}(u; k) = \sqrt{1 - k^2 y_1^2} = \sqrt{1 - k^2 \text{sn}^2(u; k)}.$$

Note that these functions are normalized by the requirements $\text{sn}(0; k) = 0$, $\text{cn}(0; k) = 1$, and $\text{dn}(0; k) = 1$. The functions $\text{cn}(\cdot; k)$ and $\text{dn}(\cdot; k)$ are even functions. These functions are periodic with $\text{sn}(u + 4K(k); k) = \text{sn}(u; k)$, $\text{cn}(u + 4K(k); k) = \text{cn}(u; k)$, $\text{dn}(u + 2K(k); k) = \text{dn}(u; k)$. Moreover, we have the relations $\text{sn}^2 u + \text{cn}^2 u = 1$, $k^2 \text{sn}^2 u + \text{dn}^2 u = 1$, $k'^2 \text{sn}^2 u + \text{cn}^2 u = \text{dn}^2 u$, $\text{sn}(u + 2K; k) = -\text{sn}(u; k)$, $\text{cn}(u + 4K; k) = -\text{cn}(u; k)$. We also have the following explicit values: $\text{sn}(0) = 0$, $\text{cn}(0) = 1$, $\text{sn}(K) = 0$, $\text{cn}(K) = 0$ and the asymptotic behaviors $\text{sn}(u; 0) = \sin u$, $\text{cn}(u; 0) = \cos u$, $\text{sn}(u; 1) = \tanh u$, $\text{cn}(u; 1) = \text{sech} u$. Finally, the formulas

$$\frac{\partial}{\partial u} \text{snu} = \text{cnudnu}, \quad \frac{\partial}{\partial u} \text{cnu} = -\text{snudnu}, \quad \frac{\partial}{\partial u} \text{dnu} = -k^2 \text{cnusnu}$$

are straightforwardly deduced from the foregoing material.

7.2. A nonlocal Poincaré–Wirtinger inequality. The following inequality was used in Proposition 5.1. Suppose $f \in H_{2L}^{\frac{1}{2}}$ such that $\int_{-L}^L f(x)dx = 0$; then for $D = \mathcal{H}\partial_x$

$$\int_{-L}^L [D^{\frac{1}{2}}f]^2 dx \geq \frac{\pi}{L} \int_{-L}^L f^2 dx.$$

Proof. Let f be in \mathcal{P} , and consider the Fourier expansion of f given by

$$f(x) = \sum_{n=-\infty}^{+\infty} a_n e^{\frac{in\pi x}{L}} = \sum_{n=-\infty}^{+\infty} \widehat{f}(n) e^{\frac{in\pi x}{L}},$$

where $\widehat{f}(n) = \frac{1}{2L} \int_{-L}^L f(x) e^{-\frac{in\pi x}{L}} dx$, with $\widehat{f}(0) = 0$. Then, $\mathcal{H}\partial_x f(x) = \frac{\pi}{L} \sum_{n=-\infty}^{+\infty} a_n |n| e^{\frac{in\pi x}{L}}$. From the Parseval theorem we obtain

$$\int_{-L}^L f \mathcal{H}\partial_x f dx = 2L \sum_{n=-\infty}^{+\infty} \frac{\pi}{L} |a_n|^2 |n| = 2L \sum_{n \neq 0} \frac{\pi}{L} |a_n|^2 |n|$$

and

$$\int_{-L}^L f^2 dx = 2L \sum_{n \neq 0} |a_n|^2.$$

Therefore,

$$\int_{-L}^L f \mathcal{H}\partial_x f dx = 2L \sum_{n \neq 0} \frac{\pi}{L} |a_n|^2 |n| \geq \frac{\pi}{L} \int_{-L}^L f^2 dx.$$

So, using density arguments, we can show that the inequality above occurs for $f \in H_{2L}^{\frac{1}{2}}$. This completes the proof of the inequality. \square

Acknowledgments. The second author would like to express his thanks to the Institute of Mathematics and Statistic (IME) of the University of São Paulo/SP-Brazil for its hospitality. The authors also thank the referees for some helpful comments.

REFERENCES

- [1] L. ABDELOUHAB, J. BONA, M. FELLAND, AND J. C. SAUT, *Nonlocal models for nonlinear, dispersive wave*, Phys. D, 40 (1989), pp. 360–392.
- [2] J. P. ALBERT, *Positivity properties and stability of solitary-wave solutions of model equations for long waves*, Comm. Partial Differential Equations, 17 (1992), pp. 1–22.
- [3] J. P. ALBERT AND J. L. BONA, *Total positivity and the stability of internal waves in fluids of finite depth*, IMA J. Appl. Math., 46 (1991), pp. 1–19.
- [4] J. P. ALBERT, J. L. BONA, AND D. HENRY, *Sufficient conditions for stability of solitary-wave equation of model equations for long waves*, Phys. D, 24 (1987), pp. 343–366.
- [5] B. ALVAREZ-SAMANIEGO AND D. LANNES, *Large time existence for 3D water-waves and asymptotics*, Invent. Math., 171 (2008), pp. 485–541.
- [6] J. ANGULO, *Existence and Stability of Solitary Wave Solutions to Non-linear Dispersive Evolution Equations*, Publicações do 24º Colóquio Brasileiro de Matemática, IMPA, Rio de Janeiro, Brazil, 2003.
- [7] J. ANGULO, *Stability of cnoidal waves to Hirota-Satsuma systems*, Mat. Contemp., 27 (2004), pp. 189–223.

- [8] J. ANGULO, *Stability of dnoidal waves to Hirota-Satsuma system*, Differential Integral Equations, 18 (2005), pp. 611–645.
- [9] J. ANGULO, *Non-linear stability of periodic travelling-wave equation for the Schrödinger and modified Korteweg-de Vries equation*, J. Differential Equations, 235 (2007), pp. 1–30.
- [10] J. ANGULO, J. L. BONA, AND M. SCIALOM, *Stability of cnoidal waves*, Adv. Differential Equations, 11 (2006), pp. 1321–1374.
- [11] J. ANGULO AND F. NATALI, *Stability and Instability of Periodic Travelling Waves Wave Solutions for the Critical Korteweg-de Vries and Non-linear Schrödinger Equations*, preprint, 2007.
- [12] T. B. BENJAMIN, *Instability of periodic wavetrain in nonlinear dispersive systems*, Proc. Roy. Soc. London Ser. A, 299 (1967), pp. 59–75.
- [13] T. B. BENJAMIN, *Internal waves of permanent form in fluids of great depth*, J. Fluid Mech., 29 (1967), pp. 559–592.
- [14] T. B. BENJAMIN, *The stability of solitary waves*, Proc. Roy. Soc. London Ser. A, 338 (1972), pp. 153–183.
- [15] T. B. BENJAMIN, *Lectures on linear wave motion*, in Nonlinear Wave Motion, A. C. Newell, ed., AMS, Providence, RI, 1974, pp. 3–47.
- [16] T. B. BENJAMIN, *Solitary and periodic waves of a new kind*, Philos. Trans. Roy. Soc. London Ser. A, 354 (1996), pp. 1775–1806.
- [17] J. L. BONA, *On the stability theory of solitary waves*, Proc. Roy. Soc. London Ser. A, 344 (1975), pp. 363–374.
- [18] F. BOWMAN, *Introduction to Elliptic Functions with Applications*, Dover, New York, 1961.
- [19] P. F. BYRD AND M. D. FRIEDMAN, *Handbook of Elliptic Integrals for Engineers and Scientists*, 2nd ed., Springer, New York, 1971.
- [20] J. COLLIANDER, M. KEEL, G. STAFILANI, H. TAKAOKA, AND T. TAO, *Sharp global well-posedness for the KDV and modified KDV on \mathbb{R} and \mathbb{T}* , J. Amer. Math. Soc., 16 (2003), pp. 705–749.
- [21] A. CONSTANTIN, *The trajectories of particles in Stokes waves*, Invent. Math., 166 (2006), pp. 523–535.
- [22] A. CONSTANTIN AND J. ESCHER, *Particle trajectories in solitary water waves*, Bull. Amer. Math. Soc. (N.S.), 44 (2007), pp. 423–431.
- [23] A. CONSTANTIN AND W. STRAUSS, *Stability properties of steady water waves with vorticity*, Comm. Pure Appl. Math., 60 (2007), pp. 911–950.
- [24] T. GALLAY AND M. HÄRÄĞUŞ, *Stability of small periodic waves for the nonlinear Schrödinger equation*, J. Differential Equations, 234 (2007), pp. 544–581.
- [25] T. GALLAY AND M. HÄRÄĞUŞ, *Orbital stability of periodic waves for the nonlinear Schrödinger equation*, J. Dynam. Differential Equations, 19 (2007), pp. 825–865.
- [26] R. A. GARDNER, *Spectral analysis of long wavelength periodic waves and applications*, J. Reine Angew. Math., 491 (1997), pp. 149–181.
- [27] M. GRILLAKIS, J. SHATAH, AND W. STRAUSS, *Stability theory of solitary waves in the presence of symmetry I*, J. Funct. Anal., 74 (1987), pp. 160–197.
- [28] M. GRILLAKIS, J. SHATAH, AND W. STRAUSS, *Stability theory of solitary waves in the presence of symmetry II*, J. Funct. Anal., 94 (1990), pp. 308–348.
- [29] E. L. INCE, *The periodic Lamé functions*, Proc. Roy. Soc. Edinburgh, 60 (1940), pp. 47–63.
- [30] R. J. IORIO, JR., AND V. M. V. IORIO, *Fourier Analysis and Partial Differential Equations*, Cambridge Stud. Adv. Math. 70, Cambridge University Press, Cambridge, UK, 2001.
- [31] S. KARLIN, *The existence of eigenvalues for integral operators*, Trans. Amer. Math. Soc., 113 (1964), pp. 1–17.
- [32] S. KARLIN *Total Positivity*, Stanford University Press, Stanford, CA, 1968.
- [33] T. KATO, *Perturbation Theory for Linear Operators*, 2nd ed., Springer, Berlin, 1976.
- [34] D. J. KORTEWEG AND G. DE VRIES, *On the change of form of long wave advancing in a rectangular canal, and on a new type of long stationary waves*, Philos. Mag., 39 (1895), pp. 422–443.
- [35] W. MAGNUS AND E. OBERHETTINGER, *Formulas and Theorems for the Special Functions of Mathematical Physics*, Springer, New York, 1986.
- [36] W. MAGNUS AND S. WINKLER, *Hill's Equation*, Tracts in Pure and Appl. Math. 20, Wiley, New York, 1976.
- [37] H. P. MCKEAN, *Stability for the Korteweg-de Vries equation*, Comm. Pure Appl. Math., 30 (1997), pp. 347–353.
- [38] H. P. MCKEAN AND H. DYM, *Fourier Series and Integrals*, Academic Press, New York, London, 1972.
- [39] L. MOLINET, *Global well-posedness in L^2 for the periodic Benjamin-Ono equation*, Amer. J. Math., 130 (2008), pp. 635–683.

- [40] L. MOLINET, *Global well-posedness in the energy space for the Benjamin-Ono equation on the circle*, Math. Ann., 337 (2007), pp. 353–383.
- [41] L. MOLINET AND F. RIBAUD, *Well-Posedness in H^1 for the (Generalized) Benjamin-Ono Equation on the Circle*, preprint, 2006.
- [42] F. OBERHETTINGER, *Fourier Expansions: A Collection of Formulas*, Academic Press, New York, London, 1973.
- [43] H. ONO, *Algebraic solitary waves in stratified fluids*, J. Phys. Soc. Japan, 39 (1975), pp. 1082–1091.
- [44] A. R. OSBORNE, M. SERIO, L. BERGAMASCO, AND L. CAVALERI, *Solitons, cnoidal waves and nonlinear interactions in shallow-water ocean surface waves*, Phys. D, 123 (1998), pp. 64–81.
- [45] E. M. STEIN, *Singular Integrals and Differentiability Properties of Functions*, Princeton University Press, Princeton, NJ, 1970.
- [46] E. M. STEIN AND G. WEISS, *Introduction to Fourier Analysis on Euclidean Spaces*, Princeton University Press, Princeton, NJ, 1970.
- [47] J. F. TOLAND, *Stokes waves*, Topol. Methods Nonlinear Anal., 7 (1996), pp. 1–48.
- [48] M. I. WEINSTEIN, *Existence and dynamic stability of solitary wave solutions of equations arising in long wave propagation*, Comm. Partial Differential Equations, 12 (1987), pp. 1133–1173.
- [49] M. I. WEINSTEIN, *Lyapunov stability of ground states of nonlinear dispersive evolution equations*, Comm. Pure Appl. Math., 39 (1986), pp. 51–67.

INHOMOGENEOUS BOUNDARY VALUE PROBLEMS FOR COMPRESSIBLE NAVIER–STOKES EQUATIONS: WELL-POSEDNESS AND SENSITIVITY ANALYSIS*

P. I. PLOTNIKOV[†], E. V. RUBAN[†], AND J. SOKOLOWSKI[‡]

Abstract. In this paper compressible, stationary Navier–Stokes equations are considered. A framework for analysis of such equations is established. In particular, the well-posedness for inhomogeneous boundary value problems of elliptic-hyperbolic type is shown. Analysis is performed for small perturbations of the so-called *approximate solutions* that take form (1.12). The approximate solutions are determined from Stokes problem (1.11). The small perturbations are given by solutions to (1.13). The uniqueness of solutions for problem (1.13) is proved, and, in addition, the differentiability of solutions with respect to the coefficients of differential operators is shown. The results on the well-posedness of nonlinear problems are interesting on their own and are used to obtain the shape differentiability of the drag functional for incompressible Navier–Stokes equations. The shape gradient of the drag functional is derived in the classical and useful for computations form; an appropriate adjoint state is introduced to this end. The *material* derivatives of solutions to the Navier–Stokes equations are given by smooth functions; however, the shape differentiability is shown in a weak norm. The method of analysis proposed in this paper is general and can be used to establish the well-posedness for distributed and boundary control problems as well as for inverse problems in the case of the state equations in the form of compressible Navier–Stokes equations. The differentiability of solutions to the Navier–Stokes equations with respect to the data leads to the first order necessary conditions for a broad class of optimization problems.

Key words. Navier–Stokes equations, compressible fluids, shape optimization

AMS subject classifications. 35Q30, 49J20, 76N10

DOI. 10.1137/070694272

1. Introduction. Shape optimization for compressible Navier–Stokes equations is important for applications [27] and is investigated from a numerical point of view; however, the mathematical analysis of such problems is not available in the existing literature. One of the reasons is the lack of existence results for inhomogeneous boundary value problems for such equations.

The results established in this paper lead in particular to the first order optimality conditions for a class of shape optimization problems for compressible Navier–Stokes equations.

1.1. Problem formulation. In this paper we prove the well-posedness and perform the sensitivity analysis for inhomogeneous boundary value problems for the compressible Navier–Stokes equations. We restrict ourselves to the case of a specific shape optimization problem for stationary motion of viscous compressible non-heat-conducting isentropic gas. However, the technique of modelling and analysis presented here is general and can be used for a broad class of optimization problems for nonlinear elliptic-hyperbolic equations. The sensitivity analysis is the necessary

*Received by the editors June 11, 2007; accepted for publication (in revised form) June 20, 2008; published electronically October 22, 2008. This paper was prepared in the fall of 2006 during a visit of Pavel I. Plotnikov and Evgenya V. Ruban to the Institute Elie Cartan of the University Henri Poincaré Nancy 1.

<http://www.siam.org/journals/sima/40-3/69427.html>

[†]Lavryentyev Institute of Hydrodynamics, Siberian Division of Russian Academy of Sciences, Lavryentyev pr. 15, Novosibirsk 630090, Russia (plotnikov@hydro.nsc.ru, zhenya.ruban@mail.ru).

[‡]Institut Elie Cartan, Laboratoire de Mathématiques, Université Henri Poincaré Nancy 1, B.P. 239, 54506 Vandoeuvre lés Nancy Cedex, France (Jan.Sokolowski@iecn.u-nancy.fr).

step for numerical methods of solution for optimization problems. In general the mathematical analysis of optimization problems includes the following steps, with the mathematical proofs of the required facts:

- existence of solutions,
- uniqueness and optimality conditions,
- numerical method of solution.

The existence of an optimal shape for the drag minimization is shown in [40] under assumptions that are compared to the assumptions in the present paper. Here we present the necessary mathematical tools required for the second step of analysis, i.e., the derivation of an optimality system. In particular, we prove the shape differentiability of solutions to (1.9) and provide the classical representation of the shape derivatives of integral shape functionals in terms of an appropriate adjoint state.

We consider in detail all questions on the existence, uniqueness, and shape differentiability of solutions to stationary boundary value problems for compressible Navier–Stokes equations. Such boundary value problems can be regarded as the mathematical models of viscous gas flow around a body tested in the wind tunnel. We assume that the viscous gas occupies the double-connected domain $\Omega = B \setminus S$, where $B \subset \mathbb{R}^3$ is a hold-all domain with the smooth boundary $\Sigma = \partial B$, and $S \subset B$ is a compact obstacle. Furthermore, we assume that the velocity of the gas coincides with a given vector field $\mathbf{U} \in C^\infty(\mathbb{R}^3)^3$ on the surface Σ . In this framework, the boundary of the flow domain Ω is divided into three subsets, inlet Σ_{in} , outgoing set Σ_{out} , and characteristic set Σ_0 , which are defined by the equalities

$$(1.1) \quad \Sigma_{\text{in}} = \{x \in \Sigma : \mathbf{U} \cdot \mathbf{n} < 0\}, \quad \Sigma_{\text{out}} = \{x \in \Sigma : \mathbf{U} \cdot \mathbf{n} > 0\},$$

$\Sigma_0 = \{x \in \partial\Omega : \mathbf{U} \cdot \mathbf{n} = 0\}$, where \mathbf{n} stands for the unit outward normal to $\partial\Omega = \Sigma \cup \partial S$. In turn the compact $\Gamma = \Sigma_0 \cap \Sigma$ splits the surface Σ into three disjoint parts $\Sigma = \Sigma_{\text{in}} \cup \Sigma_{\text{out}} \cup \Gamma$. The problem is to find the velocity field \mathbf{u} and the gas density ϱ satisfying the following equations along with the boundary conditions:

$$(1.2a) \quad \Delta \mathbf{u} + \lambda \nabla \operatorname{div} \mathbf{u} = R \varrho \mathbf{u} \cdot \nabla \mathbf{u} + \frac{R}{\epsilon^2} \nabla p(\varrho) \quad \text{in } \Omega,$$

$$(1.2b) \quad \operatorname{div}(\varrho \mathbf{u}) = 0 \quad \text{in } \Omega,$$

$$(1.2c) \quad \mathbf{u} = \mathbf{U} \quad \text{on } \Sigma, \quad \mathbf{u} = 0 \quad \text{on } \partial S,$$

$$(1.2d) \quad \varrho = \varrho_0 \quad \text{on } \Sigma_{\text{in}},$$

where the pressure $p = p(\varrho)$ is a smooth, strictly monotone function of the density, ϵ is the Mach number, R is the Reynolds number, λ is the viscosity ratio, and ϱ_0 is a positive constant.

For the derivation of equations (1.2) we refer to [22]. The general theory of compressible Navier–Stokes equations is covered by the monographs [11], [25], and [33]. In particular, the main results on the existence of global weak solutions for stationary problems with the zero velocity boundary conditions were established in [25] and sharpened in [33]. See also [14], [38], [39] for generalizations.

There are numerous papers dealing with the zero velocity boundary value problem for steady compressible Navier–Stokes equations in the context of small data. We recall only that there are three different approaches to this problem proposed in [2], [35], and [30], [32], respectively. The basic results on the local existence and uniqueness of strong solutions are assembled in [33]. For an interesting overview see [36].

The inhomogeneous boundary problems were studied in papers [20] and [21], where the local existence and uniqueness results were obtained in the two-dimensional

case under the assumption that the velocity \mathbf{u} is close to a given constant vector. The question of the existence of strong solutions to boundary value problems in three spatial dimensions with nonzero velocity boundary data in smooth domains is still an open problem. There are difficulties including the problems of the total mass control and the singularities developed by solutions at the manifold $\overline{\Sigma}_{in} \cap \overline{\Sigma}_0 \cup \overline{\Sigma}_{out}$. In this paper we prove the local existence and uniqueness of strong solutions to problem (1.2) in fractional Sobolev spaces, under the assumption that the given vector field \mathbf{U} satisfies the emergent vector field conditions **(H1)**–**(H3)** on Γ . It seems that a condition of this type is necessary for the continuity of mass density ϱ .

Shape optimization problems. Among many shape optimization problems for Navier–Stokes equations, we could list the drag minimization problem, which is investigated in this paper and in [37], [38], [39], [40]. It is important to note that the drag from one side depends on the shape of the obstacle, and from the other side on the design of its surface structure as well as on the quality of its surface. The problem of the drag reduction includes the optimal design of the shape of the obstacle and control of the flow near its surface. The last question was investigated in paper [17], which also contains an overview of the problem.

Another problem of practical interest concerns the optimal shape of tunnels [27]. In the specific problem the required mass distribution on the outlet of the tunnel is given. The associated shape optimization problem can be formulated as follows. Determine an admissible domain such that the mass distribution at the inlet is given, the velocity field is prescribed on the boundary of the domain, and the mass distribution at the outlet is as close as possible to a given function. Inlet and outlet subsets are defined by the vector field \mathbf{U} which serves as the inhomogeneous boundary condition for the law of momentum conservation in the form of a Navier–Stokes stationary system. The shape optimization problem as it is formulated in [27] enters into our framework, and the results on shape sensitivity analysis can be applied to solve the problem. Another class of problems which can be investigated using the tools proposed in the paper are optimal control problems, e.g., with the boundary controls. These topics are, however, beyond the scope of the paper, and we present the drag minimization problem as an example to the general theory.

Drag minimization. One of the main applications of the theory of compressible viscous flows is the optimal shape design in aerodynamics. The classical sample is the problem of the minimization of the drag of airfoil travelling in atmosphere with uniform speed \mathbf{U}_∞ . Recall that in our framework the hydrodynamical force acting on the body S is defined by the formula [41]

$$\mathbf{J}(S) = - \int_{\partial S} \left(\nabla \mathbf{u} + (\nabla \mathbf{u})^* + (\lambda - 1) \operatorname{div} \mathbf{u} \mathbf{I} - \frac{R}{\epsilon^2} p \mathbf{I} \right) \cdot \mathbf{n} dS.$$

In a frame attached to the moving body the drag is the component of \mathbf{J} parallel to \mathbf{U}_∞ ,

$$(1.3) \quad J_D(S) = \mathbf{U}_\infty \cdot \mathbf{J}(S),$$

and the lift is the component of \mathbf{J} in the direction orthogonal to \mathbf{U}_∞ . For the fixed data, the drag can be regarded as a functional depending on the shape of the obstacle S . The drag minimization and the lift maximization are between shape optimization problems of some practical importance. The questions of the domain dependence of solutions to nonstationary compressible Navier–Stokes equations and of the solvability of the drag optimization problem were considered in papers [12], [13]. The solvability

of the drag minimization problem for stationary equations (1.2) is shown in [37], [40]. For incompressible Navier–Stokes equations, the existence of material derivatives of solutions and the formula for the shape derivative of the drag functional and adjoint state were obtained in [4], [5], and [42]; see also [43] and [44] for some generalizations. The growing literature on numerical and applied aspects of the problem is nicely surveyed in [18] and [27]. To the best of our knowledge, the mathematical sensitivity analysis for the compressible Navier–Stokes equations has not been studied yet. We derive the formula for the shape derivatives of the drag functional which can be used, in particular, for the explicit formulation of optimality conditions. In order to define the shape derivatives of the shape functional we combine the material derivatives of the solutions to the governing PDEs with an appropriate adjoint state according to the same scheme as is proposed, e.g., in [42] for steady incompressible equations.

We start with a description of our framework for shape sensitivity analysis, or more generally, for well-posedness of compressible Navier–Stokes equations. To this end we choose the vector field $\mathbf{T} \in C^2(\mathbb{R}^3)^3$ vanishing in the vicinity of Σ and define the mapping

$$(1.4) \quad y = x + \varepsilon \mathbf{T}(x),$$

which describes the perturbation of the shape of the obstacle. We refer the reader to [45] for more general framework and results in shape optimization and to [28] for shape calculus in the framework of fluid–structure interaction. For small ε , the mapping $x \rightarrow y$ diffeomorphically takes the flow region Ω onto $\Omega_\varepsilon = B \setminus S_\varepsilon$, where the perturbed obstacle $S_\varepsilon = y(S)$. Let $(\bar{\mathbf{u}}_\varepsilon, \bar{\varrho}_\varepsilon)$ be solutions to problem (1.2) in Ω_ε . After substituting $(\bar{\mathbf{u}}_\varepsilon, \bar{\varrho}_\varepsilon)$ into the formula for \mathbf{J} , the drag becomes the function of the parameter ε . Our aim is, in fact, to prove that this function is well defined and differentiable at $\varepsilon = 0$. This leads to the first order shape sensitivity analysis for solutions to compressible Navier–Stokes equations. It is convenient to reduce such an analysis to the analysis of dependence of solutions with respect to the coefficients of the governing equations. To this end, we introduce the functions $\mathbf{u}_\varepsilon(x)$ and $\varrho_\varepsilon(x)$ defined in the unperturbed domain Ω by the formulae

$$\mathbf{u}_\varepsilon(x) = \mathbf{N}\bar{\mathbf{u}}_\varepsilon(x + \varepsilon \mathbf{T}(x)), \quad \varrho_\varepsilon(x) = \bar{\varrho}_\varepsilon(x + \varepsilon \mathbf{T}(x)),$$

where

$$(1.5) \quad \mathbf{N}(x) = \det(\mathbf{I} + \varepsilon \mathbf{T}'(x))(\mathbf{I} + \varepsilon \mathbf{T}'(x))^{-1}$$

is the adjugate matrix of the Jacobi matrix $\mathbf{I} + \varepsilon \mathbf{T}'$. Furthermore, we also use the notation $\mathbf{g}(x) = \sqrt{\det \mathbf{N}}$. It is easy to see that the matrices $\mathbf{N}(x)$ depend analytically upon the small parameter ε and

$$(1.6) \quad \mathbf{N} = \mathbf{I} + \varepsilon \mathbf{D}(x) + \varepsilon^2 \mathbf{D}_1(\varepsilon, x),$$

where $\mathbf{D} = \operatorname{div} \mathbf{T} \mathbf{I} - \mathbf{T}'$. Calculations show that, for $\mathbf{u}_\varepsilon, \varrho_\varepsilon$, the following boundary value problem is obtained:

$$(1.7a) \quad \Delta \mathbf{u}_\varepsilon + \nabla \left(\lambda \mathbf{g}^{-1} \operatorname{div} \mathbf{u}_\varepsilon - \frac{R}{\varepsilon^2} p(\varrho_\varepsilon) \right) = \mathcal{A}(\mathbf{u}_\varepsilon) + R \mathcal{B}(\varrho_\varepsilon, \mathbf{u}_\varepsilon, \mathbf{u}_\varepsilon) \quad \text{in } \Omega,$$

$$(1.7b) \quad \operatorname{div}(\varrho_\varepsilon \mathbf{u}_\varepsilon) = 0 \quad \text{in } \Omega,$$

$$(1.7c) \quad \mathbf{u}_\varepsilon = \mathbf{U} \quad \text{on } \Sigma, \quad \mathbf{u}_\varepsilon = 0 \quad \text{on } \partial S,$$

$$(1.7d) \quad \varrho_\varepsilon = \varrho_0 \quad \text{on } \Sigma_{\text{in}}.$$

Here the linear operator \mathcal{A} and the nonlinear mapping \mathcal{B} are defined in terms of \mathbf{N} :

$$(1.8) \quad \begin{aligned} \mathcal{A}(\mathbf{u}) &= \Delta \mathbf{u} - (\mathbf{N}^*)^{-1} \operatorname{div} (\mathbf{g}^{-1} \mathbf{N} \mathbf{N}^* \nabla (\mathbf{N}^{-1} \mathbf{u})), \\ \mathcal{B}(\varrho, \mathbf{u}, \mathbf{w}) &= \varrho (\mathbf{N}^*)^{-1} (\mathbf{u} \nabla (\mathbf{N}^{-1} \mathbf{w})). \end{aligned}$$

For the derivation of these equations, see Appendix C.

The specific structure of the matrix \mathbf{N} does not play any particular role in the further analysis. Therefore, we consider a general problem of the existence, uniqueness, and dependence on coefficients of the solutions to equations (1.7) under the assumption that \mathbf{N} is a given matrix-valued function which is close, in an appropriate norm, to the identity mapping \mathbf{I} and coincides with \mathbf{I} in the vicinity of Σ . By abuse of notation, we write simply \mathbf{u} and ϱ , instead of \mathbf{u}_ε and ϱ_ε , when studying the well-posedness and dependence on \mathbf{N} . Before formulation of main results we write the governing equation in more transparent form using the change of unknown functions proposed in [35], [15]. To do so we introduce *the effective viscous pressure*

$$q = \frac{R}{\epsilon^2} p(\varrho) - \lambda \mathbf{g}^{-1} \operatorname{div} \mathbf{u}$$

and rewrite equations (1.7) in the equivalent form

$$(1.9a) \quad \Delta \mathbf{u} - \nabla q = \mathcal{A}(\mathbf{u}) + R \mathcal{B}(\varrho, \mathbf{u}, \mathbf{u}) \quad \text{in } \Omega,$$

$$(1.9b) \quad \operatorname{div} \mathbf{u} = \mathbf{g} \sigma_0 p(\varrho) - \frac{\mathbf{g}q}{\lambda} \quad \text{in } \Omega,$$

$$(1.9c) \quad \mathbf{u} \cdot \nabla \varrho + \mathbf{g} \sigma_0 p(\varrho) \varrho = \frac{\mathbf{g}q}{\lambda} \varrho \quad \text{in } \Omega,$$

$$(1.9d) \quad \mathbf{u} = \mathbf{U} \quad \text{on } \Sigma, \quad \mathbf{u} = 0 \quad \text{on } \partial S,$$

$$(1.9e) \quad \varrho = \varrho_0 \quad \text{on } \Sigma_{\text{in}},$$

where $\sigma_0 = R/(\lambda \epsilon^2)$. In the new variables (\mathbf{u}, q, ϱ) the expression for the force \mathbf{J} reads

$$(1.10) \quad \mathbf{J} = - \int_{\Omega} [\mathbf{g}^{-1} (\mathbf{N}^* \nabla (\mathbf{N}^{-1} \mathbf{u}) + \nabla (\mathbf{N}^{-1} \mathbf{u})^* \mathbf{N} - \operatorname{div} \mathbf{u}) - q - R \varrho \mathbf{u} \otimes \mathbf{u}] \mathbf{N}^* \nabla \eta \, dx,$$

where $\eta \in C^\infty(\Omega)$ is an arbitrary function, which is equal to 1 in an open neighborhood of the obstacle S and 0 in a vicinity of Σ . The value of \mathbf{J} is independent of the choice of the function η .

We assume that $\lambda \gg 1$ and $R \ll 1$, which corresponds to almost incompressible flow with low Reynolds number. In such a case, the *approximate solutions* to problem (1.9) can be chosen in the form $(\varrho_0, \mathbf{u}_0, q_0)$, where ϱ_0 is a constant in boundary condition (1.9e), and (\mathbf{u}_0, q_0) is a solution to the boundary value problem for the Stokes equations,

$$(1.11) \quad \begin{aligned} \Delta \mathbf{u}_0 - \nabla q_0 &= 0, \quad \operatorname{div} \mathbf{u}_0 = 0 \quad \text{in } \Omega, \\ \mathbf{u}_0 &= \mathbf{U} \quad \text{on } \Sigma, \quad \mathbf{u}_0 = 0 \quad \text{on } \partial S, \quad \Pi q_0 = q_0. \end{aligned}$$

In our notation Π is the projector:

$$\Pi u = u - \frac{1}{\operatorname{meas} \Omega} \int_{\Omega} u \, dx.$$

Equations (1.11) can be obtained as the limit of equations (1.9) for the passage $\lambda \rightarrow \infty$, $R \rightarrow 0$. It follows from the standard elliptic theory that, for the boundary $\partial\Omega \in C^\infty$, we have $(\mathbf{u}_0, q_0) \in C^\infty(\Omega)$. We look for solutions to problem (1.9) in the form

$$(1.12) \quad \mathbf{u} = \mathbf{u}_0 + \mathbf{u}, \quad \varrho = \varrho_0 + \varphi, \quad q = q_0 + \lambda\sigma_0 p(\varrho_0) + \pi + \lambda m,$$

with the unknowns functions $\vartheta = (\mathbf{u}, \pi, \varphi)$ and the unknown constant m . Substituting (1.12) into (1.9) we obtain the following boundary problem for ϑ :

$$(1.13a) \quad \begin{aligned} \Delta \mathbf{u} - \nabla \pi &= \mathcal{A}(\mathbf{u}) + R\mathcal{B}(\varrho, \mathbf{u}, \mathbf{u}) \quad \text{in } \Omega, \\ \operatorname{div} \mathbf{u} &= \mathbf{g} \left(\frac{\sigma}{\varrho_0} \varphi - \Psi[\vartheta] - m \right) \quad \text{in } \Omega, \\ \mathbf{u} \cdot \nabla \varphi + \sigma \varphi &= \Psi_1[\vartheta] + m\mathbf{g}\varrho \quad \text{in } \Omega, \\ \mathbf{u} = 0 \quad \text{on } \partial\Omega, \quad \varphi = 0 \quad \text{on } \Sigma_{\text{in}}, \quad \Pi\pi &= \pi, \end{aligned}$$

where

$$\begin{aligned} \Psi_1[\vartheta] &= \mathbf{g} \left(\varrho \Psi[\vartheta] - \frac{\sigma}{\varrho_0} \varphi^2 \right) + \sigma \varphi (1 - \mathbf{g}), \quad \Psi[\vartheta] = \frac{q_0 + \pi}{\lambda} - \frac{\sigma}{p'(\varrho_0)\varrho_0} H(\varphi), \\ \sigma &= \sigma_0 p'(\varrho_0)\varrho_0, \quad H(\varphi) = p(\varrho_0 + \varphi) - p(\varrho_0) - p'(\varrho_0)\varphi, \end{aligned}$$

and the vector field \mathbf{u} and the function ϱ are given by (1.12). Finally, we specify the constant m . In our framework, in contrast to the case of the homogeneous boundary problem, the solution to such a problem is not trivial. Note that, since $\operatorname{div} \mathbf{u}$ is of the null mean value, the right-hand side of (1.13a)₃ must satisfy the compatibility condition

$$m \int_{\Omega} \mathbf{g} \, dx = \int_{\Omega} \mathbf{g} \left(\frac{\sigma}{\varrho_0} \varphi - \Psi[\vartheta] \right) dx,$$

which formally determines m . This choice of m leads to essential mathematical difficulties. To make this issue clear note that in the simplest case $\mathbf{g} = 1$ we have $m = \varrho_0^{-1} \sigma (\mathbf{I} - \Pi)\varphi + O(|\vartheta|^2, \lambda^{-1})$, and the principal linear part of the governing equations (1.13a) becomes

$$\begin{pmatrix} \Delta & -\nabla & 0 \\ \operatorname{div} & 0 & -\frac{\sigma}{\varrho_0} \\ 0 & 0 & \mathbf{u}\nabla + \sigma \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \pi \\ \varphi \end{pmatrix} + \begin{pmatrix} 0 \\ m \\ -m\varrho_0 \end{pmatrix} \sim \begin{pmatrix} \Delta \mathbf{u} - \nabla \pi \\ \operatorname{div} \mathbf{u} - \frac{\sigma}{\varrho_0} \Pi \varphi \\ \mathbf{u}\nabla \varphi + \sigma \Pi \varphi \end{pmatrix}.$$

Hence, the question of solvability of the linearized equations derived for (1.13) can be reduced to the question of solvability of the boundary value problem for the nonlocal transport equation

$$\mathbf{u}\nabla \varphi + \sigma \Pi \varphi = f,$$

which is very difficult because of the loss of maximum principle. In fact, this question is concerned with the problem of the control of the total gas mass in compressible flows. Recall that the absence of the mass control is the main obstacle for proving the global solvability of inhomogeneous boundary problems for compressible Navier–Stokes equations; we refer to [25] for discussion. In order to cope with this difficulty we write the compatibility condition in a sophisticated form, which allows us to control

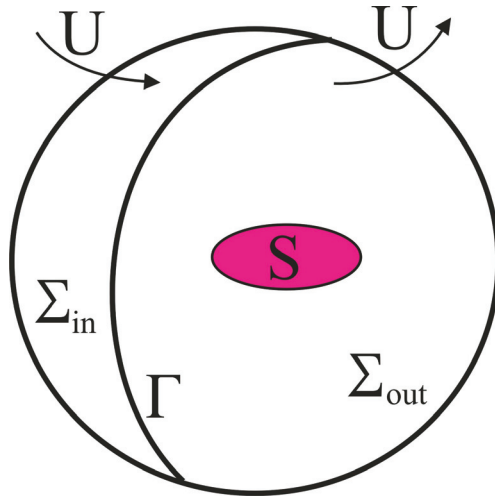


FIG. 1. Flow domain Ω with an obstacle S .

the total mass of the gas. To this end we introduce the auxiliary function ζ satisfying the equations

$$(1.13b) \quad -\operatorname{div}(\mathbf{u}\zeta) + \sigma\zeta = \sigma\mathbf{g} \text{ in } \Omega, \quad \zeta = 0 \text{ on } \Sigma_{\text{out}},$$

and fix the constant m as follows:

$$(1.13c) \quad m = \varkappa \int_{\Omega} (\varrho_0^{-1}\Psi_1[\vartheta]\zeta - \mathbf{g}\Psi[\vartheta]) dx, \quad \varkappa = \left(\int_{\Omega} \mathbf{g}(1 - \zeta - \varrho_0^{-1}\zeta\varphi) dx \right)^{-1}.$$

In this way the auxiliary function ζ becomes an integral part of the solution to problem (1.13). Now our aim is to prove the existence and uniqueness of solutions to problem (1.13) and investigate the dependence of the solutions on matrices \mathbf{N} . Before the presentation of the main results we introduce some notation and formulate preliminary results.

Geometrical conditions on the flow region. We assume that a surface $\Sigma = \Sigma_{\text{in}} \cup \Sigma_{\text{out}} \cup \Gamma$ and a given vector field \mathbf{U} satisfy the following conditions, referred to as the *emergent vector field conditions*.

CONDITION 1.1. *The set Γ is a closed $C^{2+\alpha}$ one-dimensional manifold. Moreover, there is a positive constant c such that*

$$(1.14) \quad \mathbf{U} \cdot \nabla(\mathbf{U} \cdot \mathbf{n}) > c > 0 \text{ on } \Gamma.$$

These conditions have a simple geometric interpretation, that $\mathbf{U} \cdot \mathbf{n}$ vanishes only up to the first order at Γ , and \mathbf{U} is transversal to Γ ; furthermore, for each point $P \in \Gamma$, the vector $\mathbf{U}(P)$ points to the part of Σ where \mathbf{U} is an exterior vector field (see Figure 1). Note that the emergent vector field condition plays an important role in the theory of the classical oblique derivative problem; see [16]. In the context of our problem it is equivalent to the following conditions:

- (H1) The boundary of Ω belongs to class $C^{2+\alpha}$, $\alpha \in (0, 1)$. For each point $P \in \Gamma$ there exist the local Cartesian coordinates (x_1, x_2, x_3) with the origin at P such that in the new coordinates $\mathbf{U}(P) = (U, 0, 0)$ with $U = |\mathbf{U}(P)|$, and

$\mathbf{n}(P) = (0, 0, -1)$. Moreover, there is a neighborhood $\mathcal{O} = [-k, k]^2 \times [-t, t]$ of P such that the intersections $\Sigma \cap \mathcal{O}$ and $\Gamma \cap \mathcal{O}$ are defined by the equations

$$F_0(x) \equiv x_3 - F(x_1, x_2) = 0, \quad \nabla F_0(x) \cdot \mathbf{U}(x) = 0,$$

and $\Omega \cap \mathcal{O}$ is the epigraph $\{F_0 > 0\} \cap \mathcal{O}$. The function F satisfies the conditions

$$(1.15) \quad \|F\|_{C^2([-k,k]^2)} \leq K, \quad F(0, 0) = 0, \quad \nabla F(0, 0) = 0,$$

where the constants $k, t < 1$ and $K > 1$ depend only on the curvature of Σ and are independent of the point P .

(H2) For a suitable choice of the constant k , with k independent of $P \in \Gamma$, the manifold $\Gamma \cap \mathcal{O}$ admits the parameterization

$$(1.16) \quad x = \mathbf{x}^0(x_2) := (\Upsilon(x_2), x_2, F(\Upsilon(x_2), x_2))$$

such that $\Upsilon(0) = 0$ and $\|\Upsilon\|_{C^2([-k,k])} \leq C_\Gamma$, where the constant $C_\Gamma > 1$ depends only on Σ and \mathbf{U} .

(H3) There are positive constants N^\pm independent of P such that for $x \in \Sigma$ given by the condition $F_0(x_1, x_2, x_3) = x_3 - F(x_1, x_2) = 0$ we have

$$(1.17) \quad \begin{aligned} N^-(x_1 - \Upsilon(x_2)) &\leq -\nabla F_0(x) \cdot \mathbf{U}(x) \leq N^+(x_1 - \Upsilon(x_2)) \quad \text{for } x_1 > \Upsilon(x_2), \\ -N^-(x_1 - \Upsilon(x_2)) &\leq \nabla F_0(x) \cdot \mathbf{U}(x) \leq -N^+(x_1 - \Upsilon(x_2)) \quad \text{for } x_1 < \Upsilon(x_2). \end{aligned}$$

Function spaces. In this paragraph we assemble some technical results which are used throughout the paper. Function spaces play a central role, and we recall some notation, fundamental definitions, and properties, which can be found in [1] and [6]. For the convenience of the reader we collect in Appendix B the basic facts from the theory of interpolation spaces. For our applications we need the results in three spatial dimensions; however, the results are presented for the dimension $d \geq 2$.

Let Ω be the whole space \mathbb{R}^d or a bounded domain in \mathbb{R}^d with the boundary $\partial\Omega$ of class C^1 . For an integer $l \geq 0$ and for an exponent $r \in [1, \infty)$, we denote by $W^{l,r}(\Omega)$ the Sobolev space endowed with the norm $\|u\|_{W^{l,r}(\Omega)} = \sup_{|\alpha| \leq l} \|\partial^\alpha u\|_{L^r(\Omega)}$. For real $0 < s < 1$, the fractional Sobolev space $W^{s,r}(\Omega)$ is obtained by the real interpolation method (see [46] for the proofs) between $L^r(\Omega)$ and $W^{1,r}(\Omega)$ and consists of all measurable functions with the finite norm

$$\|u\|_{W^{s,r}(\Omega)} = \|u\|_{L^r(\Omega)} + |u|_{s,r,\Omega},$$

where

$$(1.18) \quad |u|_{s,r,\Omega}^r = \int_{\Omega \times \Omega} |x - y|^{-d-rs} |u(x) - u(y)|^r dx dy.$$

In the general case, the Sobolev space $W^{l+s,r}(\Omega)$ is defined as the space of measurable functions with the finite norm $\|u\|_{W^{l+s,r}(\Omega)} = \sup_{|\alpha| \leq l} \|\partial^\alpha u\|_{W^{s,r}(\Omega)}$. For $0 < s < 1$, the Sobolev space $W^{s,r}(\Omega)$ is, in fact [6], the interpolation space $[L^r(\Omega), W^{1,r}(\Omega)]_{s,r}$.

Furthermore, the notation $W_0^{l,r}(\Omega)$, with an integer l , stands for the closed subspace of the space $W^{l,r}(\Omega)$ of all functions $u \in L^r(\Omega)$ which being extended by zero outside of Ω belong to $W^{l,r}(\mathbb{R}^d)$.

Denote by $\mathcal{W}_0^{0,r}(\Omega)$ and $\mathcal{W}_0^{1,r}(\Omega)$ the subspaces of $L^r(\mathbb{R}^d)$ and $W^{1,r}(\mathbb{R}^d)$, respectively, of all functions vanishing outside of Ω . Obviously $\mathcal{W}_0^{1,r}(\Omega)$ and $W_0^{1,r}(\Omega)$ are isomorphic topologically and algebraically and we can identify them. However, we need the interpolation spaces $\mathcal{W}_0^{s,r}(\Omega)$ for nonintegers, in particular, for $s = 1/r$.

DEFINITION 1.2. For all $0 < s \leq 1$ and $1 < r < \infty$, we denote by $\mathcal{W}_0^{s,r}(\Omega)$ the interpolation space $[\mathcal{W}_0^{0,r}(\Omega), \mathcal{W}_0^{1,r}(\Omega)]_{s,r}$ endowed with one of the equivalent norms (6.1) and (6.3) defined by the interpolation method.

It follows from the definition of interpolation spaces (see Appendix B) that $\mathcal{W}_0^{s,r}(\Omega) \subset W^{s,r}(\mathbb{R}^d)$ and, for all $u \in \mathcal{W}_0^{s,r}(\Omega)$,

$$(1.19) \quad \|u\|_{W^{s,r}(\mathbb{R}^d)} \leq c(r, s) \|u\|_{\mathcal{W}_0^{s,r}(\Omega)}, \quad u = 0 \quad \text{outside } \Omega.$$

In other words, $\mathcal{W}_0^{s,r}(\Omega)$ consists of all elements $u \in W^{s,r}(\Omega)$ such that the extension \bar{u} of u by 0 outside of Ω has the finite $[\mathcal{W}_0^{0,r}(\Omega), \mathcal{W}_0^{1,r}(\Omega)]_{s,r}$ -norm. We identify u and \bar{u} for the elements $u \in \mathcal{W}_0^{s,r}(\Omega)$. With this identification it follows that $W_0^{1,r}(\Omega) \subset \mathcal{W}_0^{s,r}(\Omega)$ and the space $C_0^\infty(\Omega)$ is dense in $\mathcal{W}_0^{s,r}(\Omega)$. It is worthy to note that for $0 < s < 1$ and for $1 < r < \infty$, the function \bar{u} belongs to the space $W^{s,r}(\mathbb{R}^d)$ if and only if $u \in W^{s,r}(\Omega)$ and $\text{dist}(x, \partial\Omega)^{-s}u \in L^r(\Omega)$. We also point out that the interpolation space $\mathcal{W}_0^{s,r}(\Omega)$ coincides with the Sobolev space $W_0^{s,r}(\Omega)$ for $s \neq 1/r$. Recall that the standard space $W_0^{r,s,r}(\Omega)$ is the completion of $C_0^\infty(\Omega)$ in the $W^{s,r}(\Omega)$ -norm.

Embedding theorems. For $sr > d$ and $0 \leq \alpha < s - r/d$, the embedding $W^{s,r}(\Omega) \hookrightarrow C^\alpha(\Omega)$ is continuous and compact. In particular, for $sr > d$, the Sobolev space $W^{s,r}(\Omega)$ is a commutative Banach algebra; i.e., for all $u, v \in W^{s,r}(\Omega)$,

$$(1.20) \quad \|uv\|_{W^{s,r}(\Omega)} \leq c(r, s) \|u\|_{W^{s,r}(\Omega)} \|v\|_{W^{s,r}(\Omega)}.$$

If $sr < d$ and $t^{-1} = r^{-1} - d^{-1}s$, then the embedding $W^{s,r}(\Omega) \hookrightarrow L^t(\Omega)$ is continuous, [1, Thm. 7.57]. We also have [1, Thm. 7.58], for $\alpha < s$, $(s - \alpha)r < d$ and $\beta^{-1} = r^{-1} - d^{-1}(s - \alpha)$,

$$(1.21) \quad \|u\|_{W^{\alpha,\beta}(\Omega)} \leq c(r, s, \alpha, \beta, \Omega) \|u\|_{W^{s,r}(\Omega)}.$$

It follows from (1.19) that all the embedding inequalities remain true for the elements of the interpolation space $\mathcal{W}_0^{s,r}(\Omega)$.

Duality. We define

$$(1.22) \quad \langle u, v \rangle = \int_{\Omega} u v \, dx$$

for any functions such that the right-hand side makes sense. For $r \in (1, \infty)$, each element $v \in L^{r'}(\Omega)$, $r' = r/(r - 1)$, determines the functional L_v of $(\mathcal{W}_0^{s,r}(\Omega))'$ by the identity $L_v(u) = \langle u, v \rangle$. We introduce the $(-s, r')$ -norm of an element $v \in L^{r'}(\Omega)$ to be by definition the norm of the functional L_v , that is

$$(1.23) \quad \|v\|_{\mathcal{W}^{-s,r'}(\Omega)} = \sup_{\substack{u \in \mathcal{W}_0^{s,r}(\Omega) \\ \|u\|_{\mathcal{W}_0^{s,r}(\Omega)} = 1}} |\langle u, v \rangle|.$$

We let $\mathcal{W}^{-s,r'}(\Omega)$ denote the completion of the space $L^{r'}(\Omega)$ with respect to the $(-s, r')$ -norm. For an integer s , $\mathcal{W}^{-s,r'}(\Omega)$ is topologically and algebraically isomorphic to $(W_0^{s,r}(\Omega))'$. Moreover, we can identify $\mathcal{W}^{-s,r'}(\Omega)$ with the interpolation space

$[L^{r'}(\Omega), W_0^{-1,r'}(\Omega)]_{s,r}$; see [6] and Appendix B. With this denotation we have the duality principle

$$(1.24) \quad \|u\|_{\mathcal{W}_0^{s,r}(\Omega)} = \sup_{\substack{v \in C_0^\infty(\Omega) \\ \|v\|_{\mathcal{W}^{-s,r'}(\Omega)}=1}} |\langle u, v \rangle|.$$

With applications to the theory of Navier–Stokes equations in mind, we introduce the smaller dual space defined as follows. We identify the function $v \in L^{r'}(\Omega)$ with the functional $L_v \in (W^{s,r}(\Omega))'$ and denote by $\mathbb{W}^{-s,r'}(\Omega)$ the completion of $L^{r'}(\Omega)$ in the norm

$$(1.25) \quad \|v\|_{\mathbb{W}^{-s,r'}(\Omega)} := \sup_{\substack{u \in W^{s,r}(\Omega) \\ \|u\|_{W^{s,r}(\Omega)}=1}} |\langle u, v \rangle|.$$

In the sense of this identification the space $C_0^\infty(\Omega)$ is dense in the interpolation space $\mathbb{W}^{-s,r}(\Omega)$. It follows immediately from the definition that

$$\mathbb{W}^{-s,r'}(\Omega) \subset (W^{s,r}(\Omega))' \subset \mathcal{W}^{-s,r'}(\Omega).$$

For an arbitrary bounded domain $\Omega \subset \mathbb{R}^3$ with a Lipschitz boundary, we introduce the Banach spaces

$$\begin{aligned} X^{s,r} &= W^{s,r}(\Omega) \cap W^{1,2}(\Omega), & Y^{s,r} &= W^{s+1,r}(\Omega) \cap W^{2,2}(\Omega), \\ Z^{s,r} &= \mathcal{W}^{s-1,r}(\Omega) \cap L^2(\Omega) \end{aligned}$$

equipped with the norms

$$\begin{aligned} \|u\|_{X^{s,r}} &= \|u\|_{W^{s,r}(\Omega)} + \|u\|_{W^{1,2}(\Omega)}, & \|u\|_{Y^{s,r}} &= \|u\|_{W^{1+s,r}(\Omega)} + \|u\|_{W^{2,2}(\Omega)}, \\ \|u\|_{Z^{s,r}} &= \|u\|_{\mathcal{W}^{s-1,r}(\Omega)} + \|u\|_{L^2(\Omega)}. \end{aligned}$$

It can be easily seen that the embeddings $Y^{s,r} \hookrightarrow X^{s,r} \hookrightarrow Z^{s,r}$ are compact and, for $sr > 3$, each of the spaces $X^{s,r}$ and $Y^{s,r}$ is a commutative Banach algebra.

Stokes equations. The following lemma is a straightforward consequence of classical results on solvability of the first boundary value problem for Stokes equations (see [9]) and the interpolation theory.

LEMMA 1.3. *Let $\Omega \subset \mathbb{R}^3$ be a bounded domain with $\partial\Omega \in C^2$ and $(F, G) \in \mathcal{W}^{s-1,r}(\Omega) \times W^{s,r}(\Omega)$ ($0 \leq s \leq 1, 1 < r < \infty$). Then the boundary value problem*

$$(1.26) \quad \begin{aligned} \Delta \mathbf{u} - \nabla \pi &= F, & \operatorname{div} \mathbf{u} &= \Pi G & \text{in } \Omega, \\ \mathbf{u} &= 0 & \text{on } \partial\Omega, & & \Pi \pi = \pi, \end{aligned}$$

has a unique solution $(\mathbf{u}, \pi) \in W^{s+1,r}(\Omega) \times W^{s,r}(\Omega)$ such that

$$(1.27) \quad \|\mathbf{u}\|_{W^{s+1,r}(\Omega)} + \|\pi\|_{W^{s,r}(\Omega)} \leq c(\Omega, r, s)(\|F\|_{\mathcal{W}^{s-1,r}(\Omega)} + \|G\|_{W^{s,r}(\Omega)}).$$

In particular, we have

$$\|\mathbf{u}\|_{Y^{s,r}} + \|\pi\|_{X^{s,r}} \leq c(\Omega, r, s)(\|F\|_{Z^{s,r}} + \|G\|_{X^{s,r}}).$$

Proof. The proof is in Appendix B. \square

Note that the lemma works for singular values $s = 1/r + \text{integer}$. But in this case the traces of solutions at the boundary are not defined.

1.2. Results. *Transport equations.* Today there exists a complete theory of generalized solutions to the class of hyperbolic-elliptic equations developed in [10] and [34] under the assumptions that the equations have C^1 coefficients and satisfy the maximum principle. The questions on smoothness properties of solutions are more difficult. We recall the classical results of [19], [34], related to the case of $\bar{\Sigma}_{\text{in}} \cap \bar{\Sigma}_{\text{out}} = \emptyset$. The particular case, with $\Sigma_{\text{in}} = \Sigma_{\text{out}} = \emptyset$, in the Sobolev spaces is covered in the papers [3] and [30], [31]. The case of the nonempty interface $\Gamma = \bar{\Sigma}_{\text{in}} \cap \bar{\Sigma}_{\text{out}}$ is still weakly investigated. The theory of boundary value problems for transport equations is an integral part of the theory of multicomponents and inhomogeneous flows. In this context this problem has been considered by many authors; see [24] and [29] for discussion.

In general the existence of strong solutions to inhomogeneous boundary value problems for transport equations is still an open problem. The following theorem, which is used throughout this paper, partially fills this gap. Let us consider the following boundary value problems for linear transport equations:

$$(1.28) \quad \mathcal{L}\varphi := \mathbf{u}\nabla\varphi + \sigma\varphi = f \text{ in } \Omega, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}},$$

$$(1.29) \quad \mathcal{L}^*\varphi^* := -\text{div}(\varphi^*\mathbf{u}) + \sigma\varphi^* = f \text{ in } \Omega, \quad \varphi^* = 0 \text{ on } \Sigma_{\text{out}}.$$

The bounded functions φ, φ^* are called the generalized solutions to problems (1.28), (1.29), respectively, if the integral identities

$$(1.30) \quad \int_{\Omega} (\varphi \mathcal{L}^* \zeta^* - f \zeta^*) dx = 0, \quad \int_{\Omega} (\varphi^* \mathcal{L} \zeta - f \zeta) dx = 0$$

hold true for all test functions $\zeta^*, \zeta \in C(\Omega) \cap W^{1,1}(\Omega)$, respectively, such that $\zeta^* = 0$ on Σ_{out} and $\zeta = 0$ on Σ_{in} .

THEOREM 1.4. *Assume that Σ and \mathbf{U} satisfy conditions (H1)–(H3), the exponents s, r satisfy the inequalities*

$$(1.31) \quad 1/2 < s \leq 1, \quad 1 < r < 3/(2s - 1),$$

and the vector field \mathbf{u} belongs to the class $C^1(\Omega)$ and satisfies the boundary condition

$$(1.32) \quad \mathbf{u} = \mathbf{U} \text{ on } \Sigma, \quad \mathbf{u} = 0 \text{ on } \partial S.$$

Then there are positive constants σ^* and δ^* depending only on Σ, \mathbf{U}, s, r , and $\|\mathbf{u}\|_{C^1(\Omega)}$ such that the following hold:

(i) For any $\sigma > \sigma^*$ and $f \in W^{s,r}(\Omega) \cap L^\infty(\Omega)$, problem (1.28) has a unique solution $\varphi \in W^{s,r}(\Omega) \cap L^\infty(\Omega)$ satisfying the inequalities

$$(1.33) \quad \|\varphi\|_{W^{s,r}(\Omega)} \leq C\|f\|_{W^{s,r}(\Omega)}, \quad \|\varphi\|_{L^\infty(\Omega)} \leq \sigma^{-1}\|f\|_{L^\infty(\Omega)}.$$

(ii) If, in addition, $\|\text{div } \mathbf{u}\|_{W^{s,r}(\Omega)} + \|\text{div } \mathbf{u}\|_{L^\infty(\Omega)} \leq \delta^*$, problem (1.29) has a unique solution $\varphi^* \in W^{s,r}(\Omega) \cap L^\infty(\Omega)$, which admits the estimates

$$(1.34) \quad \|\varphi^*\|_{W^{s,r}(\Omega)} \leq C\|f\|_{W^{s,r}(\Omega)}, \quad \|\varphi^*\|_{L^\infty(\Omega)} \leq (\sigma - \delta^*)^{-1}\|f\|_{L^\infty(\Omega)}.$$

The constant C depends only on $\|\mathbf{u}\|_{C^1(\Omega)}, r, s, \sigma, \mathbf{U}$, and Ω .

Restriction $s \geq 1/2$ in condition (1.31) is not essential and can be removed. But it plays an important role in the proof of Lemma 1.7 and the main theorem, Theorem 1.9. On the other hand, the inequality $r < 3/(2s - 1)$ is crucial for the

control of singularities of solutions at the characteristic manifold Γ . It is connected with the behavior of \mathbf{U} in the vicinity of Γ and does not depend on the dimension of the ambient space. Note also that solution is strong for $s = 1$, but not continuous since $W^{1,3}(\Omega)$ is not embedded in $C(\Omega)$. On the other hand, $sr \rightarrow \infty$ as $s \rightarrow 1/2$ and solution is continuous for s close to $1/2$. In order to obtain continuous strong solutions we have to consider the problem in the scale of spaces $X^{s,r}$

Since for $sr > 3$, the embeddings $X^{s,r} \hookrightarrow C(\Omega)$, $Y^{s,r} \hookrightarrow C^1(\Omega)$ are bounded, we have the following result on solvability of problems (1.28), (1.29) in space $X^{s,r}$.

COROLLARY 1.5. *Assume that $sr > 3$ and the vector field \mathbf{u} has the representation $\mathbf{u} = \mathbf{u}_0 + \mathbf{u}$, where $\mathbf{u}_0 \in C^\infty(\Omega)^3$ is a solution to problem (1.11). Then there exist $\tau^* \in (0, 1]$ and σ^* , depending only on Σ, \mathbf{u}_0 , and s, r , such that, for all \mathbf{u} with $\|\mathbf{u}\|_{Y^{s,r}} \leq \tau^*$, $\sigma > \sigma^*$, and $f \in X^{s,r}$, each of problems (1.28) and (1.29) has a unique solution satisfying the inequalities*

$$(1.35) \quad \|\varphi\|_{X^{s,r}} \leq C\|f\|_{X^{s,r}}, \quad \|\varphi^*\|_{X^{s,r}} \leq C\|f\|_{X^{s,r}}.$$

Existence and uniqueness theory. The second main result of this paper concerns the existence and local uniqueness of solutions to problem (1.13). Denote by E the closed subspace of the Banach space $Y^{s,r}(\Omega)^3 \times X^{s,r}(\Omega)^2$ in the form

$$(1.36) \quad E = \{\vartheta = (\mathbf{u}, \pi, \varphi) : \mathbf{u} = 0 \text{ on } \partial\Omega, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}}, \quad \Pi\pi = \pi\},$$

and denote by $\mathcal{B}_\tau \subset E$ the closed ball of radius τ centered at 0. Next, note that, for $sr > 3$, elements of the ball \mathcal{B}_τ satisfy the inequality

$$(1.37) \quad \|\mathbf{u}\|_{C^1(\Omega)} + \|\pi\|_{C(\Omega)} + \|\varphi\|_{C(\Omega)} \leq c_e(r, s, \Omega)\|\vartheta\|_E \leq c_e\tau,$$

where the norm in E is defined by

$$\|\vartheta\|_E = \|\mathbf{u}\|_{Y^{s,r}(\Omega)} + \|\pi\|_{X^{s,r}(\Omega)} + \|\varphi\|_{X^{s,r}(\Omega)}.$$

THEOREM 1.6. *Assume that the surface Σ and given vector field \mathbf{U} satisfy conditions (H1)–(H3). Furthermore, let σ^*, τ^* be constants given by Corollary 1.5, and let positive numbers r, s, σ satisfy the inequalities*

$$(1.38) \quad 1/2 < s \leq 1, \quad 1 < r < 3/(2s - 1), \quad sr > 3, \quad \sigma > \sigma^*.$$

Then there exists $\tau_0 \in (0, \tau^]$, depending only on $\mathbf{U}, \Omega, r, s, \sigma$, such that, for all*

$$(1.39) \quad \tau \in (0, \tau_0], \quad \lambda^{-1}, R \in (0, \tau^2], \quad \|\mathbf{N} - \mathbf{I}\|_{C^2(\Omega)} \leq \tau^2,$$

problem (1.13), with \mathbf{u}_0 given by (1.11), has a unique solution $\vartheta \in \mathcal{B}_\tau$. Moreover, the auxiliary function ζ and the constants \varkappa, m admit the estimates

$$(1.40) \quad \|\zeta\|_{X^{s,r}} + |\varkappa| \leq c, \quad |m| \leq c\tau < 1,$$

where the constant c depends only on \mathbf{U}, Ω, r, s , and σ .

Material derivatives of solutions. Theorem 1.6 guarantees the existence and uniqueness of solutions to problem (1.13) for all \mathbf{N} close to the identity matrix \mathbf{I} . The totality of such solutions can be regarded as the mapping from \mathbf{N} to the solution to Navier–Stokes equations. The natural question is the smoothness properties of this mapping, in particular, its differentiability. With application to shape optimization problems in mind, we consider the particular case where the matrices \mathbf{N} depend on

the small parameter ε and have representation (1.6). We assume that C^1 -norms of the matrix-valued functions \mathbf{D} and $\mathbf{D}_1(\varepsilon)$ in (1.6) have a majorant independent of ε . By virtue of Theorem 1.6, there are the positive constants ε_0 and τ such that, for all sufficiently small R, λ^{-1} and $\varepsilon \in [0, \varepsilon_0]$, problem (1.13) with $\mathbf{N} = \mathbf{N}(\varepsilon)$ has a unique solution $\vartheta(\varepsilon) = (\mathbf{u}(\varepsilon), \pi(\varepsilon), \varphi(\varepsilon)), \zeta(\varepsilon), m(\varepsilon)$, which admits the estimate

$$(1.41) \quad \|\vartheta(\varepsilon)\|_E + |m(\varepsilon)| \leq c\tau, \quad \|\zeta(\varepsilon)\|_{X^{s,r}} \leq c,$$

where the constant c is independent of ε , and the Banach space E is defined by (1.36). Denote the solution $(\vartheta(0), m(0), \zeta(0))$ for $\varepsilon = 0$ by (ϑ, m, ζ) , and define the finite differences with respect to ε ,

$$(\mathbf{w}_\varepsilon, \omega_\varepsilon, \psi_\varepsilon) = \varepsilon^{-1}(\vartheta - \vartheta(\varepsilon)), \quad \xi_\varepsilon = \varepsilon^{-1}(\zeta - \zeta(\varepsilon)), \quad n_\varepsilon = \varepsilon^{-1}(m - m(\varepsilon)).$$

Formal calculations show that the limit $(\mathbf{w}, \omega, \psi, \xi, n) = \lim_{\varepsilon \rightarrow 0} (\mathbf{w}_\varepsilon, \omega_\varepsilon, \psi_\varepsilon, \xi_\varepsilon, n_\varepsilon)$ is a solution to linearized equations

$$(1.42) \quad \begin{aligned} \Delta \mathbf{w} - \nabla \omega &= R \mathcal{C}_0(\mathbf{w}, \psi) + \mathcal{D}_0(\mathbf{D}) \quad \text{in } \Omega, \\ \operatorname{div} \mathbf{w} &= b_{21}^0 \psi - b_{22}^0 \omega + b_{23}^0 n + b_{20}^0 \mathfrak{d} \quad \text{in } \Omega, \\ \mathbf{u} \nabla \psi + \sigma \psi &= -\mathbf{w} \cdot \nabla \varphi + b_{11}^0 \psi + b_{12}^0 \omega + b_{13}^0 n + b_{10}^0 \mathfrak{d} \quad \text{in } \Omega, \\ -\operatorname{div}(\mathbf{u} \xi) + \sigma \xi &= \operatorname{div}(\zeta \mathbf{w}) + \sigma \mathfrak{d} \quad \text{in } \Omega, \\ \mathbf{w} &= 0 \quad \text{on } \partial \Omega, \quad \psi = 0 \quad \text{on } \Sigma_{\text{in}}, \quad \xi = 0 \quad \text{on } \Sigma_{\text{out}}, \\ \omega - \Pi \omega &= 0, \quad n = \varkappa \int_{\Omega} (b_{31}^0 \psi + b_{32}^0 \omega + b_{34}^0 \xi + b_{30}^0 \mathfrak{d}) \, dx, \end{aligned}$$

where $\mathfrak{d} = 1/2 \operatorname{Tr} \mathbf{D}$ and the variable coefficients b_{ij}^0 and the operators $\mathcal{C}_0, \mathcal{D}_0$, are defined by the formulae

$$(1.43) \quad \begin{aligned} b_{11}^0 &= \Psi[\vartheta] - \varrho H'(\varphi) + m - \frac{2\sigma}{\varrho_0} \varphi, \quad b_{12}^0 = \lambda^{-1} \varrho, \quad b_{13}^0 = \varrho, \\ b_{10}^0 &= \varrho \Psi[\vartheta] - \frac{\sigma}{\varrho_0} \varphi^2 - \sigma \varphi + m \varrho, \quad b_{21}^0 = \frac{\sigma}{\varrho_0} \psi_0 + H'(\varphi), \\ b_{22}^0 &= -\lambda^{-1}, \quad b_{23}^0 = -1, \quad b_{20}^0 = \sigma \varphi \varrho_0^{-1} - \Psi[\vartheta] - m, \\ b_{31}^0 &= \varrho_0^{-1} \zeta \left(\Psi[\vartheta] - \varrho H'(\varphi) - \frac{2\sigma}{\varrho_0} \varphi \right) - H'(\varphi) + m \varrho_0^{-1} \zeta, \\ b_{32}^0 &= (\lambda \varrho_0)^{-1} \varrho \zeta b_{12}^0 + \lambda^{-1}, \quad b_{34}^0 = \varrho_0^{-1} \Psi_1[\vartheta] + m(1 + \varrho_0^{-1} \varphi), \\ b_{30}^0 &= \varrho_0^{-1} \zeta (\Psi_1[\vartheta] - m \varrho) + \Psi[\vartheta] - m(1 - \zeta - \varrho_0^{-1} \zeta \varphi), \end{aligned}$$

$$(1.44) \quad \mathcal{C}_0(\psi, \mathbf{w}) = R\psi \mathbf{u} \nabla \mathbf{u} + R\varrho \mathbf{w} \nabla \mathbf{u} + R\varrho \mathbf{u} \nabla \mathbf{w},$$

$$(1.45) \quad \begin{aligned} \mathcal{D}_0(\mathbf{D}) &= R\mathbf{u} \nabla(\mathbf{D}\mathbf{u}) + R\mathbf{D}^*(\mathbf{u} \nabla \mathbf{u}) \\ &+ \operatorname{div} \left((\mathbf{D} + \mathbf{D}^*) \nabla \mathbf{u} - \frac{1}{2} \operatorname{Tr} \mathbf{D} \nabla \mathbf{u} \right) - \mathbf{D}^* \Delta \mathbf{u} - \Delta(\mathbf{D}\mathbf{u}). \end{aligned}$$

The justification of the formal procedure meets the serious problems, since the smoothness of solutions to problem (1.13) is not sufficient for the well-posedness of problem (1.42) in the standard weak formulation. In order to cope with this difficulty we define *very weak solutions* to problem (1.42). The construction of such solutions is based on

the following lemma, and the proof is given in Appendix A. The lemma is given in \mathbb{R}^d , for our application $d = 3$.

LEMMA 1.7. *Let $\Omega \subset \mathbb{R}^d$ be a bounded domain with the Lipschitz boundary, let exponents s and r satisfy the inequalities $sr > d$, $1/2 \leq s \leq 1$, and $\varphi, \varsigma \in W^{s,r}(\Omega) \cap W^{1,2}(\Omega)$, $\mathbf{w} \in \mathcal{W}_0^{1-s,r'}(\Omega) \cap W_0^{1,2}(\Omega)$. Then there is a constant c depending only on s, r , and Ω , such that the trilinear form*

$$\mathfrak{B}(\mathbf{w}, \varphi, \varsigma) = - \int_{\Omega} \varsigma \mathbf{w} \cdot \nabla \varphi \, dx$$

satisfies the inequality

$$(1.46) \quad |\mathfrak{B}(\mathbf{w}, \varphi, \varsigma)| \leq c \|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)} \|\varphi\|_{W^{s,r}(\Omega)} \|\varsigma\|_{W^{s,r}(\Omega)}$$

and can be continuously extended to $\mathfrak{B} : \mathcal{W}_0^{1-s,r'}(\Omega)^d \times W^{s,r}(\Omega)^2 \mapsto \mathbb{R}$. In particular, we have $\varsigma \nabla \varphi \in W^{s-1,r}(\Omega)$ and $\|\varsigma \nabla \varphi\|_{W^{s-1,r}(\Omega)} \leq c \|\varphi\|_{W^{s,r}(\Omega)} \|\varsigma\|_{W^{s,r}(\Omega)}$.

DEFINITION 1.8. *The vector field $\mathbf{w} \in \mathcal{W}_0^{1-s,r'}(\Omega)^3$, functionals $(\omega, \psi, \xi) \in \mathbb{W}^{-s,r'}(\Omega)^3$, and constant n are said to be a weak solution to problem (1.42) if $\langle \omega, 1 \rangle = 0$ and the identity*

$$(1.47) \quad \begin{aligned} & \int_{\Omega} \mathbf{w} \left(\mathbf{H} - R_{\varrho} \nabla \mathbf{u} \cdot \mathbf{h} + R_{\varrho} \nabla \mathbf{h}^* \mathbf{u} \right) dx - \mathfrak{B}(\mathbf{w}, \varphi, \varsigma) - \mathfrak{B}(\mathbf{w}, v, \zeta) \\ & + \langle \omega, G - b_{12}^0 \varsigma - b_{22}^0 g - \varkappa b_{32}^0 \rangle + \langle \psi, F - b_{11}^0 \varsigma - b_{21}^0 g - \varkappa b_{31}^0 - R \mathbf{u} \cdot \nabla \mathbf{u} \cdot \mathbf{h} \rangle \\ & + \langle \xi, M - \varkappa b_{34}^0 \rangle + n(1 - \langle 1, b_{13}^0 \varsigma \rangle) \\ & = \langle \mathfrak{d}, b_{10}^0 \varsigma + b_{20}^0 g + \varkappa b_{30}^0 + \sigma v \rangle + \langle \mathcal{D}_0, \mathbf{h} \rangle \end{aligned}$$

holds true for all $(\mathbf{H}, G, F, M) \in (C^\infty(\Omega))^6$ such that $G = \Pi G$. Here $\mathfrak{d} = 1/2 \operatorname{Tr} \mathbf{D}$, and the test functions $\mathbf{h}, g, \varsigma, v$ are defined by the solutions to adjoint problems

$$(1.48) \quad \mathbf{D} \mathbf{h} - \nabla g = \mathbf{H}, \quad \operatorname{div} \mathbf{h} = G, \quad \mathcal{L}^* \varsigma = F, \quad \mathcal{L} v = M \quad \text{in } \Omega,$$

$$(1.49) \quad \mathbf{h} = 0 \quad \text{on } \partial \Omega, \quad \Pi g = g, \quad \varsigma = 0 \quad \text{on } \Sigma_{\text{out}}, \quad v = 0 \quad \text{on } \Sigma_{\text{in}}.$$

We are now in a position to formulate the third main result of this paper.

THEOREM 1.9. *Under the above assumptions,*

$$(1.50) \quad \begin{aligned} & \mathbf{w}_\varepsilon \rightarrow \mathbf{w} \quad \text{weakly in } \mathcal{W}_0^{1-s,r'}(\Omega), \quad n_\varepsilon \rightarrow n \quad \text{in } \mathbb{R}, \\ & \psi_\varepsilon \rightarrow \psi, \quad \omega_\varepsilon \rightarrow \omega, \quad \xi_\varepsilon \rightarrow \xi \quad (*)\text{-weakly in } \mathbb{W}^{-s,r'}(\Omega) \quad \text{as } \varepsilon \rightarrow 0, \end{aligned}$$

where the limits, vector field \mathbf{w} , functionals ψ, ω, ξ , and the constant n are given by the weak solution to problem (1.42).

Note that the matrix $\mathbf{N}(\varepsilon)$ defined by equalities (1.5) meets all requirements of Theorem 1.9, and in the special case we have in representation (1.6)

$$(1.51) \quad \mathbf{D}(x) = \operatorname{div} \mathbf{T}(x) \mathbf{I} - \mathbf{T}'(x).$$

Therefore, Theorem 1.9, combined with formulae (1.3) and (1.10), implies the existence of the shape derivative for the drag functional at $\varepsilon = 0$. Straightforward calculations lead to the following result.

THEOREM 1.10. *Under the assumptions of Theorem 1.9, there exists the shape derivative*

$$\frac{d}{d\varepsilon} J_D(S_\varepsilon) \Big|_{\varepsilon=0} = L_e(\mathbf{T}) + L_u(\mathbf{w}, \omega, \psi),$$

where the linear forms L_e and L_u are defined by the equalities

$$\begin{aligned} L_e(\mathbf{T}) &= \int_{\Omega} \operatorname{div} \mathbf{T} (\nabla \mathbf{u} + \nabla \mathbf{u}^* - \operatorname{div} \mathbf{u} \mathbf{I}) \nabla \eta \mathbf{U}_\infty \, dx \\ &- \int_{\Omega} [\nabla \mathbf{u} + \nabla \mathbf{u}^* - \operatorname{div} \mathbf{u} - q \mathbf{I} - R \varrho \mathbf{u} \otimes \mathbf{u}] \mathbf{D} \nabla \eta \cdot \mathbf{U}_\infty \, dx \\ &- \int_{\Omega} [\mathbf{D}^* \nabla \mathbf{u} + \nabla \mathbf{u}^* \mathbf{D} - \nabla(\mathbf{D} \mathbf{u}) - \nabla(\mathbf{D} \mathbf{u})^*] \nabla \eta \cdot \mathbf{U}_\infty \, dx \end{aligned}$$

and

$$\begin{aligned} L_u(\mathbf{w}, \omega, \psi) &= \int_{\Omega} \mathbf{w} [\Delta \eta \mathbf{U}_\infty + R \varrho (\mathbf{u} \cdot \nabla \eta) \mathbf{U}_\infty + R \varrho (\mathbf{u} \cdot \mathbf{U}_\infty) \nabla \eta] \, dx \\ &+ \langle \omega, \nabla \eta \cdot \mathbf{U}_\infty \rangle + R \langle \psi, (\mathbf{u} \cdot \nabla \eta) (\mathbf{u} \cdot \mathbf{U}_\infty) \rangle. \end{aligned}$$

While L_e depends directly on the vector field \mathbf{T} , the linear form L_u depends on the weak solution $(\mathbf{w}, \psi, \omega)$ to problem (1.42), and thus depends on the *direction* \mathbf{T} in a very implicit manner, which is inconvenient for applications. In order to cope with this difficulty, we define the *adjoint state* $\mathbf{Y} = (\mathbf{h}, g, \varsigma, v, l)^\top$ given as a solution to the linear equation

$$(1.52) \quad \mathfrak{L} \mathbf{Y} - \mathfrak{U} \mathbf{Y} - \mathfrak{V} \mathbf{Y} = \Theta,$$

supplemented with boundary conditions (1.49). Here the operators \mathfrak{L} , \mathfrak{U} , \mathfrak{V} and the vector field Θ are defined by

$$\begin{aligned} \mathfrak{L} &= \begin{pmatrix} \Delta & -\nabla & 0 & 0 & 0 \\ \operatorname{div} & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathcal{L}^* & 0 & 0 \\ 0 & 0 & 0 & \mathcal{L} & 0 \\ 0 & 0 & -\mathbb{B}_{13} & 0 & 1 \end{pmatrix}, & \mathfrak{U} &= \begin{pmatrix} 0 & 0 & -\nabla \varphi & -\zeta \nabla & 0 \\ 0 & 0 & \Pi_{12} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \\ \mathfrak{V} &= \begin{pmatrix} R \varrho (\nabla \mathbf{u} - \mathbf{u} \nabla) & 0 & 0 & 0 & 0 \\ 0 & -\lambda^{-1} \Pi & 0 & 0 & \varkappa \Pi b_{32}^0 \\ R \mathbf{u} \cdot \nabla \mathbf{u} & b_{21}^0 & b_{11}^0 & 0 & \varkappa b_{31}^0 \\ 0 & 0 & 0 & 0 & \varkappa b_{34}^0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \end{aligned}$$

$$\begin{aligned} \Theta &= (\Delta \eta \mathbf{U}_\infty + R \varrho (\nabla \eta \otimes \mathbf{U}_\infty + \mathbf{U}_\infty \otimes \nabla \eta) \mathbf{u}, \quad \Pi (\nabla \eta \cdot \mathbf{U}_\infty), \quad R (\mathbf{u} \nabla \eta) (\mathbf{u} \mathbf{U}_\infty), 0, 0), \\ &\Pi_{2i}(\cdot) = \Pi(b_{2i}^0(\cdot)), \quad \mathbb{B}_{13}(\cdot) = \langle 1, b_{13}^0(\cdot) \rangle. \end{aligned}$$

The following theorem guarantees the existence of the adjoint state and gives the expression of the shape derivative for the drag functional in terms of the vector field \mathbf{T} .

THEOREM 1.11. *Let a given solution $\vartheta \in \mathcal{B}_\tau$, $(\zeta, m) \in X^{s,r} \times \mathbb{R}$, to problem (1.13) meet all requirements of Theorem 1.6. Then there exists positive constant τ_1 (depending only on \mathbf{U} , Ω , and r, s) such that if $\tau \in (0, \tau_1]$ and $R, \lambda^{-1} \leq \tau_1^2$, then*

there exists a unique solution $\mathbf{Y} \in (Y^{s,r})^3 \times (X^{s,r})^3 \times \mathbb{R}$ to problem (1.52), (1.49). The form L_u has the representation

$$(1.53) \quad L_u(\mathbf{w}, \psi, \omega) = \int_{\Omega} [\operatorname{div} \mathbf{T}(b_{10}^0 \varsigma + b_{20}^0 g + \sigma v + \varkappa b_{30}^0) + \mathcal{D}_0(\operatorname{div} \mathbf{T} - \mathbf{T}') \mathbf{h}] dx,$$

where the coefficients b_{ij}^0 and the operator \mathcal{D}_0 are defined by formulae (1.43) and (1.45).

Method and structure of the paper. The following aspects of our method deserve a brief description:

- extended form (1.13) of the governing equations which allows us to cope with the mass control problem;
- the splitting of the boundary value problem for the transport equation into two parts: the local problem in the vicinity of inlet, and the global problem with the modified vector field $\tilde{\mathbf{u}}$ and the empty inlet $\tilde{\Sigma}_{\text{in}}$;
- the estimates of solutions to the model problem (4.19) in the fractional Sobolev spaces, which cannot be obtained by the interpolation method;
- the very weak formulation of linearized equations introduced to assure the existence of material derivatives.

Now we can explain the organization of this paper. Section 2 is devoted to the proof of Theorem 1.6. First of all, we establish the existence of solutions to problem (1.13) using the Schauder fixed point theorem. Next we consider the linear equations for the difference of two solutions $(\mathbf{u}_i, \varphi_i)$, $i = 1, 2$, corresponding to arbitrary matrix-valued functions \mathbf{N}_i . Using Theorem 1.4 we deduce the weak formulation of the boundary value problem for linearized equations. The main result of this section is Theorem 2.3, which shows that solutions of the linearized problem are stable with respect to perturbations of data in the dual Sobolev space. This result implies the local uniqueness of solutions to problem (1.13). In section 3 we exploit Theorem 2.3 to prove the existence of the material derivative of solutions. The last section is devoted to the proof of Theorem 1.4.

2. Existence and uniqueness of local solutions. Proof of Theorem 1.6.

2.1. Existence theory. In this paragraph we establish the local solvability of problem (1.13) and prove the first part of Theorem 1.6. In our notation, c denotes generic constants, which are different in different places and depend only on Ω , \mathbf{U} , σ , and r, s . The proof is based on the following lemma which furnishes the regularity properties of composed functions. Let us consider functions $u, v : \Omega \mapsto B_K$, where $B_K = \{\mathbf{x} : |\mathbf{x}| \leq K\} \subset \mathbb{R}^3$ is the ball of radius K centered at 0.

LEMMA 2.1. *Assume that $u, v \in X^{s,r}$, $s \in (0, 1]$, $sr > 3$, and $f \in C^3(\Omega \times B_K)$. Then we have*

$$(2.1) \quad \|f(\cdot, u)\|_{X^{s,r}} \leq c(r, s) \|f\|_{C^1(\Omega \times B_K(0))} (1 + \|u\|_{X^{s,r}}),$$

$$(2.2) \quad \|f(\cdot, v) - f(\cdot, u)\|_{X^{s,r}} \leq c(r, s) \|f\|_{C^2(\Omega \times B_K)} (1 + \|u\|_{X^{s,r}} + \|v\|_{X^{s,r}}) \|u - v\|_{X^{s,r}}.$$

Proof. In order to prove (2.1) it suffices to note that

$$|f(x, u(x)) - f(y, u(y))|^r \leq c(r) \|f\|_{C^1(\Omega \times B_K)}^r (|x - y|^r + |u(x) - u(y)|^r),$$

which, in view of the inequality

$$\int_{\Omega \times \Omega} |x - y|^{r-3-rs} dx dy \leq c(r, s),$$

yields

$$|f(\cdot, u)|_{s,r,\Omega} \leq c(r, s) \|f\|_{C^1(\Omega \times B_K(0))} (1 + |u|_{s,r,\Omega}).$$

On the other hand, we have

$$\|\nabla f(\cdot, u)\|_{L^2(\Omega)} \leq \|f\|_{C^1(\Omega \times B_K(0))} \|\nabla u\|_{L^2(\Omega)}.$$

Combining obtained inequalities we get (2.1). It remains to note that (2.2) follows from (2.1) and the Hadamard formula for the first order expansion of f . \square

Fix sufficiently small positive τ such that

$$(2.3) \quad c_e \tau < \delta^*,$$

where δ^* is the constant determined in Corollary 1.5 and c_e is the constant from inequality (1.37). By virtue of Corollary 1.5, there is σ^* , depending only on Ω, \mathbf{U} , and r, s , such that, for all $\vartheta \in \mathcal{B}_\tau$ and $\sigma > \sigma^*$, problems (1.28) and (1.29) have solutions satisfying inequalities (1.35). Finally, fix an arbitrary $\sigma > \sigma^*$.

We solve problem (1.13) by an application of the Schauder fixed point theorem in the following framework. Choose an arbitrary element $\vartheta \in \mathcal{B}_\tau$. As mentioned above, the problem

$$(2.4) \quad \mathbf{u} \cdot \nabla \varphi_1 + \sigma \varphi_1 = \Psi_1[\vartheta] + m \mathbf{g} \varrho \text{ in } \Omega, \quad \varphi_1 = 0 \text{ on } \Sigma_{\text{in}},$$

has a unique solution satisfying the inequality

$$(2.5) \quad \|\varphi_1\|_{X^{s,r}} \leq c(\Omega, U, \sigma, r, s) (\|\Psi_1[\vartheta]\|_{X^{s,r}} + |m|).$$

Next, define \mathbf{u}_1 and π_1 to be the solutions of the boundary value problem for the Stokes equations

$$(2.6) \quad \begin{aligned} \Delta \mathbf{u}_1 - \nabla \pi_1 &= \mathcal{A}(\mathbf{u}) + R\mathcal{B}(\varrho, \mathbf{u}, \mathbf{u}) \equiv F[\vartheta] \text{ in } \Omega, \\ \varrho_0 \operatorname{div} \mathbf{u}_1 &= \Pi(\mathbf{g}\sigma\varphi_1 - \mathbf{g}\varrho_0\Psi[\vartheta] - \mathbf{g}m\varrho_0) \text{ in } \Omega, \\ \mathbf{u}_1 &= 0 \text{ on } \partial\Omega, \quad \pi_1 - \Pi\pi_1 = 0, \end{aligned}$$

where m is given by (1.13c). By Lemma 1.3, this problem admits a unique solution such that

$$(2.7) \quad \|\mathbf{u}_1\|_{Y^{s,r}} + \|\pi_1\|_{X^{s,r}} \leq c(\|F[\vartheta]\|_{Z^{s,r}} + |\Psi[\vartheta]|_{X^{s,r}} + \|\varphi_1\|_{X^{s,r}} + |m|).$$

Equations (2.4), (2.6), (1.13c) define the mapping $\Xi : \vartheta \rightarrow \vartheta_1 = (\mathbf{u}_1, \pi_1, \varphi_1)$. We claim that for a suitable choice of the constant τ , Ξ is a weakly continuous automorphism of the ball \mathcal{B}_τ . We begin with the estimates for nonlinear operators present in (2.4). Fix an arbitrary $\vartheta \in \mathcal{B}_\tau$. Applying inequality (2.2) from Lemma 2.1 to the function H which is a part of $\Psi[\vartheta]$, we obtain $\|H(\varphi)\|_{X^{s,r}} \leq c\tau^2$, which leads to the estimate

$$(2.8) \quad \|\Psi[\vartheta]\|_{X^{s,r}} \leq \frac{c}{\lambda} (\|q_0\|_{C^1(\Omega)} + \|\pi\|_{X^{s,r}}) + c\tau^2 \leq \frac{c}{\lambda} + c\tau^2 \leq c\tau^2.$$

Since, under assumptions of Theorem 1.6, $X^{s,r}(\Omega)$ is a Banach algebra and $\|\varrho\|_{X^{s,r}} \leq c + \|\varphi\|_{X^{s,r}} \leq \text{const}$, we conclude from this and (2.5) that

$$(2.9) \quad \|\varphi_1\|_{X^{s,r}} \leq c/\lambda + c\tau^2 + c|m| \leq c\tau^2 + c|m|.$$

In order to estimate the right-hand side of the first equation in (2.6) we introduce the vector function $\mathbf{z} = (\mathbf{u}, \nabla \mathbf{u}, \pi, \varphi)$ and proceed as follows. It can be easily seen that $\|\mathbf{z}\|_{X^{s,r}} \leq \|\vartheta\|_E \leq \tau$ and $|\mathbf{z}| \leq c\tau$. Recall that the operator \mathcal{B} constitutes a cubic polynomial of \mathbf{u} and ϱ . By Lemma 2.1, we have

$$(2.10) \quad R\|\mathcal{B}(\varrho, \mathbf{u}, \mathbf{u})\|_{X^{s,r}} \leq cR(1 + \|\varrho\|_{X^{s,r}} + \|\mathbf{z}\|_{X^{s,r}}) \leq c\tau^2(1 + \tau) \quad \text{in } \mathcal{B}_\tau.$$

Next note that

$$\|\mathcal{A}(\mathbf{u})\|_{Z^{s,r}} \leq c(\|\mathbf{g} - 1\|_{C^2(\Omega)} + \|\mathbf{N} - \mathbf{I}\|_{C^2(\Omega)})(1 + \|\mathbf{u}\|_{Y^{s,r}}) \leq c\tau^2\|\mathbf{u}\|_{Y^{s,r}},$$

which along with (1.39) and (2.10) implies

$$(2.11) \quad \|F[\vartheta]\|_{Z^{s,r}} \leq c\tau^2(1 + \tau) \quad \text{in } \mathcal{B}_\tau.$$

Combining inequalities (2.8) and (2.9) we get the estimate

$$\|\sigma\varphi_1 + \Psi[\vartheta]\|_{X^{s,r}} \leq c\tau^2.$$

From this, (2.11), (2.7), and Lemma 1.3 we finally obtain

$$(2.12) \quad \|\mathbf{u}_1\|_{Y^{s,r}} + |\pi_1|_{X^{s,r}} \leq c\tau^2 + c|m|.$$

It remains to estimate m . Recall that the vector field \mathbf{u} and parameter σ meet all requirements of Corollary 1.5. Therefore, problem (1.13b) has a unique solution $\zeta \in W^{s,r}(\Omega)$ for all s, r satisfying (1.38). In particular, inequalities (1.35) yield the estimate $\|\zeta\|_{X^{s,r}} \leq c$. Since, by virtue of (1.38), the pair $s = 2/3, r = 6$ is admissible and the embedding $W^{2/3,6}(\Omega) \hookrightarrow C^{1/6}(\Omega)$ is bounded, estimates (1.33) and (1.34) for $rs > 3$ yield

$$(2.13) \quad \|\zeta\|_{C^{1/6}(\Omega)} + \|\zeta\|_{W^{1,2}(\Omega)} \leq C(\mathbf{U}, \Omega, \sigma).$$

Recalling that $\operatorname{div} \mathbf{u} = \operatorname{div} \mathbf{u}$, we obtain $|\operatorname{div} \mathbf{u}| \leq c_e\tau$. From this, the inequality $|\mathbf{g}| \leq 1 + c\tau^2$, and the maximum principle (1.34), we conclude that

$$(2.14) \quad \|\zeta\|_{C(\Omega)} \leq (1 + c\tau^2)(1 - \sigma^{-1}c\tau)^{-1} \leq (1 - c\tau)^{-1},$$

which leads to the following estimate:

$$|1 - \zeta| \leq c\tau(1 - c\tau)^{-1}.$$

Now we can estimate the right-hand side of (1.13c). Rewrite the first integral in the form

$$\begin{aligned} \int_{\Omega} \mathbf{g}(1 - \zeta - \varrho_0^{-1}\zeta\varphi) \, dx &= \int_{\Omega} (1 - \zeta)^+ \, dx + \int_{\Omega} (\mathbf{g} - 1)(1 - \zeta - \varrho_0^{-1}\zeta\varphi) \, dx \\ &\quad - \int_{\Omega} ((1 - \zeta)^- + \varrho_0^{-1}\zeta\varphi) \, dx. \end{aligned}$$

We have

$$|(\mathbf{g} - 1)(1 - \zeta - \varrho_0^{-1}\zeta\varphi)| \leq c\tau^2, \quad |(1 - \zeta)^- + \varrho_0^{-1}\zeta\varphi| \leq c_e\tau + c\tau(1 - c\tau)^{-1}.$$

On the other hand, we have $\|(1 - \zeta)^+\|_{C^{1/6}(\Omega)} \leq c(\mathbf{U}, \Omega, \sigma)$ and $(1 - \zeta)^+ = 1$ on Σ_{out} . Hence

$$\int_{\Omega} (1 - \zeta)^+ dx > \kappa(\mathbf{U}, \Omega, \sigma) > 0.$$

Thus, we get

$$\varkappa^{-1} \geq \kappa(1 - c\kappa^{-1}\tau(1 - c\tau)^{-1}).$$

In particular, there is a positive τ_0 depending only on \mathbf{U} , Ω , and σ , such that

$$|\varkappa| \leq c \text{ for all } \tau \leq \tau_0.$$

Repeating these arguments and using inequalities (2.8), (1.13c), we arrive at $|m| \leq c\tau^2$. Combining this estimate with (2.9) and (2.12), we finally obtain $\|\vartheta_1\|_{X^{s,r}} \leq c\tau^2$. Choose sufficiently small $\tau_0 = \tau_0(\mathbf{U}, \Omega, \sigma)$ such that $c\tau_0^2 < \tau_0$. Thus, for all $\tau \leq \tau_0$, Ξ maps the ball \mathcal{B}_τ into itself. Let us show that Ξ is weakly continuous. Choose an arbitrary sequence $\vartheta_n \in \mathcal{B}_\tau$ such that $\vartheta_n = (\mathbf{u}_n, \pi_n, \varphi_n)$ converges weakly in E to some ϑ . Since the ball \mathcal{B}_τ is closed and convex, ϑ belongs to \mathcal{B}_τ . Let us consider the corresponding sequences of the elements $\vartheta_{1,n} = \Xi(\vartheta_n) \in \mathcal{B}_\tau$ and functions ζ_n . There are subsequences $\{\vartheta_{1,j}\} \subset \{\vartheta_{1,n}\}$ and $\{\zeta_j\} \subset \{\zeta_n\}$ such that $\vartheta_{1,j}$ converges weakly in E to some element $\vartheta_1 \in \mathcal{B}_\tau$ and ζ_j converges weakly in $X^{s,r}$ to some function $\zeta \in X^{s,r}$. Since the embedding $E \hookrightarrow C(\Omega)^5$ is compact, we have $\vartheta_n \rightarrow \vartheta$, $\vartheta_{1j} \rightarrow \vartheta_1$ in $C(\Omega)^5$, and

$$\nabla \zeta_j \rightharpoonup \nabla \zeta \text{ weakly in } L^2(\Omega), \quad \zeta_j \rightarrow \zeta \text{ in } C(\Omega).$$

Substituting ϑ_j and $\vartheta_{1,j}$ into equations (2.4), (2.6), (1.13c) and letting $j \rightarrow \infty$, we obtain that the limits ϑ and ϑ_1 also satisfy (2.4), (2.6), (1.13c). Thus, we get $\vartheta_1 = \Xi(\vartheta)$. Since for given ϑ , a solution to equations (2.4), (2.6) is unique, we conclude from this that all weakly convergent subsequences of $\vartheta_{1,n}$ have the unique limit ϑ_1 . Therefore, the whole sequence $\vartheta_{1,n} = \Xi(\vartheta_n)$ converges weakly to $\Xi(\vartheta)$. Hence the mapping $\Xi : \mathcal{B}_\tau \mapsto \mathcal{B}_\tau$ is weakly continuous, and, by virtue of the Schauder fixed point theory, there is $\vartheta \in \mathcal{B}(\tau)$ such that $\vartheta = \Xi(\vartheta)$.

It remains to prove that ϑ is given by a solution to problem (1.13a). For $\vartheta_1 = \vartheta$, the only difference between problems (1.13a) and (2.6), (1.13c) is the presence of the projector Π in the right-hand side of (2.6). Hence, it suffices to show that

$$(2.15) \quad \Pi(\varrho_0^{-1} \mathbf{g} \sigma \varphi - \mathbf{g} \Psi[\vartheta] - \mathbf{g} m) = \varrho_0^{-1} \mathbf{g} \sigma \varphi - \mathbf{g} \Psi[\vartheta] - \mathbf{g} m.$$

To this end we note that φ is a generalized solution to the transport equation

$$\mathbf{u} \cdot \nabla \varphi + \sigma \varphi = \Psi_1[\vartheta] + m \mathbf{g} \varrho.$$

Using ζ as a test function and recalling the integral identity (1.30) we obtain

$$\sigma \int_{\Omega} \varphi \mathbf{g} dx = \int_{\Omega} \zeta (\Psi_1[\vartheta] + m \mathbf{g} \varrho) dx.$$

On the other hand, equality (1.13c) reads

$$\int_{\Omega} \zeta (\varrho_0^{-1} \Psi_1[\vartheta] + m \mathbf{g} (1 + \varphi \varrho_0^{-1})) dx = \int_{\Omega} (\mathbf{g} \Psi[\vartheta] + \mathbf{g} m) dx.$$

Combining these equalities and noting that $1 + \varrho_0^{-1}\varphi = \varrho/\varrho_0$, we obtain

$$\int_{\Omega} \left(\frac{\mathfrak{g}\sigma}{\varrho_0}\varphi - \mathfrak{g}\Psi[\vartheta] - m\mathfrak{g} \right) dx = 0$$

which yields (2.15), and the proof of Theorem 1.6 is completed. \square

2.2. Uniqueness and stability. In this paragraph we prove that, under the assumptions of Theorem 1.6, a solution to problem (1.13) is unique, and we investigate in detail the dependence of the solution on the matrix function \mathbf{N} .

Weak formulation of linearized equations. Assume that matrices \mathbf{N}_i , $i = 1, 2$, satisfy conditions of Theorem 1.6, and denote by $(\vartheta_i, \zeta_i, m_i) \in E \times X^{s,r} \times \mathbb{R}$, $i = 1, 2$, the corresponding solutions to problem (1.13). Recall that the solutions $(\vartheta_i, \zeta_i, m_i)$, together with the constants \varkappa_i , satisfy the inequalities

$$(2.16) \quad |m_i| + \|\vartheta_i\|_E \leq c\tau, \quad |\varkappa_i| + \|\zeta_i\|_{X^{s,r}} \leq c,$$

where the constant c depends only on \mathbf{U}, Ω, r, s , and σ . We denote $\mathbf{u}_i = \mathbf{u}_0 + \mathbf{u}_i$, $i = 1, 2$, the solutions to (1.9) for \mathbf{N}_i , $i = 1, 2$. Now set

$$\begin{aligned} \mathfrak{d} &= \mathfrak{g}_1 - \mathfrak{g}_2 \equiv \sqrt{\det \mathbf{N}_1} - \sqrt{\det \mathbf{N}_2}, \\ \mathbf{w} &= \mathbf{u}_1 - \mathbf{u}_2, \quad \omega = \pi_1 - \pi_2, \quad \psi = \varphi_1 - \varphi_2, \quad \xi = \zeta_1 - \zeta_2, \quad n = m_1 - m_2. \end{aligned}$$

It follows from (1.13) that

$$(2.17) \quad \begin{aligned} \mathbf{u}_1 \nabla \psi + \sigma \psi &= -\mathbf{w} \cdot \nabla \varphi_2 + b_{11}\psi + b_{12}\omega + b_{13}n + b_{10}\mathfrak{d} \quad \text{in } \Omega, \\ \Delta \mathbf{w} - \nabla \omega &= \mathcal{A}_1(\mathbf{w}) + R\mathcal{C}_1(\psi, \mathbf{w}) + \mathcal{D} \quad \text{in } \Omega, \\ \operatorname{div} \mathbf{w} &= b_{21}\psi + b_{22}\omega + b_{23}n + b_{20}\mathfrak{d} \quad \text{in } \Omega, \\ -\operatorname{div}(\mathbf{u}_1 \xi) + \sigma \xi &= \operatorname{div}(\zeta_2 \mathbf{w}) + \sigma \mathfrak{d} \quad \text{in } \Omega, \\ \mathbf{w} &= 0 \quad \text{on } \partial\Omega, \quad \psi = 0 \quad \text{on } \Sigma_{\text{in}}, \quad \xi = 0 \quad \text{on } \Sigma_{\text{out}}, \\ \omega - \Pi\omega = 0, \quad n &= \varkappa \int_{\Omega} \left(b_{31}\psi + b_{32}\omega + b_{34}\xi + b_{30}\mathfrak{d} \right) dx. \end{aligned}$$

Here the coefficients are given by the formula

$$(2.18) \quad \begin{aligned} b_{11} &= \sigma(1 - \mathfrak{g}_2) + \mathfrak{g}_2\Psi[\vartheta_1] - \mathfrak{g}_2\varrho_2\Phi_1(\varphi_1, \varphi_2) + \mathfrak{g}_2m_2 - \frac{\sigma\mathfrak{g}_2}{\varrho_0}(\varphi_1 + \varphi_2), \\ b_{12} &= \lambda^{-1}\varrho_2\mathfrak{g}_2, \quad b_{13} = \mathfrak{g}_2\varrho_1, \quad b_{10} = \varrho_1\Psi[\vartheta_1] - \frac{\sigma}{\varrho_0}\varphi_1^2 - \sigma\varphi_1 + m_1\varrho_1, \\ b_{21} &= \mathfrak{g}_2 \left(\frac{\sigma}{\varrho_0} + \Phi_1(\varphi_1, \varphi_2) \right), \quad b_{22} = -\mathfrak{g}_2/\lambda, \\ b_{23} &= -\mathfrak{g}_2, \quad b_{20} = \sigma\varphi_1\varrho_0^{-1} - \Psi[\vartheta_1] - m_2, \\ b_{31} &= \varrho_0^{-1}\zeta_1 \left(\sigma(1 - \mathfrak{g}_2) + \mathfrak{g}_2\Psi[\vartheta_1] - \mathfrak{g}_2\varrho_2\Phi_1(\varphi_1, \varphi_2) \right. \\ &\quad \left. - \frac{\mathfrak{g}_2\sigma}{\varrho_0}(\varphi_1 + \varphi_2) \right) - \mathfrak{g}_2\Phi_1(\varphi_1, \varphi_2) + m_2\mathfrak{g}_1\varrho_0^{-1}\zeta_2, \\ b_{32} &= \varrho_0^{-1}\zeta_1b_{12} - b_{22}, \quad b_{34} = \varrho_0^{-1}\Psi_1[\vartheta_2] + m_2\mathfrak{g}_1(1 + \varrho_0^{-1}\varphi_1), \\ b_{30} &= \varrho_0^{-1}\zeta_1(b_{10} - m_1\varrho_1) + \Psi[\vartheta_1] - m_2(1 - \zeta_2 - \varrho_0^{-1}\zeta_2\varphi_2), \\ \Phi_1(\varphi_1, \varphi_2) &= (p'(\varrho_0)\varrho_0)^{-1}\sigma \int_0^1 H'(\varphi_1s + \varphi_2(1-s)) ds, \end{aligned}$$

and the operators \mathcal{C}_1 and \mathcal{D} are defined by the equalities

$$\begin{aligned} \mathcal{C}(\mathbf{w}) &= \mathcal{B}_1(\psi, \mathbf{u}_1, \mathbf{u}_1) + \mathcal{B}_1(\varrho_2, \mathbf{w}, \mathbf{u}_1) + \mathcal{B}_1(\varrho_2, \mathbf{u}_2, \mathbf{w}), \\ \mathcal{D} &= \mathcal{A}_1(\mathbf{u}_2) - \mathcal{A}_2(\mathbf{u}_2) + R(\mathcal{B}_1(\varrho_2, \mathbf{u}_2, \mathbf{u}_2) - \mathcal{B}_2(\varrho_2, \mathbf{u}_2, \mathbf{u}_2)), \end{aligned}$$

where \mathcal{A}_i and \mathcal{B}_i are given by (1.8) with \mathbf{N}_i instead of \mathbf{N} .

We consider \mathcal{D} and \mathfrak{d} as *given functions* and equality (2.17) as the system of equations and boundary conditions for unknowns \mathbf{w} , ψ , ξ , and n . The next step is crucial for further analysis. We replace equations (2.17) by an integral identity, which leads to the notion of a very weak solution to problem (2.17). To this end choose an arbitrary function $(\mathbf{H}, G, F, M) \in C^\infty(\Omega)^6$ such that $G - \Pi G = 0$, and consider the auxiliary boundary value problems

$$(2.19) \quad \mathcal{L}^* \zeta = F, \quad \mathcal{L}v = M \text{ in } \Omega, \quad \zeta = 0 \text{ on } \Sigma_{\text{out}}, \quad v = 0 \text{ on } \Sigma_{\text{in}}.$$

$$(2.20) \quad \Delta \mathbf{h} - \nabla g = \mathbf{H}, \quad \operatorname{div} \mathbf{h} = \Pi G \text{ in } \Omega, \quad \mathbf{h} = 0 \text{ on } \partial\Omega, \quad \Pi g = g.$$

Since, under the assumptions of Theorem 1.6, \mathbf{u} and σ meet all requirements of Corollary 1.5, each of problems (2.19) has a unique solution such that

$$(2.21) \quad \|\zeta\|_{W^{s,r}(\Omega)} \leq c\|F\|_{W^{s,r}(\Omega)}, \quad \|v\|_{W^{s,r}(\Omega)} \leq c\|M\|_{W^{s,r}(\Omega)},$$

where c depends only on \mathbf{U} , Ω , r , s , and σ . On the other hand, by virtue of Lemma 1.3, problem (2.20) has a unique solution satisfying the inequality

$$(2.22) \quad \|\mathbf{h}\|_{W^{1+s,r}(\Omega)} + \|g\|_{W^{s,r}(\Omega)} \leq c\|\mathbf{H}\|_{W^{1+s,r}(\Omega)} + c\|G\|_{W^{s,r}(\Omega)}.$$

Recall that $\mathbf{w} \in W^{2,2}(\Omega)^3 \cap C^1(\Omega)^3$ vanishes on $\partial\Omega$, and $(\omega, \psi, \xi) \in W^{1,2}(\Omega)^3 \cap C(\Omega)^3$. Multiplying both sides of the first equation in system (2.17) by ζ , both sides of the fourth equation in (2.17) by v , integrating the results over Ω , and using the Green formula for the Stokes equations, we obtain the system of integral equalities

$$(2.23) \quad \begin{aligned} \int_{\Omega} \psi F \, dx &= \int_{\Omega} (-\mathbf{w} \cdot \nabla \varphi_2 + b_{11}\psi + b_{12}\omega + b_{13}n + b_{10}\mathfrak{d})\zeta \, dx, \\ \int_{\Omega} \mathbf{w}\mathbf{H} \, dx + \int_{\Omega} \omega G \, dx &= \int_{\Omega} (b_{21}\psi + b_{22}\omega + b_{23}n + b_{20}\mathfrak{d})g \, dx \\ + \int_{\Omega} (\mathcal{A}_1(\mathbf{w}) + R\mathcal{C}(\mathbf{w}, \psi) + \mathcal{D})\mathbf{h} \, dx, & \quad \int_{\Omega} \xi M \, dx = \int_{\Omega} (\operatorname{div}(\zeta_2 \mathbf{w}) + \sigma \mathfrak{d})v \, dx. \end{aligned}$$

Next, since $\operatorname{div}(\varrho_2 \mathbf{u}_2) = 0$, we have

$$\begin{aligned} &\int_{\Omega} (\mathcal{B}_1(\varrho_2, \mathbf{w}, \mathbf{u}_1) + \mathcal{B}_1(\varrho_2, \mathbf{u}_2, \mathbf{w})) \cdot \mathbf{h} \, dx \\ &= \int_{\Omega} \varrho_2 \mathbf{w} \cdot \left(\nabla(\mathbf{N}_1^{-1} \mathbf{u}_1) \cdot (\mathbf{N}_1^{-1} \mathbf{h}) - (\mathbf{N}_1^{-1})^* \nabla(\mathbf{N}_1^{-1} \mathbf{h})^* \mathbf{u}_2 \right) \, dx. \end{aligned}$$

On the other hand, integration by parts gives

$$\int_{\Omega} \operatorname{div}(\zeta_2 \mathbf{w})v \, dx = - \int_{\Omega} \zeta_2 \mathbf{w} \nabla v \, dx.$$

Using these identities and recalling the duality pairing, we can collect relations (2.23), together with the expression for n , in one integral identity

$$\begin{aligned}
 (2.24) \quad & \int_{\Omega} \mathbf{w} \left(\mathbf{H} - R\varrho_2 \nabla(\mathbf{N}_1^{-1} \mathbf{u}_1) \cdot (\mathbf{N}_1^{-1} \mathbf{h}) + R\varrho_2 (\mathbf{N}_1^{-1})^* \nabla(\mathbf{N}_1^{-1} \mathbf{h})^* \mathbf{u}_2 \right) dx \\
 & - \mathfrak{B}(\mathbf{w}, \varphi_2, \varsigma) - \mathfrak{B}(\mathbf{w}, v, \zeta_2) - \mathfrak{A}_1(\mathbf{w}, \mathbf{h}) + \langle \omega, G - b_{12}\varsigma - b_{22}g - \varkappa b_{32} \rangle \\
 & + \langle \psi, F - b_{11}\varsigma - b_{21}g - \varkappa b_{31} - R\mathbf{u}_1 \cdot \nabla(\mathbf{N}_1^{-1} \mathbf{u}_1) \cdot \mathbf{N}_1^{-1} \mathbf{h} \rangle \\
 & + \langle \xi, M - \varkappa b_{34} \rangle + n - n \langle 1, b_{13}\varsigma + b_{23}g \rangle = \langle \mathfrak{d}, b_{10}\varsigma + b_{20}g + \varkappa b_{30} + \sigma v \rangle + \langle \mathcal{D}, \mathbf{W} \rangle.
 \end{aligned}$$

Here, the trilinear form \mathfrak{B} and the bilinear form \mathfrak{A}_1 are defined by the equalities

$$\mathfrak{B}(\mathbf{w}, \varphi_2, \varsigma) = - \int_{\Omega} \varsigma \mathbf{w} \cdot \nabla \varphi_2 dx, \quad \mathfrak{A}_1(\mathbf{w}, \mathbf{h}) = \int_{\Omega} \mathcal{A}_1(\mathbf{w}) \cdot \mathbf{h} dx.$$

Note that relations (2.24) are well defined for all $\mathbf{w} \in \mathcal{W}_0^{1-s, r'}(\Omega)$ and $\psi, \xi \in \mathbb{W}^{-s, r'}(\Omega)$. It is obviously true for all terms, with the possible exception of \mathfrak{A}_1 and \mathfrak{B} . Well-posedness of the form \mathfrak{B} follows from Lemma 1.7. The well-posedness of the form \mathfrak{A}_1 results from the following lemma; the proof is given in Appendix A.

LEMMA 2.2. *Let $sr > 3$, $1/2 \leq s \leq 1$, and $\mathbf{w} \in \mathcal{W}_0^{1-s, r'}(\Omega) \cap W_0^{1,2}(\Omega)$, $\mathbf{h} \in W^{1+s, r}(\Omega)$, and \mathbf{N}_1 satisfy (1.39). Then there is a constant c depending only on s, r , and Ω such that*

$$(2.25) \quad |\mathfrak{A}_1(\mathbf{w}, \mathbf{h})| \leq c\tau^2 \|\mathbf{w}\|_{\mathcal{W}_0^{1-s, r'}(\Omega)} \|\mathbf{h}\|_{W^{1+s, r}(\Omega)}.$$

Consequently, the form \mathfrak{A}_1 can be continuously extended to $\mathfrak{A}_1 : \mathcal{W}_0^{1-s, r'}(\Omega)^3 \times W^{1+s, r}(\Omega)^3 \mapsto \mathbb{R}$.

Thus, relations (2.24) are well defined for all $(\mathbf{w}, \psi, \omega, \xi) \in \mathcal{W}_0^{1-s, r'}(\Omega)^3 \times \mathbb{W}^{-s, r'}(\Omega)^3$. Equalities (2.24) along with (2.19), (2.20) are called *the very weak formulation* of problem (2.17). The natural question is the uniqueness of solutions to such weak formulation. The following theorem, which is the main result of this section, guarantees the uniqueness of very weak solutions for sufficiently small τ .

THEOREM 2.3. *Let s, r , and σ satisfy condition (1.38), parameters λ, R and matrices $\mathbf{N}_i, i = 1, 2$, satisfy conditions (1.39), and constants τ meet all requirements of Theorem 1.6. Also let the solutions $(\vartheta_i, \zeta_i, m_i), i = 1, 2$, to problem (1.13) with the matrices $\mathbf{N}_i, i = 1, 2$, belong to $\mathcal{B}_\tau \times X^{s, r} \times \mathbb{R}$. Furthermore, assume that, for any $(\mathbf{H}, G, F, M) \in C^\infty(\Omega)^6$ and for $(\varsigma, v, \mathbf{h}, g)$ satisfying (2.19)–(2.20), the elements $(\mathbf{w}, \omega, \psi, \xi) \in \mathcal{W}_0^{1-s, r'}(\Omega)^3 \times \mathbb{W}^{-s, r'}(\Omega)^3$ and the constant n satisfy identity (2.24).*

Then there are constants c, τ_1 depending only on s, r, σ , and Ω, \mathbf{U} such that, for $\tau \in (0, \tau_1]$, we have

$$\begin{aligned}
 (2.26) \quad & \|\mathbf{w}\|_{\mathcal{W}_0^{1-s, r'}(\Omega)} + \|\omega\|_{\mathbb{W}^{-s, r'}(\Omega)} + \|\psi\|_{\mathbb{W}^{-s, r'}(\Omega)} + \|\xi\|_{\mathbb{W}^{-s, r'}(\Omega)} + |n| \\
 & \leq c(\|\mathcal{D}\|_{L^1(\Omega)} + \|\mathfrak{d}\|_{\mathbb{W}^{-s, r'}(\Omega)}).
 \end{aligned}$$

Proof. The proof is based upon two auxiliary lemmas; the first lemma establishes the bounds for coefficients of problem (2.17).

LEMMA 2.4. *Under the assumptions of Theorem 2.3, all the coefficients of identity (2.24) satisfy the inequalities $\|b_{ij}\|_{X^{s, r}} \leq c$; furthermore*

$$\begin{aligned}
 (2.27) \quad & \|b_{12}\|_{X^{s, r}} + \|b_{22}\|_{X^{s, r}} + \|b_{11}\|_{X^{s, r}} + \|b_{10}\|_{X^{s, r}} + \|b_{20}\|_{X^{s, r}} \leq c\tau, \\
 & \|b_{31}\|_{X^{s, r}} + \|b_{32}\|_{X^{s, r}} + \|b_{34}\|_{X^{s, r}} \leq c\tau.
 \end{aligned}$$

Proof. The proof follows from Lemma 2.1 combined with formula (2.18). □

In order to formulate the second auxiliary result we introduce the following denotations:

$$\begin{aligned} \mathfrak{J}_1 &= \langle \psi, b_{11}\varsigma \rangle + \langle \omega, b_{12}\varsigma \rangle + \langle \mathfrak{d}, b_{10}\varsigma \rangle, & \mathfrak{J}_2 &= \langle \psi, b_{21}g \rangle + \langle \omega, b_{22}g \rangle + \langle \mathfrak{d}, b_{20}g \rangle, \\ \mathfrak{J}_3 &= \varkappa(\langle \psi, b_{31} \rangle + \langle \omega, b_{32} \rangle + \langle \xi, b_{34} \rangle + \langle \mathfrak{d}, b_{30} \rangle), & \mathfrak{J}_4 &= \langle \psi, \mathbf{u}_2 \nabla(\mathbf{N}_1 \mathbf{u}_1) \cdot \mathbf{N}_1^{-1} \mathbf{h} \rangle, \\ \mathfrak{J}_5 &= \int_{\Omega} \varrho_2 \mathbf{w} \cdot \left(\nabla(\mathbf{N}_1^{-1} \mathbf{u}_1) \cdot (\mathbf{N}_1^{-1} \mathbf{h}) - (\mathbf{N}_1^{-1})^* \nabla(\mathbf{N}_1^{-1} \mathbf{h})^* \mathbf{u}_2 \right) dx, \\ \mathfrak{G} &= \|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)} + \|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\omega\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\xi\|_{\mathbb{W}^{-s,r'}(\Omega)}, \\ \mathfrak{Q} &= \|\mathbf{H}\|_{\mathcal{W}^{s-1,r}(\Omega)} + \|G\|_{W^{s,r}(\Omega)} + \|F\|_{W^{s,r}(\Omega)} + \|M\|_{W^{s,r}(\Omega)}. \end{aligned}$$

LEMMA 2.5. *Under the assumptions of Theorem 2.3, there is a constant c , depending only on \mathbf{U} , Ω , s , r , and σ , such that*

$$(2.28) \quad \mathfrak{J}_1 \leq c\tau\mathfrak{Q}[\mathfrak{G} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}],$$

$$(2.29) \quad \mathfrak{J}_2 \leq c\mathfrak{Q}[\tau\mathfrak{G} + \|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}],$$

$$(2.30) \quad \mathfrak{J}_3 \leq c\tau[\mathfrak{G} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}], \quad \mathfrak{J}_4 + \mathfrak{J}_5 \leq c\mathfrak{Q}\mathfrak{G}.$$

Proof. We have

$$\begin{aligned} \langle \psi, b_{11}\varsigma \rangle + \langle \omega, b_{12}\varsigma \rangle + \langle \mathfrak{d}, b_{10}\varsigma \rangle &\leq \|b_{11}\varsigma\|_{W^{s,r}(\Omega)} \|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} \\ &\quad + \|b_{12}\varsigma\|_{W^{s,r}(\Omega)} \|\omega\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|b_{10}\varsigma\|_{W^{s,r}(\Omega)} \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}. \end{aligned}$$

Recall that, for $rs > 3$, $W^{s,r}(\Omega)$ is a Banach algebra. From this, estimate (2.21), and inequalities (2.27), we obtain

$$\begin{aligned} &\|b_{11}\varsigma\|_{W^{s,r}(\Omega)} \|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|b_{12}\varsigma\|_{W^{s,r}(\Omega)} \|\omega\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|b_{10}\varsigma\|_{W^{s,r}(\Omega)} \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)} \\ &\leq c\|\varsigma\|_{W^{s,r}(\Omega)} (\|b_{11}\varsigma\|_{W^{s,r}(\Omega)} \|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|b_{12}\varsigma\|_{W^{s,r}(\Omega)} \|\omega\|_{\mathbb{W}^{-s,r'}(\Omega)} \\ &\quad + \|b_{10}\varsigma\|_{W^{s,r}(\Omega)} \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}) \\ &\leq c\tau\|F\|_{W^{s,r}(\Omega)} (\|\psi\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\omega\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}), \end{aligned}$$

which gives (2.28). Repeating these arguments and using inequality (2.21), we obtain the estimates for \mathfrak{J}_2 and \mathfrak{J}_3 . Next we have

$$\begin{aligned} \|\mathbf{u}_2 \cdot \nabla(\mathbf{N}_1^{-1} \mathbf{u}_1) \cdot \mathbf{N}_1^{-1} \mathbf{h}\|_{W^{s,r}(\Omega)} &\leq c\|\mathbf{u}_2\|_{W^{s,r}(\Omega)} \|\mathbf{u}_1\|_{W^{1+s,r}(\Omega)} \|\mathbf{h}\|_{W^{1+s,r}(\Omega)} \\ &\leq c\|\mathbf{H}\|_{W^{1+s,r}(\Omega)}, \end{aligned}$$

which gives the estimate for \mathfrak{J}_4 . Since the embeddings $W^{s,r}(\Omega) \hookrightarrow C(\Omega)$, $W^{1+s,r}(\Omega) \hookrightarrow C^1(\Omega)$ are bounded and $\|\mathbf{N}^{\pm 1}\|_{C^1(\Omega)} \leq c$, we have

$$\varrho_2 |\nabla(\mathbf{N}_1^{-1} \mathbf{u}_1)| |\mathbf{N}_1^{-1} \mathbf{h}| + \varrho_2 |(\mathbf{N}_1^{-1} \mathbf{u}_2)| |\nabla(\mathbf{N}_1^{-1} \mathbf{h})| \leq c\|\mathbf{h}\|_{W^{1+s,r}(\Omega)},$$

which leads to the inequality

$$\mathfrak{J}_5 \leq c\|\mathbf{h}\|_{W^{1+s,r}(\Omega)} \|\mathbf{w}\|_{L^1(\Omega)} \leq c(\|\mathbf{H}\|_{\mathcal{W}^{s-1,r}(\Omega)} + \|G\|_{W^{s,r}(\Omega)}) \|\mathbf{w}\|_{\mathcal{W}^{1-s,r'}(\Omega)},$$

and the proof of Lemma 2.5 is completed. \square

Let us return to the proof of Theorem 2.3. It follows from the duality principle that the theorem is proved provided we show that, under the assumptions of Theorem 1.6, the following inequality holds:

$$(2.31) \quad \begin{aligned} & \sup_{\Omega(\mathbf{H}, G, F, M)=1} (\langle \mathbf{w}, \mathbf{H} \rangle + \langle \omega, G \rangle + \langle \psi, F \rangle + \langle \xi, M \rangle) + |n| \\ & \leq c\tau (\mathfrak{G}(\mathbf{w}, \omega, \psi, \xi) + |n|) + c(\|\mathcal{D}\|_{L^1(\Omega)} + \|\mathfrak{d}\|_{\mathbb{W}^{-s, r'}(\Omega)}), \end{aligned}$$

where the constant c depends only on Ω, \mathbf{U} and r, s, σ . Therefore, our task is to estimate step by step all terms in the left-hand side of (2.31). We begin with an estimate for the term $\langle \psi, F \rangle$. To this end, take $\mathbf{H} = \mathbf{h} = 0, G = g = 0, M = v = 0$, and rewrite identity (2.24) in the form

$$\langle \psi, F \rangle = \mathfrak{B}(\mathbf{w}, \varphi_2, \varsigma) + \mathfrak{J}_1 + \mathfrak{J}_3 + n\langle 1, b_{13}\varsigma \rangle - n.$$

By virtue of Lemma 1.7 and estimate (2.22), we have

$$(2.32) \quad \mathfrak{B}(\mathbf{w}, \varphi_2, \varsigma) \leq c\tau \|\mathbf{w}\|_{\mathcal{W}_0^{1-s, r'}(\Omega)} \|\varsigma\|_{W^{s, r}(\Omega)} \leq c\tau \|\mathbf{w}\|_{W^{1-s, r'}(\Omega)} \|F\|_{W^{s, r}(\Omega)}.$$

On the other hand, Lemma 2.4 and inequality (2.21) yield $|\langle 1, b_{13}\varsigma \rangle| \leq c\|F\|_{W^{s, r}(\Omega)}$. From this and (2.28), (2.30) we finally obtain

$$(2.33) \quad \langle \psi, F \rangle \leq |n| + c\|F\|_{W^{s, r}(\Omega)} [\tau\mathfrak{G} + \|\mathfrak{d}\|_{\mathbb{W}^{-s, r'}(\Omega)} + |n|].$$

Moreover, by virtue of the duality principle

$$\|\psi\|_{\mathbb{W}^{-s, r'}(\Omega)} = \sup_{\|F\|_{W^{s, r}(\Omega)}=1} |\langle \psi, F \rangle|,$$

we have the following estimate for ψ :

$$(2.34) \quad \|\psi\|_{\mathbb{W}^{-s, r'}(\Omega)} \leq c\tau\mathfrak{G} + c\|\mathfrak{d}\|_{\mathbb{W}^{-s, r'}(\Omega)} + c|n|.$$

Let us estimate \mathbf{w} and ω . Substituting $F = \varsigma = 0$ and $M = v = 0$ into (2.24) we obtain

$$\langle \mathbf{w}, \mathbf{H} \rangle + \langle \omega, G \rangle = \mathfrak{A}_1(\mathbf{w}, \mathbf{h}) + \mathfrak{J}_2 + \mathfrak{J}_3 + R\mathfrak{J}_4 + R\mathfrak{J}_5 + n\langle 1, b_{23}g \rangle - n + \langle \mathcal{D}, \mathbf{h} \rangle.$$

By virtue of Lemma 2.2 and (2.22), the first term in the right-hand side is bounded:

$$|\mathfrak{A}_1(\mathbf{w}, \mathbf{h})| \leq c\tau^2(\|\mathbf{H}\|_{\mathcal{W}^{s-1, r}(\Omega)} + \|G\|_{W^{s, r}(\Omega)})\|\mathbf{w}\|_{\mathcal{W}_0^{1-s, r'}(\Omega)}.$$

Next we have

$$|n\langle 1, b_{23}g \rangle| \leq c(\|\mathbf{H}\|_{\mathcal{W}^{s-1, r}(\Omega)} + \|G\|_{W^{s, r}(\Omega)})|n|.$$

Obviously

$$|\langle \mathcal{D}, \mathbf{W} \rangle| \leq c\|\mathcal{D}\|_{L^1(\Omega)}\|\mathbf{h}\|_{C(\Omega)} \leq c(\|\mathbf{H}\|_{\mathcal{W}^{s-1, r}(\Omega)} + \|G\|_{W^{s, r}(\Omega)})\|\mathcal{D}\|_{L^1(\Omega)}.$$

These inequalities together with estimates (2.29)–(2.30) and inequality $R \leq \tau^2$ imply

$$\begin{aligned} & \langle \mathbf{w}, \mathbf{H} \rangle + \langle \omega, G \rangle \leq |n| + c\tau\mathfrak{Q}\mathfrak{G} \\ & + c\mathfrak{Q}\left(\|\psi\|_{\mathbb{W}^{-s, r'}(\Omega)} + |n| + \|\mathfrak{d}\|_{\mathbb{W}^{-s, r'}(\Omega)} + \|\mathcal{D}\|_{L^1(\Omega)}\right). \end{aligned}$$

Combining this result with (2.34), we obtain

$$(2.35) \quad \langle \mathbf{w}, \mathbf{H} \rangle + \langle \omega, G \rangle \leq |n| + c\tau\mathfrak{Q}\mathfrak{G} + c\mathfrak{Q}(|n| + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\mathcal{D}\|_{L^1(\Omega)}),$$

where $\mathfrak{Q} = \mathfrak{Q}(\mathbf{H}, G, 0, 0)$. For $G = 0$ and by the duality principle

$$\|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)} = \sup_{\|\mathbf{H}\|_{\mathcal{W}^{s-1,r}(\Omega)}=1} \langle \mathbf{H}, \mathbf{w} \rangle,$$

we conclude that

$$(2.36) \quad \|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)} \leq |n| + c\tau\mathfrak{G} + c(|n| + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\mathcal{D}\|_{L^1(\Omega)}).$$

Next, substituting $\mathbf{H} = \mathbf{h} = 0$, $G = g = F = \zeta = 0$ into identity (2.24), we arrive at

$$\langle \xi, M \rangle = \mathfrak{B}(\mathbf{w}, \zeta_2, v) + \mathfrak{J}_3 + \sigma\langle \mathfrak{d}, v \rangle - n.$$

Lemma 1.7 and (2.21) give the estimate for the first term:

$$|\mathfrak{B}(\mathbf{w}, \zeta_2, v)| \leq c\|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)}\|v\|_{\mathcal{W}^{s,r}(\Omega)} \leq c\|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)}\|M\|_{\mathcal{W}^{s,r}(\Omega)}.$$

From this and estimates (2.30), (2.21), we obtain

$$\langle \xi, M \rangle \leq c\tau\mathfrak{Q}\mathfrak{G} + c\mathfrak{Q}(\|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}) + |n|.$$

Combining this result with inequality (2.36), we arrive at

$$(2.37) \quad \langle \xi, M \rangle \leq c\mathfrak{Q}(\tau\mathfrak{G} + |n| + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)} + \|\mathcal{D}\|_{L^1(\Omega)}) + c|n|.$$

Finally, choosing all test functions in (2.24) equal to 0, we obtain $n = \mathfrak{J}_3$, which together with (2.30) yields

$$(2.38) \quad |n| \leq c\tau\mathfrak{Q}(\mathfrak{G} + \|\mathfrak{d}\|_{\mathbb{W}^{-s,r'}(\Omega)}).$$

From (2.33), (2.35), (2.37), combined with (2.38), it follows (2.31), and the proof of Theorem 2.3 is completed. \square

Uniqueness of solutions. The important consequence of Theorem 2.3 is the following result on uniqueness of solutions to problem (1.13).

PROPOSITION 2.6. *Under the assumptions of Theorem 1.6, there exists a positive τ_0 such that, for all $\tau \in (0, \tau_0]$, problem (1.13) admits a unique solution in the ball \mathcal{B}_τ .*

Proof. If for some \mathbf{N} the problem has two distinct solutions $(\vartheta_i, \zeta_i, m_i)$, $i = 1, 2$, with $\vartheta_i \in \mathcal{B}_\tau$, then the corresponding *finite differences* of the solutions \mathbf{w} , ψ , ω , and ξ meet all requirements of Theorem 2.3 with $\mathfrak{d} = 0$ and $\mathcal{D} = 0$. Therefore, in view of (2.26) all the elements \mathbf{w} , ψ , ω , and ξ are equal to 0, which completes the proof. \square

3. Proofs of Theorems 1.9 and 1.11.

Proof of Theorem 1.9. Let us consider a family of matrices $\mathbf{N}(\varepsilon)$ having representation (1.6) and the sequence of corresponding solutions $(\vartheta(\varepsilon), \zeta(\varepsilon), m(\varepsilon))$ to problem (1.13), where $\vartheta(\varepsilon) = (\mathbf{u}(\varepsilon), \pi(\varepsilon), \varphi(\varepsilon))$. By virtue of (1.41), we can assume that, possibly after passing to a subsequence, the sequence $(\vartheta(\varepsilon), \zeta(\varepsilon), m(\varepsilon))$ converges weakly in $(Y^{s,r})^3 \times (X^{s,r})^3 \times \mathbb{R}$ to some element (ϑ, ζ, m) , which satisfies equations (1.13) with $\mathbf{N} = \mathbf{I}$ and meets all requirements of Theorem 1.6. Since the solution (ϑ, ζ, m) to problem (1.13) is unique, the limit is independent of the choice of a subsequence,

and the whole sequence converges to the limit (ϑ, ζ, m) . It follows from (2.17) that the differences $\vartheta - \vartheta(\varepsilon)$, $\zeta - \zeta(\varepsilon)$, $m - m(\varepsilon)$ satisfy (2.17), with the coefficients $b_{ij}^{(\varepsilon)}$ and the operator \mathcal{D}_ε given by formula (2.18) with

$$\mathbf{N}_1 = \mathbf{I}, \mathbf{N}_2 = \mathbf{N}(\varepsilon), (\vartheta_1, \zeta_1, m_1) = (\vartheta, \zeta, m), (\vartheta_2, \zeta_2, m_2) = (\vartheta(\varepsilon), \zeta(\varepsilon), m(\varepsilon)).$$

In particular, the operator \mathcal{D}_ε is defined by the equality

$$\begin{aligned} \mathcal{D}_\varepsilon &= R\varrho(\varepsilon) \left(\mathbf{u}(\varepsilon) \nabla \mathbf{u}(\varepsilon) - (\mathbf{N}(\varepsilon)^*)^{-1} (\mathbf{u}(\varepsilon) \nabla (\mathbf{N}(\varepsilon)^{-1} \mathbf{u}(\varepsilon))) \right) \\ &+ (\mathbf{N}(\varepsilon)^*)^{-1} \operatorname{div} (\mathbf{g}^{-1}(\varepsilon) \mathbf{N}(\varepsilon) \mathbf{N}(\varepsilon)^* \nabla (\mathbf{N}(\varepsilon)^{-1} \mathbf{u}(\varepsilon))) - \Delta \mathbf{u}(\varepsilon) \end{aligned}$$

and admits the representation $\mathcal{D}_\varepsilon = \varepsilon \mathcal{D}_0(\mathfrak{d}, \mathbf{D}) + \varepsilon^2 \mathcal{D}_1(\varepsilon)$, where \mathcal{D}_0 is given by (1.45). Moreover, since the norms $\|\varrho(\varepsilon)\|_{C(\Omega)}$ and $\|\mathbf{u}(\varepsilon)\|_{W^{2,2}(\Omega)}$ are uniformly bounded, we have

$$(3.1) \quad \|\mathcal{D}_1(\varepsilon)\|_{L^2(\Omega)} \leq c(\mathbf{U}, \Omega, \sigma).$$

Next note that $\mathbf{g}(\varepsilon)$ admits the decomposition $\mathbf{g}(\varepsilon) = 1 + \varepsilon \mathfrak{d} + \varepsilon^2 \mathfrak{d}_1(\varepsilon)$, where $\mathfrak{d} = \operatorname{Tr} \mathbf{D}$ and the remainder $\mathfrak{d}_1(\varepsilon)$ is uniformly bounded in $C^1(\Omega)$. Proceeding as in the previous section and recalling the equalities $\mathcal{A}_1 = \mathfrak{A}_1 = 0$, we conclude that the finite differences

$$(\mathbf{w}_\varepsilon, \omega_\varepsilon, \psi_\varepsilon) = \varepsilon^{-1}(\vartheta - \vartheta(\varepsilon)), \quad \xi_\varepsilon = \varepsilon^{-1}(\zeta - \zeta(\varepsilon)), \quad n_\varepsilon = \varepsilon^{-1}(m - m(\varepsilon))$$

satisfy the integral identity

$$\begin{aligned} (3.2) \quad & \int_{\Omega} \mathbf{w}_\varepsilon \left(\mathbf{H} - \varrho(\varepsilon) \nabla \mathbf{u} \cdot \mathbf{h} + R\varrho(\varepsilon) \nabla \mathbf{h}^* \mathbf{u}(\varepsilon) \right) dx \\ & - \mathfrak{B}(\mathbf{w}_\varepsilon, \varphi(\varepsilon), \varsigma) - \mathfrak{B}(\mathbf{w}_\varepsilon, v, \zeta(\varepsilon)) + \langle \omega_\varepsilon, G - b_{12}^{(\varepsilon)} \varsigma - b_{22}^{(\varepsilon)} g - \varkappa b_{32}^{(\varepsilon)} \rangle \\ & + \langle \psi_\varepsilon, F - b_{11}^{(\varepsilon)} \varsigma - b_{21}^{(\varepsilon)} g - \varkappa b_{31}^{(\varepsilon)} - R\mathbf{u} \cdot \nabla \mathbf{u} \cdot \mathbf{h} \rangle + \langle \xi_\varepsilon, M - \varkappa b_{34}^{(\varepsilon)} \rangle + n_\varepsilon \\ & - n_\varepsilon \langle 1, b_{13}^{(\varepsilon)} \varsigma + b_{23}^{(\varepsilon)} g \rangle = \langle \mathfrak{d} + \varepsilon \mathfrak{d}_1(\varepsilon), b_{10}^{(\varepsilon)} \varsigma + b_{20}^{(\varepsilon)} g + \varkappa b_{30}^{(\varepsilon)} + \sigma v \rangle + \langle \mathcal{D}_0 + \varepsilon \mathcal{D}_1(\varepsilon), \mathbf{h} \rangle, \end{aligned}$$

along with the orthogonality conditions $\langle \omega_\varepsilon, 1 \rangle = 0$. Here (\mathbf{H}, G, F, M) are arbitrary smooth functions such that $G = \Pi G$, and the test functions ς, v, g , and \mathbf{h} are defined by equations (2.19), (2.20); therefore, the test functions are independent of ε . Recall that the elements $\vartheta_1 = \vartheta$ and $\vartheta_2 = \vartheta(\varepsilon)$ belong to the ball \mathcal{B}_τ and meet all requirements of Theorem 2.3. Hence, there exists $\tau_1 > 0$, depending only on Ω, \mathbf{U} and s, r, σ , such that the conditions

$$\lambda^{-1}, \quad R \leq \tau^2, \quad \|\mathbf{N}(\varepsilon) - \mathbf{I}\|_{C^1(\Omega)} \leq \tau^2, \quad 0 < \tau \leq \tau_1,$$

imply

$$\begin{aligned} \|\mathbf{w}_\varepsilon\|_{\mathcal{W}_0^{1-s, r'}(\Omega)} + \|\omega_\varepsilon\|_{\mathbb{W}^{-s, r'}(\Omega)} + \|\psi_\varepsilon\|_{\mathbb{W}^{-s, r'}(\Omega)} + \|\xi_\varepsilon\|_{\mathbb{W}^{-s, r'}(\Omega)} + |n_\varepsilon| \\ \leq c(\|\mathcal{D}_0 + \varepsilon \mathcal{D}_1\|_{L^1(\Omega)} + \|\mathfrak{d} + \varepsilon \mathfrak{d}_1\|_{C^2(\Omega)}) \leq c. \end{aligned}$$

Therefore, after possibly passing to a subsequence, we can assume that the sequence \mathbf{w}_ε converges to \mathbf{w} weakly in $\mathcal{W}_0^{1-s, r'}(\Omega)$, and $(\omega_\varepsilon, \psi_\varepsilon, \xi_\varepsilon)$ converge to (ω, ψ, ξ) (*)-weakly in $\mathbb{W}^{-s, r'}(\Omega)$ as $\varepsilon \rightarrow 0$.

Next, choose $s' > s$ satisfying conditions (1.38). By virtue of Theorem 1.6, there exists $\tau'_0 > 0$ (depending only on $\Omega, \mathbf{U}, r, s, \sigma$) such that, for all $\tau \in (0, \tau'_0]$, the functions $(\vartheta(\varepsilon), \zeta(\varepsilon))$ are bounded in $(Y^{s',r})^3 \times (X^{s',r})^3$. It follows from this that the family $\mathbf{u}(\varepsilon) = \mathbf{u}_0 + \mathbf{u}(\varepsilon)$ converges to \mathbf{u} strongly in $Y^{s,r}$, and $(\varphi(\varepsilon), \pi(\varepsilon), \zeta(\varepsilon))$ converges to (φ, π, ζ) strongly in $X^{s,r}$. Therefore, by virtue of Lemma 2.1, the sequence $b_{ij}^{(\varepsilon)}$ converges strongly in $X^{s,r}$ to b_{ij}^0 . Hence, we can pass to the limit in (3.2). It is easy to see that the limits, the vector field \mathbf{w} , and the functionals ψ, ω, ξ are given by a unique weak solution to problem (1.42) and, in addition, meet all requirements of Definition 1.8. It remains to note that, by virtue of Theorem 2.3, the limit is independent of the choice of a subsequence, which completes the proof of Theorem 1.9. \square

Proof of Theorem 1.11. Assume that r, s, σ , and τ satisfy inequalities (1.31), (1.38), and that $\vartheta = (\mathbf{u}, \pi, \varphi) \in \mathcal{B}_\tau$ is a solution to problem (1.13) given by Theorem 1.6. Denote by $Y_0^{s,r}$ the subspace of the space $Y^{s,r}$ of all functions vanishing on $\partial\Omega$, by $X_{\text{in}}^{s,r}$ and $X_{\text{out}}^{s,r}$ the subspaces of $X^{s,r}$ which consist of all functions vanishing on Σ_{in} and Σ_{out} , respectively, and by $X_{\Pi}^{s,r}$ the subspace of all function in $X^{s,r}$ having the zero mean value. Introduce the Banach spaces $\mathcal{E} = (Y_0^{s,r})^3 \times X_{\Pi}^{s,r} \times X_{\text{out}}^{s,r} \times X_{\text{in}}^{s,r} \times \mathbb{R}$ and $\mathcal{F} = (Z^{s,r})^3 \times X_{\Pi}^{s,r} \times (X^{s,r})^2 \times \mathbb{R}$. Our first task is to show that, for all $\Theta \in \mathcal{F}$, problem (1.52), (1.49) has a unique solution $Y \in \mathcal{E}$. We begin with the observation that, by virtue of Lemma 1.3, the Stokes operator has the bounded inverse

$$\begin{pmatrix} \Delta & -\nabla \\ \text{div} & 0 \end{pmatrix}^{-1} : (Z^{s,r})^3 \times X_{\Pi}^{s,r} \rightarrow (Y_0^{s,r})^3 \times X_{\Pi}^{s,r}.$$

On the other hand, by virtue of Corollary 1.5, the operators \mathcal{L} and \mathcal{L}^* from (1.48) have the bounded inverses $\mathcal{L}^{-1} : X^{s,r} \rightarrow X_{\text{in}}^{s,r}$, $(\mathcal{L}^*)^{-1} : X^{s,r} \rightarrow X_{\text{out}}^{s,r}$. Therefore, there exists the bounded operator

$$\begin{pmatrix} \mathcal{L}^* & 0 & 0 \\ 0 & \mathcal{L} & 0 \\ -\mathbb{B}_{1,3} & 0 & 1 \end{pmatrix}^{-1} : (X^{s,r})^2 \times \mathbb{R} \rightarrow X_{\text{out}}^{s,r} \times X_{\text{in}}^{s,r} \times \mathbb{R}.$$

It follows from this that, for all $\Theta \in \mathcal{F}$, the equation $\mathfrak{L}\mathbf{Y} = \Theta$ has a unique solution satisfying boundary conditions (1.49) and the inequality $\|\mathbf{Y}\|_{\mathcal{E}} \leq c\|\Theta\|_{\mathcal{F}}$, where the constant c is independent of τ . Let us consider the operators \mathfrak{U} . By virtue of Lemma 1.7, we have

$$\|\varsigma \nabla \varphi\|_{W^{s-1,r}(\Omega)} \leq c\tau \|\varsigma\|_{X^{s,r}}.$$

It is easy to see that

$$\|\varsigma \nabla \varphi\|_{L^2(\Omega)} \leq c\|\varsigma\|_{X^{s,r}} \|\nabla \varphi\|_{L^2(\Omega)} \leq c\|\varsigma\|_{X^{s,r}} \|\varphi\|_{X^{s,r}} \leq c\tau \|\varsigma\|_{X^{s,r}}.$$

Combining the obtained estimates we get the inequality $\|\varsigma \nabla \varphi\|_{Z^{s,r}} \leq c\tau \|\varsigma\|_{X^{s,r}}$. Repetition of these arguments gives the inequality $\|\zeta \nabla v\|_{Z^{s,r}} \leq c\|v\|_{X^{s,r}}$. Since the norms $\|b_{ij}^0\|_{X^{s,r}}$ are uniformly bounded, we conclude from this that $\|\mathfrak{U}\mathbf{Y}\|_{\mathcal{F}} \leq c\|Y\|_{\mathcal{E}}$. Finally, let us consider the operator \mathfrak{V} . Since the space $X^{s,r}$ is the commutative Banach algebra and $\nabla \mathbf{u} \in X^{s,r}$, we have

$$\|\mathbf{R}\mathbf{u}\nabla \mathbf{u}\mathbf{h}\|_{X^{s,r}} \leq c\tau^2 \|\mathbf{h}\|_{X^{s,r}}, \quad \|\mathbf{R}\varrho(\nabla \mathbf{u} + \mathbf{u}\nabla)\mathbf{h}\|_{X^{s,r}} \leq c\tau^2 \|\mathbf{h}\|_{Y^{s,r}}.$$

On the other hand, by virtue of Lemma 2.1 and (1.43), the coefficients b_{ij}^0 in the expression for \mathfrak{V} satisfy the inequalities

$$\|b_{12}^0\|_{X^{s,r}} + \|b_{11}^0\|_{X^{s,r}} + \|b_{31}^0\|_{X^{s,r}} + \|b_{32}^0\|_{X^{s,r}} + \|b_{34}^0\|_{X^{s,r}} \leq c\tau,$$

which yield the estimate $\|\mathfrak{Y}\mathbf{Y}\|_{\mathcal{F}} \leq c\tau\|Y\|_{\mathcal{E}}$. Thus we get that the diagonal matrix operator \mathfrak{L} has the bounded inverse, \mathfrak{U} is the bounded upper triangular (with respect to \mathfrak{L}) matrix operator, and \mathfrak{V} is the small bounded operator. Hence, for all sufficiently small τ the operator $\mathfrak{L} - \mathfrak{U} - \mathfrak{V} : \mathcal{E} \rightarrow \mathcal{F}$ has the bounded inverse, which implies the existence of an adjoint state satisfying (1.52) and boundary conditions (1.49).

It remains to prove identity (1.53). Fix the adjoint state $\mathbf{Y} = (\mathbf{h}, g, \varsigma, v, l)$, and set

$$\mathbf{H} = \Delta\mathbf{h} - \nabla g, \quad G = \operatorname{div} \mathbf{h}, \quad F = \mathcal{L}^*\varsigma, \quad M = \mathcal{L}v.$$

It follows from (1.52) that

$$\begin{aligned} \mathbf{H} - R\rho(\nabla\mathbf{u} - \mathbf{u}\nabla)\mathbf{h} + \varsigma\nabla\varphi + \zeta\nabla v &= \Delta\eta\mathbf{U}_\infty + R\rho((\mathbf{u}\nabla\eta)\mathbf{U}_\infty + (\mathbf{u}\mathbf{U}_\infty)\nabla\eta), \\ (3.3) \quad G - \Pi(b_{21}^0\varsigma + b_{22}^0g - \varkappa b_{32}^0l) &= \Pi(\nabla\eta\mathbf{U}_\infty), \\ F - R\mathbf{u}\nabla\mathbf{u}\mathbf{h} - b_{12}^0g - b_{11}^0\varsigma - \varkappa b_{31}^0l &= (\mathbf{u}\nabla\eta)(\mathbf{u}\mathbf{U}_\infty), \quad M = \varkappa b_{34}^0l. \end{aligned}$$

By virtue of Theorem 1.9, the *material* derivative $(\mathbf{w}, \omega, \psi, \xi, n)$ satisfies the integral identities (1.47). On the other hand, (F, \mathbf{H}, G, M) together with the components of the adjoint state \mathbf{Y} can be regarded as a collection of test functions for this identity. Substituting these test functions into (1.47), using equalities (3.3), and recalling the identity $\langle \omega, 1 \rangle = 0$, we obtain

$$\begin{aligned} L_u(\mathbf{w}, \omega, \psi) + \varkappa(l - 1) &\left(\langle \psi, b_{31}^0 \rangle + \langle \omega, b_{32}^0 \rangle + \langle \xi, b_{34}^0 \rangle \right) \\ + n - n\langle 1, b_{13}^0\varsigma \rangle &= \langle \mathfrak{d}, b_{10}^0\varsigma + b_{20}^0g + \varkappa b_{30}^0 + \sigma v \rangle + \langle \mathcal{D}_0, \mathbf{W} \rangle. \end{aligned}$$

It follows from (1.52) that $l = \langle 1, b_{13}^0\varsigma \rangle$, which leads to

$$\begin{aligned} (3.4) \quad L_u(\mathbf{w}, \omega, \psi) + \varkappa(l - 1) &\left(\varkappa(\langle \psi, b_{31}^0 \rangle + \langle \omega, b_{32}^0 \rangle + \langle \xi, b_{34}^0 \rangle) - n \right) \\ &= \langle \mathfrak{d}, b_{10}^0\varsigma + b_{20}^0g + \varkappa b_{30}^0 + \sigma v \rangle + \langle \mathcal{D}_0, \mathbf{W} \rangle. \end{aligned}$$

Next note that the identities (1.42) imply the following expression for the constant n :

$$n = \varkappa(\langle \psi, b_{31}^0 \rangle + \langle \omega, b_{32}^0 \rangle + \langle \xi, b_{34}^0 \rangle + \langle \mathfrak{d}, b_{30}^0 \rangle).$$

Substituting this equality into (3.4) and noting that $\mathfrak{d} = \operatorname{Tr} \mathbf{D}$, we obtain (1.53), which completes the proof. \square

4. Proof of Theorem 1.4. Our strategy is the following. First we use the classical method of characteristics (see [8] for instance) to show that in the vicinity of each point $P \in \Sigma_{\text{in}} \cup \Gamma$ there exist normal coordinates (y_1, y_2, y_3) such that $\mathbf{u}\nabla_x = \mathbf{e}_1\nabla_{y_1}$. Hence the problem of existence of solutions to the transport equation in the neighborhood of $\Sigma_{\text{in}} \cap \Gamma$ is reduced to a boundary value problem for the model equation $\partial_{y_1}\varphi + \sigma\varphi = f$ in a parabolic domain. Next we prove that the boundary value problem for the model equation admits a unique solution in fractional Sobolev space, which leads to the existence and uniqueness of solutions in the neighborhood of the inlet set. Using the existence of a local solution we can reduce problem (1.28) to a boundary value problem for the modified equation, which does not require any boundary data. Application of well-known results on solvability of elliptic-hyperbolic equations in the case $\Gamma = \emptyset$ finally gives the existence and uniqueness of solutions to problems (1.28) and (1.29).

First we introduce some notation which will be used throughout this section. For any $a > 0$ we denote by Q_a the cube $[-a, a]^3$ and by Q_a^+ the slab $[-a, a]^2 \times [0, a]$ in the space of points $y = (y_1, y_2, y_3) \in \mathbb{R}^3$. We will write Y instead of (y_2, y_3) so that $y = (y_1, Y)$.

DEFINITION 4.1. A standard parabolic neighborhood associated with the constant c_0 is a compact subset of the slab Q_a^+ , defined by the inequalities

$$(4.1) \quad \mathcal{P}_a = \{y = (y_1, Y) \in Q_a^+ : a^-(Y) \leq y_1 \leq a^+(Y)\},$$

where $a^\pm : [-a, a] \times [0, a] \mapsto \mathbb{R}$ are continuous, piecewise C^1 -functions satisfying the inequalities

$$(4.2) \quad \begin{aligned} -a &\leq a^-(Y) \leq 0 \leq a^+(Y) \leq a, \\ -c_0\sqrt{y_3} &\leq a^-(Y) \leq a^+(Y) \leq c_0\sqrt{y_3}, \\ |\partial_{y_2} a^\pm(Y)| &\leq c_0, \quad |\partial_{y_3} a^\pm(Y)| \leq c_0/\sqrt{y_3}. \end{aligned}$$

Set $Q_{\text{in}} = \{Y : a^-(Y) > -a\}$ and $Q_{\text{out}} = \{Y : a^+(Y) < a\}$. Denote by Σ_{in}^y and Σ_{out}^y the surfaces determined by the relations

$$\begin{aligned} \Sigma_{\text{in}}^y &= \{y : Y \in Q_{\text{in}}, \quad y_1 = a^-(Y)\}, \\ \Sigma_{\text{out}}^y &= \{y : Y \in Q_{\text{out}}, \quad y_1 = a^+(Y)\}. \end{aligned}$$

It is clear that $\partial\mathcal{P}_a = (\partial Q_a \cap \partial\mathcal{P}_a) \cup \Sigma_{\text{in}}^y \cup \Sigma_{\text{out}}^y$.

LEMMA 4.2. Assume that the C^2 -manifold $\Sigma = \partial B$ and the vector field $\mathbf{U} \in C^2(\Sigma)^3$ satisfy conditions **(H1)**–**(H3)**. Let $\mathbf{u} \in C^1(\mathbb{R}^3)^3$ be a compactly supported vector field such that

$$\mathbf{u} = \mathbf{U} \text{ on } \Sigma, \quad \mathbf{u} = 0 \text{ on } S.$$

Denote $M = \|\mathbf{u}\|_{C^1(\mathbb{R}^3)}$. Then there are positive constants a, c, C, ρ_c , and R_c , depending only on $M, \partial\Omega$, and \mathbf{U} , with the following properties:

(P1) For any point $P \in \Gamma$ there exists a mapping $y \rightarrow \mathbf{x}(y)$ which takes diffeomorphically the cube Q_a onto a neighborhood \mathcal{O}_P of P and satisfies the equation

$$(4.3) \quad \partial_{y_1} \mathbf{x}(y) = \mathbf{u}(\mathbf{x}(y)) \text{ in } Q_a$$

and the inequalities

$$(4.4) \quad \|\mathbf{x}\|_{C^1(Q_a)} + \|\mathbf{x}^{-1}\|_{C^1(\mathcal{O}_P)} \leq C, \quad |\mathbf{x}(y)| \leq C|y|.$$

(P2) There is a standard parabolic neighborhood \mathcal{P}_a associated with the constant c such that

$$(4.5) \quad \mathbf{x}(\mathcal{P}_a) = \mathcal{O}_P \cup \Omega, \quad \mathbf{x}(\Sigma_{\text{in}}^y) = \Sigma_{\text{in}} \cap \mathcal{O}_P, \quad \mathbf{x}(\Sigma_{\text{out}}^y) = \Sigma_{\text{out}} \cap \mathcal{O}_P.$$

(P3) Denote by $G_a \subset \mathcal{P}_a$ the domain

$$(4.6) \quad G_a = \{y = (y_1, Y) \in \mathcal{P}_a : Y \in Q_{\text{in}}\},$$

and by $B_P(\rho)$ the ball $|x - P| \leq \rho$. Then we have the inclusions

$$(4.7) \quad B_P(\rho_c) \cap \Omega \subset \mathbf{x}(G_a) \subset \mathcal{O}_P \cap \Omega \subset B_P(R_c) \cap \Omega.$$

Proof. We start with the proof of **(P1)**. Recall condition **(H1)** and fix the standard Cartesian coordinate system (x_1, x_2, x_3) associated with the point $P \in \Gamma$. Let us consider the Cauchy problem.

$$(4.8) \quad \begin{aligned} &\partial_{y_1} \mathbf{x} = \mathbf{u}(\mathbf{x}(y)) \quad \text{in } Q_a, \\ &x_1(y) = \Upsilon(y_2), \quad x_2(y) = y_2 \quad \text{for } y_1 = 0, \\ &x_3 = F(\Upsilon(y_2), y_2) + y_3 \quad \text{for } y_1 = 0. \end{aligned}$$

Without any loss of generality we can assume that $0 < a < k < 1$. For any such a , problem (4.8) has a unique solution of class $C^1(Q_a)$. Denote by $\mathfrak{F}(y) = D_y \mathbf{x}(y)$ the Jacobian matrix function. The calculations show that

$$\mathfrak{F}_0 := \mathfrak{F}(y) \Big|_{y_1=0} = \begin{pmatrix} u_1 & \Upsilon'(y_2) & 0 \\ u_2 & 1 & 0 \\ u_3 & \partial_{y_2} F(\Upsilon(y_2), y_2) & 1 \end{pmatrix}, \quad \mathfrak{F}(0) = \begin{pmatrix} U & \Upsilon'(0) & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix},$$

which implies

$$(4.9) \quad \|\mathfrak{F}(0)^{\pm 1}\| \leq C/3, \quad \|\mathfrak{F}_0(y) - \mathfrak{F}(0)\| \leq ca,$$

where the constants C, c are independent of a .

Differentiation of (4.8) leads to the ordinary differential equation for \mathfrak{F}

$$\partial_{y_1} \mathfrak{F} = D_y \mathbf{u}(\mathbf{x}) \mathfrak{F}, \quad \mathfrak{F} \Big|_{y_1=0} = \mathfrak{F}_0.$$

From this we get

$$\partial_{y_1} \|\mathfrak{F} - \mathfrak{F}_0\| \leq M(\|\mathfrak{F} - \mathfrak{F}_0\| + \|\mathfrak{F}_0\|),$$

and hence $\|\mathfrak{F} - \mathfrak{F}_0\| \leq c(M)\|\mathfrak{F}_0\|a$. Combining this result with (4.9) we finally arrive at

$$(4.10) \quad \|\mathfrak{F}(y) - \mathfrak{F}(0)\| \leq ca.$$

This inequality along with the implicit function theorem implies the existence of $a > 0$, depending only on M and Ω , such that the mapping $x = \mathbf{x}(y)$ takes diffeomorphically the cube Q_a onto some neighborhood of the point P and satisfies inequalities (4.4).

Let us turn to the proof of **(P2)**. We begin with the observation that the manifold $\mathbf{x}^{-1}(\partial\Omega \cap \mathcal{O})$ is defined by the equation

$$\Phi_0(y) := x_3(y) - F(x_1(y), x_2(y)) = 0, \quad y \in Q_a.$$

Let us show that Φ_0 is strictly monotone in y_3 and has the opposite signs on the faces $y_3 = \pm a$. To this end note that the formula for $\mathfrak{F}(0)$ along with (4.10) implies the estimates

$$|\partial_{y_3} x_3(y) - 1| + |\partial_{y_3} x_1(y)| + |\partial_{y_3} x_2(y)| \leq ca \quad \text{in } Q_a.$$

Thus, we get

$$1 - ca \leq \partial_{y_3} \Phi_0(y) = \partial_{y_3} x_3(y) - \partial_{x_i} F(x_1, x_2) \partial_{y_3} x_i(y) \leq 1 + ca.$$

It follows from (4.10) that, for $y_3 = 0$, we have $|x_3(y)| \leq ca|y|$, which along with (4.4) yields the estimate

$$|\Phi_0(y)| \leq |x_3(y)| + |F(x(y))| \leq ca|y| + KC|y|^2 \leq ca^2 \quad \text{for } y_3 = 0.$$

Hence there is a positive a , depending only on M and Ω , such that the inequalities

$$1/2 \leq \partial_{y_3} \Phi_0(y) \leq 2, \quad \pm \Phi_0(y_1, y_2, \pm a) > 0$$

hold true for all $y \in Q_a$. Therefore, the equation $\Phi_0(y) = 0$ has a unique solution $y_3 = \Phi(y_1, y_2)$ in the cube Q_a . Moreover, the function $\Phi \in C^1([-a, a]^2)$ vanishes for $y_1 = y_3 = 0$. Thus, we get

$$\mathcal{P}_a := \mathbf{x}^{-1}(\mathcal{O} \cap \Omega) = \{\Phi(y_1, y_2) < y_3 < a, \quad |y_1|, |y_2| \leq a\}.$$

Note that $|\mathbf{u}(\mathbf{x}(y)) - U\mathbf{e}_1| \leq M|\mathbf{x}(y)| \leq Ca$. Therefore, we can choose $a = a(M, \Omega)$ such that $2U/3 \leq u_1 \leq 4U/3$ and $C|u_2| \leq U/3$ in Q_a . Recall that $x_1(y) - \Upsilon(x_2(y))$ vanishes at the plane $y_1 = 0$ and

$$\partial_{y_1} [x_1(y) - \mathfrak{g}(x_2(y))] = u_1(y) - \Upsilon'(x_2(y))u_2(y).$$

We obtain from this that, for a suitable choice of a ,

$$(4.11) \quad |y_1|U/3 \leq |x_1(y) - \Upsilon(x_2(y))| \leq |y_1|5U/3 \quad \text{for } y \in Q_a.$$

Equations (4.8) imply the identity

$$\partial_{y_1} \Phi_0(y) \equiv \nabla F_0(\mathbf{x}(y)) \cdot \mathbf{u}(x(y)) = \nabla F_0(\mathbf{x}(y)) \cdot \mathbf{U}(x(y)) \quad \text{for } \Phi_0(y) = 0.$$

Combining this result with (1.17) and (4.11), we obtain the estimates

$$|y_1|N^-U/3 \leq |\partial_{y_1} \Phi_0(y)| \leq |y_1|N^+U5/3,$$

which along with the identity

$$\partial_{y_1} \Phi = -\partial_{y_1} \Phi_0 (\partial_{y_3} \Phi_0)^{-1}$$

yield the inequalities

$$(4.12) \quad \begin{aligned} -c < \partial_{y_1} \Phi(y_1, y_2) &\leq cy_1 \quad \text{for } -a < y_1 < 0, \\ cy_1 < \partial_{y_1} \Phi(y_1, y_2) &\leq c \quad \text{for } 0 < y_1 < a, \\ |\partial_{y_2} \Phi(y_1, y_2)| &\leq c, 0 \leq \Phi(y_1, y_2) \leq cy_1^2. \end{aligned}$$

It is clear that for sufficiently small a , depending only on \mathbf{U} and Ω , the functions $\Phi^\pm(y_2) = \Phi(\pm a, y_2)$ admit the estimates $ca^2 \leq \Phi^\pm(y_2) < a$. Set

$$\begin{aligned} Q_{\text{in}} &= \{Y \in [-a, a] \times [0, a] : 0 < y_3 < \Phi^-(y_2)\}, \\ Q_{\text{out}} &= \{Y \in [-a, a] \times [0, a] : 0 < y_3 < \Phi^+(y_2)\}. \end{aligned}$$

It follows from (4.12) that for every $Y \in Q_{\text{in}}$ ($Y \in Q_{\text{out}}$) the equation $y_3 = \Phi(y_1, y_2)$ has a unique solution $a^-(Y) < 0$ ($a^+(Y) > 0$). We adopt the convention that $a^\pm(Y) = \pm a$ for $y_3 > \Phi^\pm(y_2)$. It remains to note that, by virtue of (4.12), the functions a^\pm meet all requirements of Lemma 4.2. \square

The next lemma shows the existence of the normal coordinates in the vicinity of points of the inlet Σ_{in} .

LEMMA 4.3. *Let vector fields \mathbf{u} and \mathbf{U} meet all requirements of Lemma 4.2 and $U_n = -\mathbf{U}(P) \cdot \mathbf{n} > N > 0$. Then there is $b > 0$, depending only on N , Σ , and $M = \|\mathbf{u}\|_{C^1(\Omega)}$, with the following properties. There exists a mapping $y \rightarrow \mathbf{x}(y)$,*

which takes diffeomorphically the cube $Q_b = [-b, b]^3$ onto a neighborhood \mathcal{O}_P of P and satisfies the equations

$$(4.13) \quad \partial_{y_3} \mathbf{x}(y) = \mathbf{u}(\mathbf{x}(y)) \text{ in } Q_b, \quad \mathbf{x}(y_1, y_2, 0) \in \Sigma \cap \mathcal{O}_P \text{ for } |y_2| \leq a,$$

and the inequalities

$$(4.14) \quad \|\mathbf{x}\|_{C^1(Q_b)} + \|\mathbf{x}^{-1}\|_{C^1(\mathcal{O}_P)} \leq C_{M,N}, \quad |\mathbf{x}(y)| \leq C_M |y|,$$

where $C_{M,N} = 3(1 + N^{-1})(M^2 + 2)^{1/2}$. The inclusions

$$(4.15) \quad B_P(\rho_i) \cap \Omega \subset \mathbf{x}(Q_b \cap \{y_3 > 0\}) \subset B_P(R_i) \cap \Omega$$

hold true for $\rho_i = C_{M,N}^{-1}b$ and $R_i = C_{M,N}b$.

Proof. The proof simulates the proof of Lemma 4.2. Choose the local Cartesian coordinates (x_1, x_2, x_3) centered at P such that in new coordinates $\mathbf{n}(P) = \mathbf{e}_3$. By the smoothness of Σ , there is a neighborhood $\mathcal{O} = [-k, k]^2 \times [-t, t]$ such that the manifold $\Sigma \cap \mathcal{O}$ is defined by the equation

$$x_3 = F(x_1, x_2), \quad F(0, 0) = 0, \quad |\nabla F(x_1, x_2)| \leq K(|x_1| + |x_2|).$$

The constants k, t , and K depend only on Σ . Let us consider the initial value problem

$$(4.16) \quad \partial_{y_3} \mathbf{x} = \mathbf{u}(\mathbf{x}(y)) \text{ in } Q_a, \quad \mathbf{x}\Big|_{y_3=0} = (y_1, y_2, F(y_1, y_2)).$$

Without any loss of generality we can assume that $0 < b < k < 1$. It follows from **(H1)** that for any such b problem (4.16) has a unique solution of class $C^1(Q_b)$. Next note that for $y_3 = 0$ we have

$$(4.17) \quad |\mathbf{x}(y)| \leq (K + 1)|y|, \quad |\mathbf{u}(\mathbf{x}(y)) - \mathbf{u}(0)| \leq M(K + 1)|y|.$$

Denote by $\mathfrak{F}(y) = D_y \mathbf{x}(y)$. The calculations show that

$$\mathfrak{F}_0 := \mathfrak{F}(y)\Big|_{y_3=0} = \begin{pmatrix} 1 & 0 & u_1 \\ 0 & 1 & u_2 \\ 0 & 0 & u_3 \end{pmatrix}, \quad \mathfrak{F}(0) = \begin{pmatrix} 1 & 0 & u_1(P) \\ 0 & 1 & u_2(P) \\ 0 & 0 & U_n \end{pmatrix},$$

which along with (4.17) implies

$$(4.18) \quad \|\mathfrak{F}(0)^{\pm 1}\| \leq C_{M,N}/3, \quad \|\mathfrak{F}_0(y) - \mathfrak{F}(0)\| \leq cb.$$

Next, differentiation of (4.16) with respect to y leads to the equation

$$\partial_{y_1} \mathfrak{F} = D_y \mathbf{u}(\mathbf{x}) \mathfrak{F}, \quad \mathfrak{F}\Big|_{y_3=0} = \mathfrak{F}_0.$$

Arguing as in the proof of Lemma 4.2 we obtain $\|\mathfrak{F} - \mathfrak{F}_0\| \leq c(M)\|\mathfrak{F}_0\|b$. Combining this result with (4.18), we finally arrive at $\|\mathfrak{F}(y) - \mathfrak{F}(0)\| \leq cb$. From this and the implicit function theorem we conclude that there is positive b , depending only on M and Σ , such that the mapping $x = \mathbf{x}(y)$ takes diffeomorphically the cube Q_b onto some neighborhood of the point P and satisfies inequalities (4.14). Inclusions (4.15) easily follow from (4.14). \square

Model equation. Let us consider the following boundary value problem in standard parabolic neighborhood \mathcal{P}_a :

$$(4.19) \quad \partial_{y_1} \varphi(y) + \sigma \varphi(y) = f(y) \text{ in } \mathcal{P}_a, \quad \varphi(y) = 0 \text{ for } y_1 = a^-(y_2, y_3).$$

LEMMA 4.4. *Assume that*

$$(4.20) \quad 1/2 < s \leq 1 \text{ and } 1 < r < 3/(2s - 1).$$

Then for any $f \in W^{s,r}(Q_a) \cap L^\infty(Q_a^\phi)$, problem (4.19) has a unique solution satisfying the inequalities

$$(4.21) \quad \begin{aligned} \|\varphi\|_{W^{s,r}(\mathcal{P}_a)} &\leq c(r, s) \left(a^{4/r-s} \|f\|_{L^\infty(\mathcal{P}_a)} + a^{1/r} \|f\|_{W^{s,r}(\mathcal{P}_a)} \right), \\ \|\varphi\|_{L^\infty(Q_a^\phi)} &\leq \sigma^{-1} \|f\|_{L^\infty(Q_a^\phi)}. \end{aligned}$$

Proof. It suffices to prove the lemma for $s < 1$. For every $y, z \in \mathbb{R}^3$, we denote by $Y = (y_2, y_3)$, $Z = (z_2, z_3)$, respectively. Obviously, we have

$$(4.22) \quad \varphi(y) = \int_{a^-(Y)}^{y_1} e^{\sigma(x_1-y_1)} f(x_1, Y) dx_1 \text{ and } \sigma \|\varphi\|_{C(\mathcal{P}_a)} \leq \|f\|_{C(\mathcal{P}_a)}.$$

Therefore, it suffices to estimate the seminorm $|\varphi|_{s,r,\mathcal{P}_a}$. Choose an arbitrary $y, z \in \mathcal{P}_a$. Without any loss of generality we can assume that $a^-(Z) \leq a^-(Y)$. The identity

$$\begin{aligned} \varphi(z) - \varphi(y) &= \varphi(z_1, Z) - \varphi(y_1, Z) + \int_{a^-(Z)}^{a^-(Y)} e^{\sigma(x_1-y_1)} f(x_1, Z) dx_1 \\ &\quad + \int_{a^-(Y)}^{y_1} e^{\sigma(x_1-y_1)} (f(x_1, Z) - f(x_1, Y)) dx_1 \end{aligned}$$

implies the estimate

$$|\varphi(z) - \varphi(y)| \leq \|f\|_{L^\infty(\mathcal{P}_a)} (2|y_1 - z_1| + |a^-(Y) - a^-(Z)|) + \int_{-a}^a |(f(x_1, Z) - f(x_1, Y))| dx_1,$$

which along with the inequality

$$\left(\int_{-a}^a |(f(x_1, Z) - f(x_1, Y))| dx_1 \right)^r \leq a^{r-1} \int_{-a}^a |(f(x_1, Z) - f(x_1, Y))|^r dx_1$$

leads to the estimate

$$(4.23) \quad |\varphi|_{s,r,\mathcal{P}_a}^r \leq 2 \|f\|_{C(\mathcal{P}_a)}^r (I_1 + I_2) + a^{(r-1)} I_3.$$

Here we denote

$$\begin{aligned} I_1 &= \int_{Q_a \times Q_a} \frac{|y_1 - z_1|^r}{|x - y|^{3+rs}} dx dy, \quad I_2 = \int_{\mathcal{P}_a \times \mathcal{P}_a} \frac{|a^-(Y) - a^-(Z)|^r}{|x - y|^{3+rs}} dx dy, \\ I_3 &= \int_{[-a,a]^4} \int_{-a}^a \frac{|f(x_1, Y) - f(x_1, Z)|^r}{|x - y|^{3+rs}} dx dy dx_1. \end{aligned}$$

Let us estimate the terms $I_j, j = 1, 2, 3$. We begin with the observation that

$$\begin{aligned} \int_{[-a,a]^2} \frac{dZ}{|x-y|^{3+rs}} &= \frac{1}{|y_1-z_1|^{3+rs}} \int_{[-a,a]^2} \frac{dZ}{(|Y-Z|^2/|y_1-z_1|^2+1)^{(3+rs)/2}} \\ &\leq \frac{1}{|y_1-z_1|^{1+rs}} \int_{\mathbb{R}^2} \frac{dZ}{(|Z|^2+1)^{(3+rs)/2}} \leq \frac{c}{|y_1-z_1|^{1+rs}}, \end{aligned}$$

and hence

$$\int_{[-a,a]^4} \frac{dY dZ}{|x-y|^{3+rs}} \leq \frac{ca^2}{|y_1-z_1|^{1+rs}}.$$

From this we obtain

$$(4.24) \quad I_1 \leq ca^2 \int_{-a}^a \left(\int_{-a}^a |y_1-z_1|^{r(1-s)-1} dz_1 \right) dy_1 \leq c(r,s)a^{3+r(1-s)}.$$

In order to estimate I_2 , note that, by Lemma 4.2,

$$|a^-(Y) - a^-(Z)| \leq c|Y - Z|^{1/2}, \quad |a^-(Y) - a^-(Z)| \leq c|Y - Z| \left(\frac{1}{\sqrt{y_3}} + \frac{1}{\sqrt{z_3}} \right).$$

Next, it follows from the assumptions of the lemma that there is $\lambda \in (0, 1)$ such that

$$(4.25) \quad \lambda < 3/r, \quad 0 < (1 + \lambda)/2 - s < 1/r.$$

Noting that

$$|a^-(Y) - a^-(Z)| \leq c|Y - Z|^{(1+\lambda)/2} ((y_3)^{-\lambda/2} + (z_3)^{-\lambda/2}),$$

we obtain

$$\begin{aligned} I_2 &\leq c \int_{\mathcal{P}_a} (y_3)^{-r\lambda/2} \left(\int_{Q_a} \frac{|Y - Z|^{r(1+\lambda)/2}}{|x-y|^{3+rs}} dz \right) dy \\ &\quad + c \int_{\mathcal{P}_a} (z_3)^{-r\lambda/2} \left(\int_{Q_a} \frac{|Y - Z|^{r(1+\lambda)/2}}{|x-y|^{3+rs}} dy \right) dz \\ &\leq 2c \int_{\mathcal{P}_a} (y_3)^{-r\lambda/2} \left(\int_{Q_a} \frac{|Y - Z|^{r(1+\lambda)/2}}{|x-y|^{3+rs}} dz \right) dy. \end{aligned}$$

Next, inequalities (4.25) imply

$$\begin{aligned} &\int_{Q_a} \frac{|Y - Z|^{r(1+\lambda)/2}}{|x-y|^{3+rs}} dz \\ &\leq \int_{-a}^a |y_1-z_1|^{-3-rs} \left(\int_{\mathbb{R}^2} \frac{|Y - Z|^{r(1+\lambda)/2}}{(|Y - Z|^2|y_1-z_1|^{-2}+1)^{(3+rs)/2}} dZ \right) \\ &\leq \int_{-a}^a |y_1-z_1|^{-rs-1+r(1+\lambda)/2} dz_1 \int_{\mathbb{R}^2} \frac{|Z|^{r(1+\lambda)/2}}{(1+|Z|^2)^{(3+rs)/2}} dZ \\ &\leq c(r,s) \int_{-a}^a |y_1-z_1|^{r(1+\lambda)/2-rs-1} dz_1 \leq c(r,s). \end{aligned}$$

From this and Lemma 4.2 we conclude that

$$\begin{aligned}
 I_2 &\leq c(r, s) \int_{\mathcal{P}_a} (y_3)^{-r\lambda/2} dy \\
 (4.26) \quad &\leq c(r, s) \int_{[-a, a]^2} \left(\int_{cy_1^2 \leq y_3 \leq a} (y_3)^{-r\lambda/2} dy_3 \right) dy_1 dy_2 \\
 &\leq ac(r, s) \int_{-a}^a |y_1|^{2-r\lambda} dy_1 \leq ac(r, s).
 \end{aligned}$$

The remaining part of the proof is based on the following proposition.

PROPOSITION 4.5. *Let $f \in W^{s,r}(Q_a)$ and $sr > 1$. Then f has an extension \bar{f} onto \mathbb{R}^3 , which vanishes outside the set Q_{3a} and satisfies*

$$(4.27) \quad \|\bar{f}\|_{W^{s,r}} \leq ca^{(3-rs)/r} \|f\|_{L^\infty(Q_a)} + |f|_{s,r,Q_a}.$$

Proof. Define an extension of f onto the slab $[-3a, 3a] \times [-a, a]^2$ by the formula

$$f(x^\pm) = f(x) \text{ for } x \in Q_a, \text{ where } x^\pm = (\pm(2a - x_1), x_2, x_3).$$

It easily follows from the definition of the seminorm $|\cdot|_{r,s,\Omega}$ that

$$\|f\|_{W^{s,r}([-3a,3a] \times [-a,a]^2)} \leq 3\|f\|_{L^r(Q_a)} + 6|f|_{s,r,Q_a} \leq 6\|f\|_{W^{s,r}(Q_a)}.$$

Proceeding in the same way as before, first we can extend f onto the plate $[-3a, 3a]^2 \times [-a, a]$ and then over the cube Q_{3a} . Obviously, the extended function, still denoted by f , satisfies the inequalities

$$(4.28) \quad \|f\|_{W^{s,r}(Q_{3a})} \leq 216\|f\|_{W^{s,r}(Q_a)}, \quad \|f\|_{C(Q_{3a})} \leq \|f\|_{C(Q_a)}.$$

Next choose $\mu \in C^\infty(\mathbb{R}^3)$ such that $0 \leq \mu \leq 1$, $\mu = 1$ in Q_1 , and $\mu = 0$ outside of Q_{22} . Set $\bar{f} = f\mu_a$, where $\mu_a(x) = \mu(x/a)$. Next, the interpolation inequality along with the estimate $|\nabla\mu_a| \leq ca^{-1}$ implies

$$\|\mu_a\|_{W^{s,r}(\mathbb{R}^3)} \leq \|\mu_a\|_{L^r(\mathbb{R}^3)}^{1-s} \|\mu_a\|_{W^{1,r}(\mathbb{R}^3)}^s \leq ca^{3(1-s)/r} a^{(3-r)s/r} = ca^{(3-rs)/r}.$$

From this and the obvious inequality $\|\mu_a f\|_{W^{s,r}(\mathbb{R}^3)} \leq \|f\|_{L^\infty(Q_{3a})} \|\mu_a\|_{W^{s,r}(\mathbb{R}^3)} + \|f\|_{W^{s,r}(Q_{3a})}$ we conclude that

$$\|\mu_a f\|_{W^{s,r}(\mathbb{R}^3)} \leq ca^{(3-rs)/r} \|f\|_{L^\infty(Q_a)} + \|f\|_{W^{s,r}(Q_a)}.$$

Hence $\bar{f} = \mu_a f$ satisfies (4.27), and the proposition follows. \square

Let us return to the proof of Lemma 4.4. We have

$$\begin{aligned}
 I_3 &= \int_{[-a,a]^4} \int_{-a}^a |Z - Y|^{-3-rs} (|z_1 - y_1|^2 |Z - Y|^{-2} + 1)^{-(3+rs)/2} \\
 &\quad \times |f(x_1, Z) - f(x_1, Y)|^r dx dy dx_1 \leq ca \int_{\mathbb{R}} (|t|^2 + 1)^{-(3+rs)/2} dt \\
 &\quad \times \int_{[-a,a]^5} |Z - Y|^{-2-rs} |f(x_1, Z) - f(x_1, Y)|^r dX dY dx_1 \\
 &= ca \int_{-a}^a |f(x_1, \cdot)|_{r,s,[-a,a]^2}^r dx_1,
 \end{aligned}$$

which yields

$$(4.29) \quad I_3 \leq ca\|f\|_{L^r(-a,a;W^{s,r}([-a,a]^2))}^r \leq ca\|\bar{f}\|_{L^r(\mathbb{R};W^{s,r}(\mathbb{R}^2))}^r.$$

Recall that, for $s = 0, 1$, the embedding operator $W^{s,r}(\mathbb{R}^3) \hookrightarrow L^r(\mathbb{R};W^{s,r}(\mathbb{R}^2))$ is bounded. By virtue of Lemma B.1, this results holds true for all $s \in [0, 1]$, which along with Proposition 4.5 and inequality (4.29) implies

$$I_3 \leq ca^{4-rs}\|f\|_{L^\infty(Q_a)}^r + a(|f|_{s,r,Q_a})^r.$$

Combining this result with (4.24), (4.26), since $3 + r(1 - s) \geq 4 - rs$, we finally obtain

$$\|f\|_{L^\infty(Q_a)}^r(I_1 + I_2) + I_3 \leq ca^{4-sr}\|f\|_{L^\infty(Q_a)}^r + ca|f|_{s,r,Q_a}^r.$$

Substituting this inequality into (4.23) gives (4.21), and the lemma follows. \square

Let us consider the following boundary value problem in the slab $Q_a^+ = [-a, a]^2 \times [0, a]$:

$$(4.30) \quad \partial_{y_3}\varphi(y) + \sigma\varphi(y) = f(y) \text{ in } Q_a^+, \quad \varphi(y) = 0 \text{ for } y_3 = 0.$$

LEMMA 4.6. *Problem (4.30) has a unique solution satisfying the inequality*

$$(4.31) \quad \|\varphi\|_{W^{s,r}(Q_a^+)} \leq c(r, s)(a^{4/r-s}\|f\|_{L^\infty(Q_a^+)} + a^{1/r}\|f\|_{W^{s,r}(Q_a^+)}).$$

Proof. The proof of Lemma 4.4 can also be used in this case. \square

Local existence results. It follows from the assumptions of Theorem 1.4 that the vector field \mathbf{u} and the manifold Σ satisfy all assumptions of Lemma 4.2. Therefore, there exist positive numbers a, ρ_c , and R_c , depending only on Σ and $\|\mathbf{u}\|_{C^1(\Omega)}$, such that for all $P \in \Gamma$, the canonical diffeomorphism $\mathbf{x} : Q_a \mapsto \mathcal{O}_P$ is well defined and meets all requirements of Lemma 4.2. Fix an arbitrary point $P \in \Gamma$ and consider the boundary value problem

$$(4.32) \quad \mathbf{u} \cdot \nabla\varphi + \sigma\varphi = f \text{ in } \mathcal{O}_P, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}} \cap \mathcal{O}_P.$$

LEMMA 4.7. *Suppose that the exponents s, r satisfy condition (1.31). Then, for any $f \in C^1(\Omega)$, problem (4.32) has a unique solution satisfying the inequalities*

$$(4.33) \quad |\varphi|_{s,r,B_P(\rho_c)} \leq c(\|f\|_{C(B_P(R_c))} + |f|_{s,r,B_P(R_c)}), \quad \|\varphi\|_{C(B_P(\rho_c))} \leq \sigma^{-1}\|f\|_{C(B_P(R_c))},$$

where the constant c depends only on Σ, M, σ, s, r and the constant ρ_c is determined by Lemma 4.2.

Proof. We transform (4.33) using the normal coordinates (y_1, y_2, y_3) given by Lemma 4.2. Set $\bar{\varphi}(y) = \varphi(\mathbf{x}(y))$ and $\bar{f}(y) = f(\mathbf{x}(y))$. Next note that (4.3) implies the identity $\mathbf{u} \cdot \nabla_x \varphi = \partial_{y_1} \bar{\varphi}(y)$. Therefore the function $\bar{\varphi}(y)$ satisfies the following equation along with the boundary conditions:

$$(4.34) \quad \partial_{y_1} \bar{\varphi} + \sigma \bar{\varphi} = \bar{f} \text{ in } Q_a \cap \{y_3 > \Phi\}, \quad \bar{\varphi} = 0 \text{ for } y_3 = \Phi(y_1, y_2), y_1 < 0.$$

It follows from Lemma 4.4 that problem (4.34) has a unique solution $\bar{\varphi} \in W^{s,r}(G_a)$ satisfying the inequality

$$(4.35) \quad |\bar{\varphi}|_{s,r,G_a} \leq c(\|\bar{f}\|_{C(Q_a)} + |\bar{f}|_{s,r,Q_a}), \quad \|\bar{\varphi}\|_{C(G_a)} \leq \sigma^{-1}\|f\|_{C(Q_a)},$$

where the domain G_α is defined by (4.6). It remains to note that, by estimate (4.4), the mappings $\mathbf{x}^{\pm 1}$ are uniformly Lipschitz, which along with inclusions (4.7) implies the estimates

$$|\varphi|_{s,r,B_P(\rho_c)} \leq c|\bar{\varphi}|_{s,r,G_\alpha}, \quad |\bar{f}|_{s,r,Q_\alpha} \leq c|f|_{s,r,B_P(R_c)}.$$

Combining these results with (4.35) we finally obtain (4.33) and the lemma follows. \square

In order to formulate the similar result for interior points of the inlet Σ_{in} we introduce the set

$$(4.36) \quad \Sigma'_{\text{in}} = \{x \in \Sigma_{\text{in}} : \text{dist}(x, \Gamma) \geq \rho_c/3\},$$

where the constant ρ_c is given by Lemma 4.2. It is clear that

$$\inf_{P \in \Sigma'_{\text{in}}} \mathbf{U}(P) \cdot \mathbf{n}(P) \geq N > 0,$$

where the constant N depends only on M and Σ . It follows from Lemma 4.3 that there are positive numbers b , ρ_i , and R_i such that for each $P \in \Sigma'_{\text{in}}$, the canonical diffeomorphism $\mathbf{x} : Q_b \mapsto \mathcal{O}_P$ is well defined and satisfies the hypotheses of Lemma 4.3. The following lemma gives the local existence and uniqueness of solutions to the boundary value problem

$$(4.37) \quad \mathbf{u} \cdot \nabla \varphi + \sigma \varphi = f \text{ in } \mathcal{O}_P, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}} \cap \mathcal{O}_P.$$

LEMMA 4.8. *Suppose that the exponents s, r satisfy condition (1.38). Then, for any $f \in C^1(\Omega)$ and $P \in \Sigma'_{\text{in}}$, problem (4.32) has a unique solution satisfying the inequalities*

$$(4.38) \quad \begin{aligned} |\varphi|_{s,r,B_P(\rho_i)} &\leq c(\|f\|_{C(B_P(R_i))} + |f|_{s,r,B_P(R_i)}), \\ \|\varphi\|_{C(B_P(R))} &\leq \sigma^{-1} \|f\|_{C(B_P(R_i))}, \end{aligned}$$

where c depends on Σ , M , σ , and exponents s, r .

Proof. Using the normal coordinates given by Lemma 4.3 we rewrite (4.37) in the form

$$\partial_{y_3} \bar{\varphi} + \sigma \bar{\varphi} = \bar{f} \text{ in } Q_b, \quad \bar{\varphi} = 0 \text{ for } y_3 = 0.$$

Applying Lemma 4.4 and arguing as in the proof of Lemma 4.7 we obtain (4.38). \square

Existence of solutions near inlet. The next step is based on the well-known geometric lemma (see [21, Chap. 3]).

LEMMA 4.9. *Suppose that a given set $A \subset \mathbb{R}^d$ is covered by balls such that each point $x \in A$ is the center of a certain ball $B_x(r(x))$ of radius $r(x)$. If $\sup r(x) < \infty$, then from the system of the balls $\{B_x(r(x))\}$ it is possible to select a countable system $B_{x_k}(r(x_k))$ covering the entire set A and having multiplicity not greater than a certain number $\mathbf{n}(d)$ depending only on the dimension d .*

The following lemma gives the dependence of the multiplicity of radii of the covering balls.

LEMMA 4.10. *Assume that a collection of balls $B_{x_k}(r) \subset \mathbb{R}^3$ of constant radius r has the multiplicity \mathbf{n}_r . Then the multiplicity of the collections of the balls $B_{x_k}(R)$, $r < R$, is bounded by the constant $27(R/r)^3 \mathbf{n}_r$.*

Proof. Let \mathbf{n}_R be a multiplicity of the system $\{B_{x_k}(R)\}$. This means that at least \mathbf{n}_R balls, say $B_{x_1}(R), \dots, B_{x_{\mathbf{n}_R}}(R)$, have the common point P . In particular, we

have $B_{x_i}(r) \subset B_P(3R)$ for all $i \leq n_R$. Introduce the counting function $\iota(x)$ for the collection of balls $B_{x_i}(r)$, defined by

$$\iota(x) = \text{card}\{i : x \in B_{x_i}(r), 1 \leq i \leq n_r\}.$$

Note that $\iota(x) \leq n_r$. We have

$$\frac{4\pi}{3} n_R r^3 = \sum_{i=1}^{n_R} \text{meas } B_{x_i}(r) = \int_{\cup_i B_{x_i}(r)} \iota(x) dx \leq n_r \int_{\cup_i B_{x_i}(r)} dx \leq \frac{4\pi}{3} (3R)^3 n_r,$$

and the lemma follows. \square

We are now in a position to prove the local existence and uniqueness of a solution to the first boundary value problem for the transport equation in the neighborhood of the inlet Σ_{in} . Let Ω_t be the t -neighborhood of the set Σ_{in} ,

$$\Omega_t = \{x \in \Omega : \text{dist}(x, \Sigma_{\text{in}}) < t\}.$$

LEMMA 4.11. *Let $t = \min\{\rho_c/2, \rho_i/2\}$ and $T = \max\{R_c, R_i\}$, where the constants ρ_α, R_α are defined by Lemmas 4.2 and 4.3. Then there exists a constant C , depending only on M, Σ , and σ , such that, for any $f \in C^1(\Omega)$, the boundary value problem*

$$(4.39) \quad \mathbf{u} \cdot \nabla \varphi + \sigma \varphi = f \text{ in } \Omega_t, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}}$$

has a unique solution satisfying the inequalities

$$(4.40) \quad |\varphi|_{s,r,\Omega_t} \leq C(\|f\|_{C(\Omega_T)} + |f|_{s,r,\Omega_T}), \quad \|\varphi\|_{C(\Omega_t)} \leq \sigma^{-1} \|f\|_{C(\Omega_T)}.$$

Proof. It follows from Lemma 4.9 that there is a covering of the characteristic manifold Γ by the finite collection of balls $B_{P_i}(\rho_c/4)$, $1 \leq i \leq \mathbf{m}$, $P_i \in \Gamma$, of the multiplicity \mathbf{n} . The cardinality \mathbf{m} of this collection does not exceed $4\mathbf{n}(\rho_c)^{-1}L$, where L is the length of Γ . Obviously, the balls $B_{P_i}(\rho_c)$ cover the set

$$\mathcal{V}_\Gamma = \{x \in \Omega : \text{dist}(x, \Gamma) < \rho_c/2\}.$$

By virtue of Lemma 4.7, in each such ball the solution to problem (4.39) satisfies inequalities (4.33), which leads to the estimate

$$(4.41) \quad |\varphi|_{s,r,\mathcal{V}_\Gamma}^r \leq \sum_i |\varphi|_{s,r,B_{P_i}(\rho_c)}^r \leq c \sum_i \|f\|_{C(B_{P_i}(R_c))}^r + c \sum_i |f|_{s,r,B_{P_i}(R_c)}^r,$$

where c depends only on M, Σ , and σ . By Lemma 4.10, the multiplicity of the system of balls $B_{P_i}(R_c)$ is bounded from above by $12^3(R_c/\rho_c)^3$, which along with the inclusion $\cup_i B_{P_i}(R_c) \subset \Omega_T$ yields

$$\sum_{i=1}^{\mathbf{m}} |f|_{s,r,B_{P_i}(R_c)}^r \leq 12^3(R_c/\rho_c)^3 |f|_{s,r,\Omega_T}^r.$$

Obviously we have

$$\sum_i \|f\|_{C(B_{P_i}(R_c))}^r \leq \mathbf{m} \|f\|_{C(\Omega_T)}^r \leq 4\mathbf{n}(\rho_c)^{-1}L \|f\|_{C(\Omega_T)}^r.$$

Combining these results with (4.41) we obtain the estimates for the solution to problem (4.39) in the neighborhood of the characteristic manifold Γ ,

$$(4.42) \quad |\varphi|_{s,r,\mathcal{V}_\Gamma} \leq c\|f\|_{C(\Omega_T)} + c|f|_{s,r,\Omega_T}.$$

Our next task is to obtain the similar estimate in the neighborhood of the compact $\Sigma'_{\text{in}} \subset \Sigma_{\text{in}}$. To this end, we introduce the set

$$\mathcal{V}_{\text{in}} = \{x \in \Omega : \text{dist}(x, \Sigma'_{\text{in}}) < \rho_i/2\},$$

where Σ'_{in} is given by (4.36). By virtue of Lemma 4.9, there exists the finite collection of balls $B_{P_k}(\rho_i/4)$, $1 \leq k \leq \mathbf{m}$, $P_k \in \Sigma'_{\text{in}}$, of the multiplicity \mathbf{n} which covers Σ'_{in} . Obviously $\mathbf{m} \leq 16n(\rho_i)^{-2} \text{meas } \Sigma_{\text{in}}$, and the balls $B_{P_k}(\rho_i)$ cover the set \mathcal{V}_{in} . From this and Lemma 4.8 we conclude that

$$|\varphi|_{s,r,\mathcal{V}_{\text{in}}}^r \leq \sum_k |\varphi|_{s,r,B_{P_k}(\rho_i)}^r \leq c \sum_k \|f\|_{C(B_{P_k}(R_i))}^r + c \sum_k |f|_{s,r,B_{P_k}(R_i)}^r.$$

By virtue of Lemma 4.10, the multiplicity of the system of balls $B_{P_i}(R_i)$ is not greater than $12^3(R_i/\rho_i)^3$, which yields

$$\sum_i |f|_{s,r,B_{P_i}(R_i)}^r \leq 12^3(R_i/\rho_i)^3 |f|_{s,r,\Omega_T}^r.$$

Obviously we have

$$\sum_k \|f\|_{C(B_{P_k}(R_i))}^r \leq \mathbf{m} \|f\|_{C(\Omega_T)}^r \leq 16\mathbf{n}(\rho_k)^{-2} \text{meas } \Sigma_{\text{in}} \|f\|_{C(\Omega_T)}^r.$$

Thus we get

$$(4.43) \quad |\varphi|_{s,r,\mathcal{V}_{\text{in}}} \leq c\|f\|_{C(\Omega_T)} + c|f|_{s,r,\Omega_T}.$$

Since \mathcal{V}_Γ and \mathcal{V}_{in} cover Ω_t , this inequality along with inequality (4.42) yields (4.40), and the lemma follows. \square

Partition of unity. Let us turn to the analysis of the general problem

$$(4.44) \quad \mathcal{L}\varphi := \mathbf{u} \cdot \nabla \varphi + \sigma \varphi = f \text{ in } \Omega, \quad \varphi = 0 \text{ on } \Sigma_{\text{in}}.$$

The next step is based on the theory of partial differential equations with nonnegative characteristic form. The following lemma is a particular case of general results of Oleinik and Radkevich; we refer to Theorems 1.5.1 and 1.6.2 in [34].

LEMMA 4.12. *Assume that Ω is a bounded domain of the class C^2 , the vector field \mathbf{u} belongs to the class $C^1(\Omega)^3$, and $\sigma - \text{div } \mathbf{u}(x) > \delta > 0$. Then, for any $f \in L^\infty(\Omega)$, problem (1.28) has a unique solution such that $\|\varphi\|_{L^\infty(\Omega)} \leq \delta^{-1} \|f\|_{L^\infty}$. Moreover, this solution is continuous at the interior points of Σ_{in} and vanishes on Σ_{in} . If, in addition, $\Gamma = \text{cl}(\Sigma_{\text{out}} \cap \Sigma_0) \cap \text{cl } \Sigma_{\text{in}}$ is a smooth one-dimensional manifold, then a bounded generalized solution to problem (4.44) is unique.*

The question of smoothness of solutions to boundary value problems for transport equations is more complicated. All known results [19], [34] are related to the case of $\Gamma = \emptyset$. The following lemma is a consequence of Theorem 1.8.1 in the monograph [34].

LEMMA 4.13. *Assume that Ω is a bounded domain of the class C^2 and $\Sigma_{\text{out}} = \emptyset$. Furthermore, let the following conditions hold.*

- (1) The vector field \mathbf{u} and the function f belong to the class $C^1(\mathbb{R}^3)$.
- (2) There is $\Omega' \ni \Omega$ such that the inequality

$$\sigma - \sup_{\Omega'} \left\{ |\operatorname{div} \mathbf{u}| - \frac{1}{2} \sup_i \sum_{j \neq i} \left| \frac{\partial u_i}{\partial x_j} \right| - \frac{1}{2} \sup_j \sum_{i \neq j} \left| \frac{\partial u_j}{\partial x_i} \right| \right\} > 0$$

is fulfilled. Then a weak solution to problem (1.28) satisfies the Lipschitz condition in $\bar{\Omega}$.

Using these results we can construct a strong solution to problem (1.28). Recall that by Lemma 4.11, for any $f \in C^1(\Omega)$, problem (4.44) has a unique strong solution defined in neighborhood Ω_t of the inlet Σ_{in} . On the other hand, Lemma 4.12 guarantees the existence and uniqueness of a bounded weak solution to problem (4.44). The following lemma shows that both solutions coincide in Ω_t .

LEMMA 4.14. *Under the assumptions of Theorem 1.4 and Lemma 4.11, each bounded generalized solution to problem (4.44) coincides in Ω_t with the local solution φ_t .*

Proof. Let $\varphi \in L^\infty(\Omega)$ be a weak solution to problem (4.44). Recall that each point $P \in \Gamma$ has a canonical neighborhood $\mathcal{O}_P := \mathbf{x}(Q_a)$, where canonical diffeomorphism $\mathbf{x} : Q_a \mapsto \mathcal{O}_P$ is defined by Lemma 4.2. Choose an arbitrary function $\zeta \in C^1(\Omega)$ vanishing on Σ_{in} and outside of \mathcal{O}_P and set

$$\bar{\varphi}(y) = \varphi(\mathbf{x}(y)), \quad \bar{f}(y) = f(\mathbf{x}(y)), \quad \bar{\zeta}(y) = \zeta(\mathbf{x}(y)), \quad y \in Q_a \cap \{y_3 > \Phi\}.$$

By the definition of a weak solution to the transport equation we have

$$\int_{\mathcal{O}_P \cap \Omega} (\sigma \varphi \zeta - \varphi \operatorname{div}(\zeta \mathbf{u}) - f \zeta) dx = 0.$$

Direct calculations lead to the identity $\operatorname{div}_x(\zeta \mathbf{u}) = \det \mathfrak{F}^{-1} \operatorname{div}_y(\bar{\zeta} \det \mathfrak{F} \mathfrak{F}^{-1} \bar{\mathbf{u}})$, in which the notation \mathfrak{F} stands for the Jacobi matrix $\mathfrak{F} = D_y \mathbf{x}(y)$. On the other hand, (4.3) implies the equality $\mathfrak{F}^{-1} \bar{\mathbf{u}} = \mathbf{e}_1$. From this we conclude that

$$\int_{Q_a \cap \{y_3 > \Phi\}} \left((\det \mathfrak{F} \bar{\zeta})(\sigma \bar{\varphi} - \bar{f}) - \bar{\varphi} \frac{\partial}{\partial y_1} (\det \mathfrak{F} \bar{\zeta}) \right) dy = 0.$$

Recall that, by Lemma 4.2, $\partial_{y_1} \mathfrak{F}$ is continuous and $\det \mathfrak{F}$ is strictly positive in the cube Q_a . Setting $\xi = \det \mathfrak{F} \bar{\zeta}$ we conclude that the integral identity

$$\int_{Q_a \cap \{y_3 > \Phi\}} \left(\xi (\sigma \bar{\varphi} - \bar{f}) - \bar{\varphi} \frac{\partial \xi}{\partial y_1} \right) dy = 0$$

holds true for all functions $\xi \in C_0(Q_a)$ having continuous derivative $\partial_{y_1} \xi \in C(Q_a)$ and vanishing for $y_3 = \Phi(y_1, y_2)$, $y_1 < 0$. Since \bar{f} is continuously differentiable, $\bar{\varphi}$ belongs to the class $C^1_{loc}(Q_a \cap \{y_3 > \Phi\})$ and satisfies (4.34). On the other hand, $\bar{\varphi}_t$ also satisfies (4.34). Obviously, all solutions to problem (4.34) coincide in the domain G_a , and hence $\bar{\varphi}_t = \bar{\varphi}$ in this domain. Recalling that $B_P(\rho_c) \subset \mathbf{x}(G_a)$ we obtain that $\varphi_t = \varphi$ in the ball $B_P(\rho_c)$. The same arguments show that, for any $P \in \Sigma'_{\text{in}}$, the function φ_t is equal to φ in the ball $B_P(\rho_i)$. It remains to note that the balls $B_P(\rho_c)$ and $B_P(\rho_i)$ cover Ω_t , and the lemma follows. \square

Furthermore, we split the weak solution $\varphi \in L^\infty(\Omega)$ to problem (4.44) into two parts, namely, the local solution φ_t and the remainder vanishing near the inlet. To this end fix a function $\Lambda \in C^\infty(\mathbb{R})$ such that

$$(4.45) \quad 0 \leq \Lambda' \leq 3, \quad \Lambda(u) = 0 \quad \text{for } u \leq 1 \quad \text{and} \quad \Lambda(u) = 1 \quad \text{for } u \geq 3/2,$$

and introduce the one-parametric family of smooth functions

$$(4.46) \quad \chi_t(x) = \frac{1}{t^3} \int_{\mathbb{R}^3} \Theta \left(\frac{2(x-y)}{t} \right) \Lambda \left(\frac{\text{dist}(y, \Sigma_{\text{in}})}{t} \right) dy,$$

where $\Theta \in C^\infty(\mathbb{R}^3)$ is a standard mollifying kernel supported in the unit ball. It follows that

$$(4.47) \quad \begin{aligned} \chi_t(x) &= 0 \quad \text{for } \text{dist}(x, \Sigma_{\text{in}}) \leq t/2, \quad \chi_t(x) = 1 \quad \text{for } \text{dist}(x, \Sigma_{\text{in}}) \geq 2t, \\ |\partial^l \chi_t(x)| &\leq \varpi(l)t^{-l} \quad \text{for all } l \geq 0, \end{aligned}$$

where $\varpi(l)$ is a constant. Now fix a number $t = t(\Sigma, M)$ satisfying all assumptions of Lemma 4.11 and set

$$(4.48) \quad \varphi(x) = (1 - \chi_{t/2}(x))\varphi_t(x) + \phi(x).$$

By virtue of (4.47) and Lemma 4.14, the function $\phi \in L^\infty(\Omega)$ vanishes in $\Omega_{t/2}$ and in a weak sense satisfies the equations

$$\mathbf{u}\nabla\phi + \sigma\phi = \chi_{t/2}f + \varphi_t\mathbf{u}\nabla\chi_{t/2} =: F \quad \text{in } \Omega, \quad \phi = 0 \quad \text{on } \Sigma_{\text{in}}.$$

Next introduce a new vector field $\tilde{\mathbf{u}}(x) = \chi_{t/8}(x)\mathbf{u}(x)$. It is easy to see that $\chi_{t/8} = 1$ on the support of ϕ , and hence the function ϕ is also a weak solution to the modified transport equation

$$(4.49) \quad \tilde{\mathcal{L}}\phi := \tilde{\mathbf{u}}\nabla\phi + \sigma\phi = F \quad \text{in } \Omega.$$

The advantage of such an approach is that the topology of integral lines of the modified vector field $\tilde{\mathbf{u}}$ drastically differs from the topology of integral lines of \mathbf{u} . The corresponding inlet, outgoing set, and characteristic set have the other structure and $\tilde{\Sigma}_{\text{in}} = \emptyset$. In particular, (4.49) does not require boundary conditions. Finally note that the C^1 -norm of the modified vector fields has the majorant

$$(4.50) \quad \|\tilde{\mathbf{u}}\|_{C^1(\Omega)} \leq M(1 + 16\varpi(1)t^{-1}),$$

where $\varpi(1)$ is a constant from (4.47). The following lemma constitutes the existence and uniqueness of solutions to the modified equation.

LEMMA 4.15. *Suppose that*

$$(4.51) \quad \sigma > \sigma^*(M, \Sigma) = 4M(1 + 16\varpi(1)t^{-1}) + 1, \quad M = \|\mathbf{u}\|_{C^1(\Omega)},$$

and $0 \leq s \leq 1, r > 1$. Then, for any $F \in W^{s,r}(\Omega) \cap L^\infty(\Omega)$, (4.49) has a unique weak solution $\phi \in W^{s,r}(\Omega) \cap L^\infty(\Omega)$ such that

$$(4.52) \quad \|\phi\|_{W^{s,r}(\Omega)} \leq c\|F\|_{W^{s,r}(\Omega)}, \quad \|\phi\|_{L^\infty(\Omega)} \leq \sigma^{-1}\|F\|_{L^\infty(\Omega)},$$

where c depends only on r .

Proof. Without any loss of generality we can assume that $F \in C^1(\Omega)$. By virtue of (4.50) and (4.51), the vector field $\tilde{\mathbf{u}}$ and σ meet all requirements of Lemma 4.13. Hence (4.49) has a unique solution $\phi \in W^{1,\infty}(\Omega)$. For $i = 1, 2, 3$ and $\tau > 0$, define the finite difference operator

$$\delta_{i\tau}\phi = \frac{1}{\tau}(\phi(x + \tau\mathbf{e}_i) - \phi(x)).$$

It is easy to see that

$$(4.53) \quad \tilde{\mathbf{u}}\nabla\delta_{i\tau}\phi + \sigma\delta_{i\tau}\phi = \delta_{i\tau}F - \delta_{i\tau}\tilde{\mathbf{u}}\nabla\phi(x + \tau\mathbf{e}_i) \quad \text{in } \Omega \cap (\Omega - \tau\mathbf{e}_i).$$

Next introduce the function $\eta \in C^\infty(\mathbb{R})$ such that $\eta' \geq 0$, $\eta(u) = 0$ for $u \leq 1$, and $\eta(u) = 1$ for $u \geq 1$, and set $\eta_h(x) = \eta(\text{dist}(x, \partial\Omega)/h)$. Since $\tilde{\Sigma}_{\text{in}} = \emptyset$, the inequality

$$(4.54) \quad \limsup_{h \rightarrow 0} \int_{\Omega} g\tilde{\mathbf{u}} \cdot \nabla\eta_h(x) \, dx \leq 0$$

holds true for all nonnegative functions $g \in L^\infty(\Omega)$. Choosing $h > \tau$, multiplying both sides of (4.53) by $\eta_h|\delta_{i\tau}\phi|^{r-2}\delta_{i\tau}\phi$, and integrating the result over $\Omega \cap (\Omega - \tau\mathbf{e}_i)$, we obtain

$$\begin{aligned} \int_{\Omega \cap (\Omega - \tau\mathbf{e}_i)} \eta_h|\delta_{i\tau}\phi|^r \left(\sigma - \frac{1}{r} \text{div } \tilde{\mathbf{u}} \right) \, dx - \int_{\Omega \cap (\Omega - \tau\mathbf{e}_i)} |\delta_{i\tau}\phi|^r \tilde{\mathbf{u}}\nabla\eta_h \, dx \\ = \int_{\Omega \cap (\Omega - \tau\mathbf{e}_i)} (\delta_{i\tau}F - \delta_{i\tau}\tilde{\mathbf{u}}\nabla\phi(x + \tau\mathbf{e}_i))\eta_h|\delta_{i\tau}\phi|^{r-2}\delta_{i\tau}\phi \, dx. \end{aligned}$$

Letting $\tau \rightarrow 0$ and then $h \rightarrow 0$ and using inequality (4.54), we obtain

$$(4.55) \quad \int_{\Omega} |\partial_{x_i}\phi|^r \left(\sigma - \frac{1}{r} \text{div } \tilde{\mathbf{u}} \right) \, dx \leq \int_{\Omega} (\partial_{x_i}F - \partial_{x_i}\tilde{\mathbf{u}}\nabla\phi)|\partial_{x_i}\phi|^{r-2}\partial_{x_i}\phi \, dx.$$

Next note that

$$\sum_i \partial_{x_i}\tilde{\mathbf{u}}\nabla\phi|\partial_{x_i}\phi|^{r-2}\partial_{x_i}\phi \leq 3\|\tilde{\mathbf{u}}\|_{C^1(\Omega)} \sum_i |\partial_{x_i}\phi|^r.$$

On the other hand, since $1/r + 3 \leq 4$, inequalities (4.50) and (4.51) imply

$$\sigma - \left(\frac{1}{r} + 3 \right) \|\tilde{\mathbf{u}}\|_{C^1(\Omega)} \geq \sigma - 4M(1 + 16\varpi(1)t^{-1}) \geq 1.$$

From this we conclude that

$$\begin{aligned} \sum_i \int_{\Omega} |\partial_{x_i}\phi|^r \, dx &\leq \sum_i \int_{\Omega} |\partial_{x_i}\phi|^{r-1} |\partial_{x_i}F| \, dx \\ &\leq \left(\sum_i \int_{\Omega} |\partial_{x_i}\phi|^{r/(r-1)} \, dx \right)^{(r-1)/r} \left(\sum_i \int_{\Omega} |\partial_{x_i}F|^r \, dx \right)^{1/r}, \end{aligned}$$

which leads to the estimate

$$(4.56) \quad \|\nabla\phi\|_{L^r(\Omega)} \leq c(r)\|\nabla F\|_{L^r(\Omega)}.$$

Next, multiplying both sides of (4.49) by $|\phi|^{r-2}\eta_h$ and integrating the result over Ω , we get the identity

$$\int_{\Omega} \left(\sigma - \frac{1}{r} \operatorname{div} \tilde{\mathbf{u}} \right) \eta_h |\phi|^r dx - \int_{\Omega} |\phi|^r \tilde{\mathbf{u}} \nabla \eta_h dx = \int_{\Omega} F \eta_h |\phi|^{r-2} \phi dx.$$

The passage $h \rightarrow 0$ gives the inequality

$$\int_{\Omega} \left(\sigma - \frac{1}{r} \operatorname{div} \tilde{\mathbf{u}} \right) |\phi|^r dx \leq \int_{\Omega} |F| |\phi|^{r-1} dx.$$

Recalling that $\sigma - 1/r \operatorname{div} \tilde{\mathbf{u}} \geq 1$, we finally obtain

$$(4.57) \quad \|\phi\|_{L^r(\Omega)} \leq c(r) \|F\|_{L^r(\Omega)}.$$

Inequalities (4.56) and (4.57) imply estimate (4.52) for $s = 0, 1$. Hence the linear operator $\tilde{\mathcal{L}}^{-1} : F \mapsto \phi$ is continuous in the Banach spaces $W^{0,r}(\Omega)$ and $W^{1,r}(\Omega)$, and its norm does not exceed $c(r)$. Recall that $W^{s,r}(\Omega)$ is the interpolation space $[L^r(\Omega), W^{1,r}(\Omega)]_{s,r}$. From this and Lemma B.1 we conclude that inequality (4.52) is fulfilled for all $s \in [0, 1]$, which completes the proof. \square

Proof of Theorem 1.4. We begin with the proof of statement (i). Fix $\sigma > \sigma^*$, where the constant σ^* depends only on Σ, \mathbf{U} , and $\|\mathbf{u}\|_{C^1(\Omega)}$ and is defined by (4.51). Without any loss of generality we can assume that $f \in C^1(\Omega)$. The existence and uniqueness of a weak bounded solution for $\sigma > \sigma^*$ follows from Lemma 4.12. Moreover, by virtue of Lemma 4.12, such a solution satisfies the second inequality in (1.33). Therefore, it suffices to prove estimate (1.33) for $\|\varphi\|_{W^{s,r}(\Omega)}$. Since $W^{s,r}(\Omega) \cap L^\infty(\Omega)$ is the Banach algebra, representation (4.48) together with inequality (4.47) implies

$$(4.58) \quad \|\varphi\|_{W^{s,r}(\Omega)} \leq c(1 + t^{-1})(\|\varphi_t\|_{W^{s,r}(\Omega_t)} + \|\varphi_t\|_{L^\infty(\Omega_t)}) + c\|\phi\|_{W^{s,r}(\Omega)}.$$

On the other hand, Lemma 4.15 along with (4.49) yields

$$\|\phi\|_{W^{s,r}(\Omega)} \leq c\|F\|_{W^{s,r}(\Omega)} \leq c\|\chi_{t/2} f\|_{W^{s,r}(\Omega)} + \|\varphi_t \mathbf{u} \nabla \chi_{t/2}\|_{W^{s,r}(\Omega)}.$$

The first term in the right-hand side is bounded:

$$\|\chi_{t/2} f\|_{W^{s,r}(\Omega)} \leq c(1 + t^{-1})(\|f\|_{W^{s,r}(\Omega)} + \|\varphi_t\|_{L^\infty(\Omega)}).$$

In order to estimate the second term we note that, by virtue of (4.47), $\|\mathbf{u} \nabla \chi_{t/2}\|_{C^1(\Omega)} \leq cM(1 + t^{-2})$, which gives

$$\|\varphi_t \mathbf{u} \nabla \chi_{t/2}\|_{W^{s,r}(\Omega)} \leq cM(1 + t^{-2})(\|\varphi_t\|_{W^{s,r}(\Omega_t)} + \|\varphi_t\|_{L^\infty(\Omega_t)}).$$

Substituting the obtained estimates into (4.58) we arrive at the inequality

$$\|\varphi\|_{W^{s,r}(\Omega)} \leq c(M + 1)(1 + t^{-2})(\|\varphi_t\|_{W^{s,r}(\Omega)} + \|\varphi_t\|_{L^\infty(\Omega_t)} + \|f\|_{W^{s,r}(\Omega_t)} + \|f\|_{L^\infty(\Omega)}),$$

which along with (4.40) leads to the estimate (1.33). In order to prove statement (ii) of Theorem 1.4, we note that the adjoint equation can be written in the form

$$-\mathbf{u} \nabla \varphi^* + \sigma \varphi^* = f + \varphi^* \operatorname{div} \mathbf{u}.$$

Since

$$\|\operatorname{div} \mathbf{u}\varphi^*\|_{W^{s,r}(\Omega)} \leq c(\|\operatorname{div} \mathbf{u}\|_{W^{s,r}(\Omega)} + \|\operatorname{div} \mathbf{u}\|_{C(\Omega)})\|\varphi^*\|_{W^{s,r}(\Omega)},$$

we have

$$\|\operatorname{div} \mathbf{u}\varphi^*\|_{W^{s,r}(\Omega)} + \|\operatorname{div} \mathbf{u}\varphi^*\|_{C(\Omega)} \leq \delta(\|\varphi^*\|_{W^{s,r}(\Omega)} + \|\varphi^*\|_{C(\Omega)}),$$

and the needed result follows from (i) and the contraction mapping principle. \square

Appendix A. Proof of Lemmas 1.7 and 2.2.

Proof of Lemma 1.7. Since $\partial\Omega$ belongs to the class C^1 , functions φ, ς have the extensions $\bar{\varphi}, \bar{\varsigma} \in W^{s,r}(\Omega) \cap W^{1,2}(\Omega)$ such that $\bar{\varphi}, \bar{\varsigma}$ are compactly supported in \mathbb{R}^d and

$$\|\bar{\varphi}\|_{W^{s,r}(\mathbb{R}^d)} \leq c\|\varphi\|_{W^{s,r}(\Omega)}, \quad \|\bar{\varsigma}\|_{W^{s,r}(\mathbb{R}^d)} \leq c\|\varsigma\|_{W^{s,r}(\Omega)}.$$

By virtue of Definition 1.2 and inequality (1.19), function \mathbf{w} has the extension by 0 outside Ω , denoted by $\bar{\mathbf{w}}$, such that

$$\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^d)} \leq c\|\mathbf{w}\|_{W_0^{1-s,r'}(\Omega)}.$$

Obviously we have

$$\mathfrak{B}(\mathbf{w}, \varphi, \varsigma) = - \int_{\mathbb{R}^d} \bar{\mathbf{w}} \cdot \nabla \bar{\varphi} \bar{\varsigma} \, dx.$$

The following multiplicative inequality is due to Maz'ya [26]. For all $s_1 > 0, r_1 > 1$, and $r_1 s_1 < d$,

$$(5.1) \quad \|uv\|_{W^{s_1,r_1}(\mathbb{R}^d)} \leq c(r_1, s_1, d)(\|v\|_{W^{s_1,d/s_1}(\mathbb{R}^d)} + \|v\|_{L^\infty(\mathbb{R}^d)})\|u\|_{W^{s_1,r_1}(\mathbb{R}^d)}.$$

By virtue of (5.1), we have

$$\|\bar{\mathbf{w}}\bar{\varsigma}\|_{W^{1-s,r'}(\mathbb{R}^d)} \leq c\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^d)}(\|\bar{\varsigma}\|_{W^{1-s,d/(1-s)}(\mathbb{R}^d)} + \|\bar{\varsigma}\|_{L^\infty(\mathbb{R}^d)}).$$

On the other hand, since $r^{-1} - (s - (1 - s))/d \leq (1 - s)/d$ for $sr > d$, embedding inequality (1.21) yields

$$\|\bar{\varsigma}\|_{W^{1-s,d/(1-s)}(\mathbb{R}^d)} \leq c\|\bar{\varsigma}\|_{W^{s,r}(\mathbb{R}^d)}, \quad \|\bar{\varsigma}\|_{L^\infty(\mathbb{R}^d)} \leq c\|\bar{\varsigma}\|_{W^{s,r}(\mathbb{R}^d)}.$$

Thus we get

$$\|\bar{\mathbf{w}}\bar{\varsigma}\|_{W^{1-s,r'}(\mathbb{R}^d)} \leq c\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^d)}\|\bar{\varsigma}\|_{W^{s,r}(\mathbb{R}^d)}.$$

Next note that the operator $\nabla : W^{1,r}(\mathbb{R}^d) \mapsto W^{0,r}(\mathbb{R}^d)$, $\nabla : W^{0,r}(\mathbb{R}^d) \mapsto W^{-1,r}(\mathbb{R}^d)$ is continuous. By virtue of the basic property of interpolation spaces, the operator $\nabla : W^{s,r}(\mathbb{R}^d) \mapsto W^{s-1,r}(\mathbb{R}^d)$ is continuous for all $s \in [0, 1]$. In particular we have $\|\nabla \bar{\varphi}\|_{W^{s-1,r}(\mathbb{R}^d)} \leq c\|\bar{\varphi}\|_{W^{s,r}(\mathbb{R}^d)}$. Since $\bar{\varphi}$ and $\bar{\mathbf{w}}\bar{\varsigma}$ are compactly supported in \mathbb{R}^d , the duality principle implies

$$\begin{aligned} \int_{\mathbb{R}^d} \bar{\mathbf{w}}\bar{\varsigma}\nabla \bar{\varphi} \, dx &\leq c\|\bar{\mathbf{w}}\bar{\varsigma}\|_{W^{1-s,r'}(\mathbb{R}^d)}\|\nabla \bar{\varphi}\|_{W^{s-1,r}(\mathbb{R}^d)} \\ &\leq c\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^d)}\|\bar{\varsigma}\|_{W^{s,r}(\mathbb{R}^d)}\|\bar{\varphi}\|_{W^{s,r}(\mathbb{R}^d)}, \end{aligned}$$

which completes the proof. \square

Proof of Lemma 2.2. By virtue of (1.19), the extension $\bar{\mathbf{w}}$ satisfies the inequalities

$$\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^3)} \leq c\|\mathbf{w}\|_{\mathcal{W}_0^{1-s,r'}(\Omega)}, \quad \|\bar{\mathbf{w}}\|_{W^{1,2}(\mathbb{R}^3)} \leq c\|\mathbf{w}\|_{W^{1,2}(\Omega)}.$$

On the other hand, the vector field \mathbf{h} has a compactly supported extension $\bar{\mathbf{h}} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ such that $\|\bar{\mathbf{h}}\|_{W^{1+s,r}(\mathbb{R}^3)} \leq c\|\mathbf{h}\|_{W^{1+s,r}(\Omega)}$; however, this extension does not vanish outside Ω . Substituting the expression for \mathcal{A} into the formula for \mathfrak{A} and integrating by parts, we conclude that $\mathfrak{A}(\mathbf{w}, \mathbf{h})$ equals

$$\int_{\mathbb{R}^3} \mathbf{g}^{-1} \left(\nabla \left((\mathbf{N}_1^*)^{-1} - \mathbf{I} \right) \bar{\mathbf{w}} : \left(\mathbf{N}_1 \mathbf{N}_1^* \nabla (\mathbf{N}_1^{-1*} \bar{\mathbf{h}}) \right) \right. \\ \left. + \nabla \bar{\mathbf{w}} : \left(\mathbf{N}_1 \mathbf{N}_1^* \nabla (\mathbf{N}_1^{-1*} \bar{\mathbf{h}}) - \mathbf{g} \nabla \bar{\mathbf{h}} \right) \right) dx.$$

Since $\|\mathbf{N}_1^{\pm 1} - \mathbf{I}\|_{C^2(\Omega)} \leq c\tau^2$, we have

$$\|((\mathbf{N}_1^*)^{-1} - \mathbf{I})\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^3)} \leq c\tau^2\|\bar{\mathbf{w}}\|_{W^{1-s,r'}(\mathbb{R}^3)}, \\ \|\mathbf{g}^{-1}(\mathbf{N}_1 \mathbf{N}_1^* \nabla (\mathbf{N}_1^{-1*} \bar{\mathbf{h}})) - \nabla \bar{\mathbf{h}}\|_{W^{s,r}(\mathbb{R}^3)} \leq c\tau^2\|\bar{\mathbf{h}}\|_{W^{1+s,r}(\mathbb{R}^3)}.$$

Application of the duality arguments, similar to that in the proof of Lemma 1.7, completes the proof. \square

Appendix B. Interpolation. In this section we recall some results from the interpolation theory; see [6] for the proofs. Let A_0 and A_1 be Banach spaces. For $t > 0$ introduce two nonnegative functions $K : A_0 + A_1 \mapsto \mathbb{R}$ and $J : A_0 \cap A_1 \mapsto \mathbb{R}$ defined by

$$K(t, u, A_0, A_1) = \inf_{\substack{u=u_0+u_1 \\ u_i \in A_i}} \|u_0\|_{A_0} + t\|u_1\|_{A_1}, \\ J(t, u, A_0, A_1) = \max\{\|u\|_{A_0}, t\|u\|_{A_1}\}.$$

For each $s \in (0, 1)$, $1 < r < \infty$, the K -interpolation space $[A_0, A_1]_{s,r,K}$ consists of all elements $u \in A_0 + A_1$ having the finite norm

$$(6.1) \quad \|u\|_{[A_0, A_1]_{s,r,K}} = \left(\int_0^\infty t^{-1-sr} K(t, u, A_0, A_1)^r dt \right)^{1/r}.$$

On the other hand, J -interpolation space $[A_0, A_1]_{s,r,J}$ consists of all elements $u \in A_0 + A_1$ which admit the representation

$$(6.2) \quad u = \int_0^\infty \frac{v(t)}{t} dt, \quad v(t) \in A_1 \cap A_0 \text{ for } t \in (0, \infty),$$

and have the finite norm

$$(6.3) \quad \|u\|_{[A_0, A_1]_{s,r,J}} = \inf_{v(t)} \left(\int_0^\infty t^{-1-sr} J(t, v(t), A_0, A_1)^r dt \right)^{1/r} < \infty,$$

where the infimum is taken over the set of all $v(t)$ satisfying (6.2). The first main result of interpolation theory reads as follows: For all $s \in (0, 1)$ and $r \in (1, \infty)$ the spaces $[A_0, A_1]_{s,r,K}$ and $[A_0, A_1]_{s,r,J}$ are isomorphic topologically and algebraically. Hence the introduced norms are equivalent, and we can omit indices J and K . The following simple properties of interpolation spaces directly follow from definitions.

(1) If $A_1 \subset A_0$ is dense in A_0 , then $[A_0, A_1]_{s,r} \subset A_0$ is dense in A_0 .

(2) If $\tilde{A}_i, i = 0, 1$, are closed subspaces of A_i , then $[\tilde{A}_0, \tilde{A}_1]_{s,r} \subset [A_0, A_1]_{s,r}$ and $\|u\|_{[A_0, A_1]_{s,r}} \leq \|u\|_{[\tilde{A}_0, \tilde{A}_1]_{s,r}}$. One of the important results of the interpolation theory is the following representation for the interpolation of dual spaces. Let A_i be Banach spaces such that $A_1 \cap A_0$ is dense in $A_0 + A_1$. Then the Banach spaces $[(A_0)', (A_1)']_{s,r'}$ and $([A_0, A_1]_{s,r})'$ are isomorphic topologically and algebraically. Hence the spaces can be identified with equivalent norms.

In particular, if $A_1 \subset A_0, A'_0 \subset A'_1$ are dense in A_0 and A'_1 , respectively, then $([A_0, A_1]_{s,r})'$ is the completion of A'_0 in the $([A_0, A_1]_{s,r})'$ -norm.

The following lemma is the central result of the interpolation theory.

LEMMA B.1. *Let $A_i, B_i, i = 0, 1$, be Banach spaces, and let $T : A_i \mapsto B_i$ be a bounded linear operator. Then, for all $s \in (0, 1)$ and $r \in (1, \infty)$, the operator $T : [A_0, A_1]_{s,r} \mapsto [B_0, B_1]_{s,r}$ is bounded and*

$$\|T\|_{\mathcal{L}([A_0, A_1]_{s,r}, [B_0, B_1]_{s,r})} \leq \|T\|_{\mathcal{L}(A_0, B_0)}^s \|T\|_{\mathcal{L}(A_1, B_1)}^{1-s}.$$

Now we show that all basic properties of spaces $\mathcal{W}_0^{s,r}$ determined by Definition 1.2 easily follow from previously mentioned results of the interpolation theory. Let Ω be a bounded domain with a boundary of the class C^1 or $\Omega = \mathbb{R}^d$. It is well known (see [46]) that, for all $s \in (0, 1)$ and $r \in (1, \infty)$, the Sobolev space $W^{s,r}(\Omega) = [L^r(\Omega), W^{1,r}(\Omega)]_{s,r}$. Since $\mathcal{W}^{0,r}(\Omega)$ and $\mathcal{W}_0^{1,r}(\Omega)$ are closed subspaces of $W^{0,r}(\mathbb{R}^d)$ and $W^{1,r}(\mathbb{R}^d)$, respectively, the interpolating space $\mathcal{W}_0^{s,r}$ determined by Definition 1.2 satisfies inequality (1.19).

Next note that, by virtue of pairing (1.22), the space $L^{r'}(\Omega)$ can be identified with $(\mathcal{W}_0^{0,r})'$, which is dense in $\mathcal{W}^{-1,r}(\Omega) = (\mathcal{W}_0^{1,r}(\Omega))'$. Therefore, the space $(\mathcal{W}_0^{s,r})'$ is the completion of $L^{r'}(\Omega)$ in the norm of $(\mathcal{W}_0^{s,r}(\Omega))'$, which is exactly equal to the norm of $\mathcal{W}^{-s,r'}(\Omega)$. Hence $(\mathcal{H}_0^{s,r}(\Omega))' = \mathcal{W}^{-s,r'}(\Omega)$, which leads to the duality principle (1.24).

Proof of Lemma 1.3. Finally we show that Lemma 1.3 is a straightforward consequence of classical results on solvability of the first boundary value problem for the Stokes equations. Note that, by virtue of Theorem 6.1 in [9], for any $\mathbf{F} \in \mathcal{H}^{s-1,r}(\Omega)$ and $G \in H^{s,r}(\Omega)$ with $s = 0, 1$, problem (1.26) has a unique solution \mathbf{u}, π satisfying inequality

$$\|\mathbf{u}\|_{H^{s+1,r}(\Omega)} + \|\pi\|_{H^{s,r}(\Omega)} \leq c(\Omega, r, s)(\|\mathbf{F}\|_{\mathcal{H}^{s-1,r}(\Omega)} + \|G\|_{H^{s,r}(\Omega)}).$$

Thus the relation $(F, G) \mapsto (\mathbf{u}, \pi)$ determines the linear operator $T : \mathcal{H}^{s-1,r}(\Omega) \times H^{s,r}(\Omega) \mapsto H^{s+1,r}(\Omega) \times H^{s,r}(\Omega)$. Therefore, Lemma 1.3 is a consequence of Lemma B.1. \square

Appendix C. Change of variables. In this section we derive the equations (1.7). We will write $\mathbf{u}(y)$ and $\varrho(y), y \in \Omega$, and set

$$y = x + \varepsilon \mathbf{T}(x), \quad \mathbf{M}(x) = \mathbf{I} + \varepsilon \mathbf{T}'(x), \quad \tilde{\mathbf{u}}(x) = \mathbf{u}(y(x)), \quad \varrho_\varepsilon(x) = \varrho(y(x)).$$

Thus we get $\mathbf{u}_\varepsilon = N\tilde{\mathbf{u}}$. The Jacobi matrix \mathbf{M} is connected with the matrix \mathbf{N} by the relations

$$(7.1) \quad \det \mathbf{M} = (\det N)^{1/2} \equiv \mathfrak{g}, \quad \mathbf{M} = \mathfrak{g} \mathbf{N}^{-1}.$$

For any function $\phi \in C^1(\Omega)$ we have $\nabla_y \phi = (\mathbf{M}^*)^{-1} \nabla_x \tilde{\phi}$, where $\tilde{\phi}(x) = \phi(y(x))$. It follows from this that the identities

$$\begin{aligned} \int_{\tilde{\Omega}} (\operatorname{div}_y \mathbf{u})(y(x)) \tilde{\phi}(x) \det \mathbf{M} \, dx &= \int_{\Omega} (\operatorname{div}_y \mathbf{u})(y) \phi(y) \det \, dy = - \int_{\Omega} \mathbf{u} \cdot \nabla_y \phi \, dy \\ &= - \int_{\tilde{\Omega}} \tilde{\mathbf{u}} \cdot (\mathbf{M}^*)^{-1} \nabla_x \tilde{\phi}(x) \det \mathbf{M} \, dx = \int_{\tilde{\Omega}} \operatorname{div}_x ((\det \mathbf{M}) \mathbf{M}^{-1} \tilde{\mathbf{u}}) \tilde{\phi}(x) \, dx \end{aligned}$$

hold true for all $\phi \in C_0^\infty(\Omega)$. On the other hand, by virtue of (7.1) we have $(\det \mathbf{M}) \mathbf{M}^{-1} \tilde{\mathbf{u}} = \mathbf{u}_\varepsilon(x)$. This leads to the equalities

$$(7.2) \quad \begin{aligned} (\operatorname{div}_y \mathbf{u})(y(x)) &= \mathbf{g}^{-1} \operatorname{div}_x (\mathbf{N} \tilde{\mathbf{u}}(x)) \equiv \mathbf{g}^{-1} \operatorname{div}_x \mathbf{u}_\varepsilon(x), \\ \operatorname{div}_y (\varrho \mathbf{u})(y(x)) &= \mathbf{g}^{-1} \operatorname{div}_x (\varrho_\varepsilon \mathbf{u}_\varepsilon), \end{aligned}$$

which imply the modified mass balance equation (1.7b). From (7.2) and the identity $(\mathbf{M}^*)^{-1} = \mathbf{g}^{-1} \mathbf{N}^*$ we obtain

$$(7.3) \quad \nabla \left(\lambda \operatorname{div} \mathbf{u} - \frac{R}{\varepsilon^2} p(\varrho) \right) = \mathbf{g}^{-1} \mathbf{N}^* \nabla \left(\lambda \mathbf{g}^{-1} \operatorname{div} \mathbf{u}_\varepsilon - \frac{R}{\varepsilon^2} p(\varrho_\varepsilon) \right).$$

Combining (7.2) with the identity $\Delta = \operatorname{div} \nabla$ we obtain

$$(7.4) \quad \begin{aligned} \Delta \mathbf{u}(y) &= \mathbf{g}^{-1} \operatorname{div} (\mathbf{N} (\mathbf{M}^*)^{-1} \nabla \tilde{\mathbf{u}}) \\ &= \mathbf{g}^{-1} \operatorname{div} (\mathbf{g}^{-1} \mathbf{N} \mathbf{N}^* \nabla (\mathbf{N}^{-1} \mathbf{u}_\varepsilon)) = \mathbf{g}^{-1} \mathbf{N}^* (\Delta \mathbf{u}_\varepsilon - \mathcal{A}(\mathbf{u}_\varepsilon)). \end{aligned}$$

Next note that the components $(\mathbf{u} \nabla \mathbf{u})_i$ of the vector $\mathbf{u} \nabla \mathbf{u}$ satisfy the equalities

$$(\mathbf{u} \nabla \mathbf{u})_i = \mathbf{u} \cdot \nabla_y u_i = \tilde{\mathbf{u}} \cdot ((\mathbf{M}^*)^{-1} \nabla \tilde{u}_i) = \mathbf{g}^{-1} \mathbf{N} \tilde{\mathbf{u}} \cdot \nabla \tilde{u}_i = \mathbf{g}^{-1} \mathbf{u}_\varepsilon \cdot \nabla (\mathbf{N}^{-1} \mathbf{u}_\varepsilon)_i.$$

This gives

$$(7.5) \quad \varrho \mathbf{u} \nabla \mathbf{u} = \mathbf{g}^{-1} \mathbf{N}^* \mathcal{B}(\varrho_\varepsilon, \mathbf{u}_\varepsilon, \mathbf{u}_\varepsilon).$$

Substituting (7.3)–(7.5) into mass balance equation (1.2a) and multiplying both sides of the resulting equality by $\mathbf{g}(\mathbf{N}^*)^{-1}$, we obtain modified equation (1.7a).

REFERENCES

[1] R. A. ADAMS, *Sobolev Spaces*, Academic Press, New York, 1975.
 [2] H. BEIRAO DA VEIGA, *Stationary motions and the incompressible limit for compressible viscous limit*, Houston J. Math., 13 (1987), pp. 527–544.
 [3] H. BEIRAO DA VEIGA, *Existence results in Sobolev spaces for a transport equation*, Ricerche Mat., 36 (1987), pp. 173–184.
 [4] J. A. BELLO, E. FERNÁNDEZ-CARA, J. LEMOINE, AND J. SIMON, *The differentiability of the drag with respect to the variations of a Lipschitz domain in a Navier–Stokes flow*, SIAM J. Control. Optim., 35 (1997), pp. 626–640.
 [5] J. A. BELLO, E. FERNÁNDEZ-CARA, AND J. SIMON, *Optimal shape design for Navier–Stokes flow*, in System Modelling and Optimization, Lecture Notes in Control and Inform. Sci. 180, D. Kall, ed., Springer-Verlag, Berlin, 1992, pp. 481–489.
 [6] J. BERGH AND J. LÖFSTRÖM, *Interpolation Spaces. An Introduction*, Springer-Verlag, Berlin, Heidelberg, New York, 1976.
 [7] R. J. DiPERNA AND P. L. LIONS, *Ordinary differential equations, transport theory and Sobolev spaces*, Invent. Math., 98 (1989), pp. 511–547.
 [8] L. C. EVANS, *Partial Differential Equations*, AMS, Providence, RI, 1998.

- [9] G. GALDI, *An Introduction to the Mathematical Theory of the Navier–Stokes Equations*, Vol. VI, Springer-Verlag, Berlin, Heidelberg, New York, 1998.
- [10] G. FICHERA, *Sulle equazioni differenziali lineari ellittico-paraboliche del secondo ordine*, Atti Accad. Naz. Lincei, Mem. Cl. Sci. Fis. Mat. Nat., Sez. I, 5 (1956), pp. 1–30.
- [11] E. FEIREISL, *Dynamics of Viscous Compressible Fluids*, Oxford University Press, Oxford, UK, 2004.
- [12] E. FEIREISL, A. H. NOVOTNÝ, AND H. PETZELTOVÁ, *On the domain dependence of solutions to the compressible Navier–Stokes equations of a barotropic fluid*, Math. Methods Appl. Sci., 25 (2002), pp. 1045–1073.
- [13] E. FEIREISL, *Shape optimization in viscous compressible fluids*, Appl. Math. Optim., 47 (2003), pp. 59–78.
- [14] J. FREHSE, S. GOJ, AND M. STEINHAUER, *L^p -estimates for the Navier–Stokes equations for steady compressible flow*, Manuscripta Math., 116 (2005), pp. 265–275.
- [15] J. G. HEYWOOD AND M. PADULA, *On the uniqueness and existence theory for steady compressible viscous flow*, in Fundamental Directions in Mathematical Fluid Mechanics, Adv. Math. Fluid Mech., Birkhäuser, Basel, 2000, pp. 171–189.
- [16] L. HÖRMANDER, *Pseudo-differential operators and non-elliptic boundary value problems*, Ann. of Math. (2), 83 (1966), pp. 129–209.
- [17] W. JAGER AND A. MIKELIC, *Couette flow over a rough boundary and drag reduction*, Comm. Math. Phys., 232 (2003), pp. 429–455.
- [18] B. KAWOHL, O. PIRONNEAU, L. TARTAR, AND J. ZOLESIO, *Optimal Shape Design*, Lecture Notes in Math. 1740, Springer-Verlag, Berlin, 2000.
- [19] J. J. KOHN AND L. NIRENBERG, *Degenerate elliptic-parabolic equations of second order*, Comm. Pure Appl. Math., 20 (1967), pp. 797–872.
- [20] J. R. KWEON AND R. B. KELLOGG, *Compressible Navier–Stokes equations in a bounded domain with inflow boundary condition*, SIAM J. Math. Anal., 28 (1997), pp. 94–108.
- [21] J. R. KWEON AND R. B. KELLOGG, *Regularity of solutions to the Navier–Stokes equations for compressible barotropic flows on a polygon*, Arch. Ration. Mech. Anal., 163 (2000), pp. 36–64.
- [22] L. D. LANDAU AND E. M. LIFSHITZ, *Course of Theoretical Physics*, Vol. 6, *Fluid Mechanics*, Pergamon Press, Oxford, UK, 1987.
- [23] N. S. LANDKOF, *Foundations of Modern Potential Theory*, Springer-Verlag, Berlin, Heidelberg, New York, 1972.
- [24] P. L. LIONS, *Mathematical Topics in Fluid Dynamics*, Volume 1, *Incompressible Models*, Clarendon Press, Oxford, UK, 1996.
- [25] P. L. LIONS, *Mathematical Topics in Fluid Dynamics*, Volume 2, *Compressible Models*, Clarendon Press, Oxford, UK, 1998.
- [26] V. G. MAZ'YA AND T. O. SHAPOSHNIKOVA, *Multipliers in Spaces of Differential Functions*, Leningrad University, Leningrad, 1986.
- [27] B. MOHAMMADI AND O. PIRONNEAU, *Shape optimization in fluid mechanics*, in Annual Review of Fluid Mechanics, Annu. Rev. Fluid Mech. 36, Annual Reviews, Palo Alto, CA, 2004, pp. 255–279.
- [28] M. MOUBACHIR AND J.-P. ZOLESIO, *Moving Shape Analysis And Control: Applications to Fluid Structure Interactions*, Chapman & Hall/CRC, Boca Raton, FL, 2006.
- [29] A. NOIRI, F. POUPAUND, AND Y. DEMAY, *An existence theorem for the multi-fluid Stokes problem*, Quart. Appl. Math., 55 (1997), pp. 421–435.
- [30] A. NOVOTNY, *About steady transport equation. I. L^p -approach in domains with smooth boundaries*, Comment. Math. Univ. Carolin., 37 (1996), pp. 43–89.
- [31] A. NOVOTNY, *About steady transport equation. II. Schauder estimates in domains with smooth boundaries*, Portugal. Math., 54 (1997), pp. 317–333.
- [32] A. NOVOTNÝ AND M. PADULA, *Existence and uniqueness of stationary solutions for viscous compressible heat conductive fluid with large potential and small non-potential external forces*, Siberian Math. J., 34 (1993), pp. 120–146.
- [33] A. NOVOTNÝ AND I. STRAŠKRABA, *Introduction to the Mathematical Theory of Compressible Flow*, Oxford Lecture Ser. Math. Appl. 27, Oxford University Press, Oxford, UK, 2004.
- [34] O. A. OLEINIK AND E. V. RADKEVICH, *Second order equations with non-negative characteristic form*, AMS, Providence, RI, Plenum Press, New York, London, 1973.
- [35] M. PADULA, *Existence and uniqueness for viscous steady compressible motions*, Arch. Rational Mech. Anal., 97 (1986), pp. 1–20.
- [36] M. PADULA, *Steady flows of barotropic viscous fluids*, in Classical Problems in Mechanics, Quad. Mat. 1, Dipartimento di Matematica, Seconda Università di Napoli, Caserta, Italy, 1997, pp. 253–345.

- [37] P. I. PLOTNIKOV AND J. SOKOLOWSKI, *On compactness, domain dependence, and existence of steady state solutions to compressible isothermal Navier-Stokes equations*, J. Math. Fluid Mech., 7 (2005), pp. 529–573.
- [38] P. I. PLOTNIKOV AND J. SOKOLOWSKI, *Concentrations of solutions to time-discretized compressible Navier-Stokes equations*, Comm. Math. Phys., 258 (2005), pp. 567–608.
- [39] P. I. PLOTNIKOV AND J. SOKOLOWSKI, *Stationary boundary value problems for Navier-Stokes equations with adiabatic index $\gamma < 3/2$* , Doklady Mathematics, 70 (2004), pp. 535–538 (in English); Dok. Akad. Nauk, 397 (2004), pp. 166–169 (in Russian).
- [40] P. I. PLOTNIKOV AND J. SOKOLOWSKI, *Domain dependence of solutions to compressible Navier-Stokes equations*, SIAM J. Control Optim., 45 (2006), pp. 1165–1197.
- [41] H. SCHLICHTING, *Boundary-Layer Theory*, McGraw-Hill Ser. Mech. Engrg., McGraw-Hill, New York, 1955.
- [42] J. SIMON, *Domain variation for drag in Stokes flow*, in Control Theory of Distributed Parameter Systems and Applications, Lecture Notes in Control and Inform. Sci. 159, Springer-Verlag, Berlin, 1991, pp. 28–42.
- [43] T. SLAWIG, *A formula for the derivative with respect to domain variations in Navier-Stokes flow based on an embedding domain method*, SIAM J. Control Optim., 42 (2003), pp. 495–512.
- [44] T. SLAWIG, *An explicit formula for the derivative of a class of cost functionals with respect to domain variations in Stokes flow*, SIAM J. Control Optim., 39 (2000), pp. 141–158.
- [45] J. SOKOLOWSKI AND J.-P. ZOLÉSIO, *Introduction to Shape Optimization. Shape Sensitivity Analysis*, Springer Ser. Comput. Math. 16, Springer-Verlag, Berlin, 1992.
- [46] L. TARTAR, *An introduction to Sobolev spaces and interpolation spaces*, Lecture Notes of the Unione Matematica Italiana 3, Springer-Verlag, Berlin, and UMI, Bologna, Italy, 2007.

QUASI-STATIC EVOLUTION FOR A MODEL IN STRAIN GRADIENT PLASTICITY*

ALESSANDRO GIACOMINI[†] AND LUCA LUSSARDI[†]

Abstract. We prove the existence of a quasi-static evolution for a model in strain gradient plasticity proposed by Gurtin and Anand concerning isotropic, plastically irrotational materials under small deformations. This is done by means of the energetic approach to rate-independent evolution problems. Finally we study the asymptotic behavior of the evolution as the strain gradient length scales tend to zero recovering in the limit a quasi-static evolution in perfect plasticity.

Key words. variational models, energy minimization, quasi-static evolution, higher-order stress, flow rule

AMS subject classifications. 74D99, 74C05, 74G65, 49J45.

DOI. 10.1137/070708202

1. Introduction. Since the early attempts of Aifantis [3], strain gradient plasticity models have been proposed in order to capture size effects in metals such as ϵ -plasticity and ϵ -viscoplasticity. These effects, which take place approximately at the scale of 500 nm–50 μ m, cannot be modeled by conventional theories of plasticity. This fact led to the development of continuum theories of plasticity that incorporate size dependence by accounting for strain gradients, namely, the gradient of plastic strains. Following the classical papers by Nye [30] and by Ashby [4, 5], strain gradients induce dislocations, and these dislocations together with grain boundaries are mainly responsible for size effects.

Several strain gradient theories, different from one another, have been recently proposed by different authors [1, 7, 10, 11, 14, 15, 16, 17, 18, 20, 12, 21]. In this paper we focus on the theory proposed by Gurtin and Anand [19]. In the context of small deformations, and in the absence of plastic rotation, the strain gradient dependence enters the model via a microstress associated to the gradient of the plastic strain and by a free energy dependent of the macroscopic Burgers tensor.

Let $\Omega \subseteq \mathbb{R}^3$ be the reference configuration of the body. The strain $(\mathbf{E}u)_{ij} := (\partial_i u_j + \partial_j u_i)/2$ of the displacement $u : \Omega \rightarrow \mathbb{R}^3$ is decomposed as usual in the form

$$(1.1) \quad \mathbf{E}u = \mathbf{E}^e + \mathbf{E}^p,$$

where $\mathbf{E}^e \in M_{\text{sym}}^{3 \times 3}$ is the elastic strain, while \mathbf{E}^p is referred to as the plastic strain. It is assumed that \mathbf{E}^p has zero trace; i.e., \mathbf{E}^p belongs to the space of deviatoric matrices $M_D^{3 \times 3}$. Besides the usual Cauchy stress \mathbf{T} , which satisfies the classical macroscopic force balance, the stress configuration of the system is described by the deviatoric stress tensor \mathbf{T}^p and a higher-order stress tensor \mathbb{K}^p which satisfy the equilibrium condition

$$(1.2) \quad \mathbf{T}_D = \mathbf{T}^p - \text{div} \mathbb{K}^p.$$

*Received by the editors November 14, 2004; accepted for publication (in revised form) April 22, 2008; published electronically October 29, 2008.

<http://www.siam.org/journals/sima/40-3/70820.html>

[†]Dipartimento di Matematica, Facoltà di Ingegneria, Università degli Studi di Brescia, Via Valotti 9, 25133 Brescia, Italy (alessandro.giacomini@ing.unibs.it, luca.lussardi@ing.unibs.it).

Here \mathbf{T}_D denotes the deviatoric part of \mathbf{T} , i.e., $\mathbf{T}_D := \mathbf{T} - \frac{1}{3}\text{tr}(\mathbf{T})Id$. The triple $(\mathbf{T}, \mathbf{T}^P, \mathbb{K}^P)$ furnishes the internal power expenditure within a subbody $\mathcal{B} \subseteq \Omega$ by means of the relation

$$\mathcal{W}_{\text{int}}(\mathcal{B}) = \int_{\mathcal{B}} (\mathbf{T} : \dot{\mathbf{E}}^e + \mathbf{T}^P : \dot{\mathbf{E}}^P + \mathbb{K}^P : \nabla \dot{\mathbf{E}}^P) dx,$$

where $\dot{\mathbf{E}}^e$ and $\dot{\mathbf{E}}^P$ derive from a virtual velocity $(\dot{u}, \dot{\mathbf{E}}^e, \dot{\mathbf{E}}^P)$ of the system. So \mathbf{T}^P and \mathbb{K}^P are stresses conjugated to the plastic strain and its gradient, \mathbb{K}^P being a third order tensor, since it is conjugated to the gradient of the plastic strain. The balance equations for \mathbf{T} , \mathbf{T}^P , and \mathbb{K}^P follow by equating the internal power expenditure to the power expenditure associated to the external loads. This entails also boundary conditions for the normal components of \mathbf{T} and \mathbb{K}^P which are connected to the imposed traction and, respectively, on parts of the boundary (see section 3 for details).

The free energy of the system is a function of the elastic strain \mathbf{E}^e and of the macroscopic Burgers tensor $\mathbf{G} = \text{curl} \mathbf{E}^P$. In the separable quadratic isotropic case, it assumes the form

$$(1.3) \quad \psi = \mu |\mathbf{E}_D^e|^2 + \frac{1}{2} k |\text{tr} \mathbf{E}^e|^2 + \frac{\mu L^2}{2} |\text{curl} \mathbf{E}^P|^2,$$

where μ and k are the elastic shear and bulk moduli and L is an energetic length scale. The presence of $\text{curl} \mathbf{E}^P$ inside the free energy accounts for the incompatibility of the tensor field \mathbf{E}^P , and so it is connected to the presence of geometrically necessary dislocations in Ω . By means of ψ , the energetic third order tensor \mathbb{K}_{en}^P is defined as the symmetric-deviatoric part (in the first two subscripts) of $\frac{\partial \psi}{\partial \mathbf{G}}$. This entails a decomposition of \mathbb{K}^P into dissipative and energetic parts $\mathbb{K}_{\text{diss}}^P$ and \mathbb{K}_{en}^P , respectively, with

$$\mathbb{K}^P = \mathbb{K}_{\text{diss}}^P + \mathbb{K}_{\text{en}}^P.$$

Let Ω be subject to body forces $f(t)$ and to traction forces $g(t)$ on a part $\partial_N \Omega$ of its boundary, with $t \in [0, T]$. Let $\partial \Omega$ be subject to prescribed displacements $w(t)$; i.e., $w(t) = 0$ on $\partial_D \Omega$ and $w(t) = w(t)$ on $\partial_N \Omega$, where $\partial \Omega = \partial_D \Omega \cup \partial_N \Omega$ occurs (see section 3 for details). Let us assume that a displacement $w(t)$ is imposed on $\partial_D \Omega := \partial \Omega \setminus \partial_N \Omega$. The laws governing the evolution $(u(t), \mathbf{E}^e(t), \mathbf{E}^P(t))$ of the system are obtained by the thermodynamical requirement

$$\dot{\psi}(\mathcal{B}) \leq \mathcal{W}_{\text{int}}(\mathcal{B}),$$

where $\psi(\mathcal{B})$ is the free energy of the subbody \mathcal{B} obtained by integrating (1.3) over \mathcal{B} and $\dot{\psi}(\mathcal{B})$ denotes its time derivative. In order to match such an inequality, Gurtin and Anand propose a flow rule involving $\dot{\mathbf{E}}^P(t)$, $\nabla \dot{\mathbf{E}}^P(t)$, $\mathbf{T}^P(t)$, $\mathbb{K}_{\text{diss}}^P(t)$, a hardening internal variable $d^P(t, x)$, $l > 0$, and a hardening internal variable. This law reduces to the usual flow rules of classical plasticity when the length scales l and L are set to zero. In the case of a hardening internal variable, and neglecting the hardening internal variable, it takes the form

$$(1.4) \quad \mathbf{T}^P(t, x) = S_Y \frac{\dot{\mathbf{E}}^P(t, x)}{d^P(t, x)}, \quad \mathbb{K}_{\text{diss}}^P(t, x) = S_Y \frac{l^2 \nabla \dot{\mathbf{E}}^P(t, x)}{d^P(t, x)}.$$

Here $\dot{\mathbf{E}}^P(t, x)$ and $\nabla \dot{\mathbf{E}}^P(t, x)$ denote the time derivative of $\mathbf{E}^P(t, x)$ and $\nabla \mathbf{E}^P(t, x)$, respectively, S_Y is the yield strength, and

$$d^P(t, x) := \sqrt{|\dot{\mathbf{E}}^P(t, x)|^2 + l^2 |\nabla \dot{\mathbf{E}}^P(t, x)|^2}$$

is an ff ... fl ... The stresses $\mathbf{T}^P(t)$ and $\mathbb{K}_{\text{diss}}^P(t)$ satisfy the

$$(1.5) \quad \sqrt{|\mathbf{T}^P(x)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^P(x)|^2} \leq S_Y,$$

and (1.4) is valid when relation (1.5) holds with equality; $(\dot{\mathbf{E}}^P(t), \nabla \dot{\mathbf{E}}^P(t)) = (0, 0)$ otherwise. Notice that, by setting $l = L = 0$, we have $\mathbb{K}^P = 0$ and $\mathbf{T}^P = \mathbf{T}_D$ and (1.4) reduces to the usual flow rule of the von Mises type.

The aim of the paper is to provide an existence result of an evolution for the Gurtin–Anand model in the rate-independent case without hardening. The case with positive hardening has been considered recently by Reddy, Ebobisse, and McBride [31]. Adopting a primal formulation, they study the problem by means of variational inequalities in abstract Hilbert spaces. In the case without hardening, coercivity estimates fail, and the use of the abstract setting is no longer possible. This fact reflects what happens also at the level of classical plasticity, where perfect plasticity deserves an “ad hoc” treatment (see [32, 8]).

Inspired by the recent paper of Dal Maso, DeSimone, and Mora [8] concerning perfect plasticity, we recast the problem of the evolution for the Gurtin–Anand model in the framework of the energetic approach to rate-independent processes developed in [24, 25, 27, 28, 29].

Let us consider $\Omega \subseteq \mathbb{R}^N$ open, bounded, and with a Lipschitz boundary ($N \geq 2$). By means of variational arguments, we first construct a discretized-in-time evolution $(u_{k,i}, \mathbf{E}_{k,i}^e, \mathbf{E}_{k,i}^p)$ relative to the nodes t_k^i of a subdivision $0 = t_k^0 < t_k^1 < \dots < t_k^k = T$ of the time interval $[0, T]$ with step T/k .

In order to enforce variationally the stress constraint (1.5), we consider the function

$$\mathbf{E}^P \mapsto S_Y \int_{\Omega} \sqrt{|\mathbf{E}^P|^2 + l^2|\nabla \mathbf{E}^P|^2} dx.$$

Since this map has linear growth in $\nabla \mathbf{E}^P$, in order to perform direct minimization, we are naturally led to consider \mathbf{E}^P as a... $BV(\Omega; M_D^{N \times N})$ and to relax the functional to the form

$$\mathcal{H}(\mathbf{E}^P) := S_Y \int_{\Omega} \sqrt{|\mathbf{E}^P|^2 + l^2|\nabla \mathbf{E}^P|^2} dx + lS_Y|D^s \mathbf{E}^P|(\Omega),$$

where $D^s \mathbf{E}^P$ denotes the singular part of the derivative of \mathbf{E}^P .

The minimization problem that we consider in order to construct $(u_{k,i+1}, \mathbf{E}_{k,i+1}^e, \mathbf{E}_{k,i+1}^p)$ relative to the boundary displacement $w(t_k^{i+1})$ once constructed $(u_{k,i}, \mathbf{E}_{k,i}^e, \mathbf{E}_{k,i}^p)$ is the following:

$$(1.6) \quad \min_{(u, \mathbf{E}^e, \mathbf{E}^p) \in \mathcal{A}(w(t_k^{i+1}))} \mathcal{Q}_1(\mathbf{E}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p) - \langle \mathcal{L}(t_k^{i+1}), u \rangle + \mathcal{H}(\mathbf{E}^p - \mathbf{E}_{k,i}^p).$$

Here $\mathcal{A}(w(t_k^{i+1}))$ is the class of admissible configurations for $w(t_k^{i+1})$,

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}^e) &:= \int_{\Omega} \left(\mu |\mathbf{E}_D^e|^2 + \frac{1}{2} k |\text{tr} \mathbf{E}^e|^2 \right) dx, & \mathcal{Q}_2(\text{curl} \mathbf{E}^p) &:= \frac{\mu L^2}{2} \int_{\Omega} |\text{curl} \mathbf{E}^p|^2 dx, \\ \langle \mathcal{L}(t), u \rangle &:= \int_{\Omega} f(t) \cdot u dx + \int_{\partial_N \Omega} g(t) \cdot u d\mathcal{H}^{N-1}, \end{aligned}$$

where \mathcal{H}^{N-1} denotes the $(N - 1)$ -dimensional Hausdorff measure.

In order to have a well-defined energy in (1.6), it suffices that the elastic strain \mathbf{E}^e and the Burgers tensor $\text{curl}\mathbf{E}^p$ belong to the space of square integrable functions. As a consequence, the class $\mathcal{A}(t_k^i)$ turns out to be defined as the triples $(u, \mathbf{E}^e, \mathbf{E}^p)$, with

$$\begin{aligned} u &\in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N), & \mathbf{E}^e &\in L^2(\Omega; M_{\text{sym}}^{N \times N}), \\ \mathbf{E}^p &\in BV(\Omega; M_D^{N \times N}), & \text{curl}\mathbf{E}^p &\in L^2(\Omega; M^{N \times N}), \end{aligned}$$

which satisfy the boundary condition $u = w(t_k^i)$ on $\partial_D\Omega$, and such that the compatibility condition (1.1) holds. Notice that the requirement $u \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ follows by (1.1) and by the assumptions on \mathbf{E}^e and \mathbf{E}^p in view of Korn’s inequality. We assume that $f(t) \in L^N(\Omega; \mathbb{R}^N)$ and $g(t) \in L^N(\partial_N\Omega; \mathbb{R}^N)$ so that the work $\mathcal{L}(t)$ of external forces turns out to be well-defined. The displacement on $\partial_D\Omega$ is assumed to be given by the trace of a map in $W^{1,2}(\Omega; \mathbb{R}^N)$.

The minimum problem (1.6) admits solutions in $\mathcal{A}(w(t_k^i))$ provided that the external loads satisfy a suitable safe load condition (see (4.13)–(4.14)) which appears also in the study of evolutions in perfect plasticity. This condition entails some coercivity in BV for \mathbf{E}^p from the interaction between $\mathcal{H}(\mathbf{E}^p - \mathbf{E}_{k,i-1}^p)$ and the linear term $\langle \mathcal{L}(t_k^i), u \rangle$. The existence of a solution for (1.6) follows by applying the direct method of the calculus of variations (Lemma 6.1).

The continuous-in-time evolution is obtained by interpolating the discrete evolution $(u_{k,i}, \mathbf{E}_{k,i}^e, \mathbf{E}_{k,i}^p)$ and sending $k \rightarrow +\infty$ (section 7). If $w \in AC(0, T; W^{1,2}(\Omega; \mathbb{R}^N))$, $f \in AC(0, T; L^N(\Omega; \mathbb{R}^N))$, $g \in AC(0, T; L^\infty(\partial_N\Omega; \mathbb{R}^N))$, and the safe load condition on f, g holds uniformly in time, we prove the convergence towards a quasi-static evolution $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}(w(t))$ which is absolutely continuous in time and which satisfies the following two conditions:

- (a) Global minimality: For every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$

$$\begin{aligned} (1.7) \quad \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(t)) - \langle \mathcal{L}(t), u(t) \rangle \\ \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl}\mathbf{p}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^p(t)); \end{aligned}$$

- (b) Energy balance:

$$\begin{aligned} (1.8) \quad \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) = \mathcal{E}(0) + \int_0^t \int_{\Omega} \mathbf{T}(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau \\ - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau, \end{aligned}$$

where $\mathbf{T}(t)$ is the Cauchy stress tensor, $\mathcal{E}(t) := \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(t)) - \langle \mathcal{L}(t), u(t) \rangle$, $\dot{\mathcal{L}}(t)$ is associated to $\dot{f}(t)$, $\dot{g}(t)$, and $\mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t)$ defined as

$$\mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; a, b) := \sup \left\{ \sum_{j=1}^k \mathcal{H}(\mathbf{E}^p(t_j) - \mathbf{E}^p(t_{j-1})) : a = t_0 < t_1 < \dots < t_k = b \right\}$$

has the role of a ... We refer to an evolution satisfying (a) and (b) as a ... for the Gurtin–Anand model (Definition 5.1).

The analysis of the global minimality condition (1.7) leads to the existence of stresses $\mathbf{T}^p(t)$, $\mathbb{K}^p(t)$, and $\mathbb{S}^p(t)$ which together with the Cauchy stress $\mathbf{T}(t)$ satisfy the balance of internal and external powers in Ω

$$(1.9) \quad \int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \int_{\Omega} \mathbf{T}^p(t) : \mathbf{p} \, dx + \int_{\Omega} \mathbb{K}^p(t) : \nabla \mathbf{p} \, dx + \langle \mathbb{S}^p(t), D^s \mathbf{p} \rangle = \langle \mathcal{L}(t), v \rangle$$

for every virtual velocity $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$ (Lemma 8.1). Notice that a new higher-order stress $\mathbb{S}^p(t)$ conjugated to $D^s \mathbf{E}^p$ appears from our approach: This is somehow natural since $D^s \mathbf{E}^p$ is treated at the same level of $\nabla \mathbf{E}^p$. The balance (1.9) entails the usual balance equation for the Cauchy stress (Proposition 8.2), the balance equation (1.2), the stress constraint (1.5), and the confinement $\|\mathbb{S}^p(t)\| \leq lS_Y$ for the singular stress $\mathbb{S}^p(t)$ (Proposition 8.3).

The flow rule (1.4) follows from the analysis of the energy balance equality (1.8) (Proposition 8.10). It is also supplemented by a weak flow rule for the singular stress $\mathbb{S}^p(t)$ (Proposition 8.9).

Concerning the uniqueness of the evolution, it turns out that the maps $t \mapsto \mathbf{E}^e(t)$ and $t \mapsto \text{curl} \mathbf{E}^p(t)$ are uniquely determined by the initial conditions (Proposition 8.8). As for the maps $t \mapsto u(t)$ and $t \mapsto \mathbf{E}^p(t)$, we suspect that nonuniqueness can occur as in the case of standard perfect plasticity (see [32, section 2.1]), even if we do not have at the moment an explicit counterexample.

In section 9, we study the asymptotic behavior of a quasi-static evolution for the Gurtin–Anand model when the length scales l and L vanish. As noted previously, by setting l and L equal to zero, the model reduces to the classical model of perfect plasticity of von Mises. Under a suitable assumption on the initial configuration, we prove (Theorem 9.2) that the quasi-static evolution for the Gurtin–Anand model converges in a suitable sense to the evolution for elastic-perfectly plastic bodies in the framework proposed by Dal Maso, DeSimone, and Mora [8]. The main difficulty we have to handle is the change in the mathematical setting of the problem, especially concerning the plastic strain. While in the strain gradient context \mathbf{E}^p is a BV function, in [8] it is modeled simply as a Radon measure.

The paper is organized as follows. In section 2 we fix the notation and recall some basic tools we need from the theory of BV functions. In section 3 we give a brief sketch of the Gurtin–Anand model, while in section 4 we settle the mathematical framework that we adopt in the analysis. The main results are stated in section 5. The existence of a quasi-static evolution is obtained in section 7 after exploiting the convergence of the discrete evolution constructed in section 6. Section 8 is devoted to the proof of the balance equations and the flow rule. Finally section 9 contains the asymptotic analysis as the strain gradient effects vanish.

2. Notation and preliminaries. In this section we recall some basic definitions and results employed in the rest of the paper.

Matrices. We will denote by $M^{N \times N}$ the space of $N \times N$ matrices $\mathbf{A} = (a_{ij})$, with $a_{ij} \in \mathbb{R}$ endowed with the scalar product

$$(2.1) \quad \mathbf{A} : \mathbf{B} := \sum_{i,j} a_{ij} b_{ij}.$$

The norm of \mathbf{A} induced by the scalar product (2.1) is denoted by $|\mathbf{A}|$.

We will denote by $M_{\text{sym}}^{N \times N}$ the subspace of symmetric matrices and by $M_D^{N \times N}$ the subspace of $M_{\text{sym}}^{N \times N}$ of matrices \mathbf{A} with zero trace, that is, such that $\text{tr} \mathbf{A} := \sum_i a_{ii} = 0$.

Given $\mathbf{A} \in M_{\text{sym}}^{N \times N}$, we denote by \mathbf{A}_D its projection on $M_D^{N \times N}$, i.e.,

$$(2.2) \quad \mathbf{A}_D := \mathbf{A} - \frac{1}{N}(\text{tr} \mathbf{A})\mathbf{Id},$$

where \mathbf{Id} is the identity matrix.

The symmetrized gradient of an \mathbb{R}^N -valued function $u(x)$ is defined as

$$\mathbf{E}u := \frac{\nabla u + \nabla u^T}{2},$$

where $(\nabla u)_{ij} = \frac{\partial u_i}{\partial x_j}$ is the gradient of u and ∇u^T denotes its transpose.

The gradient, the divergence, and the curl of a $M^{N \times N}$ -valued function $\mathbf{A}(x) = (a_{ij}(x))$ are defined as

$$(\nabla \mathbf{A})_{ijk} := \frac{\partial a_{ij}}{\partial x_k}, \quad (\text{div} \mathbf{A})_i := \sum_j \frac{\partial a_{ij}}{\partial x_j}, \quad (\text{curl} \mathbf{A})_{ij} := \sum_{p,q} \epsilon_{ipq} \frac{\partial a_{jq}}{\partial x_p},$$

respectively, where ϵ_{ipq} are the standard permutation symbols.

We will indicate by $M^{N \times N \times N}$ the space of third order tensors $\mathbb{A} = (a_{ijk})$ with scalar product

$$\mathbb{A} : \mathbb{B} := \sum_{i,j,k} a_{ijk} b_{ijk},$$

and $|\mathbb{A}|$ will denote the induced norm of \mathbb{A} .

We say that $\mathbb{A} = (a_{ijk}) \in M^{N \times N \times N}$ is *antisymmetric* if

$$a_{ijk} = a_{jik} \quad \text{and} \quad \sum_p a_{ppk} = 0.$$

We write $\mathbb{A} \in M_D^{N \times N \times N}$.

The divergence of a $M^{N \times N \times N}$ -valued function $\mathbb{A}(x) = (a_{ijk}(x))$ is given by

$$(\text{div} \mathbb{A})_{ij} := \sum_k \frac{\partial a_{ijk}}{\partial x_k}.$$

Functional spaces and measures. Given $A \subseteq \mathbb{R}^N$ open and $1 \leq p < +\infty$, we will denote by $L^p(A; \mathbb{R}^M)$ the space of p -summable functions on A with values in \mathbb{R}^M and by $W^{1,p}(A; \mathbb{R}^M)$ the usual Sobolev space of functions in $L^p(A; \mathbb{R}^M)$ whose derivatives in the sense of distributions belong to L^p . Finally, $\mathcal{M}_b(A; \mathbb{R}^M)$ will denote the space of \mathbb{R}^M -valued Radon measures on A , and for every $\mu \in \mathcal{M}_b(A; \mathbb{R}^M)$ we will indicate by $|\mu|(A)$ its total mass. We set $\|\mu\|_{\mathcal{M}_b(A; \mathbb{R}^M)} := |\mu|(A)$. We refer the reader to [9] for the main properties concerning Sobolev spaces and Radon measures.

Let us recall some results from the theory of BV functions. We refer the reader to [2] for an exhaustive treatment of the subject.

We say that $u \in BV(A; \mathbb{R}^M)$ if $u \in L^1(A; \mathbb{R}^M)$, and its distributional derivative Du is a vector-valued Radon measure on A . $BV(A; \mathbb{R}^M)$ is a Banach space with respect to the norm

$$\|u\|_{BV(A; \mathbb{R}^M)} := \|u\|_{L^1(A; \mathbb{R}^M)} + |Du|(A).$$

We will denote by $D^s u$ the singular part of Du with respect to the Lebesgue measure \mathcal{L}^N and by ∇u the density of its absolutely continuous part.

We will say that a sequence $(u_n)_{n \in \mathbb{N}}$ in $BV(A; \mathbb{R}^M)$ converges weakly* in $BV(A; \mathbb{R}^M)$ to $u \in BV(A; \mathbb{R}^M)$ if

$$(2.3) \quad \begin{aligned} u_n &\rightarrow u && \text{strongly in } L^1(A; \mathbb{R}^M), \\ Du_n &\overset{*}{\rightharpoonup} Du && \text{weakly* in } \mathcal{M}_b(A; \mathbb{R}^M). \end{aligned}$$

The following compactness result holds: If A is open bounded and with a Lipschitz boundary, every bounded sequence in $BV(A; \mathbb{R}^M)$ admits a subsequence converging weakly* in $BV(A; \mathbb{R}^M)$.

Finally we will use throughout the paper the following embedding property of BV : If A is bounded and with a Lipschitz boundary, then $BV(A; \mathbb{R}^M)$ is continuously embedded into $L^q(A; \mathbb{R}^M)$ for every $1 \leq q \leq \frac{N}{N-1}$, the embedding being compact for every $1 \leq q < \frac{N}{N-1}$.

One-dimensional AC and BV functions with values in Banach spaces.

Let X be a reflexive Banach space or the dual of a separable Banach space. We denote by $BV(a, b; X)$ and $AC(a, b; X)$ the space of functions of bounded variation and the space of absolutely continuous functions from $[a, b]$ to X , respectively. We refer the reader to [6] for the main properties of these spaces. We recall that the variation of $f \in BV(a, b; X)$ is defined as

$$(2.4) \quad \mathcal{V}(f; a, b) := \sup \left\{ \sum_{j=1}^k \|f(t_j) - f(t_{j-1})\|_X : a = t_0 < t_1 < \dots < t_k = b \right\}.$$

If X is reflexive and $f \in AC(a, b; X)$, then the time derivative $\dot{f}(t)$ exists for a.e. $t \in [a, b]$. If X is the dual of a separable Banach space (and this is interesting when we consider the plastic strains), the time derivative $\dot{f}(t)$ exists as a weak* limit of difference quotients for a.e. $t \in [a, b]$ (see [8, Theorem 7.1]).

We will often use the following generalization of Helly’s theorem [8, Lemma 7.2] (see also [23, Theorem 3.2]): If X is the dual of a separable Banach space and $(f_k)_{k \in \mathbb{N}}$ a sequence in $BV(a, b; X)$ with $\mathcal{V}(f_k; a, b)$ and $\|f_k(a)\|_X$ uniformly bounded, then there exist $f \in BV(a, b; X)$ and a subsequence $(f_{k_j})_{j \in \mathbb{N}}$ such that $f_{k_j}(t) \overset{*}{\rightharpoonup} f(t)$ weakly* in X for every $t \in [a, b]$.

3. The Gurtin–Anand model. In this section we quickly describe the Gurtin–Anand model [19] in strain gradient plasticity which describes the behavior of isotropic, plastically irrotational materials under small deformations. We present the case in which the internal hardening variable is neglected.

Let $\Omega \subseteq \mathbb{R}^N$ be the reference configuration of the body. The starting point of the theory is, as usual, the additive decomposition of the displacement strain $\mathbf{E}u = (\nabla u + \nabla u^T)/2$ into elastic and plastic parts

$$(3.1) \quad \mathbf{E}u = \mathbf{E}^e + \mathbf{E}^p.$$

The symmetric matrices \mathbf{E}^e and \mathbf{E}^p are referred to as the elastic and the plastic part, respectively. The plastic part \mathbf{E}^p is supposed to be unable to sustain volumetric changes, so that

$$\text{tr} \mathbf{E}^p = 0;$$

that is, $\mathbf{E}^p \in M_D^{N \times N}$.

Stresses and balance equations. Given a subbody $\mathcal{B} \subseteq \Omega$, besides the usual Cauchy stress $\mathbf{T} \in \mathbb{M}_{\text{sym}}^{N \times N}$ conjugate to \mathbf{E}^e , the analysis of its equilibrium involves also stresses $\mathbf{T}^p \in \mathbb{M}_D^{N \times N}$ and $\mathbb{K}^p \in \mathbb{M}_D^{N \times N \times N}$ conjugate to \mathbf{E}^p and $\nabla \mathbf{E}^p$, respectively. Given the rate-like kinematical descriptors $(\dot{u}, \dot{\mathbf{E}}^e, \dot{\mathbf{E}}^p)$, the power expenditure within \mathcal{B} is given by

$$\mathcal{W}_{\text{int}}(\mathcal{B}) = \int_{\mathcal{B}} (\mathbf{T} : \dot{\mathbf{E}}^e + \mathbf{T}^p : \dot{\mathbf{E}}^p + \mathbb{K}^p : \nabla \dot{\mathbf{E}}^p) dx.$$

$\mathcal{W}_{\text{int}}(\mathcal{B})$ is balanced by the power expenditure of external forces

$$\mathcal{W}_{\text{ext}}(\mathcal{B}) = \int_{\partial \mathcal{B}} (t(\nu) \cdot \dot{u} + \mathbf{K}(\nu) : \dot{\mathbf{E}}^p) dA + \int_{\mathcal{B}} f \cdot \dot{u} dV,$$

where f is the external body force and $t(\nu)$ is the boundary traction (ν is the outward normal to \mathcal{B}) which are associated as usual to \dot{u} , while $\mathbf{K}(\nu) \in \mathbb{M}_D^{N \times N}$ is associated to the plastic strain rate $\dot{\mathbf{E}}^p$. The balance of power expenditures (that is, $\mathcal{W}_{\text{int}}(\mathcal{B}) = \mathcal{W}_{\text{ext}}(\mathcal{B})$ for every subbody \mathcal{B}) leads to the equalities $t(\nu) = \mathbf{T}\nu$ and $\mathbf{K}(\nu) = \mathbb{K}^p \nu$ for the traction and microtraction and to the equilibrium equations

$$-\text{div} \mathbf{T} = f \quad \text{and} \quad \mathbf{T}^p = \mathbf{T}_D + \text{div} \mathbb{K}^p \quad \text{in } \Omega,$$

where \mathbf{T}_D is the deviatoric part of \mathbf{T} as defined in (2.2). These equations are supplemented by boundary conditions for \mathbf{T} and \mathbb{K}^p . If traction forces g are present on a part $\partial_N \Omega$ of the boundary of Ω , we have as usual

$$\mathbf{T}\nu = g \quad \text{on } \partial_N \Omega.$$

Concerning \mathbb{K}^p , assuming $\mathbb{K}^p \nu = 0$ at the boundary (see [19, section 8]), we are led to the condition

$$\mathbb{K}^p \nu = 0 \quad \text{on } \partial \Omega.$$

The free energy. The free energy ψ is assumed to depend on \mathbf{E}^e and $\text{curl} \mathbf{E}^p$: In the quadratic separable case ψ has the form

$$(3.2) \quad \psi = \frac{1}{2} \mathbb{C} \mathbf{E}^e : \mathbf{E}^e + \frac{\mu L^2}{2} |\text{curl} \mathbf{E}^p|^2,$$

where \mathbb{C} is the elastic tensor

$$(3.3) \quad \mathbb{C} \mathbf{E}^e := 2\mu \mathbf{E}_D^e + k(\text{tr} \mathbf{E}^e) \mathbf{I},$$

with μ and k the elastic shear and bulk moduli, respectively. The constant $L > 0$ is an elastic length scale. The elastic energy \mathbb{K}_{en}^p is then defined so that the identity

$$(3.4) \quad \mu L^2 \text{curl} \mathbf{E}^p : \text{curl} \mathbf{A} = \mathbb{K}_{\text{en}}^p : \nabla \mathbf{A}$$

holds for every $\mathbb{M}^{N \times N}$ -valued function \mathbf{A} . In components we have

$$(3.5) \quad (\mathbb{K}_{\text{en}}^p)_{jqp} := \mu L^2 \left[\frac{\partial \mathbf{E}_{jq}^p}{\partial x_p} - \frac{1}{2} \left(\frac{\partial \mathbf{E}_{jp}^p}{\partial x_q} + \frac{\partial \mathbf{E}_{qp}^p}{\partial x_j} \right) + \frac{1}{N} \delta_{jq} \sum_r \frac{\partial \mathbf{E}_{rp}^p}{\partial x_r} \right],$$

where δ_{jq} is the usual Kröner symbol. The stress \mathbb{K}^p is then additively decomposed in the following way:

$$\mathbb{K}^p = \mathbb{K}_{\text{diss}}^p + \mathbb{K}_{\text{en}}^p.$$

Admissibility of the stresses and the flow rule. Neglecting the hardening internal variable, i.e., if we are in the case without hardening nor softening, the admissibility for the stresses involved in the description of the behavior of Ω reads

$$(3.6) \quad \sqrt{|\mathbf{T}^p(x)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^p(x)|^2} \leq S_Y,$$

where $l > 0$ is a material parameter, and S_Y is a yield stress.

Assume now that body and traction forces vary with time, i.e., $f = f(t)$ and $g = g(t)$. The flow rule which drives the system requires that if $(\mathbf{T}^p(t, x), \mathbb{K}_{\text{diss}}^p(t, x))$ is at the yield surface (that is, (3.6) holds with equality), then

$$(3.7) \quad \begin{cases} \mathbf{T}^p(t, x) = S_Y \frac{\dot{\mathbf{E}}^p(t, x)}{\sqrt{|\dot{\mathbf{E}}^p(t, x)|^2 + l^2|\nabla \dot{\mathbf{E}}^p(t, x)|^2}}, \\ \mathbb{K}_{\text{diss}}^p(t, x) = S_Y \frac{l^2 \nabla \dot{\mathbf{E}}^p(t, x)}{\sqrt{|\dot{\mathbf{E}}^p(t, x)|^2 + l^2|\nabla \dot{\mathbf{E}}^p(t, x)|^2}}. \end{cases}$$

Here $\dot{\mathbf{E}}^p(t, x)$ and $\nabla \dot{\mathbf{E}}^p(t, x)$ denote the time derivative of $\mathbf{E}^p(t, x)$ and $\nabla \mathbf{E}^p(t, x)$, respectively. If $(\mathbf{T}^p(t, x), \mathbb{K}_{\text{diss}}^p(t, x))$ is well inside the yield surface, then no plastic phenomenon occurs, i.e., $(\dot{\mathbf{E}}^p, \nabla \dot{\mathbf{E}}^p) = (0, 0)$. The flow rule (3.7) is a generalization of the von Mises flow rule in perfect plasticity (set $l = L = 0$, and note that $\mathbf{T}^p = \mathbf{T}_D$). It moreover implies that

$$\int_{\mathcal{B}} \dot{\psi} \, dV \leq \mathcal{W}_{\text{ext}}(\mathcal{B}).$$

The previous inequality reflects the thermodynamical requirement that the increase in free energy of \mathcal{B} is less than or equal to the power expended on \mathcal{B} .

4. Functional setting. In this section we state the precise mathematical framework that we adopt to study quasi-static evolutions for the Gurtin–Anand model.

The reference configuration. Let the reference configuration be given by $\Omega \subseteq \mathbb{R}^N$, $N \geq 2$, a bounded open set with a Lipschitz boundary. Let $\partial\Omega$ be partitioned into two open (in the relative topology) disjoint sets $\partial_D\Omega$ and $\partial_N\Omega$ with the same boundary Γ such that $\mathcal{H}^{N-2}(\Gamma) < +\infty$.

Admissible configurations. Let the prescribed boundary displacement on $\partial_D\Omega$ be given by (the trace of) a Sobolev function $w \in W^{1,2}(\Omega; \mathbb{R}^N)$. By an admissible configuration relative to the boundary datum w , we will understand a triple $(u, \mathbf{E}^e, \mathbf{E}^p)$ such that

$$(4.1) \quad u \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N), \quad \mathbf{E}^e \in L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N}), \quad \mathbf{E}^p \in BV(\Omega; \mathbb{M}_D^{N \times N}),$$

with

$$(4.2) \quad u = w \quad \text{on } \partial_D\Omega,$$

$$(4.3) \quad \mathbf{E}u = \mathbf{E}^e + \mathbf{E}^p,$$

and

$$(4.4) \quad \text{curl} \mathbf{E}^p \in L^2(\Omega; \mathbb{M}^{N \times N}).$$

Equality (4.2) is intended in the sense of traces. Notice that, by the embedding properties of BV , (4.3) entails $\mathbf{E}u \in L^{\frac{N}{N-1}}(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$; the requirement $u \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ is then consistent with the regularity implied by Korn's inequality in view of the boundary condition (4.2). Let us denote by $\mathcal{A}(w)$ the family of admissible configurations for the boundary datum w , i.e.,

$$(4.5) \quad \mathcal{A}(w) := \{(u, \mathbf{E}^e, \mathbf{E}^p) \text{ such that (4.1)–(4.4) are satisfied}\}.$$

Notice that $\mathcal{A}(w) \neq \emptyset$ since $(w, \mathbf{E}w, 0) \in \mathcal{A}(w)$: Moreover $\mathcal{A}(w)$ contains the triples $(v, \mathbf{E}v, 0)$, with $v \in W^{1,2}(\Omega; \mathbb{R}^N)$ and $v = w$ on $\partial_D \Omega$, so that rigid deformations are admissible (and this motivates the choice $w \in W^{1,2}(\Omega; \mathbb{R}^N)$).

The free energy. The free energy of the configuration $(u, \mathbf{E}^e, \mathbf{E}^p) \in \mathcal{A}(w)$ is given according to (3.2) by

$$\Psi(\mathbf{E}^e, \text{curl} \mathbf{E}^p) := \mathcal{Q}_1(\mathbf{E}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p),$$

where

$$(4.6) \quad \mathcal{Q}_1(\mathbf{E}^e) := \frac{1}{2} \int_{\Omega} \mathbb{C} \mathbf{E}^e : \mathbf{E}^e \, dx$$

and

$$(4.7) \quad \mathcal{Q}_2(\text{curl} \mathbf{E}^p) := \frac{1}{2} \mu L^2 \int_{\Omega} |\text{curl} \mathbf{E}^p|^2 \, dx.$$

Here \mathbb{C} denotes the elasticity tensor (3.3), and $\mu > 0$ is the elastic shear modulus: We assume also that $k > 0$, so that there exist $0 < \alpha_{\mathbb{C}} \leq \beta_{\mathbb{C}} < +\infty$ such that for every $\mathbf{A} \in \mathbb{M}_{\text{sym}}^{N \times N}$ we have

$$(4.8) \quad \alpha_{\mathbb{C}} |\mathbf{A}|^2 \leq \mathbb{C} \mathbf{A} : \mathbf{A} \leq \beta_{\mathbb{C}} |\mathbf{A}|^2.$$

The yield function \mathcal{H} . In order to get variationally the constraint for the stresses conjugated to the plastic variables according to (3.6), we are led to consider the relaxation

$$(4.9) \quad \mathcal{H}(\mathbf{E}^p) := S_Y \int_{\Omega} \sqrt{|\mathbf{E}^p|^2 + l^2 |\nabla \mathbf{E}^p|^2} \, dx + l S_Y |D^s \mathbf{E}^p|(\Omega)$$

defined for every $\mathbf{E}^p \in BV(\Omega; \mathbb{M}_D^{N \times N})$. Simple arguments on subadditive and positively one-homogeneous functions on measures (see [13]) show that \mathcal{H} is the relaxation under the L^1 -norm of the map

$$\mathbf{E}^p \mapsto S_Y \int_{\Omega} \sqrt{|\mathbf{E}^p|^2 + l^2 |\nabla \mathbf{E}^p|^2} \, dx$$

defined for a regular plastic strain \mathbf{E}^p , which is connected to the effective flow rate proposed by Gurtin and Anand (see [19, section 6.3]). As a consequence, \mathcal{H} turns out to be naturally involved in an analysis which employs direct methods of the calculus of variations.

We will often use the lower semicontinuity of \mathcal{H} along weakly* converging sequences, which is a direct consequence of the relaxation process through which \mathcal{H} is obtained.

LEMMA 4.1. $(\mathbf{E}_n^p)_{n \in \mathbb{N}} \subset BV(\Omega; M_D^{N \times N})$

$$\mathbf{E}_n^p \xrightarrow{*} \mathbf{E}^p \quad \text{in } BV(\Omega; M_D^{N \times N})$$

if and only if $\mathbf{E}^p \in BV(\Omega; M_D^{N \times N})$

$$\mathcal{H}(\mathbf{E}^p) \leq \liminf_{n \rightarrow +\infty} \mathcal{H}(\mathbf{E}_n^p).$$

Prescribed boundary displacements and body/traction forces. We assume that the prescribed boundary displacement on $\partial_D \Omega$ is given by (the trace of) a function $w(t, x)$ which is absolutely continuous in time with values in the Sobolev space $W^{1,2}(\Omega; \mathbb{R}^N)$, i.e.,

$$(4.10) \quad w \in AC(0, T; W^{1,2}(\Omega; \mathbb{R}^N)).$$

Moreover we assume that the prescribed body forces in Ω and traction forces on $\partial_N \Omega$ are given by

$$(4.11) \quad f \in AC(0, T; L^N(\Omega; \mathbb{R}^N)) \quad \text{and} \quad g \in AC(0, T; L^N(\partial_N \Omega; \mathbb{R}^N)).$$

For every $t \in [0, T]$ let us consider $\mathcal{L}(t) : W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N) \rightarrow \mathbb{R}$ given by

$$(4.12) \quad \langle \mathcal{L}(t), u \rangle := \int_{\Omega} f(t) \cdot u \, dx + \int_{\partial_N \Omega} g(t) \cdot u \, d\mathcal{H}^{N-1}.$$

Here \mathcal{H}^{N-1} denotes the $(N-1)$ -dimensional Hausdorff measure, which is a generalization to arbitrary sets of the usual surface measure (see [9]). By means of the Sobolev embedding theorem it is easily seen that $\mathcal{L}(t)$ is a continuous linear functional on $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$.

Throughout the paper we will assume that the prescribed body and traction forces satisfy the following uniform condition: We assume that for every $t \in [0, T]$ there exists $\rho(t) \in L^N(\Omega; M_{\text{sym}}^{N \times N})$ such that

$$(4.13) \quad \begin{cases} -\operatorname{div} \rho(t) = f(t) & \text{in } \Omega, \\ \rho(t) \nu = g(t) & \text{on } \partial_N \Omega \end{cases}$$

and there exists $\alpha > 0$ such that for every $\mathbf{A} \in M_D^{N \times N}$, with $|\mathbf{A}| \leq \alpha$, we have

$$(4.14) \quad |\mathbf{A} + \rho_D(t)| \leq S_Y \quad \text{a.e. in } \Omega.$$

Moreover we assume that $t \mapsto \rho(t)$ and $t \mapsto \rho_D(t)$ are absolutely continuous from $[0, T]$ to $L^2(\Omega; M_{\text{sym}}^{N \times N})$ and $L^\infty(\Omega; M_D^{N \times N})$, respectively. Notice that the trace condition in (4.13) is well-defined in the dual of the traces on $\partial_N \Omega$ of $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ since ρ is an L^N -field with divergence in L^N . Moreover, for every $(u, \mathbf{E}^e, \mathbf{E}^p) \in \mathcal{A}(w)$ we have the following representation formula for $\mathcal{L}(t)$ (here we use $\mathcal{H}^{N-2}(\Gamma) < +\infty$):

$$(4.15) \quad \langle \mathcal{L}(t), u \rangle = -\langle \rho(t) \nu, w \rangle_{\partial_D \Omega} + \int_{\Omega} \rho(t) : \mathbf{E}^e \, dx + \int_{\Omega} \rho_D(t) : \mathbf{E}^p \, dx,$$

where the first term on the right-hand side should be interpreted as the pairing between $H^{-1/2}(\partial_D \Omega; \mathbb{R}^N)$ and $H^{1/2}(\partial_D \Omega; \mathbb{R}^N)$.

4.2. Notice that for $\mathcal{L}(t)$ to be well-defined in the dual of $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ it suffices to require $f(t) \in L^{N/2}(\Omega; \mathbb{R}^N)$ (assume that $N \geq 3$, the case $N = 2$ being different in view of Sobolev embedding). But in view of the safe load condition (4.13)–(4.14), $\rho(t)$ would be only an element of $L^{N/2}$ with divergence in $L^{N/2}$, so that its normal trace would be defined in the dual of the traces on $\partial\Omega$ of $W^{1, \frac{N}{N-2}}(\Omega; \mathbb{R}^N)$. Then the representation formula (4.15) would no longer be well-defined (since $w \in W^{1,2}(\Omega; \mathbb{R}^N)$).

As a consequence of the safe load condition, we have the following coercivity estimate for \mathcal{H} .

LEMMA 4.3. $\mathbf{E}^P \in BV(\Omega; M_D^{N \times N})$.

$$(4.16) \quad \mathcal{H}(\mathbf{E}^P) - \int_{\Omega} \rho_D(t) : \mathbf{E}^P \, dx \geq \frac{\alpha}{2} \|\mathbf{E}^P\|_{L^1(\Omega; M_D^{N \times N})} + \min \left\{ l \frac{\alpha}{2}, l S_Y \right\} \|D\mathbf{E}^P\|_{\mathcal{M}_b(\Omega; M_D^{N \times N \times N})}.$$

where $\alpha > 0$.

$$(4.17) \quad \mathcal{H}(\mathbf{E}^P) - \int_{\Omega} \rho_D(t) : \mathbf{E}^P \, dx \geq \alpha_l \|\mathbf{E}^P\|_{BV(\Omega; M_D^{N \times N})}.$$

Notice that by Hölder inequality we have

$$S_Y \int_{\Omega} \sqrt{|\mathbf{E}^P|^2 + l^2 |\nabla \mathbf{E}^P|^2} \, dx \geq \sup_{(\tau_1, \tau_2) \in \mathcal{K}} \int_{\Omega} [\tau_1 : \mathbf{E}^P + \tau_2 : \nabla \mathbf{E}^P] \, dx,$$

where

$$\mathcal{K} := \left\{ (\tau_1, \tau_2) \in L^\infty(\Omega; M_D^{N \times N}) \times L^\infty(\Omega; M_D^{N \times N \times N}) : \sqrt{|\tau_1|^2 + l^{-2} |\tau_2|^2} \leq S_Y \text{ a.e. in } \Omega \right\}.$$

We deduce that for every $(\tau_1, \tau_2) \in \mathcal{K}$

$$\mathcal{H}(\mathbf{E}^P) - \int_{\Omega} \rho_D(t) : \mathbf{E}^P \, dx \geq \int_{\Omega} [(\tau_1 - \rho_D(t)) : \mathbf{E}^P + \tau_2 : \nabla \mathbf{E}^P] \, dx + l S_Y |D^s \mathbf{E}^P|(\Omega),$$

so that in view of (4.14) we get

$$\mathcal{H}(\mathbf{E}^P) - \int_{\Omega} \rho_D(t) : \mathbf{E}^P \, dx \geq \int_{\Omega} [\tilde{\tau}_1 : \mathbf{E}^P + \tilde{\tau}_2 : \nabla \mathbf{E}^P] \, dx + l S_Y |D^s \mathbf{E}^P|(\Omega)$$

for every $\|\tilde{\tau}_1\|_{L^\infty(\Omega; M_D^{N \times N})} \leq \frac{\alpha}{2}$ and $\|\tilde{\tau}_2\|_{L^\infty(\Omega; M_D^{N \times N \times N})} \leq l \frac{\alpha}{2}$. We conclude that

$$\begin{aligned} \mathcal{H}(\mathbf{E}^P) - \int_{\Omega} \rho_D(t) : \mathbf{E}^P \, dx &\geq \frac{\alpha}{2} \|\mathbf{E}^P\|_{L^1(\Omega; M_D^{N \times N})} + l \frac{\alpha}{2} \|\nabla \mathbf{E}^P\|_{L^1(\Omega; M_D^{N \times N \times N})} + l S_Y |D^s \mathbf{E}^P|(\Omega), \end{aligned}$$

so that (4.16) holds. Inequality (4.17) follows by choosing $\alpha_l := \min\{\frac{\alpha}{2}, l \frac{\alpha}{2}, l S_Y\}$. \square

4.4. As mentioned in the introduction, the safe load condition entails in view of Lemma 4.3 some coercivity in BV for \mathbf{E}^P in the step-by-step minimization problems which we use in section 6 in order to construct a quasi-static evolution (see Lemma 6.1). Those problems contain two terms (dissipation and work of external loads) which have linear growth and which can compete, leading, in general, to a loss of compactness: This cannot happen if a safe load condition is assumed.

5. The main results. Let $T > 0$, and let w , f , and g be the prescribed boundary displacements, body forces, and traction forces according to (4.10) and (4.11), respectively. We assume that f and g satisfy the uniform safe load condition (4.13)–(4.14).

We will denote by $\dot{w}(t)$, $\dot{f}(t)$, and $\dot{g}(t)$ the derivative at time $t \in [0, T]$ of w , f , and g , respectively. Notice that these derivatives exist for a.e. $t \in [0, T]$ since the maps are absolutely continuous with values in a reflexive Banach space. We will denote by $\dot{\mathcal{L}}(t)$ the external work associated to $\dot{f}(t)$ and $\dot{g}(t)$.

Given \mathcal{H} as in (4.9), the \mathcal{H} -variation on $[a, b] \subseteq [0, T]$ of $t \mapsto \mathbf{E}^P(t)$ is defined as

$$(5.1) \quad \mathcal{D}_{\mathcal{H}}(\mathbf{E}^P; a, b) := \sup \left\{ \sum_{j=1}^k \mathcal{H}(\mathbf{E}^P(t_j) - \mathbf{E}^P(t_{j-1})) : a = t_0 < t_1 < \dots < t_k = b \right\}.$$

The notion of quasi-static evolution for the Gurtin–Anand model is the following.

DEFINITION 5.1 (quasi-static evolution).

$$t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^P(t))$$

$[0, T], W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N) \times L^2(\Omega; \mathbf{M}_{\text{sym}}^{N \times N}) \times BV(\Omega; \mathbf{M}_D^{N \times N}), (u(t), \mathbf{E}^e(t), \mathbf{E}^P(t)) \in \mathcal{A}(w(t))$

(a) Global stability $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$

$$(5.2) \quad \begin{aligned} \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}^P(t)) - \langle \mathcal{L}(t), u(t) \rangle \\ \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{p}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^P(t)); \end{aligned}$$

(b) Energy balance $t \mapsto \mathbf{E}^P(t) \in BV(\Omega; \mathbf{M}_D^{N \times N}), [0, T]$

$$(5.3) \quad \begin{aligned} \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^P; 0, t) = \mathcal{E}(0) + \int_0^t \int_{\Omega} \mathbf{T}(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau \\ - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau, \end{aligned}$$

$$\mathbf{T}(t) := \mathbb{C} \mathbf{E}^e(t).$$

$$(5.4) \quad \mathcal{E}(t) := \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}^P(t)) - \langle \mathcal{L}(t), u(t) \rangle,$$

$$\mathcal{D}_{\mathcal{H}}(\mathbf{E}^P; 0, t) = \mathcal{H}_t \quad (5.1)$$

Our first main result is the following existence theorem.

THEOREM 5.2. $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^P) \in \mathcal{A}(w(0))$

$$\mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^P) - \langle \mathcal{L}(0), u_0 \rangle \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{p}) - \langle \mathcal{L}(0), v \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}_0^P)$$

$(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(0)) \quad t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^P(t)) \quad (u(0), \mathbf{E}^e(0), \mathbf{E}^P(0)) = (u_0, \mathbf{E}_0^e, \mathbf{E}_0^P)$

Theorem 5.2 will be proved in section 7 by exploiting the convergence of a discrete-in-time evolution constructed through variational arguments in section 6.

Our second main result shows that a quasi-static evolution satisfies the required constitutive equations, balance equations, and the flow rule of the Gurtin–Anand model.

THEOREM 5.3. *Let $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ be a solution of the problem*

$$[0, T] \times W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N), L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N}), BV(\Omega; \mathbb{M}_D^{N \times N}), L^2(\Omega; \mathbb{M}^{N \times N}), t \mapsto \mathbf{E}^e(t), t \mapsto \mathbf{E}^p(t), t \mapsto \text{curl} \mathbf{E}^p(t)$$

with $t \in [0, T]$

(a) *Cauchy stress $\mathbf{T}(t) = \mathbb{C}\mathbf{E}^e(t)$ satisfies*

$$(5.5) \quad \begin{cases} -\text{div} \mathbf{T}(t) = f(t) & \text{in } \Omega, \\ \mathbf{T}(t)\nu = g(t) & \text{on } \partial_N \Omega. \end{cases}$$

(b) *Stresses conjugated to the plastic variables $\mathbf{T}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N \times N})$, $\mathbb{K}^p(t) \in L^2(\Omega; \mathbb{M}_D^{N \times N \times N})$, $\mathbb{K}_{\text{diss}}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N \times N})$, $\mathbb{S}^p(t) \in (\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}))^*$, $f_i, \mathbb{K}_{\text{en}}^p(t)$ (3.5), $\mathbf{E}^p(t)$, $\mathbf{T}_D(t) := (\mathbf{T}(t))_D$*

$$(5.6) \quad \begin{cases} \mathbb{K}^p(t) = \mathbb{K}_{\text{en}}^p(t) + \mathbb{K}_{\text{diss}}^p(t) & \text{in } \Omega, \\ \mathbf{T}^p(t) = \mathbf{T}_D(t) + \text{div} \mathbb{K}^p(t) & \text{in } \Omega, \\ \mathbb{K}^p(t)\nu = 0 & \text{on } \partial \Omega, \end{cases}$$

$$(5.7) \quad \begin{aligned} \sqrt{|\mathbf{T}^p(t)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^p(t)|^2} &\leq S_Y & \text{in } \Omega, \\ \|\mathbb{S}^p(t)\|_{(\mathcal{M}_b(\mathbb{M}_D^{N \times N \times N}))^*} &\leq lS_Y, \end{aligned}$$

where $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$

$$\int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \int_{\Omega} \mathbf{T}^p(t) : \mathbf{p} \, dx + \int_{\Omega} \mathbb{K}^p(t) : \nabla \mathbf{p} \, dx + \langle \mathbb{S}^p(t), D^s \mathbf{p} \rangle = \langle \mathcal{L}(t), v \rangle.$$

(c) *The flow rule $\dot{\mathbf{E}}^p(t)$ in $x \in \Omega$, $\nabla \dot{\mathbf{E}}^p(t)$, $\mathbf{T}^p(t)$, $\mathbb{K}_{\text{diss}}^p(t)$, f_i (3.7), f_i*

Notice that the normal trace which appears in (5.5) is well-defined in $H^{-1/2}(\partial \Omega; \mathbb{R}^N)$ since $\mathbf{T}(t)$ is an L^2 -field with divergence in L^2 . Similarly, the normal trace in (5.6) is well-defined in $H^{-1/2}(\partial \Omega; \mathbb{R}^{N \times N})$ because $\mathbb{K}^p(t)$ is an L^2 -field (by the definition of $\mathbb{K}_{\text{en}}^p(t)$ and by the constraint (5.7) for $\mathbb{K}_{\text{diss}}^p(t)$) with divergence in L^2 (by the balance equation (5.6) and by the constraint (5.7) for $\mathbf{T}^p(t)$).

In section 9 we will analyze the behavior of a quasi-static evolution as the length scales l and L go to zero, i.e., when the strain gradient effects vanish. We will prove (Theorem 9.2) that the quasi-static evolution converges to an evolution for perfect plasticity according to the framework recently proposed by Dal Maso, DeSimone, and Mora in [8].

6. The discrete-in-time evolution. In this section we construct a discretized-in-time evolution for the Gurtin–Anand model employing a step-by-step minimization procedure. The convergence of this approximated evolution to a quasi-static evolution for the Gurtin–Anand model as the time step discretization goes to zero will be proved in the next section.

Let $k \in \mathbb{N}$, $k \geq 1$, and let us set $t_k^i := \frac{i}{k}T$ for every $i = 0, 1, \dots, k$. Let us set

$$u_{k,0} := u_0, \quad \mathbf{E}_{k,0}^e := \mathbf{E}_0^e, \quad \mathbf{E}_{k,0}^p := \mathbf{E}_0^p,$$

where $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p) \in \mathcal{A}(w(0))$ is the initial configuration of the system given by Theorem 5.2.

Supposing to have constructed $(u_{k,i}, \mathbf{E}_{k,i}^e, \mathbf{E}_{k,i}^p) \in \mathcal{A}(w(t_k^i))$, let $(u_{k,i+1}, \mathbf{E}_{k,i+1}^e, \mathbf{E}_{k,i+1}^p) \in \mathcal{A}(w(t_k^{i+1}))$ ($i = 0, \dots, k - 1$) be a solution of the following minimization problem:

$$(6.1) \quad \min_{(u, \mathbf{E}^e, \mathbf{E}^p) \in \mathcal{A}(w(t_k^{i+1}))} \mathcal{Q}_1(\mathbf{E}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p) - \langle \mathcal{L}(t_k^{i+1}), u \rangle + \mathcal{H}(\mathbf{E}^p - \mathbf{E}_{k,i}^p).$$

The existence of a solution for problem (6.1) is established in the following lemma.

LEMMA 6.1. *Let (6.1) hold. Then there exists a solution to (6.1).*

The result follows by applying the direct method of the calculus of variations. In fact, let

$$(u_n, \mathbf{E}_n^e, \mathbf{E}_n^p) \in \mathcal{A}(w(t_k^{i+1}))$$

be a minimizing sequence for (6.1). By comparison with $(w(t_k^{i+1}), \mathbf{E}w(t_k^{i+1}), 0)$ we get

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_n^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_n^p) - \langle \mathcal{L}(t_k^{i+1}), u_n \rangle + \mathcal{H}(\mathbf{E}_n^p - \mathbf{E}_{k,i}^p) \\ \leq \mathcal{Q}_1(\mathbf{E}w(t_k^{i+1})) - \langle \mathcal{L}(t_k^{i+1}), w(t_k^{i+1}) \rangle + \mathcal{H}(\mathbf{E}_{k,i}^p) := C. \end{aligned}$$

By the representation formula (4.15) for $\mathcal{L}(t_k^{i+1})$ we deduce that

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_n^e) - \int_{\Omega} \rho(t_k^{i+1}) : \mathbf{E}_n^e \, dx + \mathcal{Q}_2(\text{curl} \mathbf{E}_n^p) \\ + \mathcal{H}(\mathbf{E}_n^p - \mathbf{E}_{k,i}^p) - \int_{\Omega} \rho_D(t_k^{i+1}) : (\mathbf{E}_n^p - \mathbf{E}_{k,i}^p) \, dx \\ \leq C + \int_{\Omega} \rho_D(t_k^{i+1}) : \mathbf{E}_{k,i}^p \, dx - \langle \rho(t_k^{i+1}) \nu, w(t_k^{i+1}) \rangle_{\partial_D \Omega}. \end{aligned}$$

By the coercivity of \mathcal{Q}_1 and \mathcal{Q}_2 in L^2 and by (4.17) we get

$$\|\mathbf{E}_n^e\|_{L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})}^2 + \|\text{curl} \mathbf{E}_n^p\|_{L^2(\Omega; \mathbb{M}^{N \times N})}^2 + \|\mathbf{E}_n^p - \mathbf{E}_{k,i}^p\|_{BV(\Omega; \mathbb{M}_D^{N \times N})} \leq C_1$$

for some $C_1 > 0$. Up to a subsequence we may assume that

$$\mathbf{E}_n^e \rightharpoonup \mathbf{E}^e \quad \text{weakly in } L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$$

and

$$\mathbf{E}_n^p \overset{*}{\rightharpoonup} \mathbf{E}^p \quad \text{weakly* in } BV(\Omega; \mathbb{M}_D^{N \times N}).$$

As a consequence we get $\text{curl} \mathbf{E}^p \in L^2(\Omega; \mathbb{M}^{N \times N})$ and that

$$\text{curl} \mathbf{E}_n^p \rightharpoonup \text{curl} \mathbf{E}^p \quad \text{weakly in } L^2(\Omega; \mathbb{M}^{N \times N}).$$

By the compatibility $\mathbf{E}u_n = \mathbf{E}_n^e + \mathbf{E}_n^p$ and by the embedding $BV(\Omega; \mathbb{M}_D^{N \times N}) \hookrightarrow L^{\frac{N}{N-1}}(\Omega; \mathbb{M}_D^{N \times N})$, we get that $(\mathbf{E}u_n)_{n \in \mathbb{N}}$ is bounded in $L^{\frac{N}{N-1}}(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$. In view of the boundary condition $u_n = w_{k,i+1}$ on $\partial_D \Omega$, Korn's inequality implies that $(u_n)_{n \in \mathbb{N}}$ is bounded in $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$. Up to a further subsequence we can thus suppose that

$$u_n \rightharpoonup u \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$$

for some $u \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$, with $u = w(t_k^{i+1})$ on $\partial_D \Omega$. We clearly have $(u, \mathbf{E}^e, \mathbf{E}^p) \in \mathcal{A}(w(t_k^{i+1}))$, and by lower semicontinuity (\mathcal{Q}_1 and \mathcal{Q}_2 are quadratic, $\mathcal{L}(t_k^{i+1})$ is linear, and \mathcal{H} is lower semicontinuous by Lemma 4.1), we deduce that

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p) - \langle \mathcal{L}(t_k^{i+1}), u \rangle + \mathcal{H}(\mathbf{E}^p - \mathbf{E}_{k,i}^p) \\ & \leq \liminf_{n \rightarrow +\infty} \left[\mathcal{Q}_1(\mathbf{E}_n^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_n^p) - \langle \mathcal{L}(t_k^{i+1}), u_n \rangle + \mathcal{H}(\mathbf{E}_n^p - \mathbf{E}_{k,i}^p) \right]. \end{aligned}$$

We conclude that $(u, \mathbf{E}^e, \mathbf{E}^p)$ is a minimizer for problem (6.1), so that the proof is concluded. \square

The discretized-in-time evolution is obtained by interpolating the data obtained by the minimization procedure described above. Let us set for $t_k^i \leq t < t_k^{i+1}$

$$w_k(t) := w(t_k^i) \quad \text{and} \quad \mathcal{L}_k(t) := \mathcal{L}(t_k^i).$$

We collect the main properties of the discretized-in-time evolution (essential for the passage to the limit as the time step discretization goes to zero) in the following proposition.

PROPOSITION 6.2. $t \mapsto (u_k(t), \mathbf{E}_k^e(t), \mathbf{E}_k^p(t)), \quad t \in [0, T]$.

- (a) $(u_k(0), \mathbf{E}_k^e(0), \mathbf{E}_k^p(0)) = (u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$
- (a) $(u_k(t), \mathbf{E}_k^e(t), \mathbf{E}_k^p(t)) \in \mathcal{A}(w_k(t)), \quad t \in [0, T]$

$$\begin{aligned} (6.2) \quad & \mathcal{Q}_1(\mathbf{E}_k^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}_k^p(t)) - \langle \mathcal{L}_k(t), u_k(t) \rangle \\ & \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{p}) - \langle \mathcal{L}_k(t), v \rangle + \mathcal{H}(\mathbf{E}_k^p(t) - \mathbf{p}). \end{aligned}$$

- (b)

$$\mathcal{E}_k(t) := \mathcal{Q}_1(\mathbf{E}_k^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}_k^p(t)) - \langle \mathcal{L}_k(t), u_k(t) \rangle,$$

$$t_k^i \leq t < t_k^{i+1}$$

$$\begin{aligned} (6.3) \quad & \mathcal{E}_k(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_k^p; 0, t) \leq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle \\ & + \int_0^{t_k^i} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^{t_k^i} \langle \dot{\mathcal{L}}(\tau), u_k(\tau) \rangle \, d\tau \\ & - \int_0^{t_k^i} \langle \mathcal{L}_k(\tau), \dot{w}(\tau) \rangle \, d\tau + e_k, \end{aligned}$$

- (c) $e_k \rightarrow 0, \quad k \rightarrow +\infty, \quad \mathcal{D}_{\mathcal{H}} \dots f_i \dots (5.1)$
- (c) $C \dots k \in \mathbb{N}, \quad t \in [0, T]$

$$\begin{aligned} (6.4) \quad & \|u_k(t)\|_{W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)} + \|\mathbf{E}_k^e(t)\|_{L^2(\Omega; M_{\text{sym}}^{N \times N})} \\ & + \|\text{curl} \mathbf{E}_k^p(t)\|_{L^2(\Omega; M^{N \times N})} + \mathcal{V}(\mathbf{E}_k^p; 0, t) \leq C, \end{aligned}$$

$$\mathcal{V}(\mathbf{E}_k^p; 0, t) \dots \mathbf{E}_k^p \dots [0, t] \dots f_i \dots (2.4)$$

For every $t_k^i \leq t < t_k^{i+1}$ let us set

$$u_k(t) := u_{k,i}, \quad \mathbf{E}_k^e(t) := \mathbf{E}_{k,i}^e, \quad \text{and} \quad \mathbf{E}_k^p(t) := \mathbf{E}_{k,i}^p,$$

where $(u_{k,j}, \mathbf{E}_{k,j}^e, \mathbf{E}_{k,j}^p) \in \mathcal{A}(w(t_k^j))$ is a solution of the minimization problem (6.1). The minimality property (6.2) follows immediately by the subadditivity of \mathcal{H} .

Let us prove (6.3). By construction, by comparing $(u_{k,j}, \mathbf{E}_{k,j}^e, \mathbf{E}_{k,j}^p)$ with

$$(u_{k,j-1} + w(t_k^j) - w(t_k^{j-1}), \mathbf{E}_{k,j-1}^e + \mathbf{E}w(t_k^j) - \mathbf{E}w(t_k^{j-1}), \mathbf{E}_{k,j-1}^p) \in \mathcal{A}(w(t_k^j)),$$

we get

$$\begin{aligned} (6.5) \quad & \mathcal{Q}_1(\mathbf{E}_{k,j}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k,j}^p) + \mathcal{H}(\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e) - \langle \mathcal{L}(t_k^j), u_{k,j} \rangle \\ & \leq \mathcal{Q}_1(\mathbf{E}_{k,j-1}^e + \mathbf{E}w(t_k^j) - \mathbf{E}w(t_k^{j-1})) + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k,j-1}^p) \\ & \quad - \langle \mathcal{L}(t_k^j), u_{k,j-1} + w(t_k^j) - w(t_k^{j-1}) \rangle \\ & = \mathcal{Q}_1(\mathbf{E}_{k,j-1}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k,j-1}^p) - \langle \mathcal{L}(t_k^{j-1}), u_{k,j-1} \rangle + \int_{t_k^{j-1}}^{t_k^j} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau \\ & \quad - \int_{t_k^{j-1}}^{t_k^j} \langle \dot{\mathcal{L}}(\tau), u_k(\tau) \rangle \, d\tau - \int_{t_k^{j-1}}^{t_k^j} \langle \mathcal{L}_k(\tau), \dot{w}(\tau) \rangle \, d\tau + \delta_{k,j}, \end{aligned}$$

where

$$\delta_{k,j} := \mathcal{Q}_1(\mathbf{E}w(t_k^j) - \mathbf{E}w(t_k^{j-1})) - \langle \mathcal{L}(t_k^j) - \mathcal{L}(t_k^{j-1}), w(t_k^j) - w(t_k^{j-1}) \rangle.$$

By summing up from $j = 1$ to $j = i$ we get

$$\begin{aligned} \mathcal{E}_k(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_k^p; 0, t) & \leq \mathcal{E}_k(0) + \int_0^{t_k^i} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau \\ & \quad - \int_0^{t_k^i} \langle \dot{\mathcal{L}}(\tau), u_k(\tau) \rangle \, d\tau - \int_0^{t_k^i} \langle \mathcal{L}_k(\tau), \dot{w}(\tau) \rangle \, d\tau + \sum_{j=1}^i \delta_{k,j}. \end{aligned}$$

Since

$$\begin{aligned} \delta_{k,j} & \leq \frac{\beta_{\mathbb{C}}}{k} \int_{t_k^{j-1}}^{t_k^j} \|\mathbf{E} \dot{w}(\tau)\|_{L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})}^2 \, d\tau \\ & \quad + \sup_j \left\| \mathcal{L}(t_k^j) - \mathcal{L}(t_k^{j-1}) \right\|_{(W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N))^*} \int_{t_k^{j-1}}^{t_k^j} \|\dot{w}(\tau)\|_{W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)} \, d\tau, \end{aligned}$$

by setting $e_k := \sum_{j=1}^k \delta_{k,j}$, we get $e_k \rightarrow 0$ as $k \rightarrow +\infty$. Since

$$\mathcal{E}_k(0) = \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle,$$

inequality (6.3) follows.

Let us prove (6.4). By using the safe load condition on f and g , by (4.15) we can rewrite the first inequality of (6.5) in the following form:

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_{k,j}^e) - \int_{\Omega} \rho(t_k^j) : \mathbf{E}_{k,j}^e \, dx + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k,j}^p) + \mathcal{H}(\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e) - \int_{\Omega} \rho_D(t_k^j) : \mathbf{E}_{k,j}^p \, dx \\ \leq \mathcal{Q}_1(\mathbf{E}_{k,j-1}^e + \mathbf{E}w(t_k^j) - \mathbf{E}w(t_k^{j-1})) + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k,j-1}^p) \\ - \int_{\Omega} \rho(t_k^j) : \mathbf{E}_{k,j-1}^e \, dx - \int_{\Omega} \rho_D(t_k^j) : \mathbf{E}_{k,j-1}^p \, dx, \end{aligned}$$

so that

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_{k,j}^e) &- \int_{\Omega} \rho(t_k^j) : \mathbf{E}_{k,j}^e \, dx + \mathcal{Q}_2(\operatorname{curl} \mathbf{E}_{k,j}^p) + \mathcal{H}(\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e) \\ &- \int_{\Omega} \rho_D(t_k^j) : (\mathbf{E}_{k,j}^p - \mathbf{E}_{k,j-1}^p) \, dx \leq \mathcal{Q}_1(\mathbf{E}_{k,j-1}^e) - \int_{\Omega} \rho(t_k^{j-1}) : \mathbf{E}_{k,j-1}^e \, dx \\ &\quad + \mathcal{Q}_2(\operatorname{curl} \mathbf{E}_{k,j-1}^p) + \int_{t_k^{j-1}}^{t_k^j} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau \\ &\quad - \int_{t_k^{j-1}}^{t_k^j} \int_{\Omega} \dot{\rho}(\tau) : \mathbf{E}_k^e(\tau) \, dx \, d\tau + \tilde{\delta}_{k,j}, \end{aligned}$$

where

$$\tilde{\delta}_{k,j} := \mathcal{Q}_1(\mathbf{E} w(t_k^j) - \mathbf{E} w(t_k^{j-1})) \leq \frac{\beta_C}{k} \int_{t_k^{j-1}}^{t_k^j} \|\mathbf{E} \dot{w}(\tau)\|_{L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})}^2 \, d\tau.$$

By summing up from 0 to i we have

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_{k,i}^e) &- \int_{\Omega} \rho(t_k^i) : (\mathbf{E}_{k,i}^e - \mathbf{E} w(t_k^i)) \, dx + \mathcal{Q}_2(\operatorname{curl} \mathbf{E}_{k,i}^p) \\ &\quad + \sum_{j=0}^i \left[\mathcal{H}(\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e) - \int_{\Omega} \rho_D(t_k^j) : (\mathbf{E}_{k,j}^p - \mathbf{E}_{k,j-1}^p) \, dx \right] \\ &\leq \mathcal{Q}_1(\mathbf{E}_{k,0}^e) - \int_{\Omega} \rho(0) : (\mathbf{E}_{k,0}^e - \mathbf{E} w_{k,0}) \, dx + \mathcal{Q}_2(\operatorname{curl} \mathbf{E}_{k,0}^p) \\ &\quad + \int_0^{t_k^i} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^{t_k^i} \int_{\Omega} \dot{\rho}(\tau) : (\mathbf{E}_k^e(\tau) - \mathbf{E} w_k(\tau)) \, dx \, d\tau + \tilde{e}_k, \end{aligned}$$

where $\tilde{e}_k := \sum_{j=0}^k \tilde{\delta}_{k,j} \rightarrow 0$ as $k \rightarrow +\infty$. Since by (4.17) we have

$$\begin{aligned} \sum_{j=0}^i \left[\mathcal{H}(\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e) - \int_{\Omega} \rho_D(t_k^j) : (\mathbf{E}_{k,j}^p - \mathbf{E}_{k,j-1}^p) \, dx \right] \\ \geq \alpha_l \sum_{j=0}^i \|\mathbf{E}_{k,j}^e - \mathbf{E}_{k,j-1}^e\|_{BV(\Omega; \mathbb{M}_D^{N \times N})}, \end{aligned}$$

we deduce that

$$\begin{aligned} (6.6) \quad \mathcal{Q}_1(\mathbf{E}_k^e(t)) &- \int_{\Omega} \rho(t_k^i) : (\mathbf{E}_k^e(t) - \mathbf{E} w_k(t)) \, dx + \mathcal{Q}_2(\operatorname{curl} \mathbf{E}_k^p(t)) + \alpha_l \mathcal{V}(\mathbf{E}_k^e; 0, t) \\ &\leq C_1 + \int_0^{t_k^i} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^{t_k^i} \int_{\Omega} \dot{\rho}(\tau) : (\mathbf{E}_k^e(\tau) - \mathbf{E} w_k(\tau)) \, dx \, d\tau \end{aligned}$$

for some $C_1 > 0$ independent of k and t . Since $\mathcal{Q}_1(\mathbf{E}_k^e(t))$ is quadratic, we get that $\|\mathbf{E}_k^e(t)\|_{L^2}$ is uniformly bounded in k and t . Hence from (6.6) we deduce also that $\|\operatorname{curl} \mathbf{E}_k^p(t)\|_{L^2}$ and $\mathcal{V}(\mathbf{E}_k^e; 0, t)$ are uniformly bounded with respect to k and t . Since $u_k(t) = w_k(t)$ on $\partial_D \Omega$, by Korn's inequality we have also that $u_k(t)$ is uniformly bounded in $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ with respect to k and t . The proof of (6.4) is thus concluded. \square

7. Existence of a quasi-static evolution and approximation results. In this section we prove that the discrete evolution $t \mapsto (u_k(t), \mathbf{E}_k^e(t), \mathbf{E}_k^p(t))$ given by Proposition 6.2 admits a subsequence converging (in a suitable sense) to a quasi-static evolution for the Gurtin–Anand model. This will be done in Lemmas 7.1, 7.2, and 7.3. Theorem 5.2 will thus follow by combining these lemmas.

LEMMA 7.1. *Let $(u_k(t), \mathbf{E}_k^e(t), \mathbf{E}_k^p(t))_{k \in \mathbb{N}}$ be a sequence of discrete evolutions satisfying (7.1)–(7.3) and (7.4). Then there exists a subsequence $(k_j)_{j \in \mathbb{N}}$ such that $(u_{k_j}(t), \mathbf{E}_{k_j}^e(t), \mathbf{E}_{k_j}^p(t))_{j \in \mathbb{N}}$ converges to a quasi-static evolution $(u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))_{t \in [0, T]}$ with $(u(0), \mathbf{E}^e(0), \mathbf{E}^p(0)) = (u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$ and*

$$(7.1) \quad (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}(w(t)),$$

$$u_{k_j}(t) \rightharpoonup u(t) \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N),$$

$$(7.2) \quad \mathbf{E}_{k_j}^e(t) \rightharpoonup \mathbf{E}^e(t) \quad \text{weakly in } L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N}),$$

$$(7.3) \quad \mathbf{E}_{k_j}^p(t) \overset{*}{\rightharpoonup} \mathbf{E}^p(t) \quad \text{weakly* in } BV(\Omega; \mathbb{M}_D^{N \times N}),$$

$$(7.4) \quad \text{curl} \mathbf{E}_{k_j}^p(t) \rightharpoonup \text{curl} \mathbf{E}^p(t) \quad \text{weakly in } L^2(\Omega; \mathbb{M}^{N \times N}).$$

Moreover, there exists a constant $C > 0$ depending only on $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$ and w such that

$$(7.5) \quad \|u(t)\|_{W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)} + \|\mathbf{E}^e(t)\|_{L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})} + \|\text{curl} \mathbf{E}^p(t)\|_{L^2(\Omega; \mathbb{M}^{N \times N})} + \mathcal{V}(\mathbf{E}^p; 0, t) \leq C$$

for all $t \in [0, T]$ and for all $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$.

$$(7.6) \quad \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p(t)) - \langle \mathcal{L}(t), u(t) \rangle \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{p}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^p(t)).$$

By Proposition 6.2 we have

$$(7.7) \quad \|u_k(t)\|_{W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)} + \|\mathbf{E}_k^e(t)\|_{L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})} + \|\text{curl} \mathbf{E}_k^p(t)\|_{L^2(\Omega; \mathbb{M}^{N \times N})} + \mathcal{V}(\mathbf{E}_k^p; 0, t) \leq C$$

for some C independent of k and t . Since $\mathbf{E}_k^p(0) = \mathbf{E}_0^p$ and $\mathcal{V}(\mathbf{E}_k^p; 0, T) \leq C$, the existence of $\mathbf{E}^p \in BV(0, T; BV(\Omega; \mathbb{M}_D^{N \times N}))$ and of a subsequence $t \mapsto (u_{k_j}(t), \mathbf{E}_{k_j}^e(t), \mathbf{E}_{k_j}^p(t))$ such that (7.3) holds follows by applying the generalized version of Helly’s theorem proved in [8, Lemma 7.2] (notice that BV can be seen as the dual of a separable Banach space in such a way that the associated convergence with respect to the weak* topology is precisely the weak* convergence defined in (2.3)).

Since weak* convergence in BV implies strong convergence in L^1 , by (7.7) we deduce that $\text{curl} \mathbf{E}^p(t) \in L^2(\Omega; \mathbb{M}^{N \times N})$ and that (7.4) holds.

Let us fix $t \in [0, T]$. In view of the coercivity estimate (7.7), we may assume that there exist $\tilde{u} \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$, $\tilde{\mathbf{E}}^e \in L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$, and a further subsequence k_{j_h} (depending a priori on t) such that

$$u_{k_{j_h}}(t) \rightharpoonup \tilde{u} \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$$

and

$$(7.8) \quad \mathbf{E}_{k_{j_h}}^e(t) \rightharpoonup \widetilde{\mathbf{E}}^e \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N}).$$

It follows easily that $(\tilde{u}, \widetilde{\mathbf{E}}^e, \mathbf{E}^P(t)) \in \mathcal{A}(w(t))$. We claim that for every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$ we have

$$(7.9) \quad \begin{aligned} \mathcal{Q}_1(\widetilde{\mathbf{E}}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^P(t)) - \langle \mathcal{L}(t), \tilde{u} \rangle \\ \leq \mathcal{Q}_1(\mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{p}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^P(t)). \end{aligned}$$

Notice that, in view of (7.9), it turns out that \tilde{u} and $\widetilde{\mathbf{E}}^e$ are uniquely determined by $\mathbf{E}^P(t)$. In fact, the pair $(\tilde{u}, \widetilde{\mathbf{E}}^e)$ minimizes the convex functional $(v, \mathbf{e}) \mapsto \mathcal{Q}_1(\mathbf{e}) - \langle \mathcal{L}(t), v \rangle$ on the convex set $K := \{(v, \mathbf{e}) : (v, \mathbf{e}, \mathbf{E}^P(t)) \in \mathcal{A}(w(t))\}$. Since the functional is strictly convex in \mathbf{e} , $\widetilde{\mathbf{E}}^e$ is uniquely determined, and so is \tilde{u} in view of Korn's inequality. By setting $u(t) := \tilde{u}$ and $\mathbf{E}^e(t) := \widetilde{\mathbf{E}}^e$, we get that (7.1) and (7.2) hold (without passing to a further subsequence).

In view of (7.7) we deduce that (7.5) holds. Finally, the global stability is given precisely by (7.9).

In order to conclude the proof, we need to prove claim (7.9). Let us set

$$v_h := v - \tilde{u} + u_{k_{j_h}}(t), \quad \mathbf{e}_h := \mathbf{e} - \widetilde{\mathbf{E}}^e + \mathbf{E}_{k_{j_h}}^e(t), \quad \text{and} \quad \mathbf{p}_h := \mathbf{p} - \mathbf{E}^P(t) + \mathbf{E}_{k_{j_h}}^P(t).$$

We have $(v_h, \mathbf{e}_h, \mathbf{p}_h) \in \mathcal{A}(w_{k_{j_h}}(t))$. By (6.2) we have

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}_{k_{j_h}}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}_{k_{j_h}}^P(t)) - \langle \mathcal{L}_{k_{j_h}}(t), u_{k_{j_h}}(t) \rangle \\ \leq \mathcal{Q}_1(\mathbf{e}_h) + \mathcal{Q}_2(\text{curl} \mathbf{p}_h) - \langle \mathcal{L}_{k_{j_h}}(t), v_h \rangle + \mathcal{H}(\mathbf{p}_h - \mathbf{E}_{k_{j_h}}^P(t)) \\ = \mathcal{Q}_1(\mathbf{e} - \widetilde{\mathbf{E}}^e + \mathbf{E}_{k_{j_h}}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{p} - \text{curl} \mathbf{E}^P(t) + \text{curl} \mathbf{E}_{k_{j_h}}^P(t)) \\ - \langle \mathcal{L}_{k_{j_h}}(t), v - \tilde{u} + u_{k_{j_h}}(t) \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^P(t)), \end{aligned}$$

so that we get

$$\begin{aligned} 0 \leq \mathcal{Q}_1(\mathbf{e} - \widetilde{\mathbf{E}}^e) + \int_{\Omega} \mathbb{C}(\mathbf{e} - \widetilde{\mathbf{E}}^e) : \mathbf{E}_{k_{j_h}}^e(t) \, dx \\ + \mathcal{Q}_2(\text{curl} \mathbf{p} - \text{curl} \mathbf{E}^P(t)) + \mu L^2 \int_{\Omega} (\text{curl} \mathbf{p} - \text{curl} \mathbf{E}^P(t)) : \text{curl} \mathbf{E}_{k_{j_h}}^P(t) \, dx \\ - \langle \mathcal{L}_{k_{j_h}}(t), v - \tilde{u} \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^P(t)). \end{aligned}$$

By letting $h \rightarrow +\infty$, in view of (7.8), (7.3), and (7.4) and since $t \mapsto \mathcal{L}(t)$ is absolutely continuous with values in $(W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N))^*$, we obtain

$$\begin{aligned} 0 \leq \mathcal{Q}_1(\mathbf{e} - \widetilde{\mathbf{E}}^e) + \int_{\Omega} \mathbb{C}(\mathbf{e} - \widetilde{\mathbf{E}}^e) : \widetilde{\mathbf{E}}^e \, dx \\ + \mathcal{Q}_2(\text{curl} \mathbf{p} - \text{curl} \mathbf{E}^P(t)) + \mu L^2 \int_{\Omega} (\text{curl} \mathbf{p} - \text{curl} \mathbf{E}^P(t)) : \text{curl} \mathbf{E}^P(t) \, dx \\ - \langle \mathcal{L}(t), v - \tilde{u} \rangle + \mathcal{H}(\mathbf{p} - \mathbf{E}^P(t)). \end{aligned}$$

By adding to both sides the term $\mathcal{Q}_1(\widetilde{\mathbf{E}}^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}^P(t)) - \langle \mathcal{L}(t), \tilde{u} \rangle$, we obtain precisely (7.9), so that the proof is concluded. \square

We have the following estimate from above for the total energy.

LEMMA 7.2. $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$... 7.1

$t \in [0, T]$...

$$(7.10) \quad \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \leq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle + \int_0^t \int_{\Omega} \mathbb{C} \mathbf{E}^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau,$$

$$\mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \leq f_1(t) + f_2(t) \quad (5.4) \quad (5.1)$$

Let us fix $t \in [0, T]$. By (6.3) we have

$$(7.11) \quad \mathcal{E}_k(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_k^p; 0, t) \leq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle + \int_0^{t_k^i} \int_{\Omega} \mathbb{C} \mathbf{E}_k^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^{t_k^i} \langle \dot{\mathcal{L}}(\tau), u_k(\tau) \rangle \, d\tau - \int_0^{t_k^i} \langle \mathcal{L}_k(\tau), \dot{w}(\tau) \rangle \, d\tau + e_k,$$

where $e_k \rightarrow 0$ as $k \rightarrow +\infty$. In view of (7.2), (7.4), and (7.1) and since $\mathcal{L}_k(t) \rightarrow \mathcal{L}(t)$ strongly in $(W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N))^*$, we get by lower semicontinuity

$$\mathcal{E}(t) \leq \liminf_{j \rightarrow +\infty} \mathcal{E}_{k_j}(t).$$

Moreover, by (7.3) and the lower semicontinuity of \mathcal{H} with respect to the weak* convergence in BV , the very definition of $\mathcal{D}_{\mathcal{H}}$ implies that

$$\mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \leq \liminf_{j \rightarrow +\infty} \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t).$$

By Lebesgue dominated convergence we get as $k \rightarrow +\infty$

$$\int_0^{t_{k_j}^i} \int_{\Omega} \mathbb{C} \mathbf{E}_{k_j}^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^{t_{k_j}^i} \langle \dot{\mathcal{L}}(\tau), u_{k_j}(\tau) \rangle \, d\tau - \int_0^{t_{k_j}^i} \langle \mathcal{L}_{k_j}(\tau), \dot{w}(\tau) \rangle \, d\tau \rightarrow \int_0^t \int_{\Omega} \mathbb{C} \mathbf{E}^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau.$$

Then (7.10) follows passing to the limit in (7.11). \square

The following estimate from below for the total energy holds.

LEMMA 7.3. $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$... 7.1

$t \in [0, T]$...

$$(7.12) \quad \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \geq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl} \mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle + \int_0^t \int_{\Omega} \mathbb{C} \mathbf{E}^e(\tau) : \mathbf{E} \dot{w}(\tau) \, dx \, d\tau - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau,$$

$$\mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \leq f_1(t) + f_2(t) \quad (5.4) \quad (5.1)$$

Let $t \in [0, T]$, $h \geq 1$, and let us set $s_h^j := \frac{j}{h}t$ for $j = 0, 1, \dots, h$. By the global stability condition (7.6), by comparing $(u(s_h^j), \mathbf{E}^e(s_h^j), \mathbf{E}^p(s_h^j))$ with

$$\begin{aligned} & (u(s_h^{j+1}) - w(s_h^{j+1}) + w(s_h^j), \\ & \quad \mathbf{E}^e(s_h^{j+1}) - \mathbf{E}w(s_h^{j+1}) + \mathbf{E}w(s_h^j), \mathbf{E}^p(s_h^{j+1})) \in \mathcal{A}(w(s_h^j)), \end{aligned}$$

we get

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^e(s_h^{j+1}) - \mathbf{E}w(s_h^{j+1}) + \mathbf{E}w(s_h^j)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(s_h^{j+1})) \\ & \quad - \langle \mathcal{L}(s_h^j), u(s_h^{j+1}) - w(s_h^{j+1}) + w(s_h^j) \rangle + \mathcal{H}(\mathbf{E}^p(s_h^{j+1}) - \mathbf{E}^p(s_h^j)) \\ & \quad \geq \mathcal{Q}_1(\mathbf{E}^e(s_h^j)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(s_h^j)) - \langle \mathcal{L}(s_h^j), u(s_h^j) \rangle, \end{aligned}$$

which can be rewritten in the following form:

$$\begin{aligned} (7.13) \quad & \mathcal{Q}_1(\mathbf{E}^e(s_h^{j+1})) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(s_h^{j+1})) - \langle \mathcal{L}(s_h^{j+1}), u(s_h^{j+1}) \rangle \\ & + \mathcal{H}(\mathbf{E}^p(s_h^{j+1}) - \mathbf{E}^p(s_h^j)) \geq \mathcal{Q}_1(\mathbf{E}^e(s_h^j)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(s_h^j)) - \langle \mathcal{L}(s_h^j), u(s_h^j) \rangle \\ & + \int_{s_h^j}^{s_h^{j+1}} \int_{\Omega} \mathbb{C}\overline{\mathbf{E}}_h^e(s) : \mathbf{E}\dot{w}(s) \, dx \, ds - \int_{s_h^j}^{s_h^{j+1}} \langle \dot{\mathcal{L}}(s), \bar{u}_h(s) \rangle \, ds - \int_{s_h^j}^{s_h^{j+1}} \langle \overline{\mathcal{L}}_h(s), \dot{w}(s) \rangle \, ds + \bar{\delta}_{h,j}, \end{aligned}$$

where for $s_h^j < s \leq s_h^{j+1}$ we set

$$\bar{u}_h(s) := u(s_h^{j+1}), \quad \overline{\mathbf{E}}_h^e(s) := \mathbf{E}^e(s_h^{j+1}), \quad \overline{\mathbf{E}}_h^p(s) := \mathbf{E}^p(s_h^{j+1}), \quad \overline{\mathcal{L}}_h(s) := \mathcal{L}(s_h^{j+1}),$$

and

$$\bar{\delta}_{h,j} := -\mathcal{Q}_1(\mathbf{E}w(s_h^{j+1}) - \mathbf{E}w(s_h^j)) - \langle \mathcal{L}(s_h^{j+1}) - \mathcal{L}(s_h^j), w(s_h^{j+1}) - w(s_h^j) \rangle.$$

By summing up in (7.13) from 0 to $h - 1$ we get

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^p(t)) - \langle \mathcal{L}(t), u(t) \rangle + \sum_{j=0}^{h-1} \mathcal{H}(\mathbf{E}^p(s_h^{j+1}) - \mathbf{E}^p(s_h^j)) \\ & \quad \geq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl}\mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle \\ & \quad + \int_0^t \int_{\Omega} \mathbb{C}\overline{\mathbf{E}}_h^e(s) : \mathbf{E}\dot{w}(s) \, dx \, ds - \int_0^t \langle \dot{\mathcal{L}}(s), \bar{u}_h(s) \rangle \, ds - \int_0^t \langle \overline{\mathcal{L}}_h(s), \dot{w}(s) \rangle \, ds + \bar{e}_h, \end{aligned}$$

where $\bar{e}_h \rightarrow 0$ as $h \rightarrow +\infty$. By the very definition of $\mathcal{D}_{\mathcal{H}}$ we get

$$\begin{aligned} (7.14) \quad & \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \geq \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl}\mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle \\ & \quad + \int_0^t \int_{\Omega} \mathbb{C}\overline{\mathbf{E}}_h^e(s) : \mathbf{E}\dot{w}(s) \, dx \, ds - \int_0^t \langle \dot{\mathcal{L}}(s), \bar{u}_h(s) \rangle \, ds \\ & \quad \quad \quad - \int_0^t \langle \overline{\mathcal{L}}_h(s), \dot{w}(s) \rangle \, ds + \bar{e}_h. \end{aligned}$$

Since $\mathbf{E}^p \in BV(0, T; BV(\Omega; M_D^{N \times N}))$, we have that $\mathbf{E}^p(t)$ is continuous in time with respect to the strong topology in $BV(\Omega; M_D^{N \times N})$ up to a countable set in $[0, T]$. Let $s \in [0, T]$ be a continuity point of \mathbf{E}^p , and let $s_n \rightarrow s$. Then

$$(7.15) \quad \mathbf{E}^e(s_n) \rightharpoonup \mathbf{E}^e(s) \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N})$$

and

$$(7.16) \quad u(s_n) \rightharpoonup u(s) \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N).$$

In fact up to a subsequence we have that

$$\mathbf{E}^e(s_n) \rightharpoonup \widetilde{\mathbf{E}}^e \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N})$$

and

$$u(s_n) \rightharpoonup \tilde{u} \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N),$$

with $(\tilde{u}, \widetilde{\mathbf{E}}^e, \mathbf{E}^p(s)) \in \mathcal{A}(w(s))$. Given $(v, \mathbf{e}, \mathbf{E}^p(s)) \in \mathcal{A}(w(s))$, by the global stability condition (7.6), comparing $(u(s_n), \mathbf{E}^e(s_n), \mathbf{E}^p(s_n))$ with $(v - w(s) + w(s_n), \mathbf{e} - \mathbf{E}w(s) + \mathbf{E}w(s_n), \mathbf{E}^p(s))$, and taking into account the continuity of \mathcal{H} with respect to the BV -norm, we obtain that $(\tilde{u}, \widetilde{\mathbf{E}}^e)$ is a minimizer of the convex functional $(v, \mathbf{e}) \mapsto \mathcal{Q}_1(\mathbf{e}) - \langle \mathcal{L}(s), v \rangle$ on the convex set $\mathcal{K} := \{(v, \mathbf{e}) : (v, \mathbf{e}, \mathbf{E}^p(s)) \in \mathcal{A}(w(s))\}$. By the uniqueness of the minimizer, we have that $\tilde{u} = u(s)$ and $\widetilde{\mathbf{E}}^e = \mathbf{E}^e(s)$, so that (7.15) and (7.16) follow.

By (7.15) and (7.16) we have that for a.e. every $s \in [0, t]$

$$(7.17) \quad \overline{\mathbf{E}}_h^e(s) \rightharpoonup \mathbf{E}^e(s) \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N})$$

and

$$(7.18) \quad \overline{u}_h(s) \rightharpoonup u(s) \quad \text{weakly in } W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N).$$

By taking into account that for every $s \in [0, T]$

$$\overline{\mathcal{L}}_h(s) \rightarrow \mathcal{L}(s) \quad \text{strongly in } \left(W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N) \right)^*,$$

in view of (7.17) and (7.18), passing to the limit in (7.14) we get by dominated convergence (take into account (7.5)) that (7.12) follows. \square

We are in a position to prove Theorem 5.2. Indeed, the evolution $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ given by Lemma 7.1 is a quasi-static evolution for the Gurtin–Anand model because it satisfies the global stability condition in view of (7.6), and it satisfies the energy balance because of (7.10) and (7.12).

The convergence of the discrete-in-time evolution to the continuous one can be improved in the following way.

PROPOSITION 7.4. *Let $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ be the evolution given by Lemma 7.1. Let $t_j \in [0, T]$ be a sequence such that $t_j \rightarrow +\infty$.*

$$(7.19) \quad \mathcal{E}_{k_j}(t) \rightarrow \mathcal{E}(t)$$

$$(7.20) \quad \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t) \rightarrow \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t).$$

$$(7.21) \quad \mathbf{E}_{k_j}^e(t) \rightarrow \mathbf{E}^e(t) \quad \text{in } L^2(\Omega; M_{\text{sym}}^{N \times N}),$$

$$(7.22) \quad \text{curl} \mathbf{E}_{k_j}^p(t) \rightarrow \text{curl} \mathbf{E}^p(t) \quad \text{in } L^2(\Omega; M^{N \times N}),$$

$$(7.23) \quad \mathcal{Q}_1(\mathbf{E}_{k_j}^e(t)) \rightarrow \mathcal{Q}_1(\mathbf{E}^e(t)), \quad \mathcal{Q}_2(\text{curl}\mathbf{E}_{k_j}^p(t)) \rightarrow \mathcal{Q}_2(\text{curl}\mathbf{E}^p(t))$$

$$(7.24) \quad \langle \mathcal{L}_{k_j}(t), u_{k_j}(t) \rangle \rightarrow \langle \mathcal{L}(t), u(t) \rangle.$$

Notice that by lower semicontinuity we have for every $t \in [0, T]$

$$(7.25) \quad \mathcal{E}(t) \leq \liminf_{j \rightarrow +\infty} \mathcal{E}_{k_j}(t).$$

Moreover, by the lower semicontinuity of \mathcal{H} with respect to the weak* convergence, and by the very definition of $\mathcal{D}_{\mathcal{H}}$, we deduce that for every $t \in [0, T]$

$$(7.26) \quad \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) \leq \liminf_{j \rightarrow +\infty} \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t).$$

By (6.3) and (7.12) we get that

$$\begin{aligned} \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) &\leq \liminf_{j \rightarrow +\infty} (\mathcal{E}_{k_j}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t)) \\ &\leq \limsup_{j \rightarrow +\infty} \left(\mathcal{E}_{k_j}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t) \right) \\ &\leq \limsup_{j \rightarrow +\infty} \left[\mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl}\mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle + \int_0^{t_{k_j}^i} \int_{\Omega} \mathbb{C}\mathbf{E}_{k_j}^e(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau \right. \\ &\quad \left. - \int_0^{t_{k_j}^i} \langle \dot{\mathcal{L}}(\tau), u_{k_j}(\tau) \rangle \, d\tau - \int_0^{t_{k_j}^i} \langle \mathcal{L}_{k_j}(\tau), \dot{w}(\tau) \rangle \, d\tau + e_{k_j} \right] \\ &= \mathcal{Q}_1(\mathbf{E}_0^e) + \mathcal{Q}_2(\text{curl}\mathbf{E}_0^p) - \langle \mathcal{L}(0), u_0 \rangle + \int_0^t \int_{\Omega} \mathbb{C}\mathbf{E}^e(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau \\ &\quad - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau \leq \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t). \end{aligned}$$

We conclude that for every $t \in [0, T]$

$$\lim_{j \rightarrow +\infty} (\mathcal{E}_{k_j}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}_{k_j}^p; 0, t)) = \mathcal{E}(t) + \mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t).$$

From (7.25) and (7.26) we deduce that (7.19) and (7.20) hold. Since by lower semicontinuity

$$\mathcal{Q}_1(\mathbf{E}^e(t)) \leq \liminf_{j \rightarrow +\infty} \mathcal{Q}_1(\mathbf{E}_{k_j}^e(t)) \quad \text{and} \quad \mathcal{Q}_2(\text{curl}\mathbf{E}^p(t)) \leq \liminf_{j \rightarrow +\infty} \mathcal{Q}_2(\text{curl}\mathbf{E}_{k_j}^p(t)),$$

while

$$\langle \mathcal{L}_{k_j}(t), u_{k_j}(t) \rangle \rightarrow \langle \mathcal{L}(t), u(t) \rangle,$$

from (7.19) we deduce that (7.23) and (7.24) hold. In particular (7.21) and (7.22) follow, and the proof is concluded. \square

7.5 (convergence for the elastic and the “energetic” plastic strains). Since the maps $t \mapsto \mathbf{E}^e(t)$ and $t \mapsto \text{curl}\mathbf{E}^p(t)$ turn out to be uniquely determined by the initial conditions (see Proposition 8.8), we infer that

$$\begin{aligned} \mathbf{E}_k^e(t) &\rightarrow \mathbf{E}^e(t) && \text{strongly in } L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N}), \\ \text{curl}\mathbf{E}_k^p(t) &\rightarrow \text{curl}\mathbf{E}^p(t) && \text{strongly in } L^2(\Omega; \mathbb{M}^{N \times N}) \end{aligned}$$

without passing to a subsequence $(k_j)_{j \in \mathbb{N}}$.

8. Balance equations and the flow rule. This section is devoted to the proof of Theorem 5.3; that is, we prove that a quasi-static evolution $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ for the Gurtin–Anand model satisfies the prescribed regularity and uniqueness properties, the balance equations, and the flow rule.

We need the following lemma.

LEMMA 8.1. *Let $t \in [0, T]$. Let $\mathbf{T}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N})$, $\mathbb{K}_{\text{diss}}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N \times N})$, $\mathbb{S}^p(t) \in (\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}))^*$, $(\mathbf{A}, \mathbb{B}, \mathbb{L}) \in L^1(\Omega; \mathbb{M}_D^{N \times N}) \times L^1(\Omega; \mathbb{M}_D^{N \times N \times N}) \times \mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N})$.*

$$(8.1) \quad \left| \int_{\Omega} \mathbf{T}^p(t) : \mathbf{A} \, dx + \int_{\Omega} \mathbb{K}_{\text{diss}}^p(t) : \mathbb{B} \, dx + \langle \mathbb{S}^p(t), \mathbb{L} \rangle \right| \leq S_Y \int_{\Omega} \sqrt{|\mathbf{A}|^2 + l^2 |\mathbb{B}|^2} \, dx + l S_Y |\mathbb{L}|(\Omega)$$

Let $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$

$$(8.2) \quad \int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle = - \int_{\Omega} \mathbf{T}^p(t) : \mathbf{p} \, dx - \int_{\Omega} \mathbb{K}_{\text{diss}}^p(t) : \nabla \mathbf{p} \, dx - \langle \mathbb{S}^p(t), D^s \mathbf{p} \rangle.$$

Let $\mathbb{K}^p(t) := \mathbb{K}_{\text{en}}^p(t) + \mathbb{K}_{\text{diss}}^p(t)$. Let f_i be the i -th component of $\mathbb{K}^p(t)$. Then $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$ implies

$$(8.3) \quad \int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \int_{\Omega} \mathbf{T}^p(t) : \mathbf{p} \, dx + \int_{\Omega} \mathbb{K}^p(t) : \nabla \mathbf{p} \, dx + \langle \mathbb{S}^p(t), D^s \mathbf{p} \rangle = \langle \mathcal{L}(t), v \rangle.$$

Let us fix $t \in [0, T]$. From the global stability condition (5.2), for every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$ and $\varepsilon \in \mathbb{R}$ we get

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p(t)) - \langle \mathcal{L}(t), u(t) \rangle \\ & \leq \mathcal{Q}_1(\mathbf{E}^e(t) + \varepsilon \mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p(t) + \varepsilon \text{curl} \mathbf{p}) - \langle \mathcal{L}(t), u(t) + \varepsilon v \rangle + \mathcal{H}(\varepsilon \mathbf{p}) \end{aligned}$$

so that

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^e(t) + \varepsilon \mathbf{e}) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p(t) + \varepsilon \text{curl} \mathbf{p}) - \varepsilon \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\varepsilon \mathbf{p}) \\ & \geq \mathcal{Q}_1(\mathbf{E}^e(t)) + \mathcal{Q}_2(\text{curl} \mathbf{E}^p(t)). \end{aligned}$$

By taking the left and right derivative for $\varepsilon = 0$ we get

$$\int_{\Omega} \mathbf{C} \mathbf{E}^e(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle + \mathcal{H}(\mathbf{p}) \geq 0$$

and

$$\int_{\Omega} \mathbf{C} \mathbf{E}^e(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle - \mathcal{H}(-\mathbf{p}) \leq 0$$

so that, since $\mathbf{T}(t) := \mathbf{C} \mathbf{E}^e(t)$,

$$\left| \int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle \right| \leq \mathcal{H}(\mathbf{p}).$$

The previous inequality shows that the linear functional on $\mathcal{A}(0)$

$$(v, \mathbf{e}, \mathbf{p}) \mapsto \int_{\Omega} \mathbb{C}\mathbf{E}^e(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \operatorname{curl} \mathbf{E}^p(t) : \operatorname{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle$$

depends indeed only on \mathbf{p} .

Let $X \subseteq L^1(\Omega; \mathbb{M}_D^{N \times N}) \times L^1(\Omega; \mathbb{M}_D^{N \times N \times N}) \times \mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N})$ be the linear subspace generated by $\{(\mathbf{p}, \nabla \mathbf{p}, D^s \mathbf{p}) : (v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0) \text{ for some } v \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N), \mathbf{e} \in L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})\}$. By applying the Hahn–Banach theorem we deduce that the linear functional

$$(8.4) \quad \varphi(\mathbf{p}, \nabla \mathbf{p}, D^s \mathbf{p}) := \int_{\Omega} \mathbf{T}(t) : \mathbf{e} \, dx + \mu L^2 \int_{\Omega} \operatorname{curl} \mathbf{E}^p(t) : \operatorname{curl} \mathbf{p} \, dx - \langle \mathcal{L}(t), v \rangle$$

on the linear space X can be extended in a continuous way to the entire space $L^1(\Omega; \mathbb{M}_D^{N \times N}) \times L^1(\Omega; \mathbb{M}_D^{N \times N \times N}) \times \mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N})$ in such a way that

$$(8.5) \quad |\varphi(\mathbf{A}, \mathbb{B}, \mathbb{L})| \leq S_Y \int_{\Omega} \sqrt{|\mathbf{A}|^2 + l^2 |\mathbb{B}|^2} \, dx + l S_Y |\mathbb{L}|(\Omega)$$

for every $(\mathbf{A}, \mathbb{B}, \mathbb{L}) \in L^1(\Omega; \mathbb{M}_D^{N \times N}) \times L^1(\Omega; \mathbb{M}_D^{N \times N \times N}) \times \mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N})$. By representing φ , in view of (8.5) and (8.4), we obtain that there exist $\mathbf{T}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N})$, $\mathbb{K}_{\text{diss}}^p(t) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N \times N})$, and $\mathbb{S}^p(t) \in (\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}))^*$ such that (8.1) and (8.2) hold. Finally, (8.3) follows by (8.2) in view of the very definition of $\mathbb{K}_{\text{en}}^p(t)$. \square

The following proposition concerns the balance equation for the Cauchy stress.

PROPOSITION 8.2 (balance equations for the Cauchy stress). $\dots, t \in [0, T]$

$$(8.6) \quad \begin{cases} -\operatorname{div} \mathbf{T}(t) = f(t) & \text{in } \Omega, \\ \mathbf{T}(t) \nu = g(t) & \text{on } \partial_N \Omega. \end{cases}$$

Let $v \in C^\infty(\overline{\Omega}, \mathbb{R}^N)$ such that $v = 0$ on $\partial_D \Omega$. By choosing $(v, \mathbf{E}v, 0) \in \mathcal{A}(0)$ in (8.3) we deduce that

$$(8.7) \quad \int_{\Omega} \mathbf{T}(t) : \mathbf{E}v \, dx = \langle \mathcal{L}(t), v \rangle.$$

Then clearly $-\operatorname{div} \mathbf{T}(t) = f(t)$ in the sense of distributions in Ω . Since $\mathbf{T}(t) \in L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$ and its divergence belongs in particular to $L^2(\Omega; \mathbb{R}^N)$, we have that the normal trace of $\mathbf{T}(t)$ on $\partial \Omega$ is well defined as an element of $H^{-1/2}(\partial \Omega; \mathbb{R}^N)$. By integrating by parts in (8.7) we get immediately the second relation of (8.6). \square

Concerning the stresses conjugated to the plastic variables, the following result holds.

PROPOSITION 8.3 (stresses conjugated to the plastic variables). $\dots, t \in [0, T], \dots \mathbf{T}^p(t), \mathbb{K}_{\text{diss}}^p(t), \mathbb{K}^p(t), \dots \mathbb{S}^p(t) \dots 8.1 \dots$

$$(8.8) \quad \begin{cases} \mathbf{T}^p(t) = \mathbf{T}_D(t) + \operatorname{div} \mathbb{K}^p(t) & \text{in } \Omega, \\ \mathbb{K}^p(t) \nu = 0 & \text{on } \partial \Omega, \end{cases}$$

$\mathbf{T}_D(t) := (\mathbf{T}(t))_D \dots \mathbf{T}^p(t) \dots \mathbb{K}_{\text{diss}}^p(t) \dots$

$$(8.9) \quad \sqrt{|\mathbf{T}^p(t, x)|^2 + l^{-2} |\mathbb{K}_{\text{diss}}^p(t, x)|^2} \leq S_Y \dots x \in \Omega,$$

$$(8.10) \quad \|\mathbb{S}^P(t)\|_{(\mathcal{M}_b(M_D^{N \times N \times N}))^*} \leq lS_Y.$$

The stress constraints (8.9) and (8.10) follow by choosing $(\mathbf{A}, \mathbb{B}, 0)$ and $(0, 0, \mathbb{L})$, respectively, in (8.1).

Let us come to (8.8). Let $\mathbf{p} \in C^\infty(\bar{\Omega}, M_D^{N \times N})$, so that in particular $(0, -\mathbf{p}, \mathbf{p}) \in \mathcal{A}(0)$. Then (8.3) yields

$$-\int_{\Omega} \mathbf{T}(t) : \mathbf{p} \, dx + \int_{\Omega} \mathbf{T}^P(t) : \mathbf{p} \, dx + \int_{\Omega} \mathbb{K}^P(t) : \nabla \mathbf{p} \, dx = 0.$$

Since \mathbf{p} takes values in the space of deviatoric matrices, we can replace $\mathbf{T}(t)$ by $\mathbf{T}_D(t)$ so that we obtain

$$(8.11) \quad \int_{\Omega} (\mathbf{T}^P(t) - \mathbf{T}_D(t)) : \mathbf{p} \, dx + \int_{\Omega} \mathbb{K}^P(t) : \nabla \mathbf{p} \, dx = 0.$$

We conclude that the first relation of (8.8) holds. As a consequence, in view of (8.9) and the definition of $\mathbb{K}_{\text{en}}^P(t)$, we have $\mathbb{K}^P(t) \in L^2(\Omega; M_D^{N \times N \times N})$ with divergence in $L^2(\Omega; M_D^{N \times N})$, so that its normal trace on $\partial\Omega$ is well defined as an element of $H^{-1/2}(\partial\Omega; \mathbb{R}^{N \times N})$. By integrating by parts in (8.11) we obtain also the second relation of (8.8), and the proof is concluded. \square

8.4. Note that relation (8.3) represents the balance of internal and external power expenditures on the whole body Ω (see section 3). Due to our variational approach which requires $\mathbf{E}^P(t) \in BV(\Omega; M_D^{N \times N})$ so that $D\mathbf{E}^P(t)$ has also a singular part, a stress $\mathbb{S}^P(t)$ associated to $D^s\mathbf{E}^P(t)$ appears in the balance. In order to get a balance equation for a subbody $\mathcal{B} \subset\subset \Omega$ with a sufficiently smooth boundary, we can reason as follows. Let us assume to be in the physical case $N = 3$. As a consequence, admissible displacements v turn out to belong to $L^3(\Omega; \mathbb{R}^3)$. Let us assume for simplicity that $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$ is such that \mathbf{p} belongs also to $L^2(\Omega; M_D^{3 \times 3})$, and let $\varphi \in C_c^\infty(\Omega)$. By (8.3) we can write

$$(8.12) \quad \begin{aligned} & \int_{\Omega} \mathbf{T}(t) : (\varphi \mathbf{e}) \, dx + \int_{\Omega} \mathbf{T}^P(t) : (\varphi \mathbf{p}) \, dx + \int_{\Omega} \mathbb{K}^P(t) : (\varphi \nabla \mathbf{p}) \, dx + \langle \mathbb{S}^P(t), \varphi D^s \mathbf{p} \rangle \\ &= \int_{\Omega} \mathbf{T}(t) : (\varphi \mathbf{e} + v \odot \nabla \varphi) \, dx + \int_{\Omega} \mathbf{T}^P(t) : (\varphi \mathbf{p}) \, dx + \int_{\Omega} \mathbb{K}^P(t) : \nabla(\varphi \mathbf{p}) \, dx \\ & \quad + \langle \mathbb{S}^P(t), D^s(\varphi \mathbf{p}) \rangle - \int_{\Omega} \mathbf{T}(t) : (v \odot \nabla \varphi) \, dx - \int_{\Omega} \mathbb{K}^P(t) : (\mathbf{p} \otimes \nabla \varphi) \, dx \\ &= \langle \mathcal{L}(t), \varphi v \rangle - \int_{\Omega} \mathbf{T}(t) : (v \odot \nabla \varphi) \, dx - \int_{\Omega} \mathbb{K}^P(t) : (\mathbf{p} \otimes \nabla \varphi) \, dx, \end{aligned}$$

where the last equality follows since $(\varphi v, \varphi \mathbf{e} + v \odot \nabla \varphi, \varphi \mathbf{p}) \in \mathcal{A}(0)$ (we use $\mathbf{p} \in L^2(\Omega; M_D^{3 \times 3})$ to ensure that $\text{curl}(\varphi \mathbf{p}) \in L^2(\Omega; M^{3 \times 3})$). Here $(v \odot \nabla \varphi)_{i,j} = (v_i \partial_j \varphi + v_j \partial_i \varphi)/2$. As a consequence we have that the distribution

$$\varphi \mapsto - \int_{\Omega} \mathbf{T}(t) : (v \odot \nabla \varphi) \, dx - \int_{\Omega} \mathbb{K}^P(t) : (\mathbf{p} \otimes \nabla \varphi) \, dx$$

turns out to be a measure $\mu \in \mathcal{M}_b(\Omega)$. Moreover, by considering the measure $\eta \in \mathcal{M}_b(\Omega; \mathbb{R}^3)$ given by

$$\int_{\Omega} \psi \, d\eta := \int_{\Omega} \mathbf{T}(t) : (v \odot \psi) \, dx + \int_{\Omega} \mathbb{K}^P(t) : (\mathbf{p} \otimes \psi) \, dx, \quad \psi \in C_c^\infty(\Omega; \mathbb{R}^3),$$

we get immediately that $\operatorname{div} \eta = \mu$. According to [22], for every subset $\mathcal{B} \subset \subset \Omega$ with a sufficiently smooth boundary we have that η admits normal trace $\eta \cdot \nu$ on $\partial \mathcal{B}$ defined as an element of the dual of $C^1(\partial \mathcal{B})$, in such a way that the following Gauss–Green formula holds:

$$\int_{\mathcal{B}} d(\operatorname{div} \eta) = \langle \eta \cdot \nu, 1_{\partial \mathcal{B}} \rangle.$$

Let us denote formally $\eta \cdot \nu$ by $[\mathbf{T}(t)\nu \cdot v + \mathbb{K}^{\mathbf{P}}(t)\nu : \mathbf{p}]$, and let $[\mathbb{K}^{\mathbf{P}}(t) : \nabla \mathbf{p} + \mathbb{S}^{\mathbf{P}}(t) : D^s \mathbf{p}]$ be the measure such that

$$\int_{\Omega} \varphi d[\mathbb{K}^{\mathbf{P}}(t) : \nabla \mathbf{p} + \mathbb{S}^{\mathbf{P}}(t) : D^s \mathbf{p}] = \int_{\Omega} \mathbb{K}^{\mathbf{P}}(t) : (\varphi \nabla \mathbf{p}) \, dx + \langle \mathbb{S}^{\mathbf{P}}(t), \varphi D^s \mathbf{p} \rangle.$$

By (8.12) we can write by choosing $\varphi = 1_{\mathcal{B}}$

$$\begin{aligned} (8.13) \quad & \int_{\mathcal{B}} \mathbf{T}(t) : \mathbf{e} \, dx + \int_{\mathcal{B}} \mathbf{T}^{\mathbf{P}}(t) : \mathbf{p} \, dx + [\mathbb{K}^{\mathbf{P}}(t) : \nabla \mathbf{p} + \mathbb{S}^{\mathbf{P}}(t) : D^s \mathbf{p}](\mathcal{B}) \\ & = \int_{\mathcal{B}} f(t) \cdot v \, dx + \langle [\mathbf{T}(t)\nu \cdot v + \mathbb{K}^{\mathbf{P}}(t)\nu : \mathbf{p}], 1_{\partial \mathcal{B}} \rangle, \end{aligned}$$

which is a weak form for the balance of power expenditures for the subbody \mathcal{B} relative to the virtual velocity $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$.

The *AC* regularity in time for the quasi-static evolution is proved in the following proposition. The proof relies heavily on [8, Theorem 5.2]. We exploit the calculations in our context since we aim to understand the precise dependence on the material length scales l and L of the norms involved in the statement (in view of the convergence result of section 9, where we let $L, l \rightarrow 0$).

PROPOSITION 8.5. *Let $w, \rho, \alpha, \mathbb{C}, l, L \in \mathbb{R}^+$ and $\mathbb{C} \in \mathbb{C}$. Let $t \mapsto u(t) \in W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$, $t \mapsto \mathbf{E}^e(t) \in L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$, $t \mapsto \mathbf{E}^{\mathbf{P}}(t) \in BV(\Omega; \mathbb{M}_D^{N \times N})$, $t \mapsto \dot{\rho}(t) \in L^2(\Omega; \mathbb{M}^{N \times N})$, $t \mapsto \dot{\rho}_D(t) \in L^\infty(\Omega; \mathbb{R})$, and $t \in [0, T]$.*

$$(8.14) \quad \|\dot{u}(t)\|_{W^{1, \frac{N}{N-1}}} \leq C_1(w, \rho, \alpha, \mathbb{C}, l) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}],$$

$$(8.15) \quad \|\dot{\mathbf{E}}^e(t)\|_{L^2} \leq C_2(w, \rho, \alpha, \mathbb{C}) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}],$$

$$(8.16) \quad \|\dot{\mathbf{E}}^{\mathbf{P}}(t)\|_{BV} \leq C_3(w, \rho, \alpha, \mathbb{C}, l) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}],$$

$$(8.17) \quad \|\operatorname{curl} \dot{\mathbf{E}}^{\mathbf{P}}(t)\|_{L^2} \leq C_4(w, \rho, \alpha, \mathbb{C}, L) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}],$$

where C_1, C_2, C_3, C_4 are constants depending on $w, \rho, \alpha, \mathbb{C}, l, L$. (4.13) (4.14)

Let $t \mapsto u(t) \in BD(\Omega)$, $t \mapsto \mathbf{E}^{\mathbf{P}}(t) \in L^1(\Omega; \mathbb{M}_D^{N \times N})$, $t \mapsto \dot{\rho}(t) \in L^2(\Omega; \mathbb{M}^{N \times N})$, $t \mapsto \dot{\rho}_D(t) \in L^\infty(\Omega; \mathbb{R})$, and $t \in [0, T]$.

$$(8.18) \quad \|\dot{u}(t)\|_{BD} \leq C_5(w, \rho, \alpha, \mathbb{C}) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}],$$

$$(8.19) \quad \|\dot{\mathbf{E}}^{\mathbf{P}}(t)\|_{L^1} \leq C_6(w, \rho, \alpha, \mathbb{C}) [\|\dot{\rho}(t)\|_{L^2} + \|\dot{\rho}_D(t)\|_{L^\infty} + \|\mathbf{E}\dot{w}(t)\|_{L^2}].$$

Let $t_1, t_2 \in [0, T]$, with $t_1 < t_2$. Since by the very definition of $\mathcal{D}_{\mathcal{H}}$ we have $\mathcal{D}_{\mathcal{H}}(\mathbf{E}^{\text{P}}; t_1, t_2) \geq \mathcal{H}(\mathbf{E}^{\text{P}}(t_2) - \mathbf{E}^{\text{P}}(t_1))$, by the energy balance (5.3) we may write

$$(8.20) \quad \begin{aligned} & \mathcal{Q}_1(\mathbf{E}^{\text{e}}(t_2)) - \mathcal{Q}_1(\mathbf{E}^{\text{e}}(t_1)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^{\text{P}}(t_2)) - \mathcal{Q}_2(\text{curl}\mathbf{E}^{\text{P}}(t_1)) \\ & \quad + \mathcal{H}(\mathbf{E}^{\text{P}}(t_2) - \mathbf{E}^{\text{P}}(t_1)) - \langle \mathcal{L}(t_2), u(t_2) \rangle + \langle \mathcal{L}(t_1), u(t_1) \rangle \\ & \leq \int_{t_1}^{t_2} \int_{\Omega} \mathbf{T}(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau - \int_{t_1}^{t_2} \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_{t_1}^{t_2} \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau. \end{aligned}$$

Let us consider $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(0)$ such that

$$v := u(t_2) - u(t_1) - (w(t_2) - w(t_1)), \quad \mathbf{e} := \mathbf{E}^{\text{e}}(t_2) - \mathbf{E}^{\text{e}}(t_1) - (\mathbf{E}w(t_2) - \mathbf{E}w(t_1)),$$

and $\mathbf{p} := \mathbf{E}^{\text{P}}(t_2) - \mathbf{E}^{\text{P}}(t_1)$. By combining (8.1) and (8.2), we deduce that

$$(8.21) \quad \begin{aligned} & - \int_{\Omega} \mathbf{T}(t_1) : (\mathbf{E}^{\text{e}}(t_2) - \mathbf{E}^{\text{e}}(t_1) - (\mathbf{E}w(t_2) - \mathbf{E}w(t_1))) \, dx \\ & \quad - \mu L^2 \int_{\Omega} \text{curl}\mathbf{E}^{\text{P}}(t_1) : (\text{curl}\mathbf{E}^{\text{P}}(t_2) - \text{curl}\mathbf{E}^{\text{P}}(t_1)) \, dx \\ & \quad + \langle \mathcal{L}(t_1), u(t_2) - u(t_1) - (w(t_2) - w(t_1)) \rangle \leq \mathcal{H}(\mathbf{E}^{\text{P}}(t_2) - \mathbf{E}^{\text{P}}(t_1)). \end{aligned}$$

By inserting (8.21) into (8.20) and taking into account (4.15), we obtain

$$\begin{aligned} & \mathcal{Q}_1(\mathbf{E}^{\text{e}}(t_2) - \mathbf{E}^{\text{e}}(t_1)) + \mathcal{Q}_2(\text{curl}\mathbf{E}^{\text{P}}(t_2) - \text{curl}\mathbf{E}^{\text{P}}(t_1)) \\ & \leq \int_{t_1}^{t_2} \int_{\Omega} \mathbf{T}(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau - \int_{t_1}^{t_2} \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_{t_1}^{t_2} \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau \\ & \quad + \langle \mathcal{L}(t_2) - \mathcal{L}(t_1), u(t_2) \rangle - \int_{\Omega} \mathbf{T}(t_1) : (\mathbf{E}w(t_2) - \mathbf{E}w(t_1)) \, dx \\ & \quad \quad + \langle \mathcal{L}(t_1), w(t_2) - w(t_1) \rangle \\ & = \int_{t_1}^{t_2} \int_{\Omega} (\mathbf{T}(\tau) - \mathbf{T}(t_1)) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau - \int_{t_1}^{t_2} \langle \dot{\mathcal{L}}(\tau), u(\tau) - u(t_2) \rangle \, d\tau \\ & \quad \quad - \int_{t_1}^{t_2} \langle \mathcal{L}(\tau) - \mathcal{L}(t_1), \dot{w}(\tau) \rangle \, d\tau \\ & = \int_{t_1}^{t_2} \int_{\Omega} (\mathbf{T}(\tau) - \mathbf{T}(t_1)) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau - \int_{t_1}^{t_2} \int_{\Omega} \dot{\rho}(\tau) : (\mathbf{E}^{\text{e}}(\tau) - \mathbf{E}^{\text{e}}(t_2)) \, dx \, d\tau \\ & \quad - \int_{t_1}^{t_2} \int_{\Omega} \dot{\rho}_D(\tau) : (\mathbf{E}^{\text{P}}(\tau) - \mathbf{E}^{\text{P}}(t_2)) \, dx \, d\tau - \int_{t_1}^{t_2} \int_{\Omega} (\rho(\tau) - \rho(t_1)) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau. \end{aligned}$$

By the coercivity estimate (4.8) for the elasticity tensor \mathbb{C} we deduce that

$$(8.22) \quad \begin{aligned} & \alpha_{\mathbb{C}} \|\mathbf{E}^{\text{e}}(t_2) - \mathbf{E}^{\text{e}}(t_1)\|_{L^2}^2 + \frac{\mu L^2}{2} \|\text{curl}\mathbf{E}^{\text{P}}(t_2) - \text{curl}\mathbf{E}^{\text{P}}(t_1)\|_{L^2}^2 \\ & \leq \beta_{\mathbb{C}} \int_{t_1}^{t_2} \|\mathbf{E}^{\text{e}}(\tau) - \mathbf{E}^{\text{e}}(t_1)\|_{L^2} \|\mathbf{E}\dot{w}(\tau)\|_{L^2} \, d\tau + \int_{t_1}^{t_2} \|\dot{\rho}\|_{L^2} \|\mathbf{E}^{\text{e}}(\tau) - \mathbf{E}^{\text{e}}(t_2)\|_{L^2} \, d\tau \\ & \quad + \int_{t_1}^{t_2} \|\dot{\rho}_D(\tau)\|_{L^\infty} \|\mathbf{E}^{\text{P}}(\tau) - \mathbf{E}^{\text{P}}(t_2)\|_{L^1} \, d\tau + \int_{t_1}^{t_2} \|\rho(\tau) - \rho(t_1)\|_{L^2} \|\mathbf{E}\dot{w}(\tau)\|_{L^2} \, d\tau. \end{aligned}$$

By (4.16) we have for $t_1 \leq s \leq t_2$

$$(8.23) \quad \frac{\alpha}{2} \|\mathbf{E}^P(t_2) - \mathbf{E}^P(s)\|_{L^1} + \alpha_l \|D\mathbf{E}^P(t_2) - D\mathbf{E}^P(s)\|_{\mathcal{M}_b} \\ \leq \mathcal{H}(\mathbf{E}^P(t_2) - \mathbf{E}^P(s)) - \int_{\Omega} \rho_D(t_2) : (\mathbf{E}^P(t_2) - \mathbf{E}^P(s)) \, dx,$$

where $\alpha_l := \min\{l\frac{\alpha}{2}, lS_Y\}$. By combining (8.23) and (8.20) with $t_1 = s$ and using (4.15), we obtain

$$(8.24) \quad \frac{\alpha}{2} \|\mathbf{E}^P(t_2) - \mathbf{E}^P(s)\|_{L^1} + \alpha_l \|D\mathbf{E}^P(t_2) - D\mathbf{E}^P(s)\|_{\mathcal{M}_b} \\ \leq \mathcal{Q}_1(\mathbf{E}^e(s)) - \mathcal{Q}_1(\mathbf{E}^e(t_2)) + \mathcal{Q}_2(\operatorname{curl}\mathbf{E}^P(s)) - \mathcal{Q}_2(\operatorname{curl}\mathbf{E}^P(t_2)) \\ + \langle \mathcal{L}(t_2), u(t_2) \rangle - \langle \mathcal{L}(s), u(s) \rangle + \int_s^{t_2} \int_{\Omega} \mathbf{T}(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau \\ - \int_s^{t_2} \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_s^{t_2} \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau \\ - \int_{\Omega} \rho_D(t_2)(\mathbf{E}^P(t_2) - \mathbf{E}^P(s)) \, dx \\ \leq \mathcal{Q}_1(\mathbf{E}^e(s)) - \mathcal{Q}_1(\mathbf{E}^e(t_2)) + \mathcal{Q}_2(\operatorname{curl}\mathbf{E}^P(s)) - \mathcal{Q}_2(\operatorname{curl}\mathbf{E}^P(t_2)) \\ + \int_{\Omega} \rho(t_2) : (\mathbf{E}^e(t_2) - \mathbf{E}^e(s)) \, dx + \int_{\Omega} (\rho(t_2) - \rho(s)) : \mathbf{E}^e(s) \, dx \\ + \int_{\Omega} (\rho_D(t_2) - \rho_D(s)) : \mathbf{E}^P(s) \, dx \\ - \int_s^{t_2} \int_{\Omega} \{\dot{\rho}(\tau) : \mathbf{E}^e(\tau) + \dot{\rho}_D(\tau) : \mathbf{E}^P(\tau) - (\mathbf{T}(\tau) - \rho(\tau)) : \mathbf{E}\dot{w}(\tau)\} \, dx \, d\tau.$$

Notice that

$$\sup_{\tau} \|\rho(\tau)\|_{L^2}, \quad \sup_{\tau} \|\rho_D(\tau)\|_{L^\infty},$$

and

$$\sup_{\tau} \|\mathbf{E}^e(\tau)\|_{L^2}, \quad \sup_{\tau} \|\mathbf{E}^P(\tau)\|_{L^1}, \quad \sup_{\tau} \|\operatorname{curl}\mathbf{E}^P(\tau)\|_{L^2}$$

are finite (in fact $t \mapsto \mathbf{E}^P(t)$ has bounded variation, while for $\mathbf{E}^e(t)$ and $\operatorname{curl}\mathbf{E}^P(t)$ we can use the energy balance (5.3)). From (8.24) we obtain for every $t_1 \leq s \leq t_2$

$$(8.25) \quad \frac{\alpha}{2} \|\mathbf{E}^P(t_2) - \mathbf{E}^P(s)\|_{L^1} + \alpha_l \|D\mathbf{E}^P(t_2) - D\mathbf{E}^P(s)\|_{\mathcal{M}_b} \\ \leq C_1 \left(\|\mathbf{E}^e(t_2) - \mathbf{E}^e(s)\|_{L^2} + \sqrt{\frac{\mu}{2}} L \|\operatorname{curl}\mathbf{E}^P(t_2) - \operatorname{curl}\mathbf{E}^P(s)\|_{L^2} \right. \\ \left. + \int_s^{t_2} \psi(\tau) \, d\tau \right),$$

where

$$(8.26) \quad \psi(\tau) := \|\dot{\rho}(\tau)\|_{L^2} + \|\dot{\rho}_D(\tau)\|_{L^\infty} + \|\mathbf{E}\dot{w}(\tau)\|_{L^2}$$

and C_1 depends on ρ , $\sup_\tau \|\mathbf{E}^e(\tau)\|_{L^2}$, $\sup_\tau \|\mathbf{E}^p(\tau)\|_{L^1}$, $\sup_\tau L\|\text{curl}\mathbf{E}^p(\tau)\|_{L^2}$, and the elasticity tensor \mathbb{C} . By (8.22) we conclude that

$$\begin{aligned} & \alpha_{\mathbb{C}}\|\mathbf{E}^e(t_2) - \mathbf{E}^e(t_1)\|_{L^2}^2 + \frac{\mu L^2}{2}\|\text{curl}\mathbf{E}^p(t_2) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2}^2 \\ & \leq C_2 \left(\|\mathbf{E}^e(t_2) - \mathbf{E}^e(t_1)\|_{L^2} + \sqrt{\frac{\mu}{2}}L\|\text{curl}\mathbf{E}^p(t_2) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2} \right) \int_{t_1}^{t_2} \psi(\tau) d\tau \\ & + C_2 \int_{t_1}^{t_2} \psi(\tau) \left(\|\mathbf{E}^e(\tau) - \mathbf{E}^e(t_1)\|_{L^2} + \sqrt{\frac{\mu}{2}}L\|\text{curl}\mathbf{E}^p(\tau) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2} \right) d\tau \\ & + C_2 \left(\int_{t_1}^{t_2} \psi(\tau) d\tau \right)^2, \end{aligned}$$

where C_2 depends also on α . By Cauchy’s inequality we obtain

$$\begin{aligned} & \|\mathbf{E}^e(t_2) - \mathbf{E}^e(t_1)\|_{L^2}^2 + \frac{\mu L^2}{2}\|\text{curl}\mathbf{E}^p(t_2) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2}^2 \\ & \leq C_3 \int_{t_1}^{t_2} \psi(\tau) \left(\|\mathbf{E}^e(\tau) - \mathbf{E}^e(t_1)\|_{L^2} + \sqrt{\frac{\mu}{2}}L\|\text{curl}\mathbf{E}^p(\tau) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2} \right) d\tau \\ & + C_3 \left(\int_{t_1}^{t_2} \psi(\tau) d\tau \right)^2. \end{aligned}$$

By means of a Gronwall-type lemma [8, Lemma 5.3] we get in particular that

$$(8.27) \quad \|\mathbf{E}^e(t_2) - \mathbf{E}^e(t_1)\|_{L^2} + \sqrt{\frac{\mu}{2}}L\|\text{curl}\mathbf{E}^p(t_2) - \text{curl}\mathbf{E}^p(t_1)\|_{L^2} \leq C_4 \int_{t_1}^{t_2} \psi(\tau) d\tau,$$

where C_4 depends on ρ , α , $\sup_\tau \|\mathbf{E}^e(\tau)\|_{L^2}$, $\sup_\tau \|\mathbf{E}^p(\tau)\|_{L^1}$, $\sup_\tau L\|\text{curl}\mathbf{E}^p(\tau)\|_{L^2}$, and the elasticity tensor \mathbb{C} . As a consequence we get that $t \mapsto \mathbf{E}^e(t)$ and $t \mapsto \text{curl}\mathbf{E}^p(t)$ are absolutely continuous from $[0, T]$ to $L^2(\Omega; \mathbb{M}_{\text{sym}}^{N \times N})$ and $L^2(\Omega; \mathbb{M}^{N \times N})$, respectively. By (8.25), we get that $t \mapsto \mathbf{E}^p(t)$ is absolutely continuous from $[0, T]$ to $BV(\Omega; \mathbb{M}_D^{N \times N})$.

Let us now come to the proof of (8.14)–(8.17). By the energy balance (5.3), and by the very definition of \mathcal{H} , we deduce that

$$\|\mathbf{E}^p(t)\|_{L^1} \leq C_5 \left(1 + \int_0^t (1 + \psi(\tau))\|\mathbf{E}^e(\tau)\|_{L^2} d\tau + \int_0^t (1 + \psi(\tau))\|\mathbf{E}^p(\tau)\|_{L^1} d\tau \right),$$

where ψ is as in (8.26) and C_5 depends only on the initial conditions and on w, ρ, α , and \mathbb{C} . By means of the classical Gronwall lemma and taking the sup in t , we obtain

$$\sup_{t \in [0, T]} \|\mathbf{E}^p(t)\|_{L^1} \leq C_6 \left(1 + \sup_{t \in [0, T]} \|\mathbf{E}^e(t)\|_{L^2} \right).$$

By the energy balance (5.3) we conclude that $\sup_{t \in [0, T]} \|\mathbf{E}^e(t)\|_{L^2}$ is bounded uniformly independently of l and L , so that the same holds for $\sup_{t \in [0, T]} \|\mathbf{E}^p(t)\|_{L^1}$ and $\sup_{t \in [0, T]} L\|\text{curl}\mathbf{E}^p(t)\|_{L^2}$. By (8.27) we conclude that (8.15) and (8.17) hold. Inequalities (8.16) and (8.19) follow by (8.25). Finally the absolute continuity of $t \mapsto u(t)$ from $[0, T]$ to $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$ and inequality (8.14) follow from the compatibility condition $\mathbf{E}u(t) := \mathbf{E}^e(t) + \mathbf{E}^p(t)$ and Korn’s inequality. Inequality (8.18) follows in a similar way. \square

8.6. Notice that the constants C_1, \dots, C_6 of Proposition 8.5 depend also on the initial condition $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$. More precisely, from the previous proof it can be evicted that they depend on $|\mathcal{E}(0)|$ and $\|\mathbf{E}_0^p\|_{L^1(\Omega; M_D^{N \times N})}$.

In order to prove the uniqueness result and the flow rule, we need the following lemma.

LEMMA 8.7. *Let $t \in [0, T]$. Then*

$$(8.28) \quad \mathcal{H}(\dot{\mathbf{E}}^p(t)) = - \int_{\Omega} \mathbf{T}(t) : (\dot{\mathbf{E}}^e(t) - \mathbf{E}\dot{w}(t)) \, dx - \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \dot{\mathbf{E}}^p(t) \, dx + \langle \mathcal{L}(t), \dot{u}(t) - \dot{w}(t) \rangle$$

$$(8.29) \quad \mathcal{H}(\dot{\mathbf{E}}^p(t)) = \int_{\Omega} \mathbf{T}^p(t) : \dot{\mathbf{E}}^p(t) \, dx + \int_{\Omega} \mathbb{K}_{\text{diss}}^p(t) : \nabla \dot{\mathbf{E}}^p(t) \, dx + \langle \mathbb{S}^p(t), D^s \dot{\mathbf{E}}^p(t) \rangle.$$

Since by Proposition 8.5 the map $t \mapsto \mathbf{E}^p(t)$ is absolutely continuous from $[0, T]$ to $BV(\Omega; M_D^{N \times N})$, by [8, Theorem 7.1] we obtain

$$\mathcal{D}_{\mathcal{H}}(\mathbf{E}^p; 0, t) = \int_0^t \mathcal{H}(\dot{\mathbf{E}}^p(\tau)) \, d\tau.$$

Then by differentiating the energy balance equation (5.3) we obtain for a.e. $t \in [0, T]$

$$\begin{aligned} & \int_{\Omega} \mathbf{T}(t) : \dot{\mathbf{E}}^e(t) \, dx + \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \dot{\mathbf{E}}^p(t) \, dx \\ & - \langle \dot{\mathcal{L}}(t), u(t) \rangle - \langle \mathcal{L}(t), \dot{u}(t) \rangle + \mathcal{H}(\dot{\mathbf{E}}^p(t)) \\ & = \int_{\Omega} \mathbf{T}(t) : \mathbf{E}\dot{w}(t) \, dx - \langle \dot{\mathcal{L}}(t), u(t) \rangle - \langle \mathcal{L}(t), \dot{w}(t) \rangle \end{aligned}$$

so that (8.28) follows. Since $(\dot{u}(t) - \dot{w}(t), \dot{\mathbf{E}}^e(t) - \mathbf{E}\dot{w}(t), \dot{\mathbf{E}}^p(t)) \in \mathcal{A}(0)$, by (8.2) we get that (8.29) holds. \square

Let us now prove the uniqueness result concerning the elastic strain $\mathbf{E}^e(t)$ and the “energetic plastic strain” $\text{curl} \mathbf{E}^p(t)$.

PROPOSITION 8.8. *Let $t \mapsto \mathbf{E}^e(t)$ and $t \mapsto \text{curl} \mathbf{E}^p(t)$ be the unique solution to the problem (8.1) with initial condition $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$.*

Let $t \mapsto (\tilde{u}(t), \tilde{\mathbf{E}}^e(t), \tilde{\mathbf{E}}^p(t))$ be another quasi-static evolution associated to the same initial condition $(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$. Let $\tilde{\mathbf{T}}^p(t)$, $\tilde{\mathbb{K}}^p(t)$, and $\tilde{\mathbb{S}}^p(t)$ be the associated stresses according to Lemma 8.1 and Proposition 8.3. By (8.28), (8.1), and (8.2) we have for a.e. $t \in [0, T]$

$$\begin{aligned} & - \int_{\Omega} \mathbf{T}(t) : (\dot{\mathbf{E}}^e(t) - \mathbf{E}\dot{w}(t)) \, dx - \mu L^2 \int_{\Omega} \text{curl} \mathbf{E}^p(t) : \text{curl} \dot{\mathbf{E}}^p(t) \, dx \\ & + \langle \mathcal{L}(t), \dot{u}(t) - \dot{w}(t) \rangle = \mathcal{H}(\dot{\mathbf{E}}^p(t)) \\ & \geq \int_{\Omega} \tilde{\mathbf{T}}^p(t) : \dot{\mathbf{E}}^p(t) \, dx + \int_{\Omega} \tilde{\mathbb{K}}_{\text{diss}}^p(t) : \nabla \dot{\mathbf{E}}^p(t) \, dx + \langle \tilde{\mathbb{S}}^p(t), D^s \dot{\mathbf{E}}^p(t) \rangle \\ & = - \int_{\Omega} \tilde{\mathbf{T}}(t) : (\dot{\mathbf{E}}^e(t) - \mathbf{E}\dot{w}(t)) \, dx - \mu L^2 \int_{\Omega} \text{curl} \tilde{\mathbf{E}}^p(t) : \text{curl} \dot{\mathbf{E}}^p(t) \, dx \\ & + \langle \mathcal{L}(t), \dot{u}(t) - \dot{w}(t) \rangle. \end{aligned}$$

We obtain

$$\int_{\Omega} [\mathbf{T}(t) - \tilde{\mathbf{T}}(t)] : [\dot{\mathbf{E}}^e(t) - \mathbf{E}\dot{w}(t)] dx + \mu L^2 \int_{\Omega} [\operatorname{curl} \mathbf{E}^P(t) - \operatorname{curl} \tilde{\mathbf{E}}^P(t)] : \operatorname{curl} \dot{\mathbf{E}}^P(t) dx \leq 0.$$

Similarly we obtain

$$\int_{\Omega} [\tilde{\mathbf{T}}(t) - \mathbf{T}(t)] : [\dot{\tilde{\mathbf{E}}}^e(t) - \mathbf{E}\dot{w}(t)] dx + \mu L^2 \int_{\Omega} [\operatorname{curl} \tilde{\mathbf{E}}^P(t) - \operatorname{curl} \mathbf{E}^P(t)] : \operatorname{curl} \dot{\tilde{\mathbf{E}}}^P(t) dx \leq 0.$$

By summing the two inequalities we get

$$\begin{aligned} & \int_{\Omega} [\mathbf{T}(t) - \tilde{\mathbf{T}}(t)] : [\dot{\mathbf{E}}^e(t) - \dot{\tilde{\mathbf{E}}}^e(t)] dx \\ & + \mu L^2 \int_{\Omega} [\operatorname{curl} \mathbf{E}^P(t) - \operatorname{curl} \tilde{\mathbf{E}}^P(t)] : [\operatorname{curl} \dot{\mathbf{E}}^P(t) - \operatorname{curl} \dot{\tilde{\mathbf{E}}}^P(t)] dx \leq 0, \end{aligned}$$

so that for a.e. $t \in [0, T]$ we have

$$\begin{aligned} & \frac{d}{dt} \left(\int_{\Omega} \mathbb{C} [\mathbf{E}^e(t) - \tilde{\mathbf{E}}^e(t)] : [\mathbf{E}^e(t) - \tilde{\mathbf{E}}^e(t)] dx \right. \\ & \left. + \mu L^2 \int_{\Omega} |\operatorname{curl} \mathbf{E}^P(t) - \operatorname{curl} \tilde{\mathbf{E}}^P(t)|^2 dx \right) \leq 0. \end{aligned}$$

Since we have $\tilde{\mathbf{E}}^e(0) = \mathbf{E}_0^e = \mathbf{E}^e(0)$ and $\operatorname{curl} \tilde{\mathbf{E}}^P(0) = \operatorname{curl} \mathbf{E}_0^P = \operatorname{curl} \mathbf{E}^P(0)$, the previous inequality and the absolute continuity of the maps involved yield that $\mathbf{E}^e(t) = \tilde{\mathbf{E}}^e(t)$ and $\operatorname{curl} \tilde{\mathbf{E}}^P(t) = \operatorname{curl} \mathbf{E}^P(t)$ for every $t \in [0, T]$, so that the proof is concluded. \square

We are now in a position to prove the flow rule for the Gurtin–Anand model. Let us start with the following weak form.

PROPOSITION 8.9 (weak form of the flow rule). $\forall t \in [0, T]$. . .

$$(a) \quad (\mathbf{A}, \mathbb{B}) \in L^\infty(\Omega; \mathbb{M}_D^{N \times N}) \times L^\infty(\Omega; \mathbb{M}_D^{N \times N \times N})$$

$$\sqrt{|\mathbf{A}(x)|^2 + l^{-2} |\mathbb{B}(x)|^2} \leq S_Y \quad x \in \Omega,$$

$$(8.30) \quad \int_{\Omega} (\mathbf{T}^P(t) - \mathbf{A}) : \dot{\mathbf{E}}^P(t) dx + \int_{\Omega} (\mathbb{K}_{\text{diss}}^P(t) - \mathbb{B}) : \nabla \dot{\mathbf{E}}^P(t) dx \geq 0.$$

$$(b) \quad \mathbb{L} \in (\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}))^* \quad \|\mathbb{L}\|_{(\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}))^*} \leq l S_Y$$

$$(8.31) \quad \langle \mathbb{S}^P(t) - \mathbb{L}, D^s \dot{\mathbf{E}}^P(t) \rangle \geq 0.$$

Recall that for a.e. $t \in [0, T]$

$$\begin{aligned} \mathcal{H}(\dot{\mathbf{E}}^P(t)) &= \mathcal{F}(\dot{\mathbf{E}}^P(t), \nabla \dot{\mathbf{E}}^P(t), D^s \dot{\mathbf{E}}^P(t)) \\ &:= S_Y \int_{\Omega} \sqrt{|\dot{\mathbf{E}}^P|^2 + l^2 |\nabla \dot{\mathbf{E}}^P|^2} dx + l S_Y |D^s \dot{\mathbf{E}}^P|(\Omega). \end{aligned}$$

Since $\mathcal{F} : L^1(\Omega; \mathbb{M}_D^{N \times N}) \times L^1(\Omega; \mathbb{M}_D^{N \times N \times N}) \times \mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N}) \rightarrow [0, +\infty[$ is continuous (with respect to the strong norm), we have $\mathcal{F}(\dot{\mathbf{E}}^P(t), \nabla \dot{\mathbf{E}}^P(t), D^s \dot{\mathbf{E}}^P(t)) =$

$\mathcal{F}^{**}(\dot{\mathbf{E}}^P(t), \nabla \dot{\mathbf{E}}^P(t), D^s \dot{\mathbf{E}}^P(t))$, where $*$ denotes the Fenchel transformation. Moreover, we have that \mathcal{F}^* is the indicator function of the set

$$\mathcal{K} := \left\{ (\mathbf{A}, \mathbb{B}, \mathbb{L}) \in L^\infty(\Omega; M_D^{N \times N}) \times L^\infty(\Omega; M_D^{N \times N \times N}) \times (\mathcal{M}_b(\Omega; M_D^{N \times N \times N}))^* : \sqrt{|\mathbf{A}|^2 + l^{-2}|\mathbb{B}|^2} \leq S_Y \text{ a.e. in } \Omega \text{ and } \|\mathbb{L}\|_{\mathcal{M}_b^*} \leq lS_Y \right\}.$$

As a consequence, by (8.29) we deduce that for every $(\mathbf{A}, \mathbb{B}, \mathbb{L}) \in \mathcal{K}$ we have

$$\int_{\Omega} (\mathbf{T}^P(t) - \mathbf{A}) : \dot{\mathbf{E}}^P(t) \, dx + \int_{\Omega} (\mathbb{K}_{\text{diss}}^P(t) - \mathbb{B}) : \nabla \dot{\mathbf{E}}^P(t) \, dx + \langle \mathbb{S}^P(t) - \mathbb{L}, D^s \dot{\mathbf{E}}^P(t) \rangle \geq 0.$$

By choosing $\mathbb{L} = \mathbb{S}^P(t)$, which is possible in view of the constraint (8.10), we obtain (8.30). Inequality (8.31) follows by choosing $\mathbf{A} = \mathbf{T}^P(t)$ and $\mathbb{B} = \mathbb{K}_{\text{diss}}^P(t)$. \square

Let us now prove that the weak flow rule (8.30) for the stresses $\mathbf{T}^P(t)$ and $\mathbb{K}_{\text{diss}}^P(t)$ reduces under suitable regularity assumptions to the usual flow rule given by Gurtin and Anand.

PROPOSITION 8.10 (flow rule). *Let $t \in [0, T]$. Let $\dot{\mathbf{E}}^P(t)$ and $\nabla \dot{\mathbf{E}}^P(t)$ be such that*

$$\sqrt{|\mathbf{T}^P(t, x)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^P(t, x)|^2} < S_Y,$$

$$(8.32) \quad (\dot{\mathbf{E}}^P(t, x), \nabla \dot{\mathbf{E}}^P(t, x)) = (0, 0),$$

$$\sqrt{|\mathbf{T}^P(t, x)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^P(t, x)|^2} = S_Y,$$

$$(8.33) \quad \begin{cases} \mathbf{T}^P(t, x) = S_Y \frac{\dot{\mathbf{E}}^P(t, x)}{\sqrt{|\dot{\mathbf{E}}^P(t, x)|^2 + l^2|\nabla \dot{\mathbf{E}}^P(t, x)|^2}}, \\ \mathbb{K}_{\text{diss}}^P(t, x) = S_Y \frac{l^2 \nabla \dot{\mathbf{E}}^P(t, x)}{\sqrt{|\dot{\mathbf{E}}^P(t, x)|^2 + l^2|\nabla \dot{\mathbf{E}}^P(t, x)|^2}}. \end{cases}$$

Let K be the convex set defined as

$$K := \left\{ (\mathbf{A}, \mathbb{B}) \in M_D^{N \times N} \times M_D^{N \times N \times N} : \sqrt{|\mathbf{A}|^2 + l^{-2}|\mathbb{B}|^2} \leq S_Y \right\}.$$

Let π_K denote the projection onto K , and let π_K^1 and π_K^2 be its components. Let $(\mathbf{A}, \mathbb{B}) \in K$, $\varepsilon > 0$, and let us set

$$\mathcal{C}_{\mathbf{A}, \mathbb{B}}^\varepsilon := (\mathbf{T}^P(t) + \varepsilon(\mathbf{A} - \mathbf{T}^P(t, x)), \mathbb{K}_{\text{diss}}^P(t) + \varepsilon(\mathbb{B} - \mathbb{K}_{\text{diss}}^P(t, x))) \in L^\infty(\Omega; M_D^{N \times N}) \times L^\infty(\Omega; M_D^{N \times N \times N}).$$

For every $r > 0$ let us set

$$\mathbf{F} := \begin{cases} \pi_K^1(\mathcal{C}_{\mathbf{A}, \mathbb{B}}^\varepsilon) & \text{in } B(x, r), \\ \mathbf{T}^P(t) & \text{outside } B(x, r) \end{cases}$$

and

$$\mathbb{G} := \begin{cases} \pi_K^2(C_{\mathbf{A},\mathbb{B}}^\varepsilon) & \text{in } B(x,r), \\ \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t) & \text{outside } B(x,r). \end{cases}$$

Since (\mathbf{F}, \mathbb{G}) are admissible for the weak flow rule (8.30), we obtain

$$\frac{1}{r^N} \left[\int_{B(x,r)} (\mathbf{T}^{\mathbb{P}}(t) - \mathbf{F}) : \dot{\mathbf{E}}^{\mathbb{P}}(t) \, dx + \int_{B(x,r)} (\mathbb{K}_{\text{diss}}^{\mathbb{P}}(t) - \mathbb{G}) : \nabla \dot{\mathbf{E}}^{\mathbb{P}}(t) \, dx \right] \geq 0.$$

Since π_K is a Lipschitz mapping, we have that x is also a Lebesgue point for $\pi_K(C_{\mathbf{A},\mathbb{B}}^\varepsilon)$ with Lebesgue value

$$\pi_K(\mathbf{T}^{\mathbb{P}}(t,x) + \varepsilon(\mathbf{A} - \mathbf{T}^{\mathbb{P}}(t,x)), \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x) + \varepsilon(\mathbb{B} - \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x))).$$

By sending $r \rightarrow 0$ and considering $0 < \varepsilon < 1$, in view of the convexity of K , we obtain

$$(\mathbf{A} - \mathbf{T}^{\mathbb{P}}(t,x)) : \dot{\mathbf{E}}^{\mathbb{P}}(t,x) + (\mathbb{B} - \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x)) : \nabla \dot{\mathbf{E}}^{\mathbb{P}}(t,x) \leq 0.$$

Since the previous inequality holds for every $(\mathbf{A}, \mathbb{B}) \in K$, we deduce that $(\dot{\mathbf{E}}^{\mathbb{P}}(t,x), \nabla \dot{\mathbf{E}}^{\mathbb{P}}(t,x))$ belongs to the normal cone to K at $(\mathbf{T}^{\mathbb{P}}(t,x), \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x))$. In particular, if $(\mathbf{T}^{\mathbb{P}}(t,x), \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x)) \in \text{int}K$, we get that (8.32) holds, while if $(\mathbf{T}^{\mathbb{P}}(t,x), \mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x)) \in \partial K$, (8.33) follows. \square

8.11. Notice that, in view of the presence of a singular part for $D\mathbf{E}^{\mathbb{P}}(t)$ and of its associated stress, plasticity can develop at a point $x \in \Omega$ also when $\|\mathbb{S}^{\mathbb{P}}\|_{\mathcal{M}_b(\Omega; \mathbb{M}_D^{N \times N \times N})} = lS_Y$ and $\sqrt{|\mathbf{T}^{\mathbb{P}}(t,x)|^2 + l^{-2}|\mathbb{K}_{\text{diss}}^{\mathbb{P}}(t,x)|^2} < S_Y$.

9. Asymptotic analysis as $l \rightarrow 0$ and $L \rightarrow 0$. In this section we want to understand the behavior of a quasi-static evolution for the Gurtin–Anand model as the length scales l and L vanish. Our goal is to prove that the quasi-static evolution converges in a suitable sense to an evolution for perfect plasticity. The result is somehow natural, since the strain gradient effects vanish.

More precisely, we prove under suitable assumptions the convergence to a quasi-static evolution for linearly elastic-perfectly plastic bodies recently proposed by Dal Maso, DeSimone, and Mora [8]. The main mathematical problem that we have to face in order to prove such a convergence is that the functional setting of the problem changes, in particular for what concerns the plastic strains. In fact in the strain gradient context, the plastic strain is a BV function (since its gradient enters in the equations), while in [8] it is modeled only as a Radon measure in $\Omega \cup \partial_D \Omega$. Similar problems occur for the displacements, in view of the compatibility condition.

In section 9.1 we briefly recall the model for quasi-static evolution in perfect plasticity recently proposed in [8]. Section 9.2 is devoted to the proof of the convergence result (Theorem 9.2).

9.1. The Dal Maso–DeSimone–Mora model for perfect plasticity. Let us briefly recall the model for quasi-static evolution in perfect plasticity recently proposed in [8]. We formulate the results in the particular form that we need for our asymptotic problem, using the notation of the previous sections.

Let $\Omega \subseteq \mathbb{R}^N$ ($N \geq 2$) be open bounded, let $\partial_D \Omega$ and $\partial_N \Omega$ have the same boundary Γ (relative to $\partial \Omega$), and let us assume that

$$(9.1) \quad \partial \Omega \text{ and } \Gamma \text{ are of class } C^2.$$

Given $w \in W^{1,2}(\Omega; \mathbb{R}^N)$, the class of admissible configurations for the boundary datum w is given by

$$\mathcal{A}_{\text{pp}}(w) := \left\{ (u, \mathbf{E}^e, \mathbf{E}^p) \in BD(\Omega) \times L^2(\Omega; M_{\text{sym}}^{N \times N}) \times \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N}) : \right. \\ \left. \mathbf{E}u = \mathbf{E}^e + \mathbf{E}^p \text{ in } \Omega, \mathbf{E}^p = (w - u) \odot \nu \, d\mathcal{H}^{N-1} \text{ on } \partial_D \Omega \right\}.$$

Here $BD(\Omega)$ denotes the space of functions with bounded deformation on Ω :

$$BD(\Omega) := \left\{ u \in L^1(\Omega; \mathbb{R}^N) : \mathbf{E}u \in \mathcal{M}_b(\Omega; M_{\text{sym}}^{N \times N}) \right\},$$

which is a Banach space with respect to the norm

$$\|u\|_{BD(\Omega)} := \|u\|_{L^1(\Omega; \mathbb{R}^N)} + \|\mathbf{E}u\|_{\mathcal{M}_b(\Omega; M_{\text{sym}}^{N \times N})}.$$

We refer the reader to [33] for the main properties of $BD(\Omega)$. The term $(w - u)$ on $\partial_D \Omega$ is intended in the sense of traces. Finally the subscripts “pp” stand for “perfect plasticity.”

Given $\mathbf{E}^p \in \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N})$, we set

$$\mathcal{H}_{\text{pp}}(\mathbf{E}^p) := S_Y |\mathbf{E}^p|(\Omega \cup \partial_D \Omega),$$

while for $\mathbf{E}^e \in L^2(\Omega; M_{\text{sym}}^{N \times N})$, we consider $\mathcal{Q}_1(\mathbf{E}^e)$ as defined in (4.6).

Let $t \in [0, T]$, and let the boundary displacement be given by

$$(9.2) \quad w \in AC(0, T; W^{1,2}(\Omega; \mathbb{R}^N)).$$

Let the body and traction forces be given by

$$(9.3) \quad f \in AC(0, T; L^N(\Omega; \mathbb{R}^N)) \quad \text{and} \quad g \in AC(0, T; L^\infty(\partial_N \Omega; \mathbb{R}^N)),$$

respectively, and let us denote by $\mathcal{L}(t)$ the associated work as in (4.12). Let us assume that f and g satisfy the uniform safe load condition (4.13)–(4.14). We can simply suppose as in [8] that $t \mapsto \rho(t)$ is absolutely continuous from $[0, T]$ to $L^2(\Omega; M_{\text{sym}}^{N \times N})$, since in view of the regularity of Ω we get $\rho(t) \in L^N(\Omega; M_{\text{sym}}^{N \times N})$ by the embedding result [22, Proposition 2.5].

Given an initial configuration

$$(u_0, \mathbf{E}_0^e, \mathbf{E}_0^p) \in \mathcal{A}_{\text{pp}}(w(0)),$$

a quasi-static evolution $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ in the sense of Dal Maso–DeSimone–Mora [8] is a map from $[0, T]$ to $BD(\Omega) \times L^2(\Omega; M_{\text{sym}}^{N \times N}) \times \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N})$ with $(u(0), \mathbf{E}^e(0), \mathbf{E}^p(0)) = (u_0, \mathbf{E}_0^e, \mathbf{E}_0^p)$ and such that for every $t \in [0, T]$ the following facts hold:

- (a) $(u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}_{\text{pp}}(w(t))$;
- (b) Global stability: For every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$

$$\mathcal{Q}_1(\mathbf{E}^e(t)) - \langle \mathcal{L}(t), u(t) \rangle \leq \mathcal{Q}_1(\mathbf{e}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}_{\text{pp}}(\mathbf{p} - \mathbf{E}^p(t));$$

- (b) Energy balance: The function $t \mapsto \mathbf{E}^p(t)$ has a bounded variation from $[0, T]$ to $\mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N})$ and

$$\mathcal{E}_{\text{pp}}(t) + \mathcal{D}_{\text{pp}}(\mathbf{E}^p; 0, t) = \mathcal{E}_{\text{pp}}(0) + \int_0^t \int_\Omega \mathbf{T}(\tau) : \mathbf{E}\dot{w}(\tau) \, dx \, d\tau \\ - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle \, d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{w}(\tau) \rangle \, d\tau,$$

where $\mathbf{T}(t) := \mathbb{C}\mathbf{E}^e(t)$,

$$\mathcal{E}_{pp}(t) := \mathcal{Q}_1(\mathbf{E}^e(t)) - \langle \mathcal{L}(t), u(t) \rangle$$

and $\mathcal{D}_{pp}(\mathbf{E}^p; 0, t) := S_Y \mathcal{V}(\mathbf{E}^p; 0, t)$.

In order to prove the convergence result of the next section, we need to recall the pairing between stress and strain which gives a useful representation of the work $\mathcal{L}(t)$ similar to (4.15). Following [8, section 2], for every $t \in [0, T]$ and for every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}(w(t))$ it is possible to define the measure $[\rho_D(t) : \mathbf{p}] \in \mathcal{M}_b(\Omega \cup \partial_D \Omega)$ such that

$$(9.4) \quad \langle \mathcal{L}(t), v \rangle = -\langle \rho(t)\nu, w(t) \rangle_{\partial_D \Omega} + \int_{\Omega} \rho(t) : \mathbf{e} \, dx + [\rho_D(t) : \mathbf{p}](\Omega \cup \partial_D \Omega)$$

and such that for every $\varphi \in C^1(\bar{\Omega})$

$$(9.5) \quad \int_{\Omega \cup \partial_D \Omega} \varphi \, d[\rho_D(t) : \mathbf{p}] = \langle \mathcal{L}(t), \varphi v \rangle + \langle \rho(t)\nu, \varphi w(t) \rangle_{\partial_D \Omega} - \int_{\Omega} \rho(t) : \varphi \mathbf{e} \, dx - \int_{\Omega} \rho(t) : [\nabla \varphi \odot v] \, dx.$$

A similar pairing $[\dot{\rho}_D(t) : \mathbf{p}] \in \mathcal{M}_b(\Omega \cup \partial_D \Omega)$ can also be defined (for a.e. $t \in [0, T]$), so that (9.4) and (9.5) hold with $\dot{\rho}_D(t)$, $\dot{\rho}(t)$, and $\dot{\mathcal{L}}(t)$ in place of $\rho_D(t)$, $\rho(t)$, and $\mathcal{L}(t)$.

9.2. The convergence result as $l, L \rightarrow 0$. Let $\Omega \subseteq \mathbb{R}^N$ satisfy (9.1), and let w, f , and g be as in (9.2) and (9.3): Notice that these data are admissible for an evolution for the Gurtin and Anand model. Let us assume that f and g satisfy the uniform safe load condition (4.13)–(4.14).

Let us consider $l_n \rightarrow 0$ and $L_n \rightarrow 0$, and let us denote by

$$t \mapsto (u_n(t), \mathbf{E}_n^e(t), \mathbf{E}_n^p(t))$$

a quasi-static evolution for the Gurtin–Anand model relative to the data w, f , and g and the material length scales $l = l_n$ and $L = L_n$. Let us denote by \mathcal{Q}_2^n , \mathcal{H}_n , and \mathcal{E}_n the energies corresponding to \mathcal{Q}_2 , \mathcal{H} , and \mathcal{E} , respectively.

Let us assume that the initial configuration $(u_n(0), \mathbf{E}_n^e(0), \mathbf{E}_n^p(0))$ is such that there exist $u_0 \in BD(\Omega)$, $\mathbf{E}_0^e \in L^2(\Omega; M_{\text{sym}}^{N \times N})$, and $\mathbf{E}_0^p \in \mathcal{M}_b(\Omega; M_D^{N \times N})$ with

$$(9.6) \quad u_n(0) \xrightarrow{*} u_0 \quad \text{weakly* in } BD(\Omega),$$

$$(9.7) \quad \mathbf{E}_n^e(0) \rightharpoonup \mathbf{E}_0^e \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N}),$$

and

$$(9.8) \quad \mathbf{E}_n^p(0) \xrightarrow{*} \mathbf{E}_0^p \quad \text{weakly* in } \mathcal{M}_b(\Omega; M_D^{N \times N}).$$

Recall that weak* convergence in $BD(\Omega)$ is given by weak convergence in L^1 for the functions and weak* convergence in the sense of measures for the symmetrized gradients.

Let us assume moreover that convergence for the initial free energies holds, that is,

$$(9.9) \quad \mathcal{Q}_1(\mathbf{E}_n^e(0)) + \mathcal{Q}_2^n(\text{curl} \mathbf{E}_n^p(0)) \rightarrow \mathcal{Q}_1(\mathbf{E}_0^e).$$

We have the following compactness result.

LEMMA 9.1. $(u_n(0), \mathbf{E}_n^e(0), \mathbf{E}_n^p(0)) \rightharpoonup^* (u(0), \mathbf{E}^e(0), \mathbf{E}^p(0))$ (9.6) (9.9)

$$u \in AC(0, T; BD(\Omega)), \quad \mathbf{E}^e \in AC(0, T; L^2(\Omega; M_{\text{sym}}^{N \times N})),$$

$$\mathbf{E}^p \in AC(0, T; \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N}))$$

for every $t \in [0, T]$

$$(9.10) \quad u_n(t) \xrightarrow{*} u(t) \quad \text{in } BD(\Omega),$$

$$(9.11) \quad \mathbf{E}_n^e(t) \rightharpoonup \mathbf{E}^e(t) \quad \text{in } L^2(\Omega; M_{\text{sym}}^{N \times N}),$$

$$\mathbf{E}_n^p(t) = 0 \quad \text{in } \partial_D \Omega,$$

$$(9.12) \quad \mathbf{E}_n^p(t) \xrightarrow{*} \mathbf{E}^p(t) \quad \text{in } \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N}).$$

for every $t \in [0, T]$

$$(9.13) \quad (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}_{\text{pp}}(w(t)).$$

Let B be an open ball in \mathbb{R}^N such that $\bar{\Omega} \subseteq B$, and let us set $\tilde{\Omega} := B \setminus \overline{\partial_N \Omega}$. For every $t \in [0, T]$ let us consider $\tilde{u}_n(t) \in W^{1, \frac{N}{N-1}}(\tilde{\Omega}; \mathbb{R}^N)$, $\tilde{\mathbf{E}}_n^e(t) \in L^2(\tilde{\Omega}; M_{\text{sym}}^{N \times N})$, and $\tilde{\mathbf{E}}_n^p(t) \in BV(\tilde{\Omega}; M_D^{N \times N})$ defined as

$$\tilde{u}_n(t) := \begin{cases} u_n(t) & \text{in } \Omega, \\ w(t) & \text{in } \tilde{\Omega} \setminus \Omega, \end{cases}$$

$$\tilde{\mathbf{E}}_n^e(t) := \begin{cases} \mathbf{E}_n^e(t) & \text{in } \Omega, \\ \mathbf{E}w(t) & \text{in } \tilde{\Omega} \setminus \Omega, \end{cases}$$

and

$$\tilde{\mathbf{E}}_n^p(t) := \begin{cases} \mathbf{E}_n^p(t) & \text{in } \Omega, \\ 0 & \text{in } \tilde{\Omega} \setminus \Omega. \end{cases}$$

By (9.6), (9.8), and (9.9) and in view of Remark 8.6 and Proposition 8.5, we deduce that $t \mapsto \tilde{u}_n(t)$, as a map from $[0, T]$ to $BD(\tilde{\Omega})$, has a variation which is uniformly bounded independently on n . More precisely, the sequence $(u_n)_{n \in \mathbb{N}}$ is equiabsolutely continuous. The same holds for $t \mapsto \tilde{\mathbf{E}}_n^e(t)$ and $t \mapsto \tilde{\mathbf{E}}_n^p(t)$ considered as maps from $[0, T]$ to $L^2(\tilde{\Omega}; M_{\text{sym}}^{N \times N})$ and $L^1(\tilde{\Omega}; M_D^{N \times N})$, respectively.

Recall that $BD(\tilde{\Omega})$ can be seen as a dual space, with associated weak* convergence given precisely by the weak* convergence in BD previously defined.

Then by considering $BD(\tilde{\Omega})$ as a dual space and $L^1(\tilde{\Omega}; M_D^{N \times N})$ as a subspace of $\mathcal{M}_b(\tilde{\Omega}; M_D^{N \times N})$, we may apply the generalized version of Helly's theorem [8, Lemma 7.2] to obtain

$$\tilde{u} \in AC(0, T; BD(\tilde{\Omega})), \quad \tilde{\mathbf{E}}^e \in AC(0, T; L^2(\tilde{\Omega}; M_{\text{sym}}^{N \times N})),$$

and

$$\tilde{\mathbf{E}}^P \in AC(0, T; \mathcal{M}_b(\tilde{\Omega}; M_D^{N \times N}))$$

such that, up to a subsequence, for every $t \in [0, T]$

$$(9.14) \quad \tilde{u}_n(t) \overset{*}{\rightharpoonup} \tilde{u}(t) \quad \text{weakly* in } BD(\tilde{\Omega}),$$

$$(9.15) \quad \tilde{\mathbf{E}}_n^e(t) \rightharpoonup \tilde{\mathbf{E}}^e(t) \quad \text{weakly in } L^2(\tilde{\Omega}; M_{\text{sym}}^{N \times N}),$$

and

$$(9.16) \quad \tilde{\mathbf{E}}_n^P(t) \overset{*}{\rightharpoonup} \tilde{\mathbf{E}}^P(t) \quad \text{weakly* in } \mathcal{M}_b(\tilde{\Omega}; M_D^{N \times N}).$$

We have clearly that for every $t \in [0, T]$

$$\tilde{u}(t) = w(t), \quad \tilde{\mathbf{E}}^e(t) = \mathbf{E}w(t), \quad \tilde{\mathbf{E}}^P(t) = 0 \quad \text{on } \tilde{\Omega} \setminus \bar{\Omega}.$$

Let us denote by $u(t)$ and $\mathbf{E}^e(t)$ the restrictions of $\tilde{u}(t)$ and $\tilde{\mathbf{E}}^e(t)$ to Ω , respectively, and let $\mathbf{E}^P(t)$ denote the restriction of $\tilde{\mathbf{E}}^P(t)$ to $\Omega \cup \partial_D \Omega$.

Relations (9.10) and (9.11) follow directly from (9.14) and (9.15). By (9.16), and by taking into account that $\tilde{\mathbf{E}}_n^P(t) = 0$ outside Ω , we obtain (9.12).

From the compatibility condition

$$\tilde{\mathbf{E}}u_n(t) = \tilde{\mathbf{E}}_n^e(t) + \tilde{\mathbf{E}}_n^P(t),$$

we deduce that in the limit we have

$$\tilde{\mathbf{E}}u(t) = \tilde{\mathbf{E}}^e(t) + \tilde{\mathbf{E}}^P(t)$$

so that

$$\mathbf{E}^P(t) \llcorner \partial_D \Omega = \tilde{\mathbf{E}}^P(t) \llcorner \partial_D \Omega = (w(t) - u(t)) \odot \nu \, d\mathcal{H}^{N-1} \llcorner \partial_D \Omega,$$

where $u(t)$ is intended in the sense of traces on $\partial_D \Omega$. We deduce that (9.13) holds, and the proof is concluded. \square

The main theorem of the section is the following asymptotic result.

THEOREM 9.2. *Let $t \mapsto (u_{l,L}(t), \mathbf{E}_{l,L}^e(t), \mathbf{E}_{l,L}^P(t))$ be a family of functions satisfying (9.6)*

$$(9.9) \quad \begin{aligned} & l, L \rightarrow 0, \quad l_n \rightarrow 0, \quad L_n \rightarrow 0, \quad (l_n, L_n)_{j \in \mathbb{N}} \rightarrow (l, L) \quad \text{in } L^2(\Omega; M_{\text{sym}}^{N \times N}), \\ & t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^P(t)) \end{aligned} \quad [8]$$

$$u_j := u_{l_n, L_n}, \quad \mathbf{E}_j^e := \mathbf{E}_{l_n, L_n}^e, \quad \mathbf{E}_j^P := \mathbf{E}_{l_n, L_n}^P$$

for every $t \in [0, T]$.

$$(9.17) \quad u_j(t) \overset{*}{\rightharpoonup} u(t) \quad \text{weakly* in } BD(\Omega),$$

$$(9.18) \quad \mathbf{E}_j^e(t) \rightharpoonup \mathbf{E}^e(t) \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N}),$$

$$(9.19) \quad \mathbf{E}_j^P(t) \overset{*}{\rightharpoonup} \mathbf{E}^P(t) \quad \text{weakly* in } \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N}).$$

for every $t \in [0, T]$

$$(9.20) \quad \mathcal{Q}_1(\mathbf{E}_j^e(t)) \rightarrow \mathcal{Q}_1(\mathbf{E}^e(t)) \quad \text{and} \quad \mathcal{Q}_2^{n_j}(\text{curl} \mathbf{E}_j^p(t)) \rightarrow 0,$$

We divide the proof into three steps.

1: By Lemma 9.1 there exist a subsequence n_j ,

$$u \in AC(0, T; BD(\Omega)), \quad \mathbf{E}^e \in AC(0, T; L^2(\Omega; M_{\text{sym}}^{N \times N})),$$

and

$$\mathbf{E}^p \in AC(0, T; \mathcal{M}_b(\Omega \cup \partial_D \Omega; M_D^{N \times N}))$$

such that by setting $u_j := u_{n_j}$, $\mathbf{E}_j^e := \mathbf{E}_{n_j}^e$, and $\mathbf{E}_j^p := \mathbf{E}_{n_j}^p$, for every $t \in [0, T]$ relations (9.17) and (9.19) hold,

$$(9.21) \quad \mathbf{E}_j^e(t) \rightharpoonup \mathbf{E}^e(t) \quad \text{weakly in } L^2(\Omega; M_{\text{sym}}^{N \times N}),$$

and $(u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}_{\text{pp}}(w(t))$, so that the triple $(u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ is admissible. Finally, from the energy balance (5.3), and by the assumptions for $t = 0$, we deduce that for every $t \in [0, T]$

$$(9.22) \quad \mathcal{Q}_2^{n_j}(\text{curl} \mathbf{E}_j^p(t)) \leq C$$

for some constant C independent of j and t .

2: Let us fix $t \in [0, T]$. In order to prove that $(u(t), \mathbf{E}^e(t), \mathbf{E}^p(t)) \in \mathcal{A}_{\text{pp}}(w(t))$ satisfies the global stability condition

$$(9.23) \quad \mathcal{Q}_1(\mathbf{E}^e(t)) - \langle \mathcal{L}(t), u(t) \rangle \leq \mathcal{Q}_1(\mathbf{e}) - \langle \mathcal{L}(t), v \rangle + \mathcal{H}_{\text{pp}}(\mathbf{p} - \mathbf{E}^p(t))$$

for every $(v, \mathbf{e}, \mathbf{p}) \in \mathcal{A}_{\text{pp}}(w(t))$, in view of [8, Theorem 3.6] it suffices to prove that the Cauchy stress $\mathbf{T}(t) = \mathbb{C} \mathbf{E}^e(t)$ satisfies the equilibrium conditions

$$(9.24) \quad \begin{cases} -\text{div} \mathbf{T}(t) = f(t) & \text{in } \Omega, \\ \mathbf{T}(t) \nu = g(t) & \text{on } \partial_N \Omega \end{cases}$$

and the constraint

$$(9.25) \quad |\mathbf{T}_D(t, x)| \leq S_Y \quad \text{for a.e. } x \in \Omega,$$

where $\mathbf{T}_D(t) := (\mathbf{T}(t))_D$.

Equation (9.24) follows from the equilibrium equation for the Cauchy stress $\mathbf{T}_j(t) = \mathbb{C} \mathbf{E}_j^e(t)$ given by (8.6) in view of the weak convergence of $\mathbf{T}_j(t)$ to $\mathbf{T}(t)$ which comes from (9.21).

In order to prove (9.25), let us consider the corresponding constraint in the strain gradient context given by (8.9). Let $\mathbb{K}_{\text{diss},j}^p(t)$, $\mathbb{K}_{\text{en},j}^p(t)$, $\mathbb{K}_j^p(t) = \mathbb{K}_{\text{diss},j}^p(t) + \mathbb{K}_{\text{en},j}^p(t)$ and $\mathbf{T}_j^p(t)$ be the stresses associated to $(u_j(t), \mathbf{E}_j^e(t), \mathbf{E}_j^p(t))$. Notice that, in view of (3.5) and of (9.22), we get

$$(9.26) \quad \mathbb{K}_{\text{en},j}^p(t) \rightarrow 0 \quad \text{strongly in } L^2(\Omega; M_D^{N \times N \times N}).$$

Moreover by (8.8) and (8.9) we have

$$(9.27) \quad \mathbf{T}_j^p(t) = (\mathbf{T}_j(t))_D + \operatorname{div} \mathbb{K}_j^p(t),$$

and

$$\sqrt{|\mathbf{T}_j^p(t, x)|^2 + l_{n_j}^{-2} |\mathbb{K}_{\operatorname{diss}, j}^p(t, x)|^2} \leq S_Y \quad \text{for a.e. } x \in \Omega.$$

In particular we have that $(\mathbf{T}_j^p(t))_{j \in \mathbb{N}}$ is uniformly bounded in $L^\infty(\Omega; M_D^{N \times N})$ and

$$\mathbb{K}_{\operatorname{diss}, j}^p(t) \rightarrow 0 \quad \text{strongly in } L^\infty(\Omega; M_D^{N \times N \times N}).$$

By (9.26) we conclude that $\mathbb{K}_j^p(t) \rightarrow 0$ strongly in $L^2(\Omega; M_D^{N \times N \times N})$. Notice that from (9.27) we deduce that $\operatorname{div} \mathbb{K}_j^p(t)$ is bounded in $L^2(\Omega; M_D^{N \times N})$. We obtain

$$\operatorname{div} \mathbb{K}_j^p(t) \rightharpoonup 0 \quad \text{weakly in } L^2(\Omega; M_D^{N \times N})$$

so that in view of (9.27) and (9.21)

$$(9.28) \quad \mathbf{T}_j^p(t) \rightharpoonup \mathbf{T}_D(t) \quad \text{weakly in } L^2(\Omega; M_D^{N \times N}).$$

Since $\mathbf{T}_j^p(t) \in \mathcal{K} := \{\mathbf{A} \in L^2(\Omega; M_D^{N \times N}) : |\mathbf{A}| \leq S_Y \text{ a.e. in } \Omega\}$, and \mathcal{K} is closed in the weak topology of $L^2(\Omega; M_D^{N \times N})$, by (9.28) we deduce that (9.25) holds. Hence (9.23) follows, and Step 2 is concluded.

3. Since $t \mapsto u_j(t)$ is absolutely continuous from $[0, T]$ to $W^{1, \frac{N}{N-1}}(\Omega; \mathbb{R}^N)$, by integrating by parts we can write the energy balance (5.3) in the following form:

$$(9.29) \quad \begin{aligned} \mathcal{Q}_1(\mathbf{E}_j^e(t)) + \mathcal{Q}_2^{n_j}(\operatorname{curl} \mathbf{E}_j^p(t)) + \mathcal{D}_{\mathcal{H}_{n_j}}(\mathbf{E}_j^p; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}_j(\tau) \rangle d\tau \\ = \mathcal{Q}_1(\mathbf{E}_j^e(0)) + \mathcal{Q}_2^{n_j}(\operatorname{curl} \mathbf{E}_j^p(0)) + \int_0^t \int_\Omega \mathbf{T}_j(\tau) : \mathbf{E} \dot{u}(\tau) dx d\tau \\ - \int_0^t \langle \mathcal{L}(\tau), \dot{u}(\tau) \rangle d\tau. \end{aligned}$$

We claim that for every $t \in [0, T]$

$$(9.30) \quad \begin{aligned} \liminf_{j \rightarrow +\infty} \left[\mathcal{D}_{\mathcal{H}_{n_j}}(\mathbf{E}_j^p; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}_j(\tau) \rangle d\tau \right] \geq \mathcal{D}_{\mathcal{H}_{pp}}(\mathbf{E}^p; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}(\tau) \rangle d\tau \\ = \mathcal{D}_{\mathcal{H}_{pp}}(\mathbf{E}^p; 0, t) - \langle \mathcal{L}(t), u(t) \rangle + \langle \mathcal{L}(0), u(0) \rangle + \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle d\tau. \end{aligned}$$

By passing to the limit in (9.29), by (9.21), (9.9), and (9.30) we get for every $t \in [0, T]$

$$\begin{aligned} \mathcal{Q}_1(\mathbf{E}^e(t)) - \langle \mathcal{L}(t), u(t) \rangle + \mathcal{D}_{\mathcal{H}_{pp}}(\mathbf{E}^p; 0, t) \leq \mathcal{Q}_1(\mathbf{E}^e(0)) - \langle \mathcal{L}(0), u(0) \rangle \\ + \int_0^t \int_\Omega \mathbf{T}(\tau) : \mathbf{E} \dot{u}(\tau) dx d\tau - \int_0^t \langle \dot{\mathcal{L}}(\tau), u(\tau) \rangle d\tau - \int_0^t \langle \mathcal{L}(\tau), \dot{u}(\tau) \rangle d\tau. \end{aligned}$$

In view of the global stability condition (9.23), by [8, Theorem 4.7] we have that also the opposite inequality holds, so that the energy balance follows. From the previous steps, we conclude that $t \mapsto (u(t), \mathbf{E}^e(t), \mathbf{E}^p(t))$ is a quasi-static evolution according to Dal Maso, DeSimone, and Mora [8].

By (9.29), (9.30), and the energy balance, we deduce that for every $t \in [0, T]$ we have

$$\mathcal{Q}_1(\mathbf{E}_j^e(t)) \rightarrow \mathcal{Q}_1(\mathbf{E}^e(t)) \quad \text{and} \quad \mathcal{Q}_2^{n_j}(\text{curl} \mathbf{E}_j^p(t)) \rightarrow 0$$

so that (9.20) holds. In view of (9.21), we conclude that (9.18) follows.

Let us prove claim (9.30). Recall that by [8, Theorem 7.1] we have the following representation of the dissipation:

$$\mathcal{D}_{\mathcal{H}_{n_j}}(\mathbf{E}_j^p; 0, t) = \int_0^t \mathcal{H}_{n_j}(\dot{\mathbf{E}}_j^p(\tau)) d\tau.$$

From the representation (4.15) we get

$$(9.31) \quad \begin{aligned} \mathcal{D}_{\mathcal{H}_{n_j}}(\mathbf{E}_j^p; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}_j(\tau) \rangle d\tau &= \int_0^t \left[\mathcal{H}_{n_j}(\dot{\mathbf{E}}_j^p(\tau)) - \int_{\Omega} \rho_D(\tau) : \dot{\mathbf{E}}_j^p(\tau) dx \right] d\tau \\ &\quad - \int_0^t \int_{\Omega} \rho(\tau) : \dot{\mathbf{E}}_j^e(\tau) dx d\tau + \int_0^t \langle \rho(\tau)\nu, \dot{w}(\tau) \rangle_{\partial_D \Omega} d\tau. \end{aligned}$$

Moreover for every $0 \leq \tau \leq t$

$$\mathcal{H}_{n_j}(\dot{\mathbf{E}}_j^p(\tau)) - \int_{\Omega} \rho_D(\tau) : \dot{\mathbf{E}}_j^p(\tau) dx \geq \int_{\Omega} \left[S_Y |\dot{\mathbf{E}}_j^p(\tau)| - \rho_D(\tau) : \dot{\mathbf{E}}_j^p(\tau) \right] dx,$$

and the integrand of the right-hand side is positive in view of the safe load condition (4.14). Let $\varphi \in C^1(\bar{\Omega})$, with $0 \leq \varphi \leq 1$ and $\varphi = 0$ near $\bar{\partial}_N \Omega$. By applying again the representation result [8, Theorem 7.1] for the dissipation $\mathcal{D}_{\mathcal{H}_{pp}}$ we conclude that

$$(9.32) \quad \begin{aligned} \liminf_{j \rightarrow +\infty} \int_0^t \left[\mathcal{H}_{n_j}(\dot{\mathbf{E}}_j^p(\tau)) - \int_{\Omega} \rho_D(\tau) : \dot{\mathbf{E}}_j^p(\tau) dx \right] \\ \geq \liminf_{j \rightarrow +\infty} \int_0^t \int_{\Omega} \left[S_Y |\dot{\mathbf{E}}_j^p(\tau)| - \rho_D(\tau) : \dot{\mathbf{E}}_j^p(\tau) \right] dx d\tau \\ \geq \liminf_{j \rightarrow +\infty} \int_0^t \int_{\Omega} \left[S_Y |\varphi \dot{\mathbf{E}}_j^p(\tau)| - \rho_D(\tau) : \varphi \dot{\mathbf{E}}_j^p(\tau) \right] dx d\tau \\ = \liminf_{j \rightarrow +\infty} \left[\mathcal{D}_{\mathcal{H}_{pp}}(\varphi \mathbf{E}_j^p; 0, t) - \int_0^t \int_{\Omega} \rho_D(\tau) : \varphi \dot{\mathbf{E}}_j^p(\tau) dx d\tau \right]. \end{aligned}$$

By the very definition of $\mathcal{D}_{\mathcal{H}_{pp}}$ and by (9.19), it is easy to see that

$$(9.33) \quad \liminf_{j \rightarrow +\infty} \mathcal{D}_{\mathcal{H}_{pp}}(\varphi \mathbf{E}_j^p; 0, t) \geq \mathcal{D}_{\mathcal{H}_{pp}}(\varphi \mathbf{E}^p; 0, t).$$

On the other hand, the absolute continuity of $t \mapsto \mathbf{E}_j^p(t)$ implies that

$$(9.34) \quad \begin{aligned} \int_0^t \int_{\Omega} \rho_D(\tau) : \varphi \dot{\mathbf{E}}_j^p(\tau) dx d\tau &= \int_{\Omega} \rho_D(t) : \varphi \mathbf{E}_j^p(t) dx - \int_{\Omega} \rho_D(0) : \varphi \mathbf{E}_j^p(0) dx \\ &\quad - \int_0^t \int_{\Omega} \dot{\rho}_D(\tau) : \varphi \mathbf{E}_j^p(\tau) dx d\tau. \end{aligned}$$

By integrating by parts, for a.e. $\tau \in [0, t]$ we have

$$\begin{aligned} \int_{\Omega} \dot{\rho}_D(\tau) : \varphi \mathbf{E}_j^P(\tau) \, dx &= \langle \dot{\mathcal{L}}(\tau), \varphi u_j(\tau) \rangle + \langle \dot{\rho}(\tau) \nu, w(\tau) \rangle_{\partial_D \Omega} \\ &\quad - \int_{\Omega} \dot{\rho}(\tau) : \varphi \mathbf{E}_j^e(\tau) \, dx - \int_{\Omega} \dot{\rho}(\tau) : [\nabla \varphi \odot u_j(\tau)] \, dx. \end{aligned}$$

In view of the embedding result [22, Proposition 2.5], we get $\dot{\rho}(\tau) \in L^N(\Omega; M_{\text{sym}}^{N \times N})$ for a.e. $\tau \in [0, t]$. By (9.21) and (9.17), and since $\varphi = 0$ near $\partial_N \Omega$, we deduce for a.e. $\tau \in [0, t]$ that

$$\begin{aligned} (9.35) \quad \lim_{j \rightarrow +\infty} \int_{\Omega} \dot{\rho}_D(\tau) : \varphi \mathbf{E}_j^P(\tau) \, dx &= \langle \dot{\mathcal{L}}(\tau), \varphi u(\tau) \rangle + \langle \dot{\rho}(\tau) \nu, w(\tau) \rangle_{\partial_D \Omega} \\ &\quad - \int_{\Omega} \dot{\rho}(\tau) : \varphi \mathbf{E}^e(\tau) \, dx - \int_{\Omega} \dot{\rho}(\tau) : [\nabla \varphi \odot u(\tau)] \, dx \\ &= \int_{\Omega \cup \partial_D \Omega} \varphi \, d[\dot{\rho}_D(\tau) : \mathbf{E}^P(\tau)], \end{aligned}$$

where $[\dot{\rho}_D(\tau) : \mathbf{E}^P(\tau)]$ is the measure defined in the previous subsection, and the last equality follows by (9.5) (with $\dot{\rho}_D$ in place of ρ_D). Similarly we obtain

$$(9.36) \quad \lim_{j \rightarrow +\infty} \int_{\Omega} \rho_D(t) : \varphi \mathbf{E}_j^P(t) \, dx = \int_{\Omega \cup \partial_D \Omega} \varphi \, d[\rho_D(t) : \mathbf{E}^P(t)]$$

and

$$(9.37) \quad \lim_{j \rightarrow +\infty} \int_{\Omega} \rho_D(0) : \varphi \mathbf{E}_j^P(0) \, dx = \int_{\Omega \cup \partial_D \Omega} \varphi \, d[\rho_D(0) : \mathbf{E}^P(0)].$$

By letting $\varphi \rightarrow 1_{\Omega \cup \partial_D \Omega}$ we obtain from (9.32), (9.33), (9.34), and (9.35)–(9.37)

$$\begin{aligned} (9.38) \quad \liminf_{j \rightarrow +\infty} \int_0^t \left[\mathcal{H}_{n_j}(\dot{\mathbf{E}}_j^P(\tau)) - \int_{\Omega} \rho_D(\tau) : \dot{\mathbf{E}}_j^P(\tau) \, dx \right] &\geq \mathcal{D}_{\mathcal{H}_{\text{pp}}}(\mathbf{E}^P; 0, t) \\ - [\rho_D(t) : \mathbf{E}^P(t)](\Omega) + [\rho_D(0) : \mathbf{E}^P(0)](\Omega \cup \partial_D \Omega) &+ \int_0^t [\dot{\rho}_D(\tau) : \mathbf{E}^P(\tau)](\Omega \cup \partial_D \Omega) \, d\tau \\ &= \mathcal{D}_{\mathcal{H}_{\text{pp}}}(\mathbf{E}^P; 0, t) - \int_0^t [\rho_D(\tau) : \dot{\mathbf{E}}^P(\tau)](\Omega \cup \partial_D \Omega) \, d\tau. \end{aligned}$$

In conclusion, by passing to the limit in (9.31), by (9.38) and (9.21) we get

$$\begin{aligned} \liminf_{j \rightarrow +\infty} \left[\mathcal{D}_{\mathcal{H}_{n_j}}(\mathbf{E}_j^P; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}_j(\tau) \rangle \, d\tau \right] \\ \geq \mathcal{D}_{\mathcal{H}_{\text{pp}}}(\mathbf{E}^P; 0, t) - \int_0^t [\rho_D(\tau) : \dot{\mathbf{E}}^P(\tau)](\Omega \cup \partial_D \Omega) \, d\tau - \int_0^t \int_{\Omega} \rho(\tau) : \dot{\mathbf{E}}^e(\tau) \, dx \, d\tau \\ + \int_0^t \langle \rho(\tau) \nu, \dot{w}(\tau) \rangle_{\partial_D \Omega} \, d\tau = \mathcal{D}_{\mathcal{H}_{\text{pp}}}(\mathbf{E}^P; 0, t) - \int_0^t \langle \mathcal{L}(\tau), \dot{u}(\tau) \rangle \, d\tau, \end{aligned}$$

where the last equality comes from the integration by parts (9.4). We deduce that claim (9.30) holds, and the proof is concluded. \square

9.3. It is interesting to compare our result with the abstract convergence result recently proposed by Mielke, Roubíček, and Stefanelli [26] concerning the asymptotic behavior of quasi-static evolutions with Γ -converging energies and dissipations. Their framework cannot be easily adapted to our context since the total energy and the dissipation functional do not satisfy trivially a sort of Γ -liminf inequality. This is due to the fact that we skip from a BV to a Radon-measure setting for the plastic strain, so that concentration at the boundary $\partial_D\Omega$ can occur, and this has to be taken into account for the dissipation \mathcal{H}_{pp} and the work of the external loads. The safe load condition is the key point to rearrange the terms in order to work out suitable lower semicontinuity inequalities (see the arguments of Step 3) which are essential for the study of the asymptotics of the problem.

REFERENCES

- [1] A. ACHARYA AND J. L. BASSANI, *Incompatibility and crystal plasticity*, J. Mech. Phys. Solids, 48 (2000), pp. 1565–1595.
- [2] L. AMBROSIO, N. FUSCO, AND D. PALLARA, *Functions of Bounded Variations and Free Discontinuity Problems*, Clarendon Press, Oxford, 2000.
- [3] E. C. AIFANTIS, *On the microstructural origin of certain inelastic models*, ASME J. Eng. Mater. Technol., 106 (1984) pp. 326–330.
- [4] M. F. ASHBY, *The deformation of plastically non-homogeneous alloys*, Philos. Mag., 21 (1970), pp. 399–424.
- [5] M. F. ASHBY, *The deformation of plastically non-homogeneous alloys*, in Strengthening Methods in Crystals, A. Kelly and R. B. Nicholson, eds., Elsevier, Amsterdam, 1971, pp. 137–192.
- [6] H. BREZIS, *Operateurs Maximaux Monotones et Semi-Groups de Contractions dans les Espaces de Hilbert*, North-Holland, Amsterdam-London; Elsevier, New York, 1973.
- [7] P. CERMELLI AND M. E. GURTIN, *On the characterization of geometrically necessary dislocations in finite plasticity*, J. Mech. Phys. Solids, 49 (2000), pp. 1539–1568.
- [8] G. DAL MASO, A. DESIMONE, AND M. G. MORA, *Quasistatic evolution problems for linearly elastic-perfectly plastic materials*, Arch. Ration. Mech. Anal., 180 (2006), pp. 237–291.
- [9] L. C. EVANS AND R. F. GARIEPY, *Measure Theory and Fine Properties of Functions*, Stud. Adv. Math., CRC Press, Boca Raton, FL, 1992.
- [10] N. A. FLECK AND J. W. HUTCHINSON, *Strain gradient plasticity*, Adv. Appl. Mech., 33 (1997), pp. 295–361.
- [11] N. A. FLECK AND J. W. HUTCHINSON, *A reformulation of strain gradient plasticity*, J. Mech. Phys. Solids, 49 (2001), pp. 2245–2271.
- [12] H. GAO, Y. HUANG, W. D. NIX, AND J. W. HUTCHINSON, *Mechanism-based strain gradient plasticity, I. Theory*, J. Mech. Phys. Solids, 47 (1999), pp. 1239–1263.
- [13] C. GOFFMAN AND J. SERRIN, *Sublinear functions of measures and variational integrals*, Duke Math. J., 31 (1964), pp. 159–178.
- [14] P. GUDMUNDSON, *A unified treatment of strain gradient plasticity*, J. Mech. Phys. Solids, 52 (2004), pp. 1379–1406.
- [15] M. E. GURTIN, *On the plasticity of single crystals: Free energy, microforces, plastic strain gradients*, J. Mech. Phys. Solids, 48 (2000), pp. 989–1036.
- [16] M. E. GURTIN, *A gradient theory of single-crystal viscoplasticity that accounts for geometrically necessary dislocations*, J. Mech. Phys. Solids, 50 (2002), pp. 5–32.
- [17] M. E. GURTIN, *On a framework for small-deformation viscoplasticity: Free energy, microforces, strain gradients*, Int. J. Plasticity, 19 (2003), pp. 47–90.
- [18] M. E. GURTIN, *A gradient theory of small-deformation isotropic plasticity that accounts for the Burgers vector and for dissipation due to plastic spin*, J. Mech. Phys. Solids, 52 (2004), pp. 2545–2568.
- [19] M. E. GURTIN AND L. ANAND, *A theory of strain-gradient plasticity for isotropic, plastically irrotational materials, I. Small deformations*, J. Mech. Phys. Solids, 53 (2005), pp. 1624–1649.
- [20] M. E. GURTIN AND A. NEEDLEMAN, *Boundary conditions in small-deformation, single-crystal plasticity that account for the Burgers vector*, J. Mech. Phys. Solids, 53 (2005), pp. 1–31.
- [21] Y. HUANG, H. GAO, W. D. NIX, AND J. W. HUTCHINSON, *Mechanism-based strain gradient plasticity-II, Analysis*, J. Mech. Phys. Solids, 48 (2000), pp. 99–128.

- [22] R. KOHN AND R. TEMAM, *Dual spaces of stresses and strains, with applications to Hencky plasticity*, Appl. Math. Optim., 10 (1983), pp. 1–35.
- [23] A. MAINIK AND A. MIELKE, *Existence results for energetic models for rate-independent systems*, Calc. Var. Partial Differential Equations, 22 (2005), pp. 73–99.
- [24] A. MIELKE, *Evolution of rate-independent systems*, in Evolutionary Equations, Handb. Differ. Equ. II, Elsevier/North-Holland, Amsterdam, 2005, pp. 461–559.
- [25] A. MIELKE, *Analysis of energetic models for rate-independent materials*, in Proceedings of the International Congress of Mathematicians, Vol. III (Beijing, 2002), Higher Ed. Press, Beijing, 2002, pp. 817–828.
- [26] A. MIELKE, T. ROUBÍČEK AND U. STEFANELLI, *Γ -limits and relaxations for rate-independent evolutionary problems*, Calc. Var. Partial Differential Equations, 31 (2008), pp. 387–416.
- [27] A. MIELKE AND F. THEIL, *A mathematical model for rate independent phase transformations with hysteresis*, in Proceedings of the Workshop on “Models of Continuum Mechanics in Analysis and Engineering,” H.-D. Alber, R. Balean, and R. Farwig, eds., Shaker-Verlag, Aachen, 1999, pp. 117–129.
- [28] A. MIELKE AND F. THEIL, *On rate-independent hysteresis models*, NoDEA Nonlinear Differential Equations Appl., 11 (2004), pp. 151–189.
- [29] A. MIELKE, F. THEIL, AND V. LEVITAS, *A variational formulation of rate-independent phase transformations using an extremum principle*, Arch. Ration. Mech. Anal., 162 (2002), pp. 137–177.
- [30] J. F. NYE, *Some geometrical relations in dislocated crystals*, Acta Metall., 1 (1953), pp. 153–162.
- [31] B. D. REDDY, F. EBOBISSE, AND A. MCBRIDE, *Well-posedness of a model of strain gradient plasticity for plastically irrotational materials*, Int. J. Plasticity, 24 (2008), pp. 55–73.
- [32] P.-M. SUQUET, *Sur les équations de la plasticité: Existence et régularité des solutions*, [On the equations of plasticity: Existence and regularity of solutions], J. Mécanique, 20 (1981), pp. 3–39.
- [33] R. TEMAM, *Problèmes Mathématiques en Plasticité*, Méthodes Mathématiques de l’Informatique 12, Gauthier-Villars, Montrouge, 1983.

ON THE NEW MULTISCALE RODLIKE MODEL OF POLYMERIC FLUIDS*

HUI ZHANG[†] AND PINGWEN ZHANG[‡]

Abstract. This paper is concerned with the well-posedness for the new rigid rodlike model in a polymeric fluid recently proposed by W.N. E and P.W. Zhang [*Meth. Appl. Anal.*, 13 (2006), pp. 181–198]. The constitutive relations considered in this work are motivated by the kinetic theory. The micro equations involve five independent spatial variables (degrees of freedom): two in the configuration domain and three in the macro flow domain. We obtain the local existence of solutions with large initial data and also global existence of solutions with small Deborah and Reynolds constants in periodic domains.

Key words. polymeric fluid, rodlike model, kinetic theory, global existence

AMS subject classifications. 76B03, 65M12, 35Q35

DOI. 10.1137/050640795

1. Introduction. The Doi kinetic theory for spatially homogeneous flow of rodlike molecules has successfully described the properties of liquid crystal polymers in a solvent [6]. One of the simplest models of polymeric fluids described by the Doi theory is the rigid rodlike model, which takes into account the macro and micro behavior of the dilute or solute polymers—the effects of flow, Brownian motion, and intermolecular forces on the molecular orientation distribution (see [6, 10, 11, 16, 17, 23]). However, it does not include the so-called distortional elasticity. The Doi theory is valid only in the limit of spatial homogeneity. For small molecule liquid crystals, distortional elasticity has been formulated in the limit of weak distribution as Frank elasticity. This is one ingredient of the Leslie–Ericksen theory. Several attempts have been made to find a theory which encompasses both the molecular visco-elasticity and the distortional elasticity. Marrucci and Greco [18] made a molecular theory of distortional elasticity. They proposed a nonlocal mean field nematic potential for LCPs (liquid crystal polymers), which accounts for spatial variations of the molecular orientation distribution. Tsuji and Rey [21, 22] add distortional elasticity to the rodlike model of the Doi theory but did not give a stress tensor. Edwards and Beris [2] give an ad hoc generalization of the Frank elasticity in tensorial form. Ericksen [9] allowed the order parameter to be a variable but still required the orientation distribution to be uniaxial. An extension of Kuzuu and Doi [13] theory to flowing systems of nonhomogeneous liquid crystalline polymers is made by Wang [25], in which the author models the LCP molecules as spheroids of equal shape and size. He derives an intermolecular potential which could be considered as an extension of Marrucci–Greco potential.

*Received by the editors September 21, 2005; accepted for publication (in revised form) July 15, 2008; published electronically October 29, 2008.

<http://www.siam.org/journals/sima/40-3/64079.html>

[†]Laboratory of Mathematics and Complex Systems, Ministry of Education, School of Mathematical Sciences, Beijing Normal University, Beijing, 100875, P.R. China (hzhang@bnu.edu.cn). This author is partially supported by an Alexander von Humboldt Fellowship and is also partially supported by the key basic research project of the Ministry of Chinese Education 107016 and the state key basic research project of China 2005CB321704.

[‡]LMAM and School of Mathematical Sciences, Peking University, Beijing, 100871, P.R. China (pzhang@pku.edu.cn). This author's research is partially supported by the state key basic research project of China 2005CB321704 and the National Science Foundation of China for Distinguished Young Scholars 10225103.

All these approaches are phenomenological in nature. They invariably contain a large number of unknown parameters which in general cannot be determined rationally. Another drawback of the phenomenological theories is the lack of consistency with existing theories and among themselves. Then E and Zhang [8] developed a new model for nonhomogeneous flows of liquid crystalline polymers with a few adjustable parameters that could model a variety of configurations of polymeric liquid crystal molecules. This new model is a combination of macroscopic partial differential equations and microscopic Fokker–Planck equations. In this model, the function $\psi(\mathbf{x}, \mathbf{m}, t)$ describes the distribution of an identical rigid rodlike molecule at (\mathbf{x}, t) with the orientation \mathbf{m} . Denoting the velocity and pressure of the fluid by \mathbf{u} and p , the new multiscale rodlike model can be expressed as

$$(1.1) \quad \frac{\partial \psi}{\partial t} + \nabla \cdot (\mathbf{u}\psi) = \frac{1}{k_B T} \nabla \cdot \{ [D_{\parallel} \mathbf{m}\mathbf{m} + D_{\perp}(\mathbf{I} - \mathbf{m}\mathbf{m})] \cdot (\psi \nabla \mu) \} + \frac{D_r}{k_B T} \mathcal{R} \cdot (\psi \mathcal{R} \mu) - \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \psi), \quad \mathbf{m} \in \mathbb{S}^2,$$

where k_B is Boltzmann constant and T is the absolute temperature, $D_{\parallel} \geq 0$ and $D_{\perp} \geq 0$ are translational diffusion coefficients parallel and normal to the orientation of the LCP molecule, $D_r = \frac{\xi_r}{k_B T}$ is the rotary diffusivity and ξ_r is the friction coefficient, ∇ is the gradient operator with respect to the spatial variables \mathbf{x} , $\mathcal{R} = \mathbf{m} \times \frac{\partial}{\partial \mathbf{m}}$ is the rotational gradient operator, and \mathbb{S}^2 is the unit sphere in \mathbb{R}^3 . The symbol μ denotes the chemical potential

$$(1.2) \quad \mu = \ln \psi + \bar{U},$$

and \bar{U} represents the excluded-volume potential [6, 11]

$$(1.3) \quad \bar{U}(\mathbf{x}, \mathbf{m}, t) = k_B T \alpha \int_{\Omega} \int_{|\mathbf{m}'|=1} B(\mathbf{x}, \mathbf{x}', \mathbf{m}, \mathbf{m}') \psi(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}'.$$

The function B in (1.3) is the interactional factor among rods. Here α denotes the intensity between particles. Now we choose $B(\mathbf{x}, \mathbf{x}', \mathbf{m}, \mathbf{m}') = \frac{1}{\varepsilon^3} \chi(\frac{\mathbf{x}-\mathbf{x}'}{\varepsilon}) |\mathbf{m} \times \mathbf{m}'|^2$, where $\chi(\mathbf{x})$ is the smooth kernel; e.g., $\chi(\mathbf{x}) = C \exp(1/(|\mathbf{x}|^2 - 1))$ as $|\mathbf{x}| < 1$, and $\chi(\mathbf{x}) = 0$ as $|\mathbf{x}| \geq 1$, where C is a constant such that $\int_{|\mathbf{x}| \leq 1} \chi(\mathbf{x}) d\mathbf{x} = 1$. $\kappa = (\nabla \mathbf{u})^T$ is the velocity gradient tensor.

Let L_0 be the typical size of the flow region, V_0 be the typical velocity scale, and $T_0 = \frac{L_0}{V_0}$ be a typical convective time scale. Further De is an important parameter called the Deborah number:

$$(1.4) \quad De = \frac{\frac{\xi_r}{k_B T}}{\frac{L_0}{V_0}} = \frac{\xi_r V_0}{k_B T L_0}.$$

It is the ratio of the orientational diffusion time scale of the rods (which is the relevant relaxation time scale) and the convective time scale of the fluid. Set

$$(1.5) \quad \varepsilon = \frac{L}{L_0},$$

where L is the length of the rods. Thus the nondimensional kinetic equation is

$$(1.6) \quad \frac{\partial \psi}{\partial t} + \nabla \cdot (\mathbf{u}\psi) = \frac{\varepsilon^2}{De} \nabla \cdot \left\{ [D_{\parallel}^* \mathbf{m}\mathbf{m} + D_{\perp}^* (\mathbf{I} - \mathbf{m}\mathbf{m})] \cdot (\psi \nabla \mu) \right\} + \frac{1}{De} \mathcal{R} \cdot (\psi \mathcal{R} \mu) - \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \psi), \quad \mathbf{m} \in \mathbb{S}^2,$$

$$(1.7) \quad \mu = \ln \psi + U,$$

$$(1.8) \quad U(\mathbf{x}, \mathbf{m}, t) = \alpha \int_{\Omega} \int_{|\mathbf{m}'|=1} B(\mathbf{x}, \mathbf{x}', \mathbf{m}, \mathbf{m}') \psi(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}'.$$

The velocity field satisfies the Navier–Stokes-like equation

$$(1.9) \quad \rho(\mathbf{u}_t + (\mathbf{u} \cdot \nabla) \mathbf{u}) + \nabla p = \nabla \cdot \tau + F,$$

$$(1.10) \quad \nabla \cdot \mathbf{u} = 0.$$

In the LCP system, the extra stress τ is given by two parts, the viscous stress τ_s and the elastic stress τ_e , namely

$$(1.11) \quad \tau = \tau_s + \tau_e.$$

The viscous stress comes from two sources, one from the solvent and the other from the constrain of rods derived in [6],

$$\tau_s = 2\eta_s \mathbf{D} + \frac{1}{2} \xi_r \mathbf{D} : \langle \mathbf{m}\mathbf{m}\mathbf{m}\mathbf{m} \rangle,$$

where $\mathbf{D} = \frac{1}{2}(\kappa + \kappa^T) = \frac{1}{2}(\nabla \mathbf{u} + (\nabla \mathbf{u})^T)$ is the strain rate tensor, η_s is the solvent viscosity, and $\langle (\cdot) \rangle$ denotes averaging with respect to the distribution ψ ; i.e., $\langle (g) \rangle = \int_{|\mathbf{m}|=1} g \psi d\mathbf{m}$. The elastic stress is derived through a generalized virtual work principle [6]. The detail can be found in [8]. Now we cite the result from [8],

$$(1.12) \quad \tau_e = -\langle (\mathbf{m} \times \mathcal{R} \mu) \mathbf{m} \rangle.$$

Meanwhile the body force can also be identified as

$$(1.13) \quad \mathbf{F} = -\langle \nabla \mu \rangle.$$

Now let $\eta_p = \xi_r, \eta = \eta_s + \eta_p, \gamma = \frac{\eta_s}{\eta}$, and Re denotes the Reynolds number. Then the nondimensional Navier–Stokes-like equation

$$(1.14) \quad \mathbf{u}_t + (\mathbf{u} \cdot \nabla) \mathbf{u} + \nabla p = \frac{\gamma}{Re} \Delta \mathbf{u} + \frac{1-\gamma}{2Re} \nabla \cdot (\mathbf{D} : \langle \mathbf{m}\mathbf{m}\mathbf{m}\mathbf{m} \rangle) + \frac{1-\gamma}{Re De} (\nabla \cdot \tau_e + \mathbf{F}) \text{ for } \mathbf{x} \in \Omega$$

$$(1.15) \quad \nabla \cdot \mathbf{u} = 0, \quad \text{for } \mathbf{x} \in \Omega.$$

In this work we mainly investigate the well-posedness of this new multiscale rodlike polymeric model. Moreover, in most cases [6] $D_{\parallel}^*/D_{\perp}^* \approx 2$, so we can set $D_{\perp}^* = 1$ and $D_{\parallel}^* = 2$ for simple and without the lost of generality. Then (1.6) can be written as

$$(1.16) \quad \frac{\partial \psi}{\partial t} + \nabla \cdot (\mathbf{u}\psi) = \frac{\varepsilon^2}{De} \nabla \cdot [(\mathbf{I} + \mathbf{m}\mathbf{m})(\psi \nabla \mu)] + \frac{1}{De} \mathcal{R} \cdot (\psi \mathcal{R} \mu) - \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \psi), \quad \mathbf{m} \in \mathbb{S}^2.$$

We can verify this system (1.14)–(1.16) satisfies the energy identity in the following way.

Multiplying \mathbf{u} to (1.14) and integrating it over Ω yields

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} \int_{\Omega} |\mathbf{u}|^2 dx + \int_{\Omega} \left[\frac{\gamma}{De} |\nabla \mathbf{u}|^2 + \frac{1-\gamma}{2Re} \langle (\mathbf{m}\mathbf{m} : \mathbf{D})^2 \rangle \right] dx \\
 &= \frac{1-\gamma}{ReDe} \int_{\Omega} [-\tau_e : \nabla \mathbf{u} + \mathbf{F} \cdot \mathbf{u}] dx \\
 &= \frac{1-\gamma}{ReDe} \int_{\Omega} [\langle (\mathbf{m} \times \mathcal{R}\mu)\mathbf{m} \rangle : \nabla \mathbf{u} - \langle \nabla \mu \rangle \cdot \mathbf{u}] dx \\
 (1.17) \quad &= \frac{1-\gamma}{ReDe} \int_{\Omega} \int_{|\mathbf{m}|=1} [(\mathbf{m} \times \mathcal{R}\mu)\mathbf{m}\psi : \nabla \mathbf{u} - \psi \nabla \mu \cdot \mathbf{u}] d\mathbf{m} dx.
 \end{aligned}$$

Multiplying μ to (1.16) as well as integrating over in Ω and the unit sphere yields

$$\begin{aligned}
 & \int_{\Omega} \int_{|\mathbf{m}|=1} \frac{\partial \psi}{\partial t} \mu d\mathbf{m} dx + \int_{\Omega} \left[\frac{\varepsilon^2}{De} \langle \nabla \mu \cdot (\mathbf{I} + \mathbf{m}\mathbf{m}) \nabla \mu \rangle + \frac{1}{De} \langle \mathcal{R}\mu \cdot \mathcal{R}\mu \rangle \right] dx \\
 (1.18) \quad &= \int_{\Omega} \int_{|\mathbf{m}|=1} \mathbf{m} \times \kappa \cdot \mathbf{m} \psi \cdot \mathcal{R}\mu d\mathbf{m} dx + \int_{\Omega} \int_{|\mathbf{m}|=1} \mathbf{u} \psi \cdot \nabla \mu d\mathbf{m} dx.
 \end{aligned}$$

Additionally, we can calculate

$$\begin{aligned}
 (1.19) \quad & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} \left[\psi \ln \psi + \frac{1}{2} \psi U \right] d\mathbf{m} dx \\
 &= \int_{\Omega} \int_{|\mathbf{m}|=1} \left[\frac{\partial \psi}{\partial t} \ln \psi + \frac{\partial \psi}{\partial t} + \frac{1}{2} \frac{\partial \psi}{\partial t} U + \frac{1}{2} \psi \frac{\partial U}{\partial t} \right] d\mathbf{m} dx \\
 &= \int_{\Omega} \int_{|\mathbf{m}|=1} \left[\frac{\partial \psi}{\partial t} (\ln \psi + U) + \frac{1}{2} \left(\psi \frac{\partial U}{\partial t} - \frac{\partial \psi}{\partial t} U \right) \right] d\mathbf{m} dx \\
 &= \int_{\Omega} \int_{|\mathbf{m}|=1} \frac{\partial \psi}{\partial t} \mu d\mathbf{m} dx.
 \end{aligned}$$

Combining (1.17)–(1.19), we obtain that the system (1.14)–(1.18) satisfies the energy law:

$$\begin{aligned}
 (1.20) \quad & \frac{d}{dt} \left[\frac{1}{2} \int_{\Omega} |\mathbf{u}|^2 dx + \frac{1-\gamma}{ReDe} E(\psi) \right] = - \int_{\Omega} \left[\frac{\gamma}{De} |\nabla \mathbf{u}|^2 + \frac{1-\gamma}{2Re} \langle (\mathbf{m}\mathbf{m} : \mathbf{D})^2 \rangle \right] dx \\
 & - \frac{1-\gamma}{ReDe} \int_{\Omega} \left[\frac{\varepsilon^2}{De} \langle \nabla \mu \cdot (\mathbf{I} + \mathbf{m}\mathbf{m}) \nabla \mu \rangle + \frac{1}{De} \langle \mathcal{R}\mu \cdot \mathcal{R}\mu \rangle \right] dx,
 \end{aligned}$$

where $E(\psi)$ is a nonlocal intermolecular potential given by

$$(1.21) \quad E(\psi) = \int_{\Omega} \int_{|\mathbf{m}|=1} \psi(\mathbf{x}, \mathbf{m}, t) \ln \psi(\mathbf{x}, \mathbf{m}, t) + \frac{1}{2} U(\mathbf{x}, \mathbf{m}, t) \psi(\mathbf{x}, \mathbf{m}, t) d\mathbf{m} dx.$$

By the same way one can see that the original system (1.6)–(1.15) also satisfies the energy law like the form (1.20). Comparing to the Doi model of rodlike polymeric fluid [6], this new model is based on the more rational assumption that the particle distribution function (pdf) ψ is possibly different in the macro variable \mathbf{x} . When the pdf is the same at every point \mathbf{x} in the domain Ω , this model is similar to the Doi model. Here the other important difference is the excluded-volume potential (1.3) if we choose $B = |\mathbf{m} \times \mathbf{m}'|^2$ or $B = |\mathbf{m} \times \mathbf{m}'|$ in (1.3), which is the well-known Maier–Saupe or Onsager potential [6]. Now in [8] E and Zhang have proved that the inhomogeneous

property reduces to the Ericksen–Leslie theory in the limit of small Debroah number. Recently numerical simulation results [26, 27] have shown that this model can really describe the anisotropic long-range elasticity of polymeric molecules, and the microstructure and defect dynamics of LCP solution. Moreover they have reported in [26] that there are seven in-plane flow modes in plane Couette flow described by this new model. Four of them have also been reported by Rey and Tsuji [22], and the other three modes are new complicated in-plane modes with inner defects. Furthermore, some significant scaling properties were verified in [26], such as the tumbling period is proportional to the inverse of the shear rate. In plane Poiseuille flow, different local states, such as flow-aligning, tumbling, or wagging, arise in different flow region. There are also some related numerical analysis results for special cases of this model of (1.14)–(1.16), e.g., [3] and references therein. These numerical results require a detailed well-posedness analysis for the system (1.14)–(1.16). This is the main objective of the present work. These related problems for the macroscopic nonlinear elasticity and viscoelasticity cases were recently studied by Sideris and Thomases [23] and Lin, Liu, and Zhang [16]. For the micro-macro model with dumbbell type of potential there are lots of works, e.g., [7, 15, 17] and references therein. In [28], we gave the globally classical existence theory and a numerical analysis for the Dirichlet initial boundary problem of the system (1.14)–(1.16) in a simple case, the 1+1-dimensional case, and the pressure-driven channel flow. More precisely, it is assumed that the rodlike particles rotate in shear plane. That work is a first step towards the better understanding for currently more sophisticated models (1.14)–(1.16).

In this paper we consider the system (1.14)–(1.16) with the initial data

$$(1.22) \quad \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0, \quad \psi(\mathbf{x}, \mathbf{m}, 0) = \psi_0(\mathbf{x}, \mathbf{m}).$$

We denote the space of functions by $H^i(\Omega), i \in \mathbb{N}$, which are in $H^i_{loc}(\mathbb{R}^3)$ (i.e., $u|_\Omega$ for every open bounded set Ω) and which are periodic with period $\Omega: \mathbf{u}(x_j + Le_j) = \mathbf{u}(x_j), \psi(x_j + Le_j) = \psi(x_j), j = 1, 2, 3$. It is easy to see that ψ is also periodic with respect to the variable \mathbf{m} since $\mathbf{m} \in \mathbb{S}^2$. Denote

$$H^i_d(\Omega) = \text{closure of } \mathcal{V} \text{ in } H^i(\Omega), \text{ where } \mathcal{V} = C^\infty_0(\Omega) \cap \{\mathbf{u} : \text{div } \mathbf{u} = 0\}.$$

We can see that $H^i_d(\Omega)$ is a Sobolev space. Moreover, we define the space

$$(1.23) \quad H^l(\Omega, \mathcal{X}_k) = \left\{ \psi : \left| \sum_{j=0}^l \sum_{i=0}^k \int_\Omega \int_{|\mathbf{m}|=1} |\nabla^j \mathcal{R}^i \psi(\mathbf{x}, \mathbf{m}, t)|^2 d\mathbf{m} d\mathbf{x} < \infty \right. \right\},$$

with the natural topology of a Banach space. Now we state our main results:

THEOREM 1.1 (local existence).

- (A1) $\mathbf{u}_0 \in H^3_d(\Omega), \mathbf{u}_0 = 0, \mathbf{u}_0 \cdot \mathbf{e}_j = 0, \mathbf{e}_j \cdot \mathbf{e}_j = 1, \mathbf{e}_j \cdot \mathbf{e}_k = 0, j, k = 1, 2, 3$.
- (A2) $\psi_0 \in \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j), \psi_0 \geq 0, \int_\Omega \int_{|\mathbf{m}|=1} \psi_0(\mathbf{x}, \mathbf{m}) d\mathbf{m} d\mathbf{x} = 1.$

$$\int_\Omega \int_{|\mathbf{m}|=1} \psi_0(\mathbf{x}, \mathbf{m}) d\mathbf{m} d\mathbf{x} = 1.$$

$$\|\mathbf{u}_0\|_{H^3_d(\Omega)}^2 + \|\psi_0\|_{\cap_{i+j=3} H^i(\Omega, \mathcal{X}_j)}^2 < \infty, T' > 0, \quad (1.14), (1.16), (1.22), (\mathbf{u}, \psi).$$

$$(1.24) \quad \mathbf{u} \in L^\infty([0, T']; H^3_d(\Omega)) \cap L^2([0, T]; H^4_d(\Omega));$$

$$(1.25) \quad \psi \in L^\infty([0, T']; \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j)) \cap L^2([0, T']; \cap_{i+j=4} H^i(\Omega, \mathcal{X}_j)),$$

$$\dots \dots f_i \dots \dots$$

THEOREM 1.2 (global existence). (A1) (A2)

$$\|\mathbf{u}_0\|_{H^3_d(\Omega)}^2 + \|\psi_0\|_{H^3(\Omega, \mathcal{X}_0) \cap H^2(\Omega, \mathcal{X}_1)}^2 \leq B,$$

$$B = C_2 \left(\|\mathbf{u}_0\|_{L^\infty(\Omega)} + \|\psi_0\|_{L^\infty(\Omega, \mathcal{X}_0) \cap L^2(\Omega, \mathcal{X}_1)} \right) \quad (1.14) \quad (1.16) \quad (1.22)$$

$$(1.26) \quad \mathbf{u} \in L^\infty([0, \infty); H^3_d(\Omega)),$$

$$(1.27) \quad \psi \in L^\infty([0, \infty); H^3(\Omega, \mathcal{X}_0) \cap H^2(\Omega, \mathcal{X}_1)),$$

$$\|\mathbf{u}\|_{L^\infty([0, \infty); H^3_d(\Omega))}^2 + \|\psi\|_{L^\infty([0, \infty); H^3(\Omega, \mathcal{X}_0) \cap H^2(\Omega, \mathcal{X}_1))}^2 \leq B$$

$$Re < \gamma/C_2, \quad De < \varepsilon^2/C_2;$$

$$C_2 = C_2(n, \Omega)$$

1.1. From the proof of Theorem 1.2 in section 4 we see that C_2 is a large positive constant. Thus Theorem 1.2 requires the Deborah and Reynolds numbers to be small enough.

1.2. From the proof of Theorem 1.1 we can also obtain the “global” result for the two-dimensional system. That is, for given $T > 0$, there exists a solution under the conditions of Theorem 1.1,

$$(1.28) \quad \mathbf{u} \in L^\infty([0, T]; H^3_d(\Omega)) \cap L^2([0, T]; H^4_d(\Omega));$$

$$(1.29) \quad \psi \in L^\infty([0, T]; \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j)) \cap L^2([0, T]; \cap_{i+j=4} H^i(\Omega, \mathcal{X}_j)).$$

But the bound of this solution depends on T .

These results are similar to that of the Navier–Stokes of traditional models of complex fluids in the case of the spatially periodic solutions [5, 24] for $n = 3$. However, in contrast to traditional models of complex fluids [5, 24] which express polymer stress τ using empirical constitutive relations; τ in (1.11) expresses the polymer stress in terms of the microscopic conformations of the polymers. So the model considered in this work is closer to the original system for polymeric fluids in kinetic theory of polymers. But, on the other hand, it causes the difficulty of well-posedness analysis and numerical simulation since we have to study the configuration equation (1.16) which involves five spatial freedom variables, two of them are in the configuration domain and the others are in the macro flow domain. We utilized the properties of the Laplace–Bertrami operator on compact Riemannian manifold to obtain the existence and the preservation of the positivity of the solution to the linearized equation of (1.16). Then the regularity of the solution was strengthened by the energy estimates method. Thus, in virtue of the properties of the distribution function ψ , we can obtain the regularity of the stress τ . Then it is finished by the local well-posedness analysis of the rigid rodlike model by utilizing the Galerkin approximation and energy methods. However, we specially point out that the nonlinear stress (1.11) and the nonlinear body force (1.13) concerned with the solution of the micro-scale model (1.16) let us only obtain the global solution for small Deborah and Reynolds numbers in virtue of the method which we choose in this paper.

The paper is organized as follows. In section 2, we give the iterative scheme of the system to obtain the existence of the local solution and the scheme alternates between

solving an equation of the same type as encountered in incompressible elasticity and solving a linear diffusion equation. Section 3 is devoted to giving the detailed proof of the main lemmas. We will investigate the global existence of the solution in section 4. In this paper C denotes different constant depending only on γ, De, Re, Ω , and ε if there are no special notations. Some times we denote $H_d^p(\Omega), L^p(\Omega)$ by H^p, L^p for brevity.

2. Local solution. In this section we will construct an iterative scheme of the system (1.14)–(1.16) with (1.22) and with which we can obtain the existence of the local solution.

Motivated by the approach [12, 20], we construct an iterative scheme of the system (1.14)–(1.16) with (1.22). Given an iteration ψ^l we determine \mathbf{u}^{l+1} by solving the equations

$$\begin{aligned} \mathbf{u}_t^{l+1} + (\mathbf{u}^{l+1} \cdot \nabla)\mathbf{u}^{l+1} + \nabla p^{l+1} \\ (2.1) \quad &= \frac{\gamma}{Re} \Delta \mathbf{u}^{l+1} + \frac{1-\gamma}{2Re} \nabla \cdot (\mathbf{D} : \langle \mathbf{m m m m} \rangle)^{l+1} + \frac{1-\gamma}{ReDe} (\nabla \cdot \tau_e^l + F^l), \\ (2.2) \quad &\nabla \cdot \mathbf{u}^{l+1} = 0 \end{aligned}$$

with the initial condition $\mathbf{u}^{l+1}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x})$, where

$$\begin{aligned} (2.3) \quad &(\mathbf{D} : \langle \mathbf{m m m m} \rangle)^{l+1} = \mathbf{D}_{ks}^{l+1} \langle m_i m_j m_k m_s \rangle^l, \\ (2.4) \quad &\langle m_i m_j m_k m_s \rangle^l = \int_{|\mathbf{m}|=1} m_i m_j m_k m_s \psi^l(\mathbf{x}, \mathbf{m}, t) d\mathbf{m}, \\ (2.5) \quad &(\tau_e^l) = -\langle (\mathbf{m} \times \mathcal{R}\mu^l) \mathbf{m} \rangle^l, \\ (2.6) \quad &\mu^l = \ln \psi^l + U^l, \\ (2.7) \quad &U^l = \alpha \int_{\Omega} \int_{|\mathbf{m}'|=1} B(\mathbf{x}, \mathbf{x}'; \mathbf{m}, \mathbf{m}') \psi^l(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}', \\ (2.8) \quad &\kappa^{l+1} = (\nabla \mathbf{u}^{l+1})^T, \quad F^l = -\langle \nabla \mu^l \rangle. \end{aligned}$$

Meanwhile, for given \mathbf{u}^l , we determine ψ^l from the following initial value problem:

$$\begin{aligned} (2.9) \quad &\frac{\partial \psi^l}{\partial t} + \nabla \cdot (\mathbf{u}^l \psi^l) = \frac{\varepsilon^2}{De} [\nabla \cdot (\mathbf{I} + \mathbf{m m}) \nabla \psi^l + \nabla \cdot (\mathbf{I} + \mathbf{m m}) (\psi^l \nabla U^l)] \\ &+ \frac{1}{De} [\mathcal{R} \cdot \mathcal{R} \psi^l + \mathcal{R} \cdot (\psi^l \mathcal{R} U^l)] - \mathcal{R} \cdot (\mathbf{m} \times \kappa^l \cdot \mathbf{m} \psi^l), \\ (2.10) \quad &\psi^l(\mathbf{x}, \mathbf{m}, 0) = \psi_0(\mathbf{x}, \mathbf{m}). \end{aligned}$$

Our eventual task is to show that the mapping $\mathcal{M} : \mathbf{u}^l \mapsto \mathbf{u}^{l+1}$ has a fixed point in an appropriate complete space of functions. The fixed point of the mapping is the solution we seek.

We will consider the mapping \mathcal{M} in the function space $S(M, T)$, which is defined as a set of all functions $\mathbf{u} : \Omega \times [0, T] \rightarrow \mathbb{R}^n (n = 2, 3)$ with the following properties:

$$\begin{aligned} (2.11) \quad &\mathbf{u} \in L^\infty([0, T]; H_d^3(\Omega)) \cap L^2([0, T]; H_d^4(\Omega)), \\ (2.12) \quad &\|\mathbf{u}\|_{L^\infty([0, T]; H_d^3(\Omega))}^2 + \|\mathbf{u}\|_{L^2([0, T]; H_d^4(\Omega))}^2 \leq M. \end{aligned}$$

On $S(M, T)$, we define the metric

$$(2.13) \quad d(\mathbf{u}_1, \mathbf{u}_2) = \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^\infty([0, T]; H_d^3(\Omega))} + \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2([0, T]; H_d^4(\Omega))}.$$

It is easy to verify that $S(M, T)$ is complete with the associated metric, and also it is nonempty for large M , as evidenced in [19]. The properties of the mapping \mathcal{M} is established by proving the next three lemmas.

LEMMA 2.1. *Let $(M, T) \in S$. Then, for any $l \in \mathbb{N}$, there exists $T' > 0$ such that*

$$(2.14) \quad \|\psi^l\|_{L^\infty([0, T]; H^3(\Omega, \mathcal{X}_0)) \cap L^2([0, T]; H^4(\Omega, \mathcal{X}_0))} \leq K.$$

Moreover, for any $T' > 0$, there exists $M > 0$ such that (2.1)–(2.8) hold for all $(M, T) \in S$ with $T \leq T'$.

$$(2.15) \quad \mathbf{u}^{l+1} \in L^\infty([0, T']; H_d^3(\Omega)) \cap L^2([0, T']; H_d^4(\Omega))$$

$$(2.16) \quad \|\mathbf{u}^{l+1}\|_{L^\infty([0, T']; H_d^3(\Omega))}^2 + \|\mathbf{u}^{l+1}\|_{L^2([0, T']; H_d^4(\Omega))}^2 \leq \phi_1(T', K),$$

$$f = \frac{1-\gamma}{ReDe} (\nabla \cdot \tau_e^l + F^l).$$

$$(2.17) \quad \phi_1(T, K) = [\|\mathbf{u}_0\|_{H_d^3(\Omega)}^2 + C\|f\|_{L^2([0, T], H^2(\Omega))}^2 + C(K)]e^{CK^2T + CK^4T} + \|\mathbf{u}_0\|_{H_d^3(\Omega)}^2 + C\|f\|_{L^2([0, T], H^2(\Omega))}^2 + C(K)T.$$

LEMMA 2.2.

$$(2.18) \quad \tau_e^l \in L^2([0, T]; H^3(\Omega)), \quad F^l \in L^2([0, T]; H^2(\Omega))$$

$$(2.19) \quad \|\tau_e^l\|_{L^2([0, T]; H^3(\Omega))} \leq C\|\psi^l\|_{L^2([0, T]; H^3(\Omega, \mathcal{X}_0))},$$

$$(2.20) \quad \|F^l\|_{L^2([0, T]; H^2(\Omega))} \leq C\|\psi^l\|_{L^2([0, T]; H^3(\Omega, \mathcal{X}_0))}$$

$$\psi^l \in L^2([0, T]; H^3(\Omega, \mathcal{X}_0))$$

LEMMA 2.3. *Let $(M, T) \in S$. Then, for any $l \in \mathbb{N}$, there exists $M > 0$ such that (2.9)–(2.10) hold for all $(M, T) \in S$ with $T \leq T'$.*

$$(2.21) \quad \psi^l \in L^\infty([0, T]; \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j)) \cap L^2([0, T]; \cap_{i+j=4} H^i(\Omega, \mathcal{X}_j)),$$

$$(2.22) \quad \|\psi^l\|_{L^2([0, T], \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j))}^2 \leq \|\psi_0\|_{\cap_{i+j=3} H^i(\Omega, \mathcal{X}_j)}^2 (CT + CMT) \cdot e^{CT + CMT}.$$

By combining Lemmas 2.1–2.3, it follows easily that the map $\mathcal{M} : S(M, T') \rightarrow S(M, T')$ is a compact operator if M is chosen sufficiently large and T' is chosen sufficiently small. In fact, by using (2.22), we know

$$\|\psi\|_{L^2([0, T], H^3(\Omega, \mathcal{X}_0))}^2 \leq \|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2 (CT + CMT)e^{CT + CMT}.$$

Thus

$$\|f\|_{L^2([0, T], H^2(\Omega))} \leq C\|\nabla \cdot \tau_e^l + F^l\|_{L^2([0, T], H^2(\Omega))} \leq C\|\psi\|_{L^2([0, T], H^3(\Omega, \mathcal{X}_0))}.$$

Then, from (2.17), we have

$$\begin{aligned} & \|\mathbf{u}^{l+1}\|_{L^\infty([0, T']; H_d^3(\Omega))}^2 + \|\mathbf{u}^{l+1}\|_{L^2([0, T']; H_d^4(\Omega))}^2 \\ & \leq [\|\mathbf{u}_0\|_{H_d^3(\Omega)}^2 + \|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2 (CT + CMT)e^{CT + CMT} + C(K)]e^{CK^2T + CK^4T} \\ & + \|\mathbf{u}_0\|_{H_d^3(\Omega)}^2 + \|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2 (CT + CMT)e^{CT + CMT} + C(K)T. \end{aligned}$$

Now we choose

$$(2.23) \quad M \geq 6\|\mathbf{u}_0\|_{H^3(\Omega)}^2 + 2C(K)(2 + T_0) + 12CT_0\|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2,$$

$$(2.24) \quad T' \leq T_0 \triangleq \min \left\{ \frac{\ln 2}{C(1 + M)}, \frac{\ln 2}{C(K^2 + K^4)}, \frac{1}{12C\|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2} \right\}.$$

Then $e^{CT'+CMT'} \leq 2, e^{CK^2T'+CK^4T'} \leq 2$, and

$$\begin{aligned} & C\|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2 e^{CT'+CMT'} MT' [1 + e^{CK^2T'+CK^4T'}] \leq \frac{M}{2}, \\ & \|\mathbf{u}_0\|_{H^3(\Omega)}^2 (1 + e^{CK^2T'+CK^4T'}) + C(K)[T' + e^{CK^2T'+CK^4T'}] \\ & + C\|\psi_0\|_{H^3(\Omega, \mathcal{X}_0)}^2 e^{CT'+CMT'} MT' [1 + e^{CK^2T'+CK^4T'}] T' \leq \frac{M}{2}. \end{aligned}$$

Thus we have (2.12).

Since $S(M, T')$ is clearly a closed, convex subset of $L^2([0, T], H_d^4(\Omega))$ and is also compact, by the fixed point theorem of Leray and Schauder [14], the conclusion of Theorem 1.1 can be obtained.

3. Proof of lemmas.

3.1. Estimates of \mathbf{u} . In this section we will give the proof of Lemma 2.1. Let $\mathbf{u}^{l+1} = w, q = p^{l+1}$, and the fourth order tensor

$$(3.1) \quad A(\mathbf{x}, t) = (a_{ijkl}), \quad \text{where } a_{ijkl}(\mathbf{x}, t) = \int_{|\mathbf{m}|=1} m_i m_j m_k m_l \psi^l(\mathbf{x}, \mathbf{m}, t) d\mathbf{m}.$$

(2.1) can be rewritten as

$$(3.2) \quad \begin{aligned} w_t + (w \cdot \nabla)w + \nabla q &= \frac{\gamma}{Re} \Delta w + f \\ &+ \frac{1-\gamma}{2Re} \nabla \cdot [\mathbf{D} : A(\mathbf{x}, t)], \end{aligned}$$

$$(3.3) \quad \nabla \cdot w = 0,$$

where $f = \frac{1-\gamma}{ReDe} (\nabla \cdot \tau_e^l + F^l)$. In the following we solve this problem to obtain the existence and uniqueness of the solution by using the Galerkin approximation similar to solving the standard Navier–Stokes equation [5, 24]. The difference here is the appearance of the term $\frac{1-\gamma}{2Re} \nabla \cdot [\mathbf{D} : A(\mathbf{x}, t)]$. We can see that it is a good term when we give a priori estimates because we will obtain $-\frac{1-\gamma}{2Re} \int_{\Omega} \langle (\mathbf{D} : \mathbf{m}\mathbf{m})^2 \rangle d\mathbf{x}$ while multiplying w to (3.2) and integrating it in Ω . The high derivatives estimates are obtained similarly. Thus we will obtain the result in Lemma 2.1 provided that

$$(3.4) \quad f \in L^2([0, T]; H^2(\Omega)),$$

$$(3.5) \quad A \in L^\infty([0, T], H^3(\Omega)) \cap L^2([0, T], H^4(\Omega)).$$

The regularity of f and A can be easily obtained from the estimates of τ_e and F in section 3.2 and ones of ψ in section 3.3, respectively. Why we need conditions (3.4) and (3.5), for the aim of self-contained, will be answered in the outline of the proof to Lemma 2.1 when $n = 2$ and $n = 3$ in the appendix. We also refer the readers to [5, 24] for further details.

3.2. Estimates of τ_e and F .

Proof of Lemma 2.2. From the definition of τ_e , it is straightforward to obtain its estimates from the assumption of ψ . In fact,

$$\begin{aligned} \tau_e &= - \int_{|\mathbf{m}|=1} (\mathbf{m} \times \mathcal{R}\mu) \mathbf{m} \psi d\mathbf{m} \\ &= - \int_{|\mathbf{m}|=1} \left[\mathbf{m} \times \left(\frac{1}{\psi} \mathcal{R}\psi + \mathcal{R}U \right) \right] \mathbf{m} \psi d\mathbf{m} \\ &= - \int_{|\mathbf{m}|=1} (\mathbf{m} \times \mathcal{R}\psi) \mathbf{m} d\mathbf{m} - \int_{|\mathbf{m}|=1} (\mathbf{m} \times \mathcal{R}U) \mathbf{m} \psi d\mathbf{m} \\ &= -\mathbf{I} + 3 \int_{|\mathbf{m}|=1} \mathbf{m} \mathbf{m} \psi d\mathbf{m} - \int_{|\mathbf{m}|=1} (\mathbf{m} \times \mathcal{R}U) \mathbf{m} \psi d\mathbf{m}, \end{aligned}$$

where we used the property of operators \mathcal{R} and $\int_{|\mathbf{m}|=1} \cdot d\mathbf{m}$ (p. 293 in [6]),

$$(3.6) \quad \int_{|\mathbf{m}|=1} G(\mathbf{m}) \mathcal{R}[F(\mathbf{m})] d\mathbf{m} = - \int_{|\mathbf{m}|=1} F(\mathbf{m}) \mathcal{R}[G(\mathbf{m})] d\mathbf{m}.$$

From the definition (2.7) of U ,

$$\begin{aligned} U &= \alpha \int_{\Omega} \int_{|\mathbf{m}'|=1} \frac{1}{\varepsilon^3} \chi \left(\frac{\mathbf{x} - \mathbf{x}'}{\varepsilon} \right) |\mathbf{m} \times \mathbf{m}'|^2 \psi^l(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}' \\ &= \alpha \int_{\Omega} \int_{|\mathbf{m}'|=1} \frac{1}{\varepsilon^3} \chi \left(\frac{\mathbf{x} - \mathbf{x}'}{\varepsilon} \right) \psi^l(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}' \\ (3.7) \quad &- \alpha \mathbf{m} \mathbf{m} : \int_{\Omega} \int_{|\mathbf{m}'|=1} \frac{1}{\varepsilon^3} \chi \left(\frac{\mathbf{x} - \mathbf{x}'}{\varepsilon} \right) \mathbf{m}' \mathbf{m}' \psi^l(\mathbf{x}', \mathbf{m}', t) d\mathbf{m}' d\mathbf{x}'; \end{aligned}$$

here $\varepsilon > 0$ is in (1.5) and $\chi(\mathbf{x})$ is a smooth kernel, and we can see that $U \in C^\infty(\Omega \times \mathbb{S}^2)$. But the bounds of the derivatives of U with respect to \mathbf{x} and \mathbf{m} depend on ε , denoted by $C(\varepsilon)$.

Therefore

$$(3.8) \quad \|\tau_e\|_{L^2(\Omega)}^2 \leq C(1 + \|\psi\|_{L^2(\Omega, \mathcal{X}_0)}^2)$$

since $\psi \in L^2([0, T], H^3(\Omega, \mathcal{X}_0))$. Higher derivatives can be obtained similarly. Thus (2.19) is obtained. By the same way we can obtain

$$(3.9) \quad F = \langle \nabla \mu \rangle = \int_{|\mathbf{m}|=1} (\nabla \psi + \psi \nabla U) d\mathbf{m} \in L^2([0, T]; H^2(\Omega))$$

since $\psi \in L^2([0, T], H^3(\Omega, \mathcal{X}_0))$. Hence (2.20) is obtained by using (3.9).

3.3. Estimates of ψ . In this section we first review some results about the Laplace–Beltrami operator on a compact Riemannian manifold which will be utilized to show the existence of the solution of (2.9)–(2.10). Then we will give the regularity estimates of the solution to (2.9)–(2.10).

3.3.1. Review about the Laplace–Beltrami operator. We recall some known results about the Laplace–Beltrami operator on a compact Riemannian manifold (M_n, g) ; see (section 4 of Chap. 4) of [1] and [14].

LEMMA 3.1. Let $v(Q, t) \in C^2(M_n \times [0, t_0])$ and $v \geq 0$ on $M_n \times \{0\} \cup \partial M_n \times [0, t_0]$. Then

$$(3.10) \quad \partial v / \partial t \geq \Delta_{M_n} v + b^i(Q, t) \partial_i v + c(Q, t)v$$

if $b^i, c \in C^0(M_n \times [0, t_0])$ and $v \geq 0$. This lemma can be proved similar to the proof of the maximum principle in p. 130 of [1]. Let $C_m = \max_{M_n \times \mathbb{R}^+} |c(Q, t)|$ and $w = e^{-(C_m+1)t}v$. Then w and v have the same sign. Since

$$\partial w / \partial t = e^{-(C_m+1)t}[\partial v / \partial t - (C_m + 1)v],$$

we have

$$(3.11) \quad \partial w / \partial t \geq \Delta_{M_n} w + b^i(Q, t) \partial_i w + [c(Q, t) - C_m - 1]w.$$

Assume w is negative somewhere and let (Q, t) be a point where w achieves its minimum. Then $\Delta_{M_n} w \geq 0$, $\partial_i w = 0$, and $\partial w / \partial t \leq 0$. Thus (3.11) implies $w(Q, t) \geq 0$, which yields a contradiction. \square

LEMMA 3.2. Let $g \in L^\infty([0, t_0], L^p(M_n))$ and

$$v \in W^{1,\infty}([0, t_0], W^{2,p}(M_n))$$

such that

$$(3.12) \quad \partial v^i / \partial t = \Delta_{M_n} v^i + a_j^i \partial_i v^j + b_j^i v^j + g^i,$$

for $1 \leq \alpha \leq k$, $v(P, 0) \equiv 0$, $P \in M_n$, $v^i(i = 1, 2, \dots, k) \in C^1(M_n \times [0, \infty))$, $g^i(i = 1, 2, \dots, k) \in C^0(M_n \times [0, \infty))$, and $a_j^i, b_j^i \in C^0(M_n \times [0, \infty))$.

3.1. In fact, the condition on the coefficients a_j^i and b_j^i in the above lemma can be weakened to be bounded in $L^\infty([0, t_0] \times M_n)$. The proof is analogous.

3.3.2. Estimate of ψ .

Proof of Lemma 2.3. In this section we simply denote

$$(3.13) \quad \phi(\mathbf{x}, \mathbf{m}, t) = \psi^l(\mathbf{x}, \mathbf{m}, t) \quad W = U^l.$$

Then (2.9)–(2.10) can be rewritten as

$$(3.14) \quad \begin{aligned} \frac{\partial \phi}{\partial t} + \mathbf{u}^l \cdot \nabla \phi &= \frac{\varepsilon^2}{De} [\nabla \cdot (\mathbf{I} + \mathbf{m}\mathbf{m}) \nabla \phi + \nabla \cdot (\mathbf{I} + \mathbf{m}\mathbf{m})(\phi \nabla W)] \\ &+ \frac{1}{De} [\mathcal{R} \cdot \mathcal{R} \phi + \mathcal{R} \cdot (\phi \mathcal{R} W)] - \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \phi), \end{aligned}$$

$$(3.15) \quad \phi(\mathbf{x}, \mathbf{m}, 0) = \psi_0(\mathbf{x}, \mathbf{m}).$$

Here we still denote $(\nabla \mathbf{u}^l(\mathbf{x}, t))^T$ by κ . Equation (3.14) is a nonlinear differential-integral equation. We first linearize (3.14) replacing W by U^{l-1} and obtaining the solution of this linear equation with the initial data (3.15). If now we transform the Cartesian coordinates \mathbf{m} to the local coordinates (θ, φ) of the unit sphere \mathbb{S}^2 in

\mathbb{R}^3 by $m_1(\theta, \varphi) = \sin \theta \cos \varphi$, $m_2(\theta, \varphi) = \sin \theta \sin \varphi$, and $m_3 = \cos \varphi$, the operator $\mathcal{R} \cdot \mathcal{R}$ is the Laplace–Beltrami operator on the unit sphere [4]. Thus the operator $\frac{\varepsilon^2}{De} \Delta + \frac{1}{De} \mathcal{R} \cdot \mathcal{R}$ is also the Laplace–Beltrami operator. Moreover, the coefficients of (3.14) are all bounded and smooth when κ is bounded, which is obtained from $\mathbf{u} \in L^\infty([0, T]; H_d^3(\Omega))$. Therefore we can utilize the above lemmas to show the existence and nonnegativeness of the solution to the linearized equation when the initial data is nonnegative. In fact, using Lemma 3.2 and Remark 3.1 for $v = \phi - \psi_0$, we obtain that (3.14)–(3.15) possesses a unique global solution $v \in W^{1,\infty}([0, T], H^2(\Omega, \mathcal{X}_0) \cap H^1(\Omega, \mathcal{X}_1) \cap L^2(\Omega, \mathcal{X}_2))$. Therefore, $\psi^l \in W^{1,\infty}([0, T]; H^2(\Omega, \mathcal{X}_0) \cap H^1(\Omega, \mathcal{X}_1) \cap L^2(\Omega, \mathcal{X}_2))$ for given $T > 0$. Now the coefficient of ϕ is $\frac{\varepsilon^2}{De} \nabla \cdot (\mathbf{I} + \mathbf{m}\mathbf{m}) \nabla W + \frac{1}{De} \mathcal{R} \cdot \mathcal{R} W - \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m})$ and it is bounded since $\psi^{l-1} \in L^\infty([0, T]; \cap_{i+j=3} H^i(\Omega, \mathcal{X}_j))$ and $\mathbf{u}^l \in L^\infty([0, T]; H_d^3(\Omega))$. Then it follows that the positivity is preserved by Lemma 3.1, and by integrating both sides of (3.14) we find that

$$\int_{\Omega} \int_{|\mathbf{m}|=1} \phi(\mathbf{x}, \mathbf{m}, t) d\mathbf{m} d\mathbf{x} = \int_{\Omega} \int_{|\mathbf{m}|=1} \phi(\mathbf{x}, \mathbf{m}, 0) d\mathbf{m} d\mathbf{x} = 1 \text{ for all } t.$$

We can obtain the solution of (3.14) with (3.15) from the solution of the linearized equation to construct a suitable Sobolev space and use the Schauder fixed point theorem in the standard argument. So here we omit the detail.

Next we will further prove the regularity of the solution $\psi^l(\mathbf{x}, \mathbf{m}, t)$ if the initial data is more regular.

I. The estimate of ϕ . Multiplying ϕ to (3.14) and integrating on S^2 with respect to \mathbf{m} and in Ω with respect to \mathbf{x} yields

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x} + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\mathbf{m} \cdot \nabla \phi|^2) d\mathbf{m} d\mathbf{x} \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R} \phi|^2 d\mathbf{m} d\mathbf{x} \\ = & -\frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (\mathbf{I} + \mathbf{m}\mathbf{m}) \phi \nabla W \cdot \nabla \phi d\mathbf{m} d\mathbf{x} - \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} \phi \mathcal{R} W \cdot \mathcal{R} \phi d\mathbf{m} d\mathbf{x} \\ & + \int_{\Omega} \int_{|\mathbf{m}|=1} \mathbf{m} \times \kappa \cdot \mathbf{m} \phi \cdot \mathcal{R} \phi d\mathbf{m} d\mathbf{x}, \end{aligned}$$

where we used the periodic condition. Thus we can obtain the estimate,

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x} + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\mathbf{m} \cdot \nabla \phi|^2) d\mathbf{m} d\mathbf{x} \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R} \phi|^2 d\mathbf{m} d\mathbf{x} \\ (3.16) \quad & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x} + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x}, \end{aligned}$$

where we used the fact that $|\kappa|$ is also bounded in $\Omega \times [0, T]$ redefined by a set of measure zero from (3.7) since $\mathbf{u}^l \in L^\infty([0, T]; H^3(\Omega))$ and $|\nabla W|, |\mathcal{R} W|$ are bounded by $C(\varepsilon)$. Therefore from (3.16) we obtain

$$(3.17) \quad \sup_{t \in [0, T]} \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x} \leq \|\psi_0\|_{L^2(\Omega, \mathcal{X}_0)}^2 e^{CT+CMT} \triangleq N_1,$$

and

$$(3.18) \quad \frac{\varepsilon^2}{De} \int_0^T \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla\phi|^2 + |\mathbf{m} \cdot \nabla\phi|^2) d\mathbf{m} dx dt + \frac{1}{De} \int_0^T \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}\phi|^2 d\mathbf{m} dx dt \leq CN_1T + CMN_1T \triangleq N_2.$$

This shows that

$$\phi \in L^\infty([0, T]; L^2(\Omega, \mathcal{X}_0)) \cap L^2([0, T]; L^2(\Omega, \mathcal{X}_1)) \cap L^2([0, T]; H^1(\Omega, \mathcal{X}_0))$$

provided that $\psi_0 \in L^2(\Omega, \mathcal{X}_0)$ and $\mathbf{u}^l \in S(M, T)$.

II. The estimate of $\mathcal{R}\phi$, ϕ . Applying the operator \mathcal{R} to (3.14) and multiplying $\mathcal{R}\phi$ to (3.14) and integrating on \mathbb{S}^2 with respect to \mathbf{m} and in Ω with respect to \mathbf{x} , respectively, yields

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}\phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla\mathcal{R}\phi|^2 + |\mathbf{m} \cdot \nabla\mathcal{R}\phi|^2) d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R} \cdot \mathcal{R}\phi|^2 d\mathbf{m} dx \\ = & -\frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} \{ \mathcal{R} \cdot (\mathbf{I} + \mathbf{m}\mathbf{m}) \cdot \nabla\phi \cdot \nabla\mathcal{R}\phi - \nabla \cdot [(\mathbf{I} + \mathbf{m}\mathbf{m})(\phi\nabla W)] \mathcal{R} \cdot \mathcal{R}\phi \} d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} \mathcal{R} \cdot (\phi\mathcal{R}W) \mathcal{R} \cdot \mathcal{R}\phi d\mathbf{m} dx + \int_{\Omega} \int_{|\mathbf{m}|=1} \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m}\phi) \mathcal{R} \cdot \mathcal{R}\phi d\mathbf{m} dx. \end{aligned}$$

Thus, by using $|\nabla W|, |\nabla^2 W|, |\nabla\mathcal{R}W|, |\mathcal{R}W|, |\mathcal{R}^2W|$ are bounded for all $\mathbf{x} \in \Omega$ and $\mathbf{m} \in \mathbb{S}^2$, and $|\kappa|$ is bounded by redefined in a set of measure zero, we have

$$(3.19) \quad \begin{aligned} & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}\phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla\mathcal{R}\phi|^2 + |\mathbf{m} \cdot \nabla\mathcal{R}\phi|^2) d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R} \cdot \mathcal{R}\phi|^2 d\mathbf{m} dx \\ & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}\phi|^2 + |\nabla_{\mathbf{x}}\phi|^2 + |\phi|^2) d\mathbf{m} dx \\ & + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx. \end{aligned}$$

Similarly, differentiating (3.14) with respect to \mathbf{x} , then multiplying $\nabla\phi$ and integrating it, we obtain

$$(3.20) \quad \begin{aligned} & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla\phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\Delta\phi|^2 + |\mathbf{m} \cdot \nabla^2\phi|^2) d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla\mathcal{R}\phi|^2 d\mathbf{m} dx \\ & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla\phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx \\ & + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla\phi|^2 + |\phi|^2) d\mathbf{m} dx. \end{aligned}$$

Here we used that $|\nabla W|, |\nabla^2 W|, |\mathcal{R}W|, |\mathcal{R}^2 W|$ and κ are bounded. Combination of (3.16), (3.19), and (3.20) and application of the Grownwall inequality yields

$$\begin{aligned} & \sup_{t \in [0, T]} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx \\ & \leq \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \psi_0|^2 + |\mathcal{R}\psi_0|^2 + |\psi_0|^2) d\mathbf{m} dx e^{CT+CM} \triangleq N_3, \\ & \frac{\varepsilon^2}{De} \int_0^T \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\Delta \phi|^2 + |\nabla \mathcal{R}\phi|^2) d\mathbf{m} dx dt \\ & + \frac{1}{De} \int_0^T \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \mathcal{R}\phi|^2 + |\mathcal{R} \cdot \mathcal{R}\phi|^2 d\mathbf{m} dx dt \leq CN_3 T + CMN_3 T. \end{aligned}$$

This shows that

$$(3.21) \quad \begin{aligned} & \phi \in L^\infty([0, t]; \cap_{i+j=1, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)) \quad \text{and} \\ & \phi \in L^2([0, t]; \cap_{i+j=2, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)) \end{aligned}$$

provided that $\psi_0 \in \cap_{i+j=1, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)$.

III. The estimate of $\mathcal{R}^2 \phi$, $\mathcal{R}^2 \phi$, and $\mathcal{R}\phi$. Similar to that in section II, we can obtain the following estimates:

$$(3.22) \quad \begin{aligned} & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \mathcal{R}\phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\Delta \mathcal{R}\phi|^2 + |\mathbf{m} \cdot \nabla^2 \mathcal{R}\phi|^2) d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \mathcal{R} \cdot \mathcal{R}\phi|^2 d\mathbf{m} dx \\ & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \mathcal{R}\phi|^2 + |\nabla_{\mathbf{x}}^2 \phi|^2 + |\nabla \phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx \\ & + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^2 \phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx, \end{aligned}$$

$$(3.23) \quad \begin{aligned} & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}^2 \phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \mathcal{R} \cdot \mathcal{R}\phi|^2 + |\mathbf{m} \cdot \nabla \mathcal{R}^2 \phi|^2) d\mathbf{m} dx \\ & + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}(\mathcal{R} \cdot \mathcal{R}\phi)|^2 d\mathbf{m} dx \\ & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^2 \phi|^2 + |\mathcal{R}\phi|^2 + |\nabla \phi|^2 + |\nabla \mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx \\ & + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^2 \phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx, \end{aligned}$$

$$\begin{aligned}
& \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \phi|^2 d\mathbf{m} d\mathbf{x} + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \Delta \phi|^2 + |\mathbf{m} \cdot \nabla^3 \phi|^2) d\mathbf{m} d\mathbf{x} \\
& \quad + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R} \Delta_{\mathbf{x}} \phi|^2 d\mathbf{m} d\mathbf{x} \\
& \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^2 \phi|^2 + |\nabla \phi|^2 + |\mathcal{R} \phi|^2 + |\nabla \mathcal{R} \phi|^2 + |\phi|^2) d\mathbf{m} d\mathbf{x} \\
& \quad + C |\nabla \kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R} \phi|^2 + |\nabla \phi|^2 + |\phi|^2) d\mathbf{m} d\mathbf{x} \\
(3.24) \quad & + C |\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \mathcal{R} \phi|^2 + |\nabla \phi|^2) d\mathbf{m} d\mathbf{x}.
\end{aligned}$$

Combining the above estimates and applying the Gronwall inequality with $|\nabla \kappa| \in L^2([0, T], H^4(\Omega))$, we can obtain

$$\begin{aligned}
\|\phi\|_{L^\infty([0, t]; \cap_{i+j=2, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))}^2 & \leq \|\psi_0\|_{\cap_{i+j=2, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)}^2 e^{CT+CM T} \triangleq N_4, \\
\|\phi\|_{L^2([0, t]; \cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))}^2 & \leq CN_4 T + CMN_4 T.
\end{aligned}$$

This implies that

$$\begin{aligned}
(3.25) \quad & \phi \in L^\infty([0, t]; \cap_{i+j=2, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)) \quad \text{and} \\
& \phi \in L^2([0, t]; \cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))
\end{aligned}$$

provided that $\psi_0 \in \cap_{i+j=2, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)$.

IV. The estimate of $\mathcal{R}^i \phi(i+j=3)$. Analogously to that in section III, we can obtain the following estimates:

$$\begin{aligned}
& \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}^3 \phi|^2 d\mathbf{m} d\mathbf{x} + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \mathcal{R}^3 \phi|^2 + |\mathbf{m} \cdot \nabla \mathcal{R}^3 \phi|^2) d\mathbf{m} d\mathbf{x} \\
& \quad + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}^4 \phi|^2 d\mathbf{m} d\mathbf{x} \\
& \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^3 \phi|^2 + |\mathcal{R}^2 \phi|^2 + |\mathcal{R} \phi|^2 + |\phi|^2 \\
& \quad + |\nabla \mathcal{R}^2 \phi|^2 + |\nabla \mathcal{R} \phi|^2 + |\nabla \phi|^2) d\mathbf{m} d\mathbf{x} \\
(3.26) \quad & + C |\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^3 \phi|^2 + |\mathcal{R}^2 \phi|^2 + |\mathcal{R} \phi|^2 + |\phi|^2) d\mathbf{m} d\mathbf{x}, \\
& \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \mathcal{R}^2 \phi|^2 d\mathbf{m} d\mathbf{x} + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^2 \mathcal{R}^2 \phi|^2 + |\mathbf{m} \cdot \nabla^2 \mathcal{R}^2 \phi|^2) d\mathbf{m} d\mathbf{x} \\
& \quad + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathcal{R}^3 \nabla \phi|^2 d\mathbf{m} d\mathbf{x} \\
& \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^2 \mathcal{R} \phi|^2 + |\nabla^2 \phi|^2 + |\nabla \mathcal{R}^2 \phi|^2 + |\nabla \mathcal{R} \phi|^2 \\
& \quad + |\nabla \phi|^2 + |\mathcal{R}^2 \phi|^2 + |\phi|^2) d\mathbf{m} d\mathbf{x} \\
(3.27) \quad & + C |\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \phi|^2 d\mathbf{m} d\mathbf{x} + C |\nabla \kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} |\phi|^2 d\mathbf{m} d\mathbf{x},
\end{aligned}$$

$$\begin{aligned}
 & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \mathcal{R}\phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^3 \mathcal{R}\phi|^2 + |\mathbf{m} \cdot \nabla^3 \mathcal{R}\phi|^2) d\mathbf{m} dx \\
 & \quad + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \mathcal{R}^2 \phi|^2 d\mathbf{m} dx \\
 & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^3 \phi|^2 + |\nabla^2 \mathcal{R}\phi|^2 + |\nabla \mathcal{R}\phi|^2 + |\mathcal{R}\phi|^2 \\
 & \quad + |\nabla^2 \phi|^2 + |\nabla \phi|^2 + |\phi|^2) d\mathbf{m} dx \\
 & \quad + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\nabla \mathcal{R}\phi|^2 + |\nabla^2 \mathcal{R}\phi|^2 + |\nabla \mathcal{R}^2 \phi|^2) d\mathbf{m} dx \\
 (3.28) \quad & + C|\nabla \kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} (|\mathcal{R}^2 \phi|^2 + |\nabla \mathcal{R}\phi|^2 + |\mathcal{R}\phi|^2 + |\phi|^2) d\mathbf{m} dx, \\
 & \frac{d}{dt} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^3 \phi|^2 d\mathbf{m} dx + \frac{\varepsilon^2}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^4 \phi|^2 + |\mathbf{m} \cdot \nabla^4 \phi|^2) d\mathbf{m} dx \\
 & \quad + \frac{1}{De} \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^3 \mathcal{R}\phi|^2 d\mathbf{m} dx \\
 & \leq C \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^3 \phi|^2 + |\nabla^2 \phi|^2 + |\nabla \phi|^2 + |\phi|^2) d\mathbf{m} dx \\
 & \quad + C|\kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^3 \phi|^2 d\mathbf{m} dx + C|\nabla \kappa| \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \phi|^2 d\mathbf{m} dx \\
 & \quad + C\|\phi\|_{\mathcal{X}_0}^2 \int_{\Omega} |\nabla^3 \kappa|^2 dx + C\|\nabla^3 \kappa\|_{L^2}^2 \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla^2 \phi|^2 + |\nabla \phi|^2) d\mathbf{m} dx \\
 (3.29) \quad & + C\|\nabla^2 \kappa\|_{L^2}^2 \int_{\Omega} \int_{|\mathbf{m}|=1} (|\nabla \phi|^2 + |\phi|^2) d\mathbf{m} dx.
 \end{aligned}$$

In the last estimate we used the following estimates:

$$\begin{aligned}
 & \int_{\Omega} \left[|\nabla^2 \kappa| \int_{|\mathbf{m}|=1} |\nabla \phi| |\nabla^3 \mathcal{R}\phi| d\mathbf{m} \right] dx \\
 & \leq \int_{\Omega} |\nabla^2 \kappa| \|\nabla \phi\|_{L^2(\mathbb{S}^2)} \|\nabla^3 \mathcal{R}\phi\|_{L^2(\mathbb{S}^2)} dx \\
 & \leq \|\nabla^2 \kappa\|_{L^3(\Omega)} \|\|\nabla \phi\|_{L^2(\mathbb{S}^2)}\|_{L^6(\Omega)} \|\|\nabla^3 \mathcal{R}\phi\|_{L^2(\mathbb{S}^2)}\|_{L^2(\Omega)} \\
 & \leq \|\nabla^2 \kappa\|_{H^{\frac{1}{2}}(\Omega)} \|\|\nabla \phi\|_{L^2(\mathbb{S}^2)}\|_{H^1(\Omega)} \|\|\nabla^3 \mathcal{R}\phi\|_{L^2(\mathbb{S}^2)}\|_{L^2(\Omega)} \\
 & \leq \|\nabla^2 \kappa\|_{H^4(\Omega)} \|\phi\|_{H^2(\Omega, \mathcal{X}_0)} \|\|\nabla^3 \mathcal{R}\phi\|_{L^2(\mathbb{S}^2)}\|_{L^2(\Omega)} \\
 (3.30) \quad & \leq \frac{1}{De} \|\nabla^3 \mathcal{R}\phi\|_{L^2(\Omega, \mathcal{X}_0)}^2 + C(De) \|\mathbf{u}\|_{H^4(\Omega)}^2 \|\phi\|_{H^2(\Omega, \mathcal{X}_0)}^2.
 \end{aligned}$$

Combining the above estimates and applying the Gronwall inequality by using $\mathbf{u} \in L^2([0, T], H^4(\Omega))$, we can obtain

$$\begin{aligned}
 \|\phi\|_{L^\infty([0, t]; \cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))} & \leq \|\psi_0\|_{\cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)} e^{CT+CM T} \triangleq N_6, \\
 \|\phi\|_{L^2([0, t]; \cap_{i+j=4, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))} & \leq CN_6 T + CMN_6 T.
 \end{aligned}$$

This implies that

$$\begin{aligned}
 (3.31) \quad & \phi \in L^\infty([0, t]; \cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)) \quad \text{and} \\
 & \phi \in L^2([0, t]; \cap_{i+j=4, i, j \geq 0} H^i(\Omega, \mathcal{X}_j))
 \end{aligned}$$

provided that $\psi_0 \in \cap_{i+j=3, i, j \geq 0} H^i(\Omega, \mathcal{X}_j)$.

Up to now, we have completed the proof of Lemma 2.3.

4. Global solution. In this section, we will show that the local solution obtained in Theorem 1.1 is actually defined for $t \in \mathbb{R}^+$ if the Deborah and Reynolds numbers are small enough. To this end we derive some a priori bounds, satisfied by that solution. Notably, C_1 in this section denotes different constants depending only on n and Ω .

4.1. Some a priori estimates. Recall that the local solution (\mathbf{u}, ψ) obtained in Theorem 1.1 satisfies (1.24) and (1.25) together with (1.14) and (1.16). Now we give the detail of a priori estimates for the case $n = 3$. For $n = 2$, it can be similarly obtained.

For (1.14), the inequality (A.36) in the appendix implies that there exists constant C such that

$$(4.1) \quad \begin{aligned} \frac{1}{2} \frac{d}{dt} \|\mathbf{u}\|_{H^3}^2 + \left(\frac{\gamma}{Re} - 5\epsilon - \frac{6\epsilon}{11} \right) \|\mathbf{u}\|_{H^4}^2 &\leq \frac{5\epsilon}{11} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} [\|\mathbf{u}\|_{H^1}^{14} + \|\mathbf{u}\|_{H^1}^4 \\ &+ \|\mathbf{u}\|_{H^2}^{24} + \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^{20} + \|\psi\|_{H^2(\Omega, \mathcal{X}_0)}^{24} + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^8 + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2], \quad (n = 3). \end{aligned}$$

Here in the left side of the above inequality we omit the term $\frac{1-\gamma}{2Re} \int_{\Omega} \langle |\mathbf{m}\mathbf{m} : \nabla^3 D|^2 \rangle dx$ since it is positive. The main difficulty to obtain (4.1) is to estimate

$$\int_{\Omega} \left| \nabla^3 \left(\mathbf{D} : \int_{|\mathbf{m}|=1} \mathbf{m}\mathbf{m}\mathbf{m}\mathbf{m} \nabla \psi d\mathbf{m} \right) \nabla^4 \mathbf{u} \right| dx.$$

Now we will use the following inequalities:

$$(4.2) \quad \begin{aligned} \int_{\Omega} |\nabla^2 \mathbf{D} \int_{|\mathbf{m}|=1} \mathbf{m}\mathbf{m}\mathbf{m}\mathbf{m} \nabla \psi d\mathbf{m} \nabla^4 \mathbf{u}| dx &\leq C_1 \|\nabla^4 \mathbf{u}\|_{L^2} \|\nabla^2 \mathbf{D}\|_{L^3} \|\nabla \psi\|_{L^6(\Omega, \mathcal{X}_0)} \\ &\leq C_1 \|\mathbf{u}\|_{H^4} \|\nabla^2 \mathbf{D}\|_{H^{\frac{1}{2}}} \|\nabla \psi\|_{H^1(\Omega, \mathcal{X}_0)} \\ &\leq C_1 \|\mathbf{u}\|_{H^4} \|\mathbf{u}\|_{H^{3+\frac{1}{2}}} \|\psi\|_{H^2(\Omega, \mathcal{X}_0)} \\ &\leq C_1 \|\mathbf{u}\|_{H^4} \|\mathbf{u}\|_{H^4}^{\frac{5}{6}} \|\mathbf{u}\|_{H^1}^{\frac{1}{6}} \|\psi\|_{H^2(\Omega, \mathcal{X}_0)} \\ &\leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^1}^2 \|\psi\|_{H^2(\Omega, \mathcal{X}_0)}^{12}, \end{aligned}$$

where ϵ is a positive constant and will be chosen later. Here we used the Hölder inequality ($\int_{\Omega} |abc| dx \leq \|a\|_{L^2} \|a\|_{L^3} \|c\|_{L^6}$), Sobolev embedding theorems ($H^1(\Omega) \subset L^6(\Omega), H^{1/2}(\Omega) \subset L^3(\Omega)$ for $n = 3$), the interpolation inequality ($\|\mathbf{u}\|_{H^{3+\frac{1}{2}}} \leq C_1 \|\mathbf{u}\|_{H^4}^{\frac{5}{6}} \|\mathbf{u}\|_{H^1}^{\frac{1}{6}}$), and Young's inequality ($ab \leq \epsilon a^{12/11} + \frac{1}{\epsilon} b^{12}$), respectively. Similarly, we have

$$\begin{aligned} \int_{\Omega} |\nabla \mathbf{D} \int_{|\mathbf{m}|=1} \mathbf{m}\mathbf{m}\mathbf{m}\mathbf{m} \nabla^2 \psi d\mathbf{m} \nabla^4 \mathbf{u}| dx &\leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^1}^2 \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^4. \end{aligned}$$

Further, a different estimate from the above two is

$$\begin{aligned}
 & \int_{\Omega} |\mathbf{D} \int_{|\mathbf{m}|=1} \mathbf{m m m m} \nabla^3 \psi d\mathbf{m} \nabla^4 \mathbf{u}| dx \\
 & \leq C_1 \|\nabla^4 \mathbf{u}\|_{L^2} \|\mathbf{D}\|_{L^6} \|\nabla^3 \psi\|_{L^3(\Omega, \mathcal{X}_0)} \\
 & \leq C_1 \|\mathbf{u}\|_{H^4} \|\mathbf{D}\|_{H^1} \|\nabla^3 \psi\|_{H^{\frac{1}{2}}(\Omega, \mathcal{X}_0)} \\
 & \leq C_1 \|\mathbf{u}\|_{H^4} \|\mathbf{u}\|_{H^2} \|\psi\|_{H^{\frac{1}{6}}(\Omega, \mathcal{X}_0)}^{\frac{1}{6}} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^{\frac{5}{6}} \\
 & \leq \epsilon \left[\|\mathbf{u}\|_{H^4} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^{\frac{5}{6}} \right]^{\frac{12}{11}} + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^2}^{12} \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^{10} \\
 & \leq \epsilon \left[\|\mathbf{u}\|_{H^4}^{\frac{12}{11}} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^{\frac{10}{11}} \right] + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^2}^{12} \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^{10} \\
 (4.3) \quad & \leq \frac{6\epsilon}{11} \|\mathbf{u}\|_{H^4}^2 + \frac{5\epsilon}{11} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^2}^{12} \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^{10}.
 \end{aligned}$$

For the convect term we used the same estimate in [24]

$$\int |(\mathbf{u} \cdot \nabla) \mathbf{u} \Delta^3 \mathbf{u}| dx \leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^1}^{14},$$

and by using the results of Lemma 2.2 we have

$$\begin{aligned}
 & \int_{\Omega} |\nabla^3 \tau_e \nabla^4 \mathbf{u}| dx \leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\tau_e\|_{H^3}^2 \leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2, \\
 & \int_{\Omega} |\nabla^2 F \nabla^4 \mathbf{u}| dx \leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|F\|_{H^2}^2 \leq \epsilon \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2.
 \end{aligned}$$

Combination of all yields (4.1). For the equation of ψ , we can obtain the following estimates:

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2 + \left(\frac{\epsilon^2}{De} - 5\epsilon - \frac{\epsilon}{11} \right) \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \left(\frac{1}{De} - \frac{5\epsilon}{11} \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 \\
 & \leq \frac{5\epsilon}{11} \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} [\|\mathbf{u}\|_{H^3}^4 + \|\mathbf{u}\|_{H^2}^{24} + \|\mathbf{u}\|_{H^2}^4 + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^4 + \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^4] \\
 (4.4) \quad & + \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^{24} + \|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^4 + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2.
 \end{aligned}$$

Here the difficulties are the estimates to the nonlinear terms $\mathbf{u} \cdot \nabla \psi$ and $\mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \psi)$ in the equation. To overcome them we use the following inequalities similar to (4.2):

$$\begin{aligned}
 & \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \mathbf{u} \nabla \psi \nabla^4 \psi| d\mathbf{m} dx \leq \epsilon \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^3}^2 \|\psi\|_{H^2(\Omega, \mathcal{X}_0)}^2, \\
 & \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \mathbf{u} \nabla^2 \psi \nabla^4 \psi| d\mathbf{m} dx \leq \epsilon \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^2}^2 \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2, \\
 & \int_{\Omega} \int_{|\mathbf{m}|=1} |\mathbf{u} \nabla^3 \psi \nabla^4 \psi| d\mathbf{m} dx \leq \epsilon \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^1}^{12} \|\psi\|_{H^1(\Omega, \mathcal{X}_0)}^2.
 \end{aligned}$$

For the other term $\int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \mathcal{R} \cdot (\mathbf{m} \times \kappa \cdot \mathbf{m} \psi) \nabla^4 \psi| d\mathbf{m} dx$, we used the estimates

obtained similar to (4.3):

$$\begin{aligned} & \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla^2 \kappa \mathcal{R} \psi \nabla^4 \psi| d\mathbf{m} d\mathbf{x} \\ & \leq \frac{6\epsilon}{11} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{5\epsilon}{11} \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^1}^2 \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^{12}, \\ & \int_{\Omega} \int_{|\mathbf{m}|=1} |\kappa \nabla^2 \mathcal{R} \psi \nabla^4 \psi| d\mathbf{m} d\mathbf{x} \\ & \leq \frac{6\epsilon}{11} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{5\epsilon}{11} \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^2}^{12} \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^2, \\ & \int_{\Omega} \int_{|\mathbf{m}|=1} |\nabla \kappa \nabla \mathcal{R} \psi \nabla^4 \psi| d\mathbf{m} d\mathbf{x} \leq \epsilon \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} \|\mathbf{u}\|_{H^3}^2 \|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^2. \end{aligned}$$

Now from (4.4), we see that it is necessary to estimate $\|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^2$. But it is not difficult to obtain the following inequality in the similar way with the estimate (4.4).

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^2 + \left(\frac{\epsilon^2}{De} - 4\epsilon \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 + \frac{1}{De} \|\psi\|_{H^2(\Omega, \mathcal{X}_2)}^2 \\ (4.5) \leq & \frac{C_1}{\epsilon} \left[\|\mathbf{u}\|_{H^3}^4 + \|\mathbf{u}\|_{H^2}^{\frac{16}{3}} + \|\mathbf{u}\|_{H^1}^{16} + \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^4 + \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^8 + \|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^2 \right]. \end{aligned}$$

For the case $n = 2$, we can obtain a priori estimates by using similar approaches, and we omit the details and give the results directly as follows:

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\mathbf{u}\|_{H^3}^2 + \left(\frac{\gamma}{Re} - 5\epsilon - \frac{9\epsilon}{16} \right) \|\mathbf{u}\|_{H^4}^2 \leq \frac{7\epsilon}{16} \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \frac{C_1}{\epsilon} [\|\mathbf{u}\|_{H^1}^4 + \|\mathbf{u}\|_{H^2}^{18} \\ (4.6) \quad & + \|\psi\|_{H^2(\Omega, \mathcal{X}_0)}^{18} + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^{\frac{9}{2}}], \quad (n = 2). \end{aligned}$$

The a priori estimates for ψ are

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2 + \left(\frac{\epsilon^2}{De} - 5\epsilon - \frac{9\epsilon}{16} \right) \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \left(\frac{1}{De} - \frac{4\epsilon}{13} \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 \\ & \leq \frac{7\epsilon}{16} \|\mathbf{u}\|_{H^4}^2 + \frac{C_1}{\epsilon} [\|\mathbf{u}\|_{H^3}^4 + \|\mathbf{u}\|_{H^1}^4 + \|\mathbf{u}\|_{H^1}^8 + \|\mathbf{u}\|_{H^1}^{18} \\ (4.7) \quad & + \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^4 + \|\psi\|_{H^1(\Omega, \mathcal{X}_1)}^{18} + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^4]. \end{aligned}$$

4.2. Global existence. In this subsection we will give the proof of Theorem 1.2. Let

$$(4.8) \quad Y(t) = \|\mathbf{u}\|_{H^3(\Omega)}^2 + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2 + \|\psi\|_{H^2(\Omega, \mathcal{X}_1)}^2.$$

Combining the above estimates (4.1), (4.4), and (4.5) we have

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} Y(t) + \left(\frac{\gamma}{Re} - 6\epsilon \right) \|\mathbf{u}\|_{H^4}^2 \\ & + \left(\frac{\epsilon^2}{De} - 5\epsilon - \frac{6\epsilon}{11} \right) \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \left(\frac{1}{De} - \frac{5\epsilon}{11} \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 \\ & + \left(\frac{\epsilon^2}{De} - 4\epsilon \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 + \frac{1}{De} \|\psi\|_{H^2(\Omega, \mathcal{X}_2)}^2 \\ (4.9) \quad & \leq \frac{C_1}{\epsilon} (Y^{12} + Y^{10} + Y^7 + Y^4 + Y^2 + Y). \end{aligned}$$

Now we choose

$$(4.10) \quad \epsilon = \min \left\{ \frac{\gamma}{7Re}, \frac{\epsilon^2}{7De} \right\}, \quad \text{and} \quad \beta = \frac{C_1}{\epsilon}.$$

Since $Y(t) \leq \|u\|_{H^4}^2 + \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2$, (4.9) implies

$$(4.11) \quad Y'(t) + 2\epsilon Y(t) \leq 2\beta(Y^{12} + Y^{10} + Y^7 + Y^4 + Y^2 + Y).$$

LEMMA 4.1.

(4.11) $\implies Y(t) \leq B$ for all $t \in \mathbb{R}^+$ if $0 < B < B_0$, where B_0 is the unique positive solution of

$$(4.12) \quad B^{11} + B^9 + B^6 + B^3 + B + 1 - \frac{\epsilon}{2\beta} = 0.$$

$Y(0) \leq B \implies Y(t) \leq B$ for all $t \geq 0$.

We will prove it by contradiction. Suppose that there exists a t such that $Y(t) > B$, and define $t^* = \inf\{t \in \mathbb{R}^+, Y(t) > B\}$, then $Y(t^*) = B$ and $Y'(t^*) \geq 0$. However from (4.11) and the hypothesis made on B we deduce

$$\begin{aligned} Y'(t^*) &\leq -2\epsilon Y(t^*) + 2\beta[Y^{11}(t^*) + Y^9(t^*) + Y^7(t^*) + Y^4(t^*) + Y^2(t^*) + Y(t^*)] \\ &\leq -2\epsilon B + 2B\beta \frac{\epsilon}{2\beta} = -\epsilon B < 0, \end{aligned}$$

which contradicts the above statement. Therefore $Y(t) \leq B$ for all $t \in \mathbb{R}^+$. \square

1.2. We have seen in subsection 4.1 that a specific norm of the local solution obtained in Theorem 1.1 satisfies an inequality with the form (4.11), where C_1 depends on the domain Ω and n while ϵ depends on $\gamma, Re, \epsilon^2, De$ from (4.10). Lemma 4.1 shows that there exists a constant B_0 , depending on initial data such that

$$Y(t) \leq B, \text{ for } t \in \mathbb{R}^+, \text{ if } Y(0) \leq B < B_0.$$

But B_0 is the unique positive solution of (4.12). From (4.12), we can easily see that (4.12) possesses a unique positive solution if and only if

$$(4.13) \quad 1 - \frac{\epsilon}{2\beta} < 0,$$

which implies from (4.10) that

$$(4.14) \quad Re < \frac{\gamma}{7\sqrt{2}C_1} \quad \text{and} \quad De < \frac{\epsilon^2}{7\sqrt{2}C_1}.$$

This completes the proof of Theorem 1.2 for $n = 3$. \square

Denoting

$$Z(t) = \|u\|_{H^3(\Omega)}^2 + \|\psi\|_{H^3(\Omega, \mathcal{X}_0)}^2,$$

the addition of (4.6) and (4.7) yields

$$\begin{aligned} &\frac{1}{2} \frac{d}{dt} Z(t) + \left(\frac{\gamma}{Re} - 6\epsilon \right) \|u\|_{H^4}^2 \\ &+ \left(\frac{\epsilon^2}{De} - 6\epsilon \right) \|\psi\|_{H^4(\Omega, \mathcal{X}_0)}^2 + \left(\frac{1}{De} - \frac{4\epsilon}{13} \right) \|\psi\|_{H^3(\Omega, \mathcal{X}_1)}^2 \\ (4.15) \quad &\leq \frac{C_1}{\epsilon} (Z^9 + Y^4 + Y^{9/4} + Y^2) \quad \text{for } n = 2. \end{aligned}$$

By the same way the result of Theorem 1.2 for $n = 2$ can be obtained.

Appendix A. Proof of Lemma 2.1. Now we solve (3.2)–(3.3) with the initial value $\mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0$ by using the Galerkin approximation. Let V be the space of all divergence-free vector in $H^1(\Omega)$, and let $\{\omega^i | i \in \mathbb{N}\}$ be a basis for V . We seek an approximation to w of the form

$$(A.1) \quad w^N(\mathbf{x}, t) = \sum_{n=1}^N \rho^n(t) \omega^n(\mathbf{x}).$$

The function w^N satisfies, instead of (3.2) and (3.3),

$$(A.2) \quad \begin{aligned} w_t^N + (w^N \cdot \nabla)w^N &= \frac{\gamma}{Re} \Delta w^N + f \\ &+ \frac{1-\gamma}{2Re} \nabla \cdot [\mathbf{D}^N : A(\mathbf{x}, t)], \end{aligned}$$

$$(A.3) \quad w^N(\mathbf{x}, 0) = P_N \mathbf{u}_0(\mathbf{x}),$$

where P_N is the orthogonal projector in H_d^1 onto $W_N = Span\{\omega^i, i = 1, \dots, N\}$. The existence and uniqueness of a solution w_N to (A.2)–(A.3) with periodic boundary conditions defined on some interval $(0, T_N), T_N > 0$, is clear; in fact, the following estimate shows that $T_N = T$ for $n = 2$ and T_N suitably small for $n = 3$.

Multiplying w^N to (A.2) and integrating it, we get

$$(A.4) \quad \begin{aligned} &\frac{1}{2} \frac{d}{dt} \|w^N\|_{L^2}^2 + \frac{\gamma}{Re} \|\nabla w^N\|_{L^2}^2 + \frac{1-\gamma}{2Re} \int_{\Omega} \langle |\mathbf{D}^N : \mathbf{m}\mathbf{m}|^2 \rangle dx \\ &\leq \frac{\gamma}{2Re} \|w^N\|_{L^2}^2 + \frac{Re}{2\gamma} \|f\|_{H^{-1}}^2. \end{aligned}$$

This implies that

$$(A.5) \quad \frac{d}{dt} \|w^N\|_{L^2}^2 + \frac{\gamma}{Re} \|\nabla w^N\|_{L^2}^2 \leq \frac{Re}{\gamma} \|f\|_{H^{-1}}^2.$$

It shows that

$$(A.6) \quad \int_0^T \|\nabla w^N(t)\|_{L^2}^2 dt \leq K_1,$$

where

$$(A.7) \quad K_1 = K_1\left(\mathbf{u}_0, f, \frac{\gamma}{Re}, T\right) = \frac{Re}{\gamma} \left(\|\mathbf{u}_0\|_{L^2}^2 + \frac{Re}{\gamma} \int_0^T \|f(t)\|_{H^{-1}}^2 dt \right).$$

For $0 < s < T$, by (A.5),

$$(A.8) \quad \|w^N(s)\|_{L^2}^2 \leq K_2,$$

where

$$(A.9) \quad K_2 = K_2\left(\mathbf{u}_0, f, \frac{\gamma}{Re}, T\right) = \frac{\gamma}{Re} K_1.$$

Multiplying Δw^N by (A.2) and integrating it, we get

$$(A.10) \quad \begin{aligned} &\frac{1}{2} \frac{d}{dt} \|\nabla w^N\|_{L^2}^2 + \frac{\gamma}{Re} \|\nabla^2 w^N\|_{L^2}^2 + \frac{1-\gamma}{2Re} \int_{\Omega} \langle |\nabla \mathbf{D}^N : \mathbf{m}\mathbf{m}|^2 \rangle dx \\ &\leq \frac{1}{4} \frac{\gamma}{Re} \|\nabla w^N\|_{L^2}^2 + \frac{Re}{\gamma} \|f\|_{L^2}^2 \\ &+ \int |(w^N \cdot \nabla)w^N \Delta w^N| dx + \frac{1-\gamma}{2Re} \int |\mathbf{D}^N \nabla A \Delta w^N| dx. \end{aligned}$$

Now the following a priori of $\int |(w^N \cdot \nabla)w^N \Delta w^N| dx$ are different, depending on the dimension. We use the relation [24]

$$(A.11) \quad \int |(w \cdot \nabla)w \Delta w| dx \leq C \|w\|_{L^2}^{1/2} \|\Delta w\|_{L^2}^{3/2} \|\nabla w\|_{L^2}, \quad (n = 2),$$

$$(A.12) \quad \int |(w \cdot \nabla)w \Delta w| dx \leq C \|w\|_{L^2}^{1/4} \|\Delta w\|_{L^2}^{7/4} \|\nabla w\|_{L^2}, \quad (n = 3),$$

and the estimates

$$(A.13) \quad \int |\mathbf{D}^N \nabla A \Delta w^N| d\mathbf{x} \leq C \|\Delta w^N\|_{L^2}^{4/3} \|\nabla^2 A\|_{L^2} \|\nabla w^N\|_{L^2}^{2/3}, \quad (n = 2),$$

$$(A.14) \quad \int |\mathbf{D}^N \nabla A \Delta w^N| d\mathbf{x} \leq C \|\Delta w^N\|_{L^2}^{3/2} \|\nabla^2 A\|_{L^2} \|\nabla w^N\|_{L^2}^{1/2}, \quad (n = 3).$$

By Young's inequality (A.10) implies

$$(A.15) \quad \begin{aligned} \frac{d}{dt} \|\nabla w^N\|_{L^2}^2 + \frac{1}{2} \frac{\gamma}{Re} \|\Delta w^N\|_{L^2}^2 &\leq \frac{2Re}{\gamma} \|f\|_{L^2}^2 \\ &+ C \|w^N\|_{L^2}^2 \|\nabla w^N\|_{L^2}^4 + C \|\nabla^2 A\|_{L^2}^4 \|\nabla w^N\|_{L^2}^2, \quad (n = 2), \end{aligned}$$

$$(A.16) \quad \begin{aligned} \frac{d}{dt} \|\nabla w^N\|_{L^2}^2 + \frac{3}{2} \frac{\gamma}{Re} \|\Delta w^N\|_{L^2}^2 &\leq \frac{2Re}{\gamma} \|f\|_{L^2}^2 \\ &+ C \|w^N\|_{L^2}^2 \|\nabla w^N\|_{L^2}^8 + C \|\nabla^2 A\|_{L^2}^4 \|\nabla w^N\|_{L^2}^2, \quad (n = 3). \end{aligned}$$

Since $\psi^l \in L^\infty([0, T], H^3(\Omega, \mathcal{X}_0))$, we can obtain

$$(A.17) \quad |\nabla A(x, t)| \in L^\infty([0, T]; H^2(\Omega)).$$

Then by the same way as the a priori estimate of Theorem 3.2 in [24] we can get for $n = 2$ in virtue of (2.14) and (3.1)

$$(A.18) \quad \sup_{t \in [0, T]} \|\nabla w^N\|_{L^2}^2 \leq K_3,$$

where

$$K_3 = K_3 \left(\mathbf{u}_0, f, \frac{\gamma}{Re}, T \right) = \left(\|\mathbf{u}_0\|_{H^1}^2 + C \int_0^T \|f(t)\|_{L^2}^2 dt \right) \exp(CK_1 K_2 + CK^3 T).$$

And

$$(A.19) \quad \int_0^T \|\Delta w^N\|_{L^2}^2 dt \leq K_4,$$

where

$$K_4 = K_4(\mathbf{u}_0, f, \frac{\gamma}{Re}, T) = C\gamma \left(\|\mathbf{u}_0\|_{H^1}^2 + C \int_0^T \|f(t)\|_{L^2}^2 dt + CK_2 K_3^2 T + CK^3 K_3 T \right).$$

For $n = 3$, we will estimate in the following from (A.16). Let $y(t) = \|\nabla w^N\|_{L^2}^2 + 1$, $a(t) = C \|w^N\|_{L^2}^2$, $b(t) = C \|\nabla^2 A\|_{L^2}^4$, and $c(t) = \frac{2Re}{\gamma} \|f\|_{L^2}^2$. Then from (A.16), we have the inequality

$$(A.20) \quad \frac{dy}{dt} \leq [a(t) + b(t) + c(t)]y^4,$$

and

$$(A.21) \quad \int_0^T [a(t) + b(t) + c(t)]dt \leq T(CK_2 + CK^6 + CK^2),$$

where we used the condition of (2.14) and the estimate (3.8) and above estimates. Then they imply that for $t \in [0, T]$ and $T < 1/[(CK_2 + CK^4 + \frac{Re}{\gamma}K^2)y^3(0)] \triangleq T_0$

$$(A.22) \quad y(t) \leq \frac{y(0)}{\sqrt[3]{1 - (CK_2 + CK^4 + \frac{Re}{\gamma}K^2)y^3(0)T}} \triangleq K_5.$$

Therefore, when $T < T_0$, we have

$$(A.23) \quad \sup_{s \in [0, T]} \|\nabla w^N\|_{L^2}^2 \leq K_5$$

and

$$(A.24) \quad \int_0^T \|\Delta w^N\|_{L^2}^2 dt \leq K_6,$$

where

$$K_6 = K_6(\mathbf{u}_0, f, \frac{\gamma}{Re}, T) = C \left(\|\mathbf{u}_0\|_{H^1}^2 + C \int_0^T \|f(t)\|_{L^2}^2 dt + CK_1K_5^4T + CK^4K_5T \right).$$

Multiplying $\Delta^2 w^N$ by (A.2) and integrating it, we get

$$(A.25) \quad \begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\nabla^2 w^N\|_{L^2}^2 + \frac{\gamma}{Re} \|\nabla^3 w^N\|_{L^2}^2 + \frac{1-\gamma}{4Re} \int_{\Omega} \langle |\nabla^2 \mathbf{D}^N : \mathbf{m}\mathbf{m}|^2 \rangle dx \\ & \leq \frac{1}{4} \frac{\gamma}{Re} \|\nabla \Delta w^N\|_{L^2}^2 + C \|\nabla f\|_{L^2}^2 \\ & \quad + \frac{1-\gamma}{2Re} \int [|\mathbf{D}^N \nabla^2 A \nabla \Delta w^N| + |\nabla \mathbf{D}^N \nabla A \nabla \Delta w^N|] dx \\ & \quad + \int |(w^N \cdot \nabla) w^N \Delta^2 w^N| dx. \end{aligned}$$

In the following we will estimate (A.25). In virtue of the inequality (p. 31 [24])

$$(A.26) \quad \int |(w^N \cdot \nabla) w^N \Delta^r w^N| dx \leq \frac{1}{4} \frac{\gamma}{Re} \|w^N\|_{H^{r+1}}^2 + C \|w^N\|_{H^1}^{2r} \quad (n = 2),$$

$$(A.27) \quad \int |(w^N \cdot \nabla) w^N \Delta^r w^N| dx \leq \frac{1}{4} \frac{\gamma}{Re} \|w^N\|_{H^{r+1}}^2 + C \|w^N\|_{H^1}^{4r+2} \quad (n = 3).$$

And the estimates

$$\int |\mathbf{D}^N \nabla^2 A \nabla^3 w^N| dx \leq C \|\nabla^3 w^N\|_{L^2} \|\nabla A\|_{L^2}^{1/3} \|\nabla^3 A\|_{L^2}^{2/3} \|\nabla^2 w^N\|_{L^2}, \quad (n = 2),$$

$$\int |\mathbf{D}^N \nabla^2 A \nabla^3 w^N| dx \leq C \|\nabla^3 w^N\|_{L^2} \|\nabla A\|_{L^2}^{1/4} \|\nabla^3 A\|_{L^2}^{3/4} \|\nabla^2 w^N\|_{L^2}, \quad (n = 3),$$

$$\int |\nabla \mathbf{D}^N \nabla A \nabla^3 w^N| dx \leq C \|\nabla^3 w^N\|_{L^2}^{5/3} \|\nabla^2 A\|_{L^2} \|\nabla w^N\|_{L^2}^{1/3}, \quad (n = 2),$$

$$\int |\nabla \mathbf{D}^N \nabla A \nabla^3 w^N| dx \leq C \|\nabla^3 w^N\|_{L^2}^{7/4} \|\nabla^2 A\|_{L^2} \|\nabla w^N\|_{L^2}^{1/4}. \quad (n = 3).$$

Thus, we have

$$(A.28) \quad \begin{aligned} \frac{d}{dt} \|w^N\|_{H^2}^2 + \frac{1}{2} \frac{\gamma}{Re} \|w^N\|_{H^3}^2 &\leq C \|w^N\|_{H^1}^4 + \frac{Re}{\gamma} \|f\|_{H^1}^2 \\ &+ C \|A\|_{H^1}^{2/3} \|A\|_{H^3}^{4/3} \|w^N\|_{H^2}^2 + C \|A\|_{H^2}^6 \|w^N\|_{H^2}^2, \quad (n = 2), \end{aligned}$$

$$(A.29) \quad \begin{aligned} \frac{d}{dt} \|w^N\|_{H^2}^2 + \frac{1}{2} \frac{\gamma}{Re} \|w^N\|_{H^3}^2 &\leq C \|w^N\|_{H^1}^{10} + \frac{Re}{\gamma} \|f\|_{H^1}^2 \\ &+ C \|A\|_{H^1}^{1/2} \|A\|_{H^3}^{3/2} \|w^N\|_{H^2}^{3/2} + C \|A\|_{H^2}^8 \|w^N\|_{H^1}^2, \quad (n = 3). \end{aligned}$$

Similarly we can obtain the following estimates: for $n = 2$,

$$(A.30) \quad \sup_{t \in [0, T]} \|w^N(t)\|_{H^2}^2 \leq [CK_2^2 T + CKK_2 T + C\|f\|_{L^2(0, T; H^1)}^2 + \|\mathbf{u}_0\|_{H^2}^2] e^{CK^2 T + CK^6 T} \triangleq K_7,$$

and

$$(A.31) \quad \int_0^T \|w^N(t)\|_{H^3}^2 dt \leq K_8,$$

where

$$K_8 = C \left(CK_3^2 T + C\|f\|_{L^2(0, T; H^1)}^2 + \|\mathbf{u}_0\|_{H^2}^2 + CK^2 K_7 K_1 T + CK^6 K_7 T \right).$$

While for $n = 3$,

$$(A.32) \quad \sup_{t \in [0, T]} \|w^N(t)\|_{H^2}^2 \leq K_9,$$

where

$$K_9 = \left[CK_5^5 T + CK^8 K_1 + \frac{Re}{\gamma} \|f\|_{L^2(0, T; H^1)}^2 + \|\mathbf{u}_0\|_{H^2}^2 \right] e^{CK^2 T}.$$

Further, we have

$$(A.33) \quad \int_0^T \|w^N(t)\|_{H^3}^2 dt \leq K_{10},$$

where

$$K_{10} = C \left(CK_5^5 T + C\|f\|_{L^2(0, T; H^1)}^2 + \|\mathbf{u}_0\|_{H^2}^2 + CK^8 K_1 T + CK^2 K_9 T \right).$$

This shows for T given,

$$(A.34) \quad w^N \in L^\infty([0, T]; H^2(\Omega)) \cap L^2([0, T]; H^3(\Omega))$$

provided that $f \in L^2([0, T]; H^1(\Omega))$ and $A \in L^\infty([0, T]; H^3(\Omega)) \cap L^2([0, T]; H^4(\Omega))$.

Similarly we can obtain the estimates

$$(A.35) \quad \begin{aligned} \frac{d}{dt} \|w^N\|_{H^3}^2 + \frac{1}{2} \frac{\gamma}{Re} \|w^N\|_{H^4}^2 &\leq C \|w^N\|_{H^1}^6 + \frac{Re}{\gamma} \|f\|_{H^2}^2 \\ &+ C \|A\|_{H^2}^9 \|w^N\|_{H^1}^2 + C \|A\|_{H^3}^4 \|w^N\|_{H^1}^2 \\ &+ C \|A\|_{H^1}^{4/9} \|A\|_{H^4}^{14/9} \|w^N\|_{H^2}^2, \quad (n = 2), \end{aligned}$$

$$(A.36) \quad \begin{aligned} \frac{d}{dt} \|w^N\|_{H^3}^2 + \frac{1}{2} \frac{\gamma}{Re} \|w^N\|_{H^4}^2 &\leq C \|w^N\|_{H^1}^{14} + \frac{Re}{\gamma} \|f\|_{H^2}^2 \\ &+ C \|A\|_{H^2}^{12} \|w^N\|_{H^1}^2 + C \|A\|_{H^3}^4 \|w^N\|_{H^1}^2 \\ &+ C \|A\|_{H^1}^{1/3} \|A\|_{H^4}^{5/3} \|w^N\|_{H^2}^2, \quad (n = 3). \end{aligned}$$

From (A.35), for $n = 2$ we have

$$(A.37) \quad \sup_{t \in [0, T]} \|w^N(t)\|_{H^3}^2 \leq K_{11},$$

where

$$K_{11} = \left[\|u_0\|_{H^3}^2 + \frac{Re}{\gamma} \|f\|_{L^2(0, T; H^2)}^2 + CK_3^3 T + CK^4 K_1 + CK^9 K_1 \right] e^{CK^4 T}.$$

Furthermore, we have

$$(A.38) \quad \int_0^T \|w^N(t)\|_{H^4}^2 dt \leq K_{12},$$

where

$$K_{12} = C \left(\|u_0\|_{H^3}^2 + C \|f\|_{L^2(0, T; H^2)}^2 + CK_3^3 T + CK^4 K_1 T + CK^2 K_{11} T + CK^2 K_7 T \right).$$

Thus we know that $\phi_1(K, T) = K_{12} + K_{11}$ for $n = 2$.

From (A.36), for $n = 3$ we have

$$(A.39) \quad \sup_{t \in [0, T]} \|w^N(t)\|_{H^3}^2 \leq K_{13},$$

where

$$K_{13} = [\|u_0\|_{H^3}^2 + C \|f\|_{L^2(0, T; H^2)}^2 + CK_3^7 T + CK^{12} K_3 T + CK^4 K_3 T] e^{CK^2 T}.$$

Moreover,

$$(A.40) \quad \int_0^T \|w^N(t)\|_{H^4}^2 dt \leq K_{14},$$

where

$$K_{14} = C \left(\|u_0\|_{H^3}^2 + C \|f\|_{L^2(0, T; H^2)}^2 + CK_3^7 T + CK^{12} K_3 T + CK^4 K_3 T + CK^2 K_{13} T \right).$$

Thus we know that $\phi_1(K, T) = K_{14} + K_{13}$ for $n = 3$.

Therefore, this implies that for T given,

$$(A.41) \quad w^N \in L^\infty([0, T]; H^3(\Omega)) \cap L^2([0, T]; H^4(\Omega))$$

provided that $f \in L^2([0, T]; H^2(\Omega))$ and $A \in L^\infty([0, T], H^3(\Omega)) \cap L^2([0, T], H^4(\Omega))$.

Passing $N \rightarrow \infty$, we can obtain the estimate (2.16).

Acknowledgments. The authors are very grateful to the referee for his many valuable suggestions. We are very grateful to Professors Weinan E and Chun Liu for their helpful discussions.

REFERENCES

- [1] T. AUBIN, *Some Nonlinear Problems in Riemannian Geometry*, Springer, New York, 1998.
- [2] A.N. BERIS AND B.J. EDWARDS, *Thermodynamics of Flowing Systems with Internal Microstructure*, Oxford Science, New York, 1994.

- [3] C.V. CHAUBAL, A. SRINIVASAN, Ö. EĞECIOĞLU, AND L.G. LEAL, *Smoothed particle hydrodynamics techniques for the solution of kinetic theory problems*, Part 1: Method, *J. Non-Newtonian Fluid Mech.*, 70 (1997), pp. 125–154.
- [4] P. CONSTANTIN, I.G. KEVREKIDS, AND E.S. TITI, *Asymptotic states of a smoluchowski equation*, *Discrete Continuous Dynamical Systems*, 11 (2004), pp. 101–112.
- [5] P. CONSTANTIN AND C. FOIAS, *Navier-Stokes Equations*, The University of Chicago Press, Chicago and London, 1988.
- [6] M. DOI AND S.F. EDWARDS, *The Theory of Polymer Dynamics*, Oxford University Press, Oxford, 1986.
- [7] Q. DU, C. LIU, AND P. YU, *FENE dumbbell model and its several linear and nonlinear closure approximations*, *Multiscale Model. Simul.*, 4 (2005), pp. 709–731.
- [8] W.N. E AND P.W. ZHANG, *A molecular kinetic theory of inhomogeneous liquid crystal flow and the small Deborah number limit*, *Meth. Appl. Anal.*, 13 (2006), pp. 181–198.
- [9] J. ERICKSEN, *Liquid crystals with variable degree of orientation*, *Arch. Rat. Mech. Anal.*, 113 (1991), pp. 97–120.
- [10] V. FARAONI, M. GROSSO, S. CRESCITELLI, AND P.L. MAFFETTONE, *The rigid-rod model for nematic polymers: An analysis of the shear flow problem*, *J. Rheol.*, 43 (1999), pp. 829–843.
- [11] P.G. DE GENNES AND J. PROST, *The Physics of Liquid Crystals*, 2nd ed., Oxford Science Publications, New York, 1993.
- [12] C. GUILLOPE AND J.C. SAUT, *Existence results for the flow of viscoelastic fluids with a differential constitutive law*, *Nonl. Anal. TMA*, 15 (1990), pp. 849–869.
- [13] N. KUZUU AND M. DOI, *Constitutive equation for nematic liquid crystals under weak velocity gradient derived from a molecular kinetic equation*, *J. Phys. Soc. Japan*, 52 (1983), pp. 3486–3494.
- [14] O. LADYZENSKAJA, V. SOLONNIKOV, AND N. URAL'CEVA, *Linear and Quasilinear Equations of Parabolic Type*, *Translations of Mathematical Monographs*, Am. Math. Soc., Providence, RI, 1968.
- [15] T.J. LI, H. ZHANG AND P.W. ZHANG, *Local existence for the dumbbell model of polymeric fluids*, *Comm. Partial Differential Equations*, 29 (2004), pp. 903–923.
- [16] F.H. LIN, C. LIU, AND P. ZHANG, *On hydrodynamics of viscoelastic fluids*, *Comm. Pure Appl. Math.*, 58 (2005), pp. 1–35.
- [17] F.H. LIN, C. LIU, AND P. ZHANG, *On a micro-macro model for polymeric fluids near equilibrium*, *Comm. Pure Appl. Math.*, 60 (2006), pp. 838–866.
- [18] G. MARRUCCI AND F. GRECO, *The elastic constants of Maier-Saupe rodlike molecule nematics*, *Mol. Cryst. Liq. Cryst.*, 206 (1991), pp. 17–30.
- [19] M. RENARDY, *Local existence of solutions of the Dirichlet initial boundary problem for incompressible hypoelastic materials*, *SIAM J. Math. Anal.*, 21 (1990), pp. 1369–1385.
- [20] M. RENARDY, *An existence theorem for model equations resulting from kinetic theories*, *SIAM J. Math. Anal.*, 22 (1991), pp. 313–327.
- [21] A.D. REY AND T. TSUJI, *Orientation mode selection mechanisms for sheared nematic liquid crystalline materials*, *Phys. Rev. E*, 57 (1998), pp. 5610–5625.
- [22] A.D. REY AND T. TSUJI, *Recent advances in theoretical liquid crystal rheology*, *Macromol. Theory Simul.*, 7 (1998), pp. 623–639.
- [23] T.C. SIDERIS AND B. THOMASES, *Global existence for 3D incompressible isotropic elastodynamics via the incompressible limit*, *Comm. Pure Appl. Math.*, 58 (2005), pp. 750–788.
- [24] R. TEMAN, *Navier-Stokes Equations and Nonlinear Functional Analysis*, SIAM, Philadelphia, 1995.
- [25] Q. WANG, *A hydrodynamic theory for solutions of nonhomogeneous nematic liquid crystalline polymers of different configurations*, *J. Chem. Phys.*, 116 (2002), pp. 9120–9136.
- [26] H.J. YU AND P.W. ZHANG, *A kinetic-hydrodynamic simulation of microstructure of liquid crystal polymers in plane shear flow*, *J. Non-Newtonian Fluid Mech.*, 141 (2007), pp. 116–127.
- [27] H.J. YU, G.H. JI, AND P.W. ZHANG, *A nonhomogeneous kinetic model of liquid crystal polymers and its thermodynamic closure approximation*, preprint.
- [28] H. ZHANG AND P.W. ZHANG, *A theoretical and numerical study for the rod-like model of a polymeric fluid*, *J. Comp. Math.*, 22 (2004), pp. 319–330.

A DESCRIPTION OF SEISMIC ACOUSTIC WAVE PROPAGATION IN POROUS MEDIA VIA HOMOGENIZATION*

ANVARBEK MEIRMANOV†

Abstract. We consider a linear system of differential equations describing the joint motion of an elastic porous body and a fluid occupying the porous space. A rigorous justification is performed for the homogenization procedures under various conditions imposed on the physical parameters as the dimensionless size of the pores tends to zero, while the porous body is geometrically periodic and the process's characteristic time is sufficiently small. Such models describe the propagation of seismic acoustic waves. In the present paper, we derive the homogenized equations, which are different types of nonstandard wave equations depending on the relations between the physical parameters. The proofs are based on Nguetseng's two-scale convergence method of homogenization in periodic structures.

Key words. Stokes equations, Lamé's equations, wave equation, two-scale convergence, homogenization of periodic structures

AMS subject classifications. 35M20, 74F10, 76S05

DOI. 10.1137/070697483

1. Introduction. In the present paper, we deal with a problem of joint motion of a deformable solid (the Ω_s) perforated by a system of channels or pores (the Ω_f) and a fluid occupying the pore space. In a domain $\Omega \subset \mathbf{R}^3$, the dimensionless displacement vector \mathbf{w} of the continuum medium in the dimensionless variables

$$\mathbf{x}' = L\mathbf{x}, \quad t' = \tau t, \quad \mathbf{w}' = \frac{L^2}{g\tau^2}\mathbf{w}, \quad \rho'_s = \rho_0\rho_s, \quad \rho'_f = \rho_0\rho_f, \quad \mathbf{F}' = g\mathbf{F}$$

satisfies the differential equation

$$(1.1) \quad \bar{\rho} \frac{\partial^2 \mathbf{w}}{\partial t'^2} = \operatorname{div} P + \bar{\rho} \mathbf{F},$$

where

$$(1.2) \quad P = \bar{\chi} P^f + (1 - \bar{\chi}) P^s,$$

$$(1.3) \quad P^f = \alpha_\mu D \left(x, \frac{\partial \mathbf{w}}{\partial t'} \right) - \left(p_f - \alpha_\nu \operatorname{div} \frac{\partial \mathbf{w}}{\partial t'} \right) I,$$

$$(1.4) \quad P^s = \alpha_\lambda D(x, \mathbf{w}) + \alpha_\eta (\operatorname{div} \mathbf{w}) I,$$

$$(1.5) \quad p_f + \bar{\chi} \alpha_p \operatorname{div} \mathbf{w} = 0.$$

Hereafter, we use the notation

$$D(x, \mathbf{u}) = (1/2) (\nabla_x \mathbf{u} + (\nabla_x \mathbf{u})^T), \quad \bar{\rho} = \bar{\chi} \rho_f + (1 - \bar{\chi}) \rho_s,$$

*Received by the editors July 17, 2007; accepted for publication (in revised form) June 24, 2008; published electronically October 31, 2008. This work was partially supported by RFBR grant 08-05-00265.

<http://www.siam.org/journals/sima/40-3/69748.html>

†Mathematics Department, Belgorod State University, ul.Pobedi 85, 308015 Belgorod, Russia (meirmanov@bsu.edu.ru).

where I is the unit tensor, the given function $\bar{\chi}(\mathbf{x})$ is the characteristic function of the pore space, the given function $\mathbf{F}(\mathbf{x}, t)$ is the dimensionless vector of distributed mass forces, P^f is the liquid stress tensor, P^s is the stress tensor in the solid skeleton, and p_f is the liquid pressure.

Equations (1.1)–(1.5) mean that the displacement vector \mathbf{w} satisfies the Stokes equations in the pore space Ω_f and the Lamé equations in the solid skeleton Ω_s .

On the “solid skeleton–pore space” common boundary Γ , the displacement vector \mathbf{w} and the liquid pressure p_f satisfy the usual continuity condition

$$(1.6) \quad [\mathbf{w}](\mathbf{x}_0, t) = 0, \quad \mathbf{x}_0 \in \Gamma, \quad t \geq 0,$$

and the momentum conservation law in the form

$$(1.7) \quad [P \cdot \mathbf{n}](\mathbf{x}_0, t) = 0, \quad \mathbf{x}_0 \in \Gamma, \quad t \geq 0,$$

where $\mathbf{n}(\mathbf{x}_0)$ is the unit normal to the boundary at the point $\mathbf{x}_0 \in \Gamma$ and

$$[\varphi](\mathbf{x}_0, t) = \varphi_{(s)}(\mathbf{x}_0, t) - \varphi_{(f)}(\mathbf{x}_0, t),$$

$$\varphi_{(s)}(\mathbf{x}_0, t) = \lim_{\substack{\mathbf{x} \rightarrow \mathbf{x}_0 \\ \mathbf{x} \in \Omega_s}} \varphi(\mathbf{x}, t), \quad \varphi_{(f)}(\mathbf{x}_0, t) = \lim_{\substack{\mathbf{x} \rightarrow \mathbf{x}_0 \\ \mathbf{x} \in \Omega_f}} \varphi(\mathbf{x}, t).$$

The problem is endowed with the homogeneous initial and boundary conditions

$$(1.8) \quad \mathbf{w}(\mathbf{x}, 0) = 0, \quad \frac{\partial \mathbf{w}}{\partial t}(\mathbf{x}, 0) = 0, \quad \mathbf{x} \in \Omega,$$

$$(1.9) \quad \mathbf{w}(\mathbf{x}, t) = 0, \quad \mathbf{x} \in S = \partial\Omega, \quad t \geq 0.$$

The dimensionless constants α_i ($i = \tau, \nu, \dots$) are defined by the formulas

$$\alpha_\mu = \frac{2\mu\tau}{L^2\rho_0}, \quad \alpha_\lambda = \frac{2\lambda\tau^2}{L^2\rho_0}, \quad \alpha_\nu = \frac{\nu\tau}{L^2\rho_0},$$

$$\alpha_p = \rho_f c_f^2 \frac{\tau^2}{L^2}, \quad \alpha_\eta = \frac{\eta\tau^2}{L^2\rho_0} = \rho_s c_s^2 \frac{\tau^2}{L^2},$$

where μ is the fluid viscosity, ν is the bulk fluid viscosity, λ and η are elastic Lamé’s constants, c_f is the speed of sound in fluids, c_s is the speed of sound in solids, L is the characteristic size of the domain under study, τ is the characteristic time of the process, ρ_f and ρ_s are the respective mean dimensionless densities of the liquid and solid phases correlated with the mean density of water ρ_0 , and g is the value of acceleration due to gravity.

The corresponding mathematical model described by system (1.1)–(1.9) is commonly used (see [2], [9]) and contains a natural small parameter ε , which is the pore characteristic size l divided by the characteristic size L of the entire porous body:

$$\varepsilon = \frac{l}{L}.$$

Our aim is to derive all possible limiting regimes (the homogenized equations) as $\varepsilon \searrow 0$. Such an approximation significantly simplifies the original problem and at the same time preserves all of its main features. But even this approach is too difficult

to be realized, and some additional simplifying assumptions are necessary. In terms of geometrical properties of the medium, it is most expedient to simplify the problem by postulating that the porous structure is periodic.

We impose the following constraints.

(1) The domain $\Omega = (0, 1)^3$ is a periodic repetition of an elementary cell $Y^\varepsilon = \varepsilon Y$, where $Y = (0, 1)^3$ and the quantity $1/\varepsilon$ is an integer so that Ω always contains an integer number of elementary cells Y^ε .

(2) Let Y_s be the “solid part” of Y , and let the “liquid part” Y_f of Y be its open complement. We write $\gamma = \partial Y_f \cap \partial Y_s$ and assume that γ is a Lipschitz continuous surface.

(3) The pore space Ω_f^ε is a periodic repetition of the elementary cell εY_f , and the solid skeleton Ω_s^ε is a periodic repetition of the elementary cell εY_s . The Lipschitz continuous boundary $\Gamma^\varepsilon = \partial \Omega_s^\varepsilon \cap \partial \Omega_f^\varepsilon$ is a periodic repetition in Ω of the boundary $\varepsilon \gamma$.

(4) Here the essential assumptions are those last three on the geometry of the elementary cells Y_s and Y_f and the domains Ω_s^ε and Ω_f^ε . As for the first assumption, we take the simplest structure of Ω (namely, the cube) just to simplify the procedure. In principle, for the domain Ω we can choose any bounded domain with a Lipschitz continuous boundary $S = \partial \Omega$.

Under these assumptions, we have

$$\bar{\chi}(\mathbf{x}) = \chi^\varepsilon(\mathbf{x}) = \chi\left(\frac{\mathbf{x}}{\varepsilon}\right),$$

$$\bar{\rho} = \rho^\varepsilon(\mathbf{x}) = \chi^\varepsilon(\mathbf{x})\rho_f + (1 - \chi^\varepsilon(\mathbf{x}))\rho_s,$$

where $\chi(\mathbf{y})$ is the characteristic function of Y_f in Y .

We assume that all dimensionless parameters depend on the small parameter ε and the (finite or infinite) limits exist:

$$\lim_{\varepsilon \searrow 0} \alpha_\mu(\varepsilon) = \mu_0, \quad \lim_{\varepsilon \searrow 0} \alpha_\lambda(\varepsilon) = \lambda_0, \quad \lim_{\varepsilon \searrow 0} \alpha_\nu(\varepsilon) = \nu_0,$$

$$\lim_{\varepsilon \searrow 0} \alpha_\eta(\varepsilon) = \eta_0, \quad \lim_{\varepsilon \searrow 0} \alpha_p(\varepsilon) = p_*,$$

$$\lim_{\varepsilon \searrow 0} \frac{\alpha_\mu}{\varepsilon^2} = \mu_1, \quad \lim_{\varepsilon \searrow 0} \frac{\alpha_\lambda}{\varepsilon^2} = \lambda_1.$$

The first research aiming to find the limiting regimes in the case where the skeleton was an absolutely rigid body was carried out by Sanchez-Palencia and Tartar. Sanchez-Palencia [9, sect. 7.2] formally obtained Darcy’s law of filtration using the method of two-scale asymptotic expansions, and Tartar [9, Appendix] rigorously justified the homogenization procedure. Using the same method of two-scale expansions, Burridge and Keller [2] formally derived a system of Biot’s equations from problem (1.1)–(1.9) in the case where the parameter α_μ was of order ε^2 and the rest of the coefficients were fixed independent of ε . Under the same assumptions as in [2], a rigorous justification of Biot’s model was given by Nguetseng [8] and later by Clopeaut et al. [3]. The most general case of problem (1.1)–(1.9) where

$$\mu_0, \lambda_0^{-1}, \nu_0, p_*^{-1}, \eta_0^{-1} < \infty$$

was studied in [6].

All these authors used Nguetseng's two-scale convergence method [7, 5].

In the present paper, we use the same method to investigate all possible limiting regimes in problem (1.1)–(1.9) in the cases where

$$\nu_0, p_*, \eta_0 < \infty; \quad \mu_0 = \lambda_0 = 0, \quad 0 < p_*, \eta_0.$$

These cases correspond to the seismic acoustic wave propagation, where all the processes on distances of tens of thousands of meters ($L \nearrow \infty$) come to an end in several seconds ($\tau \searrow 0$).

We show that the homogenized equations are different types of nonstandard wave equations for a two- or one-velocity continuum (Theorem 2.2).

This is a very interesting fact: initially a one-velocity continuum becomes a two-velocity continuum after the homogenization procedure, which appears to be the result of different smoothness of the solution in the solid and liquid components:

$$\int_{\Omega} \alpha_{\mu}(\varepsilon) \chi^{\varepsilon} |\nabla \mathbf{w}^{\varepsilon}|^2 dx \leq C_0, \quad \int_{\Omega} \alpha_{\lambda}(\varepsilon) (1 - \chi^{\varepsilon}) |\nabla \mathbf{w}^{\varepsilon}|^2 dx \leq C_0,$$

where C_0 is a constant independent of the small parameter ε . To preserve the best properties of the solution, we must use the well-known extension lemma [1, 4] and extend the solution from the solid part to the liquid part and conversely. At this stage, the criteria μ_1 and λ_1 become crucial. Namely, let $\mathbf{w}_f^{\varepsilon}$ ($\mathbf{w}_s^{\varepsilon}$) be an extension of the liquid (solid) displacements to the solid (liquid) part, and let $\mu_1 = \lambda_1 = \infty$. Then the limiting (homogenized) system describes the one-velocity continuum. This is because of the fact that each of the sequences $\{\mathbf{w}^{\varepsilon}\}$, $\{\mathbf{w}_f^{\varepsilon}\}$, and $\{\mathbf{w}_s^{\varepsilon}\}$ two-scale converges to a function independent of the fast variable. This statement easily follows from Nguetseng's theorem.

If $\mu_1 < \infty$ and $\lambda_1 = \infty$ or $\mu_1 = \infty$ and $\lambda_1 < \infty$, then the homogenized systems describe the two-velocity continuum.

Finally, we note that, in practice, to solve a real physical problem in, for example, acoustics, one does not want to carry out the limiting procedure but, instead, wants to find a simple and reliable mathematical model describing the process. But there is only one exact (sufficiently reliable) mathematical model (1.1)–(1.9) with given physical constants (densities, viscosities, etc.), the characteristic size L of the physical domain under study, and the characteristic time τ of the physical process. The small parameter ε and the dimensionless quantities α_{μ} , α_{λ} , α_p , ... are functions of them. Changing the values of L and τ within some reasonable limits, one may find some rules for the behavior of the dimensionless quantities as the small parameter tends to zero. All possible limits of these quantities are described by conditions on μ_0 , λ_0 , μ_1 , ... and, as was mentioned above, each homogenized system corresponds to a given combination of them. Thus, for a given physical situation, there exists a combination of dimensionless criteria, which would suggest the choice of the form of the homogenized system for obtaining the exact mathematical model. Therefore, to find all possible homogenized systems is very important from both mathematical and practical standpoints.

2. Main results. There are various equivalent (in the sense of distributions) forms of representation of (1.1) and boundary conditions (1.6)–(1.7). In what follows, it is convenient to write them in the form of the integral identities.

We say that four functions $(\mathbf{w}^\varepsilon, p_f^\varepsilon, p_s^\varepsilon, q^\varepsilon)$ are a generalized solution of problem (1.1)–(1.9) if they satisfy the regularity conditions

$$(2.1) \quad \mathbf{w}^\varepsilon, \nabla \mathbf{w}^\varepsilon, p_f^\varepsilon, p_s^\varepsilon, q^\varepsilon \in L^2(\Omega_T)$$

in the domain $\Omega_T = \Omega \times (0, T)$, boundary condition (1.9) in the trace sense, (1.5) and the equations

$$(2.2) \quad p_s^\varepsilon + (1 - \chi^\varepsilon)\alpha_\eta \operatorname{div} \mathbf{w}^\varepsilon = 0,$$

$$(2.3) \quad q^\varepsilon = p_f^\varepsilon + \frac{\alpha_\nu}{\alpha_p} \frac{\partial p_f^\varepsilon}{\partial t}$$

a.e. in Ω_T , and, finally, the integral identity

$$(2.4) \quad \int_{\Omega_T} \left(\rho^\varepsilon \mathbf{w}^\varepsilon \cdot \frac{\partial^2 \boldsymbol{\varphi}}{\partial t^2} - \chi^\varepsilon \alpha_\mu D(\mathbf{x}, \mathbf{w}^\varepsilon) : D \left(\mathbf{x}, \frac{\partial \boldsymbol{\varphi}}{\partial t} \right) - \rho^\varepsilon \mathbf{F} \cdot \boldsymbol{\varphi} \right. \\ \left. + \{ (1 - \chi^\varepsilon) \alpha_\lambda D(\mathbf{x}, \mathbf{w}^\varepsilon) - (q^\varepsilon + p_s^\varepsilon) I \} : D(\mathbf{x}, \boldsymbol{\varphi}) \right) dx dt = 0 \quad \Bigg\}$$

for all smooth vector-functions $\boldsymbol{\varphi} = \boldsymbol{\varphi}(\mathbf{x}, t)$ such that $\boldsymbol{\varphi}|_{\partial\Omega} = \boldsymbol{\varphi}|_{t=T} = \partial \boldsymbol{\varphi} / \partial t|_{t=T} = 0$.

In this definition, we changed the form of representation of the stress tensor P in the integral identity (2.4) by introducing two new unknown functions, q^ε and p_s^ε , which in a certain sense have the meaning of pressure. In what follows, we call functions q^ε and p_s^ε the liquid and the solid pressure, respectively, and regard (2.3) as the state equation and equations (1.5) and (2.2) as the continuity equations. This special choice of the continuity and state equations simplifies the use of the homogenization procedure.

In (2.4), by $A : B$ we denote the convolution (or, equivalently, the inner tensor product) of two second-rank tensors along the both indices, i.e., $A : B = \operatorname{tr} (B^* \circ A) = \sum_{i,j=1}^3 A_{ij} B_{ji}$.

Theorems 2.1–2.2 are the main results of the paper.

THEOREM 2.1. *Let $\mathbf{F} \in L^2(\Omega) \times L^2(\Omega) \times L^2(\Omega)$, $\varepsilon > 0$, $\mathbf{w}_f^\varepsilon, \mathbf{w}_s^\varepsilon \in L^\infty(0, T; W_2^1(\Omega))$, and (1.1)–(1.9) hold.*

$$(2.5) \quad \max_{0 \leq t \leq T} \left\| \left| \frac{\partial \mathbf{w}^\varepsilon}{\partial t} \right| + \sqrt{\alpha_\mu} \chi^\varepsilon |\nabla \mathbf{w}^\varepsilon| + (1 - \chi^\varepsilon) \sqrt{\alpha_\lambda} |\nabla \mathbf{w}^\varepsilon| \right\|_{2, \Omega} (t) \leq C_0,$$

$$(2.6) \quad \|q^\varepsilon\|_{2, \Omega_T} + \|p_f^\varepsilon\|_{2, \Omega_T} + \|p_s^\varepsilon\|_{2, \Omega_T} \leq C_0,$$

THEOREM 2.2. *Let $\mathbf{F} \in L^2(\Omega) \times L^2(\Omega) \times L^2(\Omega)$, $\varepsilon > 0$, $\mathbf{w}_f^\varepsilon, \mathbf{w}_s^\varepsilon \in L^\infty(0, T; W_2^1(\Omega))$, and (1.1)–(1.9) hold.*

$$\mathbf{w}_f^\varepsilon = \mathbf{w}^\varepsilon \quad \text{in } \Omega_f^\varepsilon \times (0, T), \quad \mathbf{w}_s^\varepsilon = \mathbf{w}^\varepsilon \quad \text{in } \Omega_s^\varepsilon \times (0, T),$$

where $\mathbf{w}^\varepsilon = \mathbf{w}^\varepsilon(\mathbf{x}, t)$ is the unique solution of the problem (1.1)–(1.9) with $\{p_f^\varepsilon\}, \{q^\varepsilon\}, \{p_s^\varepsilon\}, \{\mathbf{w}^\varepsilon\}, \{\chi^\varepsilon \mathbf{w}^\varepsilon\}, \{(1 - \chi^\varepsilon) \mathbf{w}^\varepsilon\}, \{\mathbf{w}_f^\varepsilon\}, \{\mathbf{w}_s^\varepsilon\}$ replaced by $\{p_f, q, p_s, \mathbf{w}, \mathbf{w}^f, \mathbf{w}^s, \mathbf{w}_f, \mathbf{w}_s\}$ in $L^2(\Omega_T)$.

(I) $\mu_1 = \lambda_1 = \infty$, \dots , $\mathbf{w}_f = \mathbf{w}_s = \mathbf{w}$, \dots , Ω_T, \dots , $\mathbf{w}, p_f, q, \dots$, p_s, \dots

$$(2.7) \quad \hat{\rho} \frac{\partial^2 \mathbf{w}}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F},$$

$$(2.8) \quad \frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + \operatorname{div} \mathbf{w} = 0,$$

$$(2.9) \quad q = p_f + \frac{\nu_0}{p_*} \frac{\partial p_f}{\partial t}, \quad \frac{1}{m} q = \frac{1}{1-m} p_s,$$

$$(2.10) \quad \mathbf{w}(\mathbf{x}, 0) = \frac{\partial \mathbf{w}}{\partial t}(\mathbf{x}, 0) = 0, \quad \mathbf{x} \in \Omega,$$

$$(2.11) \quad \mathbf{w}(\mathbf{x}, t) \cdot \mathbf{n}(\mathbf{x}) = 0, \quad \mathbf{x} \in S, t > 0,$$

$$\hat{\rho} = m \rho_f + (1-m) \rho_s, \quad \dots, \quad m = \int_Y \chi dy,$$

(II) $\mu_1 = \infty$, $\lambda_1 < \infty$, \dots , Ω_T, \dots , $\mathbf{w}^f = m \mathbf{w}_f, \mathbf{w}^s, p_f, q, \dots$, p_s, \dots (2.9)

$$(2.12) \quad \rho_f m \frac{\partial^2 \mathbf{w}_f}{\partial t^2} + \rho_s \frac{\partial^2 \mathbf{w}^s}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F}$$

$$(2.13) \quad \frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + m \operatorname{div} \mathbf{w}_f + \operatorname{div} \mathbf{w}^s = 0,$$

$$(2.14) \quad \frac{\partial \mathbf{w}^s}{\partial t} = (1-m) \frac{\partial \mathbf{w}_f}{\partial t} + \int_0^t B_1^s(t-\tau) \cdot \mathbf{z}^s(\mathbf{x}, \tau) d\tau,$$

$$\mathbf{z}^s(\mathbf{x}, t) = -\frac{1}{m} \nabla q(\mathbf{x}, t) + \rho_s \mathbf{F}(\mathbf{x}, t) - \rho_s \frac{\partial^2 \mathbf{w}_f}{\partial t^2}(\mathbf{x}, t),$$

$\lambda_1 > 0$, \dots

$$(2.15) \quad \rho_s \frac{\partial^2 \mathbf{w}^s}{\partial t^2} = \rho_s B_2^s \cdot \frac{\partial^2 \mathbf{w}_f}{\partial t^2} + ((1-m)I - B_2^s) \cdot \left(-\frac{1}{m} \nabla q + \rho_s \mathbf{F} \right)$$

$\lambda_1 = 0$, \dots (2.9), (2.12)–(2.15),

(2.10),

(2.11),

$$\mathbf{w} = m \mathbf{w}_f + \mathbf{w}^s$$

(2.14)–(2.15) $B_1^s(t) = B_2^s = f_1$ (5.37)
 (5.39) $((1 - m)I - B_2^s) \cdot f_1$

(III) $\mu_1 < \infty, \lambda_1 = \infty, \Omega_T, \mathbf{w}^f, \mathbf{w}^s = (1 - m)\mathbf{w}_s$
 $p_f, q = p_s$
 (2.9)

(2.16) $\rho_f \frac{\partial^2 \mathbf{w}^f}{\partial t^2} + \rho_s(1 - m) \frac{\partial^2 \mathbf{w}_s}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F}$

(2.17) $\frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + \operatorname{div} \mathbf{w}^f + (1 - m) \operatorname{div} \mathbf{w}_s = 0,$

(2.18) $\frac{\partial \mathbf{w}^f}{\partial t} = m \frac{\partial \mathbf{w}_s}{\partial t} + \int_0^t B_1^f(t - \tau) \cdot \mathbf{z}^f(\mathbf{x}, \tau) d\tau,$

$\mathbf{z}^f(\mathbf{x}, t) = -\frac{1}{m} \nabla q(\mathbf{x}, t) + \rho_f \mathbf{F}(\mathbf{x}, t) - \rho_f \frac{\partial^2 \mathbf{w}_s}{\partial t^2}(\mathbf{x}, t),$

$\mu_1 > 0$
 (2.19) $\rho_f \frac{\partial^2 \mathbf{w}^f}{\partial t^2} = \rho_f B_2^f \cdot \frac{\partial^2 \mathbf{w}_s}{\partial t^2} + (mI - B_2^f) \cdot \left(-\frac{1}{m} \nabla q + \rho_f \mathbf{F} \right)$

$\mu_1 = 0$ (2.9), (2.16)–(2.19)
 (2.10)
 (2.11)

$\mathbf{w} = \mathbf{w}^f + (1 - m)\mathbf{w}_s$
 (2.18)–(2.19) $B_1^f(t) = B_2^f$ (5.44)–
 (5.45) $(mI - B_2^f) \cdot f_1$

(IV) $\mu_1 < \infty, \lambda_1 < \infty, \Omega_T, \mathbf{w}, p_f, q = p_s$
 (2.8)
 (2.9)

(2.20) $\frac{\partial \mathbf{w}}{\partial t} = \int_0^t B(t - \tau) \cdot \nabla q(\mathbf{x}, \tau) d\tau + \mathbf{f}(\mathbf{x}, t),$

$B(t) = \mathbf{f}(\mathbf{x}, t)$ (5.58) (5.59)
 (2.8), (2.9), (2.20)
 (2.10) (2.11)

As was mentioned above, even in the most simple case (I) with $\nu_0 = 0$, Theorem 2.2 gives the standard wave equation for the solid pressure p_s but with a completely new speed of sound in the mixture, which includes the porosity, densities, and speeds of sound in the solid and liquid components.

In the next simple case (IV) with $\nu_0 = 0$, Theorem 2.2 gives a new wave equation for the solid pressure in the form

(2.21) $\frac{\partial p_s}{\partial t} = \int_0^t \operatorname{div}(\tilde{B}(t - \tau) \cdot \nabla p_s(\mathbf{x}, \tau)) d\tau.$

Here $\tilde{B}(0) = c^2I$, where the time derivative of the matrix $\tilde{B}(t)$ is generally unbounded at $t = 0$. This equation has no simple solutions like traveling waves and requires a special analysis even for the smooth matrix $\tilde{B}(t)$.

The rest of the homogenized models described by Theorem 2.2 are much more complicated than the model (2.21). This is natural, because one cannot expect that a simple model gives an “accurate” approximation of the very complicated original model (1.1)–(1.9).

3. Preliminaries.

3.1. Two-scale convergence. The justification of Theorem 2.2 is based on a systematic use of the two-scale convergence method, which was proposed by Nguetseng [7] and has been recently used in a wide range of homogenization problems (see, for example, the survey [5]).

DEFINITION 3.1. Let $\{w^\varepsilon\} \subset L^2(\Omega_T)$ and $W \in L^2(\Omega_T \times Y)$ be a sequence of functions and a function, respectively, such that

$$(3.1) \quad \lim_{\varepsilon \searrow 0} \int_{\Omega_T} w^\varepsilon(\mathbf{x}, t) \sigma\left(\mathbf{x}, t, \frac{\mathbf{x}}{\varepsilon}\right) d\mathbf{x}dt = \int_{\Omega_T} \int_Y W(\mathbf{x}, t, \mathbf{y}) \sigma(\mathbf{x}, t, \mathbf{y}) d\mathbf{y}d\mathbf{x}dt$$

where $\sigma = \sigma(\mathbf{x}, t, \mathbf{y}) \in C^\infty(\Omega_T \times Y)$, $\mathbf{y} \in Y$, $f_i \in L^2(\Omega_T)$.

The existence and the main properties of weakly convergent sequences are established by the following fundamental theorem [7, 5].

THEOREM 3.2 (Nguetseng’s theorem). 1. Let $\{w^\varepsilon\} \subset L^2(\Omega_T)$ and $W \in L^2(\Omega_T \times Y)$ be a sequence of functions and a function, respectively, such that

$$2. \quad \lim_{\varepsilon \searrow 0} \int_{\Omega_T} \{w^\varepsilon\} \{\varepsilon \nabla_x w^\varepsilon\} = \int_{\Omega_T} \int_Y W = W(\mathbf{x}, t, \mathbf{y}), \quad \mathbf{y} \in Y, \quad \{w^\varepsilon\} \subset L^2(\Omega_T), \quad \{\varepsilon \nabla_x w^\varepsilon\} \subset L^2(\Omega_T \times Y), \quad \{w^\varepsilon\} \subset L^2(\Omega_T), \quad \{\varepsilon \nabla_x w^\varepsilon\} \subset L^2(\Omega_T \times Y), \quad W = \nabla_y W.$$

COROLLARY 3.3. Let $\sigma \in L^2(Y)$, $\sigma^\varepsilon(\mathbf{x}) = \sigma(\mathbf{x}/\varepsilon)$, $\{w^\varepsilon\} \subset L^2(\Omega_T)$, $W \in L^2(\Omega_T \times Y)$, $\{\sigma^\varepsilon w^\varepsilon\} \subset L^2(\Omega_T)$, σW .

3.2. An extension lemma. A typical difficulty in homogenization problems like problem (1.1)–(1.7) arises in passing to the limit as $\varepsilon \searrow 0$ because of the fact that the bounds on the displacement gradient $\nabla \mathbf{w}^\varepsilon$ may be different in the liquid and solid components. The classical approach to overcoming this difficulty consists in constructing an extension of the displacement field defined merely on Ω_s or Ω_f to the whole Ω . The following lemma is valid due to the well-known results from [1, 4, 8]. We formulate it in the form convenient for us.

LEMMA 3.4. Let $\chi^\varepsilon \in C^\infty(\Omega)$, $\chi^\varepsilon = 1$ in Ω_f , $\chi^\varepsilon = 0$ in Ω_s , $\mathbf{w}^\varepsilon \in W_2^1(\Omega)$, $\mathbf{w}_f^\varepsilon, \mathbf{w}_s^\varepsilon \in W_2^1(\Omega)$, $\mathbf{w}_f^\varepsilon = \mathbf{w}^\varepsilon$ in Ω_f , $\mathbf{w}_s^\varepsilon = \mathbf{w}^\varepsilon$ in Ω_s .

$$(3.2) \quad \chi^\varepsilon(\mathbf{x})(\mathbf{w}_f^\varepsilon(\mathbf{x}) - \mathbf{w}^\varepsilon(\mathbf{x})) = 0, \quad (1 - \chi^\varepsilon(\mathbf{x}))(\mathbf{w}_s^\varepsilon(\mathbf{x}) - \mathbf{w}^\varepsilon(\mathbf{x})) = 0, \quad \mathbf{x} \in \Omega,$$

$$(3.3) \quad \|\mathbf{w}_i^\varepsilon\|_{2, \Omega} \leq C \|\mathbf{w}^\varepsilon\|_{2, \Omega_i}, \quad \|D(x, \mathbf{w}_i^\varepsilon)\|_{2, \Omega} \leq C \|D(x, \mathbf{w}^\varepsilon)\|_{2, \Omega_i}, \quad i = f, s,$$

where $C = C(\chi^\varepsilon, Y)$, $\chi^\varepsilon \in C^\infty(\Omega)$, $\chi^\varepsilon = 1$ in Ω_f , $\chi^\varepsilon = 0$ in Ω_s .

3.3. Some notation. Further we denote the following:

(1)

$$\langle \Phi \rangle_Y = \int_Y \Phi dy, \quad \langle \Phi \rangle_{Y_f} = \int_Y \chi \Phi dy, \quad \langle \Phi \rangle_{Y_s} = \int_Y (1 - \chi) \Phi dy.$$

(2) If \mathbf{a} and \mathbf{b} are two vectors, then the matrix $\mathbf{a} \otimes \mathbf{b}$ is defined by the formula

$$(\mathbf{a} \otimes \mathbf{b}) \cdot \mathbf{c} = \mathbf{a}(\mathbf{b} \cdot \mathbf{c})$$

for any vector \mathbf{c} .

4. Proof of Theorem 2.1. Estimates (2.5)–(2.6) follow from the energy equality in the form

$$\begin{aligned} & \frac{d}{dt} \left\{ \int_{\Omega} \rho^\varepsilon \left(\frac{\partial \mathbf{w}^\varepsilon}{\partial t} \right)^2 + \alpha_\lambda \int_{\Omega} (1 - \chi^\varepsilon) D(x, \mathbf{w}^\varepsilon) : D(x, \mathbf{w}^\varepsilon) dx \right. \\ & \left. + \alpha_p \int_{\Omega} \chi^\varepsilon (\operatorname{div} \mathbf{w}^\varepsilon)^2 dx + \alpha_\eta \int_{\Omega} (1 - \chi^\varepsilon) (\operatorname{div} \mathbf{w}^\varepsilon)^2 dx \right\} + \alpha_\nu \int_{\Omega} \chi^\varepsilon \left(\operatorname{div} \frac{\partial \mathbf{w}^\varepsilon}{\partial t} \right)^2 dx \\ (4.1) \quad & + \alpha_\mu \int_{\Omega} \chi^\varepsilon D \left(x, \frac{\partial \mathbf{w}^\varepsilon}{\partial t} \right) : D \left(x, \frac{\partial \mathbf{w}^\varepsilon}{\partial t} \right) dx = \int_{\Omega} \rho^\varepsilon \frac{\partial \mathbf{F}}{\partial t} \cdot \frac{\partial \mathbf{w}^\varepsilon}{\partial t} dx \end{aligned}$$

if we use Hölder, Gronwall, and Korn inequalities and extension Lemma 3.4. In turn, the energy equality (4.1) follows from (1.1) if we express the stress tensor P and the liquid pressure p_f using state equations (1.2)–(1.4) and continuity equation (1.5), multiply the result by $\partial \mathbf{w}^\varepsilon / \partial t$, and integrate by parts. Note that all terms on the “solid skeleton–pore space” interface Γ^ε disappear due to boundary conditions (1.6)–(1.7).

The same estimates (2.5)–(2.6) guarantee the existence and uniqueness of the generalized solution for problem (1.1)–(1.9).

5. Proof of Theorem 2.2.

5.1. Weak and two-scale limits of sequences of displacements and pressures. First, we use Lemma 3.4 and conclude that there are functions $\mathbf{w}_f^\varepsilon, \mathbf{w}_s^\varepsilon \in L^\infty(0, T; W_2^1(\Omega))$ such that

$$\mathbf{w}_f^\varepsilon = \mathbf{w}^\varepsilon \text{ in } \Omega_f^\varepsilon \times (0, T), \quad \mathbf{w}_s^\varepsilon = \mathbf{w}^\varepsilon \text{ in } \Omega_s^\varepsilon \times (0, T).$$

By Theorem 2.1, the sequences $\{p_f^\varepsilon\}, \{q^\varepsilon\}, \{p_s^\varepsilon\}, \{\mathbf{w}^\varepsilon\}, \{\mathbf{w}_f^\varepsilon\}, \{\sqrt{\alpha_\mu} \nabla \mathbf{w}_f^\varepsilon\}, \{\mathbf{w}_s^\varepsilon\}$, and $\{\sqrt{\alpha_\lambda} \nabla \mathbf{w}_s^\varepsilon\}$ are bounded in $L^2(\Omega_T)$. Hence there exists a subsequence of small parameters $\{\varepsilon > 0\}$ and functions $p_f, q, p_s, \mathbf{w}, \mathbf{w}_f$, and \mathbf{w}_s such that

$$(5.1) \quad p_f^\varepsilon \rightharpoonup p_f, \quad q^\varepsilon \rightharpoonup q, \quad p_s^\varepsilon \rightharpoonup p_s, \quad \mathbf{w}^\varepsilon \rightharpoonup \mathbf{w}, \quad \mathbf{w}_f^\varepsilon \rightharpoonup \mathbf{w}_f, \quad \mathbf{w}_s^\varepsilon \rightharpoonup \mathbf{w}_s$$

weakly in $L^2(\Omega_T)$ as $\varepsilon \searrow 0$.

Note also that

$$(5.2) \quad (1 - \chi^\varepsilon) \alpha_\lambda D(x, \mathbf{w}_s^\varepsilon) \rightarrow 0, \quad \chi^\varepsilon \alpha_\mu D(x, \mathbf{w}_f^\varepsilon) \rightarrow 0$$

strongly in $L^2(\Omega_T)$ as $\varepsilon \searrow 0$.

Relabeling if necessary, we assume that the sequences themselves converge.

By Nguetseng’s theorem, there exist functions $P_f(\mathbf{x}, t, \mathbf{y})$, $P_s(\mathbf{x}, t, \mathbf{y})$, $Q(\mathbf{x}, t, \mathbf{y})$, $\mathbf{W}(\mathbf{x}, t, \mathbf{y})$, $\mathbf{W}_f(\mathbf{x}, t, \mathbf{y})$, and $\mathbf{W}_s(\mathbf{x}, t, \mathbf{y})$ that are one-periodic in \mathbf{y} and satisfy the condition that the sequences $\{p_f^\varepsilon\}$, $\{p_s^\varepsilon\}$, $\{q^\varepsilon\}$, $\{\mathbf{w}^\varepsilon\}$, $\{\mathbf{w}_f^\varepsilon\}$, and $\{\mathbf{w}_s^\varepsilon\}$ two-scale converge to $P_f(\mathbf{x}, t, \mathbf{y})$, $P_s(\mathbf{x}, t, \mathbf{y})$, $Q(\mathbf{x}, t, \mathbf{y})$, $\mathbf{W}(\mathbf{x}, t, \mathbf{y})$, $\mathbf{W}_f(\mathbf{x}, t, \mathbf{y})$, and $\mathbf{W}_s(\mathbf{x}, t, \mathbf{y})$, respectively.

LEMMA 5.1. $\mu_1 = \infty$ ($\lambda_1 = \infty$). $\mathbf{W}_f(\mathbf{x}, t, \mathbf{y}) = \mathbf{w}_f(\mathbf{x}, t)$, $\chi(\mathbf{y})\mathbf{W}(\mathbf{x}, t, \mathbf{y}) = \chi(\mathbf{y})\mathbf{w}_f(\mathbf{x}, t)$, $\mathbf{w}^f = \langle \mathbf{W} \rangle_{Y_f} = m\mathbf{w}_f$ ($\mathbf{W}_s(\mathbf{x}, t, \mathbf{y}) = \mathbf{w}_s(\mathbf{x}, t)$), $(1 - \chi(\mathbf{y}))\mathbf{W}(\mathbf{x}, t, \mathbf{y}) = (1 - \chi(\mathbf{y}))\mathbf{w}_s(\mathbf{x}, t)$, $\mathbf{w}^s = \langle \mathbf{W} \rangle_{Y_s} = (1 - m)\mathbf{w}_s$

Suppose that $\mu_1 = \infty$, and let $\Psi(\mathbf{x}, t, \mathbf{y})$ be an arbitrary smooth scalar function periodic in \mathbf{y} . The sequence $\{\sigma_{ij}^\varepsilon\}$, where

$$\sigma_{ij}^\varepsilon = \int_{\Omega} \sqrt{\alpha_\lambda} \frac{\partial w_{f,i}^\varepsilon}{\partial x_j}(\mathbf{x}, t) \Psi(\mathbf{x}, t, \mathbf{x}/\varepsilon) dx, \quad \mathbf{w}_f^\varepsilon = (w_{f,1}^\varepsilon, w_{f,2}^\varepsilon, w_{f,3}^\varepsilon),$$

is uniformly bounded in ε . Therefore,

$$\int_{\Omega} \varepsilon \frac{\partial w_{f,i}^\varepsilon}{\partial x_j}(\mathbf{x}, t) \Psi(\mathbf{x}, t, \mathbf{x}/\varepsilon) dx = \frac{\varepsilon}{\sqrt{\alpha_\lambda}} \sigma_{ij}^\varepsilon \rightarrow 0$$

as $\varepsilon \searrow 0$, which is equivalent to

$$\int_{\Omega} \int_Y W_{f,i}(\mathbf{x}, t, \mathbf{y}) \frac{\partial \Psi}{\partial y_j}(\mathbf{x}, t, \mathbf{y}) dx dy = 0, \quad \mathbf{W}_f = (W_{f,1}, W_{f,2}, W_{f,3}),$$

or $\mathbf{W}_f(\mathbf{x}, t, \mathbf{y}) = \mathbf{w}_f(\mathbf{x}, t)$. Therefore taking the two-scale limit as $\varepsilon \searrow 0$ in the relation $\chi^\varepsilon(\mathbf{w}^\varepsilon - \mathbf{w}_f^\varepsilon) = 0$, we arrive at

$$\chi(\mathbf{y})\mathbf{W}(\mathbf{x}, t, \mathbf{y}) = \chi(\mathbf{y})\mathbf{w}_f. \quad \square$$

5.2. Micro- and macroscopic equations. We start the proof of the theorem from the macro- and microscopic equations related to the continuity equations.

LEMMA 5.2. $\mathbf{x} \in \Omega$, $\mathbf{y} \in Y$. $\{p_f^\varepsilon\}$, $\{p_s^\varepsilon\}$, $\{q^\varepsilon\}$, $\{\mathbf{w}^\varepsilon\}$, $\{\mathbf{w}_f^\varepsilon\}$, $\{\mathbf{w}_s^\varepsilon\}$

$$(5.3) \quad Q = q\chi/m, \quad P_f = p_f\chi/m, \quad P_s = p_s(1 - \chi)/(1 - m), \quad Q = P + \nu_0 p_*^{-1} \partial P / \partial t;$$

$$(5.4) \quad q/m = p_s/(1 - m), \quad q = p_f + \nu_0 p_*^{-1} \partial p_f / \partial t;$$

$$(5.5) \quad p_f/p_* + p_s/\eta_0 + \operatorname{div} \mathbf{w} = 0;$$

$$(5.6) \quad \mathbf{w}(\mathbf{x}, t) \cdot \mathbf{n}(\mathbf{x}) = 0, \quad \mathbf{x} \in S, t > 0;$$

$$(5.7) \quad \operatorname{div}_y \mathbf{W} = 0;$$

$$(5.8) \quad \mathbf{W} = \chi \mathbf{W}_f + (1 - \chi) \mathbf{W}_s.$$

In order to prove (5.3), into (2.4) we substitute the test function $\psi^\varepsilon = \varepsilon \psi(\mathbf{x}, t, \mathbf{x}/\varepsilon)$, where $\psi(\mathbf{x}, t, \mathbf{y})$ is an arbitrary one-periodic function of \mathbf{y} that is finite on Y_f (or finite on Y_s or finite on Y). Passing to the limit as $\varepsilon \searrow 0$, we obtain

$$(5.9) \quad \nabla_y Q = 0, \quad \mathbf{y} \in Y_f; \quad \nabla_y P_s = 0, \quad \mathbf{y} \in Y_s; \quad \nabla_y (Q + P_s) = 0, \quad \mathbf{y} \in Y.$$

Next, fulfilling the two-scale passage to the limit in the state equation (2.3) and in the relations

$$(1 - \chi^\varepsilon)q^\varepsilon = 0, \quad \chi^\varepsilon p_s^\varepsilon = 0$$

we arrive at the last equation in (5.3) and the relations

$$(1 - \chi)Q = 0, \quad \chi P_s = 0,$$

which together with the first two equations in (5.9) prove the first three equations in (5.3).

The second equation in (5.4) is the result of integration of the last equation in (5.3) over the domain Y_f .

The first relation in (5.4) follows from (5.3) and the last equation in (5.9): the sequence $\{(q^\varepsilon + p_s^\varepsilon)\}$ two-scale converges to $(Q + P_s) = (q + p_s)$.

Equations (5.5)–(5.7) appear as a result of the two-scale passage to the limit in (1.5) and (2.2) with the proper test functions being involved. Thus, for example, (5.5) and (5.6) arise if we consider the linear combination of (1.5) and (2.2)

$$(5.10) \quad \frac{1}{\alpha_p} p_f^\varepsilon + \frac{1}{\alpha_\eta} p_s^\varepsilon + \operatorname{div} \mathbf{w}^\varepsilon = 0,$$

multiply it by an arbitrary function independent of the “fast” variable \mathbf{x}/ε , and then pass to the limit as $\varepsilon \searrow 0$. To prove (5.7), it suffices to consider the two-scale limiting relations in (5.10) as $\varepsilon \searrow 0$ with the test functions $\varepsilon \psi(\mathbf{x}/\varepsilon) h(\mathbf{x}, t)$, where ψ and h are arbitrary smooth functions.

To prove (5.8), it suffices to consider the two-scale limiting relations in

$$\mathbf{w}^\varepsilon = \chi^\varepsilon \mathbf{w}_f^\varepsilon + (1 - \chi^\varepsilon) \mathbf{w}_s^\varepsilon. \quad \square$$

LEMMA 5.3. *Let $(\mathbf{x}, t) \in \Omega_T$. Then*

$$(5.11) \quad \rho_f \frac{\partial^2 \mathbf{w}^f}{\partial t^2} + \rho_s \frac{\partial^2 \mathbf{w}^s}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F}$$

Substituting a test function of the form $\psi = \psi(\mathbf{x}, t)$ into integral identity (2.4) and passing to the limit as $\varepsilon \searrow 0$, we arrive at (5.11). \square

LEMMA 5.4. *Let $\mu_1 = \infty$, $\lambda_1 < \infty$. Then $\mathbf{W}^s = (1 - \chi) \mathbf{W}^s$, $\mathbf{w}_f = q$, and*

$$(5.12) \quad \rho_s \frac{\partial^2 \mathbf{W}^s}{\partial t^2} = \lambda_1 \Delta_y \mathbf{W}^s - \nabla_y R^s - \frac{1}{m} \nabla q + \rho_s \mathbf{F}, \quad \mathbf{y} \in Y_s,$$

$$(5.13) \quad \mathbf{W}^s = \mathbf{w}_f, \quad \mathbf{y} \in \gamma,$$

Let $\lambda_1 > 0$. Then

$$(5.14) \quad \rho_s \frac{\partial^2 \mathbf{W}^s}{\partial t^2} = -\nabla_y R^s - \frac{1}{m} \nabla q + \rho_s \mathbf{F}, \quad \mathbf{y} \in Y_s,$$

$$(5.15) \quad (\mathbf{W}^s - \mathbf{w}_f) \cdot \mathbf{n} = 0, \quad \mathbf{y} \in \gamma,$$

Let $\lambda_1 = 0$.

$$(5.16) \quad \mathbf{W}^s(\mathbf{y}, 0) = \frac{\partial \mathbf{W}^s}{\partial t}(\mathbf{y}, 0) = 0, \quad \mathbf{y} \in Y_s.$$

(5.15) $\mathbf{n} \cdot \nabla_y \mathbf{W} = \gamma$. The differential equations (5.12) and (5.14) follow as $\varepsilon \searrow 0$ from integral equality (2.4) with the test function $\psi = \varphi(x\varepsilon^{-1}) \cdot h(\mathbf{x}, t)$, where φ is solenoidal and finite in Y_s .

The boundary condition (5.13) is a consequence of the two-scale convergence of the sequence $\{\sqrt{\alpha\lambda}\nabla_x \mathbf{w}^\varepsilon\}$ to the function $\sqrt{\lambda_1}\nabla_y \mathbf{W}(\mathbf{x}, t, \mathbf{y})$. By this convergence, the function $\nabla_y \mathbf{W}(\mathbf{x}, t, \mathbf{y})$ is L^2 -integrable in Y . The boundary condition (5.15) follows from (5.7)–(5.8) and the relation $\mathbf{W}_f = \mathbf{w}_f$. \square

In the same way, one can prove the following lemma.

LEMMA 5.5. *Let $\mu_1 < \infty$, $\lambda_1 = \infty$. Then, for $\mathbf{y} \in Y_f$, $\mathbf{W}^f = \chi \mathbf{W}$, $\mathbf{w}_s = q$.*

$$(5.17) \quad \rho_f \frac{\partial^2 \mathbf{W}^f}{\partial t^2} = \mu_1 \Delta_y \frac{\partial \mathbf{W}^f}{\partial t} - \nabla_y R^f - \frac{1}{m} \nabla q + \rho_f \mathbf{F}, \quad \mathbf{y} \in Y_f,$$

$$(5.18) \quad \mathbf{W}^f = \mathbf{w}_s, \quad \mathbf{y} \in \gamma,$$

$\mu_1 > 0$

$$(5.19) \quad \rho_f \frac{\partial^2 \mathbf{W}^f}{\partial t^2} = -\nabla_y R^f - \frac{1}{m} \nabla q + \rho_f \mathbf{F}, \quad \mathbf{y} \in Y_f,$$

$$(5.20) \quad (\mathbf{W}^f - \mathbf{w}_s) \cdot \mathbf{n} = 0, \quad \mathbf{y} \in \gamma,$$

$\mu_1 = 0$

$$(5.21) \quad \mathbf{W}^f(\mathbf{y}, 0) = \frac{\partial \mathbf{W}^f}{\partial t}(\mathbf{y}, 0) = 0, \quad \mathbf{y} \in Y_f.$$

LEMMA 5.6. *Let $\mu_1 < \infty$, $\lambda_1 < \infty$, $\tilde{\rho} = \rho_f \chi + \rho_s(1 - \chi)$. Then, for $\mathbf{y} \in Y$, $\mathbf{W} = q$.*

$$(5.22) \quad \left. \begin{aligned} &\tilde{\rho} \partial^2 \mathbf{W} / \partial t^2 + 1/m \nabla q - \tilde{\rho} \mathbf{F} \\ &= \operatorname{div}_y \{ \mu_1 \chi D(y, \partial \mathbf{W} / \partial t) + \lambda_1 (1 - \chi) D(y, \mathbf{W}) - RI \} \end{aligned} \right\}$$

$$(5.23) \quad \mathbf{W}(\mathbf{y}, 0) = \frac{\partial \mathbf{W}}{\partial t}(\mathbf{y}, 0) = 0, \quad \mathbf{y} \in Y.$$

In the proof of the last lemma, we additionally use Nguetseng’s theorem, which states that the sequence $\{\varepsilon D(x, \mathbf{w}^\varepsilon)\}$ two-scale converges to the function $D(y, \mathbf{W})$.

5.3. Homogenized equations. Lemmas 5.2 and 5.3 imply the following lemma.

LEMMA 5.7. *Let $\mu_1 = \lambda_1 = \infty$. Then, $\mathbf{w}_f = \mathbf{w}_s = \mathbf{w}$, $\Omega_T = \Omega$, $\mathbf{w} = p_f q + p_s$.*

$$(5.24) \quad \hat{\rho} \frac{\partial^2 \mathbf{w}}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F},$$

$$(5.25) \quad \frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + \operatorname{div} \mathbf{w} = 0,$$

$$(5.26) \quad q = p_f + \frac{\nu_0}{p_*} \frac{\partial p_f}{\partial t}, \quad \frac{1}{m} q = \frac{1}{1-m} p_s,$$

$$(5.27) \quad \mathbf{w}(\mathbf{x}, 0) = \frac{\partial \mathbf{w}}{\partial t}(\mathbf{x}, 0) = 0, \quad \mathbf{x} \in \Omega,$$

$$(5.28) \quad \mathbf{w}(\mathbf{x}, t) \cdot \mathbf{n}(\mathbf{x}) = 0, \quad \mathbf{x} \in S, t > 0.$$

LEMMA 5.8. *Let $\mu_1 = \infty$, $\lambda_1 < \infty$ and let Ω_T be a bounded domain. Let $\mathbf{w}_f, \mathbf{w}^s, p_f, q, p_s$ be functions satisfying (5.26)*

$$(5.29) \quad \rho_f m \frac{\partial^2 \mathbf{w}_f}{\partial t^2} + \rho_s \frac{\partial^2 \mathbf{w}^s}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F},$$

$$(5.30) \quad \frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + m \operatorname{div} \mathbf{w}_f + \operatorname{div} \mathbf{w}^s = 0,$$

$$(5.31) \quad \frac{\partial \mathbf{w}^s}{\partial t} = (1-m) \frac{\partial \mathbf{w}_f}{\partial t} + \int_0^t B_1^s(t-\tau) \cdot \mathbf{z}^s(\mathbf{x}, \tau) d\tau,$$

$$\mathbf{z}^s(\mathbf{x}, t) = -\frac{1}{m} \nabla q(\mathbf{x}, t) + \rho_s \mathbf{F}(\mathbf{x}, t) - \rho_s \frac{\partial^2 \mathbf{w}_f}{\partial t^2}(\mathbf{x}, t),$$

Let $\lambda_1 > 0$, then the system (5.26)–(5.31) has a unique solution.

$$(5.32) \quad \rho_s \frac{\partial^2 \mathbf{w}^s}{\partial t^2} = \rho_s B_2^s \cdot \frac{\partial^2 \mathbf{w}_f}{\partial t^2} + ((1-m)I - B_2^s) \cdot \left(-\frac{1}{m} \nabla q + \rho_s \mathbf{F} \right)$$

Let $\lambda_1 = 0$, then the system (5.26), (5.29)–(5.32) has a unique solution if (5.27) and (5.28) are satisfied.

$$(5.33) \quad \frac{\partial \mathbf{w}^s}{\partial t} = B_1^s(t) \cdot \frac{\partial \mathbf{w}_f}{\partial t} + B_2^s(t) \cdot \mathbf{f}_1 + ((1-m)I - B_2^s) \cdot \mathbf{f}_1. \quad (5.37)$$

Equation (5.29) follows directly from (5.11). The continuity equation (5.30) follows from (5.5) if we take into account that

$$\mathbf{w} = m\mathbf{w}_f + \mathbf{w}^s.$$

To find the last two equations (5.31) and (5.32), we just have to solve the system of microscopic equations (5.7), (5.12)–(5.16) and use the formula

$$\mathbf{w}^s = \langle \mathbf{W} \rangle_{Y_s}.$$

There are two different cases.

(a) If $\lambda_1 > 0$, then the solution of the system of microscopic equations (5.7), (5.12), and (5.13) supplemented with the homogeneous initial data (5.16) is given by the formulas

$$\mathbf{W}^s = \int_0^t \left(\mathbf{v}(\mathbf{x}, \tau) + \sum_{i=1}^3 \mathbf{W}^{s,i}(\mathbf{y}, t - \tau) z_i^s(\mathbf{x}, \tau) \right) d\tau,$$

$$R^s = \int_0^t \sum_{i=1}^3 R^{s,i}(\mathbf{y}, t - \tau) z_i^s(\mathbf{x}, \tau) d\tau, \quad \mathbf{z}^s = (z_1^s, z_2^s, z_3^s),$$

and the functions $\mathbf{W}^{s,i}(\mathbf{y}, t)$ and $R^{s,i}(\mathbf{y}, t)$ are defined by virtue of the periodic initial boundary value problem

$$(5.33) \quad \rho_s \frac{\partial^2 \mathbf{W}^{s,i}}{\partial t^2} - \lambda_1 \Delta \mathbf{W}^{s,i} + \nabla R^{s,i} = 0, \quad \mathbf{y} \in Y_s, t > 0,$$

$$(5.34) \quad \operatorname{div}_y \mathbf{W}^{s,i} = 0, \quad \mathbf{y} \in Y_s, t > 0,$$

$$(5.35) \quad \mathbf{W}^{s,i} = 0, \quad \mathbf{y} \in \gamma, t > 0,$$

$$(5.36) \quad \mathbf{W}^{s,i}(\mathbf{y}, 0) = 0, \quad \rho_s \frac{\partial \mathbf{W}^{s,i}}{\partial t}(\mathbf{y}, 0) = \mathbf{e}_i, \quad \mathbf{y} \in Y_s.$$

In (5.36), \mathbf{e}_i is the standard Cartesian basis vector.

Therefore,

$$(5.37) \quad B_1^s(t) = \left\langle \sum_{i=1}^3 \frac{\partial \mathbf{W}^{s,i}}{\partial t} \right\rangle_{Y_s} (t) \otimes \mathbf{e}_i.$$

Note that (5.33) is understood in the sense of distributions and the function $B_1^s(t)$ has no time derivative at $t = 0$.

(b) If $\lambda_1 = 0$, then in solving the system (5.7), (5.14), (5.15), and (5.16), we first find the pressure $R^s(\mathbf{x}, t, \mathbf{y})$ by solving the Neumann problem for the Laplace equation in Y_s in the form

$$R^s(\mathbf{x}, t, \mathbf{y}) = \sum_{i=1}^3 R_{s,i}(\mathbf{y}) z_i^s(\mathbf{x}, t),$$

where $R_{s,i}(\mathbf{y})$ is the solution of the problem

$$(5.38) \quad \Delta_y R_{s,i} = 0, \quad \mathbf{y} \in Y_s; \quad \nabla_y R_{s,i} \cdot \mathbf{n} = \mathbf{n} \cdot \mathbf{e}_i, \quad \mathbf{y} \in \gamma; \quad \langle R_{s,i} \rangle_{Y_s} = 0.$$

Formula (5.32) is the result of integration of (5.14) over the domain Y_s and

$$(5.39) \quad B_2^s = \sum_{i=1}^3 \langle \nabla R_{s,i} \rangle_{Y_s} \otimes \mathbf{e}_i,$$

where the matrix $B = ((1 - m)I - B_2^s)$ is symmetric and strictly positive definite. In fact, let $\tilde{R} = \sum_{i=1}^3 R_{s,i} \xi_i$ for any unit vector $\xi = (\xi_1, \xi_2, \xi_3)$. Then

$$(B \cdot \xi) \cdot \xi = \langle (\xi - \nabla \tilde{R})^2 \rangle_{Y_f}$$

and $(B \cdot \xi) \cdot \xi = 0$ if and only if \tilde{R} is a linear function in \mathbf{y} . On the other hand, it follows from the assumption about the geometry of the domain Y_s that all linear periodic functions on Y_s are constant. Finally, the normalization condition $\langle R_{s,i} \rangle_{Y_s} = 0$ yields that $\tilde{R} = 0$. However, this is impossible, because the functions $R_{s,i}$ are linearly independent. \square

LEMMA 5.9. *Let $\mu_1 < \infty$, $\lambda_1 = \infty$, $\Omega_T \subset \mathbb{R}^3$, $\mathbf{w}^f, \mathbf{w}_s, p_f, q, p_s$ satisfy (5.26)*

$$(5.40) \quad \rho_f \frac{\partial^2 \mathbf{w}^f}{\partial t^2} + \rho_s (1 - m) \frac{\partial^2 \mathbf{w}_s}{\partial t^2} = -\frac{1}{m} \nabla q + \hat{\rho} \mathbf{F},$$

$$(5.41) \quad \frac{1}{p_*} p_f + \frac{1}{\eta_0} p_s + \operatorname{div} \mathbf{w}^f + (1 - m) \operatorname{div} \mathbf{w}_s = 0,$$

$$(5.42) \quad \frac{\partial \mathbf{w}^f}{\partial t} = m \frac{\partial \mathbf{w}_s}{\partial t} + \int_0^t B_1^f(t - \tau) \cdot \mathbf{z}^f(\mathbf{x}, \tau) d\tau,$$

$$\mathbf{z}^f(\mathbf{x}, t) = -\frac{1}{m} \nabla q(\mathbf{x}, t) + \rho_f \mathbf{F}(\mathbf{x}, t) - \rho_f \frac{\partial^2 \mathbf{w}_s}{\partial t^2}(\mathbf{x}, t),$$

for $\mu_1 > 0$, then the following equations hold:

$$(5.43) \quad \rho_f \frac{\partial^2 \mathbf{w}^f}{\partial t^2} = \rho_f B_2^f \cdot \frac{\partial^2 \mathbf{w}_s}{\partial t^2} + (mI - B_2^f) \cdot \left(-\frac{1}{m} \nabla q + \rho_f \mathbf{F} \right)$$

for $\mu_1 = 0$, (5.26), (5.40)–(5.43), (5.27), (5.28), $\mathbf{w} = \mathbf{w}^f + (1 - m)\mathbf{w}_s$, (5.42)–(5.43), $B_1^f(t) = B_2^f$, (5.44)–(5.45), $(mI - B_2^f) \cdot \mathbf{f}_1$.

The proof of this lemma repeats that of the previous lemma. Here we have to solve the system of microscopic equations (5.7), (5.17)–(5.21) and use the formula

$$\mathbf{w}^f = \langle \mathbf{W} \rangle_{Y_f}.$$

Thus,

$$(5.44) \quad B_1^f(t) = \left\langle \sum_{i=1}^3 \frac{\partial \mathbf{W}^{f,i}}{\partial t} \right\rangle_{Y_f}(t) \otimes \mathbf{e}_i,$$

$$(5.45) \quad B_2^f = \sum_{i=1}^3 \langle \nabla R_{f,i} \rangle_{Y_f} \otimes \mathbf{e}_i,$$

where the functions $\mathbf{W}^{f,i}(\mathbf{y}, t)$ solve the periodic initial boundary value problem

$$(5.46) \quad \rho_f \frac{\partial^2 \mathbf{W}^{f,i}}{\partial t^2} - \mu_1 \Delta \frac{\partial \mathbf{W}^{f,i}}{\partial t} + \nabla R_{f,i} = 0, \quad \mathbf{y} \in Y_f, t > 0,$$

$$(5.47) \quad \operatorname{div}_{\mathbf{y}} \mathbf{W}^{f,i} = 0, \quad \mathbf{y} \in Y_f, t > 0,$$

$$(5.48) \quad \mathbf{W}^{f,i} = 0, \quad \mathbf{y} \in \gamma, t > 0,$$

$$(5.49) \quad \mathbf{W}^{f,i}(\mathbf{y}, 0) = 0, \quad \rho_f \frac{\partial \mathbf{W}^{f,i}}{\partial t}(\mathbf{y}, 0) = \mathbf{e}_i, \quad \mathbf{y} \in Y_f,$$

and the functions $R_{f,i}(\mathbf{y})$ solve the periodic boundary value problem

$$(5.50) \quad \Delta_{\mathbf{y}} R_{f,i} = 0, \quad \mathbf{y} \in Y_f; \quad \nabla_{\mathbf{y}} R_{f,i} \cdot \mathbf{n} = \mathbf{n} \cdot \mathbf{e}_i, \quad \mathbf{y} \in \gamma; \quad \langle R_{f,i} \rangle_{Y_f} = 0.$$

Note that, as before, the matrix $(mI - B_2^f)$ is symmetric and strictly positive definite. \square

The proof of Theorem 2.2 is completed by the following lemma.

LEMMA 5.10. $\mu_1 < \infty, \lambda_1 < \infty, \Omega_T, \mathbf{w}, p_f, q,$

$$(5.25) \quad (5.26) \quad \dots$$

$$(5.51) \quad \frac{\partial \mathbf{w}}{\partial t} = \int_0^t B(t - \tau) \cdot \nabla q(\mathbf{x}, \tau) d\tau + \mathbf{f}(\mathbf{x}, t),$$

$$B(t) \cdot \mathbf{f}(\mathbf{x}, t) \quad (5.58) \quad (5.59)$$

$$(5.25), (5.26), (5.51) \quad (5.27) \quad (5.28)$$

To derive the momentum conservation law (5.51), we must solve the system of microscopic equations (5.7), (5.22) with the initial conditions (5.23) and use the formula

$$\mathbf{w} = \langle \mathbf{W} \rangle_Y.$$

Let

$$\mathbf{W} = \int_0^t \sum_{i=1}^3 \left\{ \mathbf{W}^{q,i}(\mathbf{y}, t - \tau) \frac{\partial q}{\partial x_i}(\mathbf{x}, \tau) + \mathbf{W}^{F,i}(\mathbf{y}, t - \tau) F_i(\mathbf{x}, \tau) \right\} d\tau,$$

$$R = \int_0^t \sum_{i=1}^3 \left\{ R^{q,i}(\mathbf{y}, t - \tau) \frac{\partial q}{\partial x_i}(\mathbf{x}, \tau) + R^{F,i}(\mathbf{y}, t - \tau) F_i(\mathbf{x}, \tau) \right\} d\tau,$$

where $\mathbf{F} = \sum_{i=1}^3 F_i \mathbf{e}_i$.

Then the pair $\{\mathbf{W}, R\}$ is a solution of system (5.7), (5.22) and (5.23) if and only if the functions $\{\mathbf{W}^{q,i}(\mathbf{y}, t), R^{q,i}(\mathbf{y}, t)\}$ and $\{\mathbf{W}^{F,i}(\mathbf{y}, t), R^{F,i}(\mathbf{y}, t)\}$ are periodic in \mathbf{y} solutions of the equations

$$(5.52) \quad \operatorname{div}_y \left\{ \mu_1 \chi D \left(y, \frac{\partial \mathbf{W}^{q,i}}{\partial t} \right) + \lambda_1 (1 - \chi) D(y, \mathbf{W}^{q,i}) - R^{q,i} I \right\} = \tilde{\rho} \frac{\partial^2 \mathbf{W}^{q,i}}{\partial t^2},$$

$$(5.53) \quad \operatorname{div}_y \mathbf{W}^{q,i} = 0,$$

$$(5.54) \quad \operatorname{div}_y \left\{ \mu_1 \chi D \left(y, \frac{\partial \mathbf{W}^{F,i}}{\partial t} \right) + \lambda_1 (1 - \chi) D(y, \mathbf{W}^{F,i}) - R^{F,i} I \right\} = \tilde{\rho} \frac{\partial^2 \mathbf{W}^{F,i}}{\partial t^2},$$

$$(5.55) \quad \operatorname{div}_y \mathbf{W}^{F,i} = 0$$

in the domain Y for $t > 0$ and satisfy the initial conditions

$$(5.56) \quad \mathbf{W}^{q,i}(\mathbf{y}, 0) = 0, \quad \tilde{\rho} \frac{\partial \mathbf{W}^{q,i}}{\partial t}(\mathbf{y}, 0) = -\frac{1}{m} \mathbf{e}_i, \quad \mathbf{x} \in Y,$$

$$(5.57) \quad \mathbf{W}^{F,i}(\mathbf{y}, 0) = 0, \quad \frac{\partial \mathbf{W}^{F,i}}{\partial t}(\mathbf{y}, 0) = \mathbf{e}_i, \quad \mathbf{x} \in Y.$$

Here \mathbf{e}_i is the standard Cartesian basis vector.

Therefore,

$$(5.58) \quad B(t) = \sum_{i=1}^3 \left\langle \frac{\partial \mathbf{W}^{q,i}}{\partial t}(\mathbf{y}, t) \right\rangle_Y \otimes \mathbf{e}_i,$$

$$(5.59) \quad \mathbf{f}(\mathbf{x}, t) = \int_0^t \sum_{i=1}^3 \left\langle \frac{\partial \mathbf{W}^{F,i}}{\partial t}(\mathbf{y}, t - \tau) \right\rangle_Y F_i(\mathbf{x}, \tau) d\tau.$$

The solvability and uniqueness of problems (5.52), (5.53), (5.56) and (5.54), (5.55), (5.57) follow directly from the energy identity

$$\begin{aligned} & \frac{1}{2} \int_Y \left(\tilde{\rho} \left(\frac{\partial \mathbf{W}^{j,i}}{\partial t}(\mathbf{y}, t) \right)^2 + \lambda_1 D(y, \mathbf{W}^{j,i}(\mathbf{y}, t)) : D(y, \mathbf{W}^{j,i}(\mathbf{y}, t)) \right) dy \\ & + \int_0^t \int_Y \mu_1 D \left(y, \frac{\partial \mathbf{W}^{j,i}}{\partial \tau}(\mathbf{y}, \tau) \right) : D \left(y, \frac{\partial \mathbf{W}^{j,i}}{\partial \tau}(\mathbf{y}, \tau) \right) dy d\tau = \frac{1}{2} \beta^j \end{aligned}$$

for $i = 1, 2, 3$ and $j = q, F$.

Here

$$\beta^q = \left\langle \frac{1}{\tilde{\rho}} \right\rangle_Y, \quad \beta^F = \langle \tilde{\rho} \rangle_Y.$$

As before, equations (5.51) are understood in the sense of distributions and the function $B(t)$ has no time derivative at $t = 0$. \square

Acknowledgment. The author acknowledges the anonymous referees for kind advice and for helpful remarks.

REFERENCES

- [1] E. ACERBI, V. CHIADO PIAT, G. DAL MASO, AND D. PERCIVALE, *An extension theorem from connected sets and homogenization in general periodic domains*, *Nonlinear Anal.*, 18 (1992), pp. 481–496.
- [2] R. BURRIDGE AND J. B. KELLER, *Poroelasticity equations derived from microstructure*, *J. Acoust. Soc. Amer.*, 70 (1981), pp. 1140–1146.
- [3] T. CLOPEAU, J. L. FERRIN, R. P. GILBERT, AND A. MIKELIĆ, *Homogenizing the acoustic properties of the seabed: Part II*, *Math. Comput. Modelling*, 33 (2001), pp. 821–841.
- [4] V. V. JIKOV, S. M. KOZLOV, AND O. A. OLEINIK, *Homogenization of Differential Operators and Integral Functionals*, Springer-Verlag, New York, 1994.
- [5] D. LUKKASSEN, G. NGUETSENG, AND P. WALL, *Two-scale convergence*, *Int. J. Pure Appl. Math.*, 2 (2002), pp. 35–86.
- [6] A. MEIRMANOV, *Nguetseng’s two-scale convergence method for filtration and seismic acoustic problems in elastic porous media*, *Siberian Math. J.*, 48 (2007), pp. 519–538.
- [7] G. NGUETSENG, *A general convergence result for a functional related to the theory of homogenization*, *SIAM J. Math. Anal.*, 20 (1989), pp. 608–623.
- [8] G. NGUETSENG, *Asymptotic analysis for a stiff variational problem arising in mechanics*, *SIAM J. Math. Anal.*, 21 (1990), pp. 1394–1414.
- [9] E. SANCHEZ-PALENCIA, *Nonhomogeneous Media and Vibration Theory*, Lecture Notes in Phys. 129, Springer-Verlag, Berlin, 1980.

BOUNDARY BEHAVIOR OF SOLUTIONS OF A CLASS OF GENUINELY NONLINEAR HYPERBOLIC SYSTEMS*

JULIAN GEVIRTZ†

Abstract. We study the set of boundary singularities of arbitrary classical solutions of genuinely nonlinear 2×2 planar hyperbolic systems of the form $D_k R_k = 0$, where D_k denotes differentiation in the direction $e^{i\theta_k(R_1, R_2)}$, $k = 1, 2$, and where the defining functions θ_k satisfy (i) $\theta_2 = \theta_1 + \frac{\pi}{2}$, and (ii) $A \leq \left| \frac{\partial \theta_k(R)}{\partial R_k} \right| \leq B$, for all $R \in \mathbb{R}^2$, $k = 1, 2$, for some positive constants A, B . We show that for any system of this kind there is a $\tau < 1$ such that for any locally Lipschitz solution R in a smoothly bounded domain G , the set of points of ∂G at which R fails to have a nontangential limit has Hausdorff dimension at most τ , and, on the other hand, for any such system for which the $\theta_k \in C^\infty(\mathbb{R}^2)$, we construct a C^∞ solution R on a half-plane \mathbb{H} for which the set of points of $\partial \mathbb{H}$ at which R fails to have a nontangential limit has positive Hausdorff dimension. These results are immediately applicable to constant principal strain mappings, which are defined in terms of a system of this kind for which θ_1 is a linear function of R_1 and R_2 .

Key words. boundary behavior, constant principal strain mapping, genuinely nonlinear hyperbolic system, Hausdorff dimension, nontangential limit

AMS subject classifications. 35L40, 28A78

DOI. 10.1137/070705507

1. Introduction. For hyperbolic systems in two independent variables x and t , most often associated with space and time, one usually studies the Cauchy problem, in which one seeks a solution $u(x, t)$, $t \geq 0$, for which $u(x, 0)$ coincides with a given $u_0(x)$, the questions considered including well-posedness, global existence, blow-up, and behavior of solutions as $t \rightarrow \infty$. In the nonlinear case, discussion is often limited to initial data with a small range, and, even for such data, generalized solutions must be considered.

In this paper we concern ourselves with the following inverse question for a certain family of genuinely nonlinear 2×2 hyperbolic systems: What can be said about the boundary values of an \dots , classical solution in a plane domain G ? Here “classical” can be taken to mean C^∞ , although the treatment we give will be valid for locally Lipschitz solutions. In the first place, we are interested in systems whose formulation imposes no a priori limit on the range of characteristic directions, that is, systems such that for a characteristic given parametrically by $z(s)$, $\arg\{z'(s)\}$ can potentially cover all of \mathbb{R} , in contrast to what is implicitly the case in the standard space-time context. Second, we are interested in statements valid for \dots , classical solutions rather than ones known to arise from some form of initial value problem. Because of this generality, even in geometrically simple domains such as disks or half-planes characteristics can be quite contorted curves. Although the specific focus of this paper is the size of the set of boundary points at which classical solutions can fail to have nontangential limits, it would be reasonable to investigate other aspects of their behavior and that of the associated characteristics. In any event, given the nonstandard nature of the boundary value question and of several of the issues that arise in dealing with it, we shall begin with a somewhat detailed discussion of a system

*Received by the editors October 16, 2007; accepted for publication March 31, 2008; published electronically November 5, 2008.

<http://www.siam.org/journals/sima/40-4/70550.html>

†2005 North Winthrop Road, Muncie, IN 47304-2536 (jgevirtz@gmail.com).

for which it is physically meaningful, namely the system which describes smooth planar mappings with constant principal stretches (cps-mappings), about which we have previously written [ChG], [G1], [G2], [G3], [G4], [G5]. It is in fact the study of the boundary behavior of such mappings that is the main goal of this paper, and we have chosen to work in a wider context only because it is possible to do so with little additional effort, and because this broader approach suggests some interesting questions.

A mapping $f : G \rightarrow \mathbb{C}$ with principal stretch factors $m_1 \neq m_2$ is a mapping $f : G \rightarrow \mathbb{C}$ with locally Lipschitz continuous Jacobian

$$J_f = T(-\phi)\sigma(m_1, m_2)T(\theta),$$

where

$$T(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \quad \text{and} \quad \sigma(m_1, m_2) = \begin{bmatrix} m_1 & 0 \\ 0 & m_2 \end{bmatrix}.$$

As is explained in the cited references, apart from regularity considerations, functions θ and ϕ will correspond to such a mapping on a simply connected domain G if and only if they satisfy the autonomous quasi-linear hyperbolic system

$$(1.1) \quad D_1(m_1\theta - m_2\phi) = 0, \quad D_2(m_2\theta - m_1\phi) = 0,$$

where

$$D_1u = (\cos \theta)u_x + (\sin \theta)u_y \quad \text{and} \quad D_2u = (-\sin \theta)u_x + (\cos \theta)u_y.$$

The characteristics of a solution are the integral curves of the fields $e^{i\theta}$ and $ie^{i\theta}$. It turns out that a net \mathcal{N} made up of two mutually orthogonal families of curves covering a simply connected G is the net of characteristics of a cps-mapping if and only if for any two curves C_1, C_2 belonging to one of the families of \mathcal{N} the change in the inclination of the tangent is the same along all subarcs of curves of the family which join C_1 to C_2 . Nets with this property are known as Hencky–Prandtl (HP) nets. (See [CS], [G3], [Hem], [Hen], [Hi], [Pr].) The theory of cps-mappings we have developed is based on direct application of (1.1) together with this Hencky–Prandtl (HP) property and the equations

$$(1.2) \quad D_2D_1\theta = [D_1\theta]^2 \quad \text{and} \quad D_1D_2\theta = -[D_2\theta]^2,$$

which are also effectively equivalent to (1.1) and which are very special cases of equations derived by Lax [L] in the context of considerably more general genuinely non-linear 2×2 hyperbolic systems in the plane and used by him in connection with the inevitability of singularity formation. The equations (1.2) imply that if a characteristic C has curvature κ_0 at p , then the orthogonal characteristic arc emanating from p towards the concave side of C can have length at most $\frac{1}{\kappa_0}$, that is, that the boundary of G must be encountered after moving at most a distance of $\frac{1}{\kappa_0}$ along this orthogonal characteristic. A property of this kind is a fundamental property for the theory we are developing and plays a fundamental role in all that is to follow.

When regarded as deformations with constant principal strains, cps-mappings are of concrete interest as models in a number of physically interesting contexts (see [Y]).

Consider, for example, a thin liquid film on a plane surface which upon solidification takes on a rectangular cryptocrystalline structure; that is, at each point a suitably oriented minute square of the original liquid becomes a rectangular crystal whose side lengths are constant multiples of the side length of the square. In this light global geometric results for cps-mappings tell one about the extent to which the shape of the original film can change as a result of such solidification and about how matter is moved around in the process, and statements about the existence of boundary limits of θ (and, in light of (1.1), of ϕ , and consequently of the Jacobian of the mapping) tell one to what extent the cryptocrystalline structure is present at the very edge of the solidified lamina. Applied to the system (1.1) the main result of this paper says that $\tau < 1$ implies $G \subset \partial G$ for $\theta = \phi$. On the other hand, the construction of section 5 shows that $\tau > 1$ implies $\theta \neq \phi$.

Beyond their immediate physical significance, cps-mappings constitute a particularly important and tractable class of planar quasi-isometric mappings, for which we believe they will ultimately be shown to display extremal behavior for many of the as yet unsolved distortion questions (see [J1], [J2]). In this direction a very significant inroad was made a few years ago by Gutlyanskii and Martio [GM], who showed that the mapping given by

$$g(re^{i\psi}) = re^{i(\psi + \psi_0 \frac{\log r}{\log \rho})}$$

are extremal for the problem of determining for given $\rho > 1$ and $\psi_0 > 0$ the smallest ratio $\frac{m_1}{m_2} > 1$, such that there is a quasi-isometric mapping f of the annulus $1 < |z| < \rho$ onto itself with stretching bounds m_1, m_2 which satisfies the boundary conditions

$$f(z) = z \quad \text{and} \quad f(\rho z) = \rho e^{i\psi_0} z \quad \text{for} \quad |z| = 1.$$

It is in fact not hard to see that g is indeed a cps-mapping of the entire punctured plane $\mathbb{C} \setminus \{0\}$ with

$$m_1 = \frac{\sqrt{a^2 + 4} + a}{2} \quad \text{and} \quad m_2 = \frac{\sqrt{a^2 + 4} - a}{2},$$

where $a = \frac{\psi_0}{\log \rho}$. We call these spiral mappings because the corresponding inclination functions are of the form

$$(1.3) \quad \theta(re^{i\psi}) = \psi + \alpha,$$

where $\tan \alpha = m_2$, which is to say that the characteristics form two mutually orthogonal families of logarithmic spirals, all members of each of which are rotations of each other.

We now describe the class of 2×2 systems for which we treat the boundary limit question. Let $\theta_k = \theta_k(R_1, R_2)$, $k = 1, 2$. For given $R_1(x, y)$, $R_2(x, y)$ and any $u = u(x, y)$ we write

$$(1.4) \quad D_k u = \cos \theta_k(R_1(x, y), R_2(x, y)) \frac{\partial u}{\partial x} + \sin \theta_k(R_1(x, y), R_2(x, y)) \frac{\partial u}{\partial y}.$$

It is well known that in general an autonomous 2×2 quasi-linear homogeneous hyperbolic system for unknown functions f and g is formally equivalent to a system of

the form

$$D_k R_k = 0, \quad k = 1, 2,$$

with appropriate inclination functions $\theta_k(R_1, R_2)$. The relationship between the functions R_1, R_2 and f, g is of the form $R_k(x, y) = F_k(f(x, y), g(x, y))$. Henceforth we write $\Theta = (\theta_1, \theta_2)$ and use the term *normal system* to mean a C^∞ mapping $\Theta : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, although the smoothness requirement could be weakened substantially. Obviously, a solution of the system Θ in a domain G of the plane is then a pair of functions $R_1(x, y), R_2(x, y)$ for which R_k is constant on each integral curve (henceforth referred to as a *k-characteristic*) of the field $e^{i\theta_k(R_1(x, y), R_2(x, y))}$, $k = 1, 2$. A system is said to be *hyperbolic* if the derivatives $\frac{\partial \theta_k}{\partial R_k}$, $k = 1, 2$, never vanish. It is clear that the system (1.1) for the θ and ϕ associated with cps-mappings is already in Riemann invariant form with $R_i = m_i \theta - m_j \phi$, so that in this case the two inclination functions are given by

$$(1.5) \quad \theta_1 = \theta = \frac{m_1 R_1 - m_2 R_2}{m_1^2 - m_2^2} \quad \text{and} \quad \theta_2 = \theta + \frac{\pi}{2}.$$

(Here we have used the convention, in force throughout this paper, to the effect that $\{i, j\} = \{1, 2\}$.) This system is obviously genuinely nonlinear and in fact is the simplest possible such system in that the two families of characteristics are mutually orthogonal and Θ is a linear function of $R = (R_1, R_2)$. We now define the family of systems with which we work.

DEFINITION 1.1. *Normal system.* $\Theta = (\theta_1, \theta_2)$

(i) $\theta_2 = \theta_1 + \frac{\pi}{2}$

(ii) $A, B > 0, \dots, A \leq \left| \frac{\partial \theta_k}{\partial R_k} \right| \leq B, \dots, R \in \mathbb{R}^2$

The only hyperbolic systems with which we deal will be normal systems, for which we use the symbol θ to denote θ_1 . We shall show that

$$\tau = \tau(\Theta) < 1, \dots, G \subset \mathbb{C}$$

$$\Theta : G \rightarrow \mathbb{R}^2, \dots, \partial G \rightarrow \mathbb{R}^2$$

iff \dots, τ . This will follow as an immediate consequence (Corollary 4.3) of our principal result, Main Theorem 2.2, which deals with boundary singularities of a class of functions effectively more general than the class of solutions of normal systems. An outline of the proof, which is actually made simpler by this somewhat greater generality, is given early in section 2, just after the statement of the main theorem. Furthermore, in section 5 we shall show that

$$C^\infty \dots, \mathbb{H}, \dots, \partial \mathbb{H}$$

$$R \dots, \mathbb{H}$$

2. Quasi-HP functions, the characteristic length bound, and related matters.

Let $G \subset \mathbb{C}$ be a domain. If θ is a locally Lipschitz function on G , the integral curves of the fields $e^{i\theta(z)}$ and $ie^{i\theta(z)}$ will be called the 1- and 2-characteristics of θ , respectively. The term *characteristic* will refer to the complete integral curve, and we will use the term *arc* to refer to either of the two arcs into which a full characteristic is divided by one of its points. (Characteristics which are closed curves will not arise in this paper since we are dealing with single-valued functions θ .) As indicated in the introduction we will use the convention $\{i, j\} = \{1, 2\}$ throughout. Arcs of k -characteristics will be called *k-arcs*, or, less specifically, *characteristic arcs*. With reference to a given θ , a characteristic arc joining points $a, b \in D$ will be denoted by ab and we shall use the abbreviation

$$\Delta\theta(ab) = \theta(b) - \theta(a).$$

A domain $Q \subset G$ will be said to be a K -quasi-HP function θ , if ∂Q is a Jordan curve lying in G containing four points a, b, d, c which occur in that order when ∂Q is traversed in the positive (negative) sense and such that ab and cd are i -arcs and ac and bd are j -arcs. We say that ab and cd are K -quasi-HP of each other with respect to or along any j -characteristic passing through ab . For a curve parameterized by $z = z(s), \alpha < s < \beta$, we use the terms “to the right of C ” and “to the left of C ” in the obvious sense, so that, for example, if C is a characteristic arc and $p \in C$, we describe an orthogonal characteristic arc or a half-characteristic as emanating from p to the right or left of C .

We shall denote the 2-dimensional measure of $X \subset \mathbb{C}$ by $\mu(X)$ and the 1-dimensional measure of a set A by $\lambda(A)$. The parameter s will always refer to arc length. We use the notation $N(a, r) = \{z : |z - a| < r\}$ and denote the line segment joining a to b by \overline{ab} . The overline will also be used to denote closure, but this should cause no confusion. For $X, Y \subset \mathbb{C}, \text{dist}(X, Y) = \inf\{|y - x| : x \in X, y \in Y\}$, and for $a \in \mathbb{C}, \text{dist}(a, X) = \text{dist}(\{a\}, X)$.

DEFINITION 2.1. Let $K \geq 1$ and θ be a K -quasi-HP function in G . We say that θ has the K -quasi-HP property if for any $abcd \subset G$ as above,

$$(2.1) \quad \frac{1}{K} |\Delta\theta(ac)| \leq |\Delta\theta(bd)| \leq K |\Delta\theta(ac)|.$$

A K -quasi-HP net is a net Θ in G such that $K(\Theta) \leq K$.

A simple continuity argument shows that this definition implies that the $\Delta\theta(ac)$ and $\Delta\theta(bd)$ in (2.1) must in fact have the same sign (unless both vanish). It is also obvious that a 1-quasi-HP net is an HP net. Note that while the HP property is a local condition that implies its global counterpart, this is not the case for the K -quasi-HP property when $K > 1$.

LEMMA 2.2. Let Θ be a K -quasi-HP net, where $K = K(\Theta)$. Indeed, if A and B are as in Definition 1.1, then it is clear that we can take $K(\Theta) = B/A$.

We can now state our main theorem.

MAIN THEOREM 2.2. Let $\tau = \tau(K) < 1$ and θ be a K -quasi-HP function in $G \subset \mathbb{C}^2$. Let ∂G be a C^2 curve. Then the Hausdorff dimension of the set of k -singularities of θ is at most τ .

Because the proof of this theorem, to be given in section 4, is quite involved and depends on the prior development of a considerable amount of machinery in this and the following section, we shall briefly explain here how it proceeds. As we shall show (see Proposition 3.29) for any point $p \in \partial G$ at which θ does not have a nontangential limit either there is a nontrivial fan of characteristics emanating from p or, for $k = 1$ or 2, every neighborhood of p completely contains a full k -characteristic of θ ; in the latter case we say that p is a k -singularity. Since, as will be apparent, any quasi-HP function can have at most a countable number of fans (see Proposition 3.24), we need only show that the Hausdorff dimension, $\dim(S)$, of the set S of k -singularities satisfies $\dim(S) \leq \tau < 1$. For this to be the case it is enough that there be some $\delta = \delta(K) > 0$ such that any almost straight arc $A \subset \partial G$ has a subarc of length at least $\delta\lambda(A)$ which has at most a countable set of k -singularities of θ . If for some θ there were no such δ , then for any N there would have to be an almost straight arc $A \subset \partial G$ such that there is an N -element set S_N of k -singularities of θ which are essentially uniformly distributed along A . For sufficiently large N , starting with a

set of N small k -characteristics, one very close to each of the points of S_N , we show that there must be a k -characteristic C (where the term “ k -characteristic” is used here in an appropriate sense—see the discussion of extended characteristics between Propositions 3.13 and 3.14) whose endpoints lie on A and are at least $\delta'\lambda(A)$ apart (where $\delta' > 0$ depends solely on K) and which is tangent to A at a point $m \in A$, where m is appropriately bounded away from the endpoints of C . We then use this to obtain a subarc A' of A which contains m , whose length is bounded below by $\delta''\lambda(A)$, where δ'' , like δ' , depends only on K , and on which θ can have at most a countable number of k -singularities (see Proposition 3.33), thereby arriving at the desired contradiction. A good measure of the complexity of the proof lies in establishing the existence of the extended characteristic C , which is carried out in section 4, but which depends on properties of the net of extended characteristics of quasi-HP functions developed in section 3. We begin with the following proposition, which is immediate.

PROPOSITION 2.3 (invariance of K). Let $\theta \in K$ and let G be a domain in \mathbb{C} with $a, b \in \mathbb{C}$, $a \neq 0$. Then $\theta(az + b) - \arg a \in K$ and $\frac{1}{a}(G - b)$.

The next proposition is a special case of [G1, Lemma 2]; a somewhat simpler proof than the one given there is included for the sake of completeness. A function θ is said to be a locally L -Lipschitz function on G if each point of G has a neighborhood in which θ satisfies a Lipschitz condition with constant L .

PROPOSITION 2.4. Let G be a domain in \mathbb{C} and let $\theta \in K$ be a locally L -Lipschitz function on G . Let $Q = abcd \subset G$ be a quadrilateral with $|b - a| = l$, $|c - a| = \epsilon \leq l$, and $\text{dist}(Q, \partial G) = \eta > 0$. Then $\bar{l} = \bar{l}(L, \eta) > 0$.

$$(2.2) \quad |d - c| = l - \epsilon \Delta\theta(ab) + O(\epsilon^2 + l^3), \quad l \leq \bar{l},$$

where O is big- O of L and η . In what follows, when we say that some quantity is big- O of some expression, we mean that this is so for all l less than some positive number which depends only on L and η , and that the constant corresponding to the big- O depends only on L and η . Let ab and cd be i -arcs. Without loss of generality we can assume that $a = 0$ and $b = l$. Let $\frac{\pi}{2} + \alpha$ and $\frac{\pi}{2} + \beta$ be the inclinations of the tangents to the j -characteristics at a and b , respectively, and let $\frac{\pi}{2} + \alpha' = \arg\{c - a\}$, $\frac{\pi}{2} + \beta' = \arg\{d - b\}$. Clearly, α, β, α' , and β' are all $O(l)$ and $\beta - \alpha = \Delta\theta(ab)$. It easily follows from the Lipschitz condition that $\alpha' = \alpha + O(\epsilon)$. We have $d = \epsilon ie^{i\alpha'} + te^{i\gamma}$, where $t = |d - c| = O(l)$ and $\gamma = O(l)$. We also have $d = l + sie^{i\beta'}$, where $s = O(l)$. Just as $\alpha' = \alpha + O(\epsilon)$, one sees that $\beta' = \beta + O(s)$. From the two expressions for d we have that

$$\epsilon ie^{i\alpha'} + te^{i\gamma} = l + sie^{i\beta'},$$

so that considering real and imaginary parts we have

$$(2.3) \quad -\epsilon \sin \alpha' + t \cos \gamma = l - s \sin \beta'$$

and

$$\epsilon \cos \alpha' + t \sin \gamma = s \cos \beta'.$$

The latter equation implies that

$$s(1 + O(l^2)) = \epsilon(1 + O(l^2)) + O(l^2),$$

so that $s = \epsilon + O(l^2)$. Thus, since $\beta' = \beta + O(s)$, we have

$$(2.4) \quad \beta' = \beta + O(\epsilon + l^2).$$

From (2.3) it now follows that

$$\begin{aligned} t(1 + O(l^2)) &= l + \epsilon \sin \alpha' - s \sin \beta' \\ &= l + \epsilon(\alpha' + O(l^3)) - (\epsilon + O(l^2))(\beta' + O(l^3)), \end{aligned}$$

so that from the fact that $\alpha' = \alpha + O(\epsilon)$ and (2.4) it follows that

$$\begin{aligned} |d - c| &= t = l + \epsilon(\alpha' - \beta') + O(l^3) \\ &= l + \epsilon(\alpha + O(\epsilon)) - \epsilon(\beta + O(\epsilon + l^2)) + O(l^3) \\ &= l - \epsilon(\beta - \alpha) + O(\epsilon^2 + l^3). \end{aligned}$$

Since $\beta - \alpha = \Delta\theta(ab)$, we are done. \square

Henceforth we use the notation $Df(s)$ to denote $f'(s)$.

PROPOSITION 2.5 (length change estimate). *Let $\theta : K \rightarrow \mathbb{R}^n$ be a C^1 map, $G \subset \mathbb{R}^n$ a domain, $Q = abcd$ a quadrilateral in G , $z(s)$ a curve in G , $0 \leq s \leq \alpha$, $z(0) = a$, $z(\alpha) = b$, $E(s) = z(s)w(s)$ a curve in G , ca a curve in G , $z(s)$ a curve in G , $D\theta(z(s_0)) = \kappa_0$, $s_0 \in (0, \alpha)$, $\lambda(E(s_0)) = \lambda_0$, $\delta > 0$, $\tau > 0$, $s \in [0, \alpha]$, $|z(s) - z(s_0)| < \tau$.*

$$|w(s) - w(s_0)| = |z(s) - z(s_0)|(1 - P\lambda_0\kappa_0 + R),$$

$$\frac{1}{K+\delta} \leq P \leq K + \delta, \quad |R| \leq \delta$$

By replacing G with an appropriate subdomain which contains Q , we can assume that θ is L -Lipschitz on G for some L . Let

$$\psi_0 = \psi_0(s) = \theta(z(s)) - \theta(z(s_0)).$$

For notational convenience we assume that $s > s_0$, the opposite case being effectively the same. Let $u_0(\sigma)$, $0 \leq \sigma \leq \lambda_0$, be the arc length parameterization of $E(s_0)$ with $u_0(0) = z(s_0)$. For each $\sigma \in [0, \lambda_0]$, let $u(\sigma)$ be the point of $E(s)$ joined to $u_0(\sigma)$ by an i -arc. Let $M(s)$ and $m(s)$ denote the maximum and minimum of $|u(\sigma) - u_0(\sigma)|$ for $0 \leq \sigma \leq \lambda_0$. Obviously, $M(s) \rightarrow 0$ as $s \rightarrow 0$. Let $\bar{\tau}$ be so small that

$$(2.5) \quad M(\bar{\tau}) < \bar{l}(L, \text{dist}(Q, \partial G)),$$

where \bar{l} is as in Proposition 2.4. For $\mu \in (0, 1]$ let $\tau_0 \leq \bar{\tau}$ be such that

$$(2.6) \quad |\psi_0(s) - \kappa_0|z(s) - z(s_0)|| < \mu(|\kappa_0| + 1)|z(s) - z(s_0)| \quad \text{for } |s - s_0| < \tau_0.$$

We can define an increasing sequence of numbers in $[0, \lambda_0]$ as follows. Let $\sigma_0 = 0$, and let σ_{k+1} be the first number in $[\sigma_k, \lambda_0]$ (as long as there is one) for which $\epsilon_k = |u_0(\sigma_{k+1}) - u(\sigma_k)| = \mu l_k$, where $l_k = |u_0(\sigma_k) - u(\sigma_k)|$. By (2.6) we have

$$(2.7) \quad |\psi_0 - \kappa_0 l_0| < \mu(|\kappa_0| + 1)l_0.$$

Bearing in mind that by the K -quasi-HP property $\Delta\theta(u_0(\sigma_k)u(\sigma_k)) = P_k\psi_0$, where P_k is between $\frac{1}{K}$ and K , we may apply Proposition 2.4 (in light of (2.5)) to conclude that

$$(2.8) \quad l_{k+1} = l_k - \mu l_k P_k \psi_0 + O((\mu^2 + l_k)l_k^2).$$

We continue defining σ_{k+1} until we come to $k = N$, for which $|u_0(\sigma_N) - w(s_0)| < \mu l_N$. Such an N exists, and indeed $(N + 1)\mu m(s) \leq \lambda_0$. We denote by $T \geq 1$ the constant of the big- O in (2.8). Note that, by the preceding proposition, T depends only on L and $\text{dist}(Q, \partial G)$, and so is independent of the value of μ we ultimately decide to work with.

We now restrict μ to satisfy

$$(2.9) \quad 0 < \mu \leq \mu_1 = \frac{1}{T(1 + \Lambda_0)\Lambda_0} \leq 1,$$

where

$$\Lambda_0 = 1 + K(|\kappa_0| + 2)\lambda_0.$$

Note that μ_1 depends only on K , κ_0 , λ_0 , and T . We show that for $l_0 \leq \tau_1 = \min\{\tau_0, \mu^2\}$ we have $l_k \leq \Lambda_0 l_0$ for all k . Assume inductively that $l_j \leq \Lambda_0 l_0$, $0 \leq j \leq k$. We now use the facts that $P_k \leq K$, and $\sum \mu l_j = \sum \epsilon_j \leq \lambda_0$, and that $|\psi_0| \leq (|\kappa_0| + 1)l_0$ in light of (2.9) and (2.7). It follows from (2.8) that

$$\begin{aligned} l_{k+1} &= l_0 + \sum_{j=0}^k l_{j+1} - l_j = l_0 + \sum_{j=0}^k (-\mu l_j P_j \psi_0 + O((\mu^2 + l_j)l_j^2)) \\ &\leq l_0 + K|\psi_0|\lambda_0 + T \sum_{j=0}^k (\mu^2 + l_j)l_j^2 \leq l_0(1 + K(|\kappa_0| + 1)\lambda_0) + \sum_{j=0}^k T(\mu^2 + l_j)l_j^2. \end{aligned}$$

However,

$$(2.10) \quad \begin{aligned} \sum_{j=0}^k T(\mu^2 + l_j)l_j^2 &\leq T \sum_{j=0}^k (\mu^2 + \Lambda_0 l_0) \frac{\epsilon_j}{\mu} l_j \leq T \frac{(\mu^2 + \Lambda_0 l_0)\Lambda_0 l_0 \lambda_0}{\mu} \\ &\leq T \frac{(\mu^2 + \Lambda_0 \mu^2)\Lambda_0 l_0 \lambda_0}{\mu} \leq T\mu(1 + \Lambda_0)\Lambda_0 l_0 \lambda_0. \end{aligned}$$

But by (2.9) this last expression is at most $l_0 \lambda_0 \leq K l_0 \lambda_0$, so that

$$l_{k+1} \leq l_0(1 + K(|\kappa_0| + 1)\lambda_0) + K l_0 \lambda_0 = \Lambda_0 l_0,$$

so that indeed $l_k \leq \Lambda_0 l_0$, for all k and all $l_0 \leq \tau$, provided only that μ satisfies (2.9) and $l_0 \leq \tau_1$. Since $M(s) \rightarrow 0$ as $s \rightarrow 0$, there is a $\tau_2 = \tau_2(\mu) \leq \tau_1$ such that

$$(2.11) \quad (1 - \mu)\lambda_0 \leq \sum_{k=0}^{N-1} \epsilon_k \leq \lambda_0 \quad \text{for } l_0 \leq \tau_2.$$

Assume the sequence of σ_k stops at $k = N$. Then $|w(s_0) - u(\sigma_N)| < \mu l_N \leq \mu \Lambda_0 l_0$, and, since we are assuming $l_0 = |z(s) - z(s_0)| < \tau_2 \leq \mu^2$, Proposition 2.4 implies that

$$||w(s) - w(s_0)| - l_N| \leq \mu \Lambda_0 l_0 K \psi_0 + T(\mu^2 \Lambda_0^2 l_0^2 + \Lambda_0^3 l_0^3).$$

From this it is easy to see that there is some $\mu \leq \mu_1$, which depends only on K, κ_0, λ_0 , and T , such that

$$(2.12) \quad ||w(s) - w(s_0)| - l_N| \leq \mu l_0.$$

As previously,

$$l_N - l_0 = \sum_{k=0}^{N-1} l_{k+1} - l_k = - \sum_{k=0}^{N-1} (\mu l_k P_k \psi_0 + O((\mu^2 + l_k)l_k^2)).$$

By (2.10), the sum of the big- O terms of the immediately preceding displayed line is bounded above by $\Lambda_1 \mu l_0$, where $\Lambda_1 = T(1 + \Lambda_0)\Lambda_0$. For $\psi_0 \geq 0$ it follows from this together with (2.11) and (2.12) that

$$\begin{aligned} |w(s) - w(s_0)| - l_0 &= l_N - l_0 + (|w(s) - w(s_0)| - l_N) \\ &\geq - \sum_{k=0}^{N-1} K \mu l_k \psi_0 - (\Lambda_1 + 1)\mu l_0 \geq -K \psi_0 \lambda_0 - (\Lambda_1 + 1)\mu l_0. \end{aligned}$$

But since by (2.7), $\psi_0 \leq (\kappa_0 + (|\kappa_0| + 1)\mu)l_0$, it follows that $|w(s) - w(s_0)| - l_0 \geq Al_0$, where

$$A = -K \lambda_0 \kappa_0 - \mu(K \lambda_0 (|\kappa_0| + 1) + \Lambda_1 + 1).$$

Since (2.7) also implies that $\psi_0 \geq (\kappa_0 - (|\kappa_0| + 1)\mu)l_0$, in an analogous manner one sees that for $\psi_0 \geq 0$

$$\begin{aligned} |w(s) - w(s_0)| - l_0 &\leq - \sum_{k=0}^{N-1} \frac{\mu l_k \psi_0}{K} + (\Lambda_1 + 1)\mu l_0 \\ &\leq \left(-\frac{1}{K}(1 - \mu)\lambda_0(\kappa_0 - \mu(|\kappa_0| + 1)) + (\Lambda_1 + 1)\mu \right) l_0, \end{aligned}$$

so that $|w(s) - w(s_0)| - l_0 \leq Bl_0$, where

$$B = -\frac{1}{K}\lambda_0 \kappa_0 + \mu \left(\frac{1}{K}\lambda_0 \kappa_0 + \frac{(1 - \mu)}{K}(|\kappa_0| + 1) + \Lambda_1 + 1 \right).$$

By considering separately the cases $\kappa_0 = 0$ and $\kappa_0 \neq 0$ we obtain the desired conclusion when $\psi_0 \geq 0$ with an appropriate $\mu = \mu(K, \kappa_0, \lambda_0, T, \delta)$. The case $\psi_0 \leq 0$ is the same apart from straightforward reversal of inequalities. \square

Let ϕ be a real-valued function defined on (α, β) , and let $s_0 \in (\alpha, \beta)$. We denote by $D^-\phi(s_0)$ and $D^+\phi(s_0)$ the lower and upper limits of $\frac{\phi(s) - \phi(s_0)}{s - s_0}$ as $s \rightarrow s_0$. The remaining three propositions of this section all follow easily from Proposition 2.5 together with elementary measure theory.

PROPOSITION 2.6 (curvature-characteristic length bound). *Let $z = z(s)$, $\alpha < s < \beta$, i be an i -arc of length K , θ be a θ -arc of length J , G be a G -arc of length J , C be a C -arc of length J , ∂G be a ∂G -arc of length J . Then $D^+\theta(z(s_0)) \leq \frac{K}{\lambda(J)}$.*

In particular this says that if C is an i -arc whose curvature $\kappa(p)$ at p exists, then the j -half-characteristic emanating from p towards the j -side of C (that is, to the left or right of C according as $\kappa(p) > 0$ or $\kappa(p) < 0$) has length at most $\frac{K}{\kappa(p)}$.

PROPOSITION 2.7 (length monotonicity). Let $\theta \in \text{HP}(G, K)$, $z = z(s)$, $\alpha \leq s \leq \beta$, $z(\alpha) = a$, $z(\beta) = b$, E an arc of ∂G with $\lambda(E) > 0$. Let E' be the arc of ∂G obtained by moving E to the right of C by a distance ϵ . Then $\lambda(E') \geq \lambda(E)$.

In section 4 we shall make use of the following two lower bounds for the area of certain regions made up of families of characteristic arcs of a K -quasi-HP function θ ; the constants η and η' which appear in them depend solely on K .

PROPOSITION 2.8 (area bounds). Let C be a closed curve in \mathbb{C} with $\theta \in \text{HP}(G, K)$, $w \in C$, $C'(w)$ the characteristic curve starting at w and going to the right of C . Let $U = \cup\{C'(w) : w \in C\}$. Then

- (i) $\mu(U) \geq \eta \lambda(C) \lambda_j$
- (ii) $\mu(U) \geq \eta' \lambda_j^2 |\Delta \theta(C)|$

3. Extended characteristics and regular and singular boundary behavior. Our approach to boundary behavior requires the examination of curves which are in effect characteristics whose interiors (i.e., sets of nonendpoints) contain boundary points; the subtleties that arise in this connection require careful discussion. Hereafter the symbol \mathcal{G} will denote the family of all Jordan domains $G \subset \mathbb{C}$ for which ∂G is a C^2 curve, and $\mathcal{G}(\rho) \subset \mathcal{G}$ will denote the family of all $G \in \mathcal{G}$ such that for each $p \in \partial G$ the interior of one of the circles of radius ρ tangent to ∂G at p is contained in G and that of the other such circle is contained in $\mathbb{C} \setminus \overline{G}$. Obviously, for $G \in \mathcal{G}(\rho)$ the unsigned curvature of ∂G is everywhere bounded by $\frac{1}{\rho}$, and $\mathcal{G} = \cup\{\mathcal{G}(\rho) : \rho > 0\}$. Furthermore, $\text{HP}(G, K)$ will denote the family of K -quasi-HP functions on G . Although sometimes the hypotheses of the propositions of this section do not state so explicitly, they always deal with $G \in \mathcal{G}$ and $\theta \in \text{HP}(G, K)$. (Although the results of this paper apply to unbounded and multiply connected domains as well, the proof of the main theorem itself will entail only consideration of Jordan domains, and in fact we will be able to work largely with Jordan domains of the kind we call “characteristic subdomains,” as defined below.) By an arc C we shall henceforth mean a continuous one-to-one mapping $z = z(t)$ of a closed interval $[\alpha, \beta]$ into the closure \overline{G} of G , and $z((\alpha, \beta))$ will be referred to as the interior of C . As in section 2, when an arc is considered to be oriented, we use the term “to (towards) the right (left) of C ” to refer to the part of C (immediately) to the right (left) of C , and the term “characteristic curve” will refer to a complete integral curve of θ or $\theta + \frac{\pi}{2}$ in G . The following proposition was proved in [G5, Proposition 2.8] for HP nets. Although the proof is virtually the same in the present more general context we include it here for the sake of completeness.

PROPOSITION 3.1. Let C be an arc of ∂G with $\theta \in \text{HP}(G, K)$, $z = z(s)$, $s \in (\alpha, \beta)$, $\lim_{s \rightarrow \beta} z(s) = b \in \partial G$. Then $\text{dist}(z(s), \partial G) \rightarrow 0$ as $s \rightarrow \beta$. Clearly, the conclusion holds if $\beta \neq \infty$, so that we assume $\beta = \infty$. First of all, we show that $\text{dist}(z(s), \partial G) \rightarrow 0$ as $s \rightarrow \infty$. If this were not true, then there would be a $z_0 \in G$ and an $\epsilon > 0$ such that for some sequence $\{s_i\}$ tending to ∞ , $z(s_i) \rightarrow z_0$, but $z([s_i, s_{i+1}]) \cap \partial N(z_0, \epsilon) \neq \emptyset$. But from this it would follow that some orthogonal characteristic crosses C twice, an impossible occurrence in light of the simple connectivity of G . We can now show that, in fact, $z(s) \rightarrow b \in \partial G$ as $s \rightarrow \infty$. If this is not so, the foregoing then implies that there is an arc E of ∂G , $\lambda(E) > 0$, each point of which is an accumulation point of $C_\gamma = \{z(s) : s > \gamma\}$ for each $\gamma \in (\alpha, \infty)$. Since G is bounded and $\beta = \infty$, C cannot be a straight line, so that from the characteristic length bound it follows that there is an orthogonal

half-characteristic C' of finite length which joins some $z(\sigma)$ to a point $e \in \partial G$. Since C cannot cross C' twice in G , $C_\sigma \subset G \setminus C'$. Let z_1, z_2 be distinct points of $E \setminus \{e\}$. For each $\delta > 0$, C_σ has a subarc $pp' \subset N(\partial G, \delta) \setminus C'$, with $p, p' \in N(z_1, \delta)$ and a point $p'' \in pp' \cap N(z_2, \delta)$. For obvious topological reasons, for each sufficiently small δ , there must be a point q on pp' which is joined to a point in $N(z_1, \delta)$ by an orthogonal characteristic arc B of length at least $|z_1 - z_2| - 2\delta$ such that the curvature of C at q exists and tends to infinity as $\delta \rightarrow 0$ and C is concave towards the side from which B emanates. But this clearly violates the characteristic length bound (Proposition 2.6), as indicated in the sentence immediately following it. \square

DEFINITION 3.2. Let C be an elementary i -characteristic of G with endpoints $z(\alpha), z(\beta) \in \partial G$. Then C is called an elementary i -characteristic.

In other words, an elementary characteristic is a full characteristic together with its endpoints, which are well defined by the preceding proposition. Note that for each $p \in \partial G$, $\{p\}$ is an elementary i -characteristic. Subarcs of elementary i -characteristics will be called i -arcs (or simply i -arcs, for short).

LEMMA 3.3. Let $G \in \mathcal{G}$ and let $B' = B'(G) \in (0, 1]$. Let $a, b \in \partial G$ and let C_1, C_2 be elementary i -characteristics of G with endpoints a, b ($a = b, C_1 = \{a\}, C_2 = \partial G$).

$$\text{dist}(z, C_1) + \text{dist}(z, C_2) \geq B' \text{dist}(z, \{a, b\}) \quad z \in G.$$

This is self-evident. \square

PROPOSITION 3.4 (bounded length of characteristics). Let $G \in \mathcal{G}$ and let $M = M(G, K)$. Let $\theta \in \text{HP}(G, K)$.

Let θ be a K -quasi-HP function on G , and let C be an elementary i -characteristic of θ . We regard C as being oriented and let $A = C \cap \partial G$ (A , being the set of endpoints of C , has at most two points). Let $d_0 = \text{diam}(G)$, and let $d = \sup\{\text{dist}(z, A) : z \in C\} \leq d_0$. Clearly,

$$(3.1) \quad \mu(\{z \in G : \text{dist}(z, A) \leq 8r\}) \leq Br^2,$$

where $B = 128\pi$. For $k \geq 0$, let

$$G_k = \left\{ z \in G : \frac{d}{2^k} \leq \text{dist}(z, A) \leq \frac{d}{2^{k-1}} \right\}$$

and

$$C_k = C \cap G_k,$$

so that $C \cap G = \cup\{C_k : k \geq 1\}$. For each nonendpoint p of C , let $J(p)$ denote the elementary j -characteristic containing p . Obviously, $J(p) \cap J(p') \cap G = \emptyset$ for $p \neq p'$. For $k \geq 1$ and $p \in C_k$ let $l(p)$ and $r(p)$ be the first points encountered when moving along $J(p)$ from p to the left and right of C , respectively, which are not in the interior of $G_{k-1} \cup G_k \cup G_{k+1}$. Each $q \in \{l(p), r(p)\}$ is either on ∂G or is in G and satisfies one of $\text{dist}(q, A) = \frac{d}{2^{k+1}}$ or $\text{dist}(q, A) = \frac{d}{2^{k-2}}$. Say $q \notin \partial G$ and $\text{dist}(q, A) = \frac{d}{2^{k+1}}$. Then $|p - q| \geq \frac{d}{2^{k+1}}$, since otherwise $\text{dist}(p, A) \leq |p - q| + \text{dist}(q, A) < \frac{d}{2^k}$, which contradicts the fact that $p \in C_k$. If $q \notin \partial G$ and $\text{dist}(q, A) = \frac{d}{2^{k-2}}$, then $|p - q| \geq \frac{d}{2^{k+1}}$ since otherwise

$$\text{dist}(p, A) \geq \text{dist}(q, A) - |p - q| > \frac{d}{2^{k-2}} - \frac{d}{2^{k+1}} > \frac{d}{2^{k-1}},$$

which is inconsistent with $p \in C_k$. Thus $|p - q| \geq \frac{d}{2^{k+1}}$ if $q \notin \partial G$. Hence, if at least one of $l(p), r(p)$ is not in ∂G , we have for the open subarc $J_1(p)$ of $J(p)$ with endpoints $l(p), r(p)$ that $\lambda(J_1(p)) \geq \frac{d}{2^{k+1}}$. If, on the contrary, $\{l(p), r(p)\} \subset \partial G$, then it follows from the preceding lemma that

$$\lambda(J_1(p)) \geq |l(p) - p| + |p - r(p)| \geq \frac{B'd}{2^k} > \frac{B'd}{2^{k+1}},$$

so that this bound holds in all cases for $p \in C_k$ since $B' \leq 1$. Since $J_1(p) \setminus \{l(p), r(p)\} \subset G_{k-1} \cup G_k \cup G_{k+1}$ for $p \in C_k$, it follows from (3.1) that

$$\begin{aligned} B \left(\frac{d}{2^{k+1}} \right)^2 &\geq \mu \left(\left\{ z \in G : \frac{d}{2^{k+1}} \leq \text{dist}(z, \partial G) \leq \frac{d}{2^{k-2}} \right\} \right) \\ &= \mu(G_{k-1} \cup G_k \cup G_{k+1}) \\ &\geq \mu(\cup\{J_1(p) : p \in C_k\}) \geq \eta \lambda(C_k) \frac{B'd}{2^{k+1}}, \end{aligned}$$

by Proposition 2.8(i). From this we have that $\lambda(C_k) \leq \frac{Bd}{\eta B' 2^{k+1}}$. But since $C = \cup\{C_k : k \geq 1\}$, we conclude that $\lambda(C) \leq \frac{Bd}{2B'\eta} \leq \frac{Bd_0}{2B'\eta}$. Since $B = 128\pi$ and B' and d_0 depend only on G , and η depends only on K , we are done. \square

DEFINITION 3.5. Let C be a curve in G with endpoints $p, q \in \partial G$. Let $\theta \in \text{HP}(G, K)$ and let E be a curve in G with endpoints $a, b \in \partial G$. Let $C = ab \subset G$ and $E = \{z = z(s) : 0 \leq s \leq L\}$ be a curve in G with endpoints $a, b \in \partial G$ and $\lim_{s \rightarrow 0} \theta(z(s)) = \theta(a)$ and $\lim_{s \rightarrow L} \theta(z(s)) = \theta(b)$ regularly. Let $C \cup E$ be a curve in G with endpoints $a, b \in \partial G$ and $\lim_{s \rightarrow 0} \theta(z(s)) = \theta(a)$ and $\lim_{s \rightarrow L} \theta(z(s)) = \theta(b)$ singularly.

PROPOSITION 3.6. Let $G \in \mathcal{G}$, $\theta \in \text{HP}(G, K)$, $C = ab \subset G$ and E be a curve in G with endpoints $a, b \in \partial G$ and $\lim_{s \rightarrow 0} \theta(z(s)) = \theta(a)$ and $\lim_{s \rightarrow L} \theta(z(s)) = \theta(b)$ regularly. Let $C \cup E$ be a curve in G with endpoints $a, b \in \partial G$ and $\lim_{s \rightarrow 0} \theta(z(s)) = \theta(a)$ and $\lim_{s \rightarrow L} \theta(z(s)) = \theta(b)$ singularly. Let D be the interior of the simple closed curve $C \cup E$. Obviously, $D \subset G$. Let $p \in C$ and $q \in E$ satisfy

$$|p - q| = \sup\{|z - w| : z \in C \text{ and } w \in E\}.$$

Let C' be the elementary j -characteristic containing p . Since a j -characteristic can have at most one point in common with an i -characteristic, there is a point $q' \in E$ such that C' contains a subarc $J = pq'$ whose interior lies in D . Let $\text{diam}(C) = |z_1 - z_2|$, where $z_1, z_2 \in C$. Then

$$\text{diam}(C) = |z_1 - z_2| \leq |z_1 - a| + |a - b| + |b - z_2| \leq 2|p - q| + \lambda(E),$$

since $|z_1 - a|, |b - z_2| \leq |p - q|$. But

$$|p - q| \leq |p - q'| + |q' - q| \leq \lambda(J) + \lambda(E),$$

so that

$$\text{diam}(C) \leq 2\lambda(J) + 3\lambda(E).$$

Let J be parameterized by $w(s), 0 \leq s \leq l$, with $w(0) = p$. Let I^+ and I^- be the set $s \in (0, l)$ at which $D\theta(w(s)) \geq 0$ and $D\theta(w(s)) \leq 0$, respectively. For $s \in I^+$ ($s \in I^-$) let E^+ (E^-) be the set of points of E joined to $w(s)$ by an i -arc emanating to the

right (left) of J . It follows by a simple argument based on Proposition 2.7 (length monotonicity) that $\lambda(I^+) \leq \lambda(E^+)$ and $\lambda(I^-) \leq \lambda(E^-)$. But then

$$\lambda(J) \leq \lambda(I^+) + \lambda(I^-) \leq \lambda(E^+) + \lambda(E^-) \leq \lambda(E),$$

from which the desired bound follows immediately. \square

PROPOSITION 3.7. Let $p \in \partial G$, $q_1, q_2 \in G$, $C_1 = pq_1$, $C_2 = pq_2$ be elementary i -characteristics of K with $\theta(C_1) \cap G = \emptyset$. Let $J \subset G$ be a j -arc with endpoints $c_1, c_2 \in C_1 \cap G$. Let $T > 0$ be such that $\lambda(P) = 0$ for any $P \subset J$ with $\text{dist}(P, T) \geq T$.

Suppose not. Then, since the elementary i -characteristic passing through any point of J must exit at p , after replacing the original J by an appropriate subarc and changing C_1 and C_2 accordingly, we can assume that $\lambda(P \cap J) = \epsilon > 0$, $\lambda(N) < \frac{\epsilon}{8}$, and $\lambda(\theta(J)) < \frac{1}{100K}$, where $N = J \setminus P$. Let C_k be parameterized by $z_k(s)$, $0 \leq s \leq \lambda_k$, with $z(0) = p$. Let $\{k, l\} = \{1, 2\}$. Let $J_k(s)$ denote the j -arc joining $z_k(s)$ to a point $w_k(s) \in C_l \cap G$. Note that we know only that $J_k(s)$ is defined for $s \leq \lambda_k$ sufficiently near λ_k . It follows from length monotonicity (Proposition 2.7) that $\lambda(J_k(s)) \geq \lambda(P_k(s)) \geq \epsilon$, where $P_k(s)$ is the set of points of $J_k(s)$ joined to points of P by an i -arc. By the quasi-HP property $\lambda(\theta(J_k(s))) < \frac{1}{100}$, so that $J_k(s)$ is almost straight and in particular the distance between its endpoints is at least $\frac{1}{2}\lambda(J_k(s)) \geq \frac{1}{2}\lambda(P_k(s)) \geq \frac{\epsilon}{2}$. Let ξ_k be the infimum of all s for which $J_k(s)$ is defined. Since the distance between the endpoints of $J_k(s)$ is at least $\frac{\epsilon}{2}$, it is clear that at least one of ξ_1, ξ_2 must be positive, and for definiteness we assume that $\xi_1 > 0$. For $\sigma \in (\xi_1, \lambda_1]$ let $J_1(\sigma)$ be parameterized by $\zeta(s, \sigma)$, $0 \leq s \leq \lambda(J_1(\sigma))$, with $\zeta(0, \sigma) \in C_1$. It is clear that there are $\delta, T > 0$ such that

$$(3.2) \quad \text{dist}(\zeta(s, \sigma), C_1 \cup C_2) \geq Ts \quad \text{for } s \in (0, \delta), \quad \sigma \in (\xi_1, \lambda_1].$$

From the fact that a j -arc can intersect an i -arc at most once in G it easily follows that for each point $z \in C_1 \cap G$ there is a $\delta_1 = \delta_1(z)$ such that

$$(3.3) \quad |z - \zeta(s, \sigma)| \geq \delta_1 \quad \text{for } s \in [\delta, \lambda(J_1(\sigma))], \quad \sigma \in (\xi_1, \lambda_1].$$

From (3.2) and (3.3) together with the fact that for $\sigma \in (\xi_1, \lambda_1]$, $J_1(\sigma) = J_2(\sigma')$ for some $\sigma' \in (\xi_2, \lambda_2]$ it follows that as $\sigma \rightarrow \xi_1$, $J_1(\sigma)$ tends to an arc J_0 which contains p , which joins $z_1(\xi_1)$ to the point $z_2(\xi_2)$ of C_2 , and the distance between whose endpoints is at least $\frac{\epsilon}{2}$. Furthermore, either J_0 consists of a j -arc joining $z_1(\xi_1)$ to p or (in the case that $\xi_2 > 0$) it consists of such an arc together with another j -arc joining p to $z_2(\xi_2)$. One of the endpoints of J_0 , which we henceforth call q , is at a distance of at least $\frac{\epsilon}{4}$ from p . By renaming, if necessary, we can assume that $q = z_1(\xi_1) \in C_1$. Let J'_0 be the subarc of J_0 which joins q to p . Let A denote the arc pq of C_1 . Then J'_0 and A form the two sides of a ‘‘characteristic bilateral’’ B . Let P'_0 be the subset of points of J'_0 which correspond to (i.e., are joined by i -arcs to) points of P (that is, P'_0 is, apart from a set of linear measure 0, the set of points at which J_0 is nonconcave towards the inside of B), and let $N'_0 \subset J'_0$ be the subset corresponding to N . Since, by length monotonicity, $\lambda(N'_0) \leq \lambda(N) \leq \frac{\epsilon}{8}$, it follows that

$$(3.4) \quad \lambda(P'_0) \geq \lambda(J'_0) - \lambda(N'_0) \geq \frac{\epsilon}{4} - \frac{\epsilon}{8} = \frac{\epsilon}{8}.$$

Let J'_0 be parameterized by $z_0(s)$, $0 \leq s \leq \lambda(J'_0)$, with $z_0(0) = p$, and for $s \in (0, \lambda(J'_0))$ let $A(s)$ be the part of the elementary i -characteristic through $z_0(s)$ in \overline{B} , so that $A(s)$ joins $z_0(s)$ to p (since its interior can cross neither A nor J'_0). It follows from Proposition 3.6 that $\text{diam}(A(s)) \rightarrow 0$ as $s \rightarrow 0$. Let qp' be an arc of A for which $\lambda(\theta(qp')) < \frac{1}{100K}$. Since $\lambda(\theta(W)) < \frac{1}{100}$ for any translate W of either J'_0 or qp' by the quasi-HP property, it follows that one can translate qp' in B all the way down $J'_0 \setminus \{p\}$ from q without meeting the boundary of B , since for simple geometric reasons all these translates, being essentially perpendicular to the virtually straight arc J'_0 , must stay away from p . If $A'(s)$ denotes the translate of qp' with initial point $z_0(s)$, then $A'(s) \subset A(s)$, and therefore $\text{diam}(A'(s)) \rightarrow 0$ as $s \rightarrow \lambda(J'_0)$. This means that each point $z \in qp'$ is joined to p in B by a j -arc $J''(z)$ such that if $P''(z)$ is the set of points of $J''(z)$ joined to points of P by i arcs in T , then by (3.4) and length monotonicity $\lambda(P''(z)) \geq \frac{\epsilon}{8}$ and by the quasi-HP property $\lambda(\theta(J''(z))) < \frac{1}{100}$, so that

$$(3.5) \quad |z - p| \geq \frac{\epsilon}{16}.$$

We can now repeat this process starting with $J''(p')$ instead of $J'_0 = J''(q)$ and continue doing so to obtain in the end $J''(z)$ for all $z \in A \setminus \{p\}$. However, by the argument we just gave we now have (3.5) for all $z \in A$, which is absurd since p is an endpoint of A . Therefore $\lambda(P) = 0$. \square

PROPOSITION 3.8.

This follows easily from Proposition 3.6. \square

PROPOSITION 3.9.

Assume to the contrary that points $p_1 \neq p_2$ of ∂G are joined by distinct elementary i -characteristics C_1 and C_2 . It follows from Proposition 3.8 that the elementary i -characteristic through any point of the simply connected domain D bounded by $C_1 \cup C_2$ must also have endpoints p_1 and p_2 . But then it follows easily from Proposition 3.7 that all j -characteristic arcs in D are straight line segments. But this contradicts Proposition 3.7. \square

PROPOSITION 3.10.

The proof is an easy consequence of Proposition 3.7 and the quasi-HP property. \square

DEFINITION 3.11.

θ elementary i -characteristic C_0 of G . $a, b \in \partial G \cap B$. $\partial D = C_0 \cup B$. i -characteristic subdomain

When we wish to indicate the elementary i -characteristic involved, we will denote the i -characteristic subdomain by (D, C_0) . The arc $B = \overline{\partial D} - C_0 \subset \partial G$, called the bot of (D, C_0) and denoted by $\text{bot}(D)$, will be considered to have the order “ $<$ ” corresponding to the positive orientation of ∂D . We shall freely use interval notation as well as the terms “to the right of,” “to the left of,” “between,” etc., when dealing with $\text{bot}(D)$. Furthermore, if ab is an elementary i -characteristic joining points a, b of $\text{bot}(D)$, it will be understood that $a \leq b$, unless otherwise indicated. When dealing with a characteristic subdomain (D, C_0) , we shall work with the class $\mathcal{I}(D)$ of nontrivial elementary i -characteristics C which join points p, q of $\text{bot}(D)$. Note that $C_0 \in \mathcal{I}(D)$. If $C = pq \in \mathcal{I}(D)$, then “above” C refers to the part of D not in the closed region bounded by $C \cup [p, q]$. If $F \subset \overline{D}$ is a compact set for which $F \cap \text{bot}(D) \neq \emptyset$, we

say that an elementary i -characteristic $E = pq \in \mathcal{I}(D)$ is contained in F , and write $F \preceq E$, if F is contained in the closed set bounded by the simple closed curve $[p, q] \cup pq$. We deal only with sets F for which each component of F has points in $\text{bot}(D)$. For two such sets $F_1, F_2 \subset \overline{D}$ we say that F_1 is contained in F_2 and write $F_1 \leq F_2$ if for all f_1, f_2 with $f_k \in F_k \cap \text{bot}(D)$, $k = 1, 2$, there holds $f_1 \leq f_2$. It is clear that if $C_1, C_2 \in \mathcal{I}(D)$, then one of the following is true: $C_1 \preceq C_2$, $C_2 \preceq C_1$, $C_1 \leq C_2$, or $C_2 \leq C_1$. The first and second of these possibilities include the case $C_1 = C_2$.

In what follows, $L(G, p, d)$ will denote the segment \overline{pq} of length d which is orthogonal to ∂G at p , which emanates from p into G and which is oriented from p to q .

PROPOSITION 3.12. Let $G \in \mathcal{G}(\rho)$, $\omega_0 = \omega_0 < 1$, $\omega_1 = \omega_1(K, \rho)$, $G \in \mathcal{G}(\rho)$, $p \in \partial G$, $\theta \in \text{HP}(G, K)$, $C \in \mathcal{I}(G)$, $z(s)$, $0 \leq s \leq \lambda(C)$.

Assume that C is tangent to ∂G at p and that ρ is the radius of curvature of ∂G at p . Then

$$(3.6) \quad C \cap L(G, p, \omega_0 \rho) = \{z(0)\},$$

if C is tangent to ∂G at p and $L(G, p, \omega_0 \rho)$ is the segment of length $\omega_0 \rho$ emanating from p into G and orthogonal to ∂G at p . Then $w(s)$, $0 \leq s < \lambda(\partial G)$, $w(0) = p$, $|\arg\{z'(0)\} - \arg\{w'(0)\}| \leq \frac{\pi}{2}$.

$$\arg\{z'(0)\} \geq \arg\{w'(0)\} - \omega_1 \sqrt{|z(0) - p|}$$

if C is tangent to ∂G at p and $L(G, p, \xi_0)$ is the segment of length ξ_0 emanating from p into G and orthogonal to ∂G at p .

$$\arg\{z'(0)\} \geq \arg\{w'(0)\} + \omega_1 \sqrt{|z(0) - p|}$$

if C is tangent to ∂G at p and $L(G, p, \omega_0)$ is the segment of length ω_0 emanating from p into G and orthogonal to ∂G at p .

Clearly, it is enough to handle the case in which C emanates to the right of $L(G, p, \rho)$. Without loss of generality we can assume that $p = 0 = \arg\{w'(0)\}$, so that in particular $D = N(-\rho i, \rho) \subset \mathbb{C} \setminus G$ and $\partial N(-\rho i, \rho)$ is tangent to ∂G at 0 . Let $d = |z(0) - p|$ and let $R_d = \{di + te^{-i\alpha} : t \geq 0\}$, $\alpha = \alpha(d) \in (0, \frac{\pi}{2})$ be the ray emanating to the right of $L(G, 0, \rho)$ from the point $di \in L(G, 0, \rho)$ which is tangent to ∂D , and let the point of tangency be z_d . Then $\alpha = \cos^{-1}(\frac{\rho}{\rho+d})$, so that α is bounded below and above by $\sqrt{1 - \frac{\rho}{\rho+d}}$ and $2\sqrt{1 - \frac{\rho}{\rho+d}}$, respectively. Thus

$$(3.7) \quad 2\sqrt{d/\rho} \geq \alpha \geq \frac{\sqrt{d/\rho}}{2} \quad \text{for } d \leq \rho.$$

Let T_d be the curvilinear triangle bounded by $L(G, 0, d)$, the line segment $[di, z_d]$ and the (shorter) arc of ∂D with endpoints $0, z_d$. Then in light of (3.7) it is easy to see that

$$(3.8) \quad \text{diam}(T_d) \leq \rho\alpha + d \leq 2\rho\sqrt{d/\rho} + d = 2\sqrt{\rho d} + d \quad \text{for } d \leq \rho,$$

so that

$$(3.9) \quad \text{diam}(T_d) < \rho \quad \text{for } d \leq \frac{\rho}{16}.$$

Now assume that C is as in the hypothesis with $z(0) = di$, where $d \leq \frac{\rho}{16}$, and that

$$-\beta = \arg\{z'(0)\} \leq -N\alpha.$$

Here $N > 1$ is a number yet to be determined. Since, by (3.9), $\text{diam}(C) \geq \rho > \text{diam}(T_d)$, and since obvious angle considerations imply that C enters T_d at di , in light of assumption (3.6), C must exit T_d at a point of the segment $[di, z_d]$. Say s_0 is the smallest value of s for which $z(s) \in [di, z_d]$. Then obviously $\theta([0, s_0]) \supset [-\beta, -\alpha]$, so that $[0, s_0]$ has a subinterval $[s_1, s_2]$ for which $\theta([s_1, s_2]) = [-\beta, -\alpha]$. Since $\beta - \alpha \leq \frac{\pi}{2}$, and since the length of a curve given by $y = f(x)$, $x_1 \leq x \leq x_2$ for which $|f'(x)| \leq 1$ is at most $\sqrt{2}|x_2 - x_1|$, it follows that

$$\begin{aligned} s_2 - s_1 &= \lambda(z([s_1, s_2])) \leq \sqrt{2}|z(s_2) - z(s_1)| \\ &< 2\text{diam}(T_d) \leq 2(2\sqrt{\rho d} + d), \end{aligned}$$

by (3.8). Thus by the mean value theorem there must be some $s_0 \in [s_1, s_2]$ for which

$$(3.10) \quad D^+\theta(z(s_0)) \geq \frac{(N - 1)\alpha}{2(2\sqrt{\rho d} + d)}.$$

But the curvature bound (Proposition 2.6) together with the hypothesis regarding the j -characteristics implies that $D^+\theta(z(s_0)) \leq \frac{K}{\rho}$, so that in light of (3.10) and the lower bound in (3.7) we have

$$\frac{\rho(N - 1)\sqrt{d/\rho}}{8\sqrt{\rho d} + 4d} \leq K.$$

But the left-hand side of this inequality is $(N - 1)/(8 + 4\sqrt{d/\rho})$, which is bounded below by $\frac{N-1}{9}$ for $d \leq \frac{\rho}{16}$, so that we have a contradiction for $N = 9K + 2$. Thus if $d \leq \frac{\rho}{16}$, we must have

$$\arg\{z'(0)\} \geq -2\frac{9K + 2}{\sqrt{\rho}}\sqrt{d} \geq -\frac{22K\sqrt{d}}{\sqrt{\rho}},$$

in light of the upper bound in (3.7) and the fact that $K \geq 1$. Thus we have proved the proposition with $\omega_0 = \frac{1}{16}$ and $\omega_1 = \frac{22K}{\sqrt{\rho}}$. \square

PROPOSITION 3.13. *Let $G \in \mathcal{G}(\rho)$, $\theta \in \text{HP}(G, K)$, $C = ab$ be a curve in (D, C) with $A = \text{bot}(D)$, $z(s)$, $0 \leq s \leq \lambda(C)$, $w(s)$, $0 \leq s \leq \lambda(A)$, $z(0) = w(0) = a$, $\text{dist}(p, \partial G \setminus A) \geq \rho$, $L(G, p, \xi_0)$ be a line segment in C with endpoints $p = w(\sigma_0) \in A$ and $z(s)$. Then*

$$(3.11) \quad |\arg\{z'(s)\} - \arg\{w'(\sigma_0)\}| \leq \xi_1 \sqrt{|z(s) - p|}.$$

Without loss of generality we may assume that $p = 0 = \arg\{w'(\sigma_0)\}$. Let ω_0 and ω_1 be as in the preceding proposition. Let $L(\epsilon) = L(G, p, \epsilon)$. We assume for the moment that $z(s) \in L(\omega_0 \frac{\rho}{2})$ is the only point of C on this segment. Since immediately to the left of C there are points in the complement of D and the interior of the segment $L(|z(s)|)$ lies in D , it is clear that if we write $\arg\{z'(s_0)\} = e^{i\tau}$, with $-\pi < \tau \leq \pi$, then $|\tau| \leq \frac{\pi}{2}$. Let E be either of the two subarcs of C , one of whose endpoints is $z(s)$ and the other of which is a point $q \in G$ for which $|z(s) - q| = \frac{\rho}{2}$. It follows from the preceding proposition that if $\omega_1(K, \frac{\rho}{2})\sqrt{|z(s)|} \leq \frac{\pi}{4}$, then C cannot

be tangent to $L(\omega_0 \frac{\rho}{2})$. Thus, of the two arcs E , one moves to the right of $L(\omega_0 \frac{\rho}{2})$ as we move along it away from $z(s)$ and the other moves to the left. But then by the preceding proposition we have (3.11) with

$$(3.12) \quad \xi_1 = \omega_1 \left(K, \frac{\rho}{2} \right)$$

for any $z(s)$ for which

$$|z(s)| \leq \min \left\{ \omega_0 \frac{\rho}{2}, \left(\frac{\pi}{4\omega_1(K, \rho/2)} \right)^2 \right\}.$$

Let

$$(3.13) \quad \xi_0 = \min \left\{ \omega_0 \frac{\rho}{10}, \left(\frac{\pi}{4\omega_1(K, \rho/2)} \right)^2, \frac{\pi\rho}{21K} \right\}.$$

Assume $\epsilon \leq \xi_0$ and that $L(\epsilon)$ contains at least two points of C . Then there are $s', s'' \in (0, \lambda(C))$, $s' < s''$, such that the interior of $L(\min\{|z(s')|, |z(s'')|\})$ contains no point of C (and is therefore contained in D) and $z(s')z(s'') \cap L(\epsilon) = \{z(s'), z(s'')\}$. By Proposition 3.6, $\text{diam}(z(s')z(s'')) \leq 5\epsilon < \omega_0 \frac{\rho}{2}$. There are the following two cases.

(i) $|z(s')| < |z(s'')|$. In this case, simple topological arguments show that $z(s')$ $z(s'')$ must lie to the right of $L(\epsilon)$. Also, by the foregoing and our definition of ξ_0 , one easily has

$$|\arg\{z'(s')\} - \arg\{w'(\sigma_0)\}| \leq \frac{\pi}{4}.$$

But then since C crosses $L(\epsilon)$ again at $z(s'')$, we must have that for some $t' < t''$ in (s', s'') , $\arg\{z'(t'')\} - \arg\{z'(t')\} \geq \frac{\pi}{2} - \frac{\pi}{4} = \frac{\pi}{4}$, so that by mean value considerations as in the proof of the preceding proposition together with the curvature bound we see that there must be a point $\sigma' \in (t', t'')$ at such that

$$\frac{K}{\rho/2} \geq \frac{d \arg\{z'(s)\}}{ds} \Big|_{s=\sigma'} \geq \frac{\pi/4}{2 \text{diam}(z(s')z(s''))} \geq \frac{\pi}{10\epsilon} \geq \frac{\pi}{10\xi_0},$$

which implies that $\xi_0 \geq \frac{\pi\rho}{20K}$, a contradiction, since $\xi_0 \leq \frac{\pi\rho}{21K}$.

(ii) $|z(s')| > |z(s'')|$. In this case, simple topological arguments show that $z(s')$ $z(s'')$ must lie to the left of $L(\epsilon)$, and we proceed analogously to the way we did in case (i).

Thus, for $\epsilon \leq \xi_0$, $L(\epsilon)$ can contain at most one point of C . This completes the proof of the proposition with ξ_1 and ξ_0 as defined in (3.12) and (3.13), respectively. \square

We now extend the notion of characteristic to include certain arcs whose interiors contain points of ∂G . Although what follows does not give the most exhaustive extension possible, it is sufficient for our present needs. Let (D_0, C_0) be an i -characteristic subdomain of G , and consider a monotone decreasing sequence $\{C_k = a_k b_k : k \geq 0\}$ in $\mathcal{I}(D_0)$. In other words, $a_k \leq a_{k+1} < b_{k+1} \leq b_k$, $k \geq 0$, so that the arcs $A_k = [a_k, b_k]$ of $\text{bot}(D_0) = A_0$ are nested. Let $D_k \subset D_0$ be the i -characteristic subdomain bounded by $C_k \cup A_k$, so that $D_k \supset D_{k+1}$, $k \geq 0$. We regard C_k as being oriented from a_k to b_k . Let $a_k \rightarrow a$ and $b_k \rightarrow b$. We define C to be the set of limits of sequences $\{z_l\}$, where $z_l \in C_{k_l}$, $k_l \rightarrow \infty$. Any such C will be called an i -characteristic consisting of a sequence of arcs $a \dots b$. An extended characteristic consisting of a single point $p \in \partial G$ will be called

The following proposition, in the statement and proof of which the notation is the same as that of the immediately preceding sentences, contains the basic properties of extended characteristics. Note that $M(G, K)$ is, as in Proposition 3.4, an upper bound on the length of elementary characteristics of K -quasi-HP functions on G .

PROPOSITION 3.14. Let C be a characteristic of G with endpoints a and b , $a < b$. Let $C = \{a\}$ or $C = \{b\}$ or $C = \{a, b\}$. Let $\lambda(C) \leq M(G, K)$. Let $z(s)$, $0 \leq s \leq \lambda(C)$, be a parameterization of C with $z(0) = a$ and $C \cap \partial G \subset [a, b]$. Let $A_0 = \text{bot}(D_0)$. Let $J(t)$ be the j -half-characteristic of C starting at $(0, \lambda(C))$ and ending at $(t, \lambda(C))$, $t \in (0, \lambda(C))$, $z(t) \in \partial G \setminus (a, b)$.

$$(3.14) \quad D^+ \arg\{z'(t)\} \leq \frac{K}{\lambda(J(t))}.$$

We begin by observing that

$$(3.15) \quad A_0 \cap C \subset [a, b].$$

To see this, note that for each $w \in A_0 \setminus [a, b]$ there is an n for which $w \notin A_n$. Since $C_n \cap A_0 = \{a_n, b_n\}$, w is not in C_n either, and therefore $w \notin \overline{D_n}$. But then, by the monotonicity of $\{\overline{D_k}\}$,

$$\text{dist}(w, C_k) \geq \text{dist}(w, \overline{D_k}) \geq \text{dist}(w, \overline{D_n}) > 0 \text{ for } k \geq n,$$

from which the desired conclusion follows at once. From (3.15) it follows immediately that $C \cap \partial G \subset [a, b]$.

Next we note that if $a = b$, then Proposition 3.6 implies that $\text{diam}(C_k) \rightarrow 0$, so that $C = \{a\}$. For the remainder of the proof we therefore assume that $a < b$.

By this assumption and Proposition 3.4 on the boundedness of the lengths of characteristics there exist l_1 and l_2 such that $0 < l_1 \leq \lambda(C_k) \leq l_2 < \infty$. Let C_k be parameterized by $z = z_k(s)$, $0 \leq s \leq \lambda(C_k)$. Let $0 < \epsilon < |b - a|/3$, and let n_0 be such that

$$|a_k - a|, |b_k - b| < \epsilon \text{ for } k \geq n_0.$$

For $k \geq n_0$ let

$$\alpha_k = \sup\{s : z_k(s) \in \partial N(a, \epsilon)\}$$

and

$$\beta_k = \inf\{s : z_k(s) \in \partial N(b, \epsilon)\},$$

and let $E_k = E_k(\epsilon) = z_k([\alpha_k, \beta_k])$. Then

$$(3.16) \quad \delta = \inf\{\text{dist}(E_k, \partial G \setminus A_k) : k \geq n_0\} > 0,$$

since otherwise there would be a point of C in $A_0 \setminus [a, b]$, in contradiction of (3.15). For $k \geq n_0$ and $s \in [\alpha_k, \beta_k]$, let $J_k(s)$ denote the j -half-characteristic emanating to the left of C_k and joining $z_k(s)$ to a point $w_k(s)$ of $\partial G \setminus A_k$. In light of (3.16),

$$(3.17) \quad \lambda(J_k(s)) \geq \delta > 0 \text{ for } k \geq n_0, \quad s \in [\alpha_k, \beta_k].$$

This means that for each $\epsilon > 0$ we have an upper bound on the curvature to the left of E_k . More precisely, the curvature bound implies that

$$(3.18) \quad D^+\theta(z_k(s)) \leq \frac{K}{\lambda(J_k(s))} \leq \frac{K}{\delta} \quad \text{for } k \geq n_0, \quad s \in [\alpha_k, \beta_k].$$

In addition, the curvature bound trivially implies that

$$(3.19) \quad D^-\theta(z_k(s)) \geq -\frac{K}{\text{dist}(z_k(s), \partial G)}.$$

There is a neighborhood U of ∂G such that for each $z \in U$ there is a unique $p(z) \in \partial G$ for which $|z - p(z)| = \min\{\text{dist}(z, \zeta) : \zeta \in \partial G\}$ and for which p is continuous. For $q \in \partial G$, let $e^{i\phi(q)}$ be the positively oriented unit tangent to ∂G at q , so that, in U , $e^{i\phi(p(z))}$ is continuous. It follows from (3.18) and Proposition 3.13 that there is a $\Lambda = \Lambda(\epsilon)$ such that

$$(3.20) \quad |z'_k(s) - e^{i\phi(p(z_k(s)))}| \leq \Lambda \sqrt{|z_k(s) - p(z_k(s))|} \quad \text{for } k \geq n_0, \quad s \in [\alpha_k, \beta_k].$$

Bounds (3.18), (3.19), and (3.20) imply that the family $\{z'_k(s)\}$ is uniformly bounded and equicontinuous (that is, there is a single modulus of continuity valid for all $z'_k(s)$, $k \geq n_0$, on the respective $[\alpha_k, \beta_k]$). From this it follows that there are α, β , and a sequence $\{k_l\}$ such that $\alpha_{k_l} \rightarrow \alpha$ and $\beta_{k_l} \rightarrow \beta$ and $z_{k_l}(s) \rightarrow z_\epsilon(s)$ uniformly (in the obvious sense), where z_ϵ is continuously differentiable and parameterizes an arc $C(\epsilon)$ which joins a point of $\partial N(a, \epsilon)$ to a point of $\partial N(b, \epsilon)$ in the part of $\overline{D_0}$ lying outside of both these circles. The arc $C(\epsilon)$ is simple, since otherwise there would be some $q \in (a, b) \setminus N(\{a, b\}, \epsilon)$ such that for arbitrarily small $\tau > 0$, $L(G, q, \tau)$ intersects some E_k more than once, which would contradict Proposition 3.13. Clearly, $\lambda(C(\epsilon)) \leq M(G, K)$. The nested nature of the D_k then implies that (at least for sufficiently small ϵ) the entire sequence converges to $C(\epsilon)$. Also, it is clear that $C(\epsilon')$ is an extension of $C(\epsilon)$ for $\epsilon' < \epsilon$. We have that $C = \cup\{C(\epsilon) : \epsilon > 0\}$, since $\text{diam}(z_k([0, \alpha_k]))$ and $\text{diam}(z_k([\beta_k, \lambda(C_k)]))$ tend to 0 uniformly in k as $\epsilon \rightarrow 0$ by Proposition 3.6. It follows from $C = \cup\{C(\epsilon) : \epsilon > 0\}$ that $\lambda(C) \leq M(G, K)$ and that C is a simple arc with endpoints a, b , and furthermore that C is parameterized by a function $z(s)$ which is continuously differentiable on $(0, \lambda(C))$. That the two possible orderings of the points of $C \cap [a, b]$ coincide follows from the fact that C is a simple arc and (3.15).

The existence and uniqueness of $J(t)$ is trivial when $z(t) \in G$. When $z(t) \in (a, b)$ the existence of $J(t)$ follows from a straightforward compactness argument. In light of (3.20), for such t , any corresponding $J(t)$ must be orthogonal to ∂G at $z(t)$, so that the uniqueness of $J(t)$ follows from Proposition 3.10. Bound (3.14) follows from (3.18). \square

We will refer to an extended characteristic C joining $a, b \in \partial G$ as ab ; points of $ab \cap \partial G$ will be called a_i, \dots, a_{i-1}, a_i , and contact points other than a and b will be called $a'_k, \dots, a'_{k+1}, a'_k$. It is clear that θ can be continuously extended to $G \cup (ab \setminus \{a, b\})$. For what is to follow it is important to understand that if (D_0, C_0) is an i -characteristic subdomain, then, in addition to the extended characteristics constructed above, D_0 might contain extended characteristics C' joining points $a' < b'$ arising from a sequence of i -characteristic subdomains $\{(D'_k, C'_k)\}$, where $C'_k = a'_k b'_k$ with $a'_{k+1} \leq a'_k < b'_k \leq b'_{k+1}$ (where the order is with respect to $A_0 = \text{bot}(D_0)$, as above). Here the C'_k are contained in D_0 , but the other part of the boundary of D'_k is $\partial G \setminus (a'_k, b'_k)$ (where here again the interval notation refers to the order on

$A_0 = \text{bot}(D_0)$). Note that if for such an extended characteristic C' , $C' \cap \partial G$ has points other than a' and b' (that is, if C' is not simply an elementary characteristic), then the contact points will not occur monotonically with respect to the order on A_0 when C' is traversed from a' to b' . The extended characteristics C constructed originally will be referred to as *proper*, and this other kind of extended characteristic C' with proper contact points will be said to be *improper*. We consider that an extended characteristic ab exits G at all of its contact points, and we use the terms “exits regularly” and “exits singularly” at a or b in the obvious fashion. Clearly, ab exits regularly at its proper contact points.

PROPOSITION 3.15. *Let C be an extended characteristic in G with proper contact points a, b and p . Let $z(s)$, $0 \leq s \leq L$, be a j -characteristic with $z(0) = p$ and $\phi(s) = \arg\{z'(s)\}$. Let $E(s)$, $s \in (0, L)$, be a j -characteristic with $E(s) \cap C = \{z(s)\}$. Let $\epsilon, T > 0$, $P(\epsilon, T) = \{s < \epsilon : \phi'(s) \geq T\}$ and $N(\epsilon, T) = \{s < \epsilon : \phi'(s) \leq -T\}$. Then $\lambda(P(\epsilon, T)) > 0$, $\lambda(N(\epsilon, T)) > 0$, $\text{diam}(E(s)) \rightarrow 0$ as $s \rightarrow 0$.*

Proof. Assume to the contrary that $\lambda(P(\epsilon, T)) > 0$ for some $\epsilon_0, T_0 > 0$, $\lambda(P(\epsilon_0, T_0)) = 0$, so that $\phi'(s) \leq T_0$ a.e. on $(0, \epsilon_0]$. Then $\phi = \phi^+ + \phi^-$ on $(0, \epsilon_0]$, where ϕ^+ is Lipschitz continuous and nondecreasing and ϕ^- is continuous and nonincreasing. From this in turn it follows that either ϕ^- has a finite limit as $s \rightarrow 0$ or it tends to ∞ as $s \rightarrow 0$, so that in fact the latter is the case since C exits singularly. But this means that C spirals around p , which is clearly impossible. Thus $P(\epsilon, T)$ must have positive measure for all $\epsilon, T > 0$. One sees similarly that $N(\epsilon, T)$ has positive measure. That $\text{diam}(E(s)) \rightarrow 0$ follows immediately from the characteristic length bound. To see that the last sentence of the statement is true, assume to the contrary that $E(\sigma)$ joins $z(\sigma)$ to p . Then, since distinct j -characteristics have no common point in G and can only cross C at one point in G , $E(s)$ exits at p for $0 < s \leq \sigma$. But this is impossible in light of Proposition 3.7 and the fact that both $P(\epsilon, 1)$ and $N(\epsilon, 1)$ have positive measure. \square

Let (D, C_0) be a characteristic subdomain of G , and let $F \subset D \cup \text{bot}(D)$ be a compact set for which $F \cap \text{bot}(D) \neq \emptyset$. The family $\mathcal{C}(F) = \{C \in \mathcal{I}(D) : F \preceq C\}$ is clearly linearly ordered with respect to the relation \preceq . From this it easily follows that there is a unique monotone extended i -characteristic $E_0 = ab$ for which $F \cap \text{bot}(D) \subset [a, b]$ and $E_0 \preceq C$ for all $C \in \mathcal{C}(F)$.

NOTATION 3.16. Let (D, C_0) be a characteristic subdomain of G . Let $E_0 = ab$ be the unique extended i -characteristic such that $E_0 \preceq C$ for all $C \in \mathcal{C}(F)$. Let $\min_D(F)$ denote the set $E_0 \cap \text{bot}(D)$.

In the statement and proof of the following proposition all order relations are with respect to (D, C_0) .

PROPOSITION 3.17 (structure of $\min_D(A \cup B)$). *Let (D, C_0) be a characteristic subdomain of G . Let $a_1 < a_2 < b_1 < b_2$ be points in $\text{bot}(D)$. Let $A = a_1a_2$ and $B = b_1b_2$ be elementary characteristics in $\mathcal{I}(D)$. Let $U = \min_D(A \cup B) = ef$, $e < f$. Then:*

- (i) $U = ef \in \mathcal{C}(A \cup B)$.
- (ii) $ea \preceq L \preceq U \preceq L \preceq U$ for some $L \in \mathcal{I}(D)$.
- (iii) $bf \preceq L \preceq U \preceq L \preceq U$ for some $L \in \mathcal{I}(D)$.

Proof. It follows from the hypotheses that $e < f$ (i.e., that $e \neq f$). We show that

if (i) does not hold, then at least one of (ii) or (iii) does. Let

$$e' = \sup\{z : z \in \text{bot}(D) \cap \min_D(A \cup B), z \leq a_1\}$$

and

$$f' = \inf\{z : z \in \text{bot}(D) \cap \min_D(A \cup B), z \geq b_2\}.$$

Since (i) does not hold it is easy to see that the subarc $e'f'$ of ef is a member of $\mathcal{C}(A \cup B)$, so that in fact $e'f' = ef$. Obviously, $ef \cap (A \cup B) \subset \{e, f\}$, since otherwise ef would have to be A or B . Let w_0 be any point of ef other than e or f , and consider a small j -arc E whose initial point is w_0 , which extends to the right of ef (that is, into the simply connected domain bounded by $ef \cup [e, f]$) and which is disjoint from $A \cup B$. Let E be parameterized by $z = w(s)$, $0 \leq s \leq \lambda_0$, with $w(\lambda_0) = w_0$, and let $C(s)$ denote the elementary i -characteristic through $w(s)$. Note that

$$(3.21) \quad C(s_1) \preceq C(s_2) \quad \text{for } s_1 \leq s_2.$$

Let the left and right endpoints of $C(s)$ be $l(s)$ and $r(s)$. By the minimality of ef it must be that for no $s \in (0, \lambda_0)$ can we have both $l(s) \leq a_1$ and $r(s) \geq b_2$. From (3.21) it therefore follows that either for all $s \in (0, \lambda_0)$, $l(s) > a_1$, or for all $s \in (0, \lambda_0)$, $r(s) < b_2$. Assume the latter occurs. Then in fact $r(s) \leq b_1$ for all $s \in (0, \lambda_0)$, since otherwise $C(s)$ would cross B in D , which is impossible since they are distinct elementary i -characteristics. Consider the i -characteristic subdomain $(D', C(0))$ bounded by $C(0)$ and $\partial G \setminus (l(0), r(0))$. If we take any sequence $\{s_k\}$ in $(0, \lambda_0)$ which tends monotonically to λ_0 , and consider the corresponding sequence of elementary i -characteristics $C_k = C(s_k)$ in D' , it is clear that they will give rise to an extended i -characteristic $C' = ab$ (with, as always, $a < b$ with respect to the order for $\text{bot}(D)$), which is nonmonotone with respect to the characteristic subdomain D , which contains ef as a subarc and for which $b \leq b_1$. Clearly, $bf \cap ef = \{f\}$, so that $bf \preceq ef$. If we let

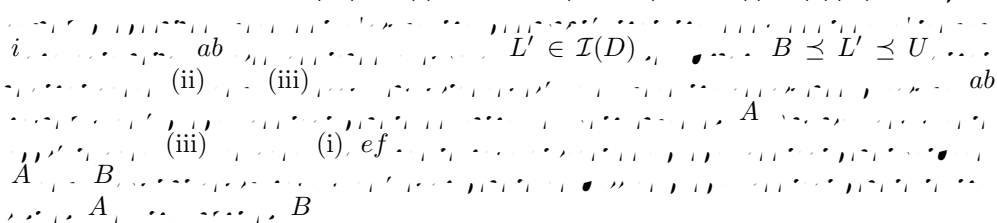
$$c' = \sup\{z : z \in \text{bot}(D) \cap bf, z \leq b_1\}$$

and

$$c'' = \inf\{z : z \in \text{bot}(D) \cap bf, z \geq b_2\},$$

then $L = c'c'' \in \mathcal{I}(D)$, so that from the minimality of ef it follows that $c' \geq a_2$ and we have conclusion (iii). In exactly the same manner one obtains conclusion (ii) in the case that $s \in (0, \lambda_0)$, $l(s) > a_1$ for all $s \in (0, \lambda_0)$. \square

COMMENT 3.18. (ii)



DEFINITION 3.19. $p \in \partial G$, regular boundary point, θ , $\theta(z)$ singular boundary point

PROPOSITION 3.20. Let $p \in \partial G$, $\phi \in (0, \pi)$.

This follows immediately from the Peano existence theorem for (local) solutions of the initial value problem $y' = F(x, y)$, $y(x_0) = y_0$ when F is continuous in a neighborhood of (x_0, y_0) . \square

DEFINITION 3.21. Let $p \in \partial G$, $\theta \in (0, \pi)$. A singularity of type 0 (singularity of type 1) is a singularity of type 2.

The following proposition is needed to make use of fans of characteristics in what is to follow.

PROPOSITION 3.22. Let $C_1 = pq_1$, $C_2 = pq_2$, $A \subset \partial G$, $T = C_1 \cup C_2 \cup A$.

- (i) $d > 0$, $z \in N(p, d) \cap T$.
- (ii) $p \in C$.

Let W be the set of all $z \in T$ for which the elementary i -characteristic $C(z)$ through z does not exit at p . For $z \in W$, $C(z)$ has no points in common with either $C_1 \setminus \{q_1\}$ or $C_2 \setminus \{q_2\}$, so that both endpoints of $C(z)$ lie on A . For each $z \in W$, let $D(z)$ denote the interior of the domain bounded by $C(z)$ and the arc of A whose endpoints are those of $C(z)$. For $w_1, w_2 \in T$, $C(w_1) \cap T$ and $C(w_2) \cap T$ are either identical or disjoint, so that for $w_1, w_2 \in W$, $D(w_1)$ and $D(w_2)$ are either nested or disjoint. Let $\xi = \frac{1}{2} \text{dist}(p, A) > 0$. If $z \in N(p, \frac{\xi}{2}) \cap W$, then it follows from Proposition 2.8(i) that $\mu(D(z)) \geq \eta \xi^2$ since $\lambda(C(z) \cap N(p, \xi)) \geq \xi$ and each of the j -arcs joining a point of $C(z) \cap N(p, \xi)$ to A in $D(z)$ has length at least ξ . Thus if $z_1, \dots, z_n \in N(p, \frac{\xi}{2}) \cap W$ are such that the corresponding $D(z_k)$ are disjoint, then $n \leq \frac{\mu(G)}{\eta \xi^2}$. If (i) is not true, a pigeonhole (area exhaustion) argument then shows that there is a sequence $\{z_k : k \geq 1\}$ of points of W tending to p such that $\{D(z_k)\}$ is an increasing sequence. But our extended characteristic construction then gives us a monotone extended characteristic C (with respect to the characteristic subdomain bounded by $C(z_1)$ and the complement in ∂G of the subarc of A joining the endpoints of $C(z_1)$) for which $p \in C$. But the endpoints of C are in A since those of all the $D(z_k)$ are in A , so that p is a proper contact point of C . Thus C satisfies (ii). \square

DEFINITION 3.23 (i -fan). Let $p \in \partial G$, $\theta \in (0, \pi)$. A fan of i -characteristics $\mathcal{F}_i(p)$ is a set of i -characteristics C with endpoints in A such that $p \in C$.

Proposition 3.22 implies that if C_1 and C_2 are distinct elementary i -characteristics in $\mathcal{F}_i(p)$, then there is a $d > 0$ such that all points of $N(p, d)$ between them belong to members of $\mathcal{F}_i(p)$, so that the interior of $\cup\{C : C \in \mathcal{F}_i(p)\}$ is nonempty. Furthermore, it follows from Proposition 3.7 that the j -characteristic through any point in the interior of $\cup\{C : C \in \mathcal{F}_i(p)\}$ is strictly concave towards the side facing p . These two facts in turn imply that if there are i -fans at $p_1 \neq p_2$, then the interiors of $\cup\{C : C \in \mathcal{F}_i(p_k)\}$, $k = 1, 2$, must be disjoint. Since $\mu(\cup\{C : C \in \mathcal{F}_i(p)\}) > 0$, we have the following.

PROPOSITION 3.24.

PROPOSITION 3.25.

Let G be a domain in the complex plane with boundary ∂G and let $p \in \partial G$ be a point where ∂G is smooth. Let C be a nontrivial arc of ∂G parameterized by $z = z(s)$, $0 \leq s \leq L$, with $z(0) = p$. Let $\theta = \arg\{z'(s)\} = \theta(z(s))$. Let $R = (R_1, R_2)$ be a sector with vertex p .

Let F be an extended i -characteristic which exits regularly at p . Without loss of generality we can assume that the positive direction along ∂G at p is that of the positive real axis. Let C be a nontrivial arc of F parameterized by $z = z(s)$, $0 \leq s \leq L$, with $z(0) = p$. Again, without loss of generality we can assume that $\gamma = \arg\{z'(s)\} = \theta(z(s))$. We consider the following three possibilities separately.

(i) $\gamma \in (0, \frac{\pi}{2}) \cup (\frac{\pi}{2}, \pi)$. In this case we can assume that C is an arc of an elementary i -characteristic exiting at p . To be specific, we assume that $\gamma \in (0, \frac{\pi}{2})$. Let q be an interior point of C . For any $\epsilon > 0$ there exist nontrivial j -characteristic arcs $E^+ = E^+(\epsilon)$ and $E^- = E^-(\epsilon)$ emanating from q to the right and left of C , respectively, such that $\lambda(\theta(E^+ \cup E^-)) < \epsilon$. Since there is no fan at p , none of the j -characteristics through any point of $E^+ \cup E^-$ other than q exits at p . Let $C(\epsilon)$ be an initial segment of C such that $\lambda(\theta(C \setminus \{p\})) < \epsilon$. Let $\alpha < \gamma$; we show that $\theta(z) \rightarrow \gamma$ as $z \rightarrow p$ between C and the ray $\arg\{z\} = \alpha$. It is easy to see that some initial segment of this ray (that is, the portion of the ray contained in $N(p, \xi)$ for some $\xi > 0$) is covered by a collection $\mathcal{Q}(\epsilon)$ of characteristic quadrilaterals Q , one of whose i -sides is a subarc of $C(\epsilon)$ and one of whose j -sides is a translate of an initial subarc of $E^+(\epsilon)$. The quasi-HP property then implies that there is a $\delta = \delta(\epsilon)$ such that if $|z - p| < \delta$ and z is between C and the ray $\arg\{z\} = \alpha$, then $|\theta(z) - \gamma| < 2K\epsilon$. Since ϵ is arbitrary, $\theta(z) \rightarrow \gamma$ as $z \rightarrow p$ between C and the ray. It is also clear that for sufficiently small ϵ the translate of E^- down to p is a nontrivial initial arc of a j -characteristic C' which exits regularly at p and which forms with ∂G an acute angle of size $\frac{\pi}{2} - \gamma$. What we have shown in regard to C now implies that if $\beta < \frac{\pi}{2} - \gamma$, then $\theta(z) + \frac{\pi}{2} \rightarrow \gamma + \frac{\pi}{2}$ as $z \rightarrow p$ between C' and the ray $\arg\{z\} = \pi - \beta$, so that $\theta(z) \rightarrow \gamma$ as $z \rightarrow p$ in that curvilinear sector. Finally, it is easy to see that $\theta(z) \rightarrow \gamma$ as $z \rightarrow p$ between C and C' . Since $\alpha < \gamma$ and $\beta < \frac{\pi}{2} - \gamma$ are arbitrary, we have the desired regularity at p . The case $\gamma \in (\frac{\pi}{2}, \pi)$ is handled in the same manner apart from minor changes of a notational nature. In the case that θ arises from the solution of a normal system it is clear from (ii) of Definition 1.1 that both $R_i(z)$ and $R_j(z)$ have limits as $z \rightarrow p$ in the sector $\alpha < \arg\{z\} < \beta$, so that the second conclusion is valid in this case.

(ii) $\gamma = \frac{\pi}{2}$. Here again we can take C to be an arc of an elementary i -characteristic exiting at p and proceed as in case (i). For any $\alpha < \frac{\pi}{2}$ it follows as in case (i) that $\theta(z) \rightarrow \gamma$ as $z \rightarrow p$ between C and the ray $\arg\{z\} = \alpha$. Again it is immediate that the same holds in the curvilinear sector between C and the ray $\arg\{z\} = \pi - \alpha$. The second conclusion is likewise immediate.

(iii) $\gamma \in \{0, \pi\}$. We deal with the case $\gamma = 0$, the case $\gamma = \pi$ being essentially the same. Here we can define $E^-(\epsilon)$ as in case (i), but because there is no fan at p the translate of $E^-(\epsilon)$ down to p is a nontrivial arc of a j -characteristic which is orthogonal to ∂G at p . This puts us in case (ii), so we are done. \square

DEFINITION 3.26.

The following is an immediate consequence of Proposition 3.15.

PROPOSITION 3.27.

PROPOSITION 3.28. Let C be any elementary i -characteristic, and let D be the characteristic subdomain bounded by C and the arc of ∂G containing p . Let $E = ef = \min_D(\{p\})$.

If $E = p$, we are done, so that $E = ef$ is a nontrivial extended characteristic. If $p \in E$, Proposition 3.25 implies that (i) $p \in \{e, f\}$, since extended characteristics exit regularly at their proper contact points. If $p \notin E$, the minimality of E implies that (ii) E is an elementary i -characteristic. We deal with possibility (i) first. Let E be parameterized by $z(s)$, $0 \leq s \leq l$, with $z(0) = p$. Since p is a singularity of type 2, for no $\sigma > 0$ is it true that $z((0, \sigma)) \subset G$. On the other hand, it follows from the regularity of ∂G and Proposition 3.25 that for no $\sigma > 0$ is it true that $z((0, \sigma)) \subset \partial G$. It is then clear that E contains a sequence $\{C_k\}$ of subarcs which are elementary i -characteristics for which $\text{diam}(C_k)$ and $\text{dist}(C_k, p)$ tend to 0 as $k \rightarrow \infty$, so that p is indeed an i -singularity. We finish by showing that possibility (ii) cannot occur; to do so we proceed as in the proof of Proposition 3.17. Let w_0 be any point of E other than e or f , and consider a small j -arc $J \subset G$, which extends from w_0 to the right of E (i.e., into D). Let J be parameterized by $w(s)$, $0 \leq s \leq \lambda_0$, with $w(\lambda_0) = w_0$, and let $C(s) = l(s)r(s)$ be the elementary i -characteristic through $w(s)$. As in the the proof of Proposition 3.17, it follows from the minimality of E that either $l(s) > p$ for all $s \in (0, \lambda_0)$ or $r(s) < p$ for all $s \in (0, \lambda_0)$. Assume, for definiteness, that the latter occurs. By considering what happens as $s \rightarrow 0$ we obtain, as in that proof, a nonmonotone extended i -characteristic ab (with $a < b$ with respect to the order on $\text{bot}(D)$) which contains ef as a proper subarc, for which $b < p$ and for which $bf \preceq ef$. But $p \notin bf$ by hypothesis, so that bf must have a subarc gh which is an elementary i -characteristic for which $g < p < h$, which contradicts the minimality of E . \square

The following is an immediate consequence of Propositions 3.25, 3.27, and 3.28 and Definitions 3.21 and 3.26.

PROPOSITION 3.29. Let $p \in \partial G$ be a type 2 singularity. Let C be an elementary i -characteristic through p .

- (i) $p \in \text{int}(C)$, $0 < s < l$, $z(s) \in C$, $k = 1, 2$
- (ii) $p \in \text{int}(C)$, $1 < s < l$, $z(s) \in C$, $k = 1, 2$
- (iii) $p \in \text{int}(C)$, $2 < s < l$, $z(s) \in C$, $k = 1, 2$

COMMENT 3.30. Let $\alpha < \pi/2$. Let C be an elementary i -characteristic through p .

PROPOSITION 3.31 (type 1 singularities with fans). Let C be an elementary i -characteristic through $p \in \partial G$. Let $q \in G$. Let $z(s)$, $0 \leq s \leq l$, be a curve in G with $z(0) = p$ and $\eta < \frac{\pi}{2}$.

$$(3.22) \quad \phi + \eta \leq \arg\{z'(s)\} \leq \phi + \pi - \eta, \quad 0 < s \leq l,$$

where $\phi = \lim_{s \rightarrow 0} \arg\{z'(s)\}$ and $\eta = \arg\{z - p\}$.

For $0 < \xi < \pi$, let $R(p, \xi)$ denote the (open) ray

$$R(p, \xi) = \{z : \arg\{z - p\} = \xi\}.$$

Without loss of generality we can assume that $\phi = 0$. It follows immediately from hypothesis (3.22) that C lies in the sector $\{z : \eta \leq \arg\{z - p\} \leq \pi - \eta\}$. The hypothesis that $\lim_{s \rightarrow 0} \arg\{z'(s)\}$ does not exist implies that there is some δ , $0 < \delta \leq \frac{\eta}{K+1}$, such

that for arbitrarily small $\sigma_1 > \sigma_2 > 0$

$$(3.23) \quad |\arg\{z'(\sigma_1)\} - \arg\{z'(\sigma_2)\}| \geq 2\delta.$$

For $0 < s \leq l$, let $F(s) = a(s)b(s)$ be the elementary j -characteristic through $z(s)$. It follows from the last sentence of Proposition 3.15 that $a(l) \neq p \neq b(l)$. Let D be the j -characteristic subdomain bounded by $F(l)$ and the arc of ∂G joining $a(l)$ to $b(l)$ and containing p . Clearly, $F(s) \cap C \cap G = \{z(s)\}$, and $a(s)$ is nonincreasing and $b(s)$ is nondecreasing with respect to the order on $\text{bot}(D)$. Since, by Proposition 3.15, $\text{diam}(F(s)) \rightarrow 0$ as $s \rightarrow 0$, by replacing C by a sufficiently short initial subarc we can assume that

$$(3.24) \quad R(p, \xi) \cap \partial D \subset F(l), \quad \frac{\delta}{2} \leq \xi \leq \pi - \frac{\delta}{2}.$$

Now consider any pair of numbers $\sigma_1 \neq \sigma_2$ in $(0, l)$ for which (3.23) holds. Obviously, for at least one of them, call it σ , there must hold $|\arg\{z'(\sigma)\} - \frac{\pi}{2}| \geq \delta$. Then $F(\sigma)$ must have a subarc $W = uv$ such that $\lambda(\theta(W)) = \delta$ and $z(\sigma) \in \{u, v\}$. To see this, say, for example, that $\arg\{z'(\sigma)\} \leq \frac{\pi}{2} - \delta$. If, along the arc of $F(\sigma)$ emanating to the left of C , the argument of the tangent were always within δ of its argument at $z(\sigma)$, then $F(\sigma)$ would never cross $R(p, \pi - \delta)$, and it would therefore not be able to exit D . An analogous argument may be used in the case that $\arg\{z'(\sigma)\} \geq \frac{\pi}{2} + \delta$.

From this it follows that there is a sequence $s_k \rightarrow 0$ of such numbers σ for all of which the corresponding j -arc $W_k = u_k v_k$ lies on one side of C or the other. To be specific, say $u_k = z(s_k)$ and W_k lies to the right of C . It follows from the quasi-HP property that if $W' = z(s)v$ is any translate of W_k along C in the direction of increasing s , then

$$\delta/K \leq \lambda(\theta(W')) \leq K\delta.$$

From this and the fact that $\delta \leq \frac{\eta}{K+1}$ it follows in turn that for $s > s_k$, along any i -characteristic arc whose initial point is in W_k and which is parallel to the subarc $z(s_k)z(s)$ of C , the inclination of the tangent is in the interval

$$[\eta - K\delta, \pi - \eta + K\delta] \subset [(K + 1)\delta - K\delta, \pi - (K + 1)\delta + K\delta] = [\delta, \pi - \delta].$$

For each k let

$$S_k = \{R(z, \xi) : z \in W_k \text{ and } \xi \in [\delta, \pi - \delta]\}.$$

Now $\text{diam}(W_k) \leq \text{diam}(F(s_k))$ and the latter tends to 0 as $k \rightarrow \infty$, so that in light of (3.24) by eliminating a finite number of elements of $\{s_k\}$ and renaming we can assume that all the $R(z, \xi)$ in each of the S_k meet ∂D a point of F before exiting G . From this it follows that the translates of $z(s_k)z(l)$ along W_k all lie in G . This means that for any $s \in [s_k, l]$ the complete translate of W_k along C from $z(s_k)$ up to $z(s)$ belongs to G . For the translate T_k of W_k up to $z(l)$ we have that $z(l) \in T_k \subset F(l)$, and by the quasi-HP property $\lambda(\theta(T_k)) \geq \delta/K$. This bound implies that there is a positive lower bound on the length of the T_k . Thus $\cap\{T_k : k \geq 1\} = T$ is a nontrivial arc of $F(l)$, one of whose endpoints is $z(l)$. However, the translate of T down to $z(s_k)$ is contained in $W_k \subset F(s_k)$. Since $\text{dist}(F(s), \{p\}) \rightarrow 0$ as $s \rightarrow 0$, it follows that all i -characteristics through T exit at p , and this establishes the existence of the desired fan. \square

We need the following elementary lemma.

LEMMA 3.32. Let f be a measurable function on $[0, T]$ such that $\int_0^T f(x)dx = A > 0$ and let $\rho \in [0, \frac{1}{2})$.

$$\lambda(\{\xi : f(\xi) \geq \rho \frac{A}{T} \text{ and } \int_0^\xi f(x) dx \geq \rho A\}) > 0.$$

Let $f(x) \leq M$ a.e. on $[0, T]$, and let $Q = \{\xi : f(\xi) \geq \rho \frac{A}{T}\}$. Obviously, $\lambda(Q) > 0$. Let Q^* be the set of density points of Q , so that $\lambda(Q^*) = \lambda(Q) > 0$. Let $\xi_0 = \sup\{\xi : \xi \in Q^*\}$. If the set $Q^* \cap \{\xi : \int_0^\xi f(x) dx \geq \rho A\}$ had measure 0, then for each $\delta > 0$ there would be a point $\xi \in (\xi_0 - \delta, \xi_0)$ for which $\int_0^\xi f(x) dx < \rho A$, so that

$$\begin{aligned} A &= \int_0^\xi f(x)dx + \int_\xi^{\xi_0} f(x)dx + \int_{\xi_0}^T f(x)dx < \rho A + \delta M + (T - \xi_0) \frac{\rho A}{T} \\ &< 2\rho A + \delta M < A, \end{aligned}$$

for δ sufficiently small, so that $\lambda(Q^* \cap \{\xi : \int_0^\xi f(x) dx \geq \rho A\})$ must be positive. □

The following proposition plays a fundamental role in the proof of Main Theorem 2.2.

PROPOSITION 3.33 (essentially singularity-free boundary arcs).

Let $\bar{\eta} = \bar{\eta}(K, \rho) < \rho$ and let $G \in \mathcal{G}(\rho)$. Let $\theta \in \text{HP}(G, K)$ and let C be a characteristic of G with $\theta \in C$. Let $p \in \partial G$ and let $E = pp'$ be a subarc of ∂G with $p' \in C$. Let D be a subarc of ∂G with $p \in D$ and $b \in A \cap C$. Let F be a subarc of ∂G with $b \in F$ and $\text{dist}(p, F) \geq \rho$. Let $[p, p'']$ be a subarc of A with $|p'' - p| = \bar{\eta}$.

Without loss of generality we may assume that $p = 0$, that the positively oriented tangent to ∂G at p has the direction of the positive real axis, and that A extends to the right of p . Let C be parameterized by $z(s) = x(s) + iy(s)$, with $z(0) = p$. We shall introduce positive constants, which will depend only on K and ρ . In most cases these constants will be denoted by the same symbol B , but the use of a single symbol to denote different constants will not cause confusion or be misleading since any statement in which B plays a role will be valid if B is taken to be any suitably large number and since this convention is used only finitely many times. The symbol B_1 , on the other hand, will refer to a specific constant, which again depends only on K and ρ and which will have the same value every time it is used.

If we give A in nonparametric form by $y = g(x)$, then

$$(3.25) \quad |g(x)| \leq \frac{x^2}{\rho} \quad \text{and} \quad |g'(x)| \leq \frac{2x}{\rho} \quad \text{for} \quad 0 \leq x \leq \frac{1}{10\rho}.$$

Clearly, Proposition 3.12 holds for extended characteristics with the obvious wording changes. Any j -arc emanating to the left of C from any point z of the subarc pb of C will exit G at a point of F , so that the length of any such j -arc is at least $\rho - |p - z|$. Since $\text{dist}(z(s), \partial G) \leq |z(s)| \leq |s|$, it follows from Proposition 3.12 that $|\arg\{z'(s)\}| \leq \frac{1}{2}$ for $0 \leq s \leq \frac{1}{B}$, so that for these values of s we have

$$\frac{1}{B} \leq x'(s) \leq 1 \quad \text{and} \quad |y'(s)| \leq 1,$$

and we can represent C nonparametrically by $y = f(x)$, where $f(x(s)) = y(s)$. For $0 \leq s \leq \frac{1}{B}$ we have

$$\begin{aligned} |f'(x(s))| &= |\tan(\arg\{z'(s)\})| \leq 2|\arg\{z'(s)\}| \leq B\sqrt{\text{dist}(z(s), A)} \\ &\leq B\sqrt{|z(s) - x(s)| + \text{dist}(x(s), A)} \leq B\sqrt{f(x(s)) + x(s)^2}, \end{aligned}$$

where the second inequality follows from Proposition 3.12 and the fourth inequality follows from the first bound of (3.25) in the form $\text{dist}(x(s), A) \leq \frac{(x(s))^2}{\rho}$. Thus there is a $B_1 \geq 1$ such that

$$(3.26) \quad |f'(x)| \leq B_1\sqrt{|f(x)| + x^2} \quad \text{and} \quad |f(x)| \leq B_1x \quad \text{for } 0 \leq x \leq \frac{1}{B_1}.$$

For $0 \leq x \leq \frac{1}{B_1}$ we then have the following. In the first place,

$$|f'(x)| \leq B_1\sqrt{B_1x + x^2} \leq 2^{1/2}B_1^{3/2}x^{1/2},$$

so that in fact $|f(x)| \leq \frac{2}{3}2^{1/2}B_1^{3/2}x^{\frac{3}{2}} < B_1^{3/2}x^{\frac{3}{2}}$. Repeating this argument we have that $|f'(x)| \leq 2^{1/2}B_1^{7/4}x^{\frac{3}{4}}$, so that in fact $|f(x)| \leq \frac{4}{7}2^{1/2}B_1^{\frac{7}{4}}x^{\frac{7}{4}} < B_1^{\frac{7}{4}}x^{\frac{7}{4}}$, so that $|f'(x)| \leq 2^{1/2}B_1^{\frac{15}{8}}x^{\frac{7}{8}}$, and so on. Thus, in fact,

$$(3.27) \quad |f(x)| \leq B_1^2x^2 \quad \text{for } 0 \leq x \leq \frac{1}{B_1}$$

and

$$(3.28) \quad |f'(x)| \leq \sqrt{2}B_1^2x \quad \text{for } 0 \leq x \leq \frac{1}{B_1}.$$

Before proceeding we remind the reader that each time the symbol B appears, it may have a value larger than it did at its previous appearance. For any $t \in [0, \frac{1}{2B}]$, let $J(t)$ be the j -half-characteristic emanating to the right of C (that is, downwards) from $t + if(t) = \alpha(t)$ and joining it to a point of A . Note that $J(t)$ will reduce to a single point if $\alpha(t) \in A$. Let $J(t)$ be given parametrically by $\zeta_t(s)$ with $\zeta_t(0) = \alpha(t)$. Let $t_0 = x_0^2$. Then, by (3.25) and (3.27), $|f(t_0) - g(t_0)| \leq Bx_0^4$ for $0 < x_0 \leq \frac{1}{B}$. For notational convenience we can assume without loss of generality that $\arg\{z'(s)\} = \theta(z(s))$. We show that, with an appropriately large value of B , $0 < x_0 \leq \frac{1}{B}$ implies that

$$(3.29) \quad |\theta(\zeta_{t_0}(s_1)) - \theta(\zeta_{t_0}(0))| \leq x_0 = \sqrt{t_0}$$

for every $\zeta_{t_0}(s_1) \in J(t_0)$. If this were not true, there would be an $s_1 \in (0, \lambda(J(t_0)))$ such that $|\theta(\zeta_{t_0}(s_1)) - \theta(\zeta_{t_0}(0))| = x_0$. Let s_2 be the smallest number in $[0, s_1]$ for which $|\theta(\zeta_{t_0}(s_2)) - \theta(\zeta_{t_0}(0))| = x_0$, so that

$$(3.30) \quad \lambda(\theta(\zeta_{t_0}([0, s_2]))) \leq 2x_0.$$

By increasing B if necessary, $t_0 = x_0^2 < \frac{1}{B^2}$ will be so small that

$$(3.31) \quad s_2 \leq 2|f(t_0) - g(t_0)| \leq Bx_0^4.$$

By the preceding lemma (with $\rho = \frac{1}{3}$, $T = s_2$, $A = x_0$) it then follows that there is an $s_3 \in (0, s_2)$ for which

$$|D\theta(\zeta_{t_0}(s_3))| \geq \frac{x_0}{3s_2} \geq \frac{1}{Bx_0^3}$$

and

$$(3.32) \quad |\theta(\zeta_{t_0}(s_3)) - \theta(\zeta_{t_0}(0))| \geq \frac{x_0}{3},$$

where the expressions inside the absolute values on the left-hand side of these last two bounds have the same sign. Let R be the i -half-characteristic through $\zeta_{t_0}(s_3)$ emanating from $J(t)$ towards its concave side (that is, to the left of $J(t_0)$ if $d \arg\{\zeta'_{t_0}(s)\}/ds$ is positive at s_3 and to the right if it is negative). It follows from the characteristic length bound (Proposition 2.5) that

$$(3.33) \quad \lambda(R) \leq Bx_0^3.$$

We concentrate on the case that R emanates to the left of $J(t_0)$; the opposite case is much easier to handle, as we indicate below. By (3.27), (3.31), and (3.33), R exits G at a point $g(\beta)$ with

$$(3.34) \quad \begin{aligned} \beta &\leq t_0 + |f(t_0)| + |\zeta_{t_0}(0) - \zeta_{t_0}(s_3)| + \lambda(R) \\ &\leq x_0^2 + O(x_0^3) \leq 2x_0^2 \quad \text{for } x_0 \leq \frac{1}{B}. \end{aligned}$$

For $t \geq t_0$ let $\bar{J}(t)$ be the j -arc joining $\alpha(t) \in C$ to a point $r(t) \in R$, so that $\bar{J}(t_0) = \zeta_{t_0}([0, s_3])$. Obviously, $\bar{J}(t)$ is not defined for all $t \in [t_0, \frac{1}{2B_1}]$, but, for any t for which it is defined, $\bar{J}(t_0)$ and $\bar{J}(t)$ are the j -sides of a characteristic quadrilateral $Q(t)$, whose i -sides are the arc $\alpha(t_0)\alpha(t)$ of C and the arc $\zeta_{t_0}(s_3)r(t)$ of R . From the K -quasi-HP property and (3.30) it follows that

$$(3.35) \quad \lambda(\theta(\bar{J}(t))) \leq K\lambda(\theta(\bar{J}(t_0))) \leq K\lambda(\theta(\zeta_{t_0}([0, s_2]))) \leq 2Kx_0,$$

and by (3.28)

$$(3.36) \quad \lambda(\theta(\zeta_{t_0}(s_3)r(t))) \leq K\lambda(\theta(\alpha(t_0)\alpha(t))) \leq 2K \tan^{-1}(Bt) \leq Bt.$$

But for any $t \leq 10x_0^2$ for which $J(t)$ is defined, we have from (3.25) and (3.28) that

$$\begin{aligned} |\theta(\zeta_t(s))| &\leq \theta(\alpha(t)) + 2Kx_0 \leq \tan^{-1}(Bt) + 2Kx_0 \\ &\leq 20Bx_0^2 + 2Kx_0 \leq \frac{1}{100} \end{aligned}$$

for $x_0 \leq \frac{1}{B}$. From this together with (3.34) it is clear that $g(\beta)$ is the endpoint of $\bar{J}(t_1)$ for some $t_1 \in [x_0^2, 10x_0^2]$. It follows from the quasi-HP property that R , when oriented from $\zeta_{t_0}(s_3)$ to $r(t_1)$, has a well-defined (one-sided) tangent, whose argument we denote by ϕ_0 . But then we have from (3.28), (3.32), (3.36), and the quasi-HP property that

$$\phi_0 \geq \frac{x_0}{3} - |\theta(\alpha(t_0))| - Bt_1 \geq \frac{x_0}{2} - \theta(\alpha(x_0^2)) - 10Bx_0^2 \geq \frac{x_0}{3} - Bx_0^2,$$

so that

$$(3.37) \quad \phi_0 \geq \frac{x_0}{4},$$

for $x_0 \leq \frac{1}{B}$. By (3.25) and (3.24)

$$|\tan^{-1}(g'(\beta))| \leq |g'(\beta)| \leq \frac{2\beta}{\rho} \leq \frac{4}{\rho}x_0^2 \quad \text{for } x_0 \leq \frac{1}{B}.$$

Taking into account the direction from which R crosses A at $g(\beta)$, one sees that $\phi_0 \leq |\tan^{-1}(g'(\beta))|$, which is a contradiction for $x_0 < \min\{\frac{1}{B}, \frac{\rho}{17}\}$. If R emanates to the right of $J(t_0)$, we get the same contradiction more easily because in that case we move back along C towards p and there are smaller bounds throughout. This establishes (3.29) for $x_0 \leq \frac{1}{B}$ (with, of course, an appropriately large B).

Thus, for $\sqrt{t} \leq \frac{1}{B}$, $\lambda(\theta(J(t))) \leq \sqrt{t}$. From this together with (3.27) it follows that

$$(3.38) \quad |\theta(z)| \leq \sqrt{t} + Bt \leq \frac{1}{100} \quad \text{for } z \in J_t, \quad 0 < t < \frac{1}{B}.$$

Finally, after increasing B if necessary, we have $|\tan^{-1}(g'(t))| \leq \frac{1}{100}$ for $|t| \leq \frac{1}{B}$. Let $l_t = \lambda(J_t)$, and let $q(t) = \zeta_t(l_t) \in A$. We now show that $\{q(t) : 0 \leq t \leq \frac{1}{B}\}$ is an arc of A , for which we have to only show that $q(t)$ is continuous on $(0, \frac{1}{B})$. Clearly, we have from (3.38)

$$(3.39) \quad \left| \arg \left\{ \frac{\partial \zeta_t(s)}{\partial s} \right\} + \frac{\pi}{2} \right| \leq \frac{1}{100} \quad \text{for } t \leq \frac{1}{B} \quad \text{and } 0 \leq s \leq \lambda(t).$$

If $l_{t_0} = 0$, then what we have already shown implies that $l_t \rightarrow 0$ as $t \rightarrow t_0$, from which it immediately follows that $q(t) \rightarrow q(t_0)$ as $t \rightarrow t_0$. Thus consider a fixed t_0 for which $l_{t_0} > 0$. Let $\delta > 0$ be small, and let $s \in (l_{t_0} - \delta, l_{t_0})$. Let the elementary i -characteristic through $\zeta_{t_0}(s)$ be parameterized by $z = w_s(\sigma)$ with $w_s(0) = \zeta_{t_0}(s)$. It is clear that there is some $\delta' \in (0, \delta]$ such that for all $\sigma \in [-\delta', \delta']$, $w_s(\sigma) \in J(t(\sigma))$ (where $t(\sigma)$ is some continuous function of σ) and $|\theta(w_s(\sigma)) - \theta(w_s(0))| \leq \frac{1}{100}$. Then it is clear from (3.39) that for $|\sigma| \leq \delta'$, $J(t(\sigma))$ intersects A at a point within 2δ of $q(t_0)$. The desired continuity follows immediately since there is a positive δ'' such that for $|t - t_0| < \delta''$, $J(t)$ intersects $w_s((-\delta', \delta'))$.

From this and (3.38) it follows that, in fact, for $0 < t < \frac{1}{B}$, $|\theta(z)| \leq 2\sqrt{t}$ for z in the domain bounded by J_t and the arcs of C and A with endpoints p and $q(t)$. The desired conclusion now follows immediately from Propositions 3.24, 3.25, and 3.31. \square

We end this section with a technical proposition to be used in the proof of the main theorem given in section 4. The numbers $\eta = \eta(K)$ and $\eta' = \eta'(K)$ arising here are those of Proposition 2.8.

PROPOSITION 3.34 (trapped area bound). $\eta_1 = \eta_1(K) \leq \min\{1, \eta(K)\}$. Let (D, V) be a domain with $U = u_1 u_2 \in \mathcal{I}(D)$ and $|p - q| \geq \xi$, $p, q \in \text{bot}(D)$, $[u_1, u_2] \subset [p, q]$, $U \subset D$, $U' \subset \overline{D}$, $U' \subset V$, $u, u', u' < u_1, u' > u_2, u = u_1, u = u_2$, $|\Delta\theta(U')| \geq \frac{\pi}{4}$, $\eta_1 \xi^2$.

For definiteness we assume that $u = u_1$, so that $u' < u_1$. Let $V = v_1 v_2$. Let U' be parameterized by $w = w(s)$, $0 \leq s \leq L$, with $w(0) = u_1$, and for $s \in [0, L)$ denote by $U(s) \preceq V$ the elementary i -characteristic containing $w(s)$, so that in particular $U(0) = U$. One of the endpoints $a(s)$ of $U(s)$ lies below U' , and the other endpoint $b(s)$ is in $[v_1, v_2] \setminus (u_1, u_2)$. There are the following two possibilities.

A. The point $b(s) \geq u_2$ (i.e., lies to the right of u_2) for $0 \leq s < L$. In this case $U \preceq U(s)$, and consequently $\lambda(U(s)) \geq \xi$, for all $s \in [0, L)$, so that by second

area bound (Proposition 2.8(ii)) the area of $\cup\{U(s) : 0 < s < L\} \subset D'$ is at least $\eta' \frac{\pi}{4} \xi^2$.

B. There is some $s \in (0, L)$ for which $b(s) \leq u_1$. Since $b(0) = u_2$, there is a decreasing sequence $\{s_k\}$ tending to some $\sigma \in [0, L)$ such that $b(\sigma) \geq u_2$, but $b(s_k) \leq u_1$ for $k \geq 1$. But then there is a nonmonotone (with respect to D) extended i -characteristic C joining $a_0 = \lim_{k \rightarrow \infty} a(s_k)$ to a point $b_0 = \lim_{k \rightarrow \infty} b(s_k) \leq a_0 \leq u_1$ for which $b(\sigma)$ is a proper contact point. Then C has subarcs $U(\sigma) = a(\sigma)b(\sigma)$ and $W = de$, which are elementary i -characteristics and for which $U(\sigma) \preceq W$ and $U(\sigma) \cap W \subset \{b(\sigma)\}$, and where $d \in [v_1, u']$, $e \in [b(\sigma), v_2]$. Since we therefore have $[u_1, u_2] \subset [d, e]$, we conclude that $\lambda(W) \geq |e - d| \geq \xi$. Regard W as oriented from e to d and consider the subarc W' of W of length $\frac{\xi}{2}$ which starts at e , and for each $p \in W'$ consider the j -arc $J(p)$ emanating to the left of W and which joins p to a point $q(p) \in \text{bot}(D)$. Because $J(p)$ cannot intersect C at any interior point of C other than p , we have $d \leq q(p) \leq u_1$. Since the length of any curve joining e to a point to the left of u_1 must be at least ξ , it follows that $\lambda(J(p)) \geq \frac{\xi}{2}$. The first area bound (Proposition 2.8(i)) then implies that

$$\mu\left(\bigcup\{J(p) : p \in W'\}\right) \geq \eta \frac{\xi^2}{4}.$$

The conclusions of Cases A and B together imply that the area between U and V is at least $\min\{\eta \frac{\xi^2}{4}, \eta' \frac{\pi}{4} \xi^2\} = \eta_1 \xi^2$, where $\eta_1 = \min\{\frac{\eta}{4}, \frac{\pi}{4} \eta'\}$. \square

4. Proof of Main Theorem 2.2. The main ingredients in the proof are Proposition 3.33 about essentially singularity-free boundary arcs, Proposition 3.29, Proposition 3.17 regarding the structure of $\min_D(A \cup B)$, and the trapped area bound (Proposition 3.34). We also use the preservation of quasi-HP conditions under linear changes of variable, which allows one to normalize the arc one is working with to have any convenient length. In reference to a quasi-HP function θ on $G \in \mathcal{G}$, we will say that $E \subset \partial G$ is *ESF- i* (abbreviated ESF- i) if it has at most countably many i -singularities of θ . As is explained in detail in the final paragraph of the proof, it is enough to prove that there is a number $\tau = \tau(K) < 1$ such that on any suitably small boundary arc the set of all i -singularities has Hausdorff dimension at most τ . In the treatment of such arcs there is considerable freedom in the choice of the various explicitly given numerical constants that we use, and for the most part they are far bigger (or smaller) than necessary and have most often been chosen either for the sake of convenience or to avoid the necessity of going into careful geometric arguments. In the same vein, many of the bounds we give are far from being sharp, so that, if we state that the value of some expression is bounded below by 1, it may be clear that it is in fact bounded below by a considerably larger number. Furthermore, in some instances we use symbols to denote numerical constants which could without much effort be determined explicitly.

Let $G \in \mathcal{G}(\rho)$, so that the unsigned curvature of ∂G is everywhere at most $\frac{1}{\rho}$. For $p \in \partial G$ let $w(p) = e^{i\phi(p)}$ be the unit tangent to ∂G (with positive orientation) at p . Let δ be any positive for which

$$(1) \delta \leq \frac{\rho}{10000}.$$

This implies in addition that

$$(2) \text{ on any arc } B \text{ of } \partial G \text{ of length } 100\delta, \lambda(w(B)) < \frac{1}{100}, \text{ and that}$$

(3) for any point $p \in \partial G$ and any $r \leq 20\delta$, $\partial N(p, r)$ meets ∂G in two points, so that for such r , $N(p, r) \cap G$ is essentially a semidisk.

The corners of $\partial N(p, 9\delta) \cap G$ can be replaced by circular arcs of radius $\frac{\delta}{90}$ and then smoothed at the joins to form a C^2 curve which is the boundary of a subdomain

$$\Sigma = \Sigma(p, \delta) \subset N(p, 10\delta) \cap G$$

of G with the following properties:

(1a) $\Sigma \in \mathcal{G}(\frac{\delta}{100})$.

(2a) The arc $N(p, 4\delta) \cap \partial\Sigma$, to be referred to as the bottom of Σ , and which is an arc of ∂G , has unsigned curvature bounded above by $\frac{1}{10000\delta}$.

(3a) $\mu(\Sigma) \leq 200\delta^2$.

Now consider the $X = X(p) = \frac{1}{\delta}(\Sigma(p, \delta) - p)$ which has the following properties:

(1b) $X \in \mathcal{G}(\frac{1}{100})$.

(2b) $B = N(p, 4) \cap \partial X$, to be referred as the bottom of X , is a C^2 arc on which the unsigned curvature is bounded above by $\frac{1}{10000}$.

(3b) $\mu(X) \leq 200$.

Note that what we call the bottom B of X is very close to a straight line segment of length 8 centered at the “midpoint” p of what is close to being the straight portion of the boundary of a semidisk of radius 9. Note also that by Proposition 2.3, $\theta(p + \delta z)$ is a K -quasi-HP function on this $X(p)$. A domain X having these three properties will be called a *quasi-HP domain*.

Until the end of the proof of the main theorem we work exclusively with quasi-HP functions θ defined on a normalized domain X . It is clear from property (1b) that there is a universal constant $\gamma_0 \in (0, 1)$ such that for any such X

$$(4.1) \quad |z_1 - z_2| \geq \gamma_0 \text{dist}_{\text{arc}}(z_1, z_2) \quad \text{for } z_1, z_2 \in \partial X,$$

where $\text{dist}_{\text{arc}}(z_1, z_2)$ denotes the length of the shorter arc of ∂X with endpoints z_1, z_2 . We shall prove the following proposition from which, as we will subsequently show, the main theorem follows almost immediately. We stress that for an arc F of the bottom of X , $\text{diam}(F)$ is the distance between its endpoints.

PROPOSITION 4.1. *Let $\Delta_0 = \Delta_0(K) > 0$ be a constant depending only on K . Let X be a normalized domain with $\mu(X) \leq \Delta_0$ and θ a K -quasi-HP function on X . Then there exists a constant $\epsilon_0 = \epsilon_0(K, \Delta_0) > 0$ such that if $\epsilon < \epsilon_0$ and θ is ϵ -essentially singular on X , then θ is ϵ_0 -essentially singular on X .*

It is clear from Proposition 2.3 that we can assume that $p = 0$ and that $w(p) = 1$. We establish the desired conclusion as follows (once again, the constants chosen are unnecessarily big or small). Let H_0 be an integer for which

$$H_0 \geq \max \left\{ \frac{10^6}{\gamma_0^2 \eta_1}, 10000 \right\},$$

and let

$$\epsilon_0 = \bar{\eta} \left(K, \frac{\gamma_0}{10H_0} \right),$$

where $\eta_1 = \eta_1(K)$ is the constant of the trapped area bound (Proposition 3.34) and $\bar{\eta}$ is the constant of Proposition 3.33 about essentially singularity-free boundary arcs; H_0, η_1 , and ϵ_0 depend solely on K . In particular we have

$$\epsilon_0 \leq \frac{\gamma_0}{10H_0} < \frac{1}{10H_0} \leq \frac{1}{10^5}.$$

Our strategy will be to show that there is some constant $\gamma_1 = \gamma_1(K)$ such that the assumption that

$$(4.2) \quad \text{the bottom } B \text{ of } X \text{ has no ESF-}i \text{ arc of diameter } \epsilon_0$$

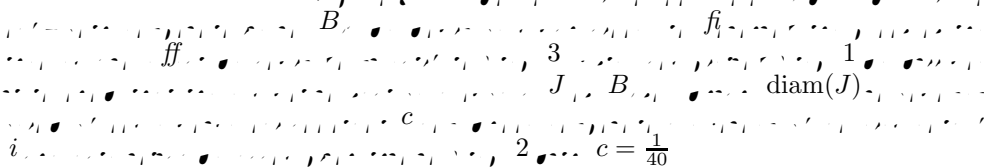
leads to the conclusion that the bottom of X has an ESF- i arc of diameter $\frac{\epsilon_0}{10^6}$ or one of diameter γ_1 , so that B has an ESF- i arc of diameter $\Delta_0 = \min\{\frac{\epsilon_0}{10^6}, \gamma_1\}$. This underlying assumption (4.2) will be in force throughout and will be used many times.

The proof the existence of such a Δ_0 is divided into three steps.

1. There is an elementary i -characteristic $C_0 = \bar{a}\bar{b}$ which has the following properties:

- (i) $|\bar{b} - \bar{a}| \geq 1$.
- (ii) At least one of \bar{a}, \bar{b} lies in B .

COMMENT 4.2.



We begin the proof of Step 1 by noting that by (4.2) there is an i -singularity on B within ϵ_0 of the left endpoint of B , so that there is an elementary i -characteristic C'_0 which joins two points of B , both of which are within $2\epsilon_0$ of the left endpoint of B . We work with the right half B_R of B ; the endpoints of B_R are 0 and a point d to the right of 0, with $|d| = 4$. In light of our assumption (4.2), it follows that there are two elementary i -arcs $F_1 = a_1a_2$ and $F_2 = b_1b_2$ joining points of B_R such that

$$|a_1 - a_2|, \quad |b_1 - b_2| < \frac{1}{H_0}$$

and

$$\frac{1}{H_0} < \text{dist}(F_1, 0), \quad \text{dist}(F_2, d) < \frac{2}{H_0}.$$

Consider the characteristic subdomain (D, C'_0) for which $\text{bot}(D)$ contains the $\partial X \setminus B$, so that the only part of X not in D is the union of C'_0 and the interior of the (small) subdomain bounded by C'_0 and the arc of B which joins its endpoints.

We apply Proposition 3.17 to $\min_D(F_1 \cup F_2)$ with $A = F_1$ and $B = F_2$ and examine in turn each of the three cases of its conclusion.

(i). If there were a contact point p of $\min_D(F_1 \cup F_2)$ between F_1 and F_2 for which

$$\text{dist}(F_1, p), \quad \text{dist}(F_2, p) > \frac{1}{H_0},$$

then Proposition 3.33 would imply that there is an ESF- i subarc of B_R of diameter ϵ_0 , contrary to (4.2). Thus, for every contact point p of $\min_D(F_1 \cup F_2)$ between F_1 and F_2 , one of

$$\text{dist}(F_1, p) \leq \frac{1}{H_0} \quad \text{or} \quad \text{dist}(F_2, p) \leq \frac{1}{H_0}$$

must hold. Say for definiteness that the former holds for some contact point p between F_1 and F_2 , and let \bar{a} be the rightmost such point. Then it is clear that $\min_D(F_1 \cup F_2)$ has an elementary subcharacteristic C_0 joining $\bar{a} \in B$ to a point $\bar{b} \in \text{bot}(D)$, where \bar{b} is either between F_1 and F_2 but within $\frac{1}{H_0}$ of F_2 or to the right of F_2 (with respect to the order on $\text{bot}(D)$). In either event, $|\bar{b} - \bar{a}| \geq 1$. The case $\text{dist}(F_2, p) \leq \frac{1}{H_0}$ is handled similarly.

For the other two cases we let a, b, e, f, U , and L be as in the conclusion of Proposition 3.17. Furthermore, let the left and right endpoints of C'_0 be l and r (with respect to the order on $\text{bot}(D)$), so that since $\text{bot}(D)$ is very close to a segment of \mathbb{R} , l lies to the right of r in the usual sense).

(ii). Let α and β be the left and right endpoints of L . If $|\beta - \alpha| \geq 1$, then we can take L for C_0 since α is to the left of F_1 and therefore is in B . Thus we may assume that $|\beta - \alpha| < 1$. Starting at a and traversing ∂X in the negative direction the indicated points occur in the following order: $a, \beta, a_2, a_1, \alpha, e, l, r, f, b$ (with some equalities possible). If $|e - f| < 1$, then the arc $e\alpha$ joins $e, \alpha \in B$ and $|e - \alpha| \geq 2$. It is then easy to see that either $e\alpha$ has a proper contact point g for which $|g - e|, |g - \alpha| > \frac{1}{H_0}$, which (in light of Proposition 3.33) would violate (4.2), or $e\alpha$ has a subarc C_0 which is an elementary i -characteristic with the desired properties. Thus we need only consider what happens if $|e - f| \geq 1$. But in this case $C_0 = ef$ has the desired properties since it is an elementary i -characteristic and $e \in B$.

(iii). Here we assume that we are not simultaneously in Case (ii), since were that so we would be done. The point b is obviously between F_1 and F_2 , but, since we are not in Case (ii), a is between l and F_1 . Let $L = pq$, with $p < q$. If $|b - p| \geq \frac{1}{H_0}$, then, again by Proposition 3.33, B would have an ESF- i arc of diameter ϵ_0 contradicting (4.2). Thus we can assume that $|b - p| < \frac{1}{H_0}$. If $|p - q| < 1$, then we will have $|b - a| > 1$ and we can take for C_0 any elementary i -characteristic closely approximating the extended characteristic ab (recall the construction of extended characteristics in section 3), since its endpoints can be made arbitrarily close to a and b , and consequently can be chosen to lie in B . On the other hand, if the distance between p and q is at least 1, then we can take C_0 to be the elementary i -characteristic pq , since $p \in B$. Thus we are done with Step 1.

Step 1 implies that there is a subarc P of ∂X and an elementary i -characteristic $C_0 = \bar{a}\bar{b}$ with $|\bar{b} - \bar{a}| \geq 1$ joining the two endpoints of P such that P contains a subarc I of B , the distance between whose endpoints is at least 1 and one of whose endpoints is an endpoint of C_0 . To be specific we assume that the left endpoint of I is the one I has in common with C_0 . We shall henceforth work with the characteristic subdomain (D_0, C_0) bounded by $P \cup C_0$, so that in particular all order related statements will refer to the order on $\text{bot}(D_0) = P$, unless otherwise indicated. We give P the usual arc length parameterization $z = z(s), 0 \leq s \leq L$, with increasing s corresponding to the positive orientation on ∂X and such that $z(0)$ is the left endpoint of I . Obviously, $z([0, 1]) \subset I$. By the bound of $\frac{1}{10000}$ on the unsigned curvature of B it is clear that

$$(4.3) \quad |z(s_2) - z(s_1)| \geq \frac{9}{10}|s_2 - s_1| \quad \text{for } z(s_1), z(s_2) \in I.$$

(Obviously, a constant considerably closer to 1 than $\frac{9}{10}$ would also work here.) We also observe that for $\alpha, \beta \in \text{bot}(D_0)$

$$(4.4) \quad |\beta - \alpha| \geq \gamma_0 \min\{1, \lambda([\alpha', \beta'])\} \quad \text{for all } [\alpha', \beta'] \subset [\alpha, \beta].$$

This follows from (4.1) and the fact that one of the arcs of ∂X joining α to β contains

$[\alpha', \beta']$ and the other contains the endpoints of C_0 , the distance between which is at least 1. It is also easy to see that

$$(4.5) \quad \text{dist}(z(0), P \setminus B) \geq 1.$$

2. There is a subarc $J = pq$ of I with $|q - p| \geq \frac{1}{40}$ such that p and q are joined by an elementary i -characteristic in D_0 .

To establish this we again use an argument based on the case that results when Proposition 3.17 is applied. However, for this step this proposition must be applied, in conjunction with the trapped area bound (Propositions 3.34), in a sequential manner. By our underlying assumption (4.2) there are elementary i -characteristics $A_1, \dots, A_{H_0} \in \mathcal{I}(D_0)$ for which the arc $I_k \subset I \subset \text{bot}(D_0)$ joining the endpoints of A_k lies in the middle third of $z([\frac{k-1}{H_0}, \frac{k}{H_0}])$ (i.e., lies in $z([\frac{k-\frac{2}{3}}{H_0}, \frac{k-\frac{1}{3}}{H_0}])$), $1 \leq k \leq H_0$. Note that by (4.3)

$$(4.6) \quad \text{dist}(I_{k_1}, I_{k_2}) \geq \frac{9}{10H_0} \left(|k_2 - k_1| - \frac{1}{3} \right).$$

Assume inductively that we have elementary i -characteristics C_k , $0 \leq k \leq t$, such that C_k envelopes only $A_{l(k)}, A_{l(k)+1}, \dots, A_{r(k)}$, but none of the other A_n , such that

$$C_0 \succeq C_1 \succeq \dots \succeq C_t,$$

$$1 = l(0) \leq l(1) \leq \dots \leq l(t) < r(t) \leq r(t-1) \leq \dots \leq r(0) = H_0$$

and such that $r(k) - l(k)$ is strictly decreasing and satisfies

$$(4.7) \quad r(k) - l(k) \geq \frac{H_0}{3}, \quad 0 \leq k \leq t.$$

We apply Proposition 3.17 with $A = A_{l(t)}$ and $B = A_{r(t)}$ to analyze $C = \min_{D_0}(A_{l(t)} \cup A_{r(t)})$ and consider separately each of the three cases in the conclusion of that proposition. More specifically, we will show three things: first, that if Case (iii) does not occur, either we will have obtained the desired J (and will therefore stop) or we will be able to produce C_{t+1} with

$$(4.8) \quad r(t+1) - l(t+1) = r(t) - l(t) - 1;$$

second, that when Case (iii) occurs we have the desired J ; and third, that Case (iii) will have to occur long before (4.7) can be violated.

(i). There can be no contact points of C between $A_{l(t)+1}$ and $A_{r(t)-1}$, since, were there to be such a point p , the distance from p to each of the endpoints of C would, by (4.6), have to be at least $\frac{9}{10H_0}(\frac{2}{3}) > \frac{1}{10H_0}$, so that we can apply Proposition 3.33 to conclude that there is an ESF- i arc of diameter ϵ_0 on B , which contradicts (4.2). If there is a contact point of C between $A_{r(t)-1}$ and $A_{r(t)}$, we will have obtained the desired conclusion because in that case C would have to have a subarc E which is an elementary i -characteristic with left endpoint to the left of $A_{l(t)+1}$ and right endpoint between $A_{r(t)-1}$ and $A_{r(t)}$, so we stop the process here, Step 2 having been established with $J = E$ in light of (4.7). If there is no contact point of C between $A_{r(t)-1}$ and $A_{r(t)}$, then there is a contact point of C between $A_{l(t)}$ and $A_{l(t)+1}$ and C has a subarc C_{t+1} which is an elementary characteristic whose left endpoint is between $A_{l(t)}$ and $A_{l(t)+1}$ and whose right endpoint lies to the right of $A_{r(t)}$. We have $l(t+1) = l(t) + 1$ and $r(t+1) = r(t)$, so that (4.8) holds.

For the other two cases we use the notation of Proposition 3.17, namely L, U, ef , and ab , where $ef = U = C = \min_{D_0}(A_{l(t)} \cup A_{r(t)})$.

(ii) Because we are not in Case (iii), b lies to the right of $A_{r(t)}$. We have that a is between $A_{l(t)}$ and $A_{l(t)+1}$, since otherwise $|e - a| \geq \frac{9}{10H_0}(\frac{2}{3}) > \frac{1}{10H_0}$, by (4.6), and then, since $[e, a] \subset [e, b]$,

$$|e - b| \geq \gamma_0 \min\{1, \lambda([e, a])\} \geq \gamma_0 \min\{1, |e - a|\} > \frac{\gamma_0}{10H_0}$$

by (4.4), so that by Proposition 3.33 applied to the contact point e there is an ESF- i arc of diameter ϵ_0 on B , which contradicts (4.2). From the definition of extended characteristics it follows that there are elementary i -characteristics $U' = a'b' \preceq U$ with a' and b' arbitrarily close to a and b , respectively. The point a' can therefore be taken to lie between $A_{l(t)}$ and $A_{l(t)+1}$, and b' can be taken to lie to the right of $A_{r(t)}$. In this case we let C_{t+1} be any such U' and set $l(t + 1) = l(t) + 1$ and $r(t + 1) = r(t)$, so that we again have (4.8).

(iii). We will show that the first time this case occurs we will have obtained the desired conclusion of Step 2. Thus we assume that when we apply Proposition 3.17 to produce $C = \min_{D_0}(A_{l(t)} \cup A_{r(t)})$ we are in Case (iii). Now we can only say that a lies somewhere to the right of the left endpoint of C_t and that f lies somewhere to the right of $A_{r(t)}$, but, since we cannot use Proposition 3.33 to any effect (because f is not necessarily in B), we cannot conclude, in analogy with the preceding case, that b lies between $A_{r(t)-1}$ and $A_{r(t)}$. If $f \in B$, then by (4.3) and (4.7) it follows that $|f - e| \geq \frac{1}{6}$. On the other hand, if $f \notin B$, then since $f \in P$, (4.5) implies that $|f - z(0)| \geq 1$. Since by (4.7), $|e - z(0)| < \frac{7}{10}$, we have that $|f - e| \geq \frac{3}{10}$. Therefore, no matter what, $|f - e| \geq \frac{1}{6}$. Let D' be the subdomain of D_0 bounded by $ab \cup [a, b]$. It follows immediately from (4.3) that $\lambda([a, b]) \leq \frac{10}{9}|b - a|$. From the definition of extended characteristics it follows that for any $\epsilon > 0$ there are elementary i -characteristics E inside D' , joining points of $a', b' \in [a, b]$ and containing points e' and f' which are within ϵ of a, b, e, f , respectively. But it follows from Proposition 3.6 that

$$\frac{1}{6} - 2\epsilon \leq |f' - e'| \leq \text{diam}(e'f') \leq \text{diam}(a'b') \leq 5\lambda([a', b']) \leq \frac{50}{9}|b' - a'|,$$

so that with ϵ appropriately small we see that $|b' - a'| > \frac{1}{40}$. But taking ϵ sufficiently small, the corresponding E will serve as the desired J .

Thus either we stop at t with the desired J or go on to $t + 1$, the latter occurring only if we are in Case (i) or if we are in Case (ii) but not Case (iii). We now show that we must actually arrive at Case (iii) long before (4.7) can be violated. Say that we have arrived at $t = T < \frac{H_0}{3}$. Let $T' = \lceil \frac{T}{2} \rceil$. At least one of the following two things must have occurred:

(A) At least T' of the times that we passed from t to $t + 1$ we will have done so because we are in Case (i) and there is a proper contact point p_t of $\min_{D_0}(A_{l(t)} \cup A_{r(t)})$ between $A_{l(t)}$ and $A_{l(t)+1}$ and p_t is the left endpoint of C_{t+1} .

(B) At least T' of the times that we passed from t to $t + 1$ we will have done so because we are in Case (ii) but not Case (iii).

We deal with possibility (A) first. Say $k_1 < k_2$ are two values of t for which we are in Case (i). Let C' be the elementary j -characteristic, one of whose endpoints is p_{k_2} . Let its second endpoint be p' . If C' were to cross C_{k_1+1} , then, since a j -characteristic can have at most one point in common with an i -characteristic, the other endpoint of

C' would lie to the left of $A_{l(k_1)}$ or to the right of $A_{r(k_1)}$, and the distance between its endpoints would, by (4.6), have to be at least $\frac{9}{10H_0}(\frac{2}{3}) > \frac{1}{2H_0}$. However (as we have seen happen before), by (4.7) together with Proposition 3.33, B would then have an ESF- i arc of diameter ϵ_0 , contradicting (4.2). Thus C' does not cross C_{k_1+1} , so that $p' \in \text{bot}(D_0)$. If α, β are points of $\text{bot}(D_0)$ with α lying to the left of C_{k_2+1} and β to its right, then $\lambda(\alpha\beta) \geq \frac{1}{4}$ since $T < \frac{H_0}{3}$, so that by (4.4)

$$(4.9) \quad |\beta - \alpha| \geq \frac{\gamma_0}{4}.$$

If p' lies to the right of C_{k_2+1} , then each point of $C_{k_2+1} \cap X$ is joined to a point in $\text{bot}(D_0)$ to the right of C_{k_2+1} by a j -half-characteristic lying between C_{k_2+1} and $C' \preceq C_{k_1+1}$. Then a simple argument (which was used in the proof of Proposition 3.34) based on Proposition 2.8(i) and the fact that for each point of the first half of C_{k_2+1} , which has length at least $\frac{2\gamma_0}{8}$, the corresponding j -half-characteristic has length at least $\frac{2\gamma_0}{8}$ shows that the area between C_{k_2+1} and C_{k_1+1} must be at least $\eta(\frac{2\gamma_0}{8})^2 = \frac{\eta\gamma_0^2}{64} \geq \frac{\eta_1\gamma_0^2}{64}$. If, on the other hand, $p' < p_{k_2}$, then $|\Delta\theta(C')| \geq \frac{\pi}{4}$. Also, as we just saw, (4.9) holds for all $\alpha, \beta \in \text{bot}(D_0)$ with α lying to the left of C_{k_2+1} and β to its right. Thus by the trapped area bound (Proposition 3.34) it follows that the area between C_{k_2+1} and C_{k_1+1} is at least $\frac{\eta_1\gamma_0^2}{16} > \frac{\eta_1\gamma_0^2}{64}$, so that this lower bound holds no matter which side of C_{k_2+1} the point p' lies on. But then $\frac{(T'-1)\eta_1\gamma_0^2}{64} \leq \mu(X) \leq 200$, so that $T' \leq \frac{12800}{\gamma_0^2\eta_1} + 1$. This in turn means that $T \leq \frac{30000}{\gamma_0^2\eta_1} < \frac{H_0}{10}$ (since γ_0 and η_1 are both in $(0, 1)$).

Possibility (B) is handled in a similar manner. Say $k_1 < k_2$ are two values of t for which we are in Case (ii) but not in Case (iii) when passing from C_t to C_{t+1} . From the construction of C_{k_1+1} we have that $U = \min_{D_0}(A_{l(k_1)} \cup A_{r(k_1)}) = u_1u_2$ is a subarc of a nonmonotone extended i -characteristic W which has another subarc S which is also an elementary i -characteristic for which $S \preceq U$ and $A_{l(k_1)} \preceq S$. It follows in addition from that construction that S precedes $C_{k_1+1} \preceq U$, so that in particular S precedes A_k for all $k \geq l(k_1) + 1$. There is an elementary j -characteristic S' which joins u_1 to a point $u' \in [u_1, u_2]$. It follows immediately from (4.7) and (4.4) that $|u_2 - u_1| \geq \frac{\gamma_0}{4}$. We claim that $S' \cap D_0 \cap C_{k_2+1} = \emptyset$. To see this, assume to the contrary that $z \in S' \cap D_0 \cap C_{k_2+1}$. It then follows from (4.4) and (4.6) that $|u' - u_1| \geq \frac{9\gamma_0}{10H_0}(\frac{2}{3}) > \frac{\gamma_0}{2H_0}$, so that by Proposition 3.33, B would have an ESF- i arc of diameter at least ϵ_0 in contradiction to our underlying assumption (4.2). Now we apply the trapped area bound to the i -characteristic subdomain (D_1, C_{k_2+1}) , where D_1 is bounded by the curve made up of C_{k_2+1} together with $\partial X \setminus (\alpha_1, \alpha_2)$, where $\alpha_1 < \alpha_2$ are the endpoints of C_{k_2+1} (that is, D_1 is the part of X that remains when the part of D_0 on and below C_{k_2+1} is removed). In particular $S' \setminus \{u_1, u'\} \subset D_1$ and $U \in \mathcal{I}(D_1)$. It is also clear from (4.7) and (4.4) that (4.9) holds for α and β in $\text{bot}(D_1)$ on opposite sides of C_{k_1+1} . For the same reason that $S' \cap D_0 \cap C_{k_2+1} = \emptyset$, we have that u' must lie between u and α_1 with respect to the order on $\text{bot}(D_0)$. This means that both endpoints of S' are in B , so that we clearly have $\Delta\theta(S') \geq \frac{\pi}{4}$. It also means that with respect to the order on $\text{bot}(D_1)$, u_1 lies between u' and u_2 . We can now apply the trapped area proposition (with (D_1, C_{k_2+1}) playing the role of (D, V)) to deduce that the area between U and C_{k_2+1} is bounded below by $\frac{\eta_1\gamma_0^2}{16}$, so that this same lower bound holds for the area between C_{k_1} and C_{k_2+1} . If the values of t in question are $t_1 < t_2 < \dots < t_{T'}$, and if R_k is the region between C_{t_k} and $C_{t_{k+1}+1}$, then R_1, R_3, \dots

are disjoint, so that

$$\left(\frac{T' - 2}{2}\right) \frac{\eta_1 \gamma_0^2}{16} \leq \mu(X) \leq 200,$$

and therefore $T' \leq \frac{6400}{\gamma_0^2 \eta_1} + 2$. Thus, as in the case that (A) holds, we again have $T < \frac{H_0}{10}$. This concludes the proof of Step 2.

3. Now we work with the $J = pq \subset I$, with $|q - p| \geq \frac{1}{40}$, whose endpoints are joined by a $C = pq \in \mathcal{I}(D_0)$. Clearly, there are two points $a_1, a_2 \in J$ such that if we consider $X_k = \Sigma(a_k, \frac{1}{10^6})$, $k = 1, 2$, where the Σ (as defined in the third paragraph of this section) are with respect to the normalized domain X with which we are now working, then

- $10^6(X_k - a_k)$ is a normalized domain, $k = 1, 2$,
- $\text{dist}(X_k, \{p, q\}) \geq \frac{1}{500}$, $k = 1, 2$,
- $\text{dist}(X_1, X_2) \geq \frac{1}{500}$, and
- the bottoms of X_1 and X_2 are in J .

By what we have shown it then follows that either the bottom one of X_1 or X_2 has an ESF- i of length $\frac{\epsilon_0}{10^6}$, in which case we have reached the desired conclusion, or for both $k = 1$ and $k = 2$ the bottom of X_k has a subarc $p_k q_k$, with

$$(4.10) \quad |q_k - p_k| \geq \frac{1}{4 \cdot 10^7}$$

whose endpoints are joined by an elementary i -characteristic C_k . In this latter case, since C joins the endpoints of J , $C_k \preceq C$, $k = 1, 2$, we can apply Proposition 3.17 to $\max_{D'}(C_1 \cup C_2)$, where D' is the characteristic subdomain bounded by $C \cup [p, q]$. But in light of the above conditions satisfied by the X_k and (4.10), and since all contact points of $\max_{D'}(C_1 \cup C_2)$ are in J , Proposition 3.33 immediately implies that there is a $\gamma_1 = \gamma_1(K)$ for which J has an ESF- i arc of J of diameter γ_1 . Thus we have proved that the bottom of a normalized domain has an ESF- i arc of length $\Delta_0 = \min\{\frac{\epsilon_0}{10^6}, \gamma_1\}$. \square

It follows from the opening discussion and the proposition we have just proved that there are $\delta_0 = \delta_0(G)$, $\alpha_0 = \alpha_0(K) \in (0, 1)$ such that any arc B of ∂G of diameter at most δ_0 has an ESF- i subarc B' , for which $\text{diam}(B') > \alpha_0 \text{diam}(B)$. It is then clear that there are numbers $\xi_0 = \xi_0(K)$, $\tau_0 = \tau_0(K)$, both in $(0, 1)$, such that every arc B with $\text{diam}(B) \leq \delta_0$ has two disjoint subarcs B_1 and B_2 for which

$$\begin{aligned} (\text{diam}(B_1))^{\tau_0} + (\text{diam}(B_2))^{\tau_0} &< \xi_0 (\text{diam}(B))^{\tau_0}, \\ \text{diam}(B_1), \text{diam}(B_2) &< (1 - \alpha_0) \text{diam}(B), \end{aligned}$$

$$B \setminus (B_1 \cup B_2) \text{ is ESF-}i.$$

Let $\text{diam}(B) \leq \delta_0$. We start with two such arcs B_1 and B_2 corresponding to B , then we apply the same fact to each of these to get four arcs of diameter $l_1, \dots, l_4 < (1 - \alpha_0)^2$ for which $\sum l_k^{\tau_0} < \xi_0^2 \delta_0^{\tau_0}$, and such that the complement of their union is ESF- i , then we do so again to get eight arcs of length at most $(1 - \alpha_0)^3$ for which the corresponding $\sum l_k^{\tau_0}$ is less than $\xi_0^3 \delta_0^{\tau_0}$ and for which the complement of their union is ESF- i , and so on. At the n th stage we have 2^n arcs $B_k^{(n)}$ for which $\sum_k (\text{diam}(B_k^{(n)}))^{\tau_0} < \xi_0^n \delta_0^{\tau_0}$, which tends to 0, and for which $B \setminus (\cup_k B_k^{(n)})$ is ESF- i . At this stage the diameter of

the largest of the 2^n arcs is at most $(1 - \alpha_0)^n$. But then B is the union of a set of τ_0 -dimensional Hausdorff measure 0 and a set that has at most countably many i -singularities. In light of Proposition 3.29 and the fact that there are at most countably many vertices of fans on ∂G (Proposition 3.24) it follows that the set of singularities on B , and therefore the set of singularities on all of ∂G , has Hausdorff dimension at most τ_0 . This finishes the proof of the main theorem.

COROLLARY 4.3. *Let Θ be a normal system with $\tau = \tau(\Theta) < 1$. Let $G \subset \mathbb{C}$ be a domain with boundary $\partial G = R \cup \Theta$, where R is a ray. Then G is τ -regular.*

This is an immediate consequence of the main theorem and Propositions 3.20 and 3.25. \square

5. Construction of solutions with singularity sets of positive Hausdorff dimension. Although we have chosen to carry out our construction in \mathbb{H} to avoid cumbersome arguments, suitable, largely straightforward changes will allow an analogous procedure to be carried out in any smoothly bounded domain. Throughout this section Θ will denote any fixed normal system. We begin with a brief discussion of characteristic initial value problems for Θ since our construction largely proceeds by joining together solutions of appropriate instances of such problems in domains sharing a boundary characteristic.

Let C_k be C^∞ curves parameterized by $z_k(s)$, $0 \leq s \leq l_k$, $k = 1, 2$, for which $z_i(0) = z_j(0)$ and for which $z'_i(0)$ and $z'_j(0)$ are mutually orthogonal. (In fact l_k can be infinite, as is the case, for example, when C_k is a ray.) We want to construct a solution of the normal system Θ for which C_k is a k -characteristic arc, $k = 1, 2$. From the definition of normal system, for any ρ_1 there is a discrete set of values ρ_2 such that $e^{i\theta(\rho_1, \rho_2)}$ and $ie^{i\theta(\rho_1, \rho_2)}$ are tangent at $z_1(0)$ to C_1 and C_2 , respectively. Given any such $\rho = (\rho_1, \rho_2)$ there are unique continuous functions $R_k(s)$ such that $e^{i\theta(\rho_1, R_2(s))}$ and $ie^{i\theta(R_1(s), \rho_2)}$ are tangent to C_1 and C_2 at the points $z_1(s)$ and $z_2(s)$, respectively. For $t = (t_1, t_2)$ we define

$$\bar{R}(t) = (R_1(t_2), R_2(t_1)), \quad 0 \leq t_k \leq l_k, \quad k = 1, 2,$$

and

$$\bar{\theta}(t) = \theta(R_1(t_2), R_2(t_1)), \quad 0 \leq t_k \leq l_k, \quad k = 1, 2.$$

We seek $\zeta : [0, l_1] \times [0, l_2] \rightarrow \mathbb{C}$, where $\zeta(t) = \xi(t) + i\eta(t)$, for which $e^{i\bar{\theta}(t)}$ and $ie^{i\bar{\theta}(t)}$ are tangent to the curves $\zeta([0, l_1], t_2)$ and $\zeta(t_1, [0, l_2])$, respectively, at the point $\zeta(t)$ and which satisfies the initial conditions

$$(5.1) \quad \zeta(t_1, 0) = z_1(t_1) \quad \text{and} \quad \zeta(t_2, 0) = z_2(t_2).$$

The tangency condition can be written as the linear hyperbolic system

$$(5.2) \quad \cos \bar{\theta} \frac{\partial \eta}{\partial t_1} - \sin \bar{\theta} \frac{\partial \xi}{\partial t_1} = 0, \quad \sin \bar{\theta} \frac{\partial \eta}{\partial t_2} + \cos \bar{\theta} \frac{\partial \xi}{\partial t_2} = 0.$$

The problem (5.2) with initial conditions (5.1) is well posed and has a C^∞ solution $\bar{\theta}$ for any C^∞ initial curves C_1 and C_2 . First consider the case in which the $\zeta(t) = \zeta(C_1, C_2, \rho, t)$ determined in this manner has an everywhere nonzero Jacobian

determinant and is globally one-to-one. We then have a solution to the system Θ in $\zeta([0, l_1] \times [0, l_2])$, which is a characteristic quadrilateral, given by

$$(5.3) \quad R(\zeta(t)) = \overline{R}(t).$$

We refer to this solution as $R(C_1, C_2, \rho, z)$ and to the characteristic quadrilateral in which it is defined as $Q(C_1, C_2, \rho)$. Even if ζ is not one-to-one on $[0, l_1] \times [0, l_2]$, then (5.3) gives a multivalued solution of Θ on any domain $\zeta(E)$, provided that ζ is a local diffeomorphism on E . Let i and j be such that $\arg\{z'_j(0)\} = \arg\{z'_i(0)\} + \frac{\pi}{2}$. It is a well-known property of the characteristic initial value problem for genuinely nonlinear systems that if $\frac{d \arg\{z'_i(t)\}}{dt} \leq 0$ on $[0, l_i]$ and $\frac{d \arg\{z'_j(s)\}}{ds} \geq 0$ on $[0, l_j]$, then ζ will be locally diffeomorphic on $[0, l_1] \times [0, l_2]$. This is simply a reflection of the fact that, in light of the quasi-HP nature of the the associated inclination function and the length monotonicity property (Proposition 2.7), under these hypotheses the curvature of j -characteristics will not blow up as one moves along i -characteristics away from the convex side of a j -characteristic C , and analogously when the roles of i and j are interchanged. We note that this property holds in the particular case in which one of the initial characteristics is a line segment or ray. We also note that it holds even if one or both of the initial curves is not a simple arc, so that ζ will be a local homeomorphism if, for example, one of these curves is a circle covered several times.

We next need to discuss briefly the smooth adjunction of line segments and circles to C^∞ curves in a specific, constructive fashion; this is another essential element of the construction process. Let C be a C^∞ curve parameterized by $Z(s)$, $0 \leq s \leq l + \delta$. We want to perform the adjunction with no change in Z on $[0, l]$. Let κ_0 be any number. It is clear that there is an operator $F(Z, l, \delta, \kappa_0, \tau)$ such that $w = F(Z, l, \delta, \kappa_0, \tau) \in C^\infty([0, \infty))$ has the following properties:

- (i) $w(s) = Z(s)$, $0 \leq s \leq l$.
- (ii) If $w'(s) = e^{i\phi(s)}$, then $\phi(s) = (s - l - \delta)\kappa_0 + \tau$, $s \geq l + \delta$.

It is clear how this can be done. Indeed, if $Z'(s) = e^{i\alpha(s)}$ on $[0, l + \delta]$, we let $\beta(s)$ be the continuous function which coincides with $\alpha(s)$ on $[0, l + \delta/3]$, which is given by $(s - l - \delta)\kappa_0 + \tau$ for $s \geq l + 2\delta/3$ and is linear on $[l + \delta/3, l + 2\delta/3]$. We let σ be a specific nonnegative C^∞ function on \mathbb{R} with support in $(-1, 1)$ and $\int_{-\infty}^\infty \sigma(s)ds = 1$ and then convolve $\beta(s)$ with $\frac{10}{\delta}\sigma(\frac{10s}{\delta})$. The desired $w = F(Z, l, \delta, \kappa_0, \tau)$ is defined by with $w'(s) = e^{i\beta(s)}$, $w(0) = Z(0)$. We note that if C is convex to the right (left) and $\kappa_0, \tau - \kappa_0\delta/3 - \alpha(l + \delta/3) \geq 0$ (≤ 0), then the curve given by $F(Z, l, \delta, \kappa_0, \tau)$ is concave towards the same side as C .

We now introduce notation and terminology to be used in our construction. We denote by \mathcal{CL} the family of C^∞ arcs with initial and terminal straight subarcs. Let $\mathcal{S}(\epsilon)$ denote the class of C^∞ curves parameterized by $z = z(s)$, $0 \leq s \leq L$, with the following properties:

- (i) $\Im\{z(0)\} = \Im\{z(L)\} = -\epsilon$.
- (ii) $\arg\{z'(0)\} = \frac{\pi}{2} + \epsilon$ and $\arg\{z'(L)\} = -(\frac{\pi}{2} + \epsilon)$.
- (iii) $\frac{d \arg\{z'(s)\}}{ds} \leq 0$, $0 \leq s \leq L$.
- (iv) $C \subset N(\partial N(0, 1) \cap \mathbb{H}, 2\epsilon)$.
- (v) $C \in \mathcal{CL}$.

For $t > 0$ and $\alpha \in \mathbb{R}$ we denote $\{tS + \alpha : S \in \mathcal{S}(\epsilon)\}$ by $\mathcal{S}(\epsilon, t, \alpha)$. If C is any arc joining $a < b$ in $\overline{\mathbb{H}}$, $B(C)$ will denote the set whose boundary is $C \cup [a, b]$. For two such arcs C_1, C_2 with $C_1 \subset B(C_2)$ we denote the closure of $B(C_2) \setminus B(C_1)$ by

$E(C_1, C_2)$. In what follows we will often deal with a C^∞ solution R defined only in a neighborhood in $B(C)$ of C , where C is an i -arc of R (by which we mean that R can be extended to a C^∞ solution in a two-sided neighborhood of C and that C is an i -arc of this extension). In this case we shall call C a i -characteristic of R . An i -characteristic C of a solution R will be said to satisfy the i -oscillation condition (OSC) if the j -characteristics passing through each point $c \in C$ all contain a straight line subarc containing c in its interior; obviously, it is enough for a single point $c \in C$ to have this property for the OSC to hold. We extend the use of the term OSC in the obvious manner to apply to one-sided characteristics. If D is a domain separated into two subdomains D_1 and D_2 by an arc C , and $R^{(1)}$ and $R^{(2)}$ are C^∞ solutions in these domains, respectively, for which C is a one-sided i -characteristic satisfying the OSC, then together $R^{(1)}$ and $R^{(2)}$ give a C^∞ solution in all of D . Our construction uses this simple fact, which allows for the smooth pasting together of solutions of characteristic initial value problems in contiguous domains.

Next we discuss a specific class of characteristic initial value problems in which one of the initial characteristics is a ray. Let C be parameterized by $z(s)$ satisfying

$$(5.4) \quad D \arg\{z'(s)\} \leq 0, \quad 0 \leq s \leq L, \quad \text{with } \Im\{z(0)\}, \quad \Im\{z(L)\} < 0,$$

and

$$(5.5) \quad \arg\{z'(0)\} \geq \frac{\pi}{2} + \epsilon \quad \text{and} \quad \arg\{z'(L)\} \leq -\left(\frac{\pi}{2} + \epsilon\right),$$

and consider the solution of the initial characteristic value problem where $C_i = C$ and C_j is a ray orthogonal to C at $z(0)$ and emanating to the left of C . Then the j -characteristics of any corresponding solution are all rays orthogonal to C_i and the i -characteristics are the orthogonal trajectories of this family of rays, and in fact are the curves $C(r)$, $r > 0$, parameterized by $z(t) + rz'(t)i$, $0 \leq t \leq L$. As r tends to ∞ the $C(r)$ tend to arcs of a circle of radius r and radian measure $\arg\{z'(0)\} - \arg\{z'(L)\}$.

We next construct two special C^∞ arcs U_i , where it is understood that the ϵ referred to below is some sufficiently small positive number. The arc U_1 , parameterized by $u_1(s)$, $0 \leq s \leq \lambda_1$ will have the following properties:

- (i) $D \arg\{u_1'(s)\} \leq 0$, $0 \leq s \leq \lambda_1$.
- (ii) $\Im\{u_1(0)\} = -\epsilon$ and $u_1(\lambda_1) \in \mathbb{R}$.
- (iii) $U_1 \in \mathcal{CL}$.
- (iv) $u_1'(0) = e^{i(\frac{\pi}{2} + \epsilon)}$.
- (v) $u_1'(s) = e^{-i\pi/4}$, $\lambda_1 - d \leq s \leq \lambda_1$ for some $d > 0$.
- (vi) $U_1 \cap N(0, 2) = \emptyset$.

(vii) Let any C' be any curve in $\mathcal{S}(\epsilon)$, and let $\alpha = \alpha(C)$ be such that the initial point of $C = C' + \alpha \in \mathcal{S}(\epsilon, 1, \alpha)$ is $p = -1 - \epsilon i$. Let $\rho \in \mathbb{R}^2$ be such that $ie^{i\theta(\rho)}$ is tangent to C at its left endpoint p . Then there is a solution to the system Θ in $E(C) = E(C \cap \mathbb{H}, U_1 \cap \mathbb{H})$ for which $R(p) = \rho$, for which C is a 2-characteristic and U_1 is an 1-characteristic, and for which both C and U_1 satisfy the OSC.

We stress that the idea is that for the appropriate translate C of C' , $C' \in \mathcal{S}(\epsilon)$ and any admissible value ρ of R at the initial point (i.e., left endpoint) c_l of C we have a solution in the domain $E(C \cap \mathbb{H}, U_1 \cap \mathbb{H})$. If we can construct such a curve, it is obvious that it will have a right-hand counterpart U_2 , having properties corresponding to (i)–(vii). In regard to (i), U_2 will terminate in a line segment of slope $\frac{\pi}{4}$ at its left end. In regard to (vii), U_2 and C will be 2- and 1-characteristics, respectively, and the initial value ρ will be the value of R at the right endpoint of C , which in any

case determines its value at the right endpoint of C since the total curvature of C is exactly $\pi + 2\epsilon$. For lack of a better term we shall refer to U_1 and U_2 as U_1, U_2, \dots, U_n for the system Θ . It is clear that translates $U_i + \alpha, \alpha \in \mathbb{R}$ of U_i have an analogous universal property.

We show that such a U_1 exists, the existence of U_2 being identical apart from trivial details. Let $A \leq |\frac{\partial \theta}{\partial R_k}| \leq B$ on \mathbb{R}^2 (as in the definition of normal system), so that the net of characteristics of any solution of the system Θ in a domain G is a K -quasi-HP net on G with $K = \frac{B}{A}$. Let W be any C^∞ arc parameterized by $w(s), 0 \leq s \leq L$, for which $D \arg\{w'(s)\} \leq 0$ and which has total curvature $-\Delta < 0$. Let $V \in \mathcal{CC}$ be parameterized by $v(s), 0 \leq s \leq 1$, with $v(0) = w(0), v'(0) = iw'(0)$, for which $0 \geq D \arg\{v'(s)\} \geq -\delta$. One easily show that for any admissible corner value ρ the characteristic coordinate mapping $\zeta(V, W, \rho, t)$ is a local diffeomorphism on all of the corresponding characteristic coordinate rectangles $S = [0, 1] \times [0, L]$ for all such W , provided that $\delta K(L + K\Delta) < 1$. To do so, one simply subdivides S into small subrectangles and solves the corresponding initial value problem piece by piece, the key point being that the length change estimate (Proposition 2.5) implies that no blow-up will occur (that is, the Jacobian determinant will never vanish). Indeed, in light of the convexity of W the lengths of its translates increase, while in light of the concavity of V the lengths of its translates decrease. The characteristic length bound as applied to the translates of W shows that they all have length at most $L + K\Delta$, so that by the characteristic length bound the lengths of the translates of V will remain bounded away from 0 if $\delta K(L + K\Delta) < 1$. We start with $W = W_0, V = V_0$ and inductively define ρ_{k+1} to be the value of $R(V_k, W_k, \rho_k, \zeta(V_k, W_k, \rho_k, (1, 0)))$, that is, the value of $R(V_k, W_k, \rho_k, p)$, where p is the the final endpoint of V_k , and we define W_{k+1} to be the final translate of W_k along V_k to the end of V_k , that is, $W_{k+1} = \zeta(V_k, W_k, \rho_k, (1, [0, L_k]))$, where $L_k = \lambda(W_k)$. By the foregoing $L_{k+1} \leq L_k + K\Delta$, and the total curvature of V_k can be any $-\delta_k$, for which $0 \leq \delta_k \leq \frac{1}{2K(L_k + K\Delta)}$. We obviously have $L_k \leq L_0 + kK\Delta$, so that the only restriction on δ_k is $0 \leq \delta_k \leq \frac{1}{2K(L_0 + (k+1)K\Delta)}$, this upper bound obviously being the general term of a divergent series for $\Delta > 0$. Let $\epsilon = \frac{1}{100}$, so that if $W_0 = C$, where C is the appropriate translate of any element of $\mathcal{S}(\epsilon)$ as indicated in (vii), then $\pi - \frac{1}{10} \leq \lambda(W_0) \leq \pi + \frac{1}{10}$ and $\Delta = \pi + \epsilon$. Clearly, there is some positive integer M such that for some sufficiently small $\eta \in (0, \frac{\pi}{4})$ there is a curve V , the union of V_0, \dots, V_M , parameterized by $v(s), 0 \leq s \leq M$, such that V terminates in a straight line segment whose initial point is in \mathbb{H} and whose slope is $-\tan \eta$. We note that by reducing the size of ϵ if necessary, we can assume that the first subarc V_0 of $V \subset V^1$ intersects \mathbb{R} at a point $v_0(s)$ for which $\arg\{v_1'(s)\} \in (-\pi, -\pi + \delta_0) \subset (-\pi, -\frac{3\pi}{4})$. In addition, we can allow the length of the final straight part of V to be as large as we want, so that we can take the imaginary part of the terminal point of V to be $-\epsilon$. We emphasize that V is independent of W_0 . We next extend V at its right end by smoothly adjoining a ray with slope exactly $-\frac{\pi}{4}$ (by means of $F(v, \lambda(V) - \frac{\epsilon}{2}, \frac{\epsilon}{2}, 0, -\frac{\pi}{4})$, as described above, where $v(s)$ parameterizes V); this extended curve will be called V^1 . We now consider the characteristic initial value problem with V^1 as the 1-arc and, as the 2-arc, a segment A of length r orthogonal to it at its initial point and emanating to the left of V^1 . Because V^1 satisfies the OSC, the solution so generated will be C^∞ in the domain made up of $E(C \cap \mathbb{H}, V^1 \cap \mathbb{H})$ together with translates of A around V^1 . Note that the corresponding solution is constant in the quarter-plane made up of all the translates of A emanating from points on the ray in which V^1 ends. From this it is clear that for $r = \lambda(A)$ sufficiently large, an appropriate subarc of the translate $V^1(r)$

of V^1 to the outer endpoint of A has properties (i), (v), (vi), and (vii). Properties (ii), (iii), and (iv) can be achieved in the following manner. First, we take r so large that at the left intersection point of $V^1(r)$ with \mathbb{R} the angle formed by $V^1(r)$ is within $\frac{\epsilon}{4}$ of $\frac{\pi}{2}$. Then we remove the arc of $V^1(r)$ joining its initial point to a point with imaginary part between $-\frac{\epsilon}{2}$ and 0 and at which the direction of the tangent vector is within $\frac{\epsilon}{2}$ of $\frac{\pi}{2}$. Finally, we can use the operator F to smoothly adjoin an arc to the remaining part of $V^1(r)$ in such a way that (ii), (iii), and (iv) hold. The resulting curve is our U_1 . Note that in (vii) we do not require the solution to bear any particular relation to the part of U_1 that lies in the lower half-plane.

The left and right endpoints of U_1 will be called e_1 and m_1 ; the left and right endpoints of U_2 will be called m_2 and e_2 (m for middle and e for end, for reasons to be made clear momentarily). It is clear that U_1 depends solely on Θ and the parameter ϵ and that once these have been fixed, the solution R generated in $E(C)$ depends on C but that its values on U_1 depend only on $\rho = R(c_l)$, c_l being the left endpoint of C ; analogous statements hold for U_2 . For $\rho \in \mathbb{R}^2$ we denote by $f_k(\rho)$ the value of this solution R at m_k . We also observe that in the case of U_1 we always have

$$(5.6) \quad \theta(f_1(\rho)) - \theta(\rho) = \theta(m_1) - \theta(c_l) = -\left(\frac{5\pi}{4} + \epsilon\right).$$

As indicated above, in the case of U_2 we still take as ρ the value of R at the left endpoint c_l of the initial curve C used in the above construction which, as we have pointed out, uniquely determines the value of R at the right endpoint of C since the total curvature of all arcs in the family $\mathcal{S}(\epsilon)$ is $\pi + 2\epsilon$. Here we have that

$$(5.7) \quad \theta(f_2(\rho)) - \theta(\rho) = \theta(m_2) - \theta(c_l) = -(\pi + 2\epsilon) + \left(\frac{5\pi}{4} + \epsilon\right) = \frac{\pi}{4} - \epsilon.$$

The f_k are continuous functions of ρ and there is a number B_0 (that depends only on the bounds on the $\frac{\partial\theta}{\partial R_k}$ associated with Θ) such that

$$(5.8) \quad \|f_k(\rho) - \rho\| \leq B_0, \quad \rho \in \mathbb{R}^2, \quad k = 1, 2,$$

where the norm is just the Euclidean norm in \mathbb{R}^2 (actually, any norm will do). Both the continuity and the bound come from the simple observation that on any i -characteristic, R_i is constant and R_j is a Lipschitz continuous function of the tangent inclination (with the Lipschitz constant depending solely on the system Θ).

For our next step we consider translates $U'_1 = U_1 - m_1$ and $U'_2 = U_2 - m_2$ which have the common point 0 . Note that U'_1 and U'_2 are orthogonal to each other at 0 and so can be used as initial curves for a characteristic initial value problem, U'_k being a k -arc. For ϵ sufficiently small (so that U_1 and U_2 are virtually orthogonal to the horizontal line $\Im\{z\} = -\epsilon$ at e_1 and e_2 , respectively) convexity considerations easily show that for any admissible corner value $\rho = R(0)$ the mapping $\zeta(U'_1, U'_2, \rho, t)$ is one-to-one on the corresponding characteristic coordinate rectangle. To see that this is indeed the case, note that since the change in θ from the left end of U'_1 to the right end of U'_2 is $-(\frac{3\pi}{2} + 2\epsilon)$ (that is, the change in θ along U'_1 + the change in θ along U'_2 , both from left to right, is $-(\frac{3\pi}{2} + 2\epsilon)$), the total change along the top two sides of the characteristic quadrilateral is also $-(\frac{3\pi}{2} + 2\epsilon)$. The resulting solution $R_\rho(z) = R(U'_1, U'_2, \rho, z)$ and the quadrilateral $Q_\rho = Q(U'_1, U'_2, \rho)$ itself depend solely on $\rho = R(0)$. Note also that because $U'_1, U'_2 \in \mathcal{CL}$, the other two sides of Q_ρ also belong to \mathcal{CL} .

Next we next show how we can produce curves $X_k \in \mathcal{S}(\epsilon, t_k, \alpha_k)$, such that $B(X_k \cap \overline{\mathbb{H}}) \supset Q_\rho \cap \mathbb{H}$, and a solution in $B(X_k \cap \overline{\mathbb{H}}) \setminus (Q_\rho \cap \mathbb{H})$ which gives a C^∞ extension of R_ρ for which X_k is a one-sided k -arc which satisfies the OSC. We show how to get X_2 , the case of X_1 being identical apart from obvious details.

First, we extend U'_2 on the right by continuously adjoining a circle of radius $\frac{\epsilon}{10}$ lying in the open lower half-plane by using $F(\bar{u}_2 - m_2, \lambda(U_2) - \delta, \delta, \frac{10}{\epsilon}, -(\frac{\pi}{2} + \epsilon))$, where \bar{u}_2 parameterizes U_2 from left to right (i.e., with $\bar{u}_2(0) = m_2$) and where δ is so small that $\bar{u}_2([\lambda(U_2) - \delta, \lambda(U_2)])$ is a line segment. This in turn allows us to extend the other 2-side of Q_ρ rightwards and eventually downwards (in light of the quasi-HP property) until it gets to a point p for which $\Im\{p\} = -\frac{\epsilon}{2}$; call this extended 2-side T . The corresponding solution is defined in the simply covered domain $E = E(U'_1 \cup U'_2, T')$, where T' is the subarc of T which joins a point of \mathbb{R} to the left of U'_1 to a point of \mathbb{R} to the right of U'_2 . This means that if we change T below \mathbb{R} to form a new curve T_2 , which we subsequently use as the 2-arc for a characteristic initial value problem with a straight initial 1-arc outside of E , then the solution is compatible with the restriction to E of the solution we have so far. Using an appropriate instance of the operator F , we obtain an extension T_2 of T by smoothly adjoining to T at its right end a segment in such a way that the argument of the tangent at the right endpoint of T_2 is less than or equal to $-(\pi/2 + \epsilon)$. Obviously, $T_2 \in \mathcal{CL}$. It is also clear from the construction that T_2 depends solely on $\rho = R(0)$ and that all derivatives of its arc length parameterization depend continuously on ρ . This outside curve T_2 clearly satisfies the hypotheses (5.4) and (5.5) of the straight line characteristic initial value construction. Using a straight 1-arc of length r_0 we therefore obtain an ‘‘almost circular’’ 2-arc S parallel T_2 which is contained in $N(0, (1 + \epsilon/2)r_0) \setminus (0, (1 - \epsilon/2)r_0)$. Let $q_1 < q_2$ be the points of $S \cap \mathbb{R}$. By choosing r_0 sufficiently large we can be certain that the interior angles of this S at q_1 and q_2 are within $\frac{\epsilon}{2}$ of $\frac{\pi}{2}$. Simple curvature considerations show that there is a single, sufficiently large value of r_0 that works for all admissible corner values $\rho = R(0)$. Let $p_1, p_2 \in T_2$ lie on the straight 1-characteristics which terminate at q_1 and q_2 , respectively. It is clear that for the solution our process has generated, $R(p_1)$ and $R(p_2)$ are continuous functions of ρ which, moreover, satisfy bounds of the form (5.8) (with $f_k(\rho)$ replaced by $R(p_k)$ and with an appropriate B_0). The solution is defined and C^∞ in the part of \mathbb{H} between $U'_1 \cup U'_2$ and S , and S is a 2-characteristic arc satisfying the OSC. Finally, we smoothly alter S below \mathbb{R} to obtain an S' such that for appropriate $\xi = \xi(\rho)$, the arc $X_2 = \xi S'$ belongs to $\mathcal{S}(\epsilon)$. If x_l is the left endpoint of X_2 , then $g_2(\rho) = R(x_l)$ is a continuous function of ρ , and similarly for the value $g_1(\rho)$ of R at the left endpoint of the analogously constructed X_1 . Here again one easily sees that in the case of X_2

$$(5.9) \quad \theta(g_2(\rho)) - \theta(\rho) = \frac{3\pi}{4} + \epsilon - \frac{\pi}{2} = \frac{\pi}{4} + \epsilon$$

and in the case of X_1 that

$$(5.10) \quad \theta(g_1(\rho)) - \theta(\rho) = \frac{3\pi}{4} + \epsilon.$$

Let $U = U'_1 \cup U'_2$. By the universal property of U'_1 , for appropriate $t^{(2)}$ and $\alpha^{(2)}$ we can use $t^{(2)}X_2 + \alpha^{(2)}$ as the 2-arc under U'_1 (see property (vii) in the discussion of the universal arcs given above), and similarly we can use some $t^{(1)}X_1 + \alpha^{(1)}$ as the 1-arc under U_2 . These $t^{(k)}X_k + \alpha^{(k)}$ in turn come from $t^{(k)}U + \alpha^{(k)}$. If we perform the construction of X_k using the value $\rho^{(k)}$ of R at the center $\alpha^{(k)}$ of $t^{(k)}U + \alpha^{(k)}$,

then the value of R at the left endpoint of X_k is $g_k(\rho^{(k)})$, so that the value of the solution R (corresponding to property (vii) of universal arcs) in the part of the $B(U'_j)$ between $t^{(i)}U + \alpha^{(i)}$ and U'_j at 0 is $f_j(g_i(\rho^{(i)}))$. We claim that there is a $\rho^{(i)}$ for which $f_j(g_i(\rho^{(i)}))$ is any given admissible corner value ρ_0 for the characteristic initial value problem with 1- and 2-arcs U_1 and U_2 . To see that such a $\rho^{(i)}$ exists in the case that $i = 2, j = 1$, for example, note that from (5.6) and (5.9) we have that

$$\begin{aligned} \theta(f_1(g_2(\rho^{(2)}))) - \theta(\rho^{(2)}) &= \theta(f_1(g_2(\rho^{(2)}))) - \theta(g_2(\rho^{(2)})) + \theta(g_2(\rho^{(2)})) - \theta(\rho^{(2)}) \\ &= -\left(\frac{5\pi}{4} + \epsilon\right) + \frac{\pi}{4} + \epsilon = -\pi. \end{aligned}$$

Thus, given ρ_0 , the point $\rho^{(2)}$ must lie on the curve

$$(5.11) \quad \theta(R) = \theta(\rho_0) + \pi$$

in the R -plane. The level curves $\theta(R) = \theta_0$ in the R -plane are the graphs $R_i = h_{i,\theta_0}(R_j)$ of monotone functions, for which h'_{i,θ_0} are uniformly bounded and uniformly bounded away from 0. Thus as $\rho^{(2)}$ moves along the curve (5.11) it follows from the continuity of f_1 and g_2 and the bound (5.8) and the corresponding bound for the g_k that $f_1(g_2(\rho^{(2)}))$ traces out the entire curve $\theta(R) = \theta(\rho_0)$, so that there is indeed a (unique) $\rho^{(1,2)}$ for which $f_1(g_2(\rho^{(2)})) = \rho_0$. The case of $f_2(g_1(\rho^{(1)})) = \rho_0$ is handled in the same manner, using (5.7) and (5.10) instead of (5.6) and (5.9).

It is now clear that we can iterate this construction to produce a solution $R^* = (R_1^*, R_2^*)$ in $B(U'_1 \cap \overline{\mathbb{H}}) \cup B(U'_2 \cap \overline{\mathbb{H}})$, and, in fact, a solution in all of \mathbb{H} , by appropriately extending the solution we have in X_1 (or in X_2 , for that matter). More specifically, the solution in $B(U'_1 \cap \overline{\mathbb{H}}) \cup B(U'_2 \cap \overline{\mathbb{H}})$ is obtained by placing a suitable similar copy of the form $tU + \alpha$ under each of U'_k in the way described and then under each of the corresponding $tU'_k + \alpha$ smaller similar copies of the form $tU + \alpha$ (with a smaller value of t , of course). At the n th stage (where we regard the initial U as corresponding to the 0th stage) we have 2^n disjoint similar copies $U_l^{(n)}, 1 \leq l \leq 2^n$, of U . Let $I_l^{(n)}$ be the interval of \mathbb{R} joining the left and right endpoints of $U_l^{(n)} \cap \overline{\mathbb{H}}$. It is an immediate consequence of the shape of the U'_k that (once ϵ has been fixed) there is a constant $\delta_1 = \delta_1(\Theta) > 0$ such that for all points $p \in I_l^{(n)}$

$$\lambda(\theta(R^*(\{z : \delta_1 < \arg\{z - p\} < \pi - \delta_1, z \in U_l^{(n)} \cap \mathbb{H}\}))) \geq \delta_1.$$

Since, on U'_i , θ and R_j are bi-Lipschitz functions of each other, it follows that there are $r_1 = r_1(\Theta), \delta_2 = \delta_2(\Theta) > 0$ such that for all $p \in I_l^{(n)}$ the range of R_k^* in

$$N(r_1^n, p) \cap \{z : \delta_1 < \arg\{z - p\} < \pi - \delta_1\}$$

is an interval of length at least $\delta_2, k = 1, 2$. Let $M^{(n)} = \cup\{I_l^{(n)} : 1 \leq l \leq 2^n\}$. Then $M = \cap\{M^{(n)} : n \geq 1\}$ consists entirely of boundary singularities of the solution R^* that we have constructed. To those who have read about Cantor sets it is probably obvious that M has positive Hausdorff dimension, but for completeness we include an appropriately modified version of the argument given by Falconer [F] for the classical excluded middle third case.

First of all, it is clear from the self-similar nature of the construction that there is a number $\gamma \in (0, 1)$ such that the minimum distance between the 2^n intervals

making up the set $M^{(n)}$ is bounded below by γ^n . We show that the τ -dimensional Hausdorff measure of M is positive, where τ is defined by

$$\gamma^\tau = \frac{1}{2}.$$

Let $\{G_k\}$ be an open cover of M , which we can assume to be finite since M is compact. Let $\max\{\text{diam}(G_k)\} < 1$. For each k there is a nonnegative integer $l = l(k)$ such that

$$\gamma^{l+1} \leq \text{diam}(G_k) < \gamma^l.$$

From the definition of γ it follows that such a G_k can have a nonempty intersection with at most one of the intervals that make up $M^{(l)}$. Consequently, for $p \geq l(k)$ at most $2^{p-l(k)}$ of the intervals making up $M^{(p)}$ can have a nonempty intersection with such a G_k . For each $p \geq l = l(k)$ we therefore have

$$2^{p-l} = \frac{2^p}{2^l} = \frac{2^p(\gamma^{l+1})^\tau}{2^l\gamma^{\tau l}\gamma^\tau} \leq \frac{2^p(\text{diam}(G_k))^\tau}{(2\gamma^\tau)^l\gamma^\tau}.$$

Let p be so large that $\min\{\text{diam}(G_k)\} \geq \gamma^{p+1}$, so that $p \geq l(k)$ for all k . Since there are 2^p intervals in $M^{(p)}$, if we denote by N_k the number of intervals touched by G_k , then

$$2^p \leq \sum_k N_k \leq \sum_k \frac{2^p(\text{diam}(G_k))^\tau}{(2\gamma^\tau)^{l(k)}\gamma^\tau} = \sum_k \frac{2^p(\text{diam}(G_k))^\tau}{\gamma^\tau},$$

so that

$$\sum_k (\text{diam}(G_k))^\tau \geq \gamma^\tau = \frac{1}{2}.$$

This means that the τ -dimensional Hausdorff measure of M is at least $\frac{1}{2}$, so that M has Hausdorff dimension at least τ .

6. A few concluding remarks. We briefly discuss some of the issues and problems suggested by the foregoing. In the first place, it would be interesting to determine whether Corollary 4.3 is true for genuinely nonlinear systems (see the definitions of the terms “system” and “genuinely nonlinear” between relations (1.4) and (1.5)) which are not necessarily normal. Probably, though, some hypothesis in the spirit of (ii) of Definition 1.1 as well as some bound like

$$\epsilon < |\theta_1(R) - \theta_2(R)| < \pi - \epsilon \quad \text{for all } R \in \mathbb{R}^2$$

is necessary, so that, using the approach of section 5, or otherwise, one might try to construct a solution of a 2×2 genuinely nonlinear hyperbolic system not satisfying one or the other or both of these conditions and which has a set of boundary singularities of Hausdorff dimension 1. In a wider context one can ask if Corollary 4.3 has any counterparts for an appropriate class of sufficiently nonlinear $n \times n$ planar hyperbolic systems. In reference to normal systems, our analysis leaves open the question of whether in a half-plane \mathbb{H} , for example, there can be a solution for which the set of boundary singularities of type 1 has positive Hausdorff dimension. Corresponding to such a solution there would have to be set $A \subset \partial G$ of positive Hausdorff dimension such that for each $a \in A$ there is a characteristic C_a exiting at a but for which

hypothesis (3.22) of Proposition 3.31 does not hold. Also open is the question of whether the word “nontangential” is necessary in the conclusion of Main Theorem 2.2.

As we shall show elsewhere, the ideas of section 5 can be used to construct in an arbitrary Jordan domain cps-mappings with arbitrary (distinct) principal stretch factors which have infinitely many isolated singularities, and in fact these singularities can be of spiral type (see [G3] for the classification of isolated singularities of cps-mappings). This raises the question of how such singularities can be distributed, and in this regard we conjecture that there is some absolute constant $\gamma > 1$ such that if $\{a_n\}$ is the sequence of isolated singularities of any cps-mapping in any smoothly bounded Jordan domain for which $\{\text{dist}(a_n, \partial G)\}$ is nonincreasing, then

$$(6.1) \quad \sum_n \text{dist}(a_n, \partial D) \gamma^n < \infty.$$

Note that this was shown with $\gamma = 1$ in [G3, Corollary 4.1]. More generally, there are other 2×2 genuinely nonlinear systems for which there exist corresponding unambiguously defined nets of characteristics which have isolated singularities, and for any such system one could attempt to obtain a classification of such singularities along the lines of [G3] and seek a bound of type (6.1) on their density.

REFERENCES

- [CS] C. CARATHÉODORY AND E. SCHMIDT, *Über die Hencky-Prandtschen Kurven*, Z. Angew. Math. Mech., 3 (1923), pp. 468–475.
- [ChG] M. CHUAQUI AND J. GEVIRTZ, *Constant principal strain mappings on 2-manifolds*, SIAM J. Math. Anal., 32 (2000), pp. 734–759.
- [F] K. FALCONER, *Fractal Geometry. Mathematical Foundations and Applications*, John Wiley and Sons, Chichester, UK, 1990.
- [G1] J. GEVIRTZ, *On planar mappings with prescribed principal strains*, Arch. Rational Mech. Anal., 117 (1992), pp. 295–320.
- [G2] J. GEVIRTZ, *A diagonal hyperbolic system for mappings with prescribed principal strains*, J. Math. Anal. Appl., 176 (1993), pp. 390–403.
- [G3] J. GEVIRTZ, *Hencky-Prandtl nets and constant principal strain mappings with isolated singularities*, Ann. Acad. Sci. Fenn. Math., 25 (2000), pp. 187–238.
- [G4] J. GEVIRTZ, *Singularity sets of constant principal strain deformations*, J. Math. Anal. Appl., 263 (2001), pp. 600–625.
- [G5] J. GEVIRTZ, *Boundary values and the transformation problem for constant principal strain mappings*, Int. J. Math. Math. Sci., 2003 (2003), pp. 739–776.
- [GM] V. GUTLYANSKII AND O. MARTIO, *Rotation estimates and spirals*, Conform. Geom. Dyn., 5 (2001), pp. 6–20.
- [Hem] W. S. HEMP, *Optimum Structures*, Clarendon Press, Oxford, UK, 1973.
- [Hen] H. HENCKY, *Über einige statisch bestimmte Fälle des Gleichgewichts in plastischen Körpern*, Z. Angew. Math. Mech., 3 (1923), pp. 241–251.
- [Hi] R. HILL, *The Mathematical Theory of Plasticity*, Clarendon Press, Oxford, UK, 1964.
- [J1] F. JOHN, *On quasi-isometric mappings*, I, Comm. Pure Appl. Math., 21 (1968), pp. 77–110.
- [J2] F. JOHN, *On quasi-isometric mappings*, II, Comm. Pure Appl. Math., 22 (1969), pp. 265–278.
- [L] P. LAX, *Development of singularities of solutions of nonlinear hyperbolic partial differential equations*, J. Mathematical Phys., 5 (1964), pp. 611–613.
- [Pr] L. PRANDTL, *Über die Eindringungsfestigkeit (Härte) plastischer Baustoffe und die Festigkeit von Schneiden*, Z. Angew. Math. Mech., 1 (1921), pp. 15–20.
- [Y] W.-L. YIN, *Two families of finite deformations with constant strain invariants*, Mech. Res. Comm., 10 (1983), pp. 127–132.

ASYMPTOTIC STABILITY OF ASCENDING SOLITARY MAGMA WAVES*

GIDEON SIMPSON[†] AND MICHAEL I. WEINSTEIN[†]

Abstract. Coherent structures, such as solitary waves, appear in many physical problems, including fluid mechanics, optics, quantum physics, and plasma physics. A less studied setting is found in geophysics, where highly viscous fluids couple to evolving material parameters, modeling partially molten rock, magma, in the Earth’s interior. Solitary waves are also found here, but the equations lack useful mathematical structures such as an inverse scattering transform or even a variational formulation. A common question in all of these applications is whether or not these structures are stable to perturbation. We prove that the solitary waves in this earth science setting are asymptotically stable and accomplish this without any preexisting Lyapunov stability. This holds true for a family of equations, extending beyond the physical parameter space. Furthermore, this extends existing results on well-posedness to data in a neighborhood of the solitary waves.

Key words. solitary waves, stability, viscously deformable porous media, magma

AMS subject classifications. 74J35, 35B40, 74L05

DOI. 10.1137/080712271

1. Introduction. Coherent structures, such as solitary waves, appear in many physical problems, including fluid mechanics, optics, quantum physics, and plasma physics. A less studied setting is found in geophysics, where highly viscous fluids couple to evolving material parameters, modeling partially molten rock, magma, in the Earth’s interior. Solitary waves are also found here, but the equations lack the useful structures such as an inverse scattering transform or even a variational formulation.

A important question in all of these applications is whether or not the coherent structures are stable to perturbation. We prove that the solitary waves in this earth science setting are asymptotically stable and accomplish this without any preexisting Lyapunov stability.

1.1. Magma: Porous flow in a viscously deformable media. Models of magma in the Earth’s interior couple Stokes flow of the viscous melt to the slow, creeping deformation of the porous rock. These force-balance equations couple to transport equations for the volume fraction of melt, the *porosity*. Formulations may be found in [5, 20, 35, 36, 41]. Nonlinearity appears in fluxes and through the material properties, the *permeability* and *viscosity* of the porous, deformable rock, which depend nonlinearly on the porosity. Consequently, such models are known, from computations, to feature localization and generate coherent structures; see [1, 2, 3, 15, 35, 36, 41, 42, 43, 44, 51]. The physical assumptions and their implications will be discussed in a forthcoming review article [39].

Under certain assumptions (small fluid fraction, absence of large-scale shear, no melting, etc.), such a system reduces to a single scalar equation for the porosity’s

*Received by the editors January 2, 2008; accepted for publication (in revised form) June 20, 2008; published electronically November 5, 2008. This work was funded in part by the U.S. National Science Foundation (NSF) Collaboration in Mathematical Geosciences (CMG), Division of Mathematical Sciences (DMS) grant DMS-05-30853, NSF Integrative Graduate Education and Research Traineeship (IGERT) grant DGE-02-21041, and NSF grants DMS-04-12305 and DMS-07-07850.

<http://www.siam.org/journals/sima/40-4/71227.html>

[†]Department of Applied Physics and Applied Mathematics, Columbia University, New York, NY 10027 (grs2103@columbia.edu, miw2103@columbia.edu).

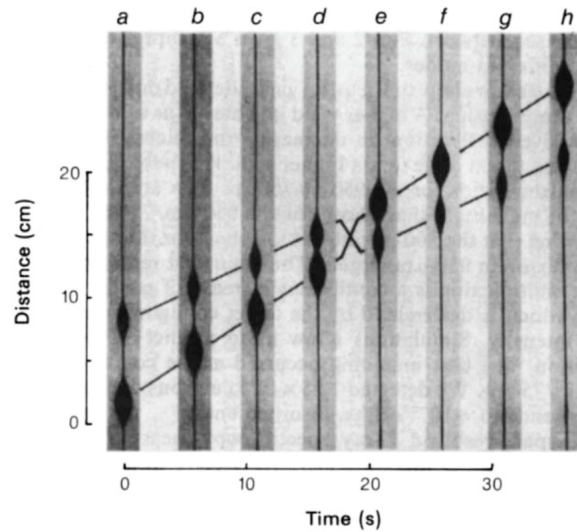


FIG. 1. *Solitary waves colliding and propagating in an experiment with honey. From Figure 1 of [37]. Reprinted by permission from Macmillan Publishers Ltd.*

evolution [2, 3, 41, 42]. The d -dimensional equation is

$$(1.1) \quad \partial_t \phi + \partial_z (\phi^n) - \nabla \cdot [\phi^n \nabla (\phi^{-m} \partial_t \phi)] = 0, \quad \mathbf{x} \in \mathbb{R}^d, t > 0,$$

with the boundary conditions that $\phi(\mathbf{x}, t) \rightarrow 1$ as $|\mathbf{x}| \rightarrow \infty$. $\nabla = (\partial_x, \partial_y, \partial_z)$ for $d = 3$ and $\nabla = (\partial_x, \partial_z)$ for $d = 2$. The nonlinearity n comes from the relationship between the permeability, K , and the porosity of the rock, $K \propto \phi^n$. m relates to the bulk viscosity, ζ , to the porosity of the rock, $\zeta \propto \phi^{-m}$. In the physical regime, these exponents have values $2 \leq n \leq 3$ and $0 \leq m \leq 1$ [12, 13, 20, 31, 35, 36, 45, 46, 52].

Equation (1.1) appears elsewhere in earth science as a model for convective mantle plumes. Manifesting themselves as *hot spots* at the surface, these plumes transport warm buoyant material. Examples include the Hawaiian Island chain and Iceland. Modeled as the flow of a viscous fluid up a conduit embedded in a higher viscosity medium, an equation of the form (1.1) was derived in [24]. There, the equation is one-dimensional ($d = 1$), the exponents are $(n, m) = (2, 1)$, and the dependent variable ϕ is the cross-sectional area of the pipe.

Numerical simulations of (1.1) in one, two, and three dimensions were performed in [2, 3, 35, 36, 51], where stable, radially symmetric, *solitary traveling waves* were observed. In [23], it was shown that in one dimension, solitary waves, $U_c(x - ct)$, in excess of the reference state, $\phi \equiv 1$, exist for $n > 1$. In the context of conduit flow, discussed in the preceding paragraph, analogue experiments using viscous syrups appear in [24, 37, 49]; robust solitary waves appeared as predicted; see Figure 1.

1.2. Stability of solitary waves. We consider (1.1) in one dimension,

$$(1.2) \quad \partial_t \phi + \partial_x (\phi^n) - \partial_x [\phi^n \partial_x (\phi^{-m} \partial_t \phi)] = 0, \quad \lim_{|x| \rightarrow \infty} \phi(x, t) = 1,$$

where the z coordinate has been relabeled x . A cursory explanation for the solitary waves' stability may be found in [50]. Under a small amplitude scaling, (1.2) is, to leading order, governed by the Korteweg–de Vries (KdV) equation. Since KdV solitons

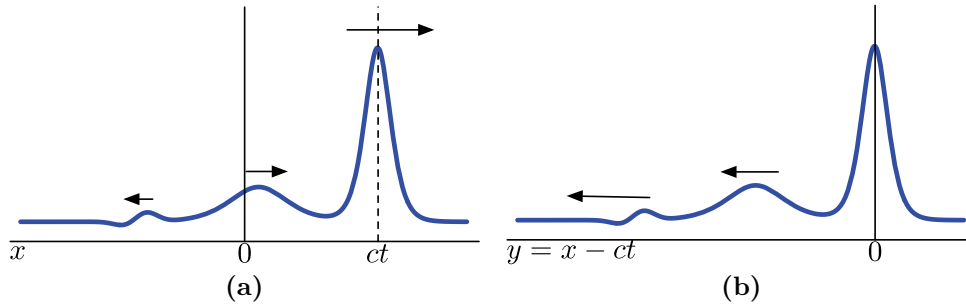


FIG. 2. The largest solitary wave in (a) travels faster than smaller solitary waves and dispersive waves behind it. In the frame of this wave, the rest of the solution appears to move leftward, as in (b).

are stable, on a timescale for which KdV approximates (1.2) its solitary waves should also be stable.

Based on observations of numerical experiments, we expect a slightly perturbed solitary wave to evolve into another wave with similar amplitude and phase. It will be accompanied by some small amplitude dispersive waves and, perhaps, another solitary wave of smaller amplitude. The leading wave will *outrun* these other disturbances, cease interacting with them, and stabilize.

Some intuition for this stability may be found in two properties. First, taller solitary waves travel with greater speed, c , than smaller ones. In the frame of the largest solitary wave, $y = x - ct$, the other waves travel leftward. Second, in the frame of the leading solitary wave, small perturbations of the reference state, $\phi(x, t) = 1 + \psi(x - ct, t)$ and $|\psi| \ll 1$, evolve under the linear flow

$$\partial_t \psi - c \partial_y \psi + n \partial_y \psi - \partial_y^2 \partial_t \psi + c \partial_y^3 \psi = 0.$$

The dispersion relation and group velocity are

$$\omega(k) = \frac{nk}{1+k^2} - ck, \quad \omega'(k) = n \frac{1-k^2}{(1+k^2)^2} - c.$$

Since the solitary waves of (1.2) that we study travel with speed $c > n > 1$, both phase and group velocities are negative for all k . Therefore, small dispersive waves also travel leftward. These two mechanisms are diagrammed in Figure 2, and we will exploit them to prove the main theorems.

As the system is conservative, perturbations, such as a small solitary wave, will not vanish in a translation invariant norm. A suitable norm will register leftward motion as decay. We will use exponentially weighted norms, in the frame of the leading solitary wave. These norms are defined in section 1.4.

Several paths to proving stability are available. One method is to seek constants of motion that can be combined into a metric centered at the solitary wave. Since the metric is time independent, if the perturbation is initially small, it will remain so. This elegant method relies on the calculus of variations, and for equations such as KdV and nonlinear Schrödinger (NLS) it may be used to prove *orbital stability* [4, 8, 47, 48]. A solitary wave, U_c , is said to be orbitally stable if for data sufficiently close to it

$$\inf_{y \in \mathbb{R}} \|u(t) - U_c(\cdot + y)\|_X < \delta \quad \text{for all } t$$

for some $\delta > 0$ and an appropriate norm $\|\cdot\|_X$. Typically, the norm is equivalent to $L^2(\mathbb{R})$ or $H^1(\mathbb{R})$.

However, in general, (1.2) lacks a sufficient number of conservation laws for this approach. Indeed, in [3], the authors searched for an additional conservation law, in hopes of proving orbital stability [7]. There are many such equations, including some of the Boussinesq systems [29] and many of the “compacton” equations [32], which also lack such structure, yet appear, in numerical experiments, to possess stable solitary waves. Note that when $n + m = 0$, (1.2) is Hamiltonian, and we have investigated this, proving orbital stability in [40]. We wish to consider the general case, which includes the physically interesting cases $(n, m) = (2, 1)$ and $(3, 0)$.

Another approach to stability is to linearize the problem $u_t = N(u, u_x, u_{xx}, \dots)$ about a solitary wave and establish linear stability. Then one seeks a way to perturbatively “boost” this to prove stability for the nonlinear flow. This may rely on direct spectral analysis of the linearized evolution operator. We employ this method, following the work of Pego and Weinstein [28, 29] and Miller and Weinstein [21]. Through this, we prove that the solitary waves are *asymptotically stable*, our main result. By this we mean that in an appropriate norm, $\|\cdot\|_Y$,

$$\|u(t) - U_c\|_Y \rightarrow 0 \quad \text{as } t \rightarrow +\infty$$

for data sufficiently close to U_c .

We note that another method recently appeared in the work of Martel and Merle [19]. Without linearizing, the authors employ a virial inequality to directly prove asymptotic stability of generalized KdV solitary waves in $H^1_{\text{loc}}(\mathbb{R})$.

Our problem is an example of an equation for which one can prove asymptotic stability in the absence of orbital stability. Upon reflection, it is clear that the asymptotic stability of generalized KdV and Benjamin–Bona–Mahony (BBM) solitary waves could have been proven without using the orbital stability results.

1.3. Main results and outline. The main results are as follows.

THEOREM 1.1. *There exists $c_* > n$ such that for all $c_0 \in (n, c_*]$, if $\phi_{c_0}(x - \theta_0)$ is a solitary wave solution of (1.2), then there exist constants $K_* > 0$, $\epsilon_* > 0$, and $a > 0$, such that for $\epsilon \leq \epsilon_*$, if*

$$\|v_0\|_{H^1} + \|e^{ax}v_0\|_{H^1} \leq \epsilon,$$

then

- (a) (1.2) has a solution with data $\phi_0(x) = \phi_{c_0}(x + \theta_0) + v_0(x)$ for all time;
- (b) there exist $c_\infty, \theta_\infty, K_*$, and $\kappa > 0$ such that

$$(1.3) \quad \|\phi(\cdot, t) - \phi_{c_\infty}(\cdot - c_\infty t + \theta_\infty)\|_{H^1} \leq K_*\epsilon,$$

$$(1.4) \quad \|e^{ax}[\phi(\cdot + c_\infty t - \theta_\infty, t) - \phi_{c_\infty}(\cdot)]\|_{H^1} \leq K_*\epsilon e^{-\kappa t},$$

$$(1.5) \quad |c_\infty - c_0| + |\theta_\infty - \theta_0| \leq K_*\epsilon.$$

COROLLARY 1.2. *Let $n + m = 0$. If $\partial_c \mathcal{N}[\phi_c] > 0$, \mathcal{N} defined in (2.8), then Theorem 1.1 holds for all $c > c_*$, except for a discrete set with no accumulation point.*

Remark 1.3. For general $n > 1$, Theorem 1.1 is limited to $c \leq c_*$ because we are only able to rigorously treat the spectrum of a linear operator perturbatively. This can be extended by a numerical computation of the spectrum. See sections 3.3–3.4.

The feature of (1.2) that allows us to prove nonlinear stability from the linear stability is a nonnegative invariant, denoted $\mathcal{N}[\phi]$, and defined in (2.8). The Taylor expansion of \mathcal{N} about a solitary wave is

$$\mathcal{N}[\phi_c + v] = \mathcal{N}[\phi_c] + \langle \delta\mathcal{N}[\phi_c], v \rangle + \langle \delta^2\mathcal{N}[\phi_c]v, v \rangle + O(\|v\|_{H^1}^3).$$

The first variation does not vanish and the second variation is not a positive definite quadratic form. However, in the frame of the solitary wave, the perturbation v is migrating to $-\infty$. Therefore

$$\langle \delta\mathcal{N}[\phi_c], v \rangle \rightarrow 0 \quad \text{as } t \rightarrow +\infty.$$

The second variation may be decomposed as $\delta^2\mathcal{N}[\phi_c] = P + Q$, P a positive quadratic form and $Q = Q(x)$ a localized function. Then since the perturbation moves leftward

$$\begin{aligned} \langle Pv, v \rangle &\geq \kappa^2 \langle v, v \rangle, \\ \langle Qv, v \rangle &\rightarrow 0 \quad \text{as } t \rightarrow +\infty. \end{aligned}$$

Asymptotically,

$$\|v\|_{H^1} \leq K\Delta\mathcal{N}, \quad \Delta\mathcal{N} = |\mathcal{N}[\phi_c + v] - \mathcal{N}[\phi_c]|,$$

giving a Lyapunov-type bound on the perturbation. Control of $\Delta\mathcal{N}$ is essential in using this estimate to prove nonlinear stability. Though we rely on the invariance of \mathcal{N} , if $\Delta\mathcal{N}$ were merely bounded in terms of the data, it would be sufficient.

However, more is needed to formalize this into a proof, notably a sense in which the perturbation recedes from the solitary wave. This is accomplished by analyzing the spectrum of the linearized evolution operator in a weighted space, in which the perturbation will decay.

The plan of the proof is as follows:

- (I) In section 2 we review properties of (1.2) and establish regularity properties of the solitary waves.
- (II) In section 3, we prove that the linearized operator, A_a , has the property that there exists $\varepsilon > 0$ such that

$$\sigma(A_a) \cap \{\Re\lambda \geq -\varepsilon\} = \{0\}$$

and zero is an eigenvalue of algebraic multiplicity two.

- (III) In section 4, we prove

$$\|w(t)\|_{H^1} = \|e^{A_a t} w_0\|_{H^1} \leq K e^{-bt} \|w_0\|_{H^1}$$

for appropriate w_0 , K , and b positive constants.

- (IV) In section 5, we make several estimates, including a formalization of the Lyapunov bound. We also formulate equations for the speed and phase parameters of the solitary wave $(c(t), \theta(t))$, coupling them to the infinite-dimensional system for the perturbation.
- (V) In section 6, we prove the main results, asymptotic stability and global existence of data near a solitary wave solution.

Some remarks are made in section 7, and additional details are located in the appendices.

1.4. Notation. Generic constants will typically be denoted by the capital letters K , M , and N , sometimes with tildes, overlines, or primes. Subscripts, such as M_γ , may appear in order to indicate that M depends on γ . We avoid using C as a generic constant, as c appears throughout the paper as the speed parameter, and an operator $C(\lambda)$ appears in section 3.

Functions will typically live in spaces $H^k(\mathbb{R}) = W^{1,k}(\mathbb{R})$, k a nonnegative integer, the spaces of square integrable functions with square integrable (weak) derivatives up to order k . We will frequently omit writing \mathbb{R} . The $L^p(\mathbb{R})$ spaces will also appear, in particular L^2 and L^∞ . While we write $\|f\|_{H^k}$ for the norm of a function in H^k , we only write $\|f\|_p$ for the norm of a function in L^p .

We will be interested in functions in the exponentially weighted space,

$$H_a^k = \{u : e^{ax}u(x) \in H^k\}$$

for $k = 0, 1, 2, \dots$, and $a > 0$ with associated norm

$$\|u\|_{H_a^k} = \|e^{ax}u\|_{H^k}.$$

We also define the norm

$$\|f\|_{H^1 \cap H_a^1} \equiv \|f\|_{H^1} + \|f\|_{H_a^1}.$$

The exponential weight will *always* be a positive number; we will often omit the assumption $a > 0$ in statements.

Frequently, we will have an operator T defined on a weighted space, H_a^k , but wish to make computations in the unweighted space. To T we associate $T_a = e^{ax}Te^{-ax}$, an operator on H^k . For the differentiation operator, $\partial_x \mapsto D_a = \partial_x - a$.

2. Preliminaries.

2.1. Properties of the equation in a weighted space. Much of the analysis involves studying (1.2) in an exponentially weighted space. We therefore state the following extension of the well-posedness results obtained in [38] for H^k spaces.

THEOREM 2.1 (local existence in time and continuous dependence upon data). *Given $0 < a < 1$, let $\phi_0(x)$ satisfy*

$$(2.1) \quad \|\phi_0 - 1\|_{H^1 \cap H_a^1} \leq R < \infty,$$

$$(2.2) \quad \inf_x \phi_0(x) \geq \alpha_0 > 0,$$

$$(2.3) \quad \inf_x \phi_0(x)^m - a^2 \phi_0(x)^n \geq \beta_0 > 0.$$

Then there exist $T_{\text{local}} > 0$ and $\phi(x, t) - 1 \in C^1([0, T_{\text{local}}), H^1 \cap H_a^1)$, a solution of (1.2) with data ϕ_0 , satisfying

$$(2.4) \quad \|\phi(\cdot, t) - 1\|_{H^1 \cap H_a^1} \leq 2R,$$

$$(2.5) \quad \inf_x \phi(x, t) \geq \frac{1}{2}\alpha_0,$$

$$(2.6) \quad \inf_x \phi(x, t)^m - a^2 \phi(x, t)^n \geq \frac{1}{2}\beta_0$$

for $t < T_{\text{local}}$.

Moreover, there is a maximal time of existence T_{\max} , such that if $T_{\max} < \infty$, then

$$(2.7) \quad \lim_{t \rightarrow T_{\max}} \|\phi(\cdot, t) - 1\|_{H^1 \cap H_a^1} + \left\| \frac{1}{\phi(\cdot, t)} \right\|_{\infty} + \left\| \frac{1}{\phi(\cdot, t)^m - a^2 \phi(\cdot, t)^n} \right\|_{\infty} = \infty.$$

Remark 2.2. When $a = 0$, (2.3) is unnecessary; this case was treated in [38]. The importance of this condition for $a > 0$ will be discussed in section 2.3.

THEOREM 2.3. *Given $0 < a < 1$, let $\phi^{(j)} - 1 \in C^1([0, T]; H^1 \cap H_a^1)$, $j = 1, 2$, be two solutions of (1.2) such that*

$$\begin{aligned} \|\phi^{(j)}(\cdot, t) - 1\|_{H^1 \cap H_a^1} &\leq R < \infty, \\ \inf_x \phi^{(j)}(x, t) &\geq \alpha_0 > 0, \\ \inf_x (\phi(x, t)^{(j)})^m - a^2 (\phi(x, t)^{(j)})^n &\geq \beta_0 > 0. \end{aligned}$$

There exists a constant $K = K(R, \alpha_0, \beta_0, a)$, such that

$$\|\phi^{(1)}(\cdot, t) - \phi^{(2)}(\cdot, t)\|_{H^1 \cap H_a^1} \leq e^{Kt} \|\phi_0^{(1)} - \phi_0^{(2)}\|_{H^1 \cap H_a^1} \quad \text{for } t \leq T.$$

Additionally, (1.2) possesses the conservation law

$$(2.8) \quad \mathcal{N}[\phi] = \begin{cases} \int \left(\frac{1}{2} \phi^{-2m} \phi_x^2 + \phi \log(\phi) - \phi + 1 \right) dx & \text{if } n + m = 1, \\ \int \left(\frac{1}{2} \phi^{-2m} \phi_x^2 + \phi - 1 - \log(\phi) \right) dx & \text{if } n + m = 2, \\ \int \left(\frac{1}{2} \phi^{-2m} \phi_x^2 + \frac{\phi^{2-n-m} - 1 + (n+m-2)(\phi-1)}{(n+m-1)(n+m-2)} \right) dx & \text{for all other } n \text{ and } m. \end{cases}$$

\mathcal{N} is well defined for ϕ bounded from below away from zero and $\|\phi - 1\|_{H^1} < \infty$. It is also locally convex about $\phi \equiv 1$. See section 3 of [38] for details.

Remark 2.4. The physical significance of \mathcal{N} remains elusive. In the Hamiltonian case, $n + m = 0$, \mathcal{N} revealed itself to be the generalized momentum of the equation [40].

2.2. Solitary waves and their analytic properties. Let us review the properties of the solitary waves associated with (1.2). In particular, we identify their decay and regularity properties and introduce the KdV scaling for later use.

Substituting the traveling wave ansatz, $\phi_c(x, t) = \phi_c(x - ct)$, into (1.2) with boundary conditions

$$(2.9) \quad \lim_{y \rightarrow \pm\infty} \phi_c(y) = 1, \quad \lim_{y \rightarrow \pm\infty} \partial_y^j \phi_c(y) = 0 \quad \text{for } j = 1, 2, \dots$$

we have, after one integration,

$$(2.10) \quad -c(\phi_c - 1) + \phi_c^n - 1 + c\phi_c^n \partial_y (\phi_c^{-m} \partial_y \phi_c) = 0.$$

Letting $u_c = 1 - \phi_c$, u_c satisfies

$$(2.11) \quad -cu_c + (u_c + 1)^n - 1 + c(u_c + 1)^n \partial_y \left((u_c + 1)^{-m} \partial_y u_c \right) = 0.$$

Equation (2.10) may also be integrated up to a first order equation,

$$(2.12) \quad \frac{1}{2} \phi_c^{-2m} (\partial_x \phi_c)^2 - F_1(\phi_c; c) = 0,$$

after applying the boundary condition $\phi_c \rightarrow 1$ at $\pm\infty$. F_1 depends on the particular exponents:

$$(2.13) \quad F_1(x; c) = \begin{cases} \frac{x^{1-n}-1}{1-n} + (1 - \frac{1}{c}) \frac{x^{-n}-1}{n} - \frac{1}{c} \log(x) & \text{if } m = 1, \\ x - 1 - (1 - \frac{1}{c}) \log(x) - \frac{1}{c} \frac{x^n-1}{c} & \text{if } n + m = 1, \\ \log(x) + (1 - \frac{1}{c}) (\frac{1}{x} - 1) - \frac{1}{c} \frac{x^{n-1}-1}{n-1} & \text{if } n + m = 2, \\ \frac{x^{2-n-m}-1}{2-n-m} - (1 - \frac{1}{c}) \frac{x^{1-n-m}-1}{1-n-m} - \frac{1}{c} \frac{x^{1-m}-1}{1-m} & \text{otherwise.} \end{cases}$$

Using (2.12), an equivalent second order, self-adjoint equation for the solitary waves is

$$(2.14) \quad F_2(\phi_c; c) = -\partial_x^2 \phi_c,$$

$$(2.15) \quad F_2(x; c) = x^{m-n} [-(x - 1) + c^{-1} (x^n - 1) - 2mx^{n-m-1} F_1(x; c)].$$

Let us introduce the KdV scaling. Define

$$(2.16) \quad \boxed{\gamma = \sqrt{\frac{c-n}{c}}}.$$

Applying the scalings,

$$(2.17) \quad \xi = \gamma(x - ct), \quad u_c(y) = \frac{\gamma^2}{n-1} U(\xi(y); \gamma),$$

(2.11) becomes

$$(2.18) \quad -U + \frac{1}{2} U^2 + \partial_\xi^2 U = O(\gamma^2).$$

Remark 2.5. The parameter γ , (2.16), will be used throughout the paper. Because it uniquely maps $c \in (n, \infty)$ onto $(0, 1)$, we will use c and γ interchangeably.

We summarize what is known about (2.11) and (2.18) in the following two results.

THEOREM 2.6. *For any $c > n > 1$, (2.11) has a unique positive, even solution u_c , going to zero at $\pm\infty$. In the KdV scaling, (2.18), U is real analytic in the arguments $(\xi, \gamma) \in \mathbb{R} \times [0, 1)$. When $\gamma = 0$*

$$U(\xi; 0) = U_\star(\xi) = 3\text{sech}^2\left(\frac{1}{2}\xi\right).$$

Furthermore, for γ in any compact subset of $[0, 1)$

$$\partial_\xi^j U(\xi; \gamma) e^{\pm\xi} (\text{sign}(\xi))^j \rightarrow K_j(\gamma) \quad \text{as } \xi \rightarrow \pm\infty \text{ for } j = 0, 1, 2.$$

COROLLARY 2.7. *Given a compact interval $[0, \gamma_0] \subset [0, 1)$, there exists a constant K such that for all $\gamma \in [0, \gamma_0]$,*

$$(2.19) \quad |\partial_\xi^j U(\xi; \gamma)| \leq K e^{-|\xi|} \quad \text{for } j = 0, 1, 2, \quad -\infty < \xi < \infty,$$

$$(2.20) \quad |\partial_y^j u_c(y)| \leq K \frac{\gamma^{2+j}}{n-1} e^{-\gamma|y|} \quad \text{for } j = 0, 1, 2, \quad -\infty < y < \infty.$$

Proof. From [23], solitary waves exist, provided $c > n > 1$ and $m \in \mathbb{R}$. Writing the problem as a two-dimensional dynamical system, we may apply the stable manifold theorem about the hyperbolic point $(0, 0)$ to deduce the exponential decay, as in [29, Theorem 2.1, Corollary 2.2]. \square

Remark 2.8. When the parameter γ is small, ϕ_c is in the regime of small amplitude, long waves, where KdV appears as the leading order equation in a perturbation expansion of (1.2), as in [50].

COROLLARY 2.9. *For each $c > n > 1$, the solitary wave solution, $\phi_c - 1$, lies in $H^\infty(\mathbb{R})$. Furthermore, there exists $\sigma_0 > 0$ such that the solitary wave ϕ_c may be analytically continued off the real axis into the strip $\{z : |\Im z| < \sigma_0\}$.*

Proof. This is a consequence of Corollary 2.7 and [9, Corollary 4.1.6]; see Appendix A.1. \square

COROLLARY 2.10. *Given a solitary wave ϕ_c , $c > n$, assume $0 < a < \gamma$. Then $\phi_c - 1 \in H_a^\infty$.*

COROLLARY 2.11. *Let $n > 1$.*

- (a) *The mapping $c \mapsto \phi_c - 1$ is $C^1((n, \infty); H^2)$. In fact the mapping is analytic.*
- (b) *The mapping $c \mapsto \phi_c - 1$ is analytic, and, for fixed x , $c \mapsto \phi_c(x)$ is an analytic function of c .*
- (c) *Given $a < \frac{1}{2}$, the mapping is also $C^1((n/(1 - 4a^2), \infty); H^2 \cap H_a^2)$.*

Proof. All parts are proved using the implicit function theorem, applied to the functional

$$\mathcal{F}[c, u] = \partial_x^2 u + F_2(1 + u; c).$$

See Appendix A.2 for details and [6] for a statement and proof of the implicit function theorem for analytic mappings. \square

2.3. Remarks and assumptions on the exponential weight. We see in Theorem 2.1 and Corollary 2.11 that the particular exponential weight restricts what data and which solitary waves will be permissible. For the solitary wave result, this restriction comes from the decay rate associated with the speed; see Corollary 2.7.

In the case of the existence theorem, (1.2) may be written as

$$\begin{aligned} \partial_t \phi &= - \{I - \partial_x [\phi^n \partial_x (\phi^{-m} \cdot)]\}^{-1} \partial_x (\phi^n) \\ (2.21) \quad &= -\phi^m \left\{ \phi^{-m} [I - \partial_x (\phi^n \partial_x (\phi^{-m} \cdot))]^{-1} \right\} \partial_x (\phi^n) \\ &= -\phi^m H_\phi^{-1} \partial_x (\phi^n). \end{aligned}$$

The operator H_ϕ is

$$(2.22) \quad H_\phi = \phi(x)^m - \partial_x (\phi(x)^n \partial_x \cdot).$$

This is a bounded operator on $L^2 \rightarrow H^1$, provided ϕ is continuous and bounded from below away from zero. However, the exponential weight introduces a second constraint. Consider solving $H_\phi u = f$, $f \in L_a^2$, for $u \in H_a^1$. Letting $g = e^{ax} f$ and $v = e^{ax} u$, this is equivalent to solving

$$[\phi^m - D_a (\phi^n D_a)] v = g, \quad v \in H^1.$$

Multiplying by v and integrating by parts,

$$\int (\phi^m - a^2 \phi^n) v^2 + \phi^n (\partial_x v)^2 dx = \int g v dx.$$

A unique solution exists, provided a and ϕ satisfy $\inf_x \phi(x)^m - a^2\phi(x)^n > 0$; this is condition (2.3).

We invert these restrictions; given a solitary wave of speed c , we will assume that a is sufficiently small so that these, and other, properties hold. There are three restrictions in what follows.

Let

$$(2.23) \quad a_1 = \frac{1}{3}\gamma(c).$$

$\phi_c - 1$ will then be in $H_{a_1}^2$ and $H_{2a_1}^2$, as will all solitary waves of nearby speed. Let

$$(2.24) \quad a_2 = \frac{1}{2} \inf_x \phi_c(x)^{(m-n)/2}.$$

Then for an $a \leq \min \{a_1, a_2\}$, ϕ_c will satisfy (2.3), with

$$\inf_x \phi_c(x)^m - a^2\phi_c(x)^n \geq \frac{3}{4} \inf_x \phi_c(x)^m > 0.$$

Hence, the solitary waves will live in a set on which the existence theorem applies. Moreover, for all data ϕ_0 sufficiently close to ϕ_c in the H^1 norm, an analogous lower bound will exist.

Remark 2.12. $\inf_x \phi_c(x)^{(m-n)/2}$ is related to a physical length scale known as the *compaction length* [20]. This length, $\delta_{\text{comp.}}$, is given by

$$\delta_{\text{comp.}}(x) = \sqrt{\phi(x)^{n-m}}.$$

It measures the distances over which there will be geometrical rearrangement of the material, appearing macroscopically as changes in ϕ , in response to viscous stresses.

The stipulation $a < \inf_x \phi(x)^{(m-n)/2}$ may be interpreted as requiring the length scale associated with the exponential weight, a^{-1} , to never be smaller than this intrinsic, spatially varying length. However, we draw no physical conclusions from this observation.

Finally, there is a constraints related to the *essential spectrum* of a linear operator, discussed in section 3. Let

$$(2.25) \quad a_3 = \sqrt{\frac{2c}{n + 2c + \sqrt{n(n + 8c)}}} \sqrt{\frac{c - n}{c}}.$$

This will ensure that for any $a \leq a_3$, the essential spectrum is located in a specific part of the complex plane. Let

$$(2.26) \quad a_*(c) = \min \left\{ a_1(c), a_2(c), \frac{1}{2}a_3(c) \right\}.$$

Then for any $a \leq a_*(c)$, all of these properties will be satisfied for ϕ sufficiently close in H^1 to ϕ_c .

2.4. Ansatz and linearization. Given a perturbed solitary wave solution, ϕ , of (1.2), assume that there exists decomposition of ϕ into a (time-dependent) solitary wave of some speed $c(t)$ and phase $\theta(t)$ and a perturbation, v ; this decomposition's

existence will be proved in section 5.1. Transforming our coordinate system into the frame of this modulating solitary wave,

$$(2.27) \quad y(x, t) = x - \int_0^t c(s)ds + \theta(t),$$

$$(2.28) \quad \phi(x, t) = \phi_{c(t)}(y(x, t)) + v(y(x, t), t) = \phi_c(y, t) + v(y, t).$$

The perturbation, v , is governed by

$$(2.29) \quad \partial_t v = A_c v - \dot{\theta} \partial_y v - \dot{c} \partial_c \phi_c - \dot{\theta} \partial_y \phi_c + \mathcal{F}_1[v; \phi_c],$$

where

$$(2.30) \quad A_c v = \phi_c^m H_{\phi_c}^{-1} \partial_y L_c v,$$

$$(2.31) \quad L_c v = -c \phi_c^n \partial_y^2 (\phi_c^{-m} v) + [c - n \phi_c^{-1} + cn(\phi_c^{-1} - 1)] v,$$

and (B.6) gives the definition of $\mathcal{F}_1[v; \phi_c]$, composed of terms nonlinear in v . We make two remarks about (2.29). First, the linear operator A_c is time dependent; $c = c(t)$, and we would prefer to work with a time-independent linear operator. Second, the appearance of the term $\dot{\theta} \partial_y v$ will prove problematic for studying the equation in H^1 .

The first problem is addressed by adding and subtracting A_{c_0} , and considering the difference $A_c - A_{c_0}$ as another perturbation of the linear flow. To remove the $\partial_y v$ term, we introduce a renormalized time,

$$(2.32) \quad \tau = c_0^{-1} \left[\int_0^t c(s)ds - \theta(t) \right].$$

The asymptotic stability proof of BBM required a similar transformation [21]. In addition, the problem will be considered in the weighted space H_a^1 , with $w(y, t) = e^{ay} v(y, t)$. All together, we have the following.

PROPOSITION 2.13 (perturbation equation). *The perturbation to a solitary wave of speed c_0 associated with the ansatz in (2.28), v , and its weighted perturbation $w = e^{ay} v$ evolves according to the equations in t - and τ -time, respectively.*

$$(2.33) \quad \partial_t v = A_c v - \dot{\theta} \partial_y v - \dot{c} \partial_c \phi_c - \dot{\theta} \partial_y \phi_c + \mathcal{F}_1[v; \phi_c],$$

$$(2.34) \quad \partial_t w = A_{c,a} w - \dot{\theta} (\partial_y - a) w - e^{ay} (\dot{c} \partial_c \phi_c + \dot{\theta} \partial_y \phi_c) + \mathcal{G}_1[w, v; \phi_c],$$

$$(2.35) \quad \partial_\tau v = A_{c_0} - \frac{c_0}{c - \theta} (c \partial_c \phi_c + \dot{\theta} \partial_y \phi_c) + S[c_0, c, \dot{\theta}] v + \frac{c_0}{c - \theta} \mathcal{F}_1[v; \phi_c],$$

$$(2.36) \quad \begin{aligned} \partial_\tau w &= A_{c_0,a} w - \frac{c_0}{c - \theta} e^{ay} (c \partial_c \phi_c + \dot{\theta} \partial_y \phi_c) + S_a[c_0, c, \dot{\theta}] w + \frac{c_0}{c - \theta} \mathcal{G}_1[w, v; \phi_c] \\ &= A_{c_0,a} w + \mathcal{G} [w, v; c_0, c, \dot{\theta}]. \end{aligned}$$

The operator S and the terms \mathcal{F}_1 and \mathcal{G}_1 are given explicitly in Appendix B.

Proof. \mathcal{G}_1 is obtained from \mathcal{F}_1 by substituting $e^{-ay} w$ for one of the v 's; $e^{ay} \mathcal{F}_1[v; \phi_c] = \mathcal{G}_1[e^{ay} v, v; \phi_c]$. The details appear in Appendix B. \square

Note: From here on, we assume c_0 to be fixed and will suppress its appearance in the linear operators A_{c_0} and $A_{c_0,a}$.

3. Spectral properties of the linearized operator. In this section we analyze the spectrum of A and A_a . We will use c in place of c_0 and x in place of y as the independent variable.

$$(3.1) \quad AY = \lambda Y,$$

$$(3.2) \quad A = \{I - \partial_x [\phi_c^n \partial_x (\phi_c^{-m} \cdot)]\}^{-1} \partial_x L_c,$$

$$(3.3) \quad L_c = -c\phi_c^n \partial_x^2 (\phi_c^{-m} \cdot) + [c - n\phi_c^{-1} + cn(\phi_c^{-1} - 1)],$$

$$(3.4) \quad A_a = e^{ax} A e^{-ax}.$$

Our goal is to prove that the linearized problem is asymptotically stable: $\|e^{A_a t} w_0\|_{H^1} \rightarrow 0$ as $t \rightarrow +\infty$. We will actually show something much stronger, that this convergence to zero happens exponentially fast. Our strategy is that of Pego and Weinstein [28, 29] and Miller and Weinstein [21]. We will

1. identify the essential spectrum of A_a by showing it to be a relatively compact perturbation of a constant coefficient operator;
2. rule out the point spectrum (eigenvalues of finite multiplicity) of A_a for $|\lambda|$ sufficiently large via an operator estimate;
3. use the *Evans function*, an infinite-dimensional analogue of the characteristic polynomial, to rule out nonzero point spectra of A_a in the set of “small” λ , which will be compact;
4. show decay in time of the C_0 -semigroup $e^{A_a t}$ associated with A_a .

The spectral analysis is handled in this section, and the semigroup theory in the following section.

The principle result of this section is as follows.

THEOREM 3.1 (spectrum of linearized operator).

- (a) *Let $a \in (0, a_*(\gamma))$. The essential spectrum of A_a denoted by $\sigma_{\text{ess}}(A_a)$ is a curve lying in the open left half-plane, with rightmost point $-\omega$,*

$$(3.5) \quad -\omega = \max\{\Re z \mid z \in \sigma_{\text{ess}}(A_a)\} < 0.$$

- (b) *There exist $\gamma_* \in (0, 1)$ and $\hat{\nu} > 0$ such that for each $\gamma \in (0, \gamma_*]$ and $a \in (\gamma\hat{\nu}, a_*(\gamma))$, there exists $\varepsilon(\gamma, a) > 0$ such that the only eigenvalue of A_a with $\Re \lambda \geq -\varepsilon$ is $\lambda = 0$, and this is an eigenvalue of algebraic multiplicity two.*

- (c) *In the Hamiltonian case, $n + m = 0$, part (b) may be extended for $\gamma \in (\gamma_*, 1)$ to all but a discrete set with no accumulation point.*

The spectrum is pictured in Figure 3.

3.1. Essential spectrum. We make use of the definition of the essential spectrum of an operator from [33, 34], which states that for a closed, densely defined operator A on a Banach space X ,

$$(3.6) \quad \sigma_{\text{ess}}(A) = \bigcap_{C \in \mathcal{K}(X)} \sigma(A + C),$$

where $\mathcal{K}(X)$ denotes the set of compact operators on X . Other definitions are possible and well known; see Chapter IX of [10] for a discussion of how these definitions relate to one another.

$\sigma(A) \setminus \sigma_{\text{ess}}(A)$ then consists of point spectra. This is so because, by Theorem 7.27 of [34], if λ is not in the essential spectrum, then $\lambda I - A$ is Fredholm with index zero.

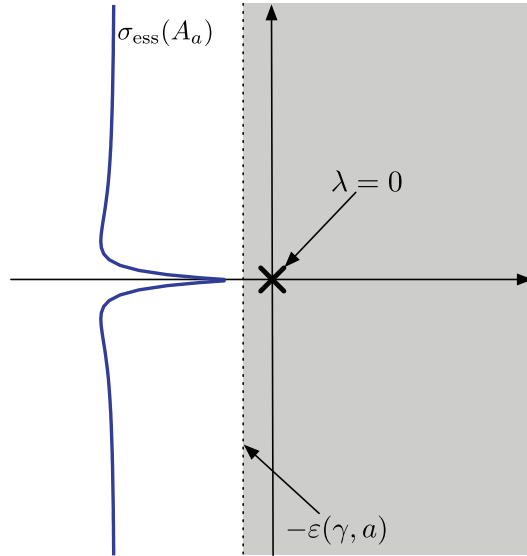


FIG. 3. The spectrum of the operator A_a . The only eigenvalue with $\Re\lambda \geq -\varepsilon$ is $\lambda = 0$.

Hence, it has a closed range and a finite kernel. Therefore, it must be an eigenvalue of finite multiplicity.

To prove Theorem 3.1 (a), we express A_a as a perturbation of a constant coefficient operator, A_a^∞ , obtained by setting $\phi_c(x)$ equal to its asymptotic state, 1:

$$(3.7) \quad A_a^\infty = (I - D_a^2)^{-1} D_a (-cD_a^2 + c - n).$$

The difference between A_a and A_a^∞ , given explicitly in (B.16), may be shown to be an A_a^∞ -compact operator. Hence, A_a is a relatively compact perturbation of A_a^∞ and

$$\sigma_{\text{ess}}(A_a) = \sigma_{\text{ess}}(A_a^\infty).$$

Upon examination of the Fourier symbol of A_a^∞ , the essential spectrum of A_a is

$$(3.8) \quad \sigma_{\text{ess}}(A_a) = \left\{ \frac{(i\ell - a)(-c(i\ell - a)^2 + c - n)}{1 - (i\ell - a)^2}, \ell \in \mathbb{R} \right\}$$

and

$$-\omega = \max \{ \Re z \mid z \in \sigma_{\text{ess}}(A_a) \} = -ac + \frac{an}{1 - a^2} < 0.$$

$\sigma_{\text{ess}}(A_a)$ lies in the open left half-plane if $0 < a < \gamma$. This set is pictured in Figure 3.

In addition, for $0 < a \leq a_* < a_3$, a_3 defined by (2.25), the spectrum moves rightward as $a \rightarrow 0$. This is because a_3 is the value for which $-\omega$ is leftmost in \mathbb{C} , maximizing the rate of decay of $e^{\alpha y}$ as $y \rightarrow -\infty$.

3.2. Large eigenvalues. As in [21], we will study the eigenvalues of A_a by considering separately a large $|\lambda|$ regime and a small $|\lambda|$ regime.

Rewriting the linear operator A_a as

$$(3.9) \quad \begin{aligned} A_a &= c\phi_c^m D_a (\phi_c^{-m} \cdot) - n\phi_c^m H_{\phi_c, a}^{-1} D_a (\phi_c^{-1} \cdot) \\ &\quad + cm\phi_c^m H_{\phi_c, a}^{-1} (\phi_c^{-1} \partial_y \phi_c \cdot) + cn\phi_c^m H_{\phi_c, a}^{-1} D_a [(\phi_c^{-1} - 1) \cdot], \end{aligned}$$

we note that λ is an L^2 eigenvalue of A_a if and only if it is also an L^2 eigenvalue of $\tilde{A}_a = \phi_c^{-m} A_a \phi_c^m$, given by

$$(3.10) \quad \begin{aligned} \tilde{A}_a &= cD_a - nH_{\phi_c, a}^{-1} D_a (\phi_c^{m-1} \cdot) \\ &\quad + cmH_{\phi_c, a}^{-1} (\phi_c^{m-1} \partial_y \phi_c \cdot) + cnH_{\phi_c, a}^{-1} D_a [\phi_c^m (\phi_c^{-1} - 1) \cdot]. \end{aligned}$$

We will rule out eigenvalues of \tilde{A}_a , thus ruling them out for A_a . This is equivalent to studying A_a in a space weighted by $\phi_c^{-m}(x)$, a strictly positive, smooth, and bounded function. \tilde{A}_a is also a relatively compact perturbation of $\tilde{A}_a^\infty = A_a^\infty$; they share the same essential spectrum.

Let the operator $C(\lambda)$ satisfy

$$C(\lambda) = (\lambda I - \tilde{A}_a^\infty)^{-1} (\tilde{A}_a - \tilde{A}_a^\infty).$$

PROPOSITION 3.2.

- (a) *The operator $C(\lambda)$ is compact for λ not in σ_{ess} . In particular, $C(\lambda)$ is compact for all λ with $\Re \lambda > -\omega$.*
- (b) *For any $\lambda \in \mathbb{C} \setminus \sigma_{\text{ess}}(A_a)$, we have that λ is an eigenvalue of A_a if and only if $1 \in \sigma(C(\lambda))$.*
- (c) *Let $\lambda \in \mathbb{C} \setminus \sigma_{\text{ess}}(A_a)$. A sufficient condition for λ not to be an eigenvalue of A_a is that $\|C(\lambda)\| < 1$, with norm either L^2 or H^1 , depending on which space is under consideration.*

Proof. Using the equivalence of eigenvalues of A_a and \tilde{A}_a , parts (b) and (c) will follow once (a) is established; see [21]. The Fourier symbol of $(\lambda I - \tilde{A}_a^\infty)^{-1}$ is

$$\frac{1 - (i\ell - a)^2}{\lambda(1 - (i\ell - a)^2) - (i\ell - a)(-c(i\ell - a)^2 + c - n)}.$$

This operator is bounded for λ not in the essential spectrum. The difference, given explicitly in (B.17), is a sum of Hilbert–Schmidt compact operators composed with bounded operators on $L^2 \rightarrow L^2$; hence $C(\lambda)$ is compact on this space.

For $C(\lambda)$ to be a compact operator on $H^1 \rightarrow H^1$, it will be sufficient to prove that

$$(I - \partial_x^2)^{1/2} C(\lambda) (I - \partial_x^2)^{-1/2} = (\lambda I - \tilde{A}_a^\infty)^{-1} (I - \partial_x^2)^{1/2} (\tilde{A}_a - \tilde{A}_a^\infty) (I - \partial_x^2)^{-1/2}$$

is compact on $L^2 \rightarrow L^2$. $(\lambda I - \tilde{A}_a^\infty)^{-1}$ is still bounded, and by commuting operators, it may be proven that $(I - \partial_x^2)^{1/2} (\tilde{A}_a - \tilde{A}_a^\infty) (I - \partial_x^2)^{-1/2}$ is compact. \square

PROPOSITION 3.3. *Let $\delta \in (0, 1)$ be fixed.*

- (a) *Let $c_\star > n$ and*

$$(3.11) \quad \hat{\vartheta} \in \left(0, \sup_{c \in (n, c_\star]} a_\star(c) / \gamma(c) \right).$$

Then there exists $M = M(c_*, \hat{\vartheta}) > 0$ such that for $c \in (n, c_*]$ and $a \in (\gamma(c)\hat{\vartheta}, a_*(c)]$, if

$$\Re\lambda \geq -\frac{1}{2}ac\gamma^2$$

and at least one of

$$|\Im\lambda| > cM\gamma^3,$$

$$\Re\lambda > cM\gamma^3$$

holds, then $\|C(\lambda)\|_{L^2 \rightarrow L^2} < 1 - \delta$.

(b) This result also holds for $\|C(\lambda)\|_{H^1 \rightarrow H^1}$.

Proof. If $M \geq 1$, then by the assumptions on $\Re\lambda$ and $\Im\lambda$, λ is not in $\sigma_{\text{ess}}(A_a)$; hence we may apply Proposition 3.2 (a) to conclude that $C(\lambda)$ is a compact operator. If the norm of $C(\lambda)$ is less than one, part (c) of that proposition will imply it is not an eigenvalue; we seek an $M \geq 1$ for which $\|C(\lambda)\|$ can be made sufficiently small.

Let $\gamma_0 = \gamma(c_*)$. For all (γ, a) in $\{(\gamma, a) \mid \gamma \in [0, \gamma_0], a \leq a_*(\gamma)\}$, there exist K_1 and K_2 such that

$$\begin{aligned} \|C(\lambda)\|_{L^2 \rightarrow L^2} &\leq \|(\lambda I - A_a^\infty)^{-1} H_{1,a}^{-1} D_a\|_{L^2 \rightarrow L^2} (cK_1\gamma^2) \\ &\quad + \|(\lambda I - A_a^\infty)^{-1} H_{1,a}^{-1}\|_{L^2 \rightarrow L^2} (cK_2\gamma^3). \end{aligned}$$

This comes from expanding the difference $\tilde{A}_a - \tilde{A}_a^\infty$ and commuting operators; see (B.17). $(\lambda I - A_a^\infty)^{-1} H_{1,a}^{-1} D_a$ and $(\lambda I - A_a^\infty)^{-1} H_{1,a}^{-1}$ are constant coefficient operators and we will treat them in Fourier space.

Thus, if we can prove that, for λ satisfying the hypotheses,

$$\begin{aligned} \sup_{\ell \in \mathbb{R}} \left| \frac{(\ell - a) cK_1\gamma^2}{\lambda \left[1 - (\ell - a)^2 \right] - (\ell - a) \left[-c(\ell - a)^2 + c\gamma^2 \right]} \right| &< \frac{1 - \delta}{2}, \\ \sup_{\ell \in \mathbb{R}} \left| \frac{cK_2\gamma^3}{\lambda \left[1 - (\ell - a)^2 \right] - (\ell - a) \left[-c(\ell - a)^2 + c\gamma^2 \right]} \right| &< \frac{1 - \delta}{2}, \end{aligned}$$

we will be done.

Introducing the scalings $\lambda = c\Lambda\gamma^3$, $\ell = \gamma\xi$, and $a = \gamma\vartheta$, this is equivalent to identifying $M \geq 1$ such that when $\Re\Lambda \geq -\frac{1}{2}\vartheta$ and $\Re\Lambda > M$ or $|\Im\Lambda| > M$ then both

$$(3.12) \quad \sup_{\xi \in \mathbb{R}} \left| \frac{(\iota\xi - \vartheta) K_1}{\Lambda \left[1 - \gamma^2 (\iota\xi - \vartheta)^2 \right] - (\iota\xi - \vartheta) \left[-(\iota\xi - \vartheta)^2 + 1 \right]} \right| < \frac{1 - \delta}{2},$$

$$(3.13) \quad \sup_{\xi \in \mathbb{R}} \left| \frac{K_2}{\Lambda \left[1 - \gamma^2 (\iota\xi - \vartheta)^2 \right] - (\iota\xi - \vartheta) \left[-(\iota\xi - \vartheta)^2 + 1 \right]} \right| < \frac{1 - \delta}{2}$$

are satisfied. By our assumption on a , $\vartheta \in (\hat{\vartheta}, 1)$.

Squaring both sides, (3.12) and (3.13) may be rewritten as two polynomial inequalities, $P_1(\Im\Lambda, \Re\Lambda, \xi) > 0$ and $P_2(\Im\Lambda, \Re\Lambda, \xi) > 0$, respectively. We will show that

for appropriately chosen Λ , the inequalities hold for all ξ . P_1 and P_2 are treated similarly. We study P_2 :

$$\begin{aligned}
 P_2(\Im\Lambda, \Re\Lambda, \xi) &= \alpha(\Im\Lambda)^2 + \beta\Im\Lambda + \eta_1(\Re\Lambda)^2 + \eta_2\Re\Lambda + \eta_3, \\
 \alpha &= \gamma^2\xi^4 + 2\xi^2\gamma^2(1 + \gamma^2\vartheta^4) + (1 - \gamma^2\vartheta^2)^2, \\
 \beta &= -2\xi^5\gamma^2 - 2\xi^3(1 + \gamma^2(1 + 2\vartheta^2)) - 2\xi(1 + (-3 + \gamma^2)\vartheta^2 + \gamma^2\vartheta^4), \\
 \eta_1 &= \alpha, \\
 \eta_2 &= 2\vartheta\xi^4\gamma^2 + 2\vartheta\xi^2(3 + \gamma^2(-1 + 2\vartheta^2)) + 2\vartheta(1 - \vartheta^2)(1 - \gamma^2\vartheta^2), \\
 \eta_3 &= \xi^6 + \vartheta^2(1 - \vartheta^2)^2 + \xi^4(2 + 3\vartheta^2) + \xi^2(1 + 3\vartheta^4) - \varpi^2, \\
 \varpi &= 2K_2/(1 - \delta).
 \end{aligned}$$

Using analysis similar to that for P_1 in [21], we first consider P_2 as quadratic in $\Im\Lambda$.¹ Examining its discriminant,

$$\begin{aligned}
 \text{discriminant} &= \gamma^4 \left[-4\xi^8 (\vartheta + \gamma^2\Re\Lambda)^2 + O(\xi^6) \right] \\
 &\quad + \gamma^2 \left[-24\vartheta^2\xi^6 + O(\xi^4) \right] \\
 &\quad - 36\xi^4\vartheta^3 + O(\xi^2).
 \end{aligned}$$

If $0 \geq \Re\Lambda \geq -\vartheta/2$, then there exists $M_0 > 0$ such that for all $\Im\Lambda$, $\gamma \in [0, \gamma_0]$, $\vartheta \in (\hat{\vartheta}, 1)$, and $|\xi| > M_0$ the discriminant is negative and $P_2 > 0$. Furthermore, since the coefficient α is always positive, there exists $R_1 > 0$ such that $P_2 > 0$ when $\gamma \in [0, \gamma_0]$, $|\Im\Lambda| > R_1$, $0 \geq \Re\Lambda \geq -\vartheta/2$, and $|\xi| \leq M_0$.

For $\gamma \in [0, \gamma_0]$ and $\vartheta \in (\hat{\vartheta}, 1)$, both η_1 and η_2 are positive; if $P_2(\Im\Lambda, \Re\Lambda, \xi) > 0$, then $P_2(\Im\Lambda, \Re\Lambda + K, \xi) > 0$ for any $K > 0$. Therefore, $P_2 > 0$ for all ξ if $\gamma \in [0, \gamma_0]$, $|\Im\Lambda| > R_1$, and $\Re\Lambda \geq -\vartheta/2$. Also, $P_2 > 0$ for $|\xi| > M_0$ and all $\Im\Lambda$ if $\Re\Lambda \geq -\vartheta/2$.

We must still treat the case of $|\Im\Lambda| \leq R_1$ and $|\xi| \leq M_0$ simultaneously. Consider P_2 as quadratic in $\Re\Lambda$. η_1 and η_2 are positive for all ξ , and η_3 is bounded from below. Therefore there is some $R_2 > 0$ such that $P_2 > 0$ for all ξ if $\gamma \in [0, \gamma_0]$, $|\Im\Lambda| \leq R_1$, and $\Re\Lambda > R_2$. Thus, for any $\vartheta \in (\hat{\vartheta}, 1)$, $\gamma_0 \in [0, \gamma_0]$, there exist $R_1 > 0$ and $R_2 > 0$ such that $P_2 > 0$ for all ξ if

$$\begin{aligned}
 &\gamma \in [0, \gamma_0], \\
 &\Re\Lambda \geq -\vartheta/2, \\
 &|\Im\Lambda| \leq R_1 \quad \text{or} \quad \Re\Lambda > R_2.
 \end{aligned}$$

For $C(\lambda) : H^1 \rightarrow H^1$, the proof is similar, with constants \tilde{K}_1 and \tilde{K}_2 in place of K_1 and K_2 . \square

¹Though this proof largely follows that for the analogous proposition in [21], we have made two modifications. The first is to introduce $\hat{\vartheta} > 0$ in (3.11). Though this restriction may not be needed, by placing a lower bound on ϑ , estimates are more obviously uniform in γ and ϑ . A second change is to take $\Re\Lambda \geq -\vartheta/2$ instead of $\Re\Lambda \geq -\vartheta/(4\gamma^2)$. Our condition is less sharp but, again, makes the intermediary bounds clearly uniform in γ and ϑ .

3.3. Small eigenvalues: The Evans function. In this section we rule out eigenvalues of A_a in the set $|\lambda| \leq M\gamma^3$. This is done using the *Evans function*, an analytic function that vanishes at eigenvalues of A_a . The Evans function, $D = D(\lambda; \gamma)$, is constructed for the eigenvalue problem using particular solutions of an associated dynamical system,

$$(3.14) \quad \dot{\mathbf{y}} = B(x, \lambda, \gamma)\mathbf{y},$$

$$(3.15) \quad \mathbf{y} = O(e^{\mu_1 x}) \quad \text{as } x \rightarrow +\infty,$$

and the adjoint system,

$$(3.16) \quad \dot{\mathbf{z}} = -\mathbf{z}B(x, \lambda, \gamma),$$

$$(3.17) \quad \mathbf{z} = O(e^{\mu_1 x}) \quad \text{as } x \rightarrow -\infty.$$

μ_1 will be the eigenvalue of smallest real part of B^∞ , the limit as $x \rightarrow \pm\infty$ of B . When certain conditions, described in Theorem 3.4, are met, the Evans function exists and may be explicitly defined as

$$(3.18) \quad D(\lambda; \gamma) = \mathbf{z}(x; \lambda, \gamma) \cdot \mathbf{y}(x; \lambda, \gamma).$$

The idea is to measure the angle between the subspace of solutions decaying at $+\infty$ with the subspace decaying at $-\infty$; hence the appearance of the dot product. The Evans function has an equivalent formulation in terms of the determinant of the fundamental solution of (3.14). For a more complete discussion of the Evans function, see [27].

THEOREM 3.4 (see [27, 29]). *Let Ω be a simply connected subset of \mathbb{C}^2 . Suppose that the system (3.14) satisfies the following hypotheses:*

- (i) $B : \mathbb{R} \times \Omega \rightarrow \mathbb{C}^{n \times n}$ is continuous in (x, λ, γ) and analytic in (λ, γ) for fixed x .
- (ii) $B^\infty(\lambda, \gamma) = \lim_{x \rightarrow \pm\infty} B(x, \lambda, \gamma)$ exists for all $(\lambda, \gamma) \in \Omega$. The limit is attained uniformly on compact subsets of Ω .
- (iii) The integral

$$\int_{-\infty}^{\infty} \|B(x, \lambda, \gamma) - B^\infty(\lambda, \gamma)\| dx$$

converges for all $(\lambda, \gamma) \in \Omega$ and the convergence is uniform on compact subsets of Ω .

- (iv) For every $(\lambda, \gamma) \in \Omega$, the matrix $B^\infty(\lambda, \gamma)$ has a unique eigenvalue of smallest real part, which is simple, denoted μ_1 .

Then $D(\lambda; \gamma)$ is well defined and analytic on Ω , such that $D(\lambda; \gamma) = 0$ if and only if 3.14 has a solution $\mathbf{y}(x)$ satisfying (3.15) and

$$(3.19) \quad \mathbf{y}(x) = o(e^{\mu_1 x}) \quad \text{as } x \rightarrow -\infty.$$

3.3.1. The KdV Evans function. In the case of the KdV equation, the eigenvalue problem may be scaled to

$$(3.20) \quad \partial_x L_{\text{KdV}} Y = \partial_x \left(-\partial_x^2 Y + Y - 3\text{sech}^2 \left(\frac{1}{2}x \right) Y \right) = \Lambda Y.$$

Because the speed parameter has been scaled out, there is only one eigenvalue parameter, Λ .

Making the identification

$$(3.21) \quad \mathbf{y} = (Y, \quad \partial_x Y, \quad L_{\text{KdV}}Y)^T,$$

we see that \mathbf{y} satisfies the dynamical system

$$(3.22) \quad \dot{\mathbf{y}} = B_{\text{KdV}}(x, \Lambda)\mathbf{y},$$

$$(3.23) \quad B_{\text{KdV}}(x, \Lambda) = \begin{pmatrix} 0 & 1 & 0 \\ 1 - 3\text{sech}(\frac{1}{2}x)^2 & 0 & -1 \\ \Lambda & 0 & 0 \end{pmatrix}.$$

A complete description of the associated Evans may be found in [28]. We summarize as follows.

THEOREM 3.5 (the KdV Evans function).

(a) *The Evans function $D_{\text{KdV}}(\Lambda)$ associated with (3.20) is given by*

$$D_{\text{KdV}}(\Lambda) = \left(\frac{\mu_1(\Lambda) + 1}{\mu_1(\Lambda) - 1} \right)^2,$$

where $\mu_1(\Lambda)$ denotes the root of $\mu^3 - \mu + \Lambda = 0$ of minimal real part.

(b) *The domain of $D_{\text{KdV}}(\Lambda)$ is the slit complex plane*

$$\Delta_{\text{KdV}} = \mathbb{C} \setminus \left(-\infty, -\sqrt{\frac{4}{27}} \right].$$

(c) *The essential spectrum of $A_{\text{KdV}} : L_a^2 \rightarrow L_a^2$ is a curve contained entirely in the domain $\{\lambda : \Re\lambda < -\epsilon\}$ for some $\epsilon > 0$. Furthermore, if $\Delta_{\text{KdV}}^+(a)$ denotes the component of $\mathbb{C} \setminus \sigma_{\text{ess}}(A_{\text{KdV}})$ that contains the right half-plane, then $D_{\text{KdV}}(\Lambda)$ has no zeros in $\Delta_{\text{KdV}}^+(a)$ except for a zero of multiplicity two at $\Lambda = 0$.*

3.3.2. The Evans function applied. The eigenvalue problem $AY = \lambda Y$,

$$\{I - \partial_x[\phi_c^n \partial_x(\phi_c^{-m} \cdot)]\}^{-1} \partial_x \{-c\phi_c^n \partial_x^2(\phi_c^{-m} \cdot) + [c - n\phi_c^{-1} + cn(\phi_c^{-1} - 1)]\} Y = \lambda Y,$$

is equivalent to

$$(3.24) \quad \boxed{\partial_x L_c Y - \lambda [I - \partial_x(\phi_c^n \partial_x(\phi_c^{-m} \cdot))] Y = 0.}$$

Defining

$$(3.25) \quad \mathbf{y} = (\phi_c^{-m} Y, \quad \partial_x(\phi_c^{-m} Y), \quad L_c Y + \lambda \phi_c^n \partial_x(\phi_c^{-m} Y))^T,$$

\mathbf{y} solves the dynamical system

$$(3.26) \quad \dot{\mathbf{y}} = B(x, \lambda, c)\mathbf{y},$$

$$(3.27) \quad B(x, \lambda, c) = \begin{pmatrix} 0 & 1 & 0 \\ c^{-1} \phi_c^{m-n} [c - n\phi_c^{-1} + cn(\phi_c^{-1} - 1)] & \lambda/c & -c^{-1} \phi_c^{-n} \\ \lambda \phi_c^m & 0 & 0 \end{pmatrix}.$$

The matrix B may be decomposed as $B = B^\infty(\lambda, c) + R(x, \lambda, c)$.

(3.28)

$$B^\infty(\lambda, c) = \begin{pmatrix} 0 & 1 & 0 \\ \gamma^2 & \lambda/c & -c^{-1} \\ \lambda & 0 & 0 \end{pmatrix},$$

(3.29)

$$R(x, \lambda, \gamma) = \begin{pmatrix} 0 & 0 & 0 \\ c^{-1}\phi_c^{m-n} [c - n\phi_c^{-1} + cn(\phi_c^{-1} - 1)] - \gamma^2 & 0 & -c^{-1}(\phi_c^{-n} - 1) \\ \lambda(\phi_c^m - 1) & 0 & 0 \end{pmatrix}.$$

Also note that for the corresponding adjoint eigenvalue problem

$$(3.30) \quad \boxed{-L_c^* \partial_x W - \lambda [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] W = 0}$$

under the identifications

$$(3.31) \quad \mathbf{z} = (-c\partial_x(\phi_c^n \partial_x W) - \lambda\phi_c^n \partial_x W, \quad c\phi_c^n \partial_x W, \quad W)$$

\mathbf{z} solves

$$(3.32) \quad \dot{\mathbf{z}} = -\mathbf{z}B(x, \lambda, \gamma).$$

THEOREM 3.6 (properties of the Evans function).

(a) *The Evans function is defined and analytic on the set $\Omega \subset \mathbb{C}^2$,*

$$(3.33) \quad \boxed{\Omega = \{(\lambda, \gamma) \mid \gamma \in (0, 1) \text{ and } \lambda \in \Omega_\gamma\}}$$

with

$$(3.34) \quad \boxed{\Omega_\gamma = \{\lambda \mid \Re \lambda > -\lambda_0\} \setminus (-\lambda_0, -\lambda_{\text{cut}}(\gamma))},$$

where

$$(3.35) \quad \begin{aligned} \lambda_{\text{cut}} &= \sqrt{\frac{1}{8} \sqrt{8c^2 + 20cn - n^2 - 8c\sqrt{n^2 + 8cn} - n\sqrt{n^2 + 8cn}}} \\ &= \frac{2}{3\sqrt{3}}n\gamma^3 + \frac{8}{9\sqrt{3}}n\gamma^5 + O(\gamma^7) \end{aligned}$$

and $\lambda_0 = n\sqrt{27/16}$.

(b) *Given $(\lambda, \gamma) \in \Omega$ and $a \leq a_*(\gamma)$, if λ is to the right of $\sigma_{\text{ess}}(A_a)$, then the following are equivalent:*

- $D(\lambda; \gamma) = 0$.
- λ is an L^2 eigenvalue of A_a .

(c) *For such zeros of D , the algebraic multiplicity of λ as an eigenvalue of A_a is equal to the order of λ as a zero of $D(\lambda; \gamma)$.*

(d) *$D(0; \gamma) = \partial_\lambda D(0; \gamma) = 0$, and hence it is an eigenvalue of algebraic multiplicity at least two.*

Remark 3.7. With this construction, for $\Re\lambda < 0$, $\lambda \in \Omega_\gamma$, the characteristic polynomial, (3.36), not only has a unique root of minimal real part, but all roots have distinct real part as well. This is stronger than is needed.

Remark 3.8. λ_{cut} is labeled as such because there is a branch cut in the Evans function there.

Remark 3.9. The eigenvalue at the origin is related to the solitary waves being a two-parameter (speed and phase) family of solutions of (1.2). The presence of the corresponding bound states, explicitly given in Proposition 3.23, could lead to a weak instability of algebraic growth in t . However, we will project out these modes by allowing our leading order solitary wave to modulate its speed and phase.

Proof of Theorem 3.6. Part (a) requires the verification of the hypotheses of Theorem 3.4 for this system. Applying the properties of the ϕ_c and Corollaries 2.9 and 2.11, and through examination of (3.27), B is clearly continuous in its three arguments for $\lambda \in \mathbb{C}$ and $c \geq n$. In addition, for fixed x , it will be analytic in (λ, c) , or, equivalently, (λ, γ) . Thus property (i) of Theorem 3.4 holds.

By Corollary 2.7, the limiting matrix B^∞ exists and $B - B^\infty$ is in L^1 . This will be uniform on compact subsets of $\mathbb{C} \times [0, 1)$, establishing properties (ii) and (iii) of Theorem 3.4.

Lastly, we must verify the existence of μ_1 , the unique eigenvalue of minimal real part. We divide this into two parts, $\Re\lambda \geq 0$ and $\Re\lambda < 0$. The characteristic polynomial, $\mathcal{P}(\mu)$, of B^∞ is

$$(3.36) \quad \boxed{c\mathcal{P}(\mu) = (\lambda - c\mu)(1 - \mu^2) + n\mu.}$$

Following the analysis in section 2(c) of [27] for a similar polynomial in the case of generalized BBM, one confirms that property (iv) holds for $\Re\lambda \geq 0$ and all γ , and hence the Evans function exists in $\{\Re\lambda \geq 0\} \times [0, 1)$.

Using the analysis in [21] for Theorem 2.7 of the polynomial, one can conclude the existence of some $\lambda_1 > 0$ and identify $\tilde{\Omega}(\gamma)$, such that for all $\gamma \in (0, 1)$, a unique root of minimal real part exists for λ in the set

$$\{\lambda : \Re\lambda < -\lambda_1\} \setminus (-\lambda_1, -\tilde{\Omega}(\gamma)].$$

Alternatively, we give a more precise analysis of (3.36) in Appendix C that yields values of λ_0 and $\tilde{\Omega}$ given in the proposition. This concludes the proof of part (a).

To prove part (b), we need a lemma regarding the location of $\sigma_{\text{ess}}(A_a)$.

LEMMA 3.10. *Let $\gamma \in (0, 1)$ and $a \leq a_*(\gamma)$. Let $\Omega_+ = \Omega_+(\gamma, a)$ denote the component of $\mathbb{C} \setminus \sigma_{\text{ess}}(A_a)$ containing the origin. Then for $\lambda \in \Omega_+ \cap \Omega_\gamma$, the roots of the characteristic polynomial satisfy the relation*

$$(3.37) \quad \Re\mu_1 < -a < \Re\mu_{j \neq 1}.$$

Proof. By inspection, if $\lambda \in \sigma_{\text{ess}}(A_a)$, there is a root μ_j of (3.36) with $\Re\mu_j = -a$. Conversely, if there is a root with real part $-a$, then λ is in the essential spectrum. Hence the characteristic polynomial has a root with real part $-a$ if and only if $\lambda \in \sigma_{\text{ess}}(A_a)$

As noted in [27, section 2 (c)], for large $|\lambda|$, the roots of the characteristic polynomial (3.36) are

$$-1 + O(|\lambda|^{-1}), \quad 1 + O(|\lambda|^{-1}), \quad \lambda/c + O(|\lambda|^{-1}).$$

So for large λ in the right half-plane, (3.37) holds because $a < 1$. Now suppose for some $\lambda \in \Omega_+ \cap \Omega_\gamma$ the inequality were false. Because the $\Re\mu_j$ depend continuously on λ , equality would have to hold for some λ , but then it must be that $\lambda \in \sigma_{\text{ess}}(A_a)$, which we have assumed is not the case. \square

We now prove part (b) of Theorem 3.6. If $D(\lambda; \gamma) = 0$, then there is a solution to the ODE, $\dot{\mathbf{y}} = B\mathbf{y}$, such that

$$\mathbf{y}(x) = O(e^{\mu_1 x}) \quad \text{as } x \rightarrow +\infty \quad \text{and} \quad \mathbf{y}(x) = o(e^{\mu_1 x}) \quad \text{as } x \rightarrow -\infty.$$

Hence, by Part 1(d) of Proposition 1.6 of [27], for sufficiently small ϵ ,

$$\mathbf{y}(x) = O(e^{\mu_* x + \epsilon|x|}) \quad \text{as } x \rightarrow -\infty.$$

Letting $W(x) = e^{ax} \phi_c(x)^m y_1(x)$,

$$W(x) = O(e^{(\mu_* + a)x + \epsilon|x|}) \quad \text{as } x \rightarrow -\infty \quad \text{and} \quad W(x) = O(e^{(\mu_1 + a)x}) \quad \text{as } x \rightarrow +\infty.$$

From (3.37), $\mu_1 + a < 0 < \mu_* + a$, so W will decay exponentially fast at $\pm\infty$. Hence it is an L^2 solution to the eigenvalue problem $A_a W = \lambda W$.

Conversely, if we have an L^2 eigenfunction, then it must satisfy

$$W(x) = O(e^{(\mu_1 + a)x}) \quad \text{as } x \rightarrow +\infty \quad \text{and} \quad W(x) = o(e^{(\mu_1 + a)x}) \quad \text{as } x \rightarrow -\infty.$$

$Y(x) = e^{-ax} W(x)$ will then satisfy (3.24) in a classical sense, although it may not be in L^2 . However it will satisfy the necessary decay estimates on \mathbf{y} , constructed from Y as in (3.25), such that $D(\lambda; \gamma) = 0$. This concludes the proof of part (b).

The proof of part (c) follows that of Lemma 2.9 from [28]. First, it is proved that if, for a given γ , λ is a zero of order k of $D(\lambda; \gamma)$, then λ is an L^2 eigenvalue of A_a of algebraic multiplicity at least k . It is then shown that it cannot have algebraic multiplicity greater than k . We omit repeating these details. Part (d) is then a consequence of (c) and the calculations in Appendix D that $D(0; \gamma) = \partial_\lambda D(0; \gamma) = 0$ for all γ . \square

Remark 3.11. For $\Re\lambda \leq 0$ with $D(\lambda; \gamma) = 0$, it is likely, but not proved, that λ is not an L^2 eigenvalue of A .

3.3.3. The Evans function in the KdV scaling. We now introduce $D_*(\Lambda; \gamma)$, the Evans function for (3.24) under the KdV scalings introduced in section 2.2:

$$(3.38) \quad \xi = \gamma x, \quad \lambda = c\Lambda\gamma^3, \quad \phi_c(x) = 1 + \frac{\gamma^2}{n-1} U(\xi(x); \gamma).$$

The eigenvalue problem is now

$$(3.39) \quad \begin{aligned} \partial_\xi L_\gamma Y &= \partial_\xi \left[- \left(1 + \frac{\gamma^2}{n-1} U \right)^n \partial_\xi^2 \left(\left(1 + \frac{\gamma^2}{n-1} U \right)^{-m} \cdot \right) \right] Y \\ &\quad + \partial_\xi \left[\left(1 + \frac{\gamma^2}{n-1} U \right)^{-1} (1-U) \right] Y \\ &= \Lambda \left[I - \gamma^2 \partial_\xi \left(\left(1 + \frac{\gamma^2}{n-1} U \right)^n \partial_\xi \left(\left(1 + \frac{\gamma^2}{n-1} U \right)^{-m} \cdot \right) \right) \right] Y. \end{aligned}$$

Recall that $U = U(\xi; \gamma)$ is the solution of (2.18), and for $\gamma = 0$, $U = U_*$, the KdV soliton.

We can construct a dynamical system formulation of (3.39), defining the vector \mathbf{Y} as

$$(3.40) \quad \mathbf{Y} = \begin{pmatrix} \left(1 + \frac{\gamma^2}{n-1}U\right)^{-m} Y \\ \partial_\xi \left[\left(1 + \frac{\gamma^2}{n-1}U\right)^{-m} Y \right] \\ L_\gamma Y + \gamma^2 \Lambda \left(1 + \frac{\gamma^2}{n-1}U\right)^n \partial_\xi \left[\left(1 + \frac{\gamma^2}{n-1}U\right)^{-m} Y \right] \end{pmatrix},$$

which satisfies

$$(3.41) \quad \dot{\mathbf{Y}} = B_\star(\xi, \Lambda, \gamma)\mathbf{Y},$$

$$(3.42) \quad B_\star(\xi, \Lambda, \gamma) = \begin{pmatrix} 0 & 1 & 0 \\ \left(1 + \frac{\gamma^2}{n-1}U(\xi; \gamma)\right)^{m-n-1} (1 - U(\xi; \gamma)) & \gamma^2 \Lambda & -\left(1 + \frac{\gamma^2}{n-1}U(\xi; \gamma)\right)^{-n} \\ \Lambda \left(1 + \frac{\gamma^2}{n-1}U(\xi; \gamma)\right)^m & 0 & 0 \end{pmatrix}.$$

As $\xi \rightarrow \infty$, the matrix is

$$(3.43) \quad B_\star^\infty(\Lambda, \gamma) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & \gamma^2 \Lambda & -1 \\ \Lambda & 0 & 0 \end{pmatrix},$$

which has the characteristic polynomial

$$(3.44) \quad P_\star(\nu; \Lambda, \gamma) = \nu^3 - \gamma^2 \Lambda \nu^2 - \nu + \Lambda.$$

A few remarks about the scaled problem. The assumptions stated in Theorem 3.4 remain the same, except now the matrix under inspection is B_\star , with eigenvalue parameters (Λ, γ) . \mathbf{Y} and \mathbf{y} are related:

$$\mathbf{y}(x) = \begin{pmatrix} Y_1(\xi(x)) \\ \gamma Y_2(\xi(x)) \\ c\gamma^2 Y_3(\xi(x)) \end{pmatrix}.$$

At $\gamma = 0$, (3.39) is

$$\partial_\xi L_0 Y = \partial_\xi [-\partial_\xi^2 Y + (1 - U_\star(\xi; 0)Y)] = \Lambda Y,$$

the KdV eigenvalue problem, (3.20), and

$$B_\star(\xi, \Lambda, 0) = \begin{pmatrix} 0 & 1 & 0 \\ 1 - U_\star(\xi) & 0 & -1 \\ \Lambda & 0 & 0 \end{pmatrix}$$

is the matrix for the KdV dynamical system.

PROPOSITION 3.12 (scaled Evans function).

(a) $D_*(\Lambda; \gamma)$ is defined and analytic on the set $\Delta \subset \mathbb{C}^2$,

$$(3.45) \quad \Delta = \{(\Lambda, \gamma) \mid \gamma \in [0, 1) \text{ and } \lambda \in \Delta_\gamma\},$$

where, for $\gamma > 0$,

$$\begin{aligned} \Delta_\gamma &= \left\{ \Lambda \mid \Re \Lambda > -\frac{\lambda_0}{c\gamma^3} \right\} \setminus \left(-\frac{\lambda_0}{c\gamma^3}, -\frac{\lambda_{\text{cut}}(\gamma)}{c\gamma^3} \right] \\ &= \left\{ \Lambda \mid \Re \Lambda > -\gamma^{-3} \sqrt{\frac{27}{16}} + O(\gamma^{-1}) \right\} \\ &\quad \setminus \left(-\gamma^{-3} \sqrt{\frac{27}{16}} + O(\gamma^{-1}), -\sqrt{\frac{4}{27}} + O(\gamma^2) \right]. \end{aligned}$$

λ_0 and $\tilde{\Omega}$ are as defined by Theorem 3.6. When $\gamma = 0$,

$$\Delta_0 = \Delta_{\text{KdV}} = \mathbb{C} \setminus \left(-\infty, -\sqrt{4/27} \right].$$

(b) For fixed $\gamma \in (0, 1)$ and $\Lambda \in \Delta_\gamma$,

$$D_*(\Lambda; \gamma) = D(c\Lambda\gamma^3; \gamma).$$

(c) For $\Lambda \in \Delta_0$,

$$D_*(\Lambda; 0) = D_{\text{KdV}}(\Lambda).$$

Proof. The proof of these statements follows that of Proposition 2.8 in [21] and Theorems 4.9–4.11 of [29].

Proof of part (a). For $\gamma > 0$, as in Theorem 3.6, we must identify a set in \mathbb{C}^2 in which the hypotheses of Theorem 3.4 are valid. Parts (i)–(iii) are obvious as the solitary wave $U(\xi; \gamma)$ decays exponentially in ξ and, for fixed ξ , will be analytic in γ . We are left to verify part (iv). The characteristic polynomial of B_\star^∞ is (3.44). As noted in [21], the roots of P_\star are related to those of P , (3.36), by

$$\mu(\lambda, \gamma) = \gamma\nu(\Lambda, \gamma).$$

So P_\star will have a unique root of minimal real part for a given Λ and γ when P has such a unique root for $\lambda = c\Lambda\gamma^3$. Therefore, for $\gamma \in (0, 1)$, if Λ is in the set

$$\Delta_\gamma = \frac{1}{c\gamma^3} \Omega_\gamma,$$

(iv) will be satisfied. If $\gamma = 0$, (3.41) and (3.42) coincide with the KdV system, for which $\Delta_0 = \mathbb{C} \setminus (-\infty, -\sqrt{4/27}]$. Clearly, as $\gamma \rightarrow 0$, Δ_γ limits to Δ_0 .

Proof of part (b). From part (a), $\lambda = c\Lambda\gamma^3 \in \Omega_\gamma$ and, by construction,

$$y_1(x, \lambda, \gamma) \sim e^{\mu_1 x} \quad \text{as } x \rightarrow +\infty,$$

$$Y_1(\xi, \Lambda, \gamma) \sim e^{\nu_1 \xi} \quad \text{as } \xi \rightarrow +\infty.$$

At $\xi = \gamma x$, $\mu_1 = \gamma\nu_1$, $\lambda = c\Lambda\gamma^3$, $y_1(x, \lambda, \gamma) = Y_1(\xi, \Lambda, \gamma)$. Using the *transmission coefficient* interpretation of the Evans function, we then have

$$\begin{aligned} \mathbf{y}(x, \lambda, \gamma) &\sim D(\lambda; \gamma)e^{\mu_1 x} \begin{pmatrix} 1 \\ \mu_1 \\ \lambda/\mu_1 \end{pmatrix}^T \quad \text{as } x \rightarrow -\infty, \\ \mathbf{Y}(\xi, \Lambda, \gamma) &\sim D_\star(\Lambda; \gamma)e^{\nu_1 \xi} \begin{pmatrix} 1 \\ \nu_1 \\ \Lambda/\nu_1 \end{pmatrix}^T \quad \text{as } \xi \rightarrow -\infty, \end{aligned}$$

implying

$$D(c\Lambda\gamma^3; \gamma) = D_\star(\Lambda; \gamma).$$

Proof of part (c). Trivially, when $\gamma = 0$, this is the KdV Evans function problem exactly. \square

PROPOSITION 3.13. *Let $\varepsilon_2 \in (0, \sqrt{4/27})$ and $M > 0$. Set*

$$\mathcal{O} = \{\Re\Lambda \geq -\varepsilon_2, \quad |\Lambda| < M\}.$$

Then for all γ sufficiently small, $\mathcal{O} \subset \Delta_\gamma$ and

$$\lim_{\gamma \rightarrow 0} D_\star(\Lambda; \gamma) = D_\star(\Lambda; 0) = D_{\text{KdV}}(\Lambda)$$

with uniform convergence for $\Lambda \in \mathcal{O}$.

Proof. Clearly, for γ sufficiently close to zero, $\mathcal{O} \subset \Delta_\gamma$. Since \mathcal{O} is compact and D_\star is analytic in both arguments,

$$\sup_{\Lambda \in \mathcal{O}} |D_\star(\Lambda; \gamma) - D_\star(\Lambda; 0)|$$

may be made arbitrary small by taking γ sufficiently close to zero. \square

THEOREM 3.14. *There exist $\gamma_\star \in (0, 1)$ and $\hat{\vartheta} > 0$ such that for all $\gamma \in (0, \gamma_\star]$ and $a \in (\gamma\hat{\vartheta}, a_\star(\gamma))$, there exists $\varepsilon = \varepsilon(\gamma, a) > 0$, such that*

- *the only eigenvalue of A_a with $\Re\lambda \geq -\varepsilon$ is $\lambda = 0$, with algebraic multiplicity two;*
- *the only zero of $D(\lambda; \gamma)$ with $\Re\lambda \geq -\varepsilon$ is $\lambda = 0$, a root of order two.*

Proof. Adapting the proof of Theorem 2.1 in [21] we break a half-plane into two parts: a bounded set about the origin and its unbounded complement. The operator estimates from Proposition 3.3 will rule out eigenvalues in the unbounded part, and the Evans function in the KdV scaling will control eigenvalues in the bounded part.

Let us apply Proposition 3.3 (a) with $\delta = 1/4$, $c_\star = 4n$, and

$$\hat{\vartheta} = \frac{1}{4} \sup_{c \in (n, c_\star]} \frac{a_\star(c)}{\gamma(c)}.$$

Then there exists $M > 0$, such that for any $c \in (n, 4n]$, $a \in (\gamma\hat{\vartheta}, a_\star(\gamma))$, and $\lambda \in \Omega_U$,

$$\Omega_U = \left\{ \lambda : \Re\lambda \geq -\frac{1}{2}ac\gamma^2, \quad |\lambda| > cM\gamma^3 \right\},$$

$\|C(\lambda)\|_{L^2 \rightarrow L^2} < 1$. By Proposition 3.2 (c), such λ cannot be L^2 eigenvalues of A_a .

Using this M , let $\varepsilon_2 = \frac{1}{4}$ and set

$$\mathcal{O} = \{\Lambda : \Re\Lambda \geq -\varepsilon_2, \quad |\Lambda| \leq M\}.$$

Define

$$m = \min\{|D_{\text{KdV}}(\Lambda)| \mid \Lambda \in \partial\mathcal{O}\}.$$

By Proposition 3.13, there exists a $\gamma_* \leq 1/2$ such that for all $\gamma \in [0, \gamma_*]$,

$$|D_*(\Lambda; \gamma) - D_*(\Lambda; 0)| = |D_*(\Lambda; \gamma) - D_{\text{KdV}}(\Lambda)| < m \quad \text{for all } \Lambda \in \partial\mathcal{O}.$$

By Theorem 3.5, the only root of D_{KdV} in \mathcal{O} is at $\Lambda = 0$, with multiplicity two. Applying Rouché’s theorem (see, for example, [16, 17]), for all $\gamma \in [0, \gamma_*]$, $D_*(\cdot, \gamma)$ and D_{KdV} have the same number of roots in \mathcal{O} . By Proposition 3.12 (b), if $\gamma \in (0, \gamma_*]$ and $\Lambda \in \mathcal{O}$, $D_*(\Lambda; \gamma) = D(c\Lambda\gamma^3; \gamma)$; therefore on the set Ω_B ,

$$\Omega_B = c\gamma^3\mathcal{O} = \{\lambda : \Re\lambda \geq -\varepsilon_2c\gamma^3, \quad |\lambda| \leq cM\gamma^3\}.$$

$D(\lambda; \gamma)$ also has only two zeros. Theorem 3.6 (d) asserts that these two roots are at the origin, corresponding to an eigenvalue of algebraic multiplicity two. By Theorem 3.6 (b), A_a then has no *nonzero* L^2 eigenvalues in Ω_B .

Combining Ω_U and Ω_B , if we set

$$(3.46) \quad \varepsilon(\gamma, a) = \min \left\{ \frac{1}{2}ac\gamma^2, \quad \varepsilon_2c\gamma^3, \quad \frac{1}{2}\omega, \quad \frac{1}{2}\lambda_{\text{cut}}(\gamma), \quad \frac{1}{2}\lambda_0 \right\},$$

then by Theorem 3.6 (b), $\lambda = 0$ is both the only eigenvalue of A_a and the only root of $D(\lambda; \gamma)$ with $\Re\lambda \geq -\varepsilon$. \square

COROLLARY 3.15. For $\gamma \in (0, \gamma_*]$, $\partial_c\mathcal{N}[\phi_c] \neq 0$.

Proof. This is a consequence the preceding theorem and the relation $\partial_\lambda^2 D(0; \gamma) \propto \partial_c\mathcal{N}[\phi_c]$; see Appendix D. \square

3.4. The Evans function beyond the KdV scaling. For $\gamma > \gamma_*$, one may compute the Evans function numerically to assert that a given λ is not an eigenvalue. Moreover, the winding number of the image of the Evans function evaluated on the line $\Re\lambda = -\lambda_1 < 0$ equals the number of zeros in the set $\Re\lambda > -\lambda_1$. Therefore, one might evaluate the Evans function on such a line and exam the plot, as we do in Figures 4 and 5.

These plots indicate that $\lambda = 0$ is the only zero in the closed right half-plane for the values of c , n , and m under consideration. Up to the acceptance of these numerics, this extends Theorem 3.14. Note that we do not evaluate out to $-\lambda_1 + i\infty$, but merely compute at sufficiently large values of λ such that we are near the asymptotic value of the Evans function. It may be proven that there exists $D_\infty(\gamma)$ such that

$$\lim_{|\lambda| \rightarrow \infty, \lambda \in \Omega} D(\lambda; \gamma) = D_\infty(\gamma).$$

A further numerical computation will reveal that this value is nonzero.

In the Hamiltonian case, $n + m = 0$, an analytical result is possible for $\gamma > \gamma_*$. The linearized operator, A , may be written as $A = J_c L_c$,

$$(3.47) \quad J_c = [I - \partial_x (\phi_c^n \partial_x (\phi_c^n))]^{-1} \partial_x, \quad J_c^* = -J_c,$$

$$(3.48) \quad L_c = -c\phi_c^n \partial_x^2 (\phi_c^n) + (c - n\phi_c^{-1} + cn(\phi_c^{-1} - 1)), \quad L_c^* = L_c.$$

This structure permits an extension of Theorem 3.14 beyond the KdV regime, given below in Theorem 3.21. However, this is absent for general n and m .

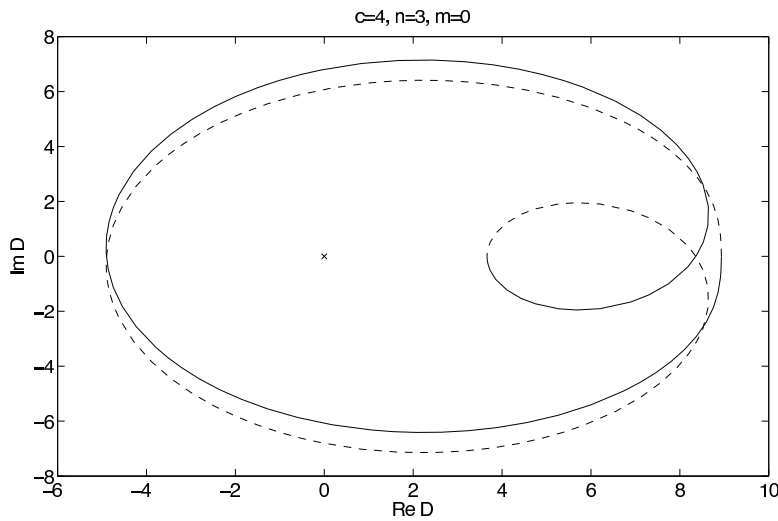


FIG. 4. $D(\cdot; c = 4)$ evaluated on a portion of the strip $\Re \lambda = -1/5$. Since $D(\bar{\lambda}) = \overline{D(\lambda)}$ we only compute $\Im \lambda > 0$, and then reflect; this is the dashed curve. The curve wraps around the origin, marked by x , twice. In this case, $n = 3$ and $m = 0$.

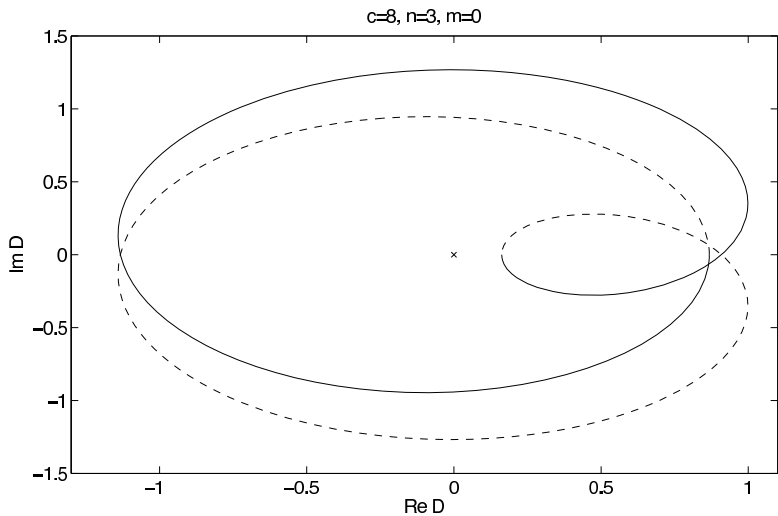


FIG. 5. $D(\cdot; c = 8)$ evaluated on a portion of the strip $\Re \lambda = -1$. Since $D(\bar{\lambda}) = \overline{D(\lambda)}$ we only compute $\Im \lambda > 0$, and then reflect; this is the dashed curve. The curve wraps around the origin, marked by x , twice. In this case, $n = 3$ and $m = 0$.

Remark 3.16. This section is the only place where the analyticity of the Evans function in the γ argument is used. In turn, this is the only place requiring analyticity of ϕ_c in c from Corollary 2.11 (b). For the results in the preceding section, *joint continuity* of $D(\lambda; \gamma)$ in its two arguments is sufficient.

LEMMA 3.17. *In the Hamiltonian case, the following are equivalent for $\Re \lambda \geq 0$, $a \leq a_*(\gamma)$:*

- λ is an L^2 eigenvalue of A .
- $D(\lambda; \gamma) = 0$.

Proof. For $\Re\lambda > 0$, we know that $\mu_\star = \min\{\Re\mu_2, \Re\mu_3\} > 0$, so an eigenfunction, having submaximal growth at $-\infty$, must decay exponentially fast. This is very similar to the relation in the case of the weighted operator from Theorem 3.6 (b).

For $\Re\lambda = 0$, the proof relies on the JL structure of the operator A . See the proof of [27, Theorem 3.6]. \square

LEMMA 3.18. *If λ is a nonzero purely imaginary eigenvalue, then $\partial_\lambda D(\lambda; \gamma) \neq 0$.*

Proof. The proof is by contradiction. Let Y^+ be the corresponding eigenfunction, $AY^+ = \lambda Y^+$. It may be proven that the subspace $\mathcal{Y} = \text{span}\{Y^+, \overline{Y^+}\}$ satisfies

$$\langle L_c u, v \rangle = 0 \quad \text{for all } u, v \in \mathcal{Y}.$$

As $\mathcal{Y} \cap \ker(L_c) = \{0\}$, we may apply Lemma 3.3 of [27] to conclude $\dim \mathcal{Y} \leq 1$, a contradiction. See Lemma 3.3 from [28] for details. \square

LEMMA 3.19 (monotonicity of functional analytically confirmed for $n = 2$). *In the Hamiltonian case, assuming $\partial_c \mathcal{N}[\phi_c] \neq 0$ for all c , then for all c , there are no eigenvalues with $\Re\lambda > 0$.*

Proof. By Theorem 3.14, the result holds for $c \leq c_\star$. We argue by contradiction to extend it beyond c_\star . Assume for some $\gamma_0 > \gamma_\star$ that there exists λ_0 , $\Re\lambda_0 > 0$, such that $D(\lambda_0; \gamma_0) = 0$. If $\Im\lambda_0 \neq 0$, then $D(\overline{\lambda_0}; \gamma_0) = 0$ and there would be two eigenvalues in the right half-plane. But by Theorem 3.1 of [27], $A = J_c L_c$ has no more eigenvalues (counting multiplicity) with $\Re\lambda > 0$ than L_c does with $\Re\lambda < 0$. As is discussed in [40], L_c has exactly one negative eigenvalue; this is a contradiction, so λ_0 is real.

Because the number of zeros in the right half-plane is at most one, counting multiplicity, we know $\partial_\lambda D(\lambda_0; \gamma_0) \neq 0$. Applying the implicit function theorem, we get an analytic function, $\lambda(\gamma)$, defined in a neighborhood of γ_0 , such that $\lambda(\gamma_0) = \lambda_0$ and $D(\lambda(\gamma); \gamma) = 0$.

Let \mathcal{O} be the maximal domain of analyticity of $\lambda(\gamma)$. In a sufficiently small neighborhood of γ_0 , $\Re\lambda(\gamma) > 0$. For real-valued γ in this neighborhood, we must have, by the above argument about complex-conjugates, that $\Im\lambda(\gamma) = 0$. Considering the power series expansion of $\lambda(\gamma)$, about γ_0 , $\lambda(\gamma)$ will be real-valued for real-valued $\gamma \in \mathcal{O}$.

Let

$$\gamma_1 = \inf \{ \gamma \in (\gamma_\star/2, \gamma_0) \cap \mathcal{O} \mid D(\lambda(\gamma); \gamma) = 0 \}.$$

For all $\gamma \in [\gamma_1, \gamma_0]$, we must have $\lambda(\gamma) > 0$. Suppose this is not the case. Then, by continuity, for some γ , we must have $\lambda(\gamma) = 0$. But this would imply that $\lambda = 0$ was a root of multiplicity three, contradicting the assumption on \mathcal{N} , which ensures it is a root of multiplicity two.

$\lambda(\gamma)$ may be analytically continued down till at least $\gamma_\star/2$. If not, then $\gamma_1 > \gamma_\star/2$ and $\partial_\lambda D(\lambda(\gamma_1); \gamma_1) \neq 0$ since this root must be simple. Therefore we could apply the implicit function theorem again and extend $\lambda(\gamma)$ below γ_1 , contradicting its minimality.

But then $D(\lambda(\gamma_\star); \gamma_\star) = 0$, and $\lambda(\gamma_\star) > 0$, contradicting Theorem 3.14. \square

Remark 3.20. An analogous result may be found in Theorem 3.4 [27] for generalized KdV and generalized BBM. However, the argument there is very different because it may be proven that $D(\lambda) \rightarrow 1$ as $|\lambda| \rightarrow \infty$. This does not hold for the Evans function associated with (1.2), due to the appearance of a nonlinearity in the dispersive term.

THEOREM 3.21 (monotonicity of functional analytically confirmed for $n = 2$). *In the Hamiltonian case, assuming $\partial_c \mathcal{N}[\phi_c] \neq 0$ for all c , Theorem 3.14 may be extended to all $\gamma \in (\gamma_*, 1)$, except for a discrete set whose only possible accumulation point is $\gamma = 1$.*

Proof. By Lemma 3.19, if Theorem 3.4 were false for some $\gamma > \gamma_*$, it would be due to a zero appearing on the imaginary axis. We will prove by contradiction that the set

$$(3.49) \quad E = \{\gamma \in [0, 1) \mid \text{there exists } \beta > 0 \text{ such that } D_\star(i\beta; \gamma) = 0\}$$

has no accumulation points. We consider only positive β 's because if $i\beta$ is a root, then so is $-i\beta$. This follows the proofs of Theorem 3.6 of [28] and Theorem 2.1 of [21].

Assuming E has a limit point, there exists a sequence, $\gamma_j \in E$, $\gamma_j \rightarrow \gamma_0 \in E$ as $j \rightarrow \infty$. Taking a subsequence if necessary, γ_j and γ_0 are bounded away from $\gamma = 1$. We will now rule out large eigenvalues, and then argue by contradiction to rule out small eigenvalues.

Applying Proposition 3.3 to this range of γ values, there will exist $M > 0$, such that the corresponding $\beta_j > 0$ must satisfy $\beta_j \leq M$. Taking a subsequence if necessary, $\beta_j \rightarrow \beta_0$, $\beta_0 \leq M$, $D_\star(i\beta_j; \gamma_j) = 0$, and, by continuity, $D_\star(i\beta_0; \gamma_0) = 0$.

Note that $\beta_0 \neq 0$. $\Lambda = 0$ is always a root of order at least two. The assumption $\partial_c \mathcal{N}[\phi_c] \propto \partial_\lambda^2 D(0; \gamma) = (c^2 \gamma^6)^{-1} \partial_\lambda^2 D_\star(0; \gamma) \neq 0$ forces it to be a root of order exactly two. But $\beta_0 = 0$ would imply that it was a zero of at least four, a contradiction.

Applying Lemma 3.18, $\partial_\Lambda D_\star(i\beta_j; \gamma_j) \neq 0$ and $\partial_\Lambda D_\star(i\beta_0; \gamma_0) \neq 0$. By the implicit function theorem, there is an analytic function $\Lambda_0(\gamma)$ defined in a neighborhood of γ_0 , such that $\Lambda_0(\gamma_0) = i\beta_0$ and $D_\star(\Lambda_0(\gamma); \gamma) = 0$. By considering the power series expansion of $\Lambda_0(\gamma)$ about γ_0 , we see, by taking γ_j sufficiently close to γ_0 , that $\Lambda_0(\gamma)$ is purely imaginary for real γ in its maximal domain of analyticity, \mathcal{O} .

Let

$$\gamma_1 = \inf \{\gamma \in [0, \gamma_0) \cap \mathcal{O} \mid D_\star(\Lambda_0(\gamma); \gamma) = 0\}.$$

For all $\gamma \in [\gamma_1, \gamma_0]$, we must have $\Im \Lambda_0(\gamma) > 0$. If not, then by continuity there would exist $\gamma \in (\gamma_1, \gamma_0)$, for which $\Lambda_0(\gamma) = 0$, yielding a contradiction as before.

Suppose $\gamma_1 > 0$. $\Im \Lambda_0(\gamma_1) > 0$ because, if not, then by continuity there would exist $\gamma \in (\gamma_1, \gamma_0)$, for which $\Lambda_0(\gamma) = 0$, leading to a contradiction again. Therefore, we may be sure that $\partial_\Lambda D_\star(\Lambda_0(\gamma_1); \gamma_1) \neq 0$. We may then apply the implicit function theorem, allowing us to continue Λ_0 below γ_1 , another contradiction. Therefore $\gamma_1 = 0$. But then $\Im \Lambda_0(0) > 0$ and $0 = D_\star(\Lambda_0(0); 0) = D_{\text{KdV}}(\Lambda_0(0))$, a contradiction. \square

Remark 3.22. This result is limited by our inability to analytically evaluate the functional $\mathcal{N}[\phi_c]$. The authors were similarly stymied in [40], where the orbital stability of the solitary waves relies on proving $\partial_c \mathcal{N}[\phi_c] > 0$. Here, as there, one may numerically evaluate the functional and observe its monotonicity in the speed argument. See [40] for the case $n = 2$.

This result, along with (3.8) and Theorem 3.14, completes the proof of Theorem 3.1.

3.5. The generalized kernel.

PROPOSITION 3.23. *Let $c > n$, $a \leq a_*(c)$ and assume $\partial_c \mathcal{N}[\phi_c] \neq 0$. Define*

$$\tilde{\xi}_1 = \partial_x \phi_c,$$

$$\tilde{\xi}_2 = \partial_c \phi_c,$$

$$\tilde{\eta}_1 = \Theta [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] \int_{-\infty}^x (L_c^*)^{-1} [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] \int_{-\infty}^x \frac{\partial_x \phi_c}{\phi_c^{n+m}} dx,$$

$$\tilde{\eta}_2 = \Theta [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] \int_{-\infty}^x \frac{\partial_x \phi_c}{\phi_c^{n+m}} dx,$$

$$\Theta = (\partial_c \mathcal{N}[\phi_c])^{-1}.$$

For $j = 1, 2$, set

$$\xi_j = e^{ay} \tilde{\xi}_j \quad \text{and} \quad \eta_j = e^{-ay} \tilde{\eta}_j.$$

Then $\{\xi_1, \xi_2\}$ and $\{\eta_1, \eta_2\}$ are biorthogonal bases for $\ker_g(A_a)$ and $\ker_g(A_a^*)$, $\langle \xi_i, \eta_j \rangle = \delta_{ij}$. They satisfy the relations

$$\begin{aligned} A_a \xi_1 &= 0, & A_a^* \eta_2 &= 0, \\ A_a \xi_2 &= -\xi_1, & A_a^* \eta_1 &= -\eta_2. \end{aligned}$$

Proof. It is easy to verify $A \tilde{\xi}_1 = 0$, $A \tilde{\xi}_2 = -\xi_1$, and $A^* \tilde{\eta}_2 = 0$.

$$\ker(L_c^*) = \text{span} \{ \phi_c^{-n-m} \partial_x \phi_c \}.$$

The kernel is orthogonal to $\tilde{\eta}_2$ because $\tilde{\eta}_2$ is an even function, while the kernel is odd. Therefore, $\tilde{\eta}_1$ is well defined and $A^* \tilde{\eta}_1 = -\tilde{\eta}_2$. \square

4. Semigroup decay. The following result proves the convective stability of solitary waves under the linearized flow. As we will do in section 6, this may be employed to prove full nonlinear stability.

THEOREM 4.1 (linearized stability). *Assume $\gamma \in (0, 1)$, $a \leq a_*(\gamma)$, and there exists $\varepsilon > 0$ such that $\lambda = 0$ is the only eigenvalue of A_a with $\Re \lambda \geq -\varepsilon$. Then the initial-value problem*

$$w_t = A_a w,$$

$$w(0) = w_0 \in H^1 \cap \ker_g(A_a^*)^\perp$$

has a unique solution $w(t) = e^{A_a t} w_0 \in C_0([0, \infty); H^1)$ with

$$(4.1) \quad \|w(t)\|_{H^1} \leq C e^{-bt} \|w_0\|_{H^1}$$

for some $C > 0$ and $b > 0$.

Remark 4.2. There exists a $b_{\max} > 0$ such that (4.1) will hold for all $b \in (0, b_{\max})$.

In particular, for $\gamma \in (0, \gamma_*)$ and $a \in (\gamma \hat{\nu}, a_*(\gamma)]$, $b_{\max} \geq \varepsilon$; see Theorem 3.1.

Proof of Theorem 4.1. This is based on a result due to Prüss [30].

THEOREM 4.3. *Let B be the infinitesimal generator of a C_0 -semigroup on a Hilbert space Z . Let $b > 0$. If there exists $M > 0$ such that*

$$\|(\lambda I - B)^{-1}\|_{Z \rightarrow Z} \leq M \quad \text{for all } \Re \lambda > -b,$$

then $\|e^{Bt}\|_{Z \rightarrow Z} \leq e^{-bt}$.

Following the approach in [21] for Proposition 3.1, we first show that A_a is the infinitesimal generator of a C_0 -semigroup on H^1 . Examining the Fourier symbol of A_a^∞ , A_a^∞ is such a semigroup. As (B.16) shows, $A_a - A_a^\infty$ is a bounded operator, so we may apply Theorem 3.1.1 of [25], which states that bounded perturbations of infinitesimal generators are also infinitesimal generators.

Consider the Hilbert space

$$Z = H^1 \cap \ker_g(A_a^*)^\perp$$

equipped with the H^1 norm and the operator

$$B = A_a|_Z$$

the restriction of A_a to Z . B inherits from A_a that it is the infinitesimal generator of a C_0 -semigroup on Z . Also note that $\sigma(B) = \sigma(A_a) \setminus \{0\}$ by Theorem III-6.17 of [14].

Recall from Theorem 3.1 that $\sigma_{\text{ess}}(A_a)$ is contained in left half-plane, and all points in $\sigma(A_a) \setminus \sigma_{\text{ess}}(A_a)$ are eigenvalues of finite multiplicity. By the assumption on the spectrum of A_a , the spectrum of B is contained in the open left half-plane.

We will now prove there exists a uniform bound on resolvent of B for $\Re\lambda > -b$ for some $b > 0$. For $\lambda \in \rho(B)$,

$$(\lambda I - B)^{-1} = (\lambda I - A_a)^{-1}|_Z,$$

so

$$\|(\lambda I - B)^{-1}\|_{Z \rightarrow Z} \leq \|(\lambda I - A_a)^{-1}\|_{H^1 \rightarrow H^1}.$$

The resolvent of A_a may be written as

$$(\lambda I - A_a)^{-1} = (I - C(\lambda))^{-1} (\lambda I - A_a^\infty)^{-1},$$

with $C(\lambda)$ as defined in Proposition 3.2. Since $\Re\sigma_{\text{ess}}(A_a) \leq -\omega < 0$, we must have $b < \omega$. Then, for $\Re\lambda \geq -b$, the Fourier symbol of $(\lambda I - A_a^\infty)^{-1}$ is uniformly bounded, so

$$\|(\lambda I - A_a^\infty)^{-1}\|_{H^1 \rightarrow H^1} \leq M' \quad \text{for some } M' > 0.$$

By Proposition 3.3, for $\Re\lambda \geq -b \geq -\frac{1}{2}ac\gamma^2$, there exists $R' > 0$ such that

$$\|(I - C(\lambda))^{-1}\|_{H^1 \rightarrow H^1} \leq 2 \quad \text{for } |\lambda| > R'.$$

Therefore,

$$\|(\lambda I - B)^{-1}\|_{Z \rightarrow Z} \leq \|(\lambda I - A_a)^{-1}\|_{H^1 \rightarrow H^1} \leq 2M' \quad \text{for } \Re\lambda \geq -b \text{ and } |\lambda| > R'.$$

For $|\lambda| \leq R'$ and $\Re\lambda \geq -\varepsilon$, B has no eigenvalues. B is a closed operator, and therefore $(\lambda I - B)^{-1}$ is holomorphic on this compact set (see Theorem III-6.7 of [14]) giving the bound

$$\|(\lambda I - B)^{-1}\|_{Z \rightarrow Z} \leq M'' \quad \text{for some } M'' > 0.$$

Hence for all $\Re\lambda \geq -\min\{b, \varepsilon\}$,

$$\|(\lambda I - B)^{-1}\|_{Z \rightarrow Z} \leq \max\{2M', M''\},$$

and we may apply Theorem 4.3. □

5. Prelude to nonlinear stability. We need three results about our system before we can prove Theorem 1.1. First, we establish a criterion for when a decomposition of ϕ into a modulating solitary wave and a perturbation is possible. Then we derive equations for the evolution of the parameters associated with this modulating solitary wave. Finally, we relate the H^1 norm to the H_a^1 norm of the perturbation.

5.1. Local existence of decomposition and continuation principles. In analyzing the weighted perturbation w , we wish to treat the nonlinear terms perturbatively, with the leading order behavior governed by the linear operator A_a . As noted in Proposition 3.23, the operator has a two-dimensional kernel. To prevent the appearance of secular terms, the perturbation must be orthogonal to $\ker_g(A_a^*)$; this reveals how the decomposition of ϕ into a perturbation and a modulated solitary wave, (2.28), occurs. This follows the strategy of [26] and [21].

PROPOSITION 5.1. *Let $c_0 > n$, $a \leq a_*(c_0)$, and $t_1 > 0$. Given $\delta_1 > 0$, there exists $\delta_0 > 0$ such that for any $\phi - 1 \in C^1([0, t_1]; H^1 \cap H_a^1)$ satisfying*

$$(5.1) \quad \sup_{t \leq t_1} \|e^{a(\cdot + \theta_0)}(\phi(\cdot, t) - \phi_{c_0}(\cdot - c_0 t + \theta_0))\|_{H^1} \leq \delta_0 \quad \text{for some } \theta_0 \in \mathbb{R}$$

there exists a unique mapping $t \mapsto (\theta(t), c(t))$ in $C^1([0, t_1]; \mathbb{R}^2)$ such that

$$(5.2) \quad \sup|\theta(t) - \theta_0| + \sup|\dot{\theta}(t)| + \sup|c(t) - c_0| + \sup|\dot{c}(t)| \leq \delta_1, \quad t \leq t_1,$$

$$(5.3) \quad \mathcal{T}_k[\phi - 1, \theta, c] = \langle \phi(x, t) - \phi_{c(t)}(y), \tilde{\eta}_j(y) \rangle = 0 \quad \text{for } j = 1, 2 \text{ and } t \in [0, t_1],$$

where $y = x - \int_0^t c(s)ds + \theta(t)$.

The number δ_0 may be chosen as a decreasing function of t_1 .

Proof. The proof is via the implicit function theorem. In this context, the Banach spaces are $C^1([0, t_1]; H_a^1)$ and $C^1([0, t_1]; \mathbb{R}^2)$, the latter space equipped with the norm

$$\sup_{t \leq t_1} |\theta(t)| + \sup_{t \leq t_1} |\dot{\theta}(t)| + \sup_{t \leq t_1} |c(t)| + \sup_{t \leq t_1} |\dot{c}(t)|.$$

The functional is

$$\mathcal{T} = (\mathcal{T}_1, \mathcal{T}_2)^T : C^1([0, t_1]; H_a^1) \times C^1([0, t_1]; \mathbb{R}^2) \rightarrow C^1([0, t_1]; \mathbb{R}^2),$$

and it is C^1 in its arguments, permitting the use of the implicit function theorem.

Setting

$$\mathbf{U}_0 = (\phi_{c_0}(x - c_0 t) - 1, \quad 0, \quad c_0)$$

we see $\mathcal{T}[\mathbf{U}_0] = 0$. The Fréchet derivative at \mathbf{U}_0 is

$$(5.4) \quad \begin{aligned} D\mathcal{T}[\mathbf{U}_0](\delta\phi, \delta\theta, \delta c) &= \left(\begin{aligned} &\langle e^{a(x - c_0 t)} \delta\phi(x, t), \eta_1(x - c_0 t) \rangle - \int_0^1 \delta c(s) ds + \delta\theta(t) \\ &\langle e^{a(x - c_0 t)} \delta\phi(x, t), \eta_2(x - c_0 t) \rangle + \delta c(t) \end{aligned} \right). \end{aligned}$$

The derivative acting on the (θ, c) is $\begin{pmatrix} I & -B \\ 0 & I \end{pmatrix}$, $(Bf)(t) = \int_0^t f(s)ds$. This block operator has a bounded inverse on $C^1([0, t_1]; \mathbb{R}^2) \rightarrow C^1([0, t_1]; \mathbb{R}^2)$. By the implicit function theorem, there exists $\delta_0 > 0$, such that if

$$\sup_{t \leq t_1} \|\phi(\cdot - \theta_0, t) - \phi_{c_0}(\cdot - c_0 t)\|_{H_a^1} < \delta_0 \quad \text{for some } \theta_0 \in \mathbb{R},$$

then there exists a mapping in $C^1([0, t_1]; \mathbb{R}^2)$, $t \mapsto (\tilde{\theta}(t), c(t))$, satisfying

$$\begin{aligned} \sup|\tilde{\theta}(t)| + \sup|\dot{\tilde{\theta}}(t)| + \sup|c(t) - c_0| + \sup|\dot{c}(t)| &< \delta_1 \quad \text{for } t \leq t_1, \\ \mathcal{T}[\phi(\cdot - \theta_0, t) - 1, \tilde{\theta}(t), c(t)] &= 0 \quad \text{for } t \leq t_1. \end{aligned}$$

Defining $\theta(t) = \tilde{\theta}(t) + \theta_0$, and applying the change of variables $x \mapsto \tilde{x} + \theta_0$, we obtain the form given in (5.1), (5.2), and (5.3).

Finally, because our norms are taken as the supremum over $t \leq t_1$, if δ_0 works for t_1 , then it will also work for any $t_2 \leq t_1$. This allows δ_0 to be constructed as a decreasing function of t_1 . \square

PROPOSITION 5.2. *Let $c_0 > n$, $a \leq a_*(c_0)$ and assume $\phi - 1 \in C^1([0, T_{\max}]; H^1 \cap H_a^1)$, $T_{\max} > 0$ solves (1.2).*

Given $\delta_1 > 0$, there exist $\delta_0 > 0$ and $\delta'_0 > 0$, such that for any $t_0 \in [0, T_{\max})$, if the decomposition of ϕ , $v(y, t) = \phi(x, t) - \phi_{c(t)}(y)$, $y = x - \int_0^t c(s)ds + \theta(t)$, and $(\theta, c) \in C^1([0, t_0]; \mathbb{R}^2)$ satisfies

$$(5.5) \quad \sup_{t \leq t_0} \|v(t)\|_{H_a^1} \leq \delta_0/3,$$

$$(5.6) \quad \sup_{t \leq t_0} |c(t) - c_0| \leq \delta'_0,$$

$$(5.7) \quad \mathcal{T}[\phi - 1, \theta, c](t) = 0 \quad \text{for } t \in [0, t_0],$$

then there is a unique extension of (θ, c) in $C^1([0, t_0 + t_]; \mathbb{R}^2)$ for some $t_* > 0$ such that*

$$(5.8) \quad \mathcal{T}[\phi - 1, \theta, c](t) = 0 \quad \text{for } t \leq t_0 + t_* \leq T_{\max},$$

$$(5.9) \quad \begin{aligned} \sup|\theta(t) - \theta(t_0)| + \sup|\dot{\theta}(t)| \\ + \sup|c(t) - c(t_0)| + \sup|\dot{c}(t)| \leq \delta_1, \quad t \in [t_0, t_0 + t_*]. \end{aligned}$$

Proof. This follows the proof of Proposition 5.2 in [28], although we are forced to modify it as we do not know a priori that a solution exists for all time.

Given δ_1 , let $\delta_0 > 0$ be the value from Proposition 5.1 with $t_1 = \frac{1}{2}(T_{\max} - t_0)$. Set $\tilde{\phi}(x, t) = \phi(x, t + t_0)$, $\theta_1 = -\int_0^{t_0} c(s)ds + \theta(t_0)$. Let δ'_0 be sufficiently small such that

$$(5.10) \quad \|\phi_{c'} - \phi_{c_0}\|_{H^1 \cap H_a^1} \leq \delta_0/3 \quad \text{for all } c' \text{ such that } |c' - c_0| \leq \delta'_0.$$

Then, since

$$\begin{aligned} \left\| e^{a(\cdot + \theta_1)} \left(\tilde{\phi}(0) - \phi_{c_0}(\cdot + \theta_1) \right) \right\|_{H^1} &\leq \left\| e^{a(\cdot + \theta_1)} \left(\tilde{\phi}(0) - \phi_{c_1}(\cdot + \theta_1) \right) \right\|_{H^1} \\ &\quad + \left\| e^{a(\cdot + \theta_1)} \left(\phi_{c_1}(\cdot + \theta_1) - \phi_{c_0}(\cdot + \theta_1) \right) \right\|_{H^1} \\ &= \|v(t_0)\|_{H_a^1} + \|\phi_{c_1} - \phi_{c_0}\|_{H_a^1} \leq \frac{2}{3}\delta_0, \end{aligned}$$

we have

$$\begin{aligned} \left\| e^{a(\cdot + \theta_1)} \left(\tilde{\phi}(t) - \phi_{c_0}(\cdot - c_0 t + \theta_1) \right) \right\|_{H^1} &\leq \left\| e^{a(\cdot + \theta_1)} \left(\tilde{\phi}(t) - \tilde{\phi}(0) \right) \right\|_{H^1} + \frac{2}{3}\delta_0 \\ &= e^{a\theta_1} \|\phi(t + t_0) - \phi(t_0)\|_{H^1 \cap H_a^1} + \frac{2}{3}\delta_0. \end{aligned}$$

As ϕ is continuous in time, there exists a $t_\star \leq t_1$, such that

$$\sup_{t \leq t_\star} \|e^{a(\cdot + \theta_1)} (\tilde{\phi}(t) - \phi_{c_0}(\cdot - c_0 t + \theta_1))\|_{H^1} \leq \delta_0.$$

Therefore, $\tilde{\phi}$ satisfies the hypotheses of Proposition 5.1. We have a unique $(\tilde{\theta}(t), \tilde{c}(t))$ with $(\tilde{\theta}(0), \tilde{c}(0)) = (\theta(t_0), c(t_0))$. This gives us the extension $(\theta(t), c(t)) = (\tilde{\theta}(t - t_0) - \theta_0 + \theta(t_0), \tilde{c}(t - t_0))$ for $t \in [t_0, t_0 + t_\star]$ and (5.9) will hold. \square

5.2. Modulation equations. Given that the perturbation must be orthogonal to $\ker_g(A_a^\star)$, the associated constraints give a pair of ODEs, coupled to the perturbation, giving a complete system of three equations for the three dependent variables.

Let P denote the projection onto $\ker_g(A_a^\star)$. Assuming this space is two-dimensional, we use the biorthogonal bases given in Proposition 3.23 to define projection onto this space and its complement,

$$(5.11) \quad P = \langle \eta_1, \cdot \rangle \xi_1 + \langle \eta_2, \cdot \rangle \xi_2,$$

$$(5.12) \quad Q = I - P.$$

The secular terms are excised from the perturbation equation, (2.36), by requiring that

$$(5.13) \quad w_\tau = A_a w + Q\mathcal{G},$$

$$(5.14) \quad P\mathcal{G} = 0,$$

$$(5.15) \quad Pw(\tau = 0) = 0.$$

Constraint (5.14) corresponds to

$$(5.16) \quad \langle \eta_j, \mathcal{G} \rangle = 0 \quad \text{for } j = 1, 2.$$

These two equations govern $c(t)$ and $\theta(t)$, completing our system.

Defining $p_1(y, t) = \partial_y \phi_{c(t)}(y) - \partial_y \phi_{c_0}(y)$ and $p_2(y, t) = \partial_c \phi_{c(t)}(y) - \partial_c \phi_{c_0}(y)$, the derived system is

$$(5.17)$$

$$\begin{pmatrix} 1 + \langle \tilde{\eta}_1, p_1 \rangle & \langle \tilde{\eta}_1, p_2 \rangle \\ \langle \tilde{\eta}_2, p_1 \rangle & 1 + \langle \tilde{\eta}_2, p_2 \rangle \end{pmatrix} \begin{pmatrix} \dot{\theta} \\ \dot{c} \end{pmatrix} = \frac{c - \dot{\theta}}{c_0} \begin{pmatrix} \langle \eta_1, S_a[c_0, c, \dot{\theta}]w \rangle \\ \langle \eta_2, S_a[c_0, c, \dot{\theta}]w \rangle \end{pmatrix} + \begin{pmatrix} \langle \eta_1, \mathcal{G}_1 \rangle \\ \langle \eta_2, \mathcal{G}_1 \rangle \end{pmatrix}.$$

However, the right-hand side still has $\dot{\theta}$ dependence. Observe that

$$(5.18)$$

$$\begin{aligned} \frac{c - \dot{\theta}}{c_0} S_a[c_0, c, \dot{\theta}]w &= cm \left[\partial_y \log \left(\frac{\phi_{c_0}}{\phi_c} \right) \right] w - \dot{\theta} m \partial_y \log(\phi_{c_0}) w \\ &+ \left[\frac{c}{c_0} n \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} D_a [(\phi_{c_0}^{-1} - c_0(\phi_{c_0}^{-1} - 1)) \cdot] \right. \\ &\quad \left. - n \phi_c^m H_{\phi_c, a}^{-1} D_a [(\phi_c^{-1} - c(\phi_c^{-1} - 1)) \cdot] \right] w \\ &- \frac{\dot{\theta}}{c_0} \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} D_a [(\phi_{c_0}^{-1} - c_0(\phi_{c_0}^{-1} - 1)) \cdot] w \\ &- cm \left[\phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} [\partial_y (\log \phi_{c_0}) \cdot] - \phi_c^m \mathcal{H}_{\phi_c, a}^{-1} [\partial_y (\log \phi_c) \cdot] \right] w \\ &+ \dot{\theta} m \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} [\partial_y (\log \phi_{c_0}) \cdot] w. \end{aligned}$$

Defining

$$(5.19) \quad \tilde{S}_a = \tilde{S}_a^1 + \tilde{S}_a^2 + \tilde{S}_a^3,$$

$$(5.20) \quad \tilde{S}_a^1 = cm (\phi_{c_0}^{-1} \partial_y \phi_{c_0} - \phi_c^{-1} \partial_y \phi_c),$$

$$(5.21) \quad \begin{aligned} \tilde{S}_a^2 = n \left\{ \frac{c}{c_0} \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} D_a [(\phi_{c_0}^{-1} - c_0(\phi_{c_0}^{-1} - 1)) \cdot] \right. \\ \left. - \phi_c^m H_{\phi_c, a}^{-1} D_a [(\phi_c^{-1} - c(\phi_c^{-1} - 1)) \cdot] \right\}, \end{aligned}$$

$$(5.22) \quad \tilde{S}_a^3 = -cm \left[\phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} (\phi_{c_0}^{-1} \partial_y \phi_{c_0} \cdot) - \phi_c^m H_{\phi_c, a}^{-1} (\phi_c^{-1} \partial_y \phi_c \cdot) \right],$$

and

$$(5.23) \quad \begin{aligned} T_a = -m \phi_{c_0}^{-1} \partial_y \phi_{c_0} - c_0^{-1} \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} D_a [(\phi_{c_0}^{-1} - c_0(\phi_{c_0}^{-1} - 1)) \cdot] \\ + m \phi_{c_0}^m H_{\phi_{c_0}, a}^{-1} (\phi_{c_0}^{-1} \partial_y \phi_{c_0} \cdot), \end{aligned}$$

the right-hand side of (5.17) may be written as

$$\begin{pmatrix} \langle \eta_1, \tilde{S}_a w \rangle \\ \langle \eta_2, \tilde{S}_a w \rangle \end{pmatrix} + \dot{\theta} \begin{pmatrix} \langle \eta_1, T_a w \rangle \\ \langle \eta_2, T_a w \rangle \end{pmatrix} + \begin{pmatrix} \langle \eta_1, \mathcal{G}_1 \rangle \\ \langle \eta_2, \mathcal{G}_1 \rangle \end{pmatrix}.$$

Equation(5.17) may be solved algebraically so that $\dot{\theta}$ appears only on the left-hand side,

$$(5.24) \quad \begin{aligned} \mathcal{B}(t) \begin{pmatrix} \dot{\theta} \\ \dot{c} \end{pmatrix} &= \begin{pmatrix} 1 + \langle \tilde{\eta}_1, p_1 \rangle - \langle \eta_1, T_a w \rangle & \langle \tilde{\eta}_1, p_2 \rangle \\ \langle \tilde{\eta}_2, p_1 \rangle - \langle \eta_2, T_a w \rangle & 1 + \langle \tilde{\eta}_2, p_2 \rangle \end{pmatrix} \begin{pmatrix} \dot{\theta} \\ \dot{c} \end{pmatrix} \\ &= \begin{pmatrix} \langle \eta_1, \tilde{S}_a w \rangle \\ \langle \eta_2, \tilde{S}_a w \rangle \end{pmatrix} + \begin{pmatrix} \langle \eta_1, \mathcal{G}_1 \rangle \\ \langle \eta_2, \mathcal{G}_1 \rangle \end{pmatrix}, \end{aligned}$$

where $\mathcal{B}(t) = I + O(|c(t) - c_0|) + O(\|w\|_{L^2})$; for sufficiently small $|c(t) - c_0| + \|w\|_{L^2}$, \mathcal{B} is invertible. Thus we have equations for $\dot{\theta}$ and \dot{c} , closing the system for (v, c, θ) .

Remark 5.3. In (5.24), we see that when $\mathcal{B}(t)$ is inverted, the right-hand side of the system is continuous in t . Therefore, provided $c(t)$, $\theta(t)$, and $w(t)$ are continuous in t , $c(t)$ and $\theta(t)$ will actually be C^1 .

5.3. Lyapunov bound. Using the functional $\mathcal{N}[\phi]$, defined in (2.8), we have the following.

PROPOSITION 5.4. *Let $c_0 > n$, $a \leq a_*(c_0)$, and let $\phi(x, t)$ be a solution to (1.2) in $C^1([0, T]; H^1 \cap H_a^1)$ with data*

$$\phi_0 = \phi_{c_0}(x + \theta_0) + v_0(x), \quad \theta_0 \in \mathbb{R}.$$

Assume the decomposition $\phi(x, t) \rightarrow (v(y, t), \theta(t), c(t))$ exists for $t \leq T$ and

$$|c(t) - c_0| + \|v(\cdot, t)\|_{H^1} \leq \delta_1 < 1 \quad \text{for } t \leq T,$$

$$\partial_c \mathcal{N}[\phi_c] \Big|_{c=c_0} \neq 0.$$

Then there exist constants K and K' such that

$$(5.25) \quad \begin{aligned} \|v\|_{H^1}^2 (1 - K'\|v\|_{H^1}) &\leq K \left(|\Delta\mathcal{N}| + \|w\|_{L^2} + \|w\|_{L^2}^2 \right. \\ &\quad \left. + |c(t) - c_0| + |c(t) - c_0|^2 + |c(t) - c_0|^3 \right), \\ \Delta\mathcal{N} &= \mathcal{N}[\phi_{c(t)} + v] - \mathcal{N}[\phi_{c_0}]. \end{aligned}$$

Remark 5.5. This proposition, which relates the unweighted norm of the perturbation to the weighted norm, is where the invariance of \mathcal{N} is particularly useful. In [29], the Boussinesq equations lacked such an invariant; thus, only linear stability was proven. Strictly speaking, it would be sufficient for $\Delta\mathcal{N}$ to be bounded in terms of the data; the conservation of \mathcal{N} is not essential, though it is helpful.

Proof. Taylor expanding \mathcal{N} about a solitary wave with perturbation z , we obtain

$$\mathcal{N}[\phi_{c_0} + z] = \mathcal{N}[\phi_{c_0}] + \langle \delta\mathcal{N}[\phi_{c_0}], z \rangle + \frac{1}{2} \langle \delta^2\mathcal{N}[\phi_{c_0}]z, z \rangle + O(\|z\|_{H^1}^3).$$

From [38], the first and second variations are

$$(5.26) \quad \begin{aligned} \langle \delta\mathcal{N}[\phi_{c_0}], z \rangle &= \int \left(\frac{1 - \phi_{c_0}^{1-n-m}}{n+m-1} + m\phi_{c_0}^{-2m-1} (\partial_x \phi_{c_0})^2 - \partial_x^2 \phi_{c_0}^{-2m} \phi_{c_0} \right) z dx \\ &= \Theta^{-1} \langle \tilde{\eta}_2, z \rangle, \end{aligned}$$

$$(5.27) \quad \begin{aligned} \langle \delta^2\mathcal{N}[\phi_{c_0}]z, z \rangle &= \int \left(\phi_{c_0}^{-n-m} - m(1+2m)\phi_{c_0}^{-2m-2} (\partial_x \phi_{c_0})^2 \right) z^2 dx \\ &\quad + \int (2m\phi_{c_0}^{-2m-1} \partial_x^2 \phi_{c_0}) z^2 dx + \int \phi_{c_0}^{-2m} (\partial_x z)^2 dx, \end{aligned}$$

where $\tilde{\eta}_2$ and Θ are as defined in Proposition 3.23.

Take $z(y, t) = \phi_{c(t)}(y) - \phi_{c_0}(y) + v(y, t) = \phi(x, t) - \phi_{c_0}(y)$. Then

$$(5.28) \quad \begin{aligned} \langle \delta\mathcal{N}[\phi_{c_0}], z \rangle &= \langle \Theta^{-1} \tilde{\eta}_2, z(y, t) \rangle \\ &= \Theta^{-1} \langle \tilde{\eta}_2, \phi_{c(t)}(y) - \phi_{c_0}(y) \rangle + \Theta^{-1} \langle \tilde{\eta}_2, v(y, t) \rangle. \end{aligned}$$

Using the continuity of $c \mapsto \phi_c - 1$,

$$(5.29) \quad \Theta^{-1} \langle \tilde{\eta}_2, \phi_{c(t)}(y) - \phi_{c_0}(y) \rangle \leq K|c(t) - c_0|.$$

The term with the perturbation, v , may be bounded by

$$(5.30) \quad \Theta^{-1} \langle \tilde{\eta}_2, v(y, t) \rangle = \Theta^{-1} \langle \eta_2, w(y, t) \rangle \leq \Theta^{-1} \|\eta_2\|_{L^2} \|w\|_{L^2} \leq K\|w\|_{L^2}.$$

Now we bound the second variation. For brevity, let

$$\Phi_{c_0} = -m(1+2m)\phi_{c_0}^{-2m-2} (\partial_x \phi_{c_0})^2 + 2m\phi_{c_0}^{-2m-1} \partial_x^2 \phi_{c_0}.$$

Then

$$(5.31) \quad \begin{aligned} \langle \delta^2\mathcal{N}[\phi_{c_0}]z, z \rangle &= \int \phi_{c_0}^{-2m} (\partial_x z)^2 + \phi_{c_0}^{-n-m} z^2 + \int \Phi_{c_0} z^2 \\ &\geq K\|z\|_{H^1}^2 + \langle \Phi_{c_0}, (\phi_{c(t)} - \phi_{c_0})^2 \rangle \\ &\quad + 2 \langle \Phi_{c_0}, (\phi_{c(t)} - \phi_{c_0}) v \rangle + \langle \Phi_{c_0}, v^2 \rangle \\ &\geq K_1\|v\|_{H^1}^2 - K_2|c(t) - c_0|^2 \\ &\quad - K_3|c(t) - c_0| - K_4 \langle \Phi_{c_0} e^{-\alpha y}, w^2 \rangle. \end{aligned}$$

We would like $\Phi_{c_0} e^{-2ay} \in L^\infty$, so that the last term may be estimated by $\|w\|_{L^2}^2$. Since $a \leq a_*(\gamma) < \frac{1}{2}\gamma$, we may do this.

Lastly, we have the remainder term $\mathcal{R}[\phi_{c_0}, z]$. Because of the a priori bound on $\|v\|_{H^1}$, this may be estimated as

$$(5.32) \quad |\mathcal{R}[\phi_{c_0}, z]| \leq K \|z\|_{H^1}^3 \leq K (|c(t) - c_0|^3 + \|v\|_{H^1}^3).$$

Combining these estimates, (5.29), (5.30), (5.31), and (5.32), gives us

$$\begin{aligned} \|v\|_{H^1}^2 (1 - D\|v\|_{H^1}) &\leq K (|\Delta\mathcal{N}| + |c(t) - c_0| + |c(t) - c_0|^2 + |c(t) - c_0|^3) \\ &\quad + K (\|w\|_{L^2} + \|w\|_{L^2}^2). \quad \square \end{aligned}$$

6. Proof of main results. Before proving Theorem 1.1, we make an a priori estimate.

6.1. A priori estimates.

PROPOSITION 6.1. *Let $c_0 > n$, $a \leq a_*(c_0)$, and assume there exists $\varepsilon > 0$ such that*

$$\sigma(A_a) \cap \{\Re\lambda \geq -\varepsilon\} = \{0\}, \quad \lambda = 0 \text{ is an eigenvalue of algebraic multiplicity two.}$$

Let $T > 0$. There exist $\delta_ \in (0, 1)$ and $K_* \geq 1$ such that if the $\phi(x, t) - 1 \in C^1([0, T]; H^1 \cap H_a^1)$ solves (1.2) and satisfies, for $t \leq T$,*

$$(6.1) \quad \inf_x \phi(x, t) \geq \alpha_0 > 0,$$

$$(6.2) \quad \inf_x \phi(x, t)^m - a^2 \phi(x, t)^n \geq \beta_0 > 0,$$

and furthermore

- (i) *the decomposition $\phi(x, t) \mapsto (v(y, t), c(t), \theta(t))$ exists for $t \leq T$;*
- (ii) *for $t \leq T$*

$$(6.3) \quad \begin{aligned} &\sqrt{|\Delta\mathcal{N}|} + \sqrt{\|w(t)\|_{H^1}} + \sqrt{|c(t) - c_0|} \\ &+ |\theta(t) - \theta_0| + \left| 1 - \frac{c_0}{c(t) - \theta(t)} \right| + \|v(t)\|_{H^1} \leq \delta_*; \end{aligned}$$

- (iii) *the data satisfies*

$$(6.4) \quad \sqrt{|c(0) - c_0|} + |\theta(0) - \theta_0| + \sqrt{|\Delta\mathcal{N}|} + \sqrt{\|w(0)\|_{H^1}} \leq \epsilon < \delta_*;$$

then for $t \in [0, T]$,

$$(6.5) \quad \sqrt{e^{\kappa t} \|w(t)\|_{H^1}} + \sqrt{|c(t) - c_0|} + |\theta(t) - \theta_0| + \left| 1 - \frac{c_0}{c(t) - \theta(t)} \right| + \|v(t)\|_{H^1} \leq K_* \epsilon,$$

with $\kappa = \kappa(\alpha_0, \beta_0, \delta_) \in (0, b_{\max})$.*

Proof. The strategy for proving this proposition is to show that the left-hand side of (6.5) may be estimated in terms of their data, (6.4), using (6.3). This largely follows the proof in [21], with a few changes. The need to estimate $|1 - c_0/(c - \theta)|$ in terms of the data will require use of the modulation equations, (5.24), to control $\dot{\theta}$, and to control $\|w\|_{H^1}$, we will need to work in τ -time. Thus we make the following estimates.

Temporal change of variables. First, let us assume that $\delta_\star \leq \frac{1}{2}$. Then, since $c_0 > n > 1$, $\frac{1}{2}c_0 \leq c(t) \leq \frac{3}{2}c_0$. Furthermore, this initial choice of δ_\star ensures

$$\frac{1}{2} \leq \frac{c_0}{c(t) - \dot{\theta}(t)} = \frac{d\tau}{dt} \leq \frac{3}{2},$$

so the change of variables $\tau = \tau(t)$ is well defined.

Time derivatives of modulation parameters. Examining (5.24), $\mathcal{B}(t) = I + O(|c(t) - c_0|) + O(\|w\|_{L^2})$; so there exists $\delta_{\mathcal{B}} > 0$ such that for $\delta_\star \leq \delta_{\mathcal{B}}$, $\mathcal{B}(t)$ will be invertible. Therefore

$$(6.6) \quad |\dot{\theta}(t)| + |\dot{c}(t)| \leq |\mathcal{B}(t)^{-1}| \left(\|\tilde{S}_a w\|_{L^2} + \|\mathcal{G}_1\|_{L^2} \right) \leq K_1 (|c - c_0| + \|v\|_{H^1}) \|w\|_{L^2},$$

permitting the estimate

$$(6.7) \quad |\dot{\theta}(t)| + |\dot{c}(t)| \leq K_1 \delta_\star^3.$$

As terms of the form $1/(c - \dot{\theta})$ will appear, we will assume that

$$\delta_\star \leq \left(\frac{1}{2K_1} \right)^{1/3} = \delta_\theta,$$

so $|\dot{\theta}| \leq \frac{1}{2}$, and this quotient will be well defined and bounded.

Weighted perturbation. In τ -time, $\|w(\tau)\|_{H^1}$ is

$$w(\tau) = e^{A_a \tau} w(\tau(t=0)) + \int_{\tau(0)}^\tau e^{A_a(\tau-s)} Q\mathcal{G}(s) ds.$$

By the semigroup decay estimate of Theorem 4.1, there exist $K_2 > 0$ and $b_{\max} > 0$ such that

$$\|w(\tau)\|_{H^1} \leq K_2 e^{-b\tau} \|w(\tau(0))\|_{H^1} + K_2 \int_{\tau(0)}^\tau e^{-b(\tau-s)} \|Q\mathcal{G}(s)\|_{H^1} ds$$

for any $b \in (0, b_{\max})$. Estimating \mathcal{G} , we obtain

$$(6.8) \quad \begin{aligned} \|\mathcal{G}\|_{H^1} &\leq K \left[\left| \frac{c_0}{c - \dot{\theta}} \right| (|\dot{\theta}| + |\dot{c}|) + \left(|c - c_0| + \left| 1 - \frac{c_0}{c - \dot{\theta}} \right| \right) \|w\|_{H^1} \right] \\ &\quad + K \left[\left| \frac{c_0}{c - \dot{\theta}} \right| \|v\|_{H^1} \|w\|_{H^1} \right] \\ &\leq K_3 \delta_\star \|w\|_{H^1}. \end{aligned}$$

We have made use of (6.6) to control $\dot{\theta}$ and \dot{c} in terms of $\|w\|_{H^1}$.

Therefore,

$$(6.9) \quad \|w(\tau)\|_{H^1} \leq K_2 e^{-b\tau} \|w(\tau(0))\|_{H^1} + K_2 K_3 \delta_\star \int_{\tau(0)}^\tau e^{-b(\tau-s)} \|w(s)\|_{H^1} ds.$$

Defining $\psi(s) = e^{bs} \|w(s)\|_{H^1}$, (6.9) becomes

$$\psi(\tau) \leq K_2 \|w(0)\|_{H^1} + K_2 K_3 \delta_\star \int_{\tau(0)}^\tau \psi(s) ds,$$

for which we may apply Gronwall’s inequality to get

$$(6.10) \quad \|w(\tau)\|_{H^1} \leq K_2 \|w(\tau(0))\|_{H^1} e^{-(b-K_2K_3\delta_\star)(\tau-\tau(0))}.$$

So for δ_\star small enough, $b - K_2K_3\delta_\star > 0$ and we induce decay in the H^1 norm of w . In particular, suppose that $\delta_\star \leq \delta_b = \frac{1}{2}b/(K_2K_3)$ and let

$$b' = b - K_2K_3\delta_\star.$$

We then return to t -time,

$$\tau - \tau(0) = \frac{1}{c_0} \int_0^\tau c(s)ds + \frac{1}{c_0} (\theta(0) - \theta(\tau)) \geq \frac{1}{c_0} (c_0 - \delta_\star) \tau - 2\frac{\delta_\star}{c_0}.$$

Therefore,

$$(6.11) \quad \|w(t)\|_{H^1} \leq \tilde{K}_2 \|w(t=0)\|_{H^1} e^{-\kappa t}$$

with $\kappa = \frac{1}{2}b'$. We now have $e^{\kappa t} \|w\|_{H^1}$ estimated in terms of the data.

Unweighted perturbation. Applying this to Proposition 5.4, we have the estimate

$$(6.12) \quad \begin{aligned} \|v(t)\|_{H^1} &\leq K \left(\sqrt{|\Delta\mathcal{N}|} + \sqrt{|c(t) - c_0|} + |c(t) - c_0| + |c(t) - c_0|^{3/2} \right) \\ &\quad + K \left(\sqrt{\|w(t)\|_{L^2}} + \|w(t)\|_{H^1} \right) \\ &\leq K \left(\sqrt{|\Delta\mathcal{N}|} + \sqrt{|c(t) - c_0|} (1 + \delta_\star + \delta_\star^2) + \sqrt{\|w(t)\|_{L^2}} (1 + \delta_\star) \right) \\ &\leq K \left(\sqrt{|\Delta\mathcal{N}|} + \sqrt{|c(t) - c_0|} + \sqrt{\|w(0)\|_{L^2}} \right). \end{aligned}$$

If we had control of $\sqrt{|c(t) - c_0|}$, then we would also control $\|v(t)\|_{H^1}$ in terms of the data.

Deviation in c from c_0 . Estimating $|c(t) - c_0|$ using (6.6) and (6.11), we get

$$\begin{aligned} |c(t) - c_0| &\leq |c(0) - c_0| + \int_0^t |\dot{c}(s)| ds \leq |c(0) - c_0| \\ &\quad + \int_0^t K_1 (|c(s) - c_0| + \|v(s)\|_{H^1}) \|w(s)\|_{L^2} ds \\ &\leq |c(0) - c_0| + K_1 \delta_\star \int_0^t \|w(s)\|_{H^1} ds \leq |c(0) - c_0| \\ &\quad + K_1 \delta_\star \int_0^t K_2 \|w(t_0)\|_{H^1} e^{-\kappa s} ds \\ &\leq |c(0) - c_0| + K_1 K_2 \delta_\star \|w(0)\|_{H^1} / \kappa. \end{aligned}$$

So we now have $|c(t) - c_0|$ in terms of data, which we rewrite as

$$(6.13) \quad \sqrt{|c(t) - c_0|} \leq K_4 \left(\sqrt{|c(0) - c_0|} + \sqrt{\|w(0)\|_{H^1}} \right),$$

which in turn gives

$$(6.14) \quad \|v(t)\|_{H^1} \leq K_5 \left(\sqrt{|\Delta\mathcal{N}|} + \sqrt{|c(0) - c_0|} + \sqrt{\|w(0)\|_{H^1}} \right).$$

Deviation in θ from θ_0 . As with the speed parameter,

$$\begin{aligned} |\theta(t) - \theta_0| &\leq |\theta(0) - \theta_0| + \int_0^t |\dot{\theta}(s)| ds \leq |\theta(0) - \theta_0| \\ &\quad + \int_0^t K_1 (|c(s) - c_0| + \|v(s)\|_{H^1}) \|w(s)\|_{L^2} ds \\ &\leq |\theta(0) - \theta_0| + K_1 \delta_* \int_0^t \|w(s)\|_{H^1} ds \leq |c(0) - c_0| \\ &\quad + K_1 \delta_* \int_0^t K_2 \|w(t_0)\|_{H^1} e^{-\kappa s} ds \\ &\leq |\theta(0) - \theta_0| + K_1 K_2 \delta_* \|w(0)\|_{H^1} / \kappa. \end{aligned}$$

This is rewritten as

$$(6.15) \quad |\theta(t) - \theta_0| \leq K_7 \left(|\theta(0) - \theta_0| + \sqrt{\|w(0)\|_{H^1}} \right).$$

Another estimate on the temporal change of variables.

$$\left| 1 - \frac{c_0}{c - \theta} \right| \leq \frac{|c - c_0| + |\theta|}{|c - \theta|} \leq 2 \left(|c - c_0| + |\theta| \right).$$

Then using (6.6) and (6.13)

$$(6.16) \quad \left| 1 - \frac{c_0}{c - \theta} \right| \leq K \left(\sqrt{|c(0) - c_0|} + \|w(t)\|_{H^1} \right) K \leq (|c(0) - c_0| + K_2 \|w(0)\|_{H^1}) \\ \leq K_6 \left(\sqrt{|c(0) - c_0|} + \sqrt{\|w(0)\|_{H^1}} \right).$$

Combining (6.11), (6.13), (6.14), (6.15), and (6.16), we have (6.5) with $\delta_* = \min\{\frac{1}{2}, \delta_A, \delta_b, \delta_\tau, \delta_\theta\}$, $K_* = \max\{\tilde{K}_2, K_4, K_5, K_6, K_7\}$. \square

6.2. Main result. We now prove Theorem 1.1. Let c_* and \hat{v} be as in Theorem 3.1 (b). Take $c_0 \in (n, c_*)$ and $a \in (\gamma \hat{v}, a_*(c_0)]$. By Theorem 3.1 (b), there exists $\varepsilon > 0$ such that the only eigenvalue of $A_{c_0, a}$ with $\Re \lambda \geq -\varepsilon$ is at the origin. Thus we satisfy the first assumption of Proposition 6.1.

Define \mathcal{T} to be the set of nonnegative numbers, T , such that, given c_0, a , and $v_0 \in H^1 \cap H_a^1$,

- a solution exists, $\phi - 1 \in C([0, T], H^1 \cap H_a^1)$, satisfying

$$(6.17) \quad \inf_x \phi(x, t) \geq \frac{1}{4} = \alpha_0 > 0,$$

$$(6.18) \quad \inf_x (\phi(x, t)^m - a^2 \phi(x, t)^n) \geq \frac{1}{4} \inf_x (\phi_{c_0}(x)^m - a^2 \phi_{c_0}(x)^n) = \beta_0 > 0;$$

- a decomposition of ϕ into $(v(y(x, t), t), \theta(t), c(t))$ exists for $t \in [0, T]$;
- (6.3) holds for $t \in [0, T]$.

Set $T_* = \sup \mathcal{T}$. We will first show that there exists $\epsilon_* > 0$, such that for $\epsilon \leq \epsilon_*$, if

$$\|v_0\|_{H^1} + \|v_0\|_{H_a^1} \leq \epsilon,$$

then $T_* > 0$. This will be proved using the continuous dependence upon the data. Using Proposition 6.1, we will then prove $T_* = \infty$.

Let δ_* and K_* be as in Proposition 6.1.

The most difficult part of the proof will be ensuring the persistence of (6.3). Consider that, at $t = 0$, the left-hand side of that equation may be estimated with

$$\begin{aligned}
 \text{LHS}(t = 0) &\leq |\Delta\mathcal{N}| + \|v_0\|_{H^1} + \|\partial_y \phi_{c_0}\|_{H^1} |\theta_0 - \theta(0)| \\
 &\quad + \sqrt{|c(0) - c_0|} + |\theta(0) - \theta_0| + \|\phi_{c_0} - \phi_{c(0)}\|_{H^1} \\
 (6.19) \quad &\quad + \left| 1 - \frac{c_0}{c(0) - \dot{\theta}(0)} \right| + \sqrt{e^{a\theta(0)} \|v_0\|_{H_a^1}} \\
 &\quad + \sqrt{\|\partial_y \phi_{c_0}\|_{H_a^1} |\theta_0 - \theta(0)|} + \sqrt{\|\phi_{c_0} - \phi_{c(0)}\|_{H_a^1}}.
 \end{aligned}$$

There exists a choice of δ' and ϵ' such that if

$$(6.20) \quad |c(0) - c_0| + |\dot{c}(0)| + |\theta(0) - \theta_0| + |\dot{\theta}(0)| \leq \delta'$$

$$(6.21) \quad \|v_0\|_{H^1 \cap H_a^1} \leq \epsilon',$$

then the right-hand side of (6.19) will be bounded by $\frac{1}{2}\delta_*$. Set $\delta_1 = \min\{\frac{1}{4}, \delta'\}$. From Propositions 5.1 and 5.2, let δ_0, δ'_0 be the corresponding values for δ_1 .

There exists $\epsilon_{\text{exist}} \in (0, 1)$ such that if $\|v_0\|_{H^1 \cap H_a^1} \leq \epsilon_{\text{exist}}$, then

$$\begin{aligned}
 \inf_x [\phi_{c_0}(x + \theta_0) + v_0(x)] &\geq 2\alpha_0, \\
 \inf_x [(\phi_{c_0}(x + \theta_0) + v_0(x))^m - a^2(\phi_{c_0}(x + \theta_0) + v_0(x))^n] &\geq 2\beta_0.
 \end{aligned}$$

By Theorem 2.1, there exist $t_1 > 0$ and a solution in $C^1([0, t_1], H^1 \cap H_a^1)$ satisfying (6.17) and (6.18). Furthermore, we will have the a priori $H^1 \cap H_a^1$ bound that

$$\sup_{t \leq t_1} \|\phi(t) - 1\|_{H^1 \cap H_a^1} \leq 2\|\phi_{c_0}(\cdot + \theta_0) + v_0 - 1\|_{H^1 \cap H_a^1} \leq 2(\|\phi_{c_0} - 1\| + 1).$$

Set

$$(6.22) \quad \epsilon_1 = \min \left\{ \epsilon_{\text{exist}}, \epsilon', \frac{1}{2}e^{-a\theta_0} \delta_0 \right\}$$

and let $\|v_0\| \leq \epsilon_1$. As noted, the solution exists, satisfying (6.17) and (6.18), until at least $t_1 > 0$. At $t = 0$,

$$\|e^{a(\cdot + \theta_0)}(\phi_0 - \phi_{c_0}(\cdot + \theta_0))\|_{H^1} \leq e^{a\theta_0} \|v_0\|_{H^1} \leq \frac{1}{2}\delta_0,$$

so by the continuity of ϕ in time, we have

$$\|e^{a(\cdot + \theta_0)}(\phi(t) - \phi_{c_0}(\cdot - c_0 t + \theta_0))\|_{H^1} \leq \delta_0 \quad \text{for some } t_2 \in (0, t_1).$$

Therefore the decomposition exists, with the δ_1 bound on the modulation parameters, up till $t_2 > 0$.

Also at $t = 0$, using the δ_1 bound on the parameters,

$$(6.23) \quad \sqrt{|\Delta\mathcal{N}|} + \sqrt{\|w(0)\|_{H^1}} + \sqrt{|c(0) - c_0|} + \left| 1 - \frac{c_0}{c(0) - \dot{\theta}(0)} \right| + \|v(0)\|_{H^1} \leq \frac{1}{2}\delta_*.$$

All of these terms are continuous in time, and there exists some $t_3 \in (0, t_2)2$, for which this remains smaller than δ_* . Therefore, for $\epsilon_* \leq \epsilon_1$, $t_3 \in \mathcal{T}$ and $T_* > 0$.

Continuing to infinity. A few more constraints on ϵ_* are needed to continue out to $t = \infty$. There exists $\epsilon_{\alpha\beta}$ such that for $\epsilon \leq \epsilon_{\alpha\beta}$, if

$$\sqrt{|c - c_0|} + \|v\|_{H^1} \leq \epsilon,$$

then

$$\begin{aligned} \inf_y [\phi_c(y) + v(y)] &\geq \alpha_0, \\ \inf_y [(\phi_c(y) + v_0(y))^m - a^2(\phi_c(y) + v_0(y))^n] &\geq \beta_0. \end{aligned}$$

Let $\epsilon_2 > 0$ be so small that

$$(6.24) \quad K_*\epsilon_2 = \min \left\{ \frac{\delta_*}{3}, \sqrt{\delta_0/3}, \sqrt{\delta'_0}, \frac{1}{2}\epsilon_{\alpha\beta} \right\}$$

and set

$$(6.25) \quad \epsilon_* = \min \{ \epsilon_1, \epsilon_2 \}.$$

Now, assume $\|v_0\|_{H^1 \cap H^1_a} \leq \epsilon \leq \epsilon_*$. As above, for this data we will have $T_* > 0$. Assume $T_* < \infty$. For any $T < T_*$, on the interval $[0, T]$, the solution exists with (6.17) and (6.18), as does the decomposition, and (6.3) holds.

Then

$$\begin{aligned} \|\phi(t) - 1\|_{H^1 \cap H^1_a} &\leq \max \left\{ 1, e^{a(\int_0^t c(s) ds - \theta(t))} \right\} (\|\phi_{c(t)} - 1\|_{H^1 \cap H^1_a} + \|v(t)\|_{H^1 \cap H^1_a}) \\ &\leq e^{a((c_0 + \delta_*)T_* + |\theta_0| + \delta_*)} \left(\sup_{|c - c_0| \leq \delta_*} \|\phi_c - 1\|_{H^1 \cap H^1_a} + \delta_* + \delta_*^2 \right) < \infty, \end{aligned}$$

and this bound is uniform in $T < T_*$. By assumption, (6.17) and (6.18) hold for $t \in [0, T]$, uniformly in $T < T_*$, which may written as

$$\left\| \frac{1}{\phi(\cdot, t)} \right\|_{L^\infty} \leq \frac{1}{\alpha_0} < \infty \quad \text{and} \quad \left\| \frac{1}{\phi(\cdot, t)^m - a^2\phi(\cdot, t)^n} \right\|_{L^\infty} \leq \frac{1}{\beta_0} < \infty.$$

Therefore, according to (2.7), $\phi(x, t)$ may be extended beyond T_* by some amount $t_2 > 0$. Hence, if $T_* \neq \infty$, it must be a failure for either the decomposition to continue to exist or for (6.3) to hold.

Again using the Proposition 6.1 and our choice of ϵ_* ,

$$\|w(\cdot, t)\|_{H^1_a} \leq \delta_0/3 \quad \text{and} \quad |c(t) - c_0| \leq \delta'_0 \quad \text{for all } t \leq T, \text{ uniformly in } T < T_*.$$

Since ϕ exists until at least $T_* + t_2$, we may apply Proposition 5.2 to extend the decomposition for some amount $t_* \leq t_2$ also beyond T_* .

By assumption,

$$\sqrt{|c(t) - c_0|} + \|v(t)\|_{H^1} \leq K_*\epsilon_* \leq \frac{1}{2}\epsilon_{\alpha\beta} \quad \text{for } t < T_*.$$

Again, by continuity, these remain bounded by $\epsilon_{\alpha\beta}$ until some time $t_3 \in (0, t_*)$ beyond T_* , so (6.17) and (6.18) also persist beyond T_* .

We may now apply Proposition 6.1 past T_* . This gives

$$\sqrt{\|w(\cdot, t)\|_{H^1}} + \sqrt{|c(t) - c_0|} + |\theta(t) - \theta_0| + \left| 1 - \frac{c_0}{c(t) - \dot{\theta}(t)} \right| + \|v(\cdot, t)\|_{H^1} \leq K_* \epsilon_* \leq \frac{1}{3} \delta_*$$

for $t \leq T < T_*$. As $\sqrt{|\Delta \mathcal{N}|}$ is time invariant and smaller than $\frac{1}{2} \delta_*$,

$$\begin{aligned} &\sqrt{|\Delta \mathcal{N}|} + \sqrt{\|w(\cdot, t)\|_{H^1}} + \sqrt{|c(t) - c_0|} + |\theta(t) - \theta_0| \\ &+ \left| 1 - \frac{c_0}{c(t) - \dot{\theta}(t)} \right| + \|v(\cdot, t)\|_{H^1} \leq K_* \epsilon_* \leq \frac{5}{6} \delta_* \end{aligned}$$

for $t \leq T$, uniformly in $T < T_*$. But all of these functions are continuous for $t \in [0, T_* + t_3]$; so for some $t_4 > 0$, this expression remains bounded by δ_* . This contradicts $T_* < \infty$. So a solution exists for all time satisfying (6.17), (6.18), along with a decomposition and (6.3).

Since we may then apply Proposition 6.1 for all time, we will always have (6.5). By virtue of $K_* \epsilon_* < \delta_* \leq \delta_{\mathcal{B}}$, we will be able to invert the matrix $\mathcal{B}(t)$ for the modulation equations, (5.24). Therefore $|\dot{c}(t)| + |\dot{\theta}(t)| \leq K \epsilon e^{-\kappa t}$ and

$$\lim_{t \rightarrow \infty} c(t) = c_\infty$$

exists, and if we define

$$\lim_{t \rightarrow \infty} \left(\theta(t) + \int_0^t (c(s) - c_\infty) ds \right) = \theta_\infty,$$

then

$$\begin{aligned} &\|\phi(\cdot, t) - \phi_{c_\infty}(\cdot - c_\infty t + \theta_\infty)\|_{H^1} \\ &\leq K_* \epsilon + \|\phi_{c_\infty}(\cdot - c_\infty t + \theta_\infty) - \phi_{c(t)}(\cdot - \int_0^t c(s) ds + \theta(t))\|_{H^1} \\ &\leq K_* \epsilon + K |c(t) - c_\infty| + \|\partial_y \phi_{c_\infty}\|_{H^1} \left| \theta(t) + \int_0^t (c(s) - c_\infty) ds - \theta_\infty \right| \\ &\leq K_* \epsilon. \end{aligned}$$

Similarly

$$\|\phi(\cdot + c_\infty t - \theta_\infty, t) - \phi_{c_\infty}\|_{H_a^1} \leq K_* \epsilon e^{-\kappa t}. \quad \square$$

6.3. Remarks. This proof is equally applicable in the Hamiltonian case, $n+m = 0$, for values of c not in the discrete set, E , of points for which A_a has an imaginary eigenvalue.

7. Summary and discussion. We have thus shown that in the space $H^1 \cap H_a^1$, the solitary waves are asymptotically stable. This dovetails with an extension of global existence to data in a neighborhood of the solitary waves. In the Hamiltonian case, we can extend it beyond c_* via analytic continuation, as was done in [21], and this is analytically verified for $n = 2$, with computations in [40] suggesting it is true for all $n > 1$. Furthermore, to the extent that we will accept a computation of the Evans function as proof, our result generalizes to large amplitude solitary waves with $c > c_*$.

To our knowledge, this is the first result for which asymptotic stability is established for a conservative PDE in the absence of a variational principle.

Open problems include a weakening of the assumption of exponential decay on the perturbation. This might be accomplished through the use of an algebraic spatial weight, which would require the perturbation to decay algebraically rapidly. Yet less restrictive would be to use the approach of Merle and his colleagues [11, 18, 19, 22]. However, there is a tradeoff in both of these approaches: weakening the assumption on the spatial decay rate of the perturbation also weakens what can be proved about the rate at which the perturbation decays in time.

Finally we remark that the multidimensional case of (1.1) is wide open. While there is no existence theory for the two- and three-dimensional problems, perhaps a similar approach, of working in a neighborhood of a solitary wave, could be applied proving both existence and stability.

Appendix A. Properties of solitary waves.

A.1. Analyticity. Here we provide a proof of Corollary 2.9. Let us restate the crucial theorem.

THEOREM A.1 (Corollary 4.1.6 of [9]). *Suppose that f is a solution of the convolution equation $f = K * G(f)$ such that $f \in L^2 \cap L^\infty$ and $\lim_{|x| \rightarrow \infty} f(x) = 0$. If the Fourier transform \hat{K} of the integral kernel K satisfies the decay condition $|\hat{K}(\xi)| \leq A_1(1 + A_2|\xi|^m)$ for some constants $A_1, A_2 > 0$ and $m \geq 1$, and G is an infinitely differentiable function whose domain contains the range $R(f)$ of f , having all its derivatives bounded on $R(f)$ and satisfying the condition $G(0) = 0$, then $f, G(f) \in H^\infty$. In addition, if G is an analytic function on an open set U containing $R(f)$, G is continuous up to the boundary of ∂U of U , and*

$$(A.1) \quad d(\partial U, R(f)) > \|K\|_{L^2},$$

then there exists a constant $\sigma_0 > 0$ such that f and $G(f)$ both have analytic extensions to the strip $\{z \in \mathbb{C} : |\Im z| < \sigma_0\}$.

Let $u_c = \phi_c - 1$. By Theorem 2.6 and Corollary 2.7, u_c is positive, in $L^\infty \cap L^2$, and decays exponentially fast at $\pm\infty$. Using (2.14) and (2.15), the equation may be written as

$$(A.2) \quad -\gamma^2 u_c + \partial_x^2 u_c + \int_0^{u_c} \partial_\tau^3 F_2(1 + \tau; c) \frac{(u_c + 1 - \tau)^2}{2} d\tau = 0.$$

Define

$$(A.3) \quad G(z) = \int_0^z \partial_\tau^3 F_2(1 + \tau; c) \frac{(z + 1 - \tau)^2}{2} d\tau.$$

Taking a Fourier transform of (A.2), the equation is

$$-\gamma^2 \widehat{u}_c(\xi) - \xi^2 \widehat{u}_c(\xi) + \widehat{G(u_c)}(\xi) = 0.$$

This becomes the nonlinear convolution equation

$$(A.4) \quad u_c(x) = K * G(u_c)(x),$$

$$(A.5) \quad \hat{K}(\xi) = \frac{1}{\gamma^2 + \xi^2}.$$

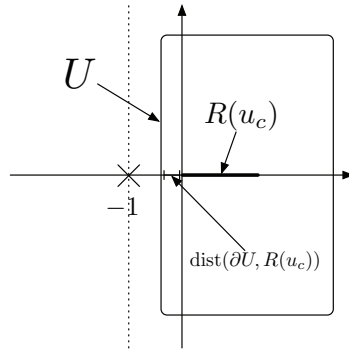


FIG. 6. A plot of a possible domain U and the range of u_c , $R(u_c)$. Note that the distance between these sets as drawn is the distance into the left-hand side of the complex plane that U extends, and that any such ovoid will be acceptable as long as it stays to the right of $\Re z = -1$.

For purposes of satisfying (A.1), let us take $\tilde{K}(x) = \alpha K(x)$ and $\tilde{G}(z) = \alpha^{-1}G(z)$ for $\alpha > 0$, where α is to be determined. Under this trivial scaling, $u_c = \tilde{K} * \tilde{G}(u_c)$.

\tilde{K} satisfies the decay estimate for Theorem A.1. $\tilde{G}(z)$ will have a singularity at $z = -1$ but is otherwise analytic. The range of u_c is the finite segment

$$R(u_c) = [0, u_{\max}],$$

and \tilde{G} is infinitely differential there, with all derivatives bounded. $G(0) = 0$. Hence, the first part of Theorem A.1 applies: u_c and $\tilde{G}(u_c)$ are in H^∞ .

Now, consider the set U in Figure 6. In this figure, $d(\partial U, R(u_c))$ is the distance U stretches into the left half-plane:

$$\|\tilde{K}\|_2 = \alpha \sqrt{\frac{\pi}{2}} \gamma^{-3/2}.$$

Picking α so small that the norm is less than 1, we can find a U such that the distance between ∂U and $R(u_c)$ exceeds $\|\tilde{K}\|_2$, satisfying (A.1) and proving analyticity in a strip. \square

A.2. Continuity as a function of speed. Consider the functional

$$(A.6) \quad \mathcal{F}[c, u] = \partial_x^2 u + F_2(1 + u; c)$$

as a mapping from $H^2 \times \mathbb{R} \rightarrow L^2$. The solitary wave $u_c = \phi_c - 1$ satisfies $\mathcal{F}[c, u_c] = 0$. Using this functional, we prove Corollary 2.11 via the implicit function theorem. Given a particular $\hat{c} > n$, set $\hat{u} = u_{\hat{c}}$.

Let H^2_{even} and L^2_{even} be the subspaces of H^2 and L^2 , respectively, of only even functions. Define the sets

$$M_0 = \left(\hat{c} - \frac{\hat{c} - n}{2}, \hat{c} + \frac{\hat{c} - n}{2} \right) \subset \mathbb{R},$$

$$N_0 = \left\{ u \in H^2_{\text{even}} : \|u - \hat{u}\|_{H^2} \leq \frac{1}{2} \right\} \subset H^2_{\text{even}},$$

$$Z = L^2_{\text{even}}.$$

Note that set M_0 is bounded away from zero and all of the functions in N_0 are uniformly bounded from below by $\frac{1}{2}$. Therefore \mathcal{F} is well defined on $M_0 \times N_0$ and will be a C^1 mapping on this set into Z .

Set

$$T = \frac{\delta \mathcal{F}}{\delta c} [\hat{c}, \hat{u}] = \partial_c F_2(1 + \hat{u}; \hat{c}),$$

$$S = \frac{\delta \mathcal{F}}{\delta u} [\hat{c}, \hat{u}] = \partial_x^2 + \partial_u F_2(1 + \hat{u}; \hat{c}).$$

T and S are bounded operators on $\mathbb{R} \rightarrow L^2_{\text{even}}$ and $H^2_{\text{even}} \rightarrow L^2_{\text{even}}$, respectively.

Let $f \in L^2_{\text{even}}$ and consider the problem $Su = f$. As an elliptic problem, this has a solution, provided $f \perp \ker(S^\dagger)$. Note that $S\partial_x \hat{u} = 0$. S is self-adjoint, has smooth coefficients, and is in one spatial dimension; this is the unique element of the kernel. But $\partial_x \hat{u}$ is an odd function, hence f is orthogonal to it and the equation has a solution u satisfying a bound

$$\|u\|_{H^2} \leq K \|u\|_{L^2}.$$

Because the coefficients in S and the right-hand side, f , are all even functions, $\tilde{u}(x) = u(-x)$ also solves $Su = f$. By the uniqueness of the solution, $u = \tilde{u}$, so u is an even function. Therefore $u \in H^2_{\text{even}}$ and the map $S : H^2_{\text{even}} \rightarrow L^2_{\text{even}}$ is onto with bounded inverse.

The kernel of S , restricted to $u \in H^2_{\text{even}}$, is trivial, so the implicit function theorem may be applied to conclude the existence of a function $\mathcal{G} : M_1 \rightarrow H^2_{\text{even}}$, $\hat{c} \in M_1 \subset M_0$, such that

$$\mathcal{F}[c, \mathcal{G}(c)] = 0$$

for all $c \in M_1$. The mapping \mathcal{G} is C^1 . For any such c ,

$$\partial_x^2 \mathcal{G}(c) + F_2(\mathcal{G}(c) + 1; c) = 0.$$

This is just the solitary wave equation. Therefore $\mathcal{G}(c) = u_c = \phi_c - 1$, and the mapping $c \mapsto \phi_c - 1$ is $C^1(\mathbb{R}; H^2)$. The analyticity of the mapping may be proven by checking the analyticity of \mathcal{F} in a neighborhood of (\hat{c}, \hat{u}) .

To prove continuity in $H^2 \cap H^2_a$, the proof is similar. Fixing $a < \frac{1}{2}$, and taking $\hat{c} \in (n/(1 - 4a^2), \infty)$, we let $M_a = M_0 \cap (n/(1 - 4a^2), \infty)$, $N_a = N_0 \cap H^2_a$, and $Z_a = Z \cap L^2_a$. We must check that $S : H^2_{\text{even}} \cap H^2_a \rightarrow L^2_{\text{even}} \cap L^2_a$ is onto with bounded inverse. This is accomplished using the previous result and studying $S_a = D^2_a + \partial_u F_2(1 + \hat{u}; \hat{c})$ on $H^2_{\text{even}} \rightarrow L^2_{\text{even}}$.

Appendix B. Perturbation expansions. Here we provide some explicit calculations, including those for Proposition 2.13.

Note the expansions

$$\phi^n = \phi_c^n + n\phi_c^{n-1}v + f_n[\phi_c, v]v,$$

$$\phi^m = \phi_c^m + m\phi_c^{m-1}v + f_m[\phi_c, v]v,$$

$$f_p[a, b] = \int_0^1 \left[p(a + \tau b)^{p-1} - pa^{p-1} \right] d\tau,$$

and

$$(B.1) \quad \begin{aligned} H_\phi^{-1} &= H_{\phi_c}^{-1} - H_{\phi_c}^{-1} B [n\phi_c^{n-1}v + f_n[\phi_c, v]v, m\phi_c^{m-1}v + f_m[\phi_c, v]v] H_{\phi_c}^{-1} \\ &\quad + H_{\phi_c}^{-1} B[\phi^n - \phi_c^n, \phi^m - \phi_c^m] \left(H_\phi^{-1} - H_{\phi_c}^{-1} \right), \end{aligned}$$

$$(B.2) \quad B[a, b]u = -\partial_x (a\partial_x u) + bu.$$

Recall (1.2),

$$\partial_t \phi = \dot{c}\partial_c \phi_c + (\dot{\theta} - c) \partial_y \phi_c + (\dot{\theta} - c) \partial_y v + v_t = -(\phi_c + v)^n H_{\phi_c+v}^{-1} \partial_y (\phi_c + v)^n.$$

Using the above expansions and the solitary wave equation, $c\partial_y \phi_c = \phi_c^m H_{\phi_c}^{-1} \partial_y (\phi_c^n)$, this may be expanded into

$$(B.3) \quad v_t = \phi_c^m H_{\phi_c}^{-1} \partial_y [-c\phi_c^n \partial_y^2 (\phi_c^{-m} v) + cv - n\phi_c^{n-1}v - cn\phi_c^{n-1} \partial_y (\phi_c^{-m} \partial_y \phi_c) v]$$

$$(B.4) \quad \begin{aligned} & - \dot{\theta} \partial_y v - \dot{c} \partial_c \phi_c - \dot{\theta} \partial_y \phi_c \\ & - f_m[\phi_c, v]v H_{\phi_c}^{-1} \partial_y (\phi^n) + m\phi_c^{m-1}v H_{\phi_c}^{-1} (H_\phi - H_{\phi_c}) H_\phi^{-1} \partial_y (\phi^n) \\ & - m\phi_c^{m-1}v H_{\phi_c}^{-1} \partial_y (\phi^n - \phi_c^n) - \phi_c^m H_{\phi_c}^{-1} \partial_y (f_n[\phi_c, v]v) \end{aligned}$$

$$(B.5) \quad \begin{aligned} & + \phi_c^m H_{\phi_c}^{-1} B[f_n[\phi_c, v]v, f_m[\phi_c, v]v] H_\phi^{-1} \partial_y (\phi^n) \\ & - \phi_c^m H_{\phi_c}^{-1} B[n\phi_c^{n-1}v, m\phi_c^{m-1}v] H_{\phi_c}^{-1} (H_\phi - H_{\phi_c}) H_\phi^{-1} \partial_y (\phi^n) \\ & + \phi_c^m H_{\phi_c}^{-1} B[n\phi_c^{n-1}v, m\phi_c^{m-1}v] H_{\phi_c}^{-1} \partial_y (\phi^n - \phi_c^n). \end{aligned}$$

(B.3) is a linear term. (B.4) will decay to zero as $\theta(t)$ and $c(t)$, the modulating parameters, approach their asymptotic limits. (B.5) is purely nonlinear in v .

We define $\mathcal{F}_1[v; \phi_c]$, the term nonlinear in v , as

$$(B.6) \quad \begin{aligned} \mathcal{F}_1[v; \phi_c] &= -f_m[\phi_c, v]v H_{\phi_c}^{-1} \partial_y (\phi^n) + m\phi_c^{m-1}v H_{\phi_c}^{-1} (H_\phi - H_{\phi_c}) H_\phi^{-1} \partial_y (\phi^n) \\ & - m\phi_c^{m-1}v H_{\phi_c}^{-1} \partial_y (\phi^n - \phi_c^n) - \phi_c^m H_{\phi_c}^{-1} \partial_y (f_n[\phi_c, v]v) \\ & + \phi_c^m H_{\phi_c}^{-1} B[f_n[\phi_c, v]v, f_m[\phi_c, v]v] H_\phi^{-1} \partial_y (\phi^n) \\ & - \phi_c^m H_{\phi_c}^{-1} B[n\phi_c^{n-1}v, m\phi_c^{m-1}v] H_{\phi_c}^{-1} (H_\phi - H_{\phi_c}) H_\phi^{-1} \partial_y (\phi^n) \\ & + \phi_c^m H_{\phi_c}^{-1} B[n\phi_c^{n-1}v, m\phi_c^{m-1}v] H_{\phi_c}^{-1} \partial_y (\phi^n - \phi_c^n). \end{aligned}$$

The operator S is given by

$$(B.7) \quad \begin{aligned} S[c_0, c, \dot{\theta}] &= mc_0 \left(\phi_{c_0}^{-1} \partial_y \phi_{c_0} - \frac{c}{c - \dot{\theta}} \phi_c^{-1} \partial_y \phi_c \right) \\ & + n\phi_{c_0}^m H_{c_0}^{-1} \partial_y [(\phi_{c_0}^{-1} - c_0(\phi_{c_0}^{-1} - 1)) \cdot] \\ & - \frac{nc_0}{c - \dot{\theta}} \phi_c^m H_{\phi_c}^{-1} \partial_y [(\phi_c^{-1} - c(\phi_c^{-1} - 1)) \cdot] \\ & - c_0 m \left\{ \phi_{c_0}^m H_{\phi_{c_0}}^{-1} [\phi_{c_0}^{-1} \partial_y \phi_{c_0} \cdot] - \frac{c}{c - \dot{\theta}} \phi_c^m H_{\phi_c}^{-1} [\phi_c^{-1} \partial_y \phi_c \cdot] \right\}. \end{aligned}$$

Finally, the terms making up \mathcal{G}_1 from (2.36):

$$(B.8) \quad \mathcal{G}_1 = \tilde{\mathcal{G}}_1 + \tilde{\mathcal{G}}_2 + \tilde{\mathcal{G}}_3 + \tilde{\mathcal{G}}_4 + \tilde{\mathcal{G}}_5 + \tilde{\mathcal{G}}_6 + \tilde{\mathcal{G}}_7,$$

$$(B.9) \quad \tilde{\mathcal{G}}_1 = -f_m[\phi_c, v]wH_{\phi_c}^{-1}\partial_y(\phi^n),$$

$$(B.10) \quad \tilde{\mathcal{G}}_2 = m\phi_c^{m-1}wH_{\phi_c}^{-1}(H_\phi - H_{\phi_c})H_{\phi_c}^{-1}\partial_y(\phi^n),$$

$$(B.11) \quad \tilde{\mathcal{G}}_3 = -m\phi_c^{m-1}wH_{\phi_c}^{-1}\partial_y(n\phi_c^{n-1}v + f_n[\phi_c, v]v),$$

$$(B.12) \quad \tilde{\mathcal{G}}_4 = -\phi_c^m H_{\phi_c, a}^{-1} D_a (f_n[\phi_c, v]w),$$

$$(B.13) \quad \tilde{\mathcal{G}}_5 = \phi_c^m H_{\phi_c, a}^{-1} B_a [f_n[\phi_c, v]w, f_m[\phi_c, v]w] H_{\phi_c}^{-1} \partial_y(\phi^n),$$

$$(B.14) \quad \tilde{\mathcal{G}}_6 = -\phi_c^m H_{\phi_c, a}^{-1} B_a [n\phi_c^{n-1}w, m\phi_c^{m-1}w] H_{\phi_c}^{-1} (H_\phi - H_{\phi_c}) H_{\phi_c}^{-1} \partial_y(\phi^n),$$

$$(B.15) \quad \tilde{\mathcal{G}}_7 = \phi_c^m H_{\phi_c, a}^{-1} B_a [n\phi_c^{n-1}w, m\phi_c^{m-1}w] H_{\phi_c}^{-1} \partial_y(\phi^n - \phi_c^n).$$

The difference between A_a and A_a^∞ may be written as

$$(B.16) \quad \begin{aligned} A_a - A_a^\infty &= -cm\phi_c^{-1}\partial_y\phi_c - (\phi_c^m - 1)H_{\phi_c, a}^{-1}D_a(n\phi_c^{n-1}\cdot) \\ &\quad + cm\phi_c^m H_{\phi_c, a}^{-1}(\phi_c^{-1}\partial_y\phi_c\cdot) - cn\phi_c^m H_{\phi_c, a}^{-1}D_a[(1 - \phi_c^{-1})\cdot] \\ &\quad + H_{1, a}^{-1}D_a[n(1 - \phi_c^{n-1})\cdot] + H_{1, a}^{-1}(\phi_c^m - 1)H_{\phi_c, a}^{-1}D_a(n\phi_c^{n-1}\cdot) \\ &\quad + H_{1, a}^{-1}D_a(1 - \phi_c^n)D_aH_{\phi_c, a}^{-1}D_a(n\phi_c^{n-1}\cdot). \end{aligned}$$

In the space weighted by $\phi_c(x)^{-m}$, this difference is

$$(B.17) \quad \begin{aligned} \tilde{A}_a - \tilde{A}_a^\infty &= nH_{1, a}^{-1}D_a[(1 - \phi_c^{m-1})\cdot] + nH_{1, a}^{-1}D_a(\phi_c^m - 1)H_{1, a}^{-1}(\phi_c^{m-1}\cdot) \\ &\quad - nH_{1, a}^{-1}[D_a(\phi_c^m - 1)]H_{1, a}^{-1}(\phi_c^{m-1}\cdot) \\ &\quad - nH_{1, a}^{-1}(\phi_c^m - 1)H_{1, a}^{-1}(\phi_c^m - 1)H_{\phi_c, a}^{-1}D_a(\phi_c^{m-1}\cdot) \\ &\quad + nH_{1, a}^{-1}(\phi_c^m - 1)H_{1, a}^{-1}D_a(\phi_c^n - 1)D_aH_{\phi_c, a}^{-1}D_a(\phi_c^{m-1}\cdot) \\ &\quad - nH_{1, a}^{-1}D_a(\phi_c^n - 1)D_aH_{\phi_c, a}^{-1}D_a(\phi_c^{m-1}\cdot) + cmH_{1, a}^{-1}(\phi_c^{m-1}\partial_y\phi_c\cdot) \\ &\quad + cmH_{1, a}^{-1}(\phi_c^m - 1)H_{\phi_c, a}^{-1}(\phi_c^{m-1}\partial_y\phi_c\cdot) \\ &\quad - cmH_{1, a}^{-1}D_a(\phi_c^n - 1)D_aH_{\phi_c, a}^{-1}(\phi_c^{m-1}\partial_y\phi_c\cdot) \\ &\quad + cnH_{1, a}^{-1}D_a[\phi_c^m(\phi_c^{-1} - 1)\cdot] \\ &\quad + cnH_{1, a}^{-1}(\phi_c^m - 1)H_{\phi_c, a}^{-1}D_a[\phi_c^m(\phi_c^{-1} - 1)\cdot] \\ &\quad - cnH_{1, a}^{-1}D_a(\phi_c^n - 1)D_aH_{\phi_c, a}^{-1}D_a[\phi_c^m(\phi_c^{-1} - 1)\cdot]. \end{aligned}$$

Appendix C. Analysis of the characteristic polynomial. In this section we prove that (3.36), $(\lambda - c\mu)(1 - \mu^2) + n\mu = 0$, has a unique root of minimal real part on a slit half-plane

$$\{\lambda : \Re\lambda > -\lambda_0\} \setminus (-\lambda_0, -\tilde{\Omega}(\gamma)].$$

There are two ways that this could be false: there could be either a multiple root or two roots with the same real part, but differing imaginary parts. As previously

noted, we will have a unique root of minimal real part for λ in the closed right half-plane, so we need only concern ourselves with $\Re\lambda < 0$.

We will identify a portion of the domain $\Re\lambda < 0$ for which there are neither multiple roots nor complex roots with the same real part. Note that this is a stricter condition than is needed, as the polynomial could have a double root for some λ , where the third root of $P(\mu)$ has a smaller real part than the multiple root.

Note that $P(\pm 1)$ never vanishes, and hence $P(\mu) = 0$ is equivalent to $R(\mu) = \lambda$, where

$$(C.1) \quad R(\mu) = c\mu + \frac{n\mu}{\mu^2 - 1}.$$

C.1. Roots of order greater than one. We start with the possibility of a double or triple root, as this is very easy to rule out. If μ is a multiple root, then in addition to $R(\mu) = \lambda$, we will also have

$$\frac{dR}{d\mu} = c - n \frac{1 + \mu^2}{(1 - \mu^2)^2} = 0,$$

which has solutions

$$\mu = -\sqrt{\frac{2c + n \pm \sqrt{8cn + n^2}}{2c}}.$$

We have ignored the roots with a + sign in front, as these will correspond to positive λ . Note that they are all real, and hence λ will also be real.

The λ one gets from the root

$$\mu_+ = -\sqrt{\frac{2c + n + \sqrt{8cn + n^2}}{2c}} \leq -\frac{3\sqrt{3}}{2}n$$

is decreasing in c . The other root,

$$\mu_- = -\sqrt{\frac{2c + n - \sqrt{8cn + n^2}}{2c}},$$

will map to λ values

$$(C.2) \quad \lambda(\mu_-) = -\sqrt{\frac{1}{8}\sqrt{8c^2 + 20cn - n^2 - 8c\sqrt{n^2 + 8cn} - n\sqrt{n^2 + 8cn}}}.$$

It can be checked that for $c > n > 1$, $R(\mu_+) < R(\mu_-)$. Therefore, for $\lambda > -\tilde{\Omega}(c)$, with

$$(C.3) \quad \tilde{\Omega}(c) = \sqrt{\frac{1}{8}\sqrt{8c^2 + 20cn - n^2 - 8c\sqrt{n^2 + 8cn} - n\sqrt{n^2 + 8cn}}},$$

$P(\mu)$ cannot have a multiple root.

C.2. Roots of differing imaginary part. If $\mu_1 = \alpha + i\beta_1$ and $\mu_2 = \alpha + i\beta_2$ are two roots of P , then

$$R(\alpha + i\beta_1) = R(\alpha + i\beta_2).$$

After matching real and imaginary parts in this expression, the three unknowns, α, β_1, β_2 , must satisfy the two equations

$$\begin{aligned} \Re\lambda &= c\alpha + \frac{n\alpha(-1 + \alpha^2 + \beta_1^2)}{(-1 + \alpha^2)^2 + 2(1 + \alpha^2)\beta_1^2 + \beta_1^4} \\ &= c\alpha + \frac{n\alpha(-1 + \alpha^2 + \beta_2^2)}{(-1 + \alpha^2)^2 + 2(1 + \alpha^2)\beta_2^2 + \beta_2^4}, \end{aligned} \tag{C.4}$$

$$\begin{aligned} \Im\lambda &= c\beta_1 - \frac{n\beta_1(1 + \alpha^2 + \beta_1^2)}{(-1 + \alpha^2)^2 + 2(1 + \alpha^2)\beta_1^2 + \beta_1^4} \\ &= c\beta_2 - \frac{n\beta_2(1 + \alpha^2 + \beta_2^2)}{(-1 + \alpha^2)^2 + 2(1 + \alpha^2)\beta_2^2 + \beta_2^4}. \end{aligned} \tag{C.5}$$

Solving the (C.4) for β_2^2 in terms of α and β_1 , there are two families of solutions:

$$\beta_2^2 = \beta_1^2, \tag{C.6}$$

$$\beta_2^2 = \frac{(1 - \alpha^2)(3 + \beta_1^2 + \alpha^2)}{\alpha^2 + \beta_1^2 - 1}. \tag{C.7}$$

Without loss of generality, we assume $\beta_1 \neq 0$.

Recall that λ is imaginary if and only if $P(\mu)$ has a purely imaginary root; hence the condition $\Re\lambda < 0$ ensures $\alpha \neq 0$. Then (C.7) implies $0 < |\alpha| \leq 1$. Furthermore, if $|\alpha| = 1$, then either the roots are conjugate or $\beta_2 = 0$. But if $\beta_2 = 0$, then $\mu_2 = 1$, which we know is not a root of $P(\mu)$.

C.2.1. Complex conjugates. When $\beta_1 = -\beta_2 = \beta$, (C.5) implies that λ is real and

$$c - \frac{n(1 + \alpha^2 + \beta^2)}{(-1 + \alpha^2)^2 + 2(1 + \alpha^2)\beta^2 + \beta^4} = 0, \tag{C.8}$$

which has roots β^2 ,

$$\beta^2 = -1 - \alpha^2 + \frac{n}{2c} \pm \frac{\sqrt{n^2 + 16c^2\alpha^2}}{2c}. \tag{C.9}$$

We may immediately rule out the negative root for β^2 . For $\beta^2 > 0$, α^2 must satisfy

$$1 + \frac{n}{2c} - \frac{\sqrt{n^2 + 24c}}{2c} < \alpha^2 < 1 + \frac{n}{2c} + \frac{\sqrt{n^2 + 24c}}{2c}. \tag{C.10}$$

Using (C.4),

$$\lambda(\alpha) = \frac{n + 8c\alpha^2 - \sqrt{n^2 + 16c^2\alpha^2}}{4\alpha}.$$

Since we are concerned only with $\lambda < 0$ here, α must, in addition to (C.10), satisfy

$$\frac{n + 8c\alpha^2 - \sqrt{n^2 + 16c^2\alpha^2}}{4\alpha} < 0. \tag{C.11}$$

When $\alpha > 0$, (C.11) requires

$$0 < \alpha^2 < \frac{c - n}{4c}.$$

But $\frac{1}{4}(c - n)/c < 1 + n/(2c) - \sqrt{n^2 + 24c}/(2c)$, so complex conjugate roots with $\alpha > 0$ are not possible with λ in the left half-plane.

When $\alpha < 0$, (C.11) requires

$$\alpha < -\frac{1}{2}\sqrt{\frac{c-n}{c}}$$

to satisfy (C.11). Consider α in the interval

$$\left[-\sqrt{1 + \frac{n}{2c} - \frac{\sqrt{n^2 + 8cn}}{2c}}, -\frac{1}{2}\sqrt{\frac{c-n}{c}} \right).$$

In this interval, there will not be complex conjugate roots, as it violates (C.10).

λ , as a function of α , is negative and increasing on this interval. For λ in the image of this interval, we may completely rule out complex conjugate roots. The image of this interval is

$$\lambda \left(\left[-\sqrt{1 + \frac{n}{2c} - \frac{\sqrt{n^2 + 8cn}}{2c}}, -\frac{1}{2}\sqrt{\frac{c-n}{c}} \right] \right) = [-\tilde{\Omega}(c), 0).$$

Hence, for $0 > \lambda > -\tilde{\Omega}$, one may rule out both multiple roots and complex conjugates.

C.2.2. Nonconjugate complex roots. Consider the case of complex roots with the same real part, but imaginary parts such that $|\beta_1| \neq |\beta_2|$. Squaring both sides of (C.5) and plugging in (C.7) for β_2^2 , we get a sixth order polynomial in β_1^2 . The roots, as functions of α , are

$$(C.12) \quad \beta_1^2 = -1 - 2\alpha - \alpha^2,$$

$$(C.13) \quad \beta_1^2 = -1 + 2\alpha - \alpha^2,$$

$$(C.14) \quad \beta_1^2 = 1 - \alpha^2 - 2\sqrt{1 - \alpha^2},$$

$$(C.15) \quad \beta_1^2 = 1 - \alpha^2 + 2\sqrt{1 - \alpha^2},$$

$$(C.16) \quad \beta_1^2 = \frac{(n/c)^2 + 4(n/c)(1 - \alpha^2) - 8(1 - \alpha^4)}{8(1 - \alpha^2)} - \frac{\sqrt{((n/c)^2 - 16(1 - \alpha^2))((n/c)^2 - 8(1 - \alpha^2)(2\alpha^2 - (n/c)))}}{8(1 - \alpha^2)},$$

$$(C.17) \quad \beta_1^2 = \frac{(n/c)^2 + 4(n/c)(1 - \alpha^2) - 8(1 - \alpha^4)}{8(1 - \alpha^2)} + \frac{\sqrt{((n/c)^2 - 16(1 - \alpha^2))((n/c)^2 - 8(1 - \alpha^2)(2\alpha^2 - (n/c)))}}{8(1 - \alpha^2)}.$$

(C.12) and (C.13) force β_1 to be imaginary, and hence they can be ruled out. Using (C.7), if β_1 is either (C.14) or (C.15), then $\beta_2 = \pm\beta_1$, conjugate roots.

In the last two cases, if β_1^2 is to be real, then

$$(C.18) \quad ((n/c)^2 - 16(1 - \alpha^2))((n/c)^2 - 8(1 - \alpha^2)(2\alpha^2 - (n/c))) \geq 0.$$

If we can find for λ in the left half-plane sufficiently close enough to the imaginary axis that it is negative, we will be done. (C.18) is negative at $\alpha = 0$, so there exists a neighborhood of the imaginary axis, such that complex nonconjugate roots may be ruled out.

The roots on the left-hand side of (C.18) are

$$(C.19) \quad \alpha^2 = 1 - \frac{1}{16} \left(\frac{n}{c}\right)^2 = \frac{1}{16}(5 - \gamma^2)(3 + \gamma^2),$$

$$(C.20) \quad \alpha^2 = \frac{1}{2} + \frac{1}{4} \frac{n}{c} + \frac{1}{2} \sqrt{1 - \frac{n}{c}} = \frac{1}{4}(3 - \gamma)(1 + \gamma),$$

$$(C.21) \quad \alpha^2 = \frac{1}{2} + \frac{1}{4} \frac{n}{c} - \frac{1}{2} \sqrt{1 - \frac{n}{c}} = \frac{1}{4}(1 - \gamma)(3 + \gamma).$$

These are positive for $\gamma \in [0, 1]$. It may be checked that (C.21) is the smallest for all γ . Hence the root of (C.18) such that μ will be closest to the imaginary axis is

$$(C.22) \quad \alpha = -\frac{1}{2} \sqrt{(1 - \gamma)(3 + \gamma)}.$$

α larger than (C.22) and less than zero will yield a λ that does not have nonconjugate complex roots.

Given λ , if $\mathcal{P}(\mu; \lambda)$ is to have two roots of same real part, α , but differing imaginary part, then, by trying any of the last four roots for β_1 (both positive and negative square roots of (C.16) and (C.17)), in (C.4) the real part of λ and α are related by

$$(C.23) \quad \Re\lambda = c\alpha \left(2 - \frac{(n/c)}{4(1 - \alpha^2) + (n/c)} \right).$$

Note that

$$(C.24) \quad \frac{d\Re\lambda}{d\alpha} = c \left(\frac{32(1 - \alpha^2)^2 + 12(n/c) - 20\alpha^2(n/c) + (n/c)^2}{(4(1 - \alpha^2) + (c/n))^2} \right)$$

and the derivative has one negative root with $|\alpha| < 1$ at

$$(C.25) \quad \alpha = -\sqrt{1 + \frac{5}{16} \frac{n}{c} - \frac{1}{16} \sqrt{64 \frac{n}{c} + 17 \left(\frac{n}{c}\right)^2}}.$$

Comparing (C.22) with (C.25), at $\gamma = 1$

$$(C.25) = -1 < 0 = (C.22),$$

and at $\gamma = 0$

$$(C.25) = -\sqrt{\frac{3}{4}} = (C.22).$$

In addition, one may check that (C.25) is increasing in $\gamma \in [0, 1]$, while (C.22) is decreasing on the same interval. Therefore

$$-\sqrt{1 + \frac{5}{16} \frac{n}{c} - \frac{1}{16} \sqrt{64 \frac{n}{c} + 17 \left(\frac{n}{c}\right)^2}} \leq -\frac{1}{2} \sqrt{(1 - \gamma)(3 + \gamma)}$$

for all $\gamma \in [0, 1]$.

For $\Re\mu = \alpha$ in

$$(C.26) \quad \left(-\frac{1}{2}\sqrt{(1-\gamma)(3+\gamma)}, 0 \right)$$

(C.23) will be an increasing function in α , and (C.24) is positive at $\alpha = 0$. On the interval (C.26), the mapping is invertible and its image is

$$(C.27) \quad \left(-\frac{1}{4}c\sqrt{1-\gamma}(\gamma+3)^{3/2}, 0 \right).$$

Therefore if λ has real part in the interval (C.27), and $\mathcal{P}(\mu; \lambda)$ is to have nonconjugate complex roots, α must lie in (C.26). But such an α violates (C.18), and we may conclude that there are no such roots. Letting $-\lambda_0$ denote the $\Re\lambda$ value at the left end point of (C.27),

$$(C.28) \quad \lambda_0 = \frac{c}{4}\sqrt{1-\gamma}(\gamma+3)^{3/2} = \frac{n}{4(1-\gamma^2)}\sqrt{1-\gamma}(\gamma+3)^{3/2}.$$

Hence for $\Re\lambda > -\lambda_0$, nonconjugate complex roots are not possible for $P(\mu)$.

Appendix D. The zero eigenvalue. In section 3.5, $\lambda = 0$ was identified as an eigenvalue of multiplicity at least two. Using the Evans function, the order of this eigenvalue may be related to the slope with respect to c of the invariant functional $\mathcal{N}[\phi_c]$.

Using the framework from section 3.3.2, set

$$(D.1) \quad \mathbf{y}^+ = \begin{pmatrix} \phi_c^{-m} \partial_x \phi_c \\ \partial_x (\phi_c^{-m} \partial_x \phi_c) \\ 0 \end{pmatrix},$$

$$(D.2) \quad \mathbf{z}^- = \left(c\partial_x (\phi_c^{-m} \partial_x \phi_c), -c\phi_c^{-m} \partial_x \phi_c, -\int_{-\infty}^x \phi_c^{-n-m} \partial_x \phi_c \right).$$

These are solutions to the dynamical systems

$$\dot{\mathbf{y}} = B(x, \lambda = 0, \gamma)\mathbf{y} \quad \text{and} \quad \dot{\mathbf{z}} = -\mathbf{z}B(x, \lambda = 0, \gamma).$$

Here, $\mu_1 = -\gamma$. Employing the notation and formulation of the Evans function of [27], by [27, Proposition 1.6, parts 2 and 3], since

$$\mathbf{y}^+ = O(e^{-\gamma x}) \quad \text{as } x \rightarrow +\infty \quad \text{and} \quad \mathbf{z}^- = O(e^{\gamma x}) \quad \text{as } x \rightarrow -\infty,$$

\mathbf{y}^+ and \mathbf{z}^- are scalar multiples of ζ^+ and η^- , respectively. ζ^+ and η^- are the solutions of the dynamical systems satisfying

$$\zeta^+ e^{\gamma x} \rightarrow \mathbf{v}^+ \quad \text{as } x \rightarrow +\infty \quad \text{and} \quad \eta^- e^{-\gamma x} \rightarrow \mathbf{w}^- \quad \text{as } x \rightarrow -\infty,$$

with

$$\begin{aligned} B^\infty \mathbf{v}^+ &= -\gamma \mathbf{v}^+, & \mathbf{w}^- B^\infty &= -\gamma \mathbf{w}^-, \\ \mathbf{v}^+ &= (1, -\gamma, 0)^T, & \mathbf{w}^- &= (2c\gamma^2)^{-1} (c\gamma^2, -c\gamma, -1), \end{aligned}$$

allowing us to define the Evans function as

$$D(\lambda = 0) = \eta^-(x, \lambda = 0) \cdot \zeta^+(x, \lambda = 0).$$

From the properties of the solitary waves, discussed in section 2.2, there exists $\beta > 0$ such that

$$\phi_c^{-m} \partial_x \phi_c e^{\gamma x} \rightarrow -\beta \text{ as } x \rightarrow +\infty,$$

and hence

$$\zeta^+ = \frac{1}{-\beta} \mathbf{y}^+ \quad \text{and} \quad \eta^- = \frac{1}{-2\gamma c \beta} \mathbf{z}^-.$$

$D(0) = 0$ by inspection. From [27], the derivative of the Evans function for a system akin to ours is²

$$(D.3) \quad \partial_\lambda D(\lambda) = - \int_{-\infty}^{\infty} \eta^-(x, \lambda) \partial_\lambda [B(x, \lambda) - \mu_1(\lambda)I] \zeta^+(x, \lambda) dx.$$

It is then trivial to compute that $\partial_\lambda D(0) = 0$, since the integrand is an odd function.

Taking the derivative of (D.3) at $\lambda = 0$ gives an equation for $\partial_\lambda^2 D(0)$,

$$(D.4) \quad \partial_\lambda^2 D(0) = - \int_{-\infty}^{\infty} \eta_\lambda^- B_\lambda \zeta^+ dx - \int_{-\infty}^{\infty} \eta^- B_\lambda \zeta_\lambda^+ dx.$$

ζ_λ^+ and η_λ^- satisfy the ODEs

$$\dot{\mathbf{y}}_\lambda = B \mathbf{y}_\lambda + B_\lambda \mathbf{y} \quad \text{and} \quad \dot{\mathbf{z}}_\lambda = -\mathbf{z}_\lambda B - \mathbf{z} B_\lambda.$$

These problems are associated with the derivatives with respect to λ of (3.24) and (3.30) at $\lambda = 0$,

$$(D.5) \quad \partial_x L_c Y_\lambda = [I - \partial_x (\phi_c^n \partial_x (\phi_c^{-m} \cdot))] Y,$$

$$(D.6) \quad -L_c^* \partial_x Z_\lambda = [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] Z,$$

through the identifications

$$\begin{aligned} y_\lambda^{(1)} &= \phi_c^{-m} Y_\lambda, & z_\lambda^{(1)} &= c \partial_x (\phi_c^n \partial_x Z_\lambda), \\ y_\lambda^{(2)} &= \partial_x (\phi_c^{-m} Y_\lambda), & z_\lambda^{(2)} &= -c \phi_c^n \partial_x Z_\lambda, \\ y_\lambda^{(3)} &= L_c Y_\lambda + \phi_c^n \partial_x (\phi_c^{-m} Y), & z_\lambda^{(3)} &= -Z_\lambda. \end{aligned}$$

For $\lambda = 0$, (3.24), (3.30), (D.5), and (D.6) are related to the generalized kernels of A and A^* :

$$\begin{aligned} Y &= \partial_x \phi_c, & Y_\lambda &= -\partial_c \phi_c, \\ Z &= \int_{-\infty}^x \frac{\partial_x \phi_c}{\phi_c^{n+m}} dx, & Z_\lambda &= - \int_{-\infty}^x (L_c^*)^{-1} [I - \phi_c^{-m} \partial_x (\phi_c^n \partial_x \cdot)] \int_{-\infty}^x \frac{\partial_x \phi_c}{\phi_c^{n+m}} dx. \end{aligned}$$

²In particular, a system for which the matrix B has the same limit at $\pm\infty$.

Then

$$\zeta_\lambda^+ = \frac{1}{-\beta} \left(-\phi_c^{-m} \partial_c \phi_c, \quad -\partial_x (\phi_c^{-m} \partial_c \phi_c), \quad -L_c \partial_c \phi_c + \phi_c^n \partial_x (\phi_c^{-m} \partial_x \phi_c) \right)^T,$$

$$\eta_\lambda^- = \frac{1}{2\gamma c \beta} (c \partial_x (\phi_c^n \partial_x Z_\lambda), \quad -c \phi_c^n \partial_x Z_\lambda, \quad -Z_\lambda).$$

Finally, we compute

$$(D.7) \quad \partial_\lambda^2 D(0) = \frac{1}{c\gamma\beta^2} \partial_c \mathcal{N}[\phi_c].$$

Acknowledgments. We thank Marc Spiegelman for his helpful comments and support, in addition to his contributions appearing in [38]. We have also benefited from discussions with Professor J. L. Bona and Professor P. Rosenau.

REFERENCES

- [1] E. AHARONOV, J. A. WHITEHEAD, P. B. KELEMEN, AND M. SPIEGELMAN, *Channeling instability of upwelling melt in the mantle*, J. Geophys. Res., 100 (1995), pp. 20433–20450.
- [2] V. BARCILON AND O. M. LOVERA, *Solitary waves in magma dynamics*, J. Fluid Mech., 204 (1989), pp. 121–133.
- [3] V. BARCILON AND F. M. RICHTER, *Non-linear waves in compacting media*, J. Fluid Mech., 165 (1986), pp. 429–448.
- [4] T. B. BENJAMIN, *The stability of solitary waves*, Proc. Roy. Soc. London Ser. A, 328 (1972), pp. 153–183.
- [5] D. BERCOVICI, Y. RICARD, AND G. SCHUBERT, *A two-phase model for compaction and damage I. General theory*, J. Geophys. Res., 106 (2001), pp. 8887–8906.
- [6] M. S. BERGER, *Nonlinearity and Functional Analysis*, Academic Press, New York, 1977.
- [7] J. L. BONA, *personal correspondence*.
- [8] J. BONA, *On the stability theory of solitary waves*, Proc. Roy. Soc. London Ser. A, 344 (1975), pp. 363–374.
- [9] J. L. BONA AND Y. A. LI, *Decay and analyticity of solitary waves*, J. Math. Pures Appl., 76 (1997), pp. 377–430.
- [10] D. E. EDMUNDS AND W. D. EVANS, *Spectral Theory and Differential Operators*, Oxford Math. Monogr., Clarendon Press, Oxford University Press, New York, 1987.
- [11] K. EL DIKA, *Asymptotic stability of solitary waves for the Benjamin-Bona-Mahony equation*, C. R. Math. Acad. Sci. Paris, 337 (2003), pp. 649–652.
- [12] G. HIRTH AND D. L. KOHLSTEDT, *Experimental constraints on the dynamics of the partially molten upper mantle 2: Deformation in the dislocation creep regime*, J. Geophys. Res., 100 (1995), pp. 15441–15052.
- [13] G. HIRTH AND D. L. KOHLSTEDT, *Experimental constraints on the dynamics of the partially molten upper mantle: Deformation in the diffusion creep regime*, J. Geophys. Res., 100 (1995), pp. 1981–2002.
- [14] T. KATO, *Perturbation Theory for Linear Operators*, Springer-Verlag, Berlin, 1995.
- [15] R. F. KATZ, M. SPIEGELMAN, AND B. HOLTZMAN, *The dynamics of melt and shear localization in partially molten aggregates*, Nature, 442 (2006), pp. 676–679.
- [16] K. KNOPP, *Theory of Functions: Part I*, Dover, New York, 1996.
- [17] K. KNOPP, *Theory of Functions: Part II*, Dover, New York, 1996.
- [18] Y. MARTEL AND F. MERLE, *A Liouville theorem for the critical generalized Korteweg-de Vries equation*, J. Math. Pures Appl., 79 (2000), pp. 339–425.
- [19] Y. MARTEL AND F. MERLE, *Asymptotic stability of solitons of the subcritical gKdV equations revisited*, Nonlinearity, 18 (2005), pp. 55–80.
- [20] D. MCKENZIE, *The generation and compaction of partially molten rock*, J. Petrology, 25 (1984), pp. 713–765.
- [21] J. R. MILLER AND M. I. WEINSTEIN, *Asymptotic stability of solitary waves for the regularized long-wave equation*, Comm. Pure Appl. Math., 49 (1996), pp. 399–441.
- [22] T. MIZUMACHI, *Asymptotic stability of solitary wave solutions to the regularized long-wave equation*, J. Differential Equations, 200 (2004), pp. 312–341.

- [23] M. NAKAYAMA AND D. P. MASON, *Rarefactive solitary waves in two-phase fluid flow of compacting media*, Wave Motion, 15 (1992), pp. 357–392.
- [24] P. OLSON AND U. CHRISTENSEN, *Solitary wave propagation in a fluid conduit within a viscous matrix*, J. Geophys. Res., 91 (1986), pp. 6367–6374.
- [25] A. PAZY, *Semigroups of Linear Operators and Applications to Partial Differential Equations*, Springer-Verlag, New York, 1983.
- [26] R. L. PEGO, P. SMEREKA, AND M. I. WEINSTEIN, *Oscillatory instability of traveling waves for a KdV-Burgers equation*, Phys. D, 67 (1993), pp. 45–65.
- [27] R. L. PEGO AND M. I. WEINSTEIN, *Eigenvalues, and instabilities of solitary waves*, Philos. Trans. Roy. Soc. London Ser. A, 340 (1992), pp. 47–94.
- [28] R. L. PEGO AND M. I. WEINSTEIN, *Asymptotic stability of solitary waves*, Comm. Math. Phys., 164 (1994), pp. 305–349.
- [29] R. L. PEGO AND M. I. WEINSTEIN, *Convective linear stability of solitary waves for Boussinesq equations*, Stud. Appl. Math., 99 (1997), pp. 311–375.
- [30] J. PRÜSS, *On the spectrum of C_0 -semigroups*, Trans. Amer. Math. Soc., 284 (1984), pp. 847–857.
- [31] J. RENNER, K. VISCKUPIC, G. HIRTH, AND B. EVANS, *Melt extraction from partially molten peridotites*, Geochemistry, Geophysics, Geosystems, 4 (2003), article 8606; doi:10.1029/2002GC000369.
- [32] P. ROSENAU, *On a model equation of traveling and stationary compactons*, Phys. Lett. A, 356 (2006), pp. 44–50.
- [33] M. SCHECHTER, *Spectra of Partial Differential Operators*, North-Holland, Amsterdam, 1971.
- [34] M. SCHECHTER, *Principles of Functional Analysis*, Grad. Stud. Math. 36, AMS, Providence, RI, 2001.
- [35] D. R. SCOTT AND D. J. STEVENSON, *Magma solitons*, Geophys. Res. Lett., 11 (1984), pp. 1161–1164.
- [36] D. R. SCOTT AND D. J. STEVENSON, *Magma ascent by porous flow*, J. Geophys. Res., 91 (1986), pp. 9283–9296.
- [37] D. R. SCOTT, D. J. STEVENSON, AND J. A. WHITEHEAD, *Observations of solitary waves in a viscously deformable pipe*, Nature, 319 (1986), pp. 759–761.
- [38] G. SIMPSON, M. SPIEGELMAN, AND M. I. WEINSTEIN, *Degenerate dispersive equations arising in the study of magma dynamics*, Nonlinearity, 20 (2007), pp. 21–49.
- [39] G. SIMPSON, M. SPIEGELMAN, AND M. I. WEINSTEIN, *Magma transport: Dynamics of fluid flow in viscously deformable porous media*, in preparation.
- [40] G. SIMPSON, M. I. WEINSTEIN, AND P. ROSENAU, *On a Hamiltonian PDE arising in magma dynamics*, Discrete Contin. Dyn. Syst. Ser. B, 10 (2008), pp. 903–924.
- [41] M. SPIEGELMAN, *Flow in deformable porous media. Part 1: Simple analysis*, J. Fluid Mech., 247 (1993), pp. 17–38.
- [42] M. SPIEGELMAN, *Flow in deformable porous media. Part 2: Numerical analysis*, J. Fluid Mech., 247 (1993), pp. 39–63.
- [43] M. SPIEGELMAN, *Linear analysis of melt band formation by simple shear*, Geochemistry, Geophysics, Geosystems, 4 (2003), pp. 1525–2027.
- [44] M. SPIEGELMAN, P. B. KELEMEN, AND E. AHARONOV, *Causes and consequences of flow organization during melt transport: The reaction infiltration instability in compactible media*, J. Geophys. Res., 106 (2001), pp. 2061–2077; available online from <http://www.ldeo.columbia.edu/~mspieg/SolFlow/>.
- [45] D. A. WARK AND E. B. WATSON, *Grain-scale permeabilities of texturally equilibrated, monomineralic rocks*, Earth and Planetary Sci. Lett., 164 (1998), pp. 591–605.
- [46] D. A. WARK, C. A. WILLIAMS, E. B. WATSON, AND J. D. PRICE, *Reassessment of pore shapes in microstructurally equilibrated rocks, with implications for permeability of the upper mantle*, J. Geophys. Res., 108 (2003).
- [47] M. I. WEINSTEIN, *Modulational stability of ground states of nonlinear Schrödinger equations*, SIAM J. Math. Anal., 16 (1985), pp. 472–491.
- [48] M. I. WEINSTEIN, *Lyapunov stability of ground states of nonlinear dispersive evolution equations*, Comm. Pure Appl. Math., 39 (1986), pp. 51–68.
- [49] J. A. WHITEHEAD, *A laboratory demonstration of solitons using a vertical watery conduit in syrup*, Amer. J. Phys., 55 (1987), pp. 998–1003.
- [50] J. A. WHITEHEAD AND K. R. HELFRICH, *The Korteweg-de Vries equation from laboratory conduit and magma migration equations*, Geophys. Res. Lett., 13 (1986), pp. 545–546.
- [51] C. WIGGINS AND M. SPIEGELMAN, *Magma migration and magmatic solitary waves in 3-d*, Geophys. Res. Lett., 22 (1995), pp. 1289–1292.
- [52] W. ZHU AND G. HIRTH, *A network model for permeability in partially molten rocks*, Earth and Planetary Sci. Lett., 212 (2003), pp. 407–416.

ENERGY TRANSPORT BY ACOUSTIC MODES OF HARMONIC LATTICES*

LISA HARRIS[†], JANI LUKKARINEN[‡], STEFAN TEUFEL[§], AND FLORIAN THEIL[†]

Abstract. We study the large scale evolution of a scalar lattice excitation u which satisfies a discrete wave equation in three dimensions, $\ddot{u}_t(\gamma) = -\sum_{\gamma'} \alpha(\gamma - \gamma') u_t(\gamma')$, where $\gamma, \gamma' \in \mathbb{Z}^3$ are lattice sites. We assume that the dispersion relation ω associated to the elastic coupling constants $\alpha(\gamma - \gamma')$ is acoustic; i.e., it has a singularity of the type $|k|$ near the vanishing wave vector, $k = 0$. To derive equations describing the macroscopic energy transport, we employ a related multiscale Wigner transform and a scale parameter $\varepsilon > 0$. The spatial and temporal scales of the Wigner transform are related to the corresponding lattice parameters via a scaling by ε . In the continuum limit, which is achieved by sending the parameter ε to 0, the Wigner transform disintegrates into three different limit objects: the Wigner transform of a rescaled weak- L^2 limit, an H-measure, and a Wigner measure. The first two provide the finer resolution of the energy concentrating at $k = 0$ so that a set of closed evolution equations may arise. We demonstrate that these three limit objects satisfy a set of decoupled transport equations: a wave equation for the weak limit, a geometric optics transport equation for the H-measure limit, and a dispersive transport equation for the standard limiting Wigner measure. This yields a complete characterization of macroscopic energy transport in harmonic lattices with regular acoustic dispersion relations.

Key words. microlocal analysis, multiscale methods, homogenization, discrete wave equations

AMS subject classifications. 70J30, 74Q15, 37K60

DOI. 10.1137/070699184

1. Introduction. The energy transport by atomic oscillations in crystalline solids is a central problem in solid state physics. To the first order approximation, the oscillations can be described by a discrete wave equation

$$(1.1) \quad \ddot{u}_t(\gamma) = -\sum_{\gamma'} \alpha(\gamma - \gamma') u_t(\gamma')$$

where $u(\gamma) \in \mathbb{R}$ is composed of the displacements of the crystal atoms from their equilibrium position, as will be discussed in section 1.1. To analyze physically relevant properties of the crystal, such as its thermal conductivity, we first need to understand how energy is transported within the crystal via purely harmonic vibrations. Such transport properties are determined by the dispersion relation ω of the crystal, here $\omega(k) = \sqrt{\hat{\alpha}(k)}$, the “hat” denoting a discrete Fourier transform. If ω is not smooth, then depending on the wavelength different types of continuum energy transport equations can arise. The different contributions to the continuum energy are described by different limit objects, so-called Wigner measures, or microlocal defect measures.

*Received by the editors August 3, 2007; accepted for publication (in revised form) June 24, 2008; published electronically November 5, 2008.

<http://www.siam.org/journals/sima/40-4/69918.html>

[†]Mathematics Institute, Warwick University, Coventry, CV4 7AL, UK (l.c.harris@warwick.ac.uk, f.theil@warwick.ac.uk).

[‡]Zentrum Mathematik, Technische Universität München, Boltzmannstr. 3, D-85747 Garching, Germany, and Department of Mathematics and Statistics, University of Helsinki, P.O. Box 68, FI-00014 Helsingin yliopisto, Finland (jani.lukkarinen@helsinki.fi). This author was supported by the Deutsche Forschungsgemeinschaft (DFG) projects Sp 181/19-1 and Sp 181/19-2 and by the Academy of Finland.

[§]Mathematisches Institut, Auf der Morgenstelle 10, 72076 Tübingen, Germany (stefan.teufel@uni-tuebingen.de).

These concepts were introduced by Tartar [23] and Gérard [10] to analyze the weak limits of certain nonlinear quantities, the energy density being one of them. Two-scale Wigner measures of the type we use here were first introduced independently by Fermanian Kammerer [5, 6] and by Nier [18].

Our main interest is to characterize the macroscopic evolution of the energy density. The starting point of our mathematical analysis is the *Wigner transform*, which can be interpreted as a “wavenumber resolved” energy density. Let us leave the details for section 2.1, and only summarize the main findings here. Let $v_t = \dot{u}_t$. We can then construct a complex field $\psi_t \in \ell_2(\mathbb{Z}^3)$, corresponding to a normal mode of the oscillations, out of the fields u_t and v_t . On the other hand, for any given $\psi \in \ell_2(\mathbb{Z}^3)$ and $\varepsilon > 0$, the lattice Wigner transform $W^\varepsilon = W^\varepsilon[\psi]$ allows defining a corresponding “energy” density, $e^\varepsilon = e^\varepsilon[\psi]$, by the formula

$$(1.2) \quad e^\varepsilon(x) = \int_{\mathbb{T}^3} dk W^\varepsilon(x, k),$$

where $\varepsilon > 0$ denotes the “lattice spacing” and $x \in \mathbb{R}^3$ is a variable which interpolates between the points of the scaled lattice $\varepsilon\mathbb{Z}^3$. Here \mathbb{T}^3 denotes the 3-torus, and we identify $\mathbb{T}^3 = \mathbb{R}^3/\mathbb{Z}^3$. Then $e^\varepsilon(x, t) = e^\varepsilon[\psi_{t/\varepsilon}](x)$ can be considered to be a proper energy density at a macroscopic time t , as it satisfies $\int_{\mathbb{R}^3} dx e^\varepsilon(x, t) = H(u_{t/\varepsilon}, v_{t/\varepsilon})$, where H denotes the Hamiltonian related to (1.1).

In general, a limit of a sequence of scaled Wigner transforms, $(W^\varepsilon[\psi^\varepsilon])$ with ε tending to 0, is given by a nonnegative Radon measure $\mu \in M_+(\mathbb{R}^3 \times \mathbb{T}^3)$. In the setup considered here, we have for a given time t a sequence of normal mode fields $\psi_{t/\varepsilon}$ which would then yield a family of nonnegative Radon measures μ_t as limit points of the corresponding Wigner transforms. However, for these limit measures to form a useful approximation of the original dynamical system, μ_t should also satisfy an autonomous evolution equation. This is typically not possible if the initial measure concentrates on the singular set of ω , i.e., to the points where ω is not smooth.

Here we will augment the standard Wigner transform scheme to encompass the most common type of singularity encountered in solid state physics: acoustic singularities, with $\omega(k)$ behaving like $|k|$ near $k = 0$. Such modes occur in general within crystal models with short range interactions, and they are particularly important as they are responsible for sound propagation in the crystal. With some effort, in the spirit of the results proven in [6], it is likely that these results could be extended to cover more complicated singular sets of ω , but we do not consider such generalizations here. Our result can also be seen as a generalization of the analysis in [17]. There it is shown that μ_t can be computed from μ_0 by solving a dispersive linear transport equation, provided that μ_0 does not concentrate on wavenumbers k where the dispersion relation ω is not C^1 .

For acoustic modes, the dispersion relation fails to be C^1 at $k = 0$, and it is necessary to resolve finer details of the solutions near the singular point, whenever a concentration at $k = 0$ is possible. This will be accomplished here by introducing a second scale to the lattice Wigner transform. The multiscale Wigner transform, defined in section 2.1, maps a field ψ into a distribution $W^\varepsilon[\psi](x, k, q)$ where the new parameter $q \in \mathbb{R}^3$ resolves the energy density for wavenumbers of the order of ε . We consider a sequence of initial conditions such that the corresponding initial fields ψ_0^ε are bounded and tight in $\ell_2(\varepsilon\mathbb{Z}^3)$. We prove in Theorem 3.2 that then for all $t \in \mathbb{R}$ there are two measures, a Wigner measure μ_t on $\mathbb{R}^3 \times \mathbb{T}_*^3$, $\mathbb{T}_*^3 = \mathbb{T}^3 \setminus \{0\}$, and an H-measure μ_t^H on $\mathbb{R}^3 \times S^2$, as well as an L^2 -function ϕ_t such that $W^\varepsilon[\psi_{t/\varepsilon}^\varepsilon]$ converges

along a subsequence to a limit determined by (μ_t, μ_t^H, ϕ_t) . The subsequence can be chosen independently of t , and it will only be relevant for determining the limit of the initial data, that is, (μ_0, μ_0^H, ϕ_0) . For all other times $t \in \mathbb{R}$, the measures μ_t, μ_t^H and the L^2 -function ϕ_t can be determined using the transport equations

$$(1.3) \quad \partial_t \mu_t(x, k) + \frac{1}{2\pi} \nabla \omega(k) \cdot \nabla_x \mu_t(x, k) = 0, \quad k \in \mathbb{T}_*^3,$$

$$(1.4) \quad \partial_t \mu_t^H(x, q) + \frac{1}{2\pi} \nabla \omega_0(q) \cdot \nabla_x \mu_t^H(x, q) = 0, \quad q \in S^2,$$

$$(1.5) \quad \partial_t^2 \phi_t(x) = \operatorname{div} \left(\frac{1}{(2\pi)^2} A_0 \nabla \phi_t(x) \right),$$

where $x \in \mathbb{R}^3$, together with the initial conditions

$$(1.6) \quad \mu_t|_{t=0} = \mu_0, \quad \mu_t^H|_{t=0} = \mu_0^H, \quad \phi_t|_{t=0} = \phi_0, \quad \partial_t \hat{\phi}_t(p)|_{t=0} = -i\omega_0(p) \hat{\phi}_0(p), \quad p \in \mathbb{R}^3,$$

given by the limit objects (μ_0, μ_0^H, ϕ_0) at $t = 0$. Here the constant matrix A_0 and the function ω_0 are determined by the Hessian of the square of the dispersion relation ω (defined in (2.6)) at $k = 0$, explicitly

$$(1.7) \quad A_0 = \frac{1}{2} D^2 \omega^2(0), \quad \omega_0(p) = \sqrt{p \cdot A_0 p}.$$

It also follows from our analysis that the sum of the energies related to μ_t, μ_t^H , and ϕ_t is a constant of motion and equals the limiting value of the total energy of the initial excitations. This shows that no energy is lost in the taking of the continuum limit.

Equation (1.3) describes the propagation of energy along the harmonic lattice with the group velocity $\nabla \omega(k)/(2\pi)$. Equation (1.5) is a wave equation describing the evolution of macroscopic fluctuations. Equation (1.4) is usually known under the name “geometric optics” and describes the evolution of macroscopic fluctuations whose wavelength is much longer than the lattice spacing ε and much smaller than 1, the wavelength of the fluctuations resolved by ϕ_t . The surprising feature about acoustic modes is the separate evolution equation for the macroscopic perturbations ϕ_t . As will be apparent in Theorem 3.2, but as also indicated by the fact that the evolution equation for ϕ_t involves only the Hessian of ω^2 at $k = 0$, its evolution is determined by the second, fine-resolution scale introduced to the Wigner transform. This is in contradistinction to the case of ω which is smooth, as then the evolution of the energy density would not depend on the second scaling parameter.

The Wigner transform, or the Wigner function, was originally introduced to study semiclassical behavior in quantum mechanics, but it has been proven to be a useful tool in studying large scale behavior of wave equations as well [12, 19]. In particular, the method of calculating continuous, macroscopic energy by finding the limit object of a sequence of energies e^ε on rescaled lattice models is one that has been widely used and justified, for example, in [9, 14, 17]. In [17], the Wigner transform of the normal modes is employed in solving the macroscopic transport of energy in the above harmonic systems for deterministic initial data. The same system is considered in [3] with random initial data and in a larger function space, excluding, however, the type of concentration effects we study here. In [6] two-scale Wigner measures of the same type we use here were introduced for the study of concentration effects near shock hypersurfaces for the heat equation. The precise connection to the results in [5, 6] will be explained in Remark 3.6 in section 3.

In [8] spatially homogeneous Wigner measures are employed to study energy asymptotics in static systems. Both the method (which was developed before the

theory of 2-microlocal measures became available) and the results are in spirit very similar to ours. Using specific profile functions, three limit objects are constructed which capture the asymptotic behavior of the energy along bounded sequences. Only Fourier methods are used in [8] and therefore the spatial distribution of the energy cannot be characterized.

The main use of the Wigner transform is that, unlike the energy density e^ε itself, it contains enough information to satisfy a closed evolution equation in the limit $\varepsilon \rightarrow 0$. Indeed, it was shown in [17] that as long as there is no concentration on the singular set of the dispersion relation faster than $\varepsilon^{1/2}$, the Wigner transform of the time-evolved state vector $\psi_{t/\varepsilon}$ converges to a limit measure $\tilde{\mu}_t$ on $\mathbb{R} \times \mathbb{K}$. Here \mathbb{K} is a suitable compactification of $\mathbb{T}^3 \setminus \mathbb{S}$, \mathbb{S} being the singular set, which allows a continuous extension of the group velocity $\nabla\omega$. The measure $\tilde{\mu}_t$ is then proven to satisfy the transport equation

$$(1.8) \quad \partial_t \tilde{\mu}_t(x, k) + \frac{1}{2\pi} \nabla\omega(k) \cdot \nabla_x \tilde{\mu}_t(x, k) = 0.$$

However, the above assumption about the rate of convergence excludes macroscopic variations of the initial data in the case of acoustic singularities, and enforces $\phi_t = 0$ in the present terminology. Separately, it was also shown in [17] that if the weak limit of the rescaled initial data exists, then the limit satisfies a continuum wave equation. However, a combination of these results was not possible, leaving open the exact effect the continuum wave equation will have on the energy density. In [15] it is shown that the Wigner measures generated by sequences of solutions of discrete wave equations satisfy transport equations; energy conservation along the sequence is not discussed.

In this paper we will show how to overcome the above difficulties in a physically relevant class of models with a singular dispersion relation. Our main additional contribution to the results of the above-mentioned references is to introduce a lattice version of the multiscale Wigner transform and to employ this in solving the macroscopic evolution of the energy density in the presence of an acoustic singularity. In section 2 we will first present the microscopic dynamical model in detail. In section 2.1 we will define the Wigner transform, and discuss its relation to the energy of the microscopic lattice model. The main results will be presented in section 3.

1.1. Relation with solid state physics. A crystal in solid state physics is a state of matter in which the atoms retain a nearly perfect periodic structure over macroscopic times. The Hamiltonian model used for the time evolution in such a crystal is, to the first order accuracy, harmonic. If we assume that each periodic cell of the idealized perfectly periodic crystal structure contains n atoms, then we can form a vector $q(\gamma) \in \mathbb{R}^{3n}$ out of the displacements of the atoms in the periodic cell labeled by $\gamma \in \mathbb{Z}^3$. The (classical) Hamiltonian equations of motion of this harmonic model are then

$$(1.9) \quad \dot{q}_i(\gamma, t) = \frac{1}{m_i} p_i(\gamma, t), \quad \dot{p}_i(\gamma, t) = - \sum_{\gamma', i'} \mathbb{A}(\gamma - \gamma')_{i, i'} q_{i'}(\gamma', t),$$

where $\gamma \in \mathbb{Z}^3$, $i = 1, \dots, 3n$, and m_i denotes the mass of the atom whose displacement q_i measures.

By the change of variables to $\tilde{q}_i(\gamma) = m_i^{1/2} q_i(\gamma)$, $\tilde{p}_i(\gamma) = m_i^{-1/2} p_i(\gamma)$, these equations can be transformed into a standard form whose force matrix is given by $\tilde{\mathbb{A}}(\gamma)_{i, i'} = m_i^{-1/2} \mathbb{A}(\gamma)_{i, i'} m_{i'}^{-1/2}$. The standard form equations can then be solved by Fourier transform, and a diagonalization of the remaining multiplicative evolution equations

decomposes the $3n$ vector degrees of freedom into independent *normal modes*, called *phonons* in solid state physics. Each normal mode is a complex scalar field on the crystal lattice, and its time evolution is unitary and uniquely determined by the corresponding dispersion relation $\omega_i(k)$ on \mathbb{T}^3 . More details on the related mathematical issues can be found in [3, 17].

The complete decoupling of the evolution equations of the normal modes allows us to study a scalar, single-mode model and still retain a straightforward applicability of the results in the above, physically more relevant, vector models. As the above procedure is fairly standard and well covered in the references, we will not present the decomposition in any more detail here. However, to illustrate the matter, we have performed the equivalent steps for the scalar model (1.1) in the beginning of section 2.

In solid state physics, the normal modes are divided into *optical* and *acoustic* depending on their regularity at $k = 0$: if the dispersion relation is regular at $k = 0$, then the mode is called optical, and if it behaves as $|k|$ at $k = 0$, the mode is called acoustic. The latter name arises as these modes are believed to be responsible for the propagation of sound waves in the crystal. Acoustic dispersion relations arise generally from atomistic Hamiltonians with short range interactions, and crystal models are typically expected to have three of them, related to the translation invariance of the microscopic forces.

Finally, let us remark that the discrete *linear* wave equation alone does not suffice to determine the physically relevant properties of the crystal, such as its thermal conductivity. However, it forms the basis for perturbative treatments which can address such questions. We refer to [1, 20, 21, 22, 25] for further details on the physical aspects of the topic, and to [2] for a review about related open problems.

2. The microscopic model. There are two mathematically equivalent descriptions of harmonic crystals. On the one hand, one can work with the Hamiltonian equations of motion and analyze the properties of the solutions. Anharmonic crystals can then be discussed in the same manner. On the other hand, one can employ the Fourier transform and the linearity to condense the Hamiltonian into the dispersion relation. Although the second approach leads to an immediate solution of the harmonic system, it cannot be used directly to analyze nonlinear models.

We will show in this chapter how the first approach reduces to the second one for harmonic lattices. Then we demonstrate how the Wigner transform can be employed in the analysis of the time evolution of the energy density of the Hamiltonian description.

We assume that the scalar excitation $u_t(\gamma)$, $\gamma \in \mathbb{Z}^3$, satisfies the discrete wave equation

$$(2.1) \quad \frac{\partial^2}{\partial t^2} u_t(\gamma) = - \sum_{\gamma' \in \mathbb{Z}^3} \alpha(\gamma - \gamma') u_t(\gamma')$$

with initial data $(u_t|_{t=0}, v_t|_{t=0}) \in X = \ell_2 \times \ell_2$, where v_t denotes the velocity field, $v_t = \partial_t u_t$. The numbers $\alpha(\gamma - \gamma')$ are the elastic coupling constants between the sites γ and γ' . We assume that α is real and symmetric ($\alpha(-\gamma) = \alpha(\gamma)$). Clearly, system (2.1) can be written in a Hamiltonian form and the energy

$$(2.2) \quad H(u, v) = \frac{1}{2} \left(\sum_{\gamma \in \mathbb{Z}^3} |v(\gamma)|^2 + \sum_{\gamma \in \mathbb{Z}^3} \left[\sum_{\gamma' \in \mathbb{Z}^3} u(\gamma) \alpha(\gamma - \gamma') u(\gamma') \right] \right)$$

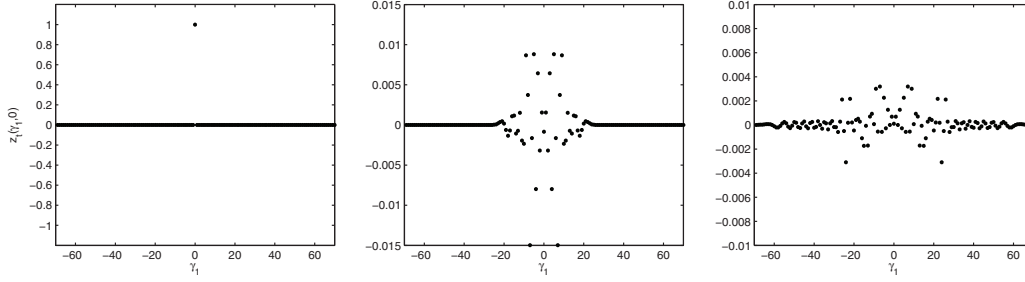


FIG. 2.1. Values of $u_t(\gamma)$ along the axis $\gamma_2 = \gamma_3 = 0$ for $t = 0, t = 0.1/\epsilon, t = 0.9/\epsilon$ with $\epsilon = \frac{1}{7} * 10^{-1}$. The evolution is given by (2.1) with the nearest neighbor elastic couplings (2.13), and the initial conditions are $u_{t=0}(0) = 1, u_{t=0}(\gamma) = 0$ for all $\gamma \neq 0$, and $\dot{u}_{t=0} \equiv 0$.

is constant along solutions. Depending on the initial conditions, the solutions of system (2.1) may develop large scale oscillations which carry a finite amount of energy; cf. Figure 2.1, where snapshots of u at several times are plotted.

Since system (2.1) is linear and invariant under discrete translations, we can write the solutions in a closed form using the Fourier transform.

DEFINITION 2.1. We define the Fourier transform $\ell_2(\mathbb{Z}^3) \rightarrow L^2(\mathbb{T}^3)$ by extending

$$(2.3) \quad (\mathcal{F}_{\gamma \rightarrow k}\psi)(k) = \hat{\psi}(k) = \sum_{\gamma \in \mathbb{Z}^3} e^{-2\pi i k \cdot \gamma} \psi(\gamma)$$

from ψ with finite support to all of $\ell_2(\mathbb{Z}^3)$. Here $\mathbb{T}^3 = \mathbb{R}^3/\mathbb{Z}^3$ denotes the unit 3-torus. The inverse transform is pointwise convergently defined by the integral

$$(2.4) \quad (\mathcal{F}_{k \rightarrow \gamma}\hat{\psi})(\gamma) = \int_{\mathbb{T}^3} dk e^{2\pi i k \cdot \gamma} \hat{\psi}(k) = \psi(\gamma),$$

where the measure dk is induced by the Lebesgue measure on the parameterization T^3 of \mathbb{T}^3 , where $T = (-\frac{1}{2}, \frac{1}{2}]$. In particular, $\int_{\mathbb{T}^3} dk = 1$.

From now on the notation T^3 will refer to the above explicit parameterization of the torus \mathbb{T}^3 . Occasionally, integration over $k \in \mathbb{T}^3$ will be applied to functions $f(k)$ which are not periodic in k . In these cases, using the parameterization $k \in T^3$ is implicitly understood. The parameterization mapping $\mathbb{T}^3 \rightarrow T^3$ will be denoted by $k \mapsto (k \bmod T^3)$ whenever the use of the parameterization needs to be explicitly stressed.

If one applies the Fourier transform to the Hamiltonian equations of motion determined by $H(u, v)$, one obtains a simpler system equivalent to (2.1):

$$(2.5) \quad \frac{\partial}{\partial t} \begin{pmatrix} \hat{u}_t(k) \\ \hat{v}_t(k) \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -\omega^2(k) & 0 \end{pmatrix} \begin{pmatrix} \hat{u}_t(k) \\ \hat{v}_t(k) \end{pmatrix}, \quad k \in \mathbb{T}^3.$$

The function $\omega : \mathbb{T}^3 \rightarrow \mathbb{R}$ is the dispersion relation and it is related to the Hamiltonian via the following formula:

$$(2.6) \quad \omega(k) = \sqrt{\hat{\alpha}(k)} = \sqrt{\sum_{\gamma \in \mathbb{Z}^3} \alpha(\gamma) \cos(2\pi \gamma \cdot k)},$$

where we have employed the assumption $\alpha(\gamma) = \alpha(-\gamma)$. Since α is real and satisfies the above symmetry property, we find that ω is also real and symmetric, i.e., $\omega(k) = \omega(-k)$.

A diagonalization of the matrix on the right-hand side of (2.5) motivates combining the real scalar fields u and v into the two complex fields $\psi_{\pm} = \psi_{\pm}[u, v] \in \ell_2(\mathbb{Z}^3, \mathbb{C})$ defined by the formula

$$(2.7) \quad \hat{\psi}_{\sigma}(k) = \frac{1}{\sqrt{2}}(\omega(k)\hat{u}(k) + i\sigma\hat{v}(k)), \quad \sigma \in \{\pm 1\}, \quad k \in \mathbb{T}^3.$$

For all $(u, v) \in X$, we clearly have $\hat{\psi}_{\sigma} \in L^2(\mathbb{T}^3)$, and thus $\psi_{\sigma} \in \ell_2(\mathbb{Z}^3)$. In addition, since $\omega(-k) = \omega(k)$, we also have $\psi_{-}(\gamma) = \psi_{+}(\gamma)$ for all γ . The transformation can always be inverted by applying

$$(2.8) \quad \hat{u} = \frac{1}{\omega\sqrt{2}}(\hat{\psi}_{+} + \hat{\psi}_{-}), \quad \hat{v} = -\frac{i}{\sqrt{2}}(\hat{\psi}_{+} - \hat{\psi}_{-}).$$

Given a solution $u_t, v_t = \dot{u}_t$ to (2.1), we then define the complex normal mode fields by the formula

$$(2.9) \quad \psi_{\sigma}(\gamma, t) = \psi_{\sigma}[u_t, v_t](\gamma), \quad \sigma \in \{\pm 1\}, \quad \gamma \in \mathbb{Z}^3, \quad t \in \mathbb{R}.$$

To see that these fields are indeed normal modes of the harmonic system, we apply (2.5) and find the evolution equations

$$(2.10) \quad \frac{\partial}{\partial t} \begin{pmatrix} \hat{\psi}_{+}(k, t) \\ \hat{\psi}_{-}(k, t) \end{pmatrix} = -i \begin{pmatrix} \omega(k) & 0 \\ 0 & -\omega(k) \end{pmatrix} \begin{pmatrix} \hat{\psi}_{+}(k, t) \\ \hat{\psi}_{-}(k, t) \end{pmatrix},$$

which are readily solved to yield for all $t \in \mathbb{R}, k \in \mathbb{R}^3$

$$(2.11) \quad \hat{\psi}_{\pm}(k, t) = e^{\mp i\omega(k)t} \hat{\psi}_{\pm}(k, 0).$$

These are exactly the two evolution equations corresponding to a “phonon” mode with a dispersion relation ω . The fields ψ_t mentioned in the introduction can now be identified with $\psi_{+}(\cdot, t)$ in the present notation.

After these reduction steps it is obvious that the dispersion relation ω fully determines the properties of the solutions. We will assume throughout this paper that ω is of the acoustic type in the following precise sense.

DEFINITION 2.2. *We call $\omega \in C(\mathbb{T}^3, [0, \infty))$ an acoustic dispersion relation if $\lambda = \omega^2$ satisfies the following:*

1. $\lambda \in C^{(3)}(\mathbb{T}^3, [0, \infty))$.
2. $\lambda(0) = 0$, and the Hessian of λ is invertible at 0.

A dispersion relation is called regular acoustic if it is acoustic and $\lambda(k) > 0$ for $k \neq 0$. The 3×3 -matrix A_0 is the Hessian of $\frac{1}{2}\lambda$ at $k = 0$ and $\omega_0(q) = \sqrt{q \cdot A_0 q}$.

Let us briefly motivate why these assumptions should be satisfied quite generally. For instance, if α is exponentially decaying, $\lambda = \hat{\alpha}$ will be analytic and thus satisfies the first requirement. The stability condition $\lambda \geq 0$ also needs to be satisfied everywhere; otherwise, the harmonic system has exponentially increasing solutions and is a poor model for oscillations in a crystal. If $\lambda(0) \neq 0$, the dispersion relation is not singular at $k = 0$. Thus in order to have a singular dispersion relation, the only real restriction for exponentially decaying stable interactions is the assumption about invertibility of the Hessian. Our analysis could likely be extended to cover also more degenerate instances of λ , although this would require some additional effort. By analogous reasoning, acoustic normal modes—in the above sense—are seen to appear commonly also in the vector models discussed in section 1.1, at least for exponentially

decaying interactions. However, in the vector models level crossings can lead to additional singularities in the dispersion relations, which would require a separate study. Finally, let us remark that exponential decay is not essential, since only sufficient regularity of the Hessian of λ will be required in the proof. For instance, if there are $C, \delta > 0$ such that $|\alpha(\gamma)| \leq C|\gamma|^{-6-\delta}$, then $\lambda = \hat{\alpha}$ satisfies item 1 in Definition 2.2. This permits us to deal also with situations where α is obtained via linearization around an equilibrium of a nonlinear system with potentials decaying slightly faster than a Lennard–Jones potential.

A prototype for the kind of dispersion relations considered here is the dispersion relation of the nearest neighbor square lattice,

$$(2.12) \quad \omega_{\text{nn}}(k) = \left[\sum_{\nu=1}^3 2(1 - \cos(2\pi k^\nu)) \right]^{\frac{1}{2}}.$$

This is clearly a regular acoustic dispersion relation, in the sense of Definition 2.2, and for it A_0 is proportional to a unit matrix and $\omega_0(q) = 2\pi|q|$. The corresponding elastic couplings are given by $\alpha_{\text{nn}}(\gamma' - \gamma) = -\Delta_{\gamma'\gamma}$, where Δ is the discrete Laplacian of the square lattice. Explicitly,

$$(2.13) \quad \alpha_{\text{nn}}(\gamma) = \begin{cases} 6 & \text{if } \gamma = 0, \\ -1 & \text{if } |\gamma| = 1, \\ 0, & \text{otherwise.} \end{cases}$$

To allow the creation of macroscopic oscillations we work with sequences of initial conditions that depend on the scaling parameter $\varepsilon > 0$ and consider the asymptotic behavior of the solutions as ε tends to 0.

2.1. Energy density and the lattice Wigner transform. From now on we will focus on analyzing asymptotic behavior of the fields ψ_\pm as ε tends to 0. As is carefully discussed in [17], generalizing the definitions of the energy density and of the Wigner transform to the discrete setting is not completely obvious. In an attempt to minimize unnecessary repetition of certain basic results related to Wigner transforms, we will resort here to the definitions used in [14], which will allow us to rely on the properties proven in Appendix B of that reference. However, we wish to keep in mind that this choice might not be optimal for all purposes, and we refer the interested reader to the discussion and references found in [16, 17, 24] for further possibilities.

We employ here the definition that for any state $(u, v) \in X$, its energy density, $e^\varepsilon = e^\varepsilon[u, v]$, scaled to a lattice spacing $\varepsilon > 0$, is the following tempered distribution defined via the complex fields $\psi_\sigma = \psi_\sigma[u, v]$ in (2.7):

$$(2.14) \quad e^\varepsilon(x) = \sum_{\gamma \in \mathbb{Z}^3} \delta(x - \varepsilon\gamma) \frac{1}{2} \sum_{\sigma=\pm 1} |\psi_\sigma(\gamma)|^2,$$

where δ denotes the Dirac delta-distribution. This is a manifestly positive distribution which is identifiable with a measure whose total mass equals the total energy,

$$(2.15) \quad \begin{aligned} \int dx e^\varepsilon[u, v](x) &= \sum_{\gamma \in \mathbb{Z}^3} \frac{1}{2} \sum_{\sigma=\pm 1} |\psi_\sigma(\gamma)|^2 = \frac{1}{2} \sum_{\sigma=\pm 1} \|\psi_\sigma\|^2 \\ &= \frac{1}{4} \sum_{\sigma=\pm 1} \int dk |\omega(k)\hat{u}(k) + i\sigma\hat{v}(k)|^2 = H(u, v) < \infty. \end{aligned}$$

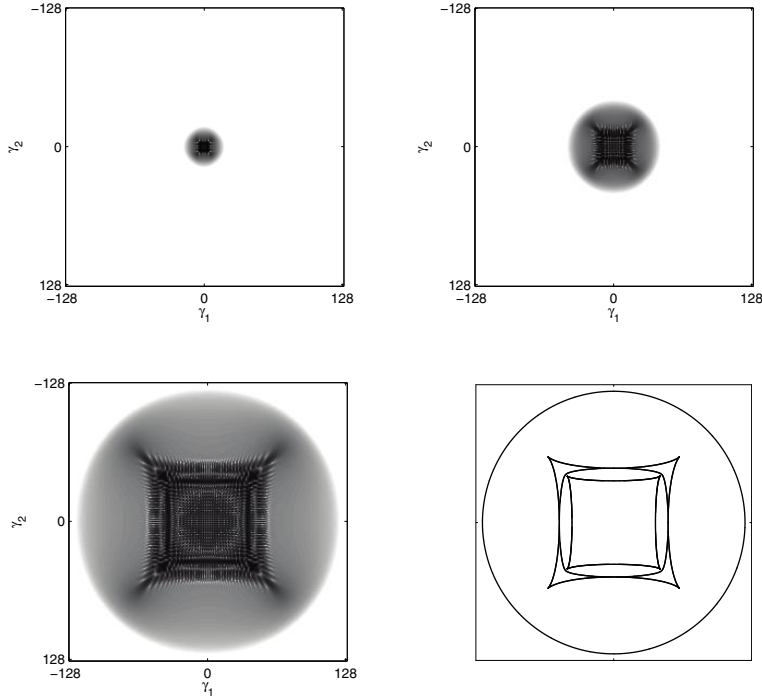


FIG. 2.2. First three panels: Snapshots of the energy density $|\psi_+(\gamma, t)|^2$ in the plane $\gamma_3 = 0$ at $t = 0.1/\varepsilon$, $t = 0.3/\varepsilon$, and $t = 0.95/\varepsilon$, $\varepsilon = \frac{1}{128}$, with initial data and elastic constants as in Figure 2.1. The plot is of the logarithm of the density, with all values less than a fixed cutoff shown white. Last panel: Plot of the restriction to the plane $x_3 = 0$ of the singular set of the energy of the solution to the transport equation (1.3), using a corresponding choice of macroscopic initial data. As the solution is scale-invariant, no explicit length scale has been denoted.

This justifies calling $e^\varepsilon[u, v]$ an energy density: it defines a distribution of the positive total energy between the lattice sites. The symmetry of ω implies that $|\psi_-(\gamma)| = |\psi_+(\gamma)|$, and thus we can also identify the energy density directly with the norm-density of $\psi_+[u, v]$,

$$(2.16) \quad e^\varepsilon(x) = \sum_{\gamma \in \mathbb{Z}^3} \delta(x - \varepsilon\gamma) |\psi_+(\gamma)|^2.$$

Then also for all $\varepsilon > 0$,

$$(2.17) \quad H(u, v) = \int dx e^\varepsilon[u, v](x) = \|\psi_+\|_{\ell^2(\mathbb{Z}^3)}^2 = \|\hat{\psi}_+\|_{L^2(\mathbb{T}^3)}^2 = \|\hat{\psi}_-\|_{L^2(\mathbb{T}^3)}^2,$$

which is independent of t for any $(u, v) = (u_t, v_t)$.

Let us now concentrate on the case when a solution u_t to (2.1) has been given, and define for any t the distribution $e^\varepsilon(x, t)$ by setting

$$(2.18) \quad e^\varepsilon(x, t) = e^\varepsilon[u_{t/\varepsilon}, v_{t/\varepsilon}](x), \quad x \in \mathbb{R}^3, t \in \mathbb{R}.$$

We have given an example of the time evolution of the energy density in Figure 2.2. The last panel in the figure contains the most obvious features which are implied

by the corresponding macroscopic evolution equation: the points of discontinuity of the solution. The macroscopic initial data is given by $d\mu_0(x, k) = \delta(x) \frac{1}{2} |\omega(k)|^2 dx dk$, which has no concentration at $k = 0$. Thus only (1.3) is relevant. It is readily solved to yield as the energy density (defined in this case simply by integrating μ_t over the k -variable)

$$(2.19) \quad e(x, t) = \int_{\mathbb{T}^3} dk \delta(x - t \frac{1}{2\pi} \nabla \omega(k)) = t^{-3} \int_{\mathbb{T}^3} dk \delta(\frac{x}{t} - \frac{1}{2\pi} \nabla \omega(k)).$$

Evaluation of such integrals has been considered, for instance, in section 6.4 of [17]. $|\nabla \omega(k)|$ has its maximum near the point of discontinuity of the gradient, at $k = 0$. This defines the outer circle outside which the solution must be zero. Inside the circle, the solution has a finite density, apart from points which correspond to values of k for which the Hessian of ω is not invertible. We have computed the positions of such points using *Mathematica* and plotted the result in the last panel in Figure 2.2. For a reader interested in the details of the computation, we point out that considering the case $x_3 = 0$ simplifies the problem, as it implies that either $k_3 = 0$ or $k_3 = \frac{1}{2}$.

We are interested in the limiting behavior of $e^\varepsilon(x, t)$ as ε tends to 0. Since the velocity of the waves with wave vector k depends on k , it is necessary to work with an object that encodes the density of waves with wave vector $k \in \mathbb{T}^3$ at $x \in \mathbb{R}^3$. This job is conveniently done by the Wigner transform. In order to avoid certain technical difficulties, we are going to define our Wigner transform only in the sense of distributions, i.e., via a duality principle.

First we introduce the space of Schwartz functions.

DEFINITION 2.3. *Let $\mathcal{S}_d = \mathcal{S}(\mathbb{R}^d)$ denote the Schwartz space and $\|\cdot\|_{\mathcal{S}_d, N}$ the corresponding N th Schwartz norm. Explicitly, with α denoting an arbitrary multi-index and with $\langle x \rangle = \sqrt{1 + x^2}$, then*

$$(2.20) \quad \|f\|_{\mathcal{S}_d, N} = \sup_{x \in \mathbb{R}^d} \max_{|\alpha| \leq N} |\langle x \rangle^N \partial^\alpha f(x)|.$$

We also employ the shorthand notation $\mathcal{S} = \mathcal{S}_3$.

To extract the relevant weak limits as ε tends to 0, we have to specify a space of suitable test functions. We need to track the evolution of three different kinds of lattice vibrations (short, medium, and long wavelength), all done compatibly with the periodicity in the Fourier variable. This leads to the following, somewhat involved notion of multiscale test functions.

DEFINITION 2.4. *We call a test function $a \in C^\infty(\mathbb{R}^3 \times \mathbb{T}^3 \times \mathbb{R}^3)$ admissible if it satisfies the following properties:*

1. $\sup_{k, q, |\alpha| \leq N} \|\partial_{k, q}^\alpha a(\cdot, k, q)\|_{\mathcal{S}, N} < \infty$ for all $N \geq 0$.
2. $q \mapsto a(x, k, q)$ is constant for all $x \in \mathbb{R}^3$ and k such that $\max_i |k_i \bmod T| \geq \frac{1}{4}$.
3. There is a function $b \in C^\infty(\mathbb{R}^3 \times \mathbb{T}^3 \times S^2)$ such that for any $N \geq 0$

$$(2.21) \quad \sup_{|q| \geq R, k \in \mathbb{T}^3} \|a(\cdot, k, q) - b(\cdot, k, \frac{q}{|q|})\|_{\mathcal{S}, N} \rightarrow 0, \quad \text{when } R \rightarrow \infty.$$

The first and third conditions are the most important and can be summarized as follows: we assume the test functions to be Schwartz in x and smooth with bounded derivatives in k and q and to have a radial limit b in q which is approached uniformly in any of the Schwartz norms. The above requirements are not minimal. The second condition is only needed in order to guarantee that $k \mapsto a(x, k, (k \bmod T^3)/\varepsilon)$ is always

smooth on \mathbb{T}^3 . Also, taking arbitrarily large N in the last step is not necessary; most likely $N = d + 3 = 6$ would suffice.

Having the notion of admissible test functions at our disposal we can define the central object of this paper: the multiscale Wigner transform.

DEFINITION 2.5. *Let $\psi \in \ell_2(\mathbb{Z}^3)$. We define the multiscale lattice Wigner transform $W^\varepsilon[\psi]$ at scale $\varepsilon > 0$ by*

$$(2.22) \quad \langle a, W^\varepsilon[\psi] \rangle = \int_{\mathbb{R}^3} dp \int_{T^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} \overline{\hat{\psi}(k - \varepsilon \frac{p}{2})} \hat{\psi}(k + \varepsilon \frac{p}{2}),$$

where a is an admissible test function and $\hat{a} = \mathcal{F}_{x \rightarrow p} a$, i.e.,

$$(2.23) \quad \hat{a}(p, k, q) = \int_{\mathbb{R}^3} dx e^{-2\pi i p \cdot x} a(x, k, q).$$

The L^2 -Wigner transform $W_{\text{cont}}^{(\varepsilon)}[\phi]$ of a function $\phi \in L^2(\mathbb{R}^3)$ at the scale $\varepsilon > 0$ is given by the distribution

$$(2.24) \quad b \mapsto \langle b, W_{\text{cont}}^{(\varepsilon)}[\phi] \rangle = \int_{\mathbb{R}^3 \times \mathbb{R}^3} dp dq \overline{\hat{b}(p, q)} \overline{\hat{\phi}(q - \varepsilon \frac{p}{2})} \hat{\phi}(q + \varepsilon \frac{p}{2})$$

for all $b \in \mathcal{S}(\mathbb{R}^3, C^\infty(\mathbb{R}^3))$, and with $\hat{b} = \mathcal{F}_{x \rightarrow p} b$.

The Wigner transform $W_{\text{cont}}^{(\varepsilon)}$ is a rescaled version of the standard Wigner transform in L^2 . The test-function space $\mathcal{S}(\mathbb{R}^3, C^\infty(\mathbb{R}^3))$ used above in its definition is obtained via the family of seminorms $p_N(b) = \sup_{|\alpha|, |x|, |q| \leq N} |\langle x \rangle^N \partial_{x,q}^\alpha b(x, q)|$ with $b \in C^\infty(\mathbb{R}^3 \times \mathbb{R}^3)$. This is a Fréchet space, and $W_{\text{cont}}^{(\varepsilon)}[\phi]$ is a continuous functional on it for any $\phi \in L^2(\mathbb{R}^3)$, as the following estimate reveals:

$$(2.25) \quad |\langle b, W_{\text{cont}}^{(\varepsilon)}[\phi] \rangle| \leq \sup_{p,q} |\langle p \rangle^4 \hat{b}(p, q)| \|\phi\|^2 \int_{\mathbb{R}^3} dp \langle p \rangle^{-4}.$$

Although \mathcal{S}_6 is not dense in this test-function space, it is nevertheless enough to know how $W_{\text{cont}}^{(\varepsilon)}$ acts on it. More precisely, if $W_i = W_{\text{cont}}^{(\varepsilon)}[\phi_i]$, $i = 1, 2$, and $\langle b, W_1 \rangle = \langle b, W_2 \rangle$ for all $b \in \mathcal{S}_6$, then $W_1 = W_2$. This follows straightforwardly from an estimate similar to (2.25) using smooth cutoff functions to cut out the infinity of the q -variable. In addition, we will also need the property that if $b(x, q) = f(x)$, with $f \in \mathcal{S}_3$, then

$$(2.26) \quad \langle b, W_{\text{cont}}^{(\varepsilon)}[\phi] \rangle = \int_{\mathbb{R}^3} dx \overline{f(x)} |\phi(x)|^2.$$

That is, $\int_{\mathbb{R}^3} dq W_{\text{cont}}^{(\varepsilon)}[\phi](x, q) = |\phi(x)|^2$. For a more careful analysis of the properties of the standard Wigner transform, see [12].

In [14] the Wigner transform of a lattice state was defined as a distribution $W_{\text{latt}}^\varepsilon \in \mathcal{S}'(\mathbb{R}^3 \times \mathbb{T}^3)$. The above definition is simply a refinement of this definition: formally for any ψ

$$(2.27) \quad W^\varepsilon(x, k, q) = \delta(q - \frac{k}{\varepsilon}) W_{\text{latt}}^\varepsilon(x, k), \quad x \in \mathbb{R}^3, k \in T^3, q \in \mathbb{R}^3.$$

This follows immediately from equation (B.6) of [14], after one realizes that if a is an admissible test function, then $(x, k) \mapsto a(x, k, (k \bmod T^3)/\varepsilon)$ belongs to $\mathcal{S}(\mathbb{R}^3 \times \mathbb{T}^3)$ for any $\varepsilon > 0$. This identification immediately allows us to use the results in [14] and

to prove that many of the basic properties of the usual Wigner transform carry over to the multiscale Wigner transform. Particularly important for us is the following relation.

PROPOSITION 2.6. *For any $f \in \mathcal{S}(\mathbb{R}^3)$, the test function $a_f(x, k, q) = f(x)$ is admissible, and for all $\varepsilon > 0$ and for all $(u, v) \in X$,*

$$(2.28) \quad \langle f, e^\varepsilon \rangle = \langle a_f, W^\varepsilon[\psi] \rangle,$$

where $e^\varepsilon = e^\varepsilon[u, v]$ and $\psi = \psi_+[u, v]$.

Proof. By inspection, we find that a is a well-defined test function and $\psi \in \ell_2$. Then, by the above-mentioned relation, $\langle a_f, W^\varepsilon[\psi] \rangle = \langle J_f, W_{\text{latt}}^\varepsilon[\psi] \rangle$ with $J_f(x, k) = f(x)$. On the other hand, it follows directly from the definition of $W_{\text{latt}}^\varepsilon$ (equation (B.2) in [14]) that

$$(2.29) \quad \begin{aligned} \langle J_f, W_{\text{latt}}^\varepsilon[\psi] \rangle &= \sum_{\gamma, \gamma' \in \mathbb{Z}^3} \psi(\gamma) \overline{\psi(\gamma')} \int_{\mathbb{T}^3} dk e^{2\pi i k \cdot (\gamma' - \gamma)} \overline{f(\varepsilon(\gamma' + \gamma)/2)} \\ &= \sum_{\gamma \in \mathbb{Z}^3} \overline{f(\varepsilon\gamma)} |\psi(\gamma)|^2 = \langle f, e^\varepsilon \rangle. \end{aligned}$$

This proves (2.28). \square

The above proposition can be formally summarized by the formula

$$(2.30) \quad e^\varepsilon(x) = \int_{\mathbb{T}^3} dk \int_{\mathbb{R}^3} dq W^\varepsilon[\psi](x, k, q),$$

which implies also (in the sense of choosing any suitable test-function sequence approaching 1)

$$(2.31) \quad \int_{\mathbb{R}^3 \times \mathbb{T}^3 \times \mathbb{R}^3} dx dk dq W^\varepsilon[\psi](x, k, q) = \|\psi\|_{\ell_2(\mathbb{Z}^3)}^2.$$

As noted earlier, analogous results hold for the L^2 -Wigner transform.

Several definitions have been used in the literature to study homogenization limits of lattice systems. We follow [14] mainly for convenience. The definition in [14] is based on the ‘‘Weyl quantization rule’’ of symbols, and a similar definition of the Wigner transform as distributions with a ‘‘classical quantization rule’’ was proposed earlier in [16]. Other choices include using Husimi transforms or considering $L^2(\mathbb{T}^3) = L^2(T^3)$ as a subspace of $L^2(\mathbb{R}^3)$ and then relying on the standard Wigner transform [17]. Finally, one can consider also interpolations of the fields between lattice sites and then using the Wigner series defined in [12]. The relations between the various definitions are discussed in [11, 16, 17, 24]. As is apparent from the results presented in these references, the various 1-microlocal definitions tend to lead to the same limit Wigner measures. This is not the case for the above 2-microlocal measures, since at least the information about the quantization rule will be carried over to the homogenization limits. However, a more systematic study would be required to settle the issue.

3. Main results. The macroscopic evolution is obtained by sending ε to 0. Our objective is to characterize the asymptotic behavior of the Wigner transform $W^\varepsilon[\psi_{t/\varepsilon}^\varepsilon]$ of the solution $\psi_{t/\varepsilon}^\varepsilon$ of (2.10). The limit strongly depends on the dispersion relation ω . We will consider here regular acoustic dispersion relations, keeping in

mind that, for instance, in the case of a scalar field and nearest neighbor interactions in \mathbb{Z}^3 the dispersion relation ω is given by (2.12), which is regular acoustic with $\omega_0(q) = 2\pi|q|$. The main achievement of this paper as compared to [17] is that complicated assumptions concerning the concentrations of the Wigner transform W^ε in wavenumber space as ε tends to 0 are no longer needed. The only remaining requirements are boundedness and tightness of the sequence of initial excitations.

ASSUMPTION 3.1. *We consider a sequence of values $\varepsilon > 0$ such that $\varepsilon \rightarrow 0$. For each ε in the sequence we assume that there is given an initial data vector $\psi_0^\varepsilon \in \ell_2(\mathbb{Z}^3)$ such that*

1. $\sup_\varepsilon \|\psi_0^\varepsilon\| < \infty$;
2. *the sequence ψ_0^ε is tight on the scale ε^{-1} :*

$$(3.1) \quad \lim_{R \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} \sum_{|\gamma| > R/\varepsilon} |\psi_0^\varepsilon(\gamma)|^2 = 0.$$

After these preparations we are in a position to state our result. The main point is that if ω is a regular acoustic dispersion relation, the asymptotic behavior of the energy density is characterized by precisely three different objects: a Wigner transform of a weak limit (macroscopic waves), an H-measure (short macroscopic waves), and a Wigner measure (microscopic waves). No assumptions concerning energy concentrations except those stated in Assumption 3.1 are required.

THEOREM 3.2. *Let $\psi_0^\varepsilon \in \ell_2(\mathbb{Z}^3)$ be a sequence which satisfies Assumption 3.1. Let ω be a regular acoustic dispersion relation and define $\psi_t^\varepsilon \in \ell_2(\mathbb{Z}^3)$ for all $t \in \mathbb{R}$ by the formula*

$$(3.2) \quad \hat{\psi}_t^\varepsilon(k) = e^{-it\omega(k)} \hat{\psi}_0^\varepsilon(k).$$

Let also $\mathbb{T}_*^3 = \mathbb{T}^3 \setminus \{0\}$. Then there are positive, bounded Radon measures μ_0, μ_0^H on $\mathbb{R}^3 \times \mathbb{T}_*^3$ and $\mathbb{R}^3 \times S^2$, respectively, a function $\phi_0 \in L^2(\mathbb{R}^3)$, and a subsequence (not relabeled) such that for all admissible test functions a and $t \in \mathbb{R}$,

$$(3.3) \quad \begin{aligned} \lim_{\varepsilon \rightarrow 0} \langle a, W^\varepsilon[\psi_{t/\varepsilon}^\varepsilon] \rangle &= \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} d\mu_t(x, k) \overline{b(x, k, \frac{k}{|k|})} \\ &+ \int_{\mathbb{R}^3 \times S^2} d\mu_t^H(x, q) \overline{b(x, 0, q)} + \langle a_0, W_{cont}^{(1)}[\phi_t] \rangle, \end{aligned}$$

where $b(x, k, q) = \lim_{R \rightarrow \infty} a(x, k, Rq)$ for $|q| = 1$, $a_0(x, q) = a(x, 0, q)$, and ϕ_t, μ_t and μ_t^H are determined for any $f \in C(\mathbb{R}^3 \times \mathbb{T}_*^3)$ and $g \in C(\mathbb{R}^3 \times S^2)$ by

$$(3.4) \quad \hat{\phi}_t(q) := e^{-it\omega_0(q)} \hat{\phi}_0(q),$$

$$(3.5) \quad \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} f(x, k) d\mu_t(x, k) := \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} f(x + t\frac{1}{2\pi}\nabla\omega(k), k) d\mu_0(x, k),$$

$$(3.6) \quad \int_{\mathbb{R}^3 \times S^2} g(x, q) d\mu_t^H(x, q) := \int_{\mathbb{R}^3 \times S^2} g(x + t\frac{1}{2\pi}\nabla\omega_0(q), q) d\mu_0^H(x, q).$$

Moreover, for all t the energy equality holds:

$$(3.7) \quad \lim_{\varepsilon \rightarrow 0} \|\psi_0^\varepsilon\|^2 = \mu_t(\mathbb{R}^3 \times \mathbb{T}_*^3) + \mu_t^H(\mathbb{R}^3 \times S^2) + \|\phi_t\|_{L^2(\mathbb{R}^3)}^2.$$

Remark 3.3. It is immediate from the definition that ϕ, μ , and μ^H are weak solutions of the set of decoupled linear transport equations (1.3)–(1.5).

The proof of Theorem 3.2 also shows that the subsequences which are extracted in the statement of the theorem can be characterized by a simple condition. In particular, the initial state of the wave equation, ϕ_0 , is determined as the weak- $L^2(\mathbb{R}^3)$ limit of the sequence of the functions (ϕ_0^ε) with Fourier transforms

$$(3.8) \quad \hat{\phi}_0^\varepsilon(q) = \begin{cases} \varepsilon^{\frac{3}{2}} \hat{\psi}_0^\varepsilon(\varepsilon q) & \text{if } \|q\|_\infty \leq \frac{1}{2\varepsilon}, \\ 0, & \text{otherwise.} \end{cases}$$

The exact characterization is contained in the following corollary, whose proof will be given in section 5.

COROLLARY 3.4. *Let (ψ_0^ε) be a sequence which satisfies Assumption 3.1. Suppose that ϕ_0^ε converges weakly to ϕ_0 , and that $\lim_{\varepsilon \rightarrow 0} \langle a, W^\varepsilon[\psi_0^\varepsilon] \rangle$ exists for every admissible test function a . Then there are unique positive, bounded Radon measures μ_0, μ_0^H on $\mathbb{R}^3 \times \mathbb{T}_*^3$ and $\mathbb{R}^3 \times S^2$, respectively, such that for every admissible test function a*

$$(3.9) \quad \begin{aligned} \lim_{\varepsilon \rightarrow 0} \langle a, W^\varepsilon[\psi_0^\varepsilon] \rangle &= \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} d\mu_0(x, k) \overline{b(x, k, \frac{k}{|k|})} \\ &+ \int_{\mathbb{R}^3 \times S^2} d\mu_0^H(x, q) \overline{b(x, 0, q)} + \langle a_0, W_{cont}^{(1)}[\phi_0] \rangle, \end{aligned}$$

where a_0 and b are defined as in Theorem 3.2. In addition, then (3.3) holds for all $t \in \mathbb{R}$ along the original sequence ε with the initial macroscopic data determined by the triplet (μ_0, μ_0^H, ϕ_0) .

Remark 3.5. According to (2.27) for $a \in \mathcal{S}(\mathbb{R}^3 \times \mathbb{T}^3)$ the lattice Wigner transform $W_{latt}^\varepsilon[\psi_0^\varepsilon]$ satisfies

$$(3.10) \quad \langle a, W_{latt}^\varepsilon[\psi_0^\varepsilon] \rangle = \langle a, W^\varepsilon[\psi_0^\varepsilon] \rangle.$$

It is well known (see, e.g., [14]) that $W_{latt}^\varepsilon[\psi_0^\varepsilon]$ has a limit μ_{latt} ,

$$(3.11) \quad \lim_{\varepsilon \rightarrow 0} \langle a, W_{latt}^\varepsilon[\psi_0^\varepsilon] \rangle = \int_{\mathbb{R}^3 \times \mathbb{T}^3} d\mu_{latt}(x, k) \overline{a(x, k)},$$

which is a positive Radon measure on $\mathbb{R}^3 \times \mathbb{T}^3$. Evaluating (3.9) for $a \in \mathcal{S}(\mathbb{R}^3 \times \mathbb{T}^3)$ shows that the “standard Wigner measure” μ_{latt} can be expressed in terms of the two-scale Wigner measure according to

$$(3.12) \quad d\mu_{latt}(x, k) = d\mu_{latt}(x, k) \mathbb{1}(k \neq 0) + \int_{q \in S^2} d\mu_0^H(x, q) \otimes \delta(k) dk + |\phi_0(x)|^2 dx \otimes \delta(k) dk,$$

where dx denotes Lebesgue measure on \mathbb{R}^3 , δ the Dirac δ -function, and $\mathbb{1}$ the characteristic function.

Remark 3.6. Our proof of the existence of two-scale Wigner measures is similar to the strategy proposed by Fermanian Kammerer in [4, 5]. For a general result on existence of two-scale Wigner measures concentrating on hypersurfaces, see [6]. (See [6, Theorem 1.4] for a simple formulation in one dimension and [6, Theorem 1.6] for the general case.) Concerning the existence part the only novelty in our result is thus the use of the lattice Wigner transform.

However, in our approach we obtain the existence of the measures together with their time evolution. While the strategy for obtaining transport equations for two-scale Wigner measures outlined in section 6 of [7] is not directly applicable in our

case due to the singular nature of the dispersion relation ω , it could certainly be adapted. The existence part and the transport part of the result could thus be separated. However, our proof is technically simpler and we do not rely on any advanced pseudodifferential calculus.

4. Proof of Theorem 3.2. Let a be an admissible test function and consider a fixed $t \in \mathbb{R}$ when we need to inspect the $\varepsilon \rightarrow 0$ limit of

$$(4.1) \quad \langle a, W^\varepsilon[\psi_{t/\varepsilon}^\varepsilon] \rangle = \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} e^{-i\frac{t}{\varepsilon}(\omega(k+\varepsilon\frac{p}{2})-\omega(k-\varepsilon\frac{p}{2}))} \overline{\hat{\psi}_0^\varepsilon(k-\varepsilon\frac{p}{2})} \hat{\psi}_0^\varepsilon(k+\varepsilon\frac{p}{2}).$$

First we identify the function ϕ_0 which contains the contributions of the long-wave excitations. Let $\hat{\phi}_0^\varepsilon$ be defined by (3.8). Since $\limsup_{\varepsilon \rightarrow 0} \|\hat{\phi}_0^\varepsilon\|_{L^2} = \limsup_{\varepsilon \rightarrow 0} \|\psi_0^\varepsilon\|_{\ell_2}$ is bounded by Assumption 3.1 there exist a subsequence and a function $\phi_0 \in L^2(\mathbb{R}^3)$ such that $\hat{\phi}_0^\varepsilon$ converges weakly to ϕ_0 in $L^2(\mathbb{R}^3)$. We then define ϕ_t by (3.4).

Using a localization function χ , which will be specified later, we now split the integral I on the right-hand side of (4.1) into three parts and an error term so that the contributions of short-, medium-, and long-wave excitation can be analyzed separately. We define

$$\begin{aligned} I_{>} &= \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} e^{-i\frac{t}{\varepsilon}(\omega(k+\varepsilon\frac{p}{2})-\omega(k-\varepsilon\frac{p}{2}))} \overline{\hat{\psi}_0^\varepsilon(k-\varepsilon\frac{p}{2})} \hat{\psi}_0^\varepsilon(k+\varepsilon\frac{p}{2})(1-\chi), \\ I_{<}^H &= \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} e^{-i\frac{t}{\varepsilon}(\omega(k+\varepsilon\frac{p}{2})-\omega(k-\varepsilon\frac{p}{2}))} \overline{(\hat{\psi}_0^\varepsilon(k-\varepsilon\frac{p}{2}) - \varepsilon^{-\frac{3}{2}}\hat{\phi}_0(\frac{k}{\varepsilon} - \frac{p}{2}))} \\ &\quad \times (\hat{\psi}_0^\varepsilon(k+\varepsilon\frac{p}{2}) - \varepsilon^{-\frac{3}{2}}\hat{\phi}_0(\frac{k}{\varepsilon} + \frac{p}{2}))\chi, \\ I_{<}^{\text{wv}} &= \varepsilon^{-3} \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} e^{-i\frac{t}{\varepsilon}(\omega(k+\varepsilon\frac{p}{2})-\omega(k-\varepsilon\frac{p}{2}))} \overline{\hat{\phi}_0(\frac{k}{\varepsilon} - \frac{p}{2})} \hat{\phi}_0(\frac{k}{\varepsilon} + \frac{p}{2})\chi, \\ R &= I - I_{>} - I_{<}^H - I_{<}^{\text{wv}}. \end{aligned}$$

The definition of R implies that

$$(4.2) \quad \langle a, W^\varepsilon[\psi_{t/\varepsilon}^\varepsilon] \rangle = I_{>} + I_{<}^H + I_{<}^{\text{wv}} + R.$$

To localize the oscillations in Fourier space we need smooth cutoff functions.

DEFINITION 4.1. Let $f \in C^\infty(\mathbb{R}, [0, 1])$ denote a fixed function which is symmetric, $f(-x) = f(x)$, strictly monotonically decreasing on $[1, 2]$, and

$$(4.3) \quad f(x) = \begin{cases} 1 & \text{if } |x| \leq 1, \\ 0 & \text{if } |x| \geq 2. \end{cases}$$

We further define $\varphi \in C^\infty(\mathbb{R}^3, [0, 1])$ by $\varphi(k) = f(|k|)$.

Let $0 < \rho \leq \frac{1}{4}$ be arbitrary and set $\chi = \chi(k, p; \varepsilon, \rho) = \varphi(\frac{k_+}{\rho})\varphi(\frac{k_-}{\rho})$, where $k_\pm = k \pm \varepsilon\frac{p}{2}$. We will continue to use this shorthand notation in the following under the tacit assumption that k_\pm is always really a function of both k and εp . Let us now also point out that all four terms in the decomposition (4.2) depend on ε and ρ , via χ , even though we have not denoted this dependence explicitly.

The first term containing $1 - \chi$ in (4.2) is zero if $|k_\pm| \leq \rho$, while the remainder is zero if $|k_+|$ or $|k_-| \geq 2\rho$. Thus the chosen decomposition splits the integration over k and p into “large,” “intermediate,” and “small” wavenumbers. We will demonstrate

that the following convergences hold: there is a sequence of ρ , a subsequence of ε , and measures μ_t, μ_t^H satisfying the statements made in the theorem, such that

$$(4.4) \quad \lim_{\rho \rightarrow 0} \lim_{\varepsilon \rightarrow 0} I_{>} = \int_{\mathbb{R}^3 \times \mathbb{T}_+^3} d\mu_t(x, k) \overline{b(x, k, \frac{k}{|k|})},$$

$$(4.5) \quad \lim_{\rho \rightarrow 0} \lim_{\varepsilon \rightarrow 0} I_{<}^H = \int_{\mathbb{R}^3 \times S^2} d\mu_t^H(x, q) \overline{b(x, 0, q)},$$

$$(4.6) \quad \lim_{\rho \rightarrow 0} \limsup_{\varepsilon \rightarrow 0} \left| I_{<}^{wv} - \langle a_0, W_{\text{cont}}^{(1)}[\phi_t] \rangle \right| = 0,$$

$$(4.7) \quad \lim_{\rho \rightarrow 0} \limsup_{\varepsilon \rightarrow 0} |R| = 0.$$

Clearly, (4.4)–(4.7) then imply (3.3).

Large wavenumbers. We split $I_{>}$ further into two parts using

$$(4.8) \quad \begin{aligned} 1 - \varphi\left(\frac{k_-}{\rho}\right) \varphi\left(\frac{k_+}{\rho}\right) &= 1 - \varphi\left(\frac{k}{\rho}\right)^2 + \left(\varphi\left(\frac{k}{\rho}\right) - \varphi\left(\frac{k_+}{\rho}\right)\right) \varphi\left(\frac{k}{\rho}\right) \\ &+ \left(\varphi\left(\frac{k}{\rho}\right) - \varphi\left(\frac{k_-}{\rho}\right)\right) \varphi\left(\frac{k_+}{\rho}\right). \end{aligned}$$

Let $h_\rho(k) = 1 - \varphi(\frac{k}{\rho})^2$, which is a smooth function. The integral then becomes $I_{>} = I_{>}^1 + R_1$, where

$$(4.9) \quad I_{>}^1 = \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{a}(p, k, \frac{k}{\varepsilon})} h_\rho(k) e^{-i\frac{t}{\varepsilon}(\omega(k+\varepsilon\frac{p}{2})-\omega(k-\varepsilon\frac{p}{2}))} \overline{\hat{\psi}_0^\varepsilon(k-\varepsilon\frac{p}{2})} \hat{\psi}_0^\varepsilon(k+\varepsilon\frac{p}{2}).$$

The remainder R_1 can be estimated using $|\varphi| \leq 1$ and

$$(4.10) \quad \varphi\left(\frac{k_\pm}{\rho}\right) - \varphi\left(\frac{k}{\rho}\right) = \pm \frac{\varepsilon}{2\rho} \int_0^1 ds p \cdot \nabla \varphi\left(\frac{1}{\rho}(k \pm s\frac{\varepsilon}{2}p)\right),$$

which yield the bound, with a universal constant C ,

$$(4.11) \quad |R_1| \leq C \frac{\varepsilon}{\rho} \sup_{k, q} \|a(\cdot, k, q)\|_{S, d+2} \|\hat{\psi}_0^\varepsilon\|^2 \|\nabla \varphi\|_\infty.$$

Therefore, there is a constant c' such that $|R_1| \leq c'\varepsilon/\rho$ and thus $R_1 \rightarrow 0$ when $\varepsilon \rightarrow 0$ for all ρ .

We then consider $I_{>}^1$. The presence of h_ρ guarantees that the integrand is zero unless $|k| \geq \rho$. Thus we can change $\hat{a}(p, k, k/\varepsilon)$ to $\hat{b}(p, k, k/|k|)$ in the integrand with an error R_2 bounded by

$$(4.12) \quad |R_2| \leq C \sup_{k, |q| \geq \rho/\varepsilon} \|a(\cdot, k, q) - b(\cdot, k, \frac{q}{|q|})\|_{S, d+1} \|\hat{\psi}_0^\varepsilon\|^2,$$

where C is a universal constant. Therefore, the assumptions imply that $R_2 \rightarrow 0$ when $\varepsilon \rightarrow 0$ for all ρ . On the other hand, for $|k| \geq \rho$ and $|p| < \rho/\varepsilon$, inequality (A.2) derived in the appendix implies

$$(4.13) \quad \left| \frac{1}{\varepsilon} (\omega(k + \varepsilon\frac{p}{2}) - \omega(k - \varepsilon\frac{p}{2})) - p \cdot \nabla \omega(k) \right| \leq C_3 \varepsilon \frac{|p|^2}{|k|} \leq C_3 \frac{\varepsilon}{\rho} |p|^2.$$

Therefore, using the estimate $|e^{ix} - e^{iy}| \leq \min(|x - y|, 2)$, valid for all $x, y \in \mathbb{R}$, we find that we can further change the t -dependent exponential in the integrand to $e^{-itp \cdot \nabla \omega(k)}$ with an error R_3 satisfying

$$(4.14) \quad |R_3| \leq C \sup_{k,q} \|b(\cdot, k, q)\|_{\mathcal{S}, d+3} \left(\frac{\varepsilon}{\rho} |t| + \int_{|p| \geq \rho/\varepsilon} dp \langle p \rangle^{-d-3} \right)$$

for some constant C . Therefore, also $\lim_{\varepsilon \rightarrow 0} R_3 = 0$ for all ρ . In summary, $I_{>}^1 = I_{>}^2 + R_2 + R_3$, where R_2, R_3 are negligible, and

$$(4.15) \quad I_{>}^2 = I_{>}^2(\varepsilon, \rho) = \int_{\mathbb{R}^3} dp \int_{\mathbb{T}^3} dk \overline{\hat{b}(p, k, \frac{k}{|k|})} h_\rho(k) e^{-itp \cdot \nabla \omega(k)} \overline{\hat{\psi}_0^\varepsilon(k - \varepsilon \frac{p}{2})} \hat{\psi}_0^\varepsilon(k + \varepsilon \frac{p}{2}).$$

Let us for a moment consider the lattice Wigner transform $W_{\text{latt}}^\varepsilon$ of ψ_0^ε as defined in [14]. As pointed out after Definition 2.5, then for any test function $f \in \mathcal{S}(\mathbb{R}^3 \times \mathbb{T}^3)$, we get an admissible test function by the formula $a_f(x, k, q) = f(x, k)$ and $\langle f, W_{\text{latt}}^\varepsilon \rangle = \langle a_f, W^\varepsilon[\psi_0^\varepsilon] \rangle$. Since ψ_0^ε is a norm-bounded sequence, the sequence $W_{\text{latt}}^\varepsilon$ is weak-* bounded, and thus there are $W_{\text{latt}}^0 \in \mathcal{S}'(\mathbb{R}^3 \times \mathbb{T}^3)$ and a subsequence along which $W_{\text{latt}}^\varepsilon \xrightarrow{*} W_{\text{latt}}^0$.

Since the sequence ψ_0^ε is by assumption also tight on the scale ε^{-1} , we can then apply Theorems B.4 and B.5 of [14] and conclude that W_{latt}^0 is given by a positive, bounded Radon measure μ on $\mathbb{R}^3 \times \mathbb{T}^3$ such that for all *continuous* functions $f \in C(\mathbb{T}^3)$ and $p \in \mathbb{R}^3$,

$$(4.16) \quad \lim_{\varepsilon \rightarrow 0} \int_{\mathbb{T}^3} dk f(k) \overline{\hat{\psi}_0^\varepsilon(k - \varepsilon \frac{p}{2})} \hat{\psi}_0^\varepsilon(k + \varepsilon \frac{p}{2}) = \int_{\mathbb{R}^3 \times \mathbb{T}^3} d\mu(x, k) f(k) e^{-2\pi i p \cdot x}.$$

Because $\nabla \omega(k)$ and $\frac{k}{|k|}$ are continuous apart from $k = 0$, the function $k \mapsto h_\rho(k) \overline{\hat{b}(p, k, \frac{k}{|k|})} e^{-itp \cdot \nabla \omega(k)}$ is everywhere continuous for all $\rho > 0$ and $p \in \mathbb{R}^3$. Therefore, by the dominated convergence theorem, for all ρ we find

$$(4.17) \quad \begin{aligned} \lim_{\varepsilon \rightarrow 0} I_{>}^2 &= \int_{\mathbb{R}^3} dp \int_{\mathbb{R}^3 \times \mathbb{T}^3} d\mu(x, k) \overline{\hat{b}(p, k, \frac{k}{|k|})} h_\rho(k) e^{-2\pi i p \cdot (x + t \nabla \omega(k) (2\pi)^{-1})} \\ &= \int_{\mathbb{R}^3 \times \mathbb{T}^3} d\mu(x, k) h_\rho(k) \overline{b(x + t \frac{1}{2\pi} \nabla \omega(k), k, \frac{k}{|k|})}. \end{aligned}$$

When $\rho \rightarrow 0$, the integrand approaches pointwise $\overline{b(x + t \frac{1}{2\pi} \nabla \omega(k), k, \frac{k}{|k|})}$ apart from $k = 0$ when the limit is 0. Therefore, by the dominated convergence theorem

$$(4.18) \quad \begin{aligned} &\lim_{\rho \rightarrow 0} \int_{\mathbb{R}^3 \times \mathbb{T}^3} d\mu(x, k) h_\rho(k) \overline{b(x + t \frac{1}{2\pi} \nabla \omega(k), k, \frac{k}{|k|})} \\ &= \int_{\mathbb{R}^3 \times \mathbb{T}^3_*} d\mu(x, k) \overline{b(x + t \frac{1}{2\pi} \nabla \omega(k), k, \frac{k}{|k|})} = \int_{\mathbb{R}^3 \times \mathbb{T}^3_*} d\mu_t(x, k) \overline{b(x, k, \frac{k}{|k|})}, \end{aligned}$$

where we have defined the bounded, positive Radon measure μ_t using $\mu_0 = \mu|_{\mathbb{R}^3 \times \mathbb{T}^3_*}$ in the formula (3.5). We have shown that (4.4) holds.

Small wavenumbers and the remainder. After a change of variables $q = \frac{k}{\varepsilon}$, one obtains that

$$(4.19) \quad I_{<}^{\text{wv}} = \int_{\mathbb{R}^3} dp \int_{T^3/\varepsilon} dq \overline{\hat{a}(p, \varepsilon q, q)} e^{-i \frac{t}{\varepsilon} (\omega(\varepsilon q_+) - \omega(\varepsilon q_-))} \overline{\hat{\phi}_0(q_-)} \hat{\phi}_0(q_+) \varphi(\frac{\varepsilon q_+}{\rho}) \varphi(\frac{\varepsilon q_-}{\rho}),$$

where $q_{\pm} = q \pm \frac{\rho}{2}$. We can immediately replace the integration region for the q -integral by \mathbb{R}^3 . To see this, note that the integrand is zero, unless $|q \pm \frac{\rho}{2}| \leq \frac{2\rho}{\varepsilon}$ for both signs. Since $2q = q_+ + q_-$, this can happen only if also $|q| \leq \frac{2\rho}{\varepsilon} \leq \frac{1}{2\varepsilon}$, which implies that the integrand is zero if $\|q\|_{\infty} > \frac{1}{2\varepsilon}$.

However, if $|\varepsilon q| \leq 2\rho$, then

$$(4.20) \quad |\hat{a}(p, \varepsilon q, q) - \hat{a}(p, 0, q)| \leq \sup_k |\nabla_k \hat{a}(p, k, q)| 2\rho,$$

and thus we can replace in the integrand the function $\hat{a}(p, \varepsilon q, q)$ by $\hat{a}(p, 0, q)$, with an error R'_1 which is bounded by $C\rho$ with a constant C independent of ε and ρ . Thus we only need to consider the integral

$$(4.21) \quad I_{<}^1 = \int_{\mathbb{R}^3} dp \int_{\mathbb{R}^3} dq \overline{\hat{a}(p, 0, q)} e^{-i\frac{t}{\varepsilon}(\omega(\varepsilon q_+) - \omega(\varepsilon q_-))} \overline{F^{\varepsilon, \rho}(q_-)} F^{\varepsilon, \rho}(q_+),$$

where $F^{\varepsilon, \rho}(q) = \hat{\phi}_0(q) \varphi(\frac{\varepsilon q}{\rho})$. Next we apply the estimate (A.3) derived in the appendix, proving that if now $|p| \leq \frac{1}{2\varepsilon}$, then

$$(4.22) \quad \left| \frac{1}{\varepsilon}(\omega(\varepsilon q_+) - \omega(\varepsilon q_-)) - \omega_0(q_+) + \omega_0(q_-) \right| \leq C_4 \varepsilon |p| |q|.$$

Following the same argument as earlier, we can then conclude that the t -dependent exponential can be changed to $e^{-it(\omega_0(q_+) - \omega_0(q_-))}$, with an error R'_2 which satisfies the estimate

$$(4.23) \quad |R'_2| \leq C \sup_q \|a(\cdot, 0, q)\|_{S, d+2} \left(|t| \rho + \int_{|p| \geq (2\varepsilon)^{-1}} dp \langle p \rangle^{-d-2} \right).$$

Thus $\lim_{\rho \rightarrow 0} \limsup_{\varepsilon \rightarrow 0} |R'_2| = 0$ for all t . Finally, we need to change $F^{\varepsilon, \rho}$ to $\hat{\phi}_0$, with an error R'_3 which can be bounded by $C \|F^{\varepsilon, \rho} - \hat{\phi}_0\|$. Since the bound goes to zero when $\varepsilon \rightarrow 0$ and

$$(4.24) \quad I_{<}^{wv} = \int_{\mathbb{R}^3} dp \int_{\mathbb{R}^3} dq \overline{\hat{a}(p, 0, q)} e^{-it(\omega_0(q_+) - \omega_0(q_-))} \overline{\hat{\phi}_0(q_-)} \hat{\phi}_0(q_+) + R'_1 + R'_2 + R'_3,$$

equation (4.6) has been established.

Similar estimates can be employed to demonstrate the vanishing of the remainder, (4.7). From the definition of R we get

$$(4.25) \quad \begin{aligned} R &= \int_{\mathbb{R}^3} dp \int_{T^3/\varepsilon} dq \overline{\hat{a}(p, \varepsilon q, q)} e^{-i\frac{t}{\varepsilon}(\omega(\varepsilon q_+) - \omega(\varepsilon q_-))} \overline{(\hat{\phi}_0^\varepsilon(q_-) - \hat{\phi}_0(q_-))} \hat{\phi}_0(q_+) \varphi\left(\frac{\varepsilon q_+}{\rho}\right) \varphi\left(\frac{\varepsilon q_-}{\rho}\right) \\ &+ \int_{\mathbb{R}^3} dp \int_{T^3/\varepsilon} dq \overline{\hat{a}(p, \varepsilon q, q)} e^{-i\frac{t}{\varepsilon}(\omega(\varepsilon q_+) - \omega(\varepsilon q_-))} \overline{\hat{\phi}_0(q_-)} (\hat{\phi}_0^\varepsilon(q_+) - \hat{\phi}_0(q_+)) \varphi\left(\frac{\varepsilon q_+}{\rho}\right) \varphi\left(\frac{\varepsilon q_-}{\rho}\right). \end{aligned}$$

We then apply the above estimates to remove the ε -dependence from all other terms in the integrands, apart from the differences $\hat{\phi}_0^\varepsilon - \hat{\phi}_0$. The error has a bound which vanishes when $\rho \rightarrow 0$. We are then left with

$$(4.26) \quad \begin{aligned} &\int_{\mathbb{R}^3} d\eta \int_{\mathbb{R}^3} d\xi \overline{\hat{a}(\eta - \xi, 0, \frac{1}{2}(\eta + \xi))} e^{-it(\omega_0(\eta) - \omega_0(\xi))} \overline{(\hat{\phi}_0^\varepsilon(\xi) - \hat{\phi}_0(\xi))} \hat{\phi}_0(\eta) \\ &+ \int_{\mathbb{R}^3} d\eta \int_{\mathbb{R}^3} d\xi \overline{\hat{a}(\eta - \xi, 0, \frac{1}{2}(\eta + \xi))} e^{-it(\omega_0(\eta) - \omega_0(\xi))} \overline{\hat{\phi}_0(\xi)} (\hat{\phi}_0^\varepsilon(\eta) - \hat{\phi}_0(\eta)), \end{aligned}$$

which vanishes as $\varepsilon \rightarrow 0$, since $\hat{\phi}_0^\varepsilon$ converges weakly to $\hat{\phi}_0$. This establishes (4.7).

Intermediate wavenumbers. Changing coordinates to $q = \frac{k}{\varepsilon}$ yields

$$I_{<}^H = \int_{\mathbb{R}^3} dp \int_{T^3/\varepsilon} dq \overline{\hat{a}(p, \varepsilon q, q)} e^{-i\frac{t}{\varepsilon}(\omega(\varepsilon q_+) - \omega(\varepsilon q_-))} \overline{\hat{f}^{\varepsilon, \rho}(q_-)} \hat{f}^{\varepsilon, \rho}(q_+),$$

where $\hat{f}^{\varepsilon, \rho}(q) = \varphi(\frac{\varepsilon q}{\rho})(\hat{\phi}_0^\varepsilon(q) - \hat{\phi}_0(q))$. Let $M > 0$ be arbitrary. We split off the values $|p| > M$ from the integral defining $I_{<}^3$. The difference $R'_4 = R'_4(\varepsilon, \rho, M)$ can be bounded by $C \int_{|p|>M} dp \langle p \rangle^{-d-1}$ and thus $\lim_{M \rightarrow \infty} \sup_{\rho, \varepsilon} |R'_4| = 0$. We divide the remaining integral over $|p| \leq M$ further into two parts using the identity

$$(4.27) \quad 1 = \tilde{\varphi}\left(\frac{q_-}{2M}\right) \tilde{\varphi}\left(\frac{q_+}{2M}\right) + 1 - \tilde{\varphi}\left(\frac{q_-}{2M}\right) \tilde{\varphi}\left(\frac{q_+}{2M}\right),$$

where $\tilde{\varphi} = 1 - \varphi$. If $|q| \geq 5M$, then $|q_\pm| \geq 4M$ and the second part is zero. It can be checked by inspection that the sequence $(f^{\varepsilon, \rho})_\varepsilon$ is bounded and tight, and it has a weak limit zero. Of these properties only the tightness is nonobvious, but this can also be easily deduced from the formula

$$(4.28) \quad f^{\varepsilon, \rho}(x) = \rho^3 \varepsilon^{-3/2} \sum_{\gamma \in \mathbb{Z}^3} \psi_0^\varepsilon(\gamma) \hat{\varphi}(\rho(\gamma - \frac{1}{\varepsilon}x)) - \int_{\mathbb{R}^3} dy \phi_0(x + \frac{\varepsilon}{\rho}y) \hat{\varphi}(y).$$

Therefore, by Lemma A.2, $\lim_{\varepsilon \rightarrow 0} \|\hat{f}^{\varepsilon, \rho}\|_{L^2(B_{6M})} = 0$ for all M, ρ . This implies that the contribution of the second term, denoted by R'_5 , satisfies $\lim_{\varepsilon \rightarrow 0} |R'_5| = 0$ for all M, ρ .

We are thus left with

$$(4.29) \quad I_{<}^4 = \int_{|p| \leq M} dp \int_{\mathbb{R}^3} dq \overline{\hat{a}(p, 0, q)} e^{-it(\omega_0(q_+) - \omega_0(q_-))} \overline{\hat{g}^{\varepsilon, \rho, M}(q_-)} \hat{g}^{\varepsilon, \rho, M}(q_+),$$

where $\hat{g}^{\varepsilon, \rho, M}(q) = \tilde{\varphi}\left(\frac{q}{2M}\right) \hat{f}^{\varepsilon, \rho}(q)$ and the integrand can be nonzero only for $|q| \geq M$. We thus only need to consider $|q| \geq M$ and $|p| \leq M$. First we replace in the integrand $\hat{a}(p, 0, q)$ by $\hat{b}(p, 0, \hat{q})$, $\hat{q} = \frac{q}{|q|}$, with an error R'_6 which is bounded by

$$(4.30) \quad |R'_6| \leq C \sup_{k, |q| \geq M} \|a(\cdot, k, q) - b(\cdot, k, \frac{q}{|q|})\|_{S, d+1}$$

for some constant C . Then we change $e^{-it(\omega_0(q_+) - \omega_0(q_-))}$ to $e^{-itp \cdot \nabla \omega_0(\hat{q})}$ with an error R'_7 which can be estimated using the inequality (A.4) of the appendix. This proves that there is a constant C such that

$$(4.31) \quad |R'_7| \leq C \frac{|t|}{M} \sup_{k, q} \|b(\cdot, k, q)\|_{S, d+3}.$$

Therefore, $I_{<}^4 = I_{<}^5 + R'_6 + R'_7$, where

$$(4.32) \quad I_{<}^5 = \int_{|p| \leq M} dp \int_{\mathbb{R}^3} dq \overline{\hat{b}(p, 0, \hat{q})} e^{-itp \cdot \nabla \omega_0(\hat{q})} \overline{\hat{g}^{\varepsilon, \rho, M}(q_-)} \hat{g}^{\varepsilon, \rho, M}(q_+)$$

and $\lim_{M \rightarrow \infty} \sup_{\rho, \varepsilon} |R'_6 + R'_7| = 0$ for all t .

Let

$$(4.33) \quad b_t(x, \hat{q}) = b(x + t\frac{1}{2\pi} \nabla \omega_0(\hat{q}), 0, \hat{q}).$$

Then $b_t \in \mathcal{S}(\mathbb{R}^3 \times S^2)$, and it has an extension to a function $J_t \in \mathcal{S}_6$; i.e., there is J_t such that $J_t(x, q) = b_t(x, q)$ for all $|q| = 1$. On the other hand, then

$$(4.34) \quad I_{<}^5 = \int_{|p| \leq M} dp \int_{\mathbb{R}^3} dq \widehat{b}_t(p, \hat{q}) \overline{\widehat{g}^{\varepsilon, \rho, M}(q_-)} \widehat{g}^{\varepsilon, \rho, M}(q_+) = \langle J_t, \Lambda^{\varepsilon, \rho, M} \rangle,$$

where $\Lambda^{\varepsilon, \rho, M} \in \mathcal{S}'_6$ denotes the distribution

$$(4.35) \quad J \mapsto \langle J, \Lambda^{\varepsilon, \rho, M} \rangle = \int_{|p| \leq M} dp \int_{\mathbb{R}^3} dq \widehat{J}(p, \hat{q}) \overline{\widehat{g}^{\varepsilon, \rho, M}(q_-)} \widehat{g}^{\varepsilon, \rho, M}(q_+).$$

Clearly, each $\Lambda^{\varepsilon, \rho, M}$ has support in $\mathbb{R}^3 \times S^2$, and there is a constant C such that for all J, ε, ρ, M

$$(4.36) \quad |\langle J, \Lambda^{\varepsilon, \rho, M} \rangle| \leq C \|\widehat{g}^{\varepsilon, \rho, M}\|^2 \|J\|_{\mathcal{S}, d+1},$$

with $\|\widehat{g}^{\varepsilon, \rho, M}\|$ denoting the $L^2(\mathbb{R}^3)$ -norm. However, since we have $\|\widehat{g}^{\varepsilon, \rho, M}\| \leq \|\widehat{f}^{\varepsilon, \rho}\|$, where $\|\widehat{f}^{\varepsilon, \rho}\|$ is bounded in ε , the Banach–Alaoglu theorem implies that the family $(\Lambda^{\varepsilon, \rho, M})_{\varepsilon, \rho, M}$ belongs to a weak-* sequentially compact set. Therefore, for every ρ, M there is a subsequence of (ε) and $\Lambda^{\rho, M}$ such that $\Lambda^{\varepsilon, \rho, M} \overset{*}{\rightharpoonup} \Lambda^{\rho, M}$ along this subsequence. In addition, for every ρ there is Λ^ρ and a sequence of integers M such that $\Lambda^{\rho, M} \overset{*}{\rightharpoonup} \Lambda^\rho$ along this sequence. Finally, there are $\Lambda \in \mathcal{S}'_6$ and a sequence of integers $N \geq 4$ such that for $\rho = \frac{1}{N}$, $\Lambda^\rho \overset{*}{\rightharpoonup} \Lambda$.

All of the above distributions clearly must have support on $\mathbb{R}^3 \times S^2$. We will soon prove that, in addition, for all $f \in \mathcal{S}_6$ and ρ

$$(4.37) \quad \langle |f|^2, \Lambda^\rho \rangle \geq 0.$$

This implies that then also $\langle |f|^2, \Lambda \rangle \geq 0$. Therefore, by the Bochner–Schwartz theorem, there is a positive Radon measure μ^H on \mathbb{R}^6 such that for all test functions J , $\langle J, \Lambda \rangle = \int \mu^H(dx, dk) J(x, k)$. Since also μ^H must have support on $\mathbb{R}^3 \times S^2$, we can thus identify it with a positive Radon measure μ_0^H on $\mathbb{R}^3 \times S^2$. By considering test functions $J(x, k) = e^{-\delta^2 x^2}$ in the limit $\delta \rightarrow 0$, it is also clear that μ_0^H must be bounded. We then define the positive, bounded Radon measures μ_t^H , $t \in \mathbb{R}$, by the formula (3.6). It follows from the construction of μ_t^H that (4.5) holds along the above sequences ρ, ε .

The main missing ingredient is provided by the following lemma.

LEMMA 4.1. *For $p, q \in \mathbb{R}^3$, let us denote $q_\pm = q \pm \frac{p}{2}$, $\hat{q}_\pm = q_\pm / |q_\pm|$, and $\hat{q} = q / |q|$. There is a constant C such that for all $q, q' \in \mathbb{R}^3$ and $f \in \mathcal{S}_6$*

$$(4.38) \quad \left| \int_{\mathbb{R}^3} dx e^{-2\pi i p \cdot x} \overline{f(x, q')} f(x, q) \right| \leq C \langle p \rangle^{-d-1} \|f\|_{\mathcal{S}, d+1}^2.$$

If, in addition, $|q| \geq M$ and $|p| \leq M$, then also

$$(4.39) \quad \left| \mathcal{F}_{x \rightarrow p}(|f|^2)(p, \hat{q}) - \int_{\mathbb{R}^3} dx e^{-2\pi i p \cdot x} \overline{f(x, \hat{q}_+)} f(x, \hat{q}_-) \right| \leq \frac{1}{M} C \langle p \rangle^{-d-1} \|f\|_{\mathcal{S}, d+2}^2.$$

Before proving the lemma we demonstrate that it implies the inequality (4.37). Let $f \in \mathcal{S}_6$ be arbitrary. Define q_\pm, \hat{q}_\pm as in Lemma 4.1, except that here also let $\hat{0} = 0$, and let

$$(4.40) \quad I^{\varepsilon, \rho, M, f} = \int_{\mathbb{R}^3} dp \int_{\mathbb{R}^3} dq \left[\int_{\mathbb{R}^3} dx e^{2\pi i p \cdot x} f(x, \hat{q}_+) \overline{f(x, \hat{q}_-)} \right] \overline{\widehat{g}^{\varepsilon, \rho, M}(q_-)} \widehat{g}^{\varepsilon, \rho, M}(q_+).$$

Note that this integral is well-defined by (4.38). Then, by the estimate (4.39) and uniform boundedness of $g^{\varepsilon,\rho,M}$, there is a constant c such that

$$(4.41) \quad |\langle |f|^2, \Lambda^{\varepsilon,\rho,M} \rangle - I^{\varepsilon,\rho,M,f}| \leq c \|f\|_{\mathcal{S},d+2}^2 \left(\int_{|p| \geq M} dp \langle p \rangle^{-d-1} + \frac{1}{M} \right).$$

On the other hand, $I^{\varepsilon,\rho,M,f} \geq 0$ always. To see this, consider first the case when $g^{\varepsilon,\rho,R} \in \mathcal{S}$. Then $\hat{g}^{\varepsilon,\rho,R} \in \mathcal{S}$, and by changing variables from (q, p) to (q_+, q_-) in (4.40) and then using Fubini's theorem to reorder the integrals, we find that

$$(4.42) \quad I^{\varepsilon,\rho,M,f} = \int_{\mathbb{R}^3} dx |G(x)|^2, \quad \text{with } G(x) = \int_{\mathbb{R}^3} dq e^{2\pi i q \cdot x} f(x, \hat{q}) \hat{g}^{\varepsilon,\rho,R}(q) \in L^2.$$

Since $I^{\varepsilon,\rho,M,f}$ depends L^2 -continuously on $g^{\varepsilon,\rho,R}$ this implies that also for general $g^{\varepsilon,\rho,R} \in L^2$ we have $I^{\varepsilon,\rho,M,f} \geq 0$. Since the right-hand side of (4.41) vanishes if first $\varepsilon \rightarrow 0$ and then $M \rightarrow \infty$, we must thus also have $\langle |f|^2, \Lambda^\rho \rangle \geq 0$. This proves (4.37).

The only remaining task is to prove Lemma 4.1. Consider first (4.38). If $|p| \leq 1$, we have trivially a bound

$$(4.43) \quad \int_{\mathbb{R}^3} dx |f(x, q')| |f(x, q)| \leq \|f\|_{\mathcal{S},d+1}^2 \int_{\mathbb{R}^3} dx \langle x \rangle^{-2d-2} \leq c \|f\|_{\mathcal{S},d+1}^2.$$

If $|p| \geq 1$, we perform N partial integrations in the direction of p , i.e., in the direction $\hat{p} = \frac{p}{|p|}$, yielding

$$(4.44) \quad \int_{\mathbb{R}^3} dx e^{-2\pi i p \cdot x} \overline{f(x, q')} f(x, q) = \frac{1}{(2\pi i |p|)^N} \int_{\mathbb{R}^3} dx e^{-2\pi i p \cdot x} (\hat{p} \cdot \nabla)^N (\overline{f(x, q')} f(x, q)).$$

By the Leibniz rule,

$$(4.45) \quad (\hat{p} \cdot \nabla)^N (\overline{f(x, q')} f(x, q)) = \sum_{n=0}^N \binom{N}{n} (\hat{p} \cdot \nabla)^n \overline{f(x, q')} (\hat{p} \cdot \nabla)^{N-n} f(x, q),$$

which is bounded by $c'_N \langle x \rangle^{-N} \|f\|_{\mathcal{S},N}^2$. Choosing $N = d + 1$ then yields (4.38) for some constant. Adjusting the constant C so that the bound is true also for $|p| \leq 1$ proves that (4.38) is valid.

To prove (4.39), consider q, p as required in the second part of the lemma. Also let $q(s) = q + \frac{s}{2}p$. Then $|q(s)| \geq |q|/2 > 0$ for all $|s| \leq 1$, and $\hat{q}(s)$ is thus well-defined and smooth. Therefore, for any $g \in \mathcal{S}_6$ and $s_0 \in [-1, 1]$,

$$(4.46) \quad \begin{aligned} g(x, \hat{q}(s_0)) - g(x, \hat{q}) &= \int_0^{s_0} ds \frac{d}{ds} g(x, \hat{q}(s)) \\ &= \int_0^{s_0} ds \left(\frac{1}{2|q(s)|} p - \frac{p \cdot \hat{q}(s)}{2|q(s)|} \hat{q}(s) \right) \cdot \nabla_q g(x, q)|_{q=\hat{q}(s)}, \end{aligned}$$

implying

$$(4.47) \quad |g(x, \hat{q}_\pm) - g(x, \hat{q})| \leq 2 \frac{|p|}{|q|} \sup_{|q|=1} |\nabla_q g(x, q)|.$$

Also for all $g_1, g_2 \in \mathcal{S}_6$,

$$(4.48) \quad \begin{aligned} & \overline{g_1(x, \hat{q})g_2(x, \hat{q})} - \overline{g_1(x, \hat{q}_-)g_2(x, \hat{q}_+)} \\ &= \overline{(g_1(x, \hat{q}) - g_1(x, \hat{q}_-))g_2(x, \hat{q})} + \overline{g_1(x, \hat{q}_-)}(g_2(x, \hat{q}) - g_2(x, \hat{q}_+)). \end{aligned}$$

Following the steps made in the first part of the proof, and replacing the earlier estimates with the above more accurate ones when necessary, we can conclude that the constant C can be adjusted so that for these values of q, p (4.39) also holds. This completes the proof of (4.5).

Energy equality. The energy equality (3.7) follows by considering a sequence of test functions $a^\delta = e^{-\delta^2 x^2}$ and taking $\delta \rightarrow 0$. To see this, first note that the right-hand side of (3.7) is clearly independent of t , and thus it is enough to consider $t = 0$. Thanks to (2.28) and to the tightness of the sequence ψ_0^ε , we obtain for this particular test function that

$$\lim_{\delta \rightarrow 0} \lim_{\varepsilon \rightarrow 0} \langle a^\delta, W^\varepsilon[\psi_0^\varepsilon] \rangle = \lim_{\delta \rightarrow 0} \lim_{\varepsilon \rightarrow 0} \langle a^\delta, e^\varepsilon \rangle = \lim_{\varepsilon \rightarrow 0} \|\psi_0^\varepsilon\|^2.$$

Equation (3.3) implies that

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \langle a^\delta, W^\varepsilon[\psi_0^\varepsilon] \rangle &= \int_{\mathbb{R}^3} dx a^\delta(x) |\phi_0(x)|^2 + \int_{\mathbb{R}^3} \int_{\mathbb{T}^3} a^\delta(x) d\mu_0(x, k) \\ &\quad + \int_{\mathbb{R}^3} \int_{S^2} a^\delta(x) d\mu_0^H(x, k). \end{aligned}$$

Sending δ to 0 yields that

$$\lim_{\varepsilon \rightarrow 0} \|\hat{\psi}_0^\varepsilon\|^2 = \|\phi_0\|_{L^2(\mathbb{R}^3)}^2 + \mu_0(\mathbb{R}^3 \times \mathbb{T}^3) + \mu_0^H(\mathbb{R}^3 \times S^2),$$

and the energy equality has been established. This finishes the proof of Theorem 3.2.

5. Proof of Corollary 3.4. Let I_0 denote the original sequence of ε , and consider an arbitrary subsequence I of I_0 . Since $(\psi_0^\varepsilon)_{\varepsilon \in I}$ then also satisfies Assumption 3.1, we can conclude from Theorem 3.2 that for every I there is a subsequence I' such that (3.3) holds for all t with the initial conditions given by some triplet (μ_I, μ_I^H, ϕ_I) . From the construction of the subsequence in the proof of Theorem 3.2, we know that ϕ_I can be chosen as the weak limit of ϕ_0^ε along the subsequence I' . The first assumption thus implies that we can always choose $\phi_I = \phi_0$. Let us also denote $\mu_0 = \mu_{I_0}$ and $\mu_0^H = \mu_{I_0}^H$, and to prove the stated uniqueness, we will prove that $\mu_I = \mu_0$ and $\mu_I^H = \mu_0^H$ for all I .

Consider first an arbitrary $\tilde{a} \in C_c^\infty(\mathbb{R}^3 \times \mathbb{T}_*^3)$. Let $a(x, k, q) = \tilde{a}(x, k)$ for $k \neq 0$ and define $a(x, 0, q) = 0$. Then a is an admissible test function with $a_0 = 0 = b(k, 0, q)$, and for any subsequence I we thus obtain, using the second assumption,

$$(5.1) \quad \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} d\mu_I(x, k) \overline{\tilde{a}(x, k)} = \lim_{\varepsilon \in I} \langle a, W^\varepsilon[\psi_0^\varepsilon] \rangle = \int_{\mathbb{R}^3 \times \mathbb{T}_*^3} d\mu_0(x, k) \overline{\tilde{a}(x, k)}.$$

Such \tilde{a} are dense in $C_c(\mathbb{R}^3 \times \mathbb{T}_*^3)$, and thus $\mu_I = \mu_0$ on $\mathbb{R}^3 \times \mathbb{T}_*^3$.

Consider then an arbitrary $b \in C_c^\infty(\mathbb{R} \times S^2)$. Let φ be a smooth cutoff function as in Definition 4.1, and define $a(x, k, q) = b(x, q/|q|)(1 - \varphi(2q))$. Then a is an admissible

test function and indeed $\lim_{R \rightarrow \infty} a(x, k, Rq) = b(x, q)$ for all $|q| = 1$. By the already proven results we then find that for any subsequence I

$$\begin{aligned} \int_{\mathbb{R}^3 \times S^2} d\mu_I^H(x, q) \overline{b(x, q)} &= \lim_{\varepsilon \rightarrow 0} \langle a, W^\varepsilon[\psi_0^\varepsilon] \rangle - \int_{\mathbb{R}^3 \times \mathbb{T}_\pm^3} d\mu_0(x, k) \overline{b(x, \frac{k}{|k|})} - \langle a_0, W_{\text{cont}}^{(1)}[\phi_0] \rangle \\ (5.2) \qquad \qquad \qquad &= \int_{\mathbb{R}^3 \times S^2} d\mu_0^H(x, q) \overline{b(x, q)}. \end{aligned}$$

Therefore, also $\mu_I^H = \mu_0^H$, which concludes the uniqueness part of the proof of the corollary.

Finally, define for $t \neq 0$ the triplet (μ_t, μ_t^H, ϕ_t) using (μ_0, μ_0^H, ϕ_0) as the initial data. By the uniqueness proved above, for any subsequence I , (3.3) holds along the subsubsequence I' . As the right-hand side is thus independent of I , this proves that the limit also holds along the original sequence I_0 . This completes the proof of the corollary. \square

Appendix. The proof of Theorem 3.2 uses two simple lemmas, which are provided here. The first lemma summarizes several properties of regular acoustic dispersion relations.

LEMMA A.1. *Let ω be a regular acoustic dispersion relation, as in Definition 2.2, and let λ , A_0 , and ω_0 be related to ω as in the definition. Then the following assertions are true.*

1. $\omega \in C^{(3)}(\mathbb{T}^3 \setminus \{0\}, \mathbb{R})$.
2. $\nabla \lambda(0) = 0$ and $A_0 > 0$.
3. There are constants $C_1, C_2 > 0$ such that for all $\|k\|_\infty \leq \frac{3}{4}$,

$$(A.1) \qquad \qquad \omega(k) \geq C_1|k|, \qquad |\nabla \lambda(k)| \leq C_2|k|.$$

In addition, $\|\nabla \omega\|_\infty < \infty$.

4. There is C_3 such that if $\varepsilon > 0$, $p \in \mathbb{R}^3$, and $\|k\|_\infty \leq \frac{1}{2}$, with $|k| > \varepsilon|p|$, then

$$(A.2) \qquad |\omega(k + \frac{1}{2}\varepsilon p) - \omega(k - \frac{1}{2}\varepsilon p) - \varepsilon p \cdot \nabla \omega(k)| \leq C_3 \varepsilon^2 \frac{|p|^2}{|k|}.$$

5. There is C_4 such that if $\varepsilon > 0$ and $p, q \in \mathbb{R}^3$ with $|p|, |q| \leq \frac{1}{2}\varepsilon^{-1}$, then for $q_\pm = q \pm \frac{1}{2}p$,

$$(A.3) \qquad |\omega(\varepsilon q_+) - \omega(\varepsilon q_-) - \omega_0(\varepsilon q_+) + \omega_0(\varepsilon q_-)| \leq C_4 \varepsilon^2 |p| |q|.$$

6. There is C_5 such that if $p, q \in \mathbb{R}^3$ with $q \neq 0$ and $|p| \leq |q|$, then for $q_\pm = q \pm \frac{1}{2}p$,

$$(A.4) \qquad |\omega_0(q_+) - \omega_0(q_-) - p \cdot \nabla \omega_0(\frac{q}{|q|})| \leq C_5 \frac{|p|^2}{|q|}.$$

Proof. From now on, we consider ω , λ , A_0 , and ω_0 satisfying Definition 2.2. The first item is then obvious, and the second one follows from the assumptions, since 0 is a minimum. The second inequality in (A.1) follows by using item 2 and $\|D^2 \lambda\|_\infty < \infty$, when by Taylor expansion $|\nabla \lambda(k)| \leq C_2|k|$ for all $k \in \mathbb{R}$. To prove the first inequality, we first note that, by continuity, also $\|D^3 \lambda\|_\infty < \infty$. Thus there is c' such that for all $k \in \mathbb{R}^3$

$$(A.5) \qquad \qquad \qquad |\lambda(k) - \lambda_0(k)| \leq c'|k|^3,$$

where $\lambda_0(k) = \frac{1}{2}k \cdot A_0 k$. Since $A_0 > 0$, there is $c > 0$ such that $\lambda_0(k) \geq c^2 k^2$ for all k . Thus there is $\delta' > 0$ such that for all $|k| < \delta'$, we have $|1 - \lambda(k)/\lambda_0(k)| \leq \frac{3}{4}$, and for these k therefore also $\omega(k) = \sqrt{\lambda(k) - \lambda_0(k) + \lambda_0(k)} \geq \frac{1}{2}\sqrt{\lambda_0(k)} \geq \frac{c}{2}|k|$. Since λ has no zeroes in the complement set, the constant can then be adjusted so that (A.1) holds for all k with $\|k\|_\infty \leq \frac{3}{4}$.

We still need to prove the third property, boundedness of $\nabla\omega$. For later use, let us, more generally, consider a nonnegative $f \in C^{(2)}(\mathbb{R}^3)$, $q \in \mathbb{R}^3$, and k such that $f(k) \neq 0$. Then we have

$$(A.6) \quad (q \cdot \nabla)\sqrt{f(k)} = \frac{1}{2\sqrt{f(k)}}q \cdot \nabla f(k),$$

$$(A.7) \quad (q \cdot \nabla)^2\sqrt{f(k)} = -\frac{1}{4f(k)^{3/2}}(q \cdot \nabla f(k))^2 + \frac{1}{2\sqrt{f(k)}}(q \cdot \nabla)^2 f(k).$$

This implies that there is a constant C such that for all $q \in \mathbb{R}^3$ and $k \neq 0$, with $\|k\|_\infty \leq \frac{3}{4}$,

$$(A.8) \quad |q \cdot \nabla\omega(k)|, |q \cdot \nabla\omega_0(k)| \leq C|q|,$$

$$(A.9) \quad |(q \cdot \nabla)^2\omega(k)|, |(q \cdot \nabla)^2\omega_0(k)| \leq C\frac{|q|^2}{|k|}.$$

In particular, by periodicity therefore $\|\nabla\omega\|_\infty < \infty$.

To prove item 4, consider k, p, ε as in the claim, and define a function $f(s) = \omega(k_+) - \omega(k_-)$ with $k_\pm = k \pm s\frac{1}{2}p$ and $s \in [-\varepsilon, \varepsilon]$. Then $\|k_\pm\|_\infty \leq \frac{3}{2}\|k\|_\infty$ and $|k_\pm| \geq |k| - \varepsilon\frac{1}{2}|p| \geq \frac{1}{2}|k| > 0$, and thus f belongs to $C^{(3)}$. In particular, $f(0) = 0$, $f'(0) = p \cdot \nabla\omega(k)$, and

$$(A.10) \quad f''(s) = \frac{1}{2} \left[\frac{1}{\omega(k_+)} \left(\frac{p}{2} \cdot \nabla\right)^2 \lambda(k_+) - \frac{1}{\omega(k_-)} \left(\frac{p}{2} \cdot \nabla\right)^2 \lambda(k_-) - \frac{1}{2\omega(k_+)^3} \left(\frac{p}{2} \cdot \nabla\lambda(k_+)\right)^2 + \frac{1}{2\omega(k_-)^3} \left(\frac{p}{2} \cdot \nabla\lambda(k_-)\right)^2 \right].$$

Using item 3, we find that there is C such that this is uniformly bounded by $C|p|^2/|k|$. Then a Taylor expansion at the origin proves item 4.

Since for all $k \neq 0$

$$(A.11) \quad \omega(k) - \omega_0(k) = \sqrt{\lambda(k)} - \sqrt{\lambda_0(k)} = \frac{\lambda(k) - \lambda_0(k)}{\sqrt{\lambda(k)} + \sqrt{\lambda_0(k)}},$$

we then also have

$$(A.12) \quad q \cdot \nabla(\omega(k) - \omega_0(k)) = \frac{q \cdot \nabla\lambda(k) - q \cdot \nabla\lambda_0(k)}{\omega(k) + \omega_0(k)} - (\lambda(k) - \lambda_0(k)) \frac{q \cdot \nabla\omega(k) + q \cdot \nabla\omega_0(k)}{(\omega(k) + \omega_0(k))^2}.$$

By Taylor expansion at the origin, we find that there is C' such that for all $q \in \mathbb{R}^3$,

$$(A.13) \quad |\lambda(k) - \lambda_0(k)| \leq C'|k|^3, \quad |q \cdot \nabla\lambda(k) - q \cdot \nabla\lambda_0(k)| \leq C'|k|^2|q|,$$

$$(A.14) \quad |(q \cdot \nabla)^2\lambda(k) - (q \cdot \nabla)^2\lambda_0(k)| \leq C'|k||q|^2.$$

Thus there is also C such that for all $q \in \mathbb{R}^3$, and for $\|k\|_\infty \leq \frac{3}{4}$ with $k \neq 0$,

$$(A.15) \quad |q \cdot \nabla(\omega(k) - \omega_0(k))| \leq C|k||q|, \quad |(q \cdot \nabla)^2(\omega(k) - \omega_0(k))| \leq C|q|^2.$$

Let us then consider q, p, ε satisfying the assumptions made in the final item. Since then $\|\varepsilon q_{\pm}\|_{\infty} \leq \frac{3}{4}$, we can apply the previous estimates. If $p = 0$, then $q_- = q_+$ and the bound in (A.3) is trivially valid for any C_4 . Consider thus $p \neq 0$, and assume first that p is not proportional to q . Since then the line segment $[0, 1] \ni s \mapsto \varepsilon q + s\varepsilon p/2$ does not pass through the origin, the function $s \mapsto \omega(\varepsilon q + s\varepsilon p/2) - \omega_0(\varepsilon q + s\varepsilon p/2)$ is in $C^{(3)}([0, 1])$. We make a Taylor expansion of this function at $s = 0$, yielding

$$(A.16) \quad \omega(\varepsilon q + \varepsilon \frac{p}{2}) - \omega_0(\varepsilon q + \varepsilon \frac{p}{2}) = \omega(\varepsilon q) - \omega_0(\varepsilon q) + \varepsilon \frac{p}{2} \cdot (\nabla \omega(\varepsilon q) - \nabla \omega_0(\varepsilon q)) + R.$$

Here (A.15) implies $|R| \leq C\varepsilon^2|p|^2$, since

$$(A.17) \quad R = \int_0^1 ds (1-s) \left(\varepsilon \frac{p}{2} \cdot \nabla \right)^2 (\omega(\varepsilon q + s\varepsilon \frac{p}{2}) - \omega_0(\varepsilon q + s\varepsilon \frac{p}{2})).$$

This proves that (A.3) holds in this case for some constant C_4 . In the final case $p \propto q$, we choose a direction u orthogonal to q , and use the previous estimate with $p + \delta u$ instead of p for an arbitrary $0 < \delta \leq 1$. Since the left-hand side of (A.3) is continuous in δ , the bound must then hold also for $\delta = 0$, proving the validity of the estimate also in this case.

To prove (A.4), we use the fact that by assumption $|q_{\pm}| \geq \frac{1}{2}|q| > 0$, and thus denoting $\hat{q} = q/|q|$, we get

$$(A.18) \quad \omega_0(q_+) - \omega_0(q_-) = |q|(\omega_0(\hat{q} + \frac{p}{2|q|}) - \omega_0(\hat{q} - \frac{p}{2|q|})) = p \cdot \nabla \omega_0(\hat{q}) + R,$$

where $|R| \leq C|p|^2/|q|$. We have thus completed the proof of the lemma. \square

The second lemma recalls a well-known fact in Fourier analysis: the Fourier transform of a bounded and tight sequence of L^2 functions converges strongly on compact sets. For the convenience of the reader we give a proof.

LEMMA A.2. *Let $f^\varepsilon \in L^2(\mathbb{R}^d)$ be a bounded and tight sequence of functions such that $f^\varepsilon \rightharpoonup 0$ as $\varepsilon \rightarrow 0$. Then for every $\Omega \subset \mathbb{R}^d$ with a finite measure*

$$(A.19) \quad \lim_{\varepsilon \rightarrow 0} \|f^\varepsilon\|_{L^2(\Omega)} = 0.$$

Proof. Let $\beta > 0$ be arbitrary. By tightness of f^ε there is $R_\beta > 0$ such that

$$(A.20) \quad \limsup_{\varepsilon \rightarrow 0} \int_{|x| \geq R_\beta} dx |f^\varepsilon(x)|^2 \leq \beta^2.$$

Define now

$$(A.21) \quad f^{\varepsilon, \beta}(x) = \begin{cases} f^\varepsilon(x) & \text{if } |x| \leq R_\beta, \\ 0, & \text{otherwise.} \end{cases}$$

By the boundedness of $f^{\varepsilon, \beta}$ in $L^2(\mathbb{R}^d)$ there exist g^β and a subsequence (ε') such that $f^{\varepsilon', \beta} \rightharpoonup g^\beta$ in $L^2(\mathbb{R}^d)$ as $\varepsilon' \rightarrow 0$. Estimate (A.20) implies that $\limsup_{\varepsilon \rightarrow 0} \|f^\varepsilon - f^{\varepsilon, \beta}\|_{L^2(\mathbb{R}^d)}^2 \leq \beta^2$. Since also $f^\varepsilon \rightharpoonup 0$, we have

$$(A.22) \quad \|g^\beta\|^2 = \lim_{\varepsilon'} |\langle g^\beta, f^{\varepsilon', \beta} \rangle| \leq \|g^\beta\| \limsup_{\varepsilon \rightarrow 0} \|f^\varepsilon - f^{\varepsilon, \beta}\|$$

and, therefore, $\|g^\beta\|_{L^2(\mathbb{R}^d)} \leq \beta$.

Since g^β has support in a ball of radius R_β , also $g_1^\beta(x) = -2\pi i x g^\beta(x) \in L^2(\mathbb{R}^d)$. Therefore, $\hat{g}^\beta \in H^1(\mathbb{R}^d)$ and $\nabla \hat{g}^\beta = \hat{g}_1^\beta$. Similarly, also $\hat{f}^{\varepsilon, \beta} \in H^1(\mathbb{R}^d)$ and it is straightforward to check that $\nabla \hat{f}^{\varepsilon', \beta} \rightharpoonup \nabla \hat{g}^\beta$ in L^2 . Since Ω is assumed to be a set of finite measure, we then have $\lim_{\varepsilon' \rightarrow 0} \|\hat{f}^{\varepsilon', \beta} - \hat{g}^\beta\|_{L^2(\Omega)} = 0$ (for a proof, see, for instance, Theorem 8.6. in [13]). Collecting all the above estimates together proves that

$$\begin{aligned} & \limsup_{\varepsilon' \rightarrow 0} \|\hat{f}^{\varepsilon'}\|_{L^2(\Omega)} \\ & \leq \limsup_{\varepsilon' \rightarrow 0} \|\hat{f}^{\varepsilon'} - \hat{f}^{\varepsilon', \beta}\|_{L^2(\Omega)} + \limsup_{\varepsilon' \rightarrow 0} \|\hat{f}^{\varepsilon', \beta} - \hat{g}^\beta\|_{L^2(\Omega)} + \|\hat{g}^\beta\|_{L^2(\Omega)} \leq 2\beta. \end{aligned}$$

Since β can be arbitrarily small, there must be a subsequence (ε'') such that $\lim_{\varepsilon'' \rightarrow 0} \|\hat{f}^{\varepsilon''}\|_{L^2(\Omega)} = 0$.

Since the assumptions on the sequence f^ε are preserved for subsequences, we can consider an arbitrary subsequence and apply the above result to it. Then we can conclude that for every subsequence there is a subsubsequence along which (A.19) holds. This implies that the limit (A.19) actually holds also along the original sequence. \square

Acknowledgments. The authors would like to thank Clotilde Fermanian Kammerer, Alexander Mielke, and Herbert Spohn for several instructive discussions. We are also grateful to the anonymous referees for a careful reading of the manuscript and for suggesting several useful references.

REFERENCES

- [1] N. W. ASHCROFT AND N. D. MERMIN, *Solid State Physics*, Holt, Rinehart and Winston, New York, 1976.
- [2] F. BONETTO, J. L. LEBOWITZ, AND L. REY-BELLET, *Fourier's law: A challenge to theorists*, in *Mathematical Physics 2000*, A. Fokas, A. Grigoryan, T. Kibble, and B. Zegarlinski, eds., Imperial College Press, London, 2000, pp. 128–150.
- [3] T. V. DUDNIKOVA AND H. SPOHN, *Local stationarity for lattice dynamics in the harmonic approximation*, *Markov Process. Related Fields*, 12 (2006), pp. 645–678.
- [4] C. FERMANIAN KAMMERER, *Measure semi-classiques et équation de la chaleur*, Ph.D. thesis, Université Paris-Sud, Orsay, 1995.
- [5] C. FERMANIAN KAMMERER, *Mesures semi-classiques 2-microlocales*, *C. R. Acad. Sci. Paris*, 331 (2000), pp. 515–518.
- [6] C. FERMANIAN KAMMERER, *Propagation and absorption of concentration effects near shock hypersurfaces for the heat equation*, *Asymptot. Anal.*, 24 (2000), pp. 107–141.
- [7] C. FERMANIAN KAMMERER, *Two scale analysis of a bounded family in L^2 on a submanifold of the phase space*, *C. R. Acad. Sci. Paris*, 340 (2005), pp. 269–274.
- [8] N. B. FIROOZY, *Homogenization on Lattices: Small Parameter Limits, H-Measures, and Discrete Wigner Measures*, IMA Preprint 1177, University of Minnesota, Minneapolis, MN, 1993; available online from <http://www.ima.umn.edu/preprints/OCTOBER1993/1177.pdf>.
- [9] G. FRANCFORT AND P. GÉRARD, *The wave equation on a thin domain: Energy density and observability*, *J. Hyperbolic Differential Equations*, 1 (2004), pp. 351–366.
- [10] P. GÉRARD, *Compacité par compensation et régularité 2-microlocale*, in *Séminar Equation aux Dérivées Partielles*, Exp. No. VII, Ecole Polytechnique, Palaiseau, 1988–1989.
- [11] P. GÉRARD AND E. LEICHTNAM, *Ergodic properties of eigenfunctions for the Dirichlet problem*, *Duke Math. J.*, 71 (1993), pp. 559–607.
- [12] P. GÉRARD, P. A. MARKOWICH, N. J. MAUSER, AND F. PAUPAUD, *Homogenization limits and Wigner transforms*, *Commun. Pure Appl. Math.*, 50 (1997), pp. 323–379.
- [13] E. H. LIEB AND M. LOSS, *Analysis*, 2nd ed., AMS, Providence, RI, 2001.
- [14] J. LUKKARINEN AND H. SPOHN, *Kinetic limit for wave propagation in a random medium*, *Arch. Ration. Mech. Anal.*, 183 (2007), pp. 93–162.
- [15] F. MACIÀ, *Propagación y control de vibraciones en medios discretos y continuous*, Ph.D. thesis, Universidad Complutense de Madrid, Madrid, Spain, 2002.

- [16] F. MACIÀ, *Wigner measures in the discrete setting: High-frequency analysis of sampling and reconstruction operators*, SIAM J. Math. Anal., 36 (2004), pp. 347–383.
- [17] A. MIELKE, *Macroscopic behavior of microscopic oscillations in harmonic lattices via Wigner-Husimi transforms*, Arch. Ration. Mech. Anal., 181 (2006), pp. 401–448.
- [18] F. NIER, *A semi-classical picture of quantum scattering*, Ann. Sci. École Norm. Sup., 29 (1996), pp. 149–183.
- [19] L. RYZHIK, G. PAPANICOLAOU, AND J. B. KELLER, *Transport equations for elastic and other waves in random media*, Wave Motion, 24 (1996), pp. 327–370.
- [20] H. SPOHN, *The phonon Boltzmann equation, properties and link to weakly anharmonic lattice dynamics*, J. Statist. Phys., 124 (2006), pp. 1041–1104.
- [21] H. SPOHN, *Erratum on “The phonon Boltzmann equation, properties and link to weakly anharmonic lattice dynamics,”* J. Statist. Phys., 123 (2006), p. 707.
- [22] H. SPOHN, *Collisional invariants for the phonon Boltzmann equation*, J. Statist. Phys., 124 (2006), pp. 1131–1135.
- [23] L. TARTAR, *H-measures, a new approach for studying homogenisation, oscillations and concentration effects in partial differential equations*, Proc. Roy. Soc. Edinburgh Sect. A, 115 (1990), pp. 193–230.
- [24] S. TEUFEL AND G. PANATI, *Propagation of Wigner functions for the Schrödinger equation with a perturbed periodic potential*, in Multiscale Methods in Quantum Mechanics, Ph. Blanchard and G. Dell’Antonio, eds., Birkhäuser, Boston, 2004, pp. 207–220.
- [25] J. M. ZIMAN, *Electrons and Phonons: The Theory of Transport Phenomena in Solids*, Oxford University Press, London, 1967.

CONCENTRATION PHENOMENA IN A NONLOCAL QUASI-LINEAR PROBLEM MODELLING PHYTOPLANKTON I: EXISTENCE*

YIHONG DU[†] AND SZE-BI HSU[‡]

Abstract. We study the positive steady state of a quasi-linear reaction-diffusion system in one space dimension introduced by Klausmeier and Litchman for the modelling of the distributions of phytoplankton biomass and its nutrient. The system has nonlocal dependence on the biomass function, and it has a biomass-dependent drifting term describing the active movement of the biomass towards the location of the optimal growth condition. We obtain complete descriptions of the profile of the solutions when the coefficient of the drifting term is large, rigorously proving the numerically observed phenomenon of concentration of biomass for this model. Our theoretical results reveal four critical numbers for the model not observed before and offer several further insights into the problem being modelled. This is Part I of a two-part series, where we obtain nearly optimal existence and nonexistence results. The asymptotic profile of the solutions is studied in the separate Part II.

Key words. quasi-linear, nonlocal dependence, phytoplankton, concentration phenomenon, reaction-diffusion equation

AMS subject classifications. 35J55, 35J65, 92D25

DOI. 10.1137/07070663X

1. Introduction. In this paper, we study the problem

$$(1.1) \quad \begin{cases} -[d_1 u_x + \sigma c(x)u]_x = [g(x) - m]u, & 0 < x < 1, \\ -d_2 v_{xx} = -g(x)u, & 0 < x < 1, \\ d_1 u_x + \sigma c(x)u = 0, & x = 0, 1, \\ v_x(0) = 0, \quad v_x(1) = \beta[v_0 - v(1)], \end{cases}$$

where d_1, d_2, σ, m, v_0 , and β are positive constants,

$$g(x) = f(\min\{\alpha v(x), w(x)\}), \quad f(s) = \frac{rs}{K_I + s},$$

and

$$w(x) = w_0 \exp \left[-A_0 x - A \int_0^x u(s) ds \right],$$

with α, r, K_I, w_0, A , and A_0 positive constants. We note that the right-hand sides of the differential equations in (1.1) depend on the unknown functions u and v in a nonlocal manner. Moreover, the positive function $c(x)$ is determined by u and v in a rather unconventional way to be explained below. We are interested in positive

*Received by the editors October 29, 2007; accepted for publication (in revised form) June 24, 2008; published electronically November 5, 2008. This research was partially supported by the Australian Research Council, the NCTS, and the National Science Council of Taiwan.

<http://www.siam.org/journals/sima/40-4/70663.html>

[†]Department of Mathematics, School of Science and Technology, University of New England, Armidale, NSW2351, Australia (ydu@turing.une.edu.au), and Department of Mathematics, Qufu Normal University, People's Republic of China.

[‡]Department of Mathematics, National Tsing-Hua University, Hsinchu, Taiwan 300, Republic of China (sbhsu@math.nthu.edu.tw).

solutions of (1.1), namely $u > 0$ and $v > 0$ in $[0, 1]$. From (1.1) it is easy to see that for any such solution, v is an increasing function. Clearly w is a decreasing function. The function $c(x)$ is defined by

$$c(x) = \frac{x - x_0}{\delta + |x - x_0|},$$

where $\delta > 0$ is a small constant and $x_0 \in [0, 1]$ is the intersection point of the functions $\alpha u(x)$ and $w(x)$ whenever such an intersection occurs in $[0, 1]$; if $\alpha u(x)$ and $w(x)$ do not intersect, then $x_0 = 0$ if $\alpha v > w$ on $[0, 1]$, and $x_0 = 1$ if $\alpha v < w$ on $[0, 1]$. In other words, x_0 is given by the following description:

$$\min\{\alpha v(x), w(x)\} = \alpha v(x) \quad \forall x \in [0, x_0]; \quad \min\{\alpha v(x), w(x)\} = w(x) \quad \forall x \in (x_0, 1].$$

This unconventional dependence of c on the unknown solution (u, v) makes (1.1) a very special quasi-linear problem.

Such a system arises in the mathematical modelling of phytoplankton in a one-dimensional water column, where $u(x)$ represents the distribution of phytoplankton biomass, $v(x)$ stands for the distribution of nutrient, and x denotes the depth in the water column, with $x = 0$ at the surface and $x = 1$ at the bottom. The term $\sigma c(x)$ is used to describe the active movement of the biomass towards the spatial location with the optimal growth condition. Klausmeier and Litchman [KL] propose using this model to study the concentration phenomenon widely observed for phytoplankton in lakes and oceans. Their numerical analysis in [KL] demonstrates that for large σ , the biomass function $u(x)$ concentrates at a certain level $x = x_*$ while the nutrient function $v(x)$ is close to a piecewise linear function. They then treat u as a constant multiple of the δ -function concentrating at x_* and propose a game theoretical model to determine the location of x_* .

In this paper, we rigorously prove the existence of such a concentration phenomenon and obtain accurate formulas for the determination of x_* and the total biomass. Our theoretical results offer several further insights into the model besides those obtained through numerical analysis in [KL]; for example, we show the existence of four critical values $v_{**} < v_* < v^* < v^{**}$ for v_0 (the nutrient level at the sediment), such that

- (i) $x_* = 0$ when $v_0 \geq v^*$, $x_* \in (0, 1)$ when $v_0 \in (v_*, v^*)$, and $x_* = 1$ when $v_0 \leq v_*$;
- (ii) the total biomass increases with v_0 in the range $v_{**} < v_0 < v^{**}$, but it stays constant for $v_0 \geq v^{**}$ or $v_0 \leq v_{**}$ (and with v_0 above a certain level so that the biomass can survive).

It turns out that the game theoretical model of [KL] is a simplified version of our equations governing x_* and the total phytoplankton biomass for the case $v_* \leq v_0 \leq v^*$. A more detailed description of these results is given in the introduction of Part II, with their biological interpretations given in section 4 there.

To explain this model more precisely, we start by a brief description of the background and motivation of this research. Phytoplankton, the generic name of microorganisms living in lakes and oceans, is the basis of the aquatic food chain. Its importance for the proper functioning of the aquatic ecosystem has long been recognized, and its behavior has been widely studied. The distribution of phytoplankton in lakes and oceans is highly heterogeneous. To better understand this property of the phytoplankton, various mathematical models have been proposed and numerically analyzed; see, for example, [EATSH, KL, PT, PTHS, HTKS]. However, little

rigorous mathematical analysis is available. In [YN], an ordinary differential equation model for the vertical distributions of phytoplankton is theoretically analyzed; see also [IT, BFH, BFHK] for earlier related research. It is our hope that the current paper may induce further rigorous mathematical research in this direction and that the techniques developed here may find more applications.

We now describe the model in more detail. In poorly mixed water columns, it has been observed that algae can be heterogeneously distributed, with thin layers of biomass on the surface, at depth, or on the sediment surface; examples for each of these cases can be found in [KL]. To model these phenomena, [KL] proposes a reaction-diffusion-taxis model of phytoplankton, nutrients, and light, based on the principle of light and nutrient competition. They use the following system to describe the distribution of phytoplankton in a one-dimensional water column, with depth represented by $0 \leq z \leq z_b$; $z = 0$ at the surface and $z = z_b$ at the bottom:

$$(1.2) \quad \begin{cases} b_t = D_b b_{zz} + [\nu(g_z^0)b]_z + [g^0 - m]b, & 0 < z < z_b, t > 0, \\ R_t = D_R R_{zz} - bg^0/Y + \epsilon mb/Y, & 0 < z < z_b, t > 0, \\ D_b b_z + \nu(g_z^0)b = 0, & z = 0, z_b, t > 0, \\ R_z(t, 0) = 0, R_z(t, z_b) = h[R_{in} - R(t, z_b)], & t > 0, \\ I(t, z) = I_{in} \exp\left(-\int_0^z [ab(t, s) + a_{bg}]ds\right), & 0 < z < z_b, t > 0, \end{cases}$$

where $a, a_{bg}, h, D_b, D_R, m, I_{in}, R_{in}$, and Y are positive constants, $\epsilon \in [0, 1)$, $\nu(s)$ is an odd decreasing function that approaches $\nu_{max} > 0$ as $s \rightarrow -\infty$, and

$$g^0(t, z) = \min\{f_I(I(t, z)), f_R(R(t, z))\},$$

with

$$f_I(s) = r \frac{s}{s + K_I}, \quad f_R(s) = r \frac{s}{s + K_R}, \quad r, K_I, K_R > 0.$$

In (1.2), $b(t, z)$ denotes the distribution of the phytoplankton biomass, $R(t, z)$ represents the nutrient distribution, and $I(t, z)$ stands for the distribution of light. The constant I_{in} is the light distribution at the surface, and, by the Lambert–Beer law, light at depth z is given by

$$I(t, z) = I_{in} \exp\left(-\int_0^z [ab(t, s) + a_{bg}]ds\right),$$

where a and a_{bg} are, respectively, the phytoplankton and background attenuation coefficients. In this model, it is assumed that the change in phytoplankton biomass at depth z results from three processes: growth, loss, and movement. The functions $f_I(I)$ and $f_R(R)$ are the phytoplankton growth rate when only one of the resources I and R is limited (the other being regarded as sufficient). By Liebig’s law of the minimum for essential resources, the gross phytoplankton growth rate is given by $g^0(t, z) = \min\{f_I(I(t, z)), f_R(R(t, z))\}$. Biomass is lost at density-dependent rate m , representing respiration, death, and grazing. D_b is the passive diffusion rate of the biomass, while $[\nu(g_z^0)b]_z$ describes active movement of the biomass towards a spatial location (i.e., depth) with a better growth condition. The no-flux boundary condition for b means that no phytoplankton enters or leaves the water column at $z = 0$ and $z = z_b$. The equation for R is based on the assumption that nutrients in

the water column are mixed with eddy diffusion with diffusion coefficient D_R and are consumed by phytoplankton at the rate $-bg^0/Y$, and the term $\epsilon mb/Y$ means that ϵ proportion of the nutrients in dead phytoplankton is immediately recycled. Here Y describes the yield of phytoplankton biomass per unit nutrient consumed. The boundary condition for R means that nutrients do not leave or enter the top of the water column but are supplied at the bottom, with nutrients in the sediments fixed at constant concentration R_{in} , which diffuse across the sediment-water interface at a rate proportional to the concentration difference across the interface; the parameter h describes the permeability of the interface.

In [KL], taking $\epsilon = 0$ and $\nu(s) = \nu_0(s) := -\nu_{max}\text{sgn}(s)$ (where $\text{sgn}(s)$ is the sign function, which equals 1, -1 , or 0 according to whether $s > 0$, $s < 0$, or $s = 0$), the equilibrium distributions of b , R , and I are calculated numerically for various parameter values (see Table 1 and Figure 1 in [KL]). The numerical simulation in [KL] shows that as ν_{max} increases, the biomass distribution concentrates at a certain depth $z = z^*$. Further, based on intuition and formal analysis, a game theoretical approach is proposed in [KL], which can be used to calculate z^* . Though the connection between (1.2) and the simplified game theoretical approach is not rigorously established, the predictions deduced from the game theoretical model agree well with the numerical results based on (1.2); see details in [KL].

In this paper, we theoretically analyze the equilibrium solutions of (1.2). So $b = b(z)$, $R = R(z)$, and $I = I(z)$. Naturally, only positive solutions are of interest to us.

As in [KL], we assume that $\epsilon = 0$. We denote $f(s) = rs/(s+K_I)$ and $\alpha = K_I/K_R$. Then

$$f_I(s) = f(s), \quad f_R(s) = f(\alpha s).$$

Since $f'(s) > 0$ we find that

$$g^0(z) = \min\{f(I(z)), f(\alpha R(z))\} = f(\min\{I(z), \alpha R(z)\}).$$

Clearly $I(z)$ is a decreasing function. Since $D_R R'' = bg^0/Y > 0$ and $R'(0) = 0$, we find that $R'(z) > 0$ for $z \in (0, z_b]$. Therefore $R(z)$ is increasing, and we can always find a unique $z_0 \in [0, z_b]$ such that

$$g^0(z) = f(\alpha R(z)) \quad \forall z \in [0, z_0], \quad g^0(z) = f(I(z)) \quad \forall z \in (z_0, z_b].$$

Evidently z_0 depends on I and R , and $z_0 = 0$ if $\alpha R(z) \geq I(z)$ on $[0, z_b]$, and $z_0 = z_b$ when $\alpha R(z) \leq I(z)$ on $[0, z_b]$.

In view of the above discussions for $g^0(z)$, we see that

$$\begin{aligned} \nu_0(g_z^0(z)) &= -\nu_{max}\text{sgn}(g_z^0(z)) = -\nu_{max} \quad \forall z \in [0, z_0), \\ \nu_0(g_z^0(z)) &= -\nu_{max}\text{sgn}(g_z^0(z)) = \nu_{max} \quad \forall z \in (z_0, z_b]. \end{aligned}$$

In this paper, we use a continuous approximation of the above step function used in [KL]; namely, we take

$$\nu(g_z^0(z)) = \nu_{\delta'}(z) := \nu_{max} \frac{z - z_0}{\delta' + |z - z_0|}.$$

It is easily seen that $\nu_{\delta'}(z) \rightarrow \nu_0(g_z^0(z))$ as $\delta' \rightarrow 0$. We stress again that $\nu_{\delta'}$ depends on I and R through the definition of z_0 .

Next we normalize the functions in (1.2) by

$$u(x) = b(z_b x)/Y, \quad v(x) = R(z_b x), \quad w(x) = I(z_b x), \quad 0 \leq x \leq 1,$$

and define

$$d_1 = z_b^2 D_b, \quad d_2 = z_b^2 D_R, \quad \sigma = \nu_{max} z_b, \quad \delta = \delta'/z_b,$$

$$A = a z_b Y, \quad A_0 = a_{bg} z_b, \quad \beta = h z_b, v_0 = R_{in}, \quad w_0 = I_{in}.$$

We denote

$$c(x) = c_{v,w}(x) = \frac{x - x_0}{\delta + |x - x_0|}, \quad x_0 = \frac{z_0}{z_b}.$$

Then after some simple calculations we find that the steady-state version of (1.2) becomes (1.1), or, written in a more comprehensive form,

$$(1.3) \quad \begin{cases} -[d_1 u_x + \sigma c(x)u]_x = [f(\min\{\alpha v, w\}) - m]u, & 0 < x < 1, \\ -d_2 v_{xx} = -f(\min\{\alpha v, w\})u, & 0 < x < 1, \\ d_1 u_x(0) + \sigma c(0)u(0) = d_1 u_x(1) + \sigma c(1)u(1) = 0, \\ v_x(0) = 0, \quad v_x(1) = \beta[v_0 - v(1)], \\ w(x) = w_0 \exp \left[-A_0 x - A \int_0^x u(s) ds \right], & 0 \leq x \leq 1. \end{cases}$$

Let us note that from the equation for u and the strong maximum principle, if u is nonnegative on $[0, 1]$, then it is either identically 0 or positive everywhere in $[0, 1]$. It is also easy to see that whenever u is positive, $0 < v < v_0$ in $[0, 1]$.

Since this paper is very long, and the techniques used in the first half of the paper are rather different from those in the second half, we divide it into two separate parts. Part I here is mainly concerned with the existence and nonexistence problem, and Part II studies the asymptotic behavior of the positive solutions as $\sigma \rightarrow \infty$.

In section 2, treating m as a parameter, we make use of a bifurcation argument to obtain two critical numbers $0 < m_* \leq m^*$ such that (1.3) has no positive solution when $m \geq m^*$ and it has at least one positive solution when $0 < m < m_*$. We also show that as m (the death rate of the biomass) decreases to 0, the biomass blows up everywhere; the exact limiting profiles of the biomass function and the nutrient function as $m \rightarrow 0$ are also obtained. In section 3, we show that as $\sigma \rightarrow \infty$, m_* and m^* converge to the same limit $f(\min\{\alpha v_0, w_0\})$. This demonstrates that our existence and nonexistence results are sharp for large σ .

The asymptotic behavior of the positive solutions when $\sigma \rightarrow \infty$ is investigated separately in Part II (see [DH]), where we fix $0 < m < f(\min\{\alpha v_0, w_0\})$ and study the asymptotic behavior of a positive solution (u_n, v_n) of (1.3) with $\sigma = \sigma_n \rightarrow \infty$.

2. Existence and nonexistence results. The function $c(x)$ appearing in (1.3) plays a very important role in our analysis. From the definition of $c(x) = c_{v,w}(x)$ we find that it is well defined if $v(x)$ is increasing in $[0, 1]$ and $w(x)$ is decreasing in $[0, 1]$, and

$$c(x) = \frac{x - x_0}{\delta + |x - x_0|},$$

where $x_0 \in [0, 1]$ is uniquely determined by the following:

$$\min\{\alpha v(x), w(x)\} = \alpha v(x) \quad \forall x \in [0, x_0]; \quad \min\{\alpha v(x), w(x)\} = w(x) \quad \forall x \in (x_0, 1].$$

It is easily seen that for the definition of $c_{v,w}(x)$, the requirement that v is increasing and w is decreasing can be relaxed; we can allow one (but not both) of the following:

- (i) v is nondecreasing, (ii) w is nonincreasing.

Let us also observe that $c(x)$ is a C^1 function, with $c'(x) = \delta(\delta + |x - x_0|)^{-2}$.

With this in mind, we find that $(u, v, w) = (0, v_0, w_*)$ solves (1.3), where $w_*(x) = w_0 e^{-A_0 x}$. We will call this the trivial solution. To find nontrivial solutions, we now treat m as a parameter and look for special values of m so that positive solutions of (1.3) may bifurcate from this trivial solution. If $m_* \geq 0$ is such a value, then there exist $m_n \rightarrow m_*$ and (u_n, v_n, w_n) solving (1.3) with $m = m_n$ such that $u_n > 0$, $v_n > 0$, and $u_n \rightarrow 0$, $v_n \rightarrow v_0$, and $w_n \rightarrow w_*$ and in $C^1[0, 1]$ as $n \rightarrow \infty$. Now $v_n(x)$ is increasing and $w_n(x)$ is decreasing for $x \in [0, 1]$; therefore $c_{v_n, w_n}(x)$ is well defined. Moreover, it is easily checked that $c_{v_n, w_n} \rightarrow c_{v_0, w_*}$ in $C([0, 1])$ as $n \rightarrow \infty$. To simplify the notation, we write $c^n(x) = c_{v_n, w_n}(x)$ and $c^0(x) = c_{v_0, w_*}(x)$. Therefore u_n satisfies

$$\begin{cases} -[d_1 u_n' + \sigma c^n(x) u_n]' = [f(\min\{\alpha v_n, w_n\}) - m_n] u_n & \text{in } (0, 1), \\ d_1 u_n' + \sigma c^n(x) u_n = 0 & \text{for } x = 0, 1. \end{cases}$$

Here and in what follows, we use the notation $u' = u_x$, etc. To determine the value of m_* , we first deduce a useful equation from the equation for u_n . So we define $\hat{u}_n = u_n / \|u_n\|_\infty$. Then we have

$$(2.1) \quad \begin{cases} -(d_1 \hat{u}_n' + \sigma c^n \hat{u}_n)' + m_n \hat{u}_n = f(\min\{\alpha v_n, w_n\}) \hat{u}_n & \text{in } (0, 1), \\ d_1 \hat{u}_n' + \sigma c^n \hat{u}_n = 0 & \text{for } x = 0, 1. \end{cases}$$

Since the right-hand side of the first equation in (2.1) and $\{\hat{u}_n\}$ are both bounded in $L^\infty([0, 1])$, and since $m_n, c^n, (c^n)'$ are bounded in $L^\infty([0, 1])$, we can use standard L^p theory for elliptic operators (see [GT]) to conclude that $\{\hat{u}_n\}$ is a bounded sequence in $W^{2,p}([0, 1])$ for any $p > 1$. By the Sobolev embedding theorem, we see that $\{\hat{u}_n\}$ is compact in $C^1([0, 1])$. By passing to a subsequence, we may assume that $\hat{u}_n \rightarrow \hat{u}$ in $C^1([0, 1])$, and then we easily see that \hat{u} satisfies (in the weak sense)

$$(2.2) \quad \begin{cases} -(d_1 \hat{u}' + \sigma c^0 \hat{u})' + m_* \hat{u} = f(\min\{\alpha v_0, w_*\}) \hat{u} & \text{in } (0, 1), \\ d_1 \hat{u}' + \sigma c^0 \hat{u} = 0 & \text{for } x = 0, 1. \end{cases}$$

Since $\hat{u} \geq 0$ and $\|\hat{u}\|_\infty = 1$, we necessarily have, by applying the strong maximum principle to (2.2), that $\hat{u} > 0$. This implies that $-m_*$ is the principal eigenvalue of the problem

$$(2.3) \quad \begin{cases} -(d_1 u' + \sigma c^0 u)' - f(\min\{\alpha v_0, w_*\}) u = \lambda u & \text{in } (0, 1), \\ d_1 u' + \sigma c^0 u = 0 & \text{for } x = 0, 1. \end{cases}$$

One easily checks that 0 is the first eigenvalue of the problem

$$\begin{cases} -(d_1 u' + \sigma c^0 u)' = \lambda u & \text{in } (0, 1), \\ d_1 u' + \sigma c^0 u = 0 & \text{for } x = 0, 1. \end{cases}$$

Since $-f(\min\{\alpha v_0, w_*\}) < 0$, by the characterization of the first eigenvalues (see, for example, Theorems 2.4 and 2.8 of [D]), the first eigenvalue of (2.3) is less than 0, and hence $m_* > 0$. On the other hand, if (u, v, w) is a positive solution to (1.3), then rewriting the equation for u in the form

$$\begin{cases} -(d_1 u' + \sigma c u)' - f(\min\{\alpha v, w\})u = -mu & \text{in } (0, 1), \\ d_1 u' + \sigma c u = 0 & \text{for } x = 0, 1, \end{cases}$$

we find that $-m$ is the first eigenvalue of the problem

$$\begin{cases} -(d_1 u' + \sigma c u)' - f(\min\{\alpha v, w\})u = \lambda u & \text{in } (0, 1), \\ d_1 u' + \sigma c u = 0 & \text{for } x = 0, 1, \end{cases}$$

which we denote by $\lambda_1(c, v, w)$. Clearly $v < v_0$ and $w < w_*$ in $(0, 1)$. It follows that $-f(\min\{\alpha v, w\}) > -f(\min\{\alpha v_0, w_*\})$ in $(0, 1)$ and hence

$$(2.4) \quad \lambda_1(c, v, w) > \lambda_1(c, v_0, w_*).$$

In this notation, clearly $-m_* = \lambda_1(c^0, v_0, w_*)$. Since $c(x) = (x - x_0)/(\delta + |x - x_0|)$ is determined completely by x_0 , it is convenient to introduce the notation

$$C_{x_0} = \frac{x - x_0}{\delta + |x - x_0|}$$

and

$$m^* = - \inf_{x_0 \in [0,1]} \lambda_1(C_{x_0}, v_0, w_*).$$

It is easy to show that m^* is achieved by some $x_0 \in [0, 1]$, and $m^* \geq m_*$. Moreover, from the above discussion, we have the following result.

PROPOSITION 2.1. *If (1.3) has a positive solution, then necessarily $m < m^*$.*

We will show that (1.3) has a positive solution for every $0 < m < m_*$. Before that we briefly discuss some further simple estimates for the values of m so that (1.3) has a positive solution. So suppose that (1.3) has a positive solution (u, v, w) . Integrating the equation for u over $[0, 1]$, we obtain

$$\int_0^1 [f(\min\{\alpha v, w\}) - m] u dx = 0.$$

Since $u > 0$, $[f(\min\{\alpha v, w\}) - m]$ must change sign over $(0, 1)$, and therefore

$$\min_{[0,1]} f(\min\{\alpha v(x), w(x)\}) < m < \max_{[0,1]} f(\min\{\alpha v(x), w(x)\}).$$

It follows that

$$(2.5) \quad f(\min\{\alpha v(0), w(1)\}) < m < f(\min\{\alpha v_0, w_0\}).$$

From a similar consideration, we have

$$(2.6) \quad m^*, m_* \in (f(\min\{\alpha v_0, w_*(1)\}), f(\min\{\alpha v_0, w_0\})).$$

We now use a global bifurcation argument to show that (1.3) has a positive solution for every $m \in (0, m_*)$. First, we transform (1.3) into an abstract nonlinear

equation. Due to its unconventional nature, we cannot use a simple inverse operator trick to do this. In fact, to cope with the rather implicit dependence of $c_{v,w}$ on (v, w) , in the following, we have to choose the function spaces for the abstract setting very carefully and then analyze the properties of the abstract operator mostly by definitions.

Fix $\gamma \in (0, 1)$ and set

$$K := \{\phi \in C^{1,\gamma}([0, 1]) : \phi \text{ is nondecreasing in } [0, 1]\},$$

$$P := \{\phi \in C^{1,\gamma}([0, 1]) : \phi \text{ is nonnegative in } [0, 1]\}.$$

Clearly they are closed convex sets in $C^{1,\gamma}([0, 1])$, and, moreover, P is a positive cone.

For given $(u, v) \in P \times K$ and $m \geq 0$, we define

$$w = w_0 \exp \left[-A_0 x - A \int_0^x u(s) ds \right],$$

$$c(x) = c_{v,w}(x), \quad v_+ = \max\{v, 0\}$$

and will use the solutions of the following problems to define an abstract operator $T(m, u, v)$ such that $T(m, u, v) = (u, v)$ for $(u, v) \in P \times K$ if and only if (u, v) is a nonnegative solution of (1.3). So we consider the problems

$$(2.7) \quad \begin{cases} -(d_1 \phi' + \sigma c \phi)' + (m + 1)\phi = f(\min\{\alpha v_+, w\})u + u, & 0 < x < 1, \\ d_1 \phi' + \sigma c \phi = 0, & x = 0, 1, \end{cases}$$

and

$$(2.8) \quad \begin{cases} -d_2 \psi'' = -f(\min\{\alpha v_+, w\})u, & 0 < x < 1, \\ \psi'(0) = 0, \quad \psi'(1) = \beta[v_0 - \psi(1)]. \end{cases}$$

Clearly (2.7) has a unique solution ϕ , and it is nonnegative. Let $\zeta(x) = v_0 - \psi(x)$. Then (2.8) becomes

$$(2.9) \quad \begin{cases} -d_2 \zeta'' = f(\min\{\alpha v_+, w\})u, & 0 < x < 1, \\ \zeta'(0) = 0, \quad \zeta'(1) + \beta \zeta(1) = 0. \end{cases}$$

It is easily seen that (2.9) has a unique solution ζ , and it is nonnegative. Moreover, from $\zeta'' \leq 0$, and $\zeta'(0) = 0$ we deduce that ζ is nonincreasing. Hence ψ is nondecreasing and $\psi(x) \leq v_0$ in $[0, 1]$. Thus the solution operator $(\phi, \psi) = T(m, u, v)$ is well defined for $(u, v) \in P \times K$ and $m \geq 0$, and $T(m, \cdot, \cdot)$ maps $P \times K$ into itself.

We show next that T is continuous. Suppose that $(m_n, u_n, v_n) \in [0, \infty) \times P \times K$, $m_n \rightarrow m$, and $(u_n, v_n) \rightarrow (u, v)$ in $C^{1,\gamma}([0, 1]) \times C^{1,\gamma}([0, 1])$. Then it is easily checked that $c^n := c_{v_n, w_n} \rightarrow c := c_{v, w}$ in $C^1([0, 1])$. Denote $(\phi_n, \psi_n) = T(m_n, u_n, v_n)$ and $\zeta_n = v_0 - \psi_n$. We have

$$(2.10) \quad \begin{cases} -(d_1 \phi'_n + \sigma c^n \phi_n)' + (m_n + 1)\phi_n = f(\min\{\alpha(v_n)_+, w_n\})u_n + u_n, & 0 < x < 1, \\ d_1 \phi'_n + \sigma c^n \phi_n = 0, & x = 0, 1, \end{cases}$$

and

$$(2.11) \quad \begin{cases} -d_2 \zeta_n'' = f(\min\{\alpha(v_n)_+, w_n\})u_n, & 0 < x < 1, \\ \zeta_n'(0) = 0, \quad \zeta_n'(1) + \beta \zeta_n(1) = 0. \end{cases}$$

Applying standard L^p theory to both (2.10) and (2.11), we find that $\{\phi_n\}$ and $\{\zeta_n\}$ are bounded in $W^{2,p}([0, 1])$ for all $p > 1$. Hence they are precompact in $C^{1,\gamma}([0, 1])$. This implies that by passing to a subsequence, $\phi_n \rightarrow \phi$ and $\psi_n \rightarrow \psi$ in $C^{1,\gamma}([0, 1])$. Moreover, letting $n \rightarrow \infty$ in (2.10) and (2.11) we find that necessarily $(\phi, \psi) = T(m, u, v)$. Therefore the entire original sequence converges with limit (ϕ, ψ) . This proves the continuity of T .

We further show that T is compact. Suppose that $\{(m_n, u_n, v_n)\} \subset [0, \infty) \times P \times K$ is bounded. Then along some subsequence n_k , (u_n, v_n) converges in the $C^1([0, 1]) \times C^1([0, 1])$ norm to some (u, v) , and, by passing to a further subsequence, we may assume that $m_{n_k} \rightarrow m$. We may now repeat the arguments in the above continuity proof to conclude that $T(m_{n_k}, u_{n_k}, v_{n_k}) \rightarrow T(m, u, v)$ in $C^{1,\gamma}([0, 1])$. Therefore T is a compact operator on $[0, \infty) \times P \times K$.

Suppose that $m \geq 0$ and $(u, v) \in P \times K$ satisfies $(u, v) = T(m, u, v)$. We claim that $v \geq 0$ in $[0, 1]$. Otherwise, due to the monotonicity of v there exists $x_0 \in (0, 1]$ such that $v < 0$ in $[0, x_0)$ and $v \geq 0$ in $(x_0, 1]$. Therefore, by (2.8),

$$-v'' = 0 \text{ in } (0, x_0).$$

Since $v'(0) = 0$, we deduce that $v'(x) = 0$ in $(0, x_0)$, and hence v is a negative constant in $(0, x_0)$, say $v = -c$. This is possible only if $x_0 = 1$ (otherwise, v is discontinuous at $x = x_0$ since $v \geq 0$ in $(x_0, 1]$), but then from $v'(1) = \beta[v_0 - v(1)]$ we deduce that $-c = v_0 > 0$, a contradiction. Therefore we must have $v \geq 0$ in $[0, 1]$, as claimed. Thus (u, v) is a nonnegative solution of (1.3) if and only if $(u, v) \in P \times K$ and $T(m, u, v) = (u, v)$.

In order to apply the global bifurcation theory to the operator equation

$$(u, v) - T(m, u, v) = 0,$$

we now calculate the Fréchet derivative of T with respect to (u, v) at $(m, 0, v_0)$, in the convex set $P \times K$ of $C^{1,\gamma}([0, 1]) \times C^{1,\gamma}([0, 1])$, where $m \geq 0$. From (2.7) and (2.8) we easily see that $T(m, 0, v_0) = (0, v_0)$. For $(u, v) \in C([0, 1]) \times C([0, 1])$, we define $(\xi, \eta) = L_m(u, v)$ to be the unique solution of the following linear problems:

$$(2.12) \quad \begin{cases} -(d_1 \xi' + \sigma c^0 \xi)' + (m + 1)\xi = f(\min\{\alpha v_0, w_*\})u + u, & 0 < x < 1, \\ d_1 \xi' + \sigma c^0 \xi = 0, & x = 0, 1, \end{cases}$$

$$(2.13) \quad \begin{cases} -d_2 \eta'' = -f(\min\{\alpha v_0, w_*\})u, & 0 < x < 1, \\ \eta'(0) = 0, \quad \eta'(1) + \beta \eta(1) = 0, \end{cases}$$

where c^0 and w_* are defined as at the beginning of this section.

Suppose $(u_n, v_n) \rightarrow (0, v_0)$ in $P \times K$. Denote

$$(\phi_n, \psi_n) = T(m, u_n, v_n) \text{ and } (\tau_n, \theta_n) = L_m(u_n, v_n - v_0).$$

We want to show that

$$(2.14) \quad \|(\phi_n, \psi_n) - (0, v_0) - (\tau_n, \theta_n)\| = o(\|(u_n, v_n - v_0)\|),$$

where $\|(u, v)\| = \max\{\|u\|, \|v\|\}$, and $\|u\| = \|u\|_{C^{1,\gamma}([0,1])}$. This would imply that the Fréchet derivative of T with respect to (u, v) at $(m, 0, v_0)$, in the convex set $P \times K$, is the linear operator L_m .

Suppose $(u_n, v_n) \rightarrow (0, v_0)$ in $P \times K$. Without loss of generality we assume that $u_n \not\equiv 0$. We define w_n and $c^n = c_{v_n, w_n}$ as before, and let $\hat{\phi}_n = \phi_n/\|u_n\|$, $\hat{u}_n = u_n/\|u_n\|$. Then

$$(2.15) \quad \begin{cases} -(d_1 \hat{\phi}'_n + \sigma c^n \hat{\phi}_n)' + (m+1)\hat{\phi}_n = f(\min\{\alpha(v_n)_+, w_n\})\hat{u}_n + \hat{u}_n, & 0 < x < 1, \\ d_1 \hat{\phi}'_n + \sigma c^n \hat{\phi}_n = 0, & x = 0, 1. \end{cases}$$

Since the right-hand side of the first equation in (2.15) is bounded in $L^\infty([0, 1])$, much as before we deduce from the L^p theory and the Sobolev imbedding theorem that there exists some positive constant C independent of n such that

$$\|\hat{\phi}_n\| \leq C \quad \forall n \geq 1.$$

We now define $\Phi_n = (\phi_n - \tau_n)/\|u_n\|$, and from the equations for ϕ_n and τ_n we obtain

$$(2.16) \quad \begin{cases} -(d_1 \Phi'_n + \sigma c^0 \Phi_n)' + (m+1)\Phi_n = f_n, & 0 < x < 1, \\ d_1 \Phi'_n + \sigma c^0 \Phi_n = g_n, & x = 0, 1, \end{cases}$$

where

$$f_n := [f(\min\{\alpha(v_n)_+, w_n\}) - f(\min\{\alpha v_0, w_*\})]\hat{u}_n + [\sigma(c^n - c^0)\hat{\phi}_n]'$$

and

$$g_n = \sigma(c^0 - c^n)\hat{\phi}_n.$$

It is easy to check that f_n converges to 0 in $L^\infty([0, 1])$, and g_n converges to 0 in $C^1([0, 1])$. Since $(m+1) \geq 1$, we may apply the L^p estimate to (2.16) to conclude that $\|\Phi_n\| \rightarrow 0$ as $n \rightarrow \infty$. Hence

$$\|\phi_n - \tau_n\| = o(\|u_n\|).$$

Define $\Psi_n = [(v_0 - \psi_n) + \theta_n]/\|u_n\|$. Then from the equations for ψ_n and θ_n we deduce that

$$(2.17) \quad \begin{cases} -d_2 \Psi''_n = [f(\min\{\alpha(v_n)_+, w_n\}) - f(\min\{\alpha v_0, w_*\})]\hat{u}_n, & 0 < x < 1, \\ \Psi'_n(0) = 0, \quad \Psi'_n(1) + \beta \Psi_n(1) = 0. \end{cases}$$

Since the right-hand side of the first equation in (2.17) converges to 0 in $L^\infty([0, 1])$, we can apply the L^p theory to (2.17) to conclude that $\|\Psi_n\| \rightarrow 0$ as $n \rightarrow \infty$. Therefore

$$\|\psi_n - v_0 - \theta_n\| = o(\|u_n\|).$$

Thus we have

$$\|(\phi_n, \psi_n) - (0, v_0) - (\tau_n, \theta_n)\| = o(\|(u_n, v_n - v_0)\|).$$

Summarizing the above discussions, we have the following result.

PROPOSITION 2.2. *The operator $T : [0, \infty) \times P \times K \rightarrow P \times K$ is completely continuous, and it is Fréchet differentiable at $(m, 0, v_0)$ with respect to (u, v) in the convex set $P \times K$, with derivative operator L_m . Moreover, $(u, v) = T(m, u, v)$ implies that $v \in P$; (u, v) is a nonnegative solution of (1.3) if and only if $(u, v) = T(m, u, v)$.*

We are now ready to prove the main result of this section.

THEOREM 2.3. *For every $m \in (0, m_*)$, problem (1.3) has at least one positive solution. Moreover, if m_n decreases to 0 and (u_n, v_n) is a positive solution of (1.3) with $m = m_n$, then $u_n \rightarrow \infty$ uniformly in $[0, 1]$ and there exists a unique $\tau \in (0, \frac{v_0}{1+\beta^{-1}})$ (determined by (2.27) below) such that*

$$\frac{u_n(x)}{\|u_n\|_\infty} \rightarrow \left(1 + \frac{x}{\delta}\right)^{\frac{\sigma}{d_1} \delta} e^{-\frac{\sigma}{d_1} x}, \quad v_n(x) \rightarrow \tau x + v_0 - \tau(1 + \beta^{-1})$$

uniformly in $[0, 1]$. Furthermore, for each $m \in (0, m_*)$, there is a positive solution (m, u, v) lying on the global bifurcation branch, $\Gamma = \{(m, u, v)\} \subset (0, \infty) \times C^{1,\gamma}([0, 1]) \times C^{1,\gamma}([0, 1])$, bifurcating from the trivial solution branch $\Gamma_0 := \{(m, 0, v_0) : m \in (-\infty, \infty)\}$ at $m = m_*$.

Proof. Since the proof is rather long, we divide it into several steps.

Step 1: Existence of an unbounded global solution branch.

We observe that 1 is an eigenvalue of L_{m_*} with eigenvector (ϕ_1, ψ_1) satisfying $\phi_1 > 0$ and $\psi_1 < 0$. Indeed, $\phi_1 > 0$ is a principal eigenfunction of (2.3) with $\lambda = -m_*$, and ψ_1 is the unique solution of (2.13) with $u = \phi_1$, and hence $\psi_1 < 0$. In order to apply the abstract global bifurcation theory in positive cones, we define $S : [0, \infty) \times P \times (-K) \rightarrow P \times (-K)$ by

$$S(m, u, \xi) = (\phi, \zeta) \text{ if and only if } T(m, u, v_0 - \xi) = (\phi, v_0 - \zeta).$$

Then from the properties of T we find that S is completely continuous. Moreover, if we denote by $DS(m, 0, 0)$ the Fréchet derivative of S with respect to (u, ξ) in $P \times (-K)$ at $(u, \xi) = (0, 0)$, then 1 is an eigenvalue of $DS(m_*, 0, 0)$ with eigenvector $(\phi_1, -\psi_1)$, where ϕ_1 and ψ_1 are as given above. Let us denote by P_0 the nonnegative functions in $(-K)$. Clearly P_0 is a cone in $C^{1,\gamma}([0, 1])$, and hence $P \times P_0$ is a cone in $C^{1,\gamma}([0, 1]) \times C^{1,\gamma}([0, 1])$. Moreover, it is easy to check through the definition of T that $S(m, \cdot, \cdot)$ maps $P \times P_0$ into itself, and 1 is the only eigenvalue of $DS(m_*, 0, 0)$ with an eigenvector in $P \times P_0$, and for any $m \geq 0, m \neq m_*$, 1 is not an eigenvalue of $DS(m, 0, 0)$ corresponding to an eigenvector in $P \times P_0$. These properties allow us to apply Corollary 18.4 of [A] to conclude that there exists a global unbounded branch of solutions of $(u, \zeta) = S(m, u, \zeta)$ in $R^1 \times (P \times P_0 \setminus \{(0, 0)\})$. We denote this global branch by $\tilde{\Gamma} = \{(m, u, \zeta)\}$ and define

$$\Gamma := \{(m, u, v_0 - \zeta) : (m, u, \zeta) \in \tilde{\Gamma}\}.$$

Clearly Γ is a global branch of solutions to $(u, v) = T(m, u, v)$ with $u \geq 0$ and $v = v_0 - \zeta \leq v_0$. We claim that $u \neq 0$. Otherwise, $u = 0$ and we deduce that $v = v_0 - \zeta = v_0$. Hence $\zeta = 0$, contradicting the fact that $(m, u, \zeta) \in R^1 \times (P \times P_0 \setminus \{(0, 0)\})$. Hence $u \geq 0$ and $u \neq 0$. It then follows from Proposition 2.2 that $v_0 - \zeta \geq 0$. But $u \neq 0$ implies that $v_0 - \zeta \neq 0$, and hence $(u, v_0 - \zeta)$ is a positive solution of (1.3). Thus we have proved that Γ is an unbounded branch of positive solutions of (1.3).

Step 2: We show that the m -range of Γ covers $(0, m_)$.*

If $(m, u, v) \in \Gamma$, then from Proposition 2.1 and (2.5) we deduce that $0 < m < m^*$. Therefore we can find a sequence $(m_n, u_n, v_n) \in \Gamma$ such that $m_n \rightarrow m_0 \in [0, m^*]$ and $\|(u_n, v_n)\| \rightarrow \infty$. Note that $c_{v_n, w_n} = C_{x_n}$ for some $x_n \in [0, 1]$ uniquely determined by v_n and w_n . By passing to a subsequence we may assume that $x_n \rightarrow x_0 \in [0, 1]$. Then it is easily seen that $C_{x_n} \rightarrow C_{x_0}$ in $C^1([0, 1])$. We necessarily have, by passing to a subsequence, that $\|u_n\|_\infty \rightarrow \infty$, for otherwise from the equation for u_n we can deduce that $\|u_n\|$ is bounded, which in turn implies that $\|v_n\|$ is bounded, contradicting our assumption that $\|(u_n, v_n)\| \rightarrow \infty$. Therefore we may assume that $\|u_n\|_\infty \rightarrow \infty$. Denote $\hat{u}_n = u_n/\|u_n\|_\infty$. Then we can use the L^p estimate to the equation for \hat{u}_n to deduce that $\{\hat{u}_n\}$ is precompact in $C^{1,\gamma}([0, 1])$. Hence we may assume that $\hat{u}_n \rightarrow \hat{u}$ in $C^{1,\gamma}([0, 1])$. Since $0 \leq f(\min\{\alpha v_n, w_n\}) \leq f(w_*)$, we may assume that $f(\min\{\alpha v_n, w_n\})$ converges to f_0 weakly in $L^2([0, 1])$. Clearly we also have $0 \leq f_0 \leq f(w_*)$. Passing to the weak limit in the equation for \hat{u}_n we deduce that \hat{u} is a weak solution of

$$(2.18) \quad \begin{cases} -(d_1 \hat{u}' + \sigma C_{x_0} \hat{u})' = (f_0 - m_0) \hat{u}, & 0 < x < 1, \\ d_1 \hat{u}' + \sigma C_{x_0} \hat{u} = 0, & x = 0, 1. \end{cases}$$

Since $\hat{u} \geq 0$ and $\|\hat{u}\|_\infty = 1$, we can apply the Harnack inequality and the strong maximum principle to (2.18) to conclude that $\hat{u} > 0$ in $[0, 1]$. This implies that $u_n = \|u_n\|_\infty \hat{u}_n \rightarrow \infty$ uniformly in $[0, 1]$. Therefore $w_n \rightarrow 0$ uniformly on any compact subset of $(0, 1]$. This implies that $f_0 = 0$, and hence we deduce from (2.18) that $-m_0$ is the first eigenvalue of

$$(2.19) \quad \begin{cases} -(d_1 u' + \sigma C_{x_0} u)' = \lambda u, & 0 < x < 1, \\ d_1 u' + \sigma C_{x_0} u = 0, & x = 0, 1. \end{cases}$$

Hence $m_0 = 0$ and

$$\hat{u} = \exp \left[-\frac{\sigma}{d_1} \int_0^x C_{x_0}(s) ds \right].$$

This implies that the entire original sequence $\{m_n\}$ converges to 0. By the connectedness of Γ , we can conclude that for every $m \in (0, m_*)$, (1.3) has at least one positive solution lying on Γ . Moreover, when $x_n \rightarrow x_0$, we have

$$\frac{u_n}{\|u_n\|_\infty} \rightarrow \exp \left[-\frac{\sigma}{d_1} \int_0^x C_{x_0}(s) ds \right].$$

Step 3: The limiting profile of u_n .

We will show in a moment that $x_0 = 0$ and hence the entire original sequence $u_n/\|u_n\|_\infty$ converges in $C^{1,\gamma}([0, 1])$ to $\exp \left[-\frac{\sigma}{d_1} \int_0^x C_0(s) ds \right]$.

Let $\zeta_n = v_0 - v_n$. Then

$$(2.20) \quad \begin{cases} -d_2 \zeta_n'' = f(\min\{\alpha v_n, w_n\}) u_n, & 0 < x < 1, \\ \zeta_n'(0) = 0, \quad \zeta_n'(1) + \beta \zeta_n(1) = 0. \end{cases}$$

We have

$$0 \leq f(\min\{\alpha v_n, w_n\}) u_n \leq f(w_n(x)) u_n(x).$$

Moreover,

$$\begin{aligned} f(w_n(x))u_n(x) &\leq (r/K_I)w_n(x)u_n(x) \\ &= (r/K_I)e^{-A_0x}e^{-A\int_0^x u_n(s)ds}u_n(x) \\ &\leq Ce^{-A\|u_n\|_\infty\int_0^x \hat{u}_n(s)ds}\|u_n\|_\infty\hat{u}_n(x). \end{aligned}$$

Since $\hat{u}_n \rightarrow \hat{u} > 0$ uniformly in $[0, 1]$, there exist $c_1, c_2 > 0$ such that $c_1 \leq \hat{u}_n \leq c_2$ and hence

$$\begin{aligned} g_n &:= Ce^{-A\|u_n\|_\infty\int_0^x \hat{u}_n(s)ds}\|u_n\|_\infty\hat{u}_n(x) \\ &\leq Cc_2e^{-Ac_1\|u_n\|_\infty x}\|u_n\|_\infty. \end{aligned}$$

One easily sees from the above inequality that $g_n \rightarrow 0$ uniformly in compact subsets of $(0, 1]$. It follows that

$$f_n := f(\min\{\alpha v_n, w_n\})u_n \rightarrow 0$$

uniformly in compact subsets of $(0, 1]$.

We can now prove that $x_0 = 0$. Otherwise, $x_0 \in (0, 1]$ and thus $x_n > x_0/2$ for all large n . Since $v_n(x)$ is increasing in x , and $\alpha v_n(x) < w_n(x)$ in $[0, x_n]$, it follows that

$$f(\min\{\alpha v_n, w_n\})u_n \leq f(\alpha v_n(x_0/2))u_n \leq f(w_n(x_0/2))u_n \text{ in } [0, x_0/2].$$

By our earlier estimates for g_n , we find that

$$f(w_n(x_0/2))u_n \leq Cc_2e^{-Ac_1\|u_n\|_\infty(x_0/2)}\|u_n\|_\infty \rightarrow 0.$$

Hence $f_n \rightarrow 0$ uniformly in $[0, 1]$, which implies, by (2.20), that $\zeta_n \rightarrow 0$ in $C^{1,\gamma}([0, 1])$. In particular, $\zeta_n \rightarrow 0$ uniformly in $[0, 1]$; but this leads to a contradiction:

$$v_0 - \zeta_n = v_n \leq w_n/\alpha \leq w_n(x_0/2)/\alpha \rightarrow 0 \text{ uniformly in } [x_0/2, x_n].$$

Hence $x_0 = 0$. Therefore we have

$$(2.21) \quad x_n \rightarrow 0 \text{ and } u_n/\|u_n\|_\infty \rightarrow \phi_0,$$

where

$$\phi_0(x) = \exp\left[-\frac{\sigma}{d_1}\int_0^x C_0(s)ds\right] = e^{-\frac{\sigma}{d_1}x}\left(1 + \frac{x}{\delta}\right)^{\frac{\sigma}{d_1}\delta}.$$

Step 4: The limiting profile of v_n .

Since $v_n \geq 0$ we have $0 \leq \zeta_n \leq v_0$. Moreover, due to (2.20) and the fact that $f_n \rightarrow 0$ uniformly in compact subsets of $(0, 1]$, we can use standard elliptic regularity theory and a diagonal process to find a subsequence of $\{\zeta_n\}$, still denoted by ζ_n , such that $\zeta_n \rightarrow \zeta$ in $C^1([\epsilon, 1])$ for every $\epsilon \in (0, 1)$, and ζ satisfies

$$-d_2\zeta'' = 0 \text{ in } (0, 1], \quad \zeta'(1) + \beta\zeta(1) = 0, \quad 0 \leq \zeta \leq v_0.$$

It follows that

$\zeta(x) = \tau(1 + \beta^{-1} - x)$ for some $\tau \geq 0$ to be determined below.

On the other hand, ζ_n can be explicitly expressed as

$$\zeta_n(x) = d_2^{-1}(1 + \beta^{-1} - x) \int_0^1 f_n(s)ds + d_2^{-1} \int_x^1 (x - s)f_n(s)ds.$$

Since

$$\int_0^1 f_n(s)ds = d_2\beta\zeta_n(1) \in (0, d_2\beta v_0(1)],$$

and $f_n \rightarrow 0$ uniformly on any compact subset of $(0, 1]$, one easily sees that

$$\int_x^1 (x - s)f_n(s)ds \rightarrow 0$$

uniformly for $x \in [0, 1]$. Thus

$$(2.22) \quad \tau = \beta\zeta(1) = \lim_{n \rightarrow \infty} \beta\zeta_n(1) = d_2^{-1} \lim_{n \rightarrow \infty} \int_0^1 f_n(s)ds = d_2^{-1} \lim_{n \rightarrow \infty} \int_0^\epsilon f_n(s)ds$$

for any $\epsilon \in (0, 1)$.

Since $v_n(x)$ is monotone increasing in x and $v_n''(x) \geq 0$, by an elementary argument we see that the fact that $v_n \rightarrow v_0 - \zeta$ in $C^1([\epsilon, 1])$ for every $\epsilon \in (0, 1)$ implies that $v_n \rightarrow v_0 - \zeta$ uniformly in $[0, 1]$.

We now show that $\tau > 0$. Suppose $\tau = 0$. Then $\zeta_n \rightarrow 0$ and hence $v_n \rightarrow v_0$ in $C^1([0, 1])$. Therefore, due to

$$w_n(0) = w_0 \text{ and } w_n \rightarrow 0 \text{ uniformly in } [\epsilon, 1] \forall \text{ small } \epsilon > 0,$$

when $w_0 > \alpha v_0$, we have $0 < x_n < 1$ for all large n , and

$$w_n(x_n) = \alpha v_n(x_n) \rightarrow \alpha v_0 \text{ as } n \rightarrow \infty;$$

if $w_0 < \alpha v_0$, then $x_n = 0$ and $w_n(x_n) = w_0$ for all large n ; if $w_0 = \alpha v_0$, then either $0 < x_n < 1$ and $w_n(x_n) = \alpha v_n(x_n)$, or $x_n = 0$ and $w_n(x_n) = w_0$; in either case we can conclude that $w_n(x_n) \rightarrow w_0$ as $n \rightarrow \infty$.

Summarizing, we find that we always have

$$w_0 e^{-A_0 x_n} e^{-A \int_0^{x_n} u_n(s)ds} = w_n(x_n) \rightarrow \sigma_0 := \min\{w_0, \alpha v_0\}.$$

Since $x_n \rightarrow 0$ and $\hat{u}_n \rightarrow \phi_0$, we have $e^{-A_0 x_n} = 1 + o(1)$, and

$$\int_0^{x_n} u_n(s)ds = \|u_n\|_\infty \int_0^{x_n} \hat{u}_n(s)ds = \|u_n\|_\infty x_n \phi_0(0)[1 + o(1)] = \|u_n\|_\infty x_n [1 + o(1)].$$

Here $o(1)$ denotes a generic sequence converging to 0 as $n \rightarrow \infty$. This implies that

$$w_0 e^{-A\|u_n\|_\infty x_n} \rightarrow \sigma_0.$$

Hence

$$(2.23) \quad \|u_n\|_\infty x_n \rightarrow \tau_0 := A^{-1} \ln \left(\frac{w_0}{\sigma_0} \right).$$

Since we now assume that $\tau = 0$, by (2.22) we must have $\lim_{n \rightarrow \infty} \int_0^\epsilon f_n(x) dx = 0$. On the other hand, making use of $x_n \rightarrow 0$, $\|u_n\|_\infty x_n \rightarrow \tau_0$, $\hat{u}_n \rightarrow \phi_0$, and $\phi_0(0) = 1$, we have, for all large n and small ϵ ,

$$\begin{aligned} \int_0^\epsilon f_n(x) dx &\geq \int_{x_n}^\epsilon f(w_n(x)) u_n(x) dx \\ &\geq \int_{x_n}^\epsilon f \left(w_0 e^{-A\epsilon} \exp[-A\|u_n\|_\infty 2\phi_0(0)x] \right) \frac{\phi_0(0)}{2} \|u_n\|_\infty dx \\ &= \int_{\|u_n\|_\infty x_n}^{\|u_n\|_\infty \epsilon} f \left(w_0 e^{-A\epsilon} e^{-2Ay} \right) \left(\frac{1}{2} \right) dy \\ &\rightarrow \left(\frac{1}{2} \right) \int_{\tau_0}^\infty f \left(w_0 e^{-A\epsilon} e^{-2Ay} \right) dy > 0. \end{aligned}$$

This contradiction shows that we must have $\tau > 0$.

In order to find the asymptotic limit of the entire sequence $\{v_n\}$, we need to determine the value of τ . Recall that by passing to a subsequence, $v_n \rightarrow v_0 - \zeta$ uniformly in $[0, 1]$. Since $x_n \rightarrow 0$, we have either

$$w_n(x_n) = \alpha v_n(x_n) \rightarrow \alpha[v_0 - \zeta(0)] = \alpha[v_0 - \tau(1 + \beta^{-1})] := \xi_\tau,$$

which is the case when $\xi_\tau < w_0$, or $w_n(x_n) \rightarrow w_0$ when $\xi_\tau \geq w_0$. By the expression of $w_n(x_n)$, as before, we deduce that

$$(2.24) \quad w_0 e^{-A\|u_n\|_\infty x_n} \rightarrow \tilde{\xi}_\tau := \min\{w_0, \xi_\tau\}, \quad \|u_n\|_\infty x_n \rightarrow \sigma_\tau := A^{-1} \ln(w_0/\tilde{\xi}_\tau).$$

For fixed $\epsilon \in (0, 1)$, we have

$$\int_0^\epsilon f_n(x) dx = \int_0^{x_n} f_n(x) dx + \int_{x_n}^\epsilon f(w_n(x)) u_n(x) dx.$$

It is easy to check that, in every possible case, we have

$$\int_0^{x_n} f_n(x) dx = f(\tilde{\xi}_\tau) \phi_0(0) [1 + o(1)] \|u_n\|_\infty x_n = f(\tilde{\xi}_\tau) \sigma_\tau [1 + o(1)].$$

Since we already know that $f(w_n(x)) u_n(x) \rightarrow 0$ uniformly on any compact subset of $(0, 1]$, we find that, for any fixed $\epsilon_1 \in (0, \epsilon)$ and all large n ,

(2.25)

$$\begin{aligned}
 & \int_{x_n}^\epsilon f(w_n(x))u_n(x)dx \\
 &= \int_{x_n}^{\epsilon_1} f(w_n(x))u_n(x)dx + o_n(1) \\
 &= \int_{x_n}^{\epsilon_1} f(w_0e^{-A_0x}e^{-A\int_0^x u_n(s)ds})u_n(x)dx + o_n(1) \\
 &= \int_{x_n}^{\epsilon_1} f(w_0[1 + o_{\epsilon_1}(1)]e^{-A\|u_n\|_\infty \hat{u}_n(0)x^{1+o_{\epsilon_1}(1)}})\|u_n\|_\infty \hat{u}_n(0)[1 + o_{\epsilon_1}(1)]dx + o_n(1) \\
 &= [1 + o_{\epsilon_1}(1)] \int_{x_n}^{\epsilon_1} f(w_0e^{-A\|u_n\|_\infty \hat{u}_n(0)x})\|u_n\|_\infty \hat{u}_n(0)dx + o_n(1) \\
 &= [1 + o_{\epsilon_1}(1)] \int_{\|u_n\|_\infty \hat{u}_n(0)x_n}^{\epsilon_1\|u_n\|_\infty \hat{u}_n(0)} f(w_0e^{-Ay})dy + o_n(1) \\
 &= [1 + o_{\epsilon_1}(1)] \left[\int_{\sigma_\tau}^\infty f(w_0e^{-Ay})dy + o_n(1) \right] + o_n(1),
 \end{aligned}$$

where $o_n(1) \rightarrow 0$ as $n \rightarrow \infty$ for fixed ϵ_1 , and $o_{\epsilon_1}(1) \rightarrow 0$ as $\epsilon_1 \rightarrow 0$ uniformly in n . For arbitrary $\epsilon_1 \in (0, \epsilon)$, we first let $n \rightarrow \infty$ and then let $\epsilon_1 \rightarrow 0$, and we obtain from (2.25) that

$$\lim_{n \rightarrow \infty} \int_{x_n}^\epsilon f(w_n(x))u_n(x)dx = \int_{\sigma_\tau}^\infty f(w_0e^{-Ay})dy.$$

Therefore

$$(2.26) \quad \lim_{n \rightarrow \infty} \int_0^\epsilon f_n(x)dx = f(\tilde{\xi}_\tau)\sigma_\tau + \int_{\sigma_\tau}^\infty f(w_0e^{-Ay})dy.$$

Making use of (2.24) and (2.26), we can rewrite (2.22) as

$$(2.27) \quad d_2\tau = f(w_0e^{-A\sigma_\tau})\sigma_\tau + \int_{\sigma_\tau}^\infty f(w_0e^{-Ay})dy.$$

It can be easily checked that the function

$$F(\theta) := f(w_0e^{-A\theta})\theta + \int_\theta^\infty f(w_0e^{-Ay})dy$$

satisfies $F'(\theta) < 0$ and hence it is decreasing in $[0, \infty)$, with

$$F(0) = \int_0^\infty f(w_0e^{-Ay})dy = \frac{r}{A} \ln \left(\frac{w_0 + K_I}{K_I} \right), \quad F(\infty) = 0.$$

From the definition of σ_τ , we find that $\tau \rightarrow \sigma_\tau$ is nondecreasing, with

$$\sigma_0 = A^{-1} \ln \left(\frac{w_0}{\min\{w_0, \alpha v_0\}} \right), \quad \sigma_{v_0/(1+\beta^{-1})} = \infty.$$

Therefore

$\tau \rightarrow F(\sigma_\tau)$ is nonincreasing in $[0, v_0/(1 + \beta^{-1})]$ with $F(\sigma_0) > 0, F(\sigma_{v_0/(1+\beta^{-1})}) = 0$.

This implies that (2.27) has a unique solution $\tau \in (0, v_0/(1 + \beta^{-1}))$. Thus,

$$v_n \rightarrow v_0 - \tau(1 + \beta^{-1} - x) \text{ uniformly in } [0, 1],$$

with τ uniquely determined by (2.27). Since $\tau > 0$ is uniquely determined, the above convergence is true for the entire original sequence $\{v_n\}$.

The proof is complete. \square

3. The limit of m_* and m^* as $\sigma \rightarrow \infty$. In order to investigate the asymptotic behavior of the positive solutions of (1.3) as $\sigma \rightarrow \infty$, we need to first understand the limits of m_* and m^* as $\sigma \rightarrow \infty$. To stress their dependence on σ , we write $m_* = m_*(\sigma)$ and $m^* = m^*(\sigma)$. Let us recall that (2.6) holds; that is, $m_*(\sigma)$ and $m^*(\sigma)$ are always between the positive numbers $f(\min\{\alpha v_0, w_*(1)\})$ and $f(\min\{\alpha v_0, w_0\})$. We now prove the following result.

THEOREM 3.1.

$$(3.1) \quad \lim_{\sigma \rightarrow \infty} m_*(\sigma) = \lim_{\sigma \rightarrow \infty} m^*(\sigma) = f(\min\{\alpha v_0, w_0\}).$$

Proof. Since

$$m_*(\sigma) \leq m^*(\sigma) \leq f(\min\{\alpha v_0, w_0\}),$$

we need only show that

$$\lim_{\sigma \rightarrow \infty} m_*(\sigma) = f(\min\{\alpha v_0, w_0\}).$$

Moreover, it suffices to prove this along an arbitrary sequence of positive numbers increasing to ∞ . Let $\{\sigma_n\}$ be such a sequence, and denote $m_n = m_*(\sigma_n)$. By definition, there exists $u_n > 0$ in $[0, 1]$ such that

$$(3.2) \quad \begin{cases} -[d_1 u'_n + \sigma_n c^0(x) u_n]' = [f(\min\{\alpha v_0, w_*(x)\}) - m_n] u_n & \text{in } (0, 1), \\ d_1 u'_n + \sigma_n c^0(x) u_n = 0 & \text{for } x = 0, 1. \end{cases}$$

To simplify the notation, we will write

$$f_0(x) = f(\min\{\alpha v_0, w_*(x)\}).$$

Moreover, we define x_0^* by

$$c^0(x) = C_{x_0^*}(x) = \frac{x - x_0^*}{\delta + |x - x_0^*|}.$$

If $x_0^* = 1$, i.e., $\alpha v_0 \leq w_*(x)$ in $[0, 1]$ and hence $f_0(x) \equiv f(\alpha v_0)$, then clearly

$$m_n \equiv f(\alpha v_0) = f(\min\{\alpha v_0, w_0\}).$$

So (3.1) holds trivially in this case.

Suppose from now on that $x_0^* \in [0, 1)$. Therefore $f_0(x)$ is a constant in $[0, x_0^*)$ and is decreasing in $[x_0^*, 1]$. As before, integrating the first equation of (3.2) we find that $(f_0(x) - m_n)$ must change sign in $(0, 1)$, and therefore there exists a unique $x_n \in (0, 1)$ such that $f_0(x) > m_n$ in $[0, x_n)$, and $f_0(x) < m_n$ in $(x_n, 1]$.

By (2.6), $\{m_n\}$ is a bounded sequence, and, by passing to a subsequence, we may assume that

$$m_n \rightarrow m_0 \leq f(\min\{\alpha v_0, w_0\}) = f_0(x_0^*).$$

To determine the value of m_0 , we use several steps.

Step 1: Change of variables.

Let

$$\phi_n(x) := \exp \left[-\frac{\sigma_n}{d_1} \int_{x_0^*}^x c^0(x) dx \right].$$

Clearly

$$d_1 \phi_n' + \sigma_n c^0 \phi_n = 0, \quad \phi_n(x_0^*) = 1, \quad 0 < \phi_n(x) \leq 1 \text{ in } [0, 1],$$

and $\phi_n \rightarrow 0$ uniformly on compact subsets of $[0, 1] \setminus \{x_0^*\}$.

Define

$$\psi_n(x) = u_n(x) / \phi_n(x).$$

Then (3.2) becomes

$$(3.3) \quad \begin{cases} -(d_1 \phi_n \psi_n')' = [f_0(x) - m_n] \phi_n \psi_n & \text{in } (0, 1), \\ \psi_n' = 0 & \text{for } x = 0, 1. \end{cases}$$

Define

$$\xi_n(x) = \sigma_n^{-1/2} \int_0^x \frac{1}{\phi_n(s)} ds.$$

Then ξ_n is an increasing function in $[0, 1]$, with $\xi_n(0) = 0$ and

$$\xi_n(1) = y_n := \sigma_n^{-1/2} \int_0^1 \frac{1}{\phi_n(x)} dx \rightarrow \infty \text{ as } n \rightarrow \infty.$$

Let $\eta_n : [0, y_n] \rightarrow [0, 1]$ be the inverse function of $\xi_n(x)$, and define

$$U_n(y) = \psi_n(\eta_n(y)) = \psi_n(x).$$

From (3.3) a simple calculation shows that

$$(3.4) \quad \begin{cases} -U_n'' = d_1^{-1} \sigma_n \phi_n^2(\eta_n(y)) [f_0(\eta_n(y)) - m_n] U_n & \text{in } (0, y_n), \\ U_n'(0) = U_n'(y_n) = 0. \end{cases}$$

Step 2: Estimates of $\sigma_n \phi_n^2(\eta_n(y))$.

For our later estimates, we need to analyze the function $\tilde{\phi}_n(y) := \sigma_n \phi_n^2(\eta_n(y))$. To this end, for some $\tau > 0$ small to be determined later, we define \hat{y}_n and Δ_n by

$$\hat{y}_n := \xi_n(x_0^*), \quad \sigma_n^{-\tau} := \sigma_n^{-1/2} \int_{x_0^*}^{x_0^* + \Delta_n} \frac{1}{\phi_n(x)} dx,$$

so that

$$\hat{y}_n \pm \sigma_n^{-\tau} = \xi_n(x_0^* \pm \Delta_n), \quad x_0^* \pm \Delta_n = \eta_n(\hat{y}_n \pm \sigma_n^{-\tau}).$$

We will show that $\tilde{\phi}_n(y)$ behaves like a δ -function concentrating at $y = \hat{y}_n$. For definiteness, we assume that $x_0^* > 0$; the case $x_0^* = 0$ can be treated by a simple modification of the arguments below. Then it is easily seen that as $n \rightarrow \infty$,

$$\hat{y}_n \rightarrow \infty, \quad y_n - \hat{y}_n \rightarrow \infty, \quad \Delta_n \rightarrow 0.$$

Using $c^0(x) = \frac{x-x_0^*}{\delta+|x-x_0^*|}$, we can easily check that, for any given small $\epsilon > 0$, there exists $\delta_0 = \delta_0(\epsilon) > 0$ small so that, when $|x - x_0^*| \leq \delta_0$,

$$(3.5) \quad \exp \left[-\frac{\sigma_n}{2\delta d_1}(x - x_0^*)^2 \right] \leq \phi_n(x) \leq \exp \left[-\frac{\sigma_n(1-\epsilon)}{2\delta d_1}(x - x_0^*)^2 \right].$$

Therefore, for all large n ,

$$\begin{aligned} \sigma_n^{1/2-\tau} &= \int_{x_0^*}^{x_0^*+\Delta_n} \frac{1}{\phi_n(x)} dx \\ &\leq \int_{x_0^*}^{x_0^*+\Delta_n} \exp \left[\frac{\sigma_n}{2\delta d_1}(x - x_0^*)^2 \right] dx \\ &= \int_0^{\Delta_n} \exp \left(\frac{\sigma_n}{2\delta d_1}x^2 \right) dx \\ &\leq \Delta_n \exp \left(\frac{\sigma_n \Delta_n^2}{2\delta d_1} \right), \\ \sigma_n^{1/2-\tau} &= \int_{x_0^*}^{x_0^*+\Delta_n} \frac{1}{\phi_n(x)} dx \\ &\geq \int_{x_0^*}^{x_0^*+\Delta_n} \exp \left[\frac{\sigma_n(1-\epsilon)}{2\delta d_1}(x - x_0^*)^2 \right] dx \\ &= \int_0^{\Delta_n} \exp \left[\frac{\sigma_n(1-\epsilon)}{2\delta d_1}x^2 \right] dx \\ &\geq \int_{(1-\epsilon)\Delta_n}^{\Delta_n} \exp \left[\frac{\sigma_n(1-\epsilon)}{2\delta d_1}x^2 \right] dx \\ &\geq \epsilon \Delta_n \exp \left[\frac{\sigma_n \Delta_n^2 (1-\epsilon)^3}{2\delta d_1} \right]. \end{aligned}$$

It follows that $\sigma_n \Delta_n^2 \rightarrow \infty$, and

$$(3.6) \quad \exp \left(-\frac{\sigma_n \Delta_n^2}{2\delta d_1} \right) \leq \Delta_n \sigma_n^{\tau-1/2},$$

and

$$\sigma_n^{1/2} \Delta_n \exp \left[\frac{\sigma_n \Delta_n^2 (1-\epsilon)^3}{2\delta d_1} \right] \leq \epsilon^{-1} \sigma_n^{1-\tau},$$

which gives

$$\left(\frac{1}{2} \right) \ln(\sigma_n \Delta_n^2) + \frac{(1-\epsilon)^3}{2\delta d_1} \sigma_n \Delta_n^2 \leq \ln \left(\epsilon^{-1} \sigma_n^{1-\tau} \right).$$

Since $\ln(\sigma_n \Delta_n^2) = o(\sigma_n \Delta_n^2)$, and

$$\ln(\epsilon^{-1} \sigma_n^{1-\tau}) = (1 - \tau) \ln \sigma_n + o(\ln \sigma_n),$$

the last inequality above implies that

$$(3.7) \quad \sigma_n \Delta_n^2 \leq C_\epsilon \ln \sigma_n$$

for some $C_\epsilon > 0$ and all large n .

As a consequence of (3.5), (3.6), and (3.7) we have

$$\begin{aligned} \tilde{\phi}_n(\hat{y}_n \pm \sigma_n^{-\tau}) &= \sigma_n \phi_n^2(x_0^* \pm \Delta_n) \\ &\leq \sigma_n \exp \left[-\frac{\sigma_n(1-\epsilon)}{\delta d_1} \Delta_n^2 \right] \\ &\leq \sigma_n [\Delta_n^2 \sigma_n^{2\tau-1}]^{(1-\epsilon)} \\ &\leq [C_\epsilon \ln \sigma_n]^{(1-\epsilon)} \sigma_n^{1+2(\tau-1)(1-\epsilon)} \\ &\rightarrow 0 \text{ as } n \rightarrow \infty, \end{aligned}$$

provided that τ is chosen in the interval $(0, 1/2)$ and $\epsilon > 0$ is small enough.

From the property of $\phi_n(x)$, we see that the above estimates imply that

$$(3.8) \quad \tilde{\phi}_n(y) \rightarrow 0 \text{ uniformly in } [0, y_n] \setminus [\hat{y}_n - \sigma_n^{-\tau}, \hat{y}_n + \sigma_n^{-\tau}],$$

and for any $M > 0$,

$$\begin{aligned} \int_{\hat{y}_n-M}^{\hat{y}_n+M} \tilde{\phi}_n(y) dy &= \int_{\hat{y}_n-\sigma_n^{-\tau}}^{\hat{y}_n+\sigma_n^{-\tau}} \tilde{\phi}_n(y) dy + o(1) \\ &= \int_{\eta_n(\hat{y}_n-\sigma_n^{-\tau})}^{\eta_n(\hat{y}_n+\sigma_n^{-\tau})} \sigma_n \phi_n^2(x) \xi_n'(x) dx + o(1) \\ &= \int_{x_0^*-\Delta_n}^{x_0^*+\Delta_n} \sigma_n^{1/2} \phi_n(x) dx + o(1) \\ &= 2\sigma_n^{1/2} \int_0^{\Delta_n} \exp \left(-\frac{\sigma_n x^2}{2\delta d_1} [1 + o(1)] \right) dx + o(1) \\ &= [2 + o(1)] \int_0^{\sigma_n^{1/2} \Delta_n} \exp \left(-\frac{x^2}{2\delta d_1} \right) dx + o(1) \\ &= 2 \int_0^\infty \exp \left(-\frac{x^2}{2\delta d_1} \right) dx + o(1). \end{aligned}$$

In other words, for any $M > 0$,

$$(3.9) \quad \lim_{n \rightarrow \infty} \int_{\hat{y}_n-M}^{\hat{y}_n+M} \tilde{\phi}_n(y) dy = c_0 := 2 \int_0^\infty \exp \left(-\frac{x^2}{2\delta d_1} \right) dx.$$

Step 3: The limiting profile of $U_n(y)$.

We now define

$$\begin{aligned} \hat{U}_n(y) &= U_n(y + \hat{y}_n - 1) / U_n(\hat{y}_n - 1), \\ \hat{f}_n &= d_1^{-1} \tilde{\phi}_n(y + \hat{y}_n - 1) [f_0(\eta_n(y + \hat{y}_n - 1)) - m_n]. \end{aligned}$$

Clearly $\hat{U}_n(0) = 1$, and by (3.4) we have

$$(3.10) \quad \begin{cases} -\hat{U}_n'' = \hat{f}_n \hat{U}_n \text{ in } (1 - \hat{y}_n, y_n + 1 - \hat{y}_n), \\ \hat{U}_n'(1 - \hat{y}_n) = \hat{U}_n'(y_n + 1 - \hat{y}_n) = 0. \end{cases}$$

From (3.8) we see that

$$\hat{f}_n(y) \rightarrow 0 \text{ uniformly in } [1 - \hat{y}_n, 1 - \sigma_n^{-\tau}] \cup [1 + \sigma_n^{-\tau}, y_n + 1 - \hat{y}_n].$$

Since $\hat{U}_n(0) = 1$, the boundedness of \hat{f}_n over $[1 - \hat{y}_n, 1 - \sigma_n^{-\tau}]$ allows us to apply the Harnack inequality to conclude that \hat{U}_n has a bound C_J independent of n over any bounded interval $J \subset [1 - \hat{y}_n, 1 - \sigma_n^{-\tau}]$ with $0 \in J$. We can now apply to (3.10) the L^p theory, the Sobolev imbedding theorem, and a standard diagonal argument, to obtain a subsequence of $\{\hat{U}_n\}$, still denoted by \hat{U}_n , such that $\hat{U}_n \rightarrow \hat{U}$ in $C^1(J)$ for any bounded interval $J \subset (-\infty, 1)$, and \hat{U} satisfies

$$(3.11) \quad \hat{U}'' = 0 \text{ in } (-\infty, 1), \quad \hat{U}(0) = 1.$$

Since \hat{U} is nonnegative in $(-\infty, 1)$, we deduce from (3.11) that

$$\hat{U}(y) = 1 + ay, \quad a \in [-1, 0].$$

Now consider the sequence $\{\hat{U}_n(2)\}$. We claim that this is a bounded sequence. Indeed, from our earlier observation for the sign of $[f_0(x) - m_n]$, we know that the right-hand side of the first equation in (3.4) changes sign from positive to negative when y increases across $\tilde{y}_n := \xi_n(x_n)$. It follows that $U_n''(y)$ changes sign from negative to positive as y increases across \tilde{y}_n . Since $U_n'(0) = U_n'(y_n) = 0$, we find that $U_n' \leq 0$ in $[0, y_n]$ and hence $U_n(y)$ is nonincreasing in y , which implies that $\hat{U}_n(y)$ is nonincreasing in y and hence $0 \leq \hat{U}_n(2) \leq \hat{U}_n(0) = 1$. We can now use the fact that $\hat{f}_n \rightarrow 0$ uniformly in $[1 + \sigma_n^{-\tau}, y_n + 1 - \hat{y}_n]$, as above, to conclude that, subject to passing to a further subsequence, $\hat{U}_n \rightarrow \hat{U}_*$ in $C^1(J)$ for any bounded interval $J \subset (1, \infty)$, and \hat{U}_* satisfies

$$\hat{U}_*'' = 0, \quad 0 \leq \hat{U}_* \leq 1 \text{ in } (1, \infty).$$

Therefore \hat{U}_* must be a constant, say $\hat{U}_* \equiv b$.

Using (3.9) we find that \hat{f}_n is a bounded sequence in $L^1([0, 2])$. By (3.10), we have

$$\hat{U}_n'(y) = \hat{U}_n'(0) - \int_0^y \hat{f}_n(y) \hat{U}_n(y) dy \quad \forall y \in [0, 2].$$

Since $\hat{U}_n'(0) \rightarrow \hat{U}'(0) = a$ and $0 \leq \hat{U}_n(y) \leq \hat{U}_n(0) = 1$, the above identity implies that $|\hat{U}_n'(y)| \leq C$ for some $C > 0$ and all $n \geq 1$ and $y \in [0, 2]$. Therefore $\{\hat{U}_n(y)\}$ is equicontinuous in $[0, 2]$. It follows that \hat{U}_* must be a continuous extension of \hat{U} . Therefore $b = 1 + a$.

Step 4: We show that $m_0 = f(x_0^)$.*

We are now ready to determine the value of m_0 by using our estimates for $\tilde{\phi}_n$ and \hat{U}_n . We note that (3.8) and (3.9) imply that, for large n , $\tilde{\phi}_n(y + \hat{y}_n - 1)$ behaves like the δ -function concentrating at $y = 1$. We now use these properties of $\tilde{\phi}_n$ and (3.10) to obtain

$$\hat{U}_n'(2) = \hat{U}_n'(0) - \int_0^2 \hat{f}_n(y) \hat{U}_n(y) dy \rightarrow a - [f_0(x_0^*) - m_0](1 + a)c_0.$$

But since \hat{U}_* is a constant function over $(1, \infty)$, we have $\hat{U}'_n(2) \rightarrow \hat{U}'_*(2) = 0$. Therefore

$$a = [f_0(x_0^*) - m_0](1 + a)c_0.$$

The right-hand side of the above identity is nonnegative, but $a \leq 0$. Therefore we must have $a = 0$ and $m_0 = f_0(x_0^*)$. This implies that the entire original sequence $\{m_n\}$ converges to $f_0(x_0^*)$. Hence (3.1) holds, and the proof is complete. \square

If we fix m such that $0 < m < f(\min\{\alpha v_0, w_0\})$ and let σ_n be an increasing sequence of positive numbers converging to ∞ , then by Theorems 2.3 and 3.1, for all large n , (1.3) with $\sigma = \sigma_n$ has at least one positive solution. Suppose that (u_n, v_n) is such a solution. We will analyze the behavior of (u_n, v_n) as $n \rightarrow \infty$. This will be done in Part II; see [DH].

Acknowledgments. We thank the referees for their useful suggestions and comments. One referee pointed us to extra references ([IT] and [HTKS]) and provided detailed suggestions on improving the exposition of the paper (for both Part I and Part II); these have been very helpful. Y. Du thanks NCTS for the hospitality during his visit, when a major part of this paper was being written.

REFERENCES

- [A] H. AMANN, *Fixed point equations and nonlinear eigenvalue problems in ordered Banach spaces*, SIAM Rev., 18 (1976), pp. 620–709.
- [BFH] E. BERETTA, A. FASANO, AND Y. HOSONO, *Equilibrium of a phytoplankton population in a laboratory controlled system*, Surveys Math. Indust., 1 (1992), pp. 283–336.
- [BFHK] E. BERETTA, A. FASANO, Y. HOSONO, AND V. B. KOLMANOVSKIĬ, *Stability analysis of the phytoplankton vertical steady states in a laboratory test tube*, Math. Methods Appl. Sci., 17 (1994), pp. 551–575.
- [D] Y. DU, *Order Structure and Topological Methods in Nonlinear Partial Differential Equations: Maximal Principles and Applications*, Vol. 1, World Scientific, Hackensack, NJ, 2006.
- [DH] Y. DU AND S.-B. HSU, *Concentration phenomena in a nonlocal quasi-linear problem modelling phytoplankton II: Limiting profile*, SIAM J. Math. Anal., 40 (2008), pp. 1441–1470.
- [EATSH] U. EBERT, M. ARRAYAS, N. M. TEMME, B. P. SOMMEIJER, AND J. HUISMAN, *Critical conditions for phytoplankton blooms*, Bull. Math. Biol., 63 (2001), pp. 1095–1124.
- [GT] D. GILBARG AND N. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, Springer-Verlag, Berlin, 1983.
- [HTKS] J. HUISMAN, N. N. P. THI, D. M. KARL, AND B. SOMMEIJER, *Reduced mixing generates oscillations and chaos in the oceanic deep chlorophyll maximum*, Nature, 439 (2006), pp. 322–325.
- [IT] H. ISHII AND I. TAKAGI, *Global stability of stationary solutions to a nonlinear diffusion equation in phytoplankton dynamics*, J. Math. Biol., 16 (1982), pp. 1–24.
- [KL] C. A. KLAUSMEIER AND E. LITCHMAN, *Algal games: The vertical distribution of phytoplankton in poorly mixed water columns*, Limnol. Oceanogr., 46 (2001), pp. 1998–2007.
- [PT] N. N. PHAM THI, *On positive solutions in a phytoplankton-nutrient model*, J. Comput. Appl. Math., 177 (2005), pp. 467–473.
- [PTHS] N. N. PHAM THI, J. HUISMAN, AND B. P. SOMMEIJER, *Simulation of three-dimensional phytoplankton dynamics: Competition in light-limited environments*, J. Comput. Appl. Math., 174 (2005), pp. 57–77.
- [YN] K. YOSHIYAMA AND H. NAKAJIMA, *Catastrophic shifts in vertical distributions of phytoplankton, the existence of a bifurcation set*, J. Math. Biol., 52 (2006), pp. 235–276.

CONCENTRATION PHENOMENA IN A NONLOCAL QUASI-LINEAR PROBLEM MODELLING PHYTOPLANKTON II: LIMITING PROFILE*

YIHONG DU[†] AND SZE-BI HSU[‡]

Abstract. This is Part II of our study on the positive steady state of a quasi-linear reaction-diffusion system in one space dimension introduced by Klausmeier and Litchman for the modelling of the distributions of phytoplankton biomass and its nutrient. In Part I, we proved nearly optimal existence and nonexistence results. In Part II, we obtain complete descriptions of the profile of the solutions when the coefficient of the drifting term is large, rigorously proving the numerically observed phenomenon of concentration of biomass for this model. Moreover, we reveal four critical numbers for the model and provide further insights to the problem being modelled.

Key words. quasi-linear, nonlocal dependence, phytoplankton, concentration phenomenon, reaction-diffusion equation

AMS subject classifications. 35J55, 35J65, 92D25

DOI. 10.1137/070706641

1. Introduction. We continue our investigation in [DH] on the problem

$$(1.1) \quad \begin{cases} -[d_1 u_x + \sigma c(x)u]_x = [g(x) - m]u, & 0 < x < 1, \\ -d_2 v_{xx} = -g(x)u, & 0 < x < 1, \\ d_1 u_x + \sigma c(x)u = 0, & x = 0, 1, \\ v_x(0) = 0, v_x(1) = \beta[v_0 - v(1)], \end{cases}$$

where d_1, d_2, σ, m, v_0 , and β are positive constants,

$$g(x) = f(\min\{\alpha v(x), w(x)\}), \quad f(s) = \frac{rs}{K_I + s},$$

and

$$w(x) = w_0 \exp \left[-A_0 x - A \int_0^x u(s) ds \right],$$

with α, r, K_I, w_0, A , and A_0 positive constants. We are interested in positive solutions of (1.1), namely, $u > 0$ and $v > 0$ in $[0, 1]$. From (1.1) it is easy to see that for any such solution v is an increasing function. Clearly w is a decreasing function. The function $c(x)$ is defined by

$$c(x) = \frac{x - x_0}{\delta + |x - x_0|},$$

*Received by the editors October 29, 2007; accepted for publication (in revised form) July 2, 2008; published electronically November 5, 2008. This research was partially supported by the Australian Research Council, the NCTS, and the National Science Council of Taiwan.

<http://www.siam.org/journals/sima/40-4/70664.html>

[†]Department of Mathematics, School of Science and Technology, University of New England, Armidale, NSW2351, Australia (ydu@turing.une.edu.au), and Department of Mathematics, Qufu Normal University, Qufu, Shandong 273165, People's Republic of China.

[‡]Department of Mathematics, National Tsing-Hua University, Hsinchu, Taiwan 300, Republic of China (sbhsu@math.nthu.edu.tw).

where $\delta > 0$ is a small constant and $x_0 \in [0, 1]$ is uniquely determined by the following description:

$$\min\{\alpha v(x), w(x)\} = \alpha v(x) \quad \forall x \in [0, x_0]; \quad \min\{\alpha v(x), w(x)\} = w(x) \quad \forall x \in (x_0, 1].$$

(Due to the monotonicity of $v(x)$ and $w(x)$, such x_0 always exists.)

Problem (1.1) is a rescaled version of a model proposed by Klausmeier and Litchman in [KL] for the study of phytoplankton in a one-dimensional water column, where $u(x)$ represents the distribution of phytoplankton biomass, $v(x)$ stands for the distribution of nutrient, and x denotes the depth in the water column, with $x = 0$ at the surface and $x = 1$ at the bottom. The term $\sigma c(x)$ is used to describe the active movement of the biomass towards spatial location with optimal growth condition. Klausmeier and Litchman [KL] use this system to model the concentration phenomenon of phytoplankton in lakes and oceans, and the numerical analysis in [KL] demonstrates that, for large σ , the biomass function $u(x)$ concentrates at a certain level $x = x_*$, while the nutrient function $v(x)$ is close to a piecewise linear function. They then treat u as a constant multiple of the δ -function concentrating at x_* and propose a game theoretical model to determine the location of x_* . We refer the reader to Part I [DH] for further details regarding the background of (1.1).

Here we rigorously prove the existence of such a concentration phenomenon and obtain exact formulas for the determination of x_* and the total biomass. In doing so, we reveal the existence of four critical values $v_{**} < v_* < v^* < v^{**}$ for v_0 (the nutrient level at the sediment) such that

- (i) $x_* = 0$ when $v_0 \geq v^*$, $x_* \in (0, 1)$ when $v_0 \in (v_*, v^*)$, and $x_* = 1$ when $v_0 \leq v_*$;
- (ii) the total biomass increases with v_0 in the range $v_{**} < v_0 < v^{**}$, but stays constant for $v_0 \geq v^{**}$ or $v_0 \leq v_{**}$ (but with v_0 above a certain level so that the biomass can survive).

In order to give a more detailed description of these results, we first recall the main results of Part I [DH], where we proved the following two theorems.

THEOREM 1.1. *There exist $0 < m_* \leq m^* < \infty$ such that (1.1) has a positive solution for $m \in (0, m_*)$ and has no positive solution for $m > m^*$.*

The values of m_* and m^* depend on the parameters in (1.1). To stress their dependence on σ , we write $m_* = m_*(\sigma)$, $m^* = m^*(\sigma)$.

THEOREM 1.2.

$$\lim_{\sigma \rightarrow \infty} m_*(\sigma) = \lim_{\sigma \rightarrow \infty} m^*(\sigma) = f(\min\{\alpha v_0, w_0\}).$$

To investigate the limiting profile of the positive solutions of (1.1) as $\sigma \rightarrow \infty$, we will fix m such that $0 < m < f(\min\{\alpha v_0, w_0\})$ and let σ_n be an increasing sequence of positive numbers converging to ∞ . By Theorems 1.1 and 1.2, for all large n , (1.1) with $\sigma = \sigma_n$ has at least one positive solution. Suppose that (u_n, v_n) is such a solution. We will analyze the behavior of (u_n, v_n) as $n \rightarrow \infty$. This will be done in the following two sections.

In section 2, we find all the possible limiting profiles that a subsequence of $\{(u_n, v_n)\}$ can have; in particular, we find the limiting equations governing these possible limiting profiles. More precisely, let $x_n \in [0, 1]$ be uniquely determined by

$$\begin{cases} \min\{\alpha v_n(x), w_n(x)\} = \alpha v_n(x) & \text{for } x \in [0, x_n), \\ \min\{\alpha v_n(x), w_n(x)\} = w_n(x) & \text{for } x \in (x_n, 1], \end{cases}$$

where

$$w_n(x) = w_0 \exp \left[-A_0 x - A \int_0^x u_n(s) ds \right].$$

We consider the following possibilities:

$$(i) \ x_n \rightarrow x_* \in (0, 1), \quad (ii) \ x_n \rightarrow 0, \quad (iii) \ x_n \rightarrow 1.$$

The cases (ii) and (iii) are each further divided into two subcases, namely, for case (ii),

$$(a1) \ \sigma_n^{1/2} x_n \rightarrow \infty, \quad (a2) \ \sigma_n^{1/2} x_n \rightarrow a_* \in [0, \infty);$$

for case (iii),

$$(b1) \ \sigma_n^{1/2} (1 - x_n) \rightarrow \infty, \quad (b2) \ \sigma_n^{1/2} (1 - x_n) \rightarrow b_* \in [0, \infty).$$

One easily sees that, subject to a subsequence, the above are all the possible behaviors of the sequence $\{x_n\}$. Eventually we will show in section 3 that the limit of the entire sequence $\{x_n\}$ always exists and that this limit is completely determined by the value of v_0 , which in turn allows us to completely determine the profiles of u_n and v_n for large n . But in order to prove these facts, we need to first find all the possible limiting profiles of $\{(u_n, v_n)\}$ and the limiting equations that govern these profiles for each of the above listed cases. The main results of section 2 are summarized below.

If case (i) occurs, we show (see Lemma 2.3) that as $n \rightarrow \infty$, subject to a subsequence, $u_n \rightarrow 0$ uniformly in $[0, x_* - \epsilon] \cup [x_* + \epsilon, 1]$ for any small $\epsilon > 0$ and

$$\int_0^1 u_n(x) dx \rightarrow \tau_* C_0, \quad C_0 := \int_{-\infty}^{\infty} e^{-x^2/(2\delta d_1)} dx = \sqrt{2\delta d_1 \pi},$$

$$v_n \rightarrow v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - \max\{x, x_*\})$$

uniformly in $[0, 1]$, where $x_* \in (0, 1)$ and $\tau_* > 0$ are determined by

$$(1.2) \quad \begin{cases} w_0 e^{-A_0 x_* - A \tau_* (C_0/2)} = \alpha \left[v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - x_*) \right], \\ m = \int_0^1 f(w_0 e^{-A_0 x_* - A \tau_* \max\{C_0/2, C_0 y\}}) dy. \end{cases}$$

If case (ii)(a1) occurs, we show (see Lemma 2.4) that the above conclusions hold with $x_* = 0$; in particular,

$$(1.3) \quad w_0 e^{-A \tau_* (C_0/2)} = \alpha \left[v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1}) \right],$$

and

$$(1.4) \quad m = \int_0^1 f(w_0 e^{-A \tau_* \max\{C_0/2, C_0 y\}}) dy.$$

Since (1.4) uniquely determines $\tau_* > 0$, we can substitute this τ_* into (1.3) to obtain a special value for v_0 , say $v_0 = v_*$.

Similarly, if case (iii)(b1) occurs, we can show (Lemma 2.5) that the conclusions of case (i) hold except that $x_* = 1$; in particular,

$$(1.5) \quad w_0 e^{-A_0 - A\tau_*(C_0/2)} = \alpha \left[v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - 1) \right],$$

and

$$(1.6) \quad m = \int_0^1 f(w_0 e^{-A_0 - A\tau_* \max\{C_0/2, C_0 y\}}) dy.$$

Analogously, $\tau_* > 0$ is uniquely determined by (1.6), and one can then use (1.5) to obtain a special value for v_0 , say $v_0 = v^*$.

If case (ii)(a2) occurs, we show (Lemma 2.4) that as $n \rightarrow \infty$, subject to a subsequence, $u_n \rightarrow 0$ uniformly in $[\epsilon, 1]$ for any small $\epsilon > 0$,

$$\int_0^1 u_n(x) dx \rightarrow \tau_* C(a_*), \quad C(a_*) = \int_{-a_*}^\infty e^{-x^2/(2\delta d_1)} dx,$$

$$v_n \rightarrow v_0 - \frac{\tau_*}{d_2} m C(a_*) (1 + \beta^{-1} - x)$$

uniformly in $[0, 1]$, where $a_* \in [0, \infty)$ and $\tau_* > 0$ are determined by

$$(1.7) \quad m = \int_0^1 f(w_0 e^{-A\tau_* \max\{C(a_*) - C_0/2, C(a_*)y\}}) dy,$$

and

$$(1.8) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2} m C(a_*) (1 + \beta^{-1}) \right) = w_0 e^{-A\tau_* [C(a_*) - C_0/2]} \quad \text{if } a_* > 0,$$

$$(1.9) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2} m \left(\frac{C_0}{2} \right) (1 + \beta^{-1}) \right) \geq w_0 \quad \text{if } a_* = 0.$$

If case (iii)(b2) occurs, we show (Lemma 2.5) that as $n \rightarrow \infty$, subject to a subsequence, $u_n \rightarrow 0$ uniformly in $[0, 1 - \epsilon]$ for any small $\epsilon > 0$,

$$\int_0^1 u_n(x) dx \rightarrow \tau_* \tilde{C}(b_*), \quad \tilde{C}(b_*) = \int_{-\infty}^{b_*} e^{-x^2/(2\delta d_1)} dx = C(-b_*),$$

and

$$v_n \rightarrow v_0 - \frac{\tau_*}{d_2} \beta^{-1} \tilde{C}(b_*),$$

where $b_* \in [0, \infty)$ and $\tau_* > 0$ are determined by

$$(1.10) \quad m = \int_0^1 f(w_0 e^{-A_0 - A\tau_* \max\{C_0/2, \tilde{C}(b_*)y\}}) dy$$

and

$$(1.11) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2 \beta} C(b_*) \right) = w_0 e^{-A_0 - A\tau_* C_0/2} \quad \text{if } b_* > 0,$$

$$(1.12) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2 \beta} \left(\frac{C_0}{2} \right) \right) \leq w_0 e^{-A_0 - A\tau_* C_0/2} \quad \text{if } b_* = 0.$$

In section 3, through careful analysis of the limiting equations (1.2)–(1.12), we show that the entire sequence $\{x_n\}$ always converges to a point $x_* \in [0, 1]$, that exactly one of the cases considered in section 2 occurs, and that in each case the limit of the entire sequence in the conclusion exists. More precisely, if $v_* < v_0 < v^*$ (recall that v_* and v^* are defined above in cases (ii)(a1) and (iii)(b1), respectively), then case (i) must occur, and (1.2) uniquely determines x_* and τ_* . If $v_0 = v_*$, then case (ii)(a1) occurs; if $v_0 = v^*$, then case (iii)(b1) occurs. If $v_0 > v^*$, then case (ii)(a2) must happen, and if $v_0 < v_*$, then case (iii)(b2) must happen. Moreover, our analysis on the limiting total biomass $\lim_{n \rightarrow \infty} \int_0^1 u_n(x) dx$ reveals two further critical values of v_0 , $v_{**} < v_*$ and $v^{**} > v^*$ such that this limiting total biomass is strictly increasing with v_0 for v_0 in the range $v_{**} \leq v_0 \leq v^{**}$ but remains constant (i.e., no longer changes with v_0) when $v_0 \geq v^{**}$ or when $v_0 \leq v_{**}$. See Theorems 3.1–3.3 for more accurate descriptions of these results.

In section 4 we give biological interpretations of our main results and compare our rigorous limiting equations with the game theoretical model of [KL].

Though the proofs are rather involved, they consist mainly of elementary mathematical analysis; most of the proofs in section 2 and all of the arguments in section 3 can be understood with sound knowledge of calculus and real analysis.

2. The limiting equations. We will keep using the notation of Part I [DH]. It turns out that the techniques used in the proof of Theorem 3.1 in Part I are not quite suitable for our purpose here. We will introduce some different techniques.

Suppose that $0 < m < f(\min\{\alpha v_0, w_0\})$ and $\sigma_n, (u_n, v_n)$ are as given in the introduction above. Suppose $c_{v_n, w_n}(x) = C_{x_n}(x)$, $x_n \in [0, 1]$. By passing to a subsequence we may assume that $x_n \rightarrow x_* \in [0, 1]$. Then

$$C_{x_n} = \frac{x - x_n}{\delta + |x - x_n|} \rightarrow C_{x_*}$$

in $C^1([0, 1])$.

In order to obtain useful equations to determine the profiles of u_n and v_n , we need to stretch the variable x appropriately. We define

$$\Phi_n(x) = \exp \left[-\frac{\sigma_n}{2d_1} \int_{x_n}^x C_{x_n}(s) ds \right]$$

and

$$\Psi_n(x) = u_n(x) / \Phi_n(x).$$

By a direct computation we obtain

$$\begin{cases} -d_1 \Psi_n'' + \sigma_n \Gamma_n(x) \Psi_n = [f(\min\{\alpha v_n, w_n\}) - m] \Psi_n, & x \in (0, 1), \\ d_1 \Psi_n' + (\sigma_n/2) C_{x_n} \Psi_n = 0, & x = 0, 1, \end{cases}$$

where

$$\Gamma_n(x) := \frac{\sigma_n(x - x_n)^2 - 2d_1\delta}{4d_1(\delta + |x - x_n|)^2}.$$

Let us introduce the stretched variable $y = \sigma_n^{1/2}(x - x_n)$ and define

$$V_n(y) := \Psi_n(\sigma_n^{-1/2}y + x_n), \quad C_n(y) := \sigma_n^{1/2} C_{x_n}(\sigma_n^{-1/2}y + x_n) = \frac{y}{\delta + \sigma_n^{-1/2}|y|},$$

$$a_n := -\sigma_n^{1/2}x_n, \quad b_n := \sigma_n^{1/2}(1 - x_n),$$

and

$$F_n(y) := f(\min\{\alpha v_n(\sigma_n^{-1/2}y + x_n), w_n(\sigma_n^{-1/2}y + x_n)\}).$$

Then

$$(2.1) \quad \begin{cases} -d_1 V_n'' + \frac{y^2 - 2d_1\delta}{4d_1(\delta + \sigma_n^{-1/2}|y|)^2} V_n = \sigma_n^{-1} [F_n(y) - m] V_n, & y \in (a_n, b_n), \\ d_1 V_n' + (1/2)C_n V_n = 0, & y = a_n, b_n. \end{cases}$$

In the discussions below, we will consider the cases $x_* \in (0, 1)$, $x_* = 0$, and $x_* = 1$ separately.

LEMMA 2.1. *Suppose $x_n \rightarrow x_* \in (0, 1)$, and set $\tilde{V}_n(y) = V_n(y)/\|V_n\|_{L^\infty([a_n, b_n])}$. Then*

$$\tilde{V}_n \rightarrow V_0 \quad \text{in } C^1(J) \text{ for any finite interval } J \subset (-\infty, \infty),$$

where $V_0(y) = \exp[-\frac{y^2}{4d_1\delta}]$ is the unique solution of

$$-d_1 V'' = \frac{2d_1\delta - y^2}{4d_1\delta^2} V, \quad 0 < V \leq 1, \quad V(0) = 1, \quad V'(0) = 0.$$

Proof. Since $x_* \in (0, 1)$, we have $a_n \rightarrow -\infty$ and $b_n \rightarrow \infty$ as $n \rightarrow \infty$. Let us note that, for $y \in [a_n, -(2d_1\delta)^{1/2} - \epsilon]$ with $\epsilon > 0$ sufficiently small and all large n , the first equation in (2.1) implies that $V_n''(y) > 0$. Since $d_1 V_n'(a_n) = -(1/2)C_n(a_n)V_n(a_n) \geq 0$, we deduce that $V_n'(y) > 0$ in $(a_n, -(2d_1\delta)^{1/2} - \epsilon]$ for all large n . Hence V_n is increasing in this range. Similarly, we can see that $V_n(y)$ is decreasing in the range $y \in [(2d_1\delta)^{1/2} + \epsilon, b_n]$ for all large n . Therefore $\max V_n = V_n(y_n)$ for some $y_n \in [-(2d_1\delta)^{1/2} - \epsilon, (2d_1\delta)^{1/2} + \epsilon]$, and $\tilde{V}_n(y) = V_n(y)/V_n(y_n)$. We may assume that $y_n \rightarrow y^*$ as $n \rightarrow \infty$. We now define

$$\tilde{F}_n(y) := \frac{2d_1\delta - y^2}{4d_1(\delta + \sigma_n^{-1/2}|y|)^2} + \sigma_n^{-1} [F_n(y) - m].$$

Then $\tilde{V}_n(y_n) = 1$, and

$$(2.2) \quad \begin{cases} -d_1 \tilde{V}_n'' = \tilde{F}_n \tilde{V}_n, & 0 < \tilde{V}_n \leq 1, \quad y \in (a_n, b_n), \\ d_1 \tilde{V}_n' + (1/2)C_n \tilde{V}_n = 0, & y = a_n, b_n. \end{cases}$$

Since $\{\tilde{F}_n\}$ is uniformly bounded over any bounded interval and $0 \leq \tilde{V}_n \leq 1$, we may apply the interior L^p theory (see [GT]) to (2.2) and use the Sobolev imbedding theorem and a standard diagonal argument to conclude that, by passing to a subsequence, $\tilde{V}_n \rightarrow \tilde{V}$ in $C^1(J)$ for any bounded interval J , and \tilde{V} satisfies

$$(2.3) \quad -d_1 \tilde{V}'' = \frac{2d_1\delta - y^2}{4d_1\delta^2} \tilde{V}, \quad 0 < \tilde{V} \leq 1 \text{ in } (-\infty, \infty), \quad \tilde{V}(y^*) = 1, \quad \tilde{V}'(y^*) = 0.$$

By the monotonicity property of $V_n(y)$ observed earlier, we know that $\tilde{V}(y)$ is nondecreasing in $(-\infty, -(2d_1\delta)^{1/2})$ and is nonincreasing in $((2d_1\delta)^{1/2}, \infty)$. We can now use (2.3) to conclude that $\tilde{V}'(y)$ is positive and increasing in $(-\infty, -(2d_1\delta)^{1/2})$, reaching a positive maximum at $y = -(2d_1\delta)^{1/2}$; then is decreasing in $(-(2d_1\delta)^{1/2},$

$(2d_1\delta)^{1/2}$), reaching a negative minimum at $y = (2d_1\delta)^{1/2}$; and for $y > (2d_1\delta)^{1/2}$, is increasing and stays negative. Therefore $V'(y)$ has a unique zero at some $y_0 \in (-(2d_1\delta)^{1/2}, (2d_1\delta)^{1/2})$, which is the unique maximum point of \tilde{V} . Thus $y_0 = y^*$. In other words, $\tilde{V}(y)$ is increasing in $(-\infty, y^*)$ and is decreasing in $(y^*, 0)$. It then follows from an elementary analysis that \tilde{V} decays to 0 as $|y| \rightarrow \infty$, and there exists $C_1, C_2 > 0$ such that

$$\tilde{V}(y), |\tilde{V}'(y)| \leq C_1 e^{-C_2|y|} \quad \forall y \in (-\infty, \infty).$$

We now multiply $\tilde{V}(-y)$ to (2.3), integrate over $[y^*, \infty)$, and then apply integration by parts. Since $\tilde{V}(-y)$ satisfies the differential equation in (2.3), we deduce

$$\tilde{V}'(-y^*)\tilde{V}(y^*) + \tilde{V}'(y^*)\tilde{V}(-y^*) = 0.$$

It follows that $\tilde{V}'(-y^*) = 0$. Since y^* is the only zero of \tilde{V}' , we must have $y^* = -y^*$, that is, $y^* = 0$. By the uniqueness theorem of initial value problems of ordinary differential equations, we must have $\tilde{V} = V_0$, the unique solution of (2.3) with $y^* = 0$. A simple calculation confirms that the function $\exp[-\frac{y^2}{4d_1\delta}]$ solves the equation for V_0 . Hence, by uniqueness,

$$V_0(y) = \exp\left[-\frac{y^2}{4d_1\delta}\right].$$

Since V_0 is uniquely determined, it follows that the entire original sequence $\{\tilde{V}_n\}$ converges to V_0 . \square

Using the monotonicity of \tilde{V}_n and the fact that $V_0(y) \rightarrow 0$ as $|y| \rightarrow \infty$, we see that Lemma 2.1 implies

$$(2.4) \quad \|\Psi_n(\cdot)/\|\Psi_n\|_\infty - V_0(\sigma_n^{1/2}(\cdot - x_n))\|_{L^\infty([0,1])} \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

since for large n the function

$$\Psi_n(x)/\|\Psi_n\|_\infty - V_0(\sigma_n^{1/2}(x - x_n))$$

is uniformly small at those values of $x \in [0, 1]$ such that $\sigma_n^{1/2}(x - x_n)$ stays bounded (by Lemma 2.1), and, by the properties of \tilde{V}_n and V_0 , the values of the function at the remaining $x \in [0, 1]$ are also small.

We now denote $\tilde{\Psi}_n(x) = \Psi_n(x)/\|\Psi_n\|_\infty$ and consider the function

$$\tilde{u}_n(x) := \sigma_n^{1/2}\Phi_n(x)\tilde{\Psi}_n(x) = \left(\frac{\sigma_n^{1/2}}{\|\Psi_n\|_\infty}\right)u_n.$$

We will show that, for large n , \tilde{u}_n behaves like the δ -function concentrating at x_* . Indeed, we have the following result.

LEMMA 2.2. *For any given small $\epsilon > 0$, $|x - x_n| \geq \epsilon$ implies*

$$(2.5) \quad 0 < \tilde{u}_n(x) \leq \sigma_n^{1/2} \exp\left[-\frac{\sigma_n}{4(\delta + 1)d_1}\epsilon^2\right] \rightarrow 0.$$

Moreover, when $x_n \rightarrow x_* \in (0, 1)$,

$$(2.6) \quad \lim_{n \rightarrow \infty} \int_0^1 \tilde{u}_n(x)dx = C_0 := \int_{-\infty}^\infty e^{-x^2/(2\delta d_1)}dx = \sqrt{2\delta d_1\pi}.$$

Proof. For any given small $\epsilon > 0$, there exists $\delta_0 = \delta_0(\epsilon) > 0$ small so that, when $|x - x_n| \leq \delta_0$,

$$\exp \left[-\frac{\sigma_n}{4\delta d_1}(x - x_n)^2 \right] \leq \Phi_n(x) \leq \exp \left[-\frac{\sigma_n(1 - \epsilon)}{4\delta d_1}(x - x_n)^2 \right].$$

For any $x \in [0, 1]$, we have

$$\exp \left[-\frac{\sigma_n}{4\delta d_1}(x - x_n)^2 \right] \leq \Phi_n(x) \leq \exp \left[-\frac{\sigma_n}{4(\delta + 1)d_1}(x - x_n)^2 \right].$$

Since $\tilde{\Psi}_n \leq 1$, for $|x - x_n| \geq \epsilon$, we have

$$\tilde{u}_n(x) \leq \sigma_n^{1/2} \exp \left[-\frac{\sigma_n}{4(\delta + 1)d_1}\epsilon^2 \right] \rightarrow 0.$$

This proves (2.5). Moreover, we have

$$\begin{aligned} \int_0^1 \tilde{u}_n(x) dx &= \int_{x_n - \epsilon}^{x_n + \epsilon} \sigma_n^{1/2} \Phi_n(x) \tilde{\Psi}_n(x) dx + o(1) \\ &= \int_{-\epsilon \sigma_n^{1/2}}^{\epsilon \sigma_n^{1/2}} \Phi_n(x_n + \sigma_n^{-1/2} y) \tilde{V}_n(y) dy + o(1) \\ &= \int_{-\infty}^{\infty} \exp \left[-\frac{y^2}{4d_1\delta} \right] V_0(y) dy + o(1) \\ &= \int_{-\infty}^{\infty} \exp \left[-\frac{y^2}{2d_1\delta} \right] dy + o(1). \end{aligned}$$

Hence (2.6) holds. For later application, let us also note from the above argument that

$$(2.7) \quad \lim_{n \rightarrow \infty} \int_{x_n}^1 \tilde{u}_n(x) dx = \lim_{n \rightarrow \infty} \int_0^{x_n} \tilde{u}_n(x) dx = C_0/2. \quad \square$$

Denote $\tau_n := \|\Psi_n\|_{\infty} \sigma_n^{-1/2}$. We find that

$$u_n(x) = \tau_n \tilde{u}_n(x).$$

LEMMA 2.3. *Suppose that $x_n \rightarrow x_* \in (0, 1)$. Then $\{\tau_n\}$ has a subsequence, still denoted by itself, such that $\tau_n \rightarrow \tau_* > 0$. Moreover, τ_* and x_* must satisfy*

$$(2.8) \quad w_0 e^{-A_0 x_* - A \tau_* (C_0/2)} = \alpha \left[v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - x_*) \right]$$

and

$$(2.9) \quad m = \int_0^1 f(w_0 e^{-A_0 x_* - A \tau_* \max\{C_0/2, C_0 y\}}) dy.$$

Furthermore, by possibly passing to a further subsequence, $u_n \rightarrow 0$ in $C([0, 1] \setminus [x_* - \epsilon, x_* + \epsilon])$, for all $\epsilon > 0$, and

$$(2.10) \quad v_n(x) \rightarrow v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - \max\{x, x_*\})$$

uniformly in $[0, 1]$.

Proof. By passing to a subsequence, we have two possible cases:

$$(i) \tau_n \rightarrow \infty, \quad (ii) \tau_n \rightarrow \tau_* \in [0, \infty).$$

Step 1. Case (i) cannot happen.

Suppose $\tau_n \rightarrow \infty$; we are going to derive a contradiction. Denote

$$f_n = f(\min\{\alpha v_n, w_n\}).$$

Since

$$w_n(x_n) \leq w_0 e^{-A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds},$$

and by (2.7)

$$\int_0^{x_n} \tilde{u}_n(s) ds \rightarrow C_0/2 > 0,$$

we easily see that $w_n(x_n) \rightarrow 0$. It follows that

$$\|f_n\|_\infty = f_n(x_n) = f(w_n(x_n)) \rightarrow 0.$$

This implies that

$$\int_0^1 f_n \tilde{u}_n dx \rightarrow 0.$$

On the other hand, we may integrate the equation for u_n to obtain

$$\int_0^1 [f_n(x) - m] u_n dx = 0,$$

which implies that

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0.$$

Letting $n \rightarrow \infty$ and using (2.6), we obtain

$$mC_0 = \lim_{n \rightarrow \infty} \int_0^1 f_n \tilde{u}_n dx = 0,$$

which contradicts our assumption that $m > 0$. Therefore case (i) cannot happen.

Step 2. The limiting profile of u_n and v_n .

We next consider case (ii), namely, $\tau_n \rightarrow \tau_* \in [0, \infty)$. In this case, due to (2.5), $u_n = \tau_n \tilde{u}_n \rightarrow 0$ in $C([0, 1] \setminus [x_* - \epsilon, x_* + \epsilon])$, for all $\epsilon > 0$, and hence

$$(2.11) \quad \tau_n f_n \tilde{u}_n \rightarrow 0 \quad \text{uniformly in } [0, x_* - \epsilon] \cup [x_* + \epsilon, 1] \quad \forall \epsilon > 0.$$

Let $\zeta_n = v_0 - v_n$. Then

$$(2.12) \quad -d_2 \zeta_n'' = \tau_n f_n \tilde{u}_n \text{ in } (0, 1), \quad \zeta_n'(0) = 0, \quad \zeta_n'(1) + \beta \zeta_n(1) = 0.$$

Since $v_n \geq 0$, we have $\zeta_n \leq v_0$. Since $\tau_n f_n \tilde{u}_n > 0$, from (2.12) and the maximum principle, we deduce that $\zeta_n > 0$. Hence we always have $0 < \zeta_n \leq v_0$. Therefore we can integrate (2.12) to obtain

$$\eta_n := \tau_n \int_0^1 f_n \tilde{u}_n dx = d_2[\zeta'_n(0) - \zeta'_n(1)] = d_2\beta\zeta_n(1) \in [0, d_2\beta v_0].$$

This implies that, by passing to a subsequence, we may assume that $\eta_n \rightarrow \eta_* \in [0, d_2\beta v_0]$.

Moreover, using (2.11), (2.12), and $\eta_n \rightarrow \eta_*$, we find that

$$\begin{aligned} \{\zeta'_n\} &\text{ is a bounded sequence in } L^\infty([0, 1]), \\ \zeta'_n(x) &\rightarrow 0 \text{ uniformly in } [0, x_* - \epsilon] \quad \forall \epsilon > 0, \\ \zeta'_n(x) &\rightarrow -\eta_*/d_2 \text{ uniformly in } [x_* + \epsilon, 1] \quad \forall \epsilon > 0. \end{aligned}$$

Since, moreover, $0 \leq \zeta_n \leq v_0$, we conclude that $\{\zeta_n\}$ is precompact in $C([0, 1])$. Hence, by passing to a subsequence, we may assume that $\zeta_n \rightarrow \zeta$ in $C([0, 1])$.

On the other hand, we may apply the L^p theory to (2.12) and the Sobolev imbedding theorem to find a further subsequence, still denoted by ζ_n , such that $\zeta_n \rightarrow \tilde{\zeta}$ in $C^1(J)$ for any compact interval $J \subset [0, x_*) \cup (x_*, 1]$, and $\tilde{\zeta}$ satisfies (in the weak sense)

$$-d_2\tilde{\zeta}'' = 0 \text{ in } [0, x_*) \cup (x_*, 1], \quad \tilde{\zeta}'(0) = 0, \quad \tilde{\zeta}'(1) + \beta\tilde{\zeta}(1) = 0.$$

Clearly we must have $\tilde{\zeta} = \zeta$. Moreover, our earlier analysis on ζ_n implies that $\zeta'(x) = 0$ in $[0, x_*)$ and $\zeta'(x) = -\eta_*/d_2$ in $(x_*, 1]$. These properties uniquely determine ζ :

$$(2.13) \quad \zeta(x) = (\eta_*/d_2)(1 + \beta^{-1} - \max\{x_*, x\}).$$

Step 3. $\tau_* > 0$.

Otherwise, $\tau_* = 0$ and hence $\eta_* = 0$. It follows that $\zeta = 0$ and $v_n \rightarrow v_0$ uniformly in $[0, 1]$, and that

$$w_n(x) = w_0 e^{-A_0 x} e^{-A\tau_n \int_0^x \tilde{u}_n(s) ds} \rightarrow w_0 e^{-A_0 x} = w_*(x)$$

uniformly in $[0, 1]$. This implies that

$$x_* = x_0^* \quad \text{and} \quad f_n(x) \rightarrow f_0(x) := f(\min\{\alpha v_0, w_*\}) \text{ uniformly in } [0, 1].$$

We may now integrate the equation for u_n to obtain, as before,

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0.$$

Letting $n \rightarrow \infty$, we deduce

$$[f_0(x_0^*) - m]C_0 = 0,$$

which contradicts our assumption that $m < f(\min\{\alpha v_0, w_0\}) = f_0(x_0^*)$. Hence $\tau_* > 0$.

Step 4. The equations for x_* and τ_* .

We now set out to find the equations that determine x_* and τ_* . By (2.7),

$$w_n(x_n) = w_0 e^{-A_0 x_n} e^{-A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds} \rightarrow w_0 e^{-A_0 x_*} e^{-A\tau_*(C_0/2)}.$$

On the other hand,

$$w_n(x_n) = \alpha v_n(x_n) \rightarrow \alpha [v_0 - \zeta(x_*)].$$

Thus we necessarily have

$$(2.14) \quad w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} = \alpha [v_0 - \zeta(x_*)] = \alpha [v_0 - (\eta_*/d_2)(1 + \beta^{-1} - x_*)].$$

Moreover, using (2.5), (2.7), and the fact that $\alpha v_n \rightarrow \alpha(v_0 - \zeta)$ uniformly in $[0, 1]$, we deduce

$$(2.15) \quad \int_0^{x_n} f(\alpha v_n) \tilde{u}_n dx \rightarrow (C_0/2) f(\alpha v_0 - \alpha \zeta(x_*)).$$

Using

$$w_n(x) = w_0 e^{-A_0 x} e^{-A\tau_n \int_0^x \tilde{u}_n(s) ds}$$

and the property of \tilde{u}_n , we obtain, for any small $\epsilon > 0$,

$$\begin{aligned} & \int_{x_n}^1 f(w_n) \tilde{u}_n dx \\ &= \int_{x_n}^1 f(w_0 e^{-A_0 x - A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds - A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx \\ &= \int_{x_n}^{x_* + \epsilon} f(w_0 e^{-A_0 x - A\tau_n(C_0/2)} e^{-A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= [1 + o_\epsilon(1)] \int_{x_n}^{x_* + \epsilon} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= [1 + o_\epsilon(1)] \int_0^{[\int_{x_n}^1 \tilde{u}_n(s) ds]} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_n y}) dy + o(1) \\ &= [1 + o_\epsilon(1)] \int_0^{C_0/2} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_* y}) dy + o(1), \end{aligned}$$

where $o_\epsilon(1)$ represents a quantity that converges to 0 as $\epsilon \rightarrow 0$.

Thus

$$(2.16) \quad \int_{x_n}^1 f(w_n) \tilde{u}_n(x) dx \rightarrow \int_0^{C_0/2} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_* y}) dy$$

as $n \rightarrow \infty$.

Combining (2.15) and (2.16), we obtain

$$(2.17) \quad \begin{aligned} \eta_* &= \lim_{n \rightarrow \infty} \tau_n \int_0^1 f_n \tilde{u}_n dx \\ &= \tau_* \left[(C_0/2) f(\alpha v_0 - \alpha \zeta(x_*)) + \int_0^{C_0/2} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_* y}) dy \right]. \end{aligned}$$

Moreover, we may integrate the equation for u_n to obtain

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0.$$

Letting $n \rightarrow \infty$ and using (2.15), (2.16), we obtain

$$mC_0 = (C_0/2)f(\alpha v_0 - \alpha\zeta(x_*)) + \int_0^{C_0/2} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_* y}) dy.$$

This combined with (2.17) yields

$$(2.18) \quad \eta_* = \tau_* m C_0$$

and combined with (2.14) gives

$$\begin{aligned} m &= (1/2)f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)}) + C_0^{-1} \int_0^{C_0/2} f(w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} e^{-A\tau_* y}) dy \\ &= C_0^{-1} \int_0^{C_0} f(w_0 e^{-A_0 x_* - A\tau_* \max\{C_0/2, y\}}) dy \\ &= \int_0^1 f(w_0 e^{-A_0 x_* - A\tau_* \max\{C_0/2, C_0 y\}}) dy; \end{aligned}$$

thus (2.9) is proved. Equation (2.8) and (2.10) clearly follow from (2.13), (2.14), and (2.18). \square

We now consider the case $x_* = 0$. By passing to a subsequence, we have two subcases:

$$(a1) \ a_n := \sigma_n^{1/2} x_n \rightarrow \infty, \quad (a2) \ a_n \rightarrow a_* \in [0, \infty).$$

LEMMA 2.4. *In subcase (a1), all of the conclusions in Lemmas 2.2 and 2.3 hold. In subcase (a2), $\{\tau_n\}$ has a subsequence, still denoted by itself, such that $\tau_n \rightarrow \tau_* > 0$. Moreover, τ_* and a_* must satisfy*

$$(2.19) \quad m = \int_0^1 f(w_0 e^{-A\tau_* \max\{C(a_*) - C_0/2, C(a_*)y\}}) dy$$

and

$$(2.20) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2} m C(a_*) (1 + \beta^{-1}) \right) = w_0 e^{-A\tau_* [C(a_*) - C_0/2]} \quad \text{if } a_* > 0,$$

$$(2.21) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2} m \left(\frac{C_0}{2} \right) (1 + \beta^{-1}) \right) \geq w_0 \quad \text{if } a_* = 0,$$

where

$$C(a_*) := \int_{-a_*}^\infty \exp \left[-\frac{y^2}{2d_1 \delta} \right] dy.$$

Furthermore, by possibly passing to a further subsequence, $u_n \rightarrow 0$ in $C([\epsilon, 1])$, for all $\epsilon \in (0, 1)$,

$$(2.22) \quad \lim_{n \rightarrow \infty} \int_0^{x_n} \tilde{u}_n(x) dx = C(a_*) - C_0/2, \quad \lim_{n \rightarrow \infty} \int_0^1 \tilde{u}_n(x) dx = C(a_*),$$

and

$$(2.23) \quad v_n(x) \rightarrow v_0 - \frac{\tau_*}{d_2} m C(a_*) (1 + \beta^{-1} - x)$$

uniformly in $[0, 1]$.

Proof. In subcase (a1), we may repeat the arguments used for the case $x_* \in (0, 1)$ above to see that all the conclusions there (with x_* replaced by 0) remain valid; the proofs carry over with minor modifications.

Consider now subcase (a2). In this case, we may use interior and boundary L^p estimates and the Sobolev imbedding theorem to conclude that, by passing to a subsequence, $\|\tilde{V}_n - \tilde{V}\|_{C^1([a_n, M])} \rightarrow 0$ for all $M > 0$, where \tilde{V} satisfies, instead of (2.3),

$$(2.24) \quad \begin{cases} -d_1 \tilde{V}'' = \frac{2d_1 \delta - y^2}{4d_1 \delta^2} \tilde{V}, & 0 < \tilde{V} \leq 1 \text{ in } (-a_*, \infty), \\ d_1 \tilde{V}'(-a_*) - \frac{a_*}{2\delta} \tilde{V}(-a_*) = 0, & \tilde{V}(y^*) = 1, \tilde{V}'(y^*) = 0. \end{cases}$$

Note that as before \tilde{V} is decreasing in $[(2d_1 \delta)^{1/2}, \infty)$. This and (2.24) imply that \tilde{V} converges to 0 as $y \rightarrow \infty$. Moreover, an elementary consideration shows that

$$|\tilde{V}'(y)|, \tilde{V}(y) \leq C_1 e^{-C_2 y}$$

for some $C_1, C_2 > 0$, and all $y > 0$.

We will show that $y^* = 0$ and \tilde{V} is again the unique solution of (2.3) with $y^* = 0$, namely V_0 . Since V_0 and $|V_0'|$ are bounded from above by a function of the form $C_1 e^{-C_2 |y|}$, we can multiply the first equation in (2.24) by V_0 , integrate over $[y^*, \infty)$, and use integration by parts to deduce

$$d_1 [\tilde{V} V_0' - \tilde{V}' V_0] \Big|_{y^*}^\infty = 0.$$

It follows that $V_0'(y^*) = 0$, which implies that $y^* = 0$. Therefore, by the uniqueness of initial value problems of the ordinary differential equations, we deduce $\tilde{V} \equiv V_0$. Let us note that a direct calculation shows

$$d_1 V_0'(y) + \frac{y}{2\delta} V_0(y) = 0 \quad \text{for every } y \in (-\infty, \infty).$$

Therefore (2.24) does not introduce any restriction for a_* .

Since now $\sigma_n^{1/2} x_n \rightarrow a_*$, instead of (2.6), we have

$$(2.25) \quad \lim_{n \rightarrow \infty} \int_0^{x_n} \tilde{u}_n(x) dx = C(a_*) - C_0/2, \quad \lim_{n \rightarrow \infty} \int_0^1 \tilde{u}_n(x) dx = C(a_*),$$

where

$$C(a_*) := \int_{-a_*}^\infty \exp\left[-\frac{y^2}{4d_1 \delta}\right] V_0(y) dy = \int_{-a_*}^\infty \exp\left[-\frac{y^2}{2d_1 \delta}\right] dy.$$

We proceed as in the case $x_* \in (0, 1)$ and have two possibilities for τ_n as before. We show that, in the current case, we still cannot have $\tau_n \rightarrow \infty$. Arguing indirectly, we assume that $\tau_n \rightarrow \infty$.

Then in the case $a_* > 0$, we have $C(a_*) - C_0/2 > 0$, and hence

$$w_n(x_n) \leq w_0 e^{-A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds} \rightarrow 0.$$

It follows that

$$\|f_n\|_\infty = f_n(x_n) = f(w_n(x_n)) \rightarrow 0$$

and

$$\int_0^1 f_n \tilde{u}_n dx \rightarrow 0.$$

If $a_* = 0$, then $C(a_*) - C_0/2 = 0$ and

$$\begin{aligned} \int_0^1 f_n(x) \tilde{u}_n(x) dx &= \int_{x_n}^1 f(w_n(x)) \tilde{u}_n(x) dx + o(1) \\ &\leq \int_{x_n}^1 f(w_0 e^{-A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= \int_0^{[\int_{x_n}^1 \tilde{u}_n(s) ds]} f(w_0 e^{-A\tau_n y}) dy + o(1) \\ &\leq \epsilon f(w_0) + \int_\epsilon^{C_0/2} f(w_0 e^{-A\tau_n y}) dy + o(1) \\ &= \epsilon f(w_0) + o(1) \quad \forall \epsilon \in (0, C_0/2). \end{aligned}$$

Therefore we always have

$$\int_0^1 f_n \tilde{u}_n dx \rightarrow 0 \text{ as } n \rightarrow \infty.$$

As before, we may integrate the equation for u_n to obtain

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0.$$

Letting $n \rightarrow \infty$ and using the above estimate, we deduce

$$-mC(a_*) = 0,$$

a contradiction to our assumption that $m > 0$. Therefore we cannot have $\tau_n \rightarrow \infty$.

Thus we can only have the case $\tau_n \rightarrow \tau_*$. Then much as before we deduce $u_n \rightarrow 0$ in $C([\epsilon, 1])$ for all $\epsilon \in (0, 1)$, and

$$\zeta_n \rightarrow \zeta := (\eta_*/d_2)(1 + \beta^{-1} - x)$$

in $C([0, 1]) \cap C^1([\epsilon, 1])$ for all $\epsilon \in (0, 1)$. If $\tau_* = 0$, we can deduce as before that $m = f_0(x_0^*)$, a contradiction to our initial assumption on m . Therefore $\tau_* > 0$.

If $a_* = 0$, we first choose $y_n \in (x_n, 1)$ such that $y_n \rightarrow 0$ and $\int_{y_n}^1 \tilde{u}_n(x) dx \rightarrow 0$, and then we have

$$\begin{aligned} \int_0^1 f_n(x) \tilde{u}_n(x) dx &= \int_{x_n}^{y_n} f(w_n(x)) \tilde{u}_n(x) dx + o(1) \\ &= \int_{x_n}^{y_n} f(w_0 e^{-A_0 x - A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds - A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= \int_{x_n}^{y_n} f(w_0 e^{-A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= \int_{x_n}^1 f(w_0 e^{-A\tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= \int_0^{C_0/2} f(w_0 e^{-A\tau_* y}) dy + o(1). \end{aligned}$$

If $a_* > 0$, then $x_n > 0$ and $w_n(x_n) = \alpha v_n(x_n)$. From

$$v_n(x_n) \rightarrow v_0 - \zeta(0)$$

and

$$w_n(x_n) = w_0 e^{-A_0 x_n - A \tau_n \int_0^{x_n} \tilde{u}_n dx} \rightarrow w_0 e^{-A \tau_* [C(a_*) - C_0/2]}$$

we obtain

$$\alpha[v_0 - \zeta(0)] = w_0 e^{-A \tau_* [C(a_*) - C_0/2]}.$$

Moreover, similar to the above,

$$\begin{aligned} \int_{x_n}^1 f_n(x) \tilde{u}_n(x) dx &= \int_{x_n}^{y_n} f(w_n(x)) \tilde{u}_n(x) dx + o(1) \\ &= \int_{x_n}^{y_n} f(w_0 e^{-A \tau_n \int_0^{x_n} \tilde{u}_n(s) ds - A \tau_n \int_{x_n}^x \tilde{u}_n(s) ds}) \tilde{u}_n dx + o(1) \\ &= \int_0^{C_0/2} f(w_0 e^{-A \tau_* [C(a_*) - C_0/2] - A \tau_* y}) dy + o(1), \end{aligned}$$

and

$$\begin{aligned} \int_0^{x_n} f_n(x) \tilde{u}_n(x) dx &= \int_0^{x_n} f(\alpha v_n(x)) \tilde{u}_n(x) dx \\ &= f(\alpha[v_0 - \zeta(0)]) [C(a_*) - C_0/2] + o(1) \\ &= [C(a_*) - C_0/2] f(w_0 e^{-A \tau_* [C(a_*) - C_0/2]}) + o(1). \end{aligned}$$

Therefore we always have

$$(2.26) \quad \int_0^1 f_n \tilde{u}_n dx \rightarrow \int_0^{C(a_*)} f(w_0 e^{-A \tau_* \max\{[C(a_*) - C_0/2], y\}}) dy.$$

We may now use

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0$$

to obtain

$$mC(a_*) = \int_0^{C(a_*)} f(w_0 e^{-A \tau_* \max\{C(a_*) - C_0/2, y\}}) dy.$$

Therefore

$$m = \int_0^1 f(w_0 e^{-A \tau_* \max\{C(a_*) - C_0/2, C(a_*)y\}}) dy,$$

and (2.19) is proved.

We thus obtain

$$\eta_* = \tau_* \lim_{n \rightarrow \infty} \int_0^1 f_n \tilde{u}_n dx = \tau_* mC(a_*).$$

Therefore,

$$v_n(x) \rightarrow v_0 - \zeta = v_0 - \frac{\tau_*}{d_2} mC(a_*)(1 + \beta^{-1} - x)$$

uniformly in $[0, 1]$; that is, (2.23) holds.

Let us note that (2.22) was already proved in (2.25). So it remains to prove (2.20) and (2.21). If $a_* > 0$, then $x_n > 0$, and we necessarily have $\alpha v_n(x_n) = w_n(x_n)$. Recall that

$$w_n(x_n) \rightarrow w_0 e^{-A\tau_*[C(a_*)-C_0/2]}, \quad v_n(x_n) \rightarrow v_0 - \zeta(0).$$

Hence

$$\alpha \left(v_0 - \frac{\tau_*}{d_2} mC(a_*)(1 + \beta^{-1}) \right) = w_0 e^{-A\tau_*[C(a_*)-C_0/2]}.$$

If $a_* = 0$, then $x_n = 0$ is possible, and so we have $\alpha v_n(x_n) \geq w_n(x_n)$ in general, and instead of the above identity we should have

$$\alpha \left(v_0 - \frac{\tau_*}{d_2} m \left(\frac{C_0}{2} \right) (1 + \beta^{-1}) \right) \geq w_0.$$

Thus (2.20) and (2.21) are established. The proof is now complete. \square

Finally we consider the case $x_* = 1$. By passing to a subsequence, we have two subcases:

$$(b1) \ b_n := \sigma_n^{1/2}(1 - x_n) \rightarrow \infty, \quad (b2) \ b_n \rightarrow b_* \in [0, \infty).$$

LEMMA 2.5. *In subcase (b1), all of the conclusions in Lemmas 2.2 and 2.3 hold. In subcase (b2), $\{\tau_n\}$ has a subsequence, still denoted by itself, such that $\tau_n \rightarrow \tau_* > 0$. Moreover, τ_* and b_* must satisfy*

$$(2.27) \quad m = \int_0^1 f(w_0 e^{-A_0 - A\tau_* \max\{C_0/2, \tilde{C}(b_*)y\}}) dy$$

and

$$(2.28) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2\beta} \tilde{C}(b_*) \right) = w_0 e^{-A_0 - A\tau_* C_0/2} \quad \text{if } b_* > 0,$$

$$(2.29) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2\beta} \left(\frac{C_0}{2} \right) \right) \leq w_0 e^{-A_0 - A\tau_* C_0/2} \quad \text{if } b_* = 0,$$

where

$$\tilde{C}(b_*) := \int_{-\infty}^{b_*} \exp \left[-\frac{y^2}{2d_1\delta} \right] dy = C(-b_*).$$

Furthermore, by possibly passing to a further subsequence, $u_n \rightarrow 0$ in $C([0, 1 - \epsilon])$ for every $\epsilon \in (0, 1)$,

$$(2.30) \quad \lim_{n \rightarrow \infty} \int_{x_n}^1 \tilde{u}_n(x) dx = \tilde{C}(b_*) - C_0/2, \quad \lim_{n \rightarrow \infty} \int_0^1 \tilde{u}_n(x) dx = \tilde{C}(b_*),$$

$$(2.31) \quad v_n(x) \rightarrow v_0 - \zeta = v_0 - \frac{\tau_*}{d_2\beta} \tilde{C}(b_*)$$

uniformly in $[0, 1]$.

Proof. In subcase (b1), we may repeat the arguments used in Lemmas 2.2 and 2.3 for the case $x_* \in (0, 1)$ to see that all the conclusions there (with x_* replaced by 1) remain valid; the proofs need only minor modifications.

We now consider subcase (b2). Then instead of (2.3) we have

$$(2.32) \quad \begin{cases} -d_1 \tilde{V}'' = \frac{2d_1\delta - y^2}{4d_1\delta} \tilde{V}, & 0 < \tilde{V} \leq 1 \text{ in } (-\infty, b_*), \\ \tilde{V}'(b_*) + \frac{b_*}{2\delta} \tilde{V}(b_*) = 0, & \tilde{V}(y^*) = 1, \tilde{V}'(y^*) = 0. \end{cases}$$

Note that as before \tilde{V} is increasing in $(-\infty, -(2d_1\delta)^{1/2}]$. This and (2.32) imply that \tilde{V} converges to 0 as $y \rightarrow -\infty$. Moreover, an elementary consideration shows that

$$|\tilde{V}'(y)|, \tilde{V}(y) \leq C_1 e^{-C_2|y|}$$

for some $C_1, C_2 > 0$, and all $y < 0$.

As in the case for (2.24), we can similarly show that $y^* = 0$ and $\tilde{V} \equiv V_0$, the unique solution of (2.3) with $y^* = 0$. Moreover, (2.32) introduces no restriction for b_* .

Since $\sigma_n^{1/2}(1 - x_n) \rightarrow b_*$, instead of (2.6), we have

$$\lim_{n \rightarrow \infty} \int_{x_n}^1 \tilde{u}_n(x) dx = \tilde{C}(b_*) - C_0/2, \quad \lim_{n \rightarrow \infty} \int_0^1 \tilde{u}_n(x) dx = \tilde{C}(b_*),$$

where

$$\tilde{C}(b_*) := \int_{-\infty}^{b_*} \exp\left[-\frac{y^2}{4d_1\delta}\right] V_0(y) dy = C(-b_*).$$

This establishes (2.30).

We proceed as in the case $x_* \in (0, 1)$ and have two possibilities for τ_n as before. We show that in the current case, we still cannot have $\tau_n \rightarrow \infty$. Arguing indirectly, we assume that $\tau_n \rightarrow \infty$.

Since $\int_0^{x_n} \tilde{u}_n dx \rightarrow C_0/2$, we have

$$w_n(x_n) \leq w_0 e^{-A\tau_n \int_0^{x_n} \tilde{u}_n(s) ds} \rightarrow 0.$$

It follows that

$$\|f_n\|_\infty = f_n(x_n) \leq f(w_n(x_n)) \rightarrow 0,$$

and

$$\int_0^1 f_n \tilde{u}_n dx \rightarrow 0.$$

As before, we may integrate the equation for u_n to obtain

$$\int_0^1 [f_n(x) - m] \tilde{u}_n dx = 0.$$

Letting $n \rightarrow \infty$ and using the above estimate, we deduce

$$-m\tilde{C}(b_*) = 0,$$

a contradiction to our assumption that $m > 0$. Therefore we cannot have $\tau_n \rightarrow \infty$.

Thus we can have only the case $\tau_n \rightarrow \tau_*$. Then much as before we deduce $u_n \rightarrow 0$ in $C([0, 1 - \epsilon])$ for each $\epsilon \in (0, 1)$ and $\zeta_n \rightarrow \zeta$ in $C([0, 1]) \cap C^1([0, 1 - \epsilon])$, for all $\epsilon \in (0, 1)$, with ζ satisfying

$$\zeta'' = 0 \text{ in } [0, 1), \quad \zeta' = 0 \text{ in } [0, 1).$$

Hence ζ is a constant. To determine its value, we use

$$-d_2 \zeta'_n(1) = \int_0^1 \tau_n f_n \tilde{u}_n dx \rightarrow \tau_* \tilde{C}(b_*)$$

and

$$\zeta'_n(1) + \beta \zeta_n(1) = 0$$

to deduce

$$-\frac{\tau_*}{d_2} \tilde{C}(b_*) + \beta \zeta = 0,$$

and hence

$$(2.33) \quad \zeta = \frac{\tau_*}{d_2 \beta} \tilde{C}(b_*).$$

If $\tau_* = 0$, then $\zeta \equiv 0$, and hence $v_n \rightarrow v_0$ uniformly in $[0, 1]$ and

$$w_n(x) = w_0 e^{-A_0 x - A \tau_n \int_0^x \tilde{u}_n dx} \rightarrow w_0 e^{-A_0 x}$$

uniformly in $[0, 1]$. Then we can deduce as before that $m = f_0(x_0^*)$, a contradiction to our initial assumption on m . Therefore $\tau_* > 0$.

We have

$$\begin{aligned} \int_0^{x_n} f_n(x) \tilde{u}_n(x) dx &= \int_0^{x_n} f(\alpha v_n(x)) \tilde{u}_n(x) dx \\ &= (C_0/2) f(\alpha(v_0 - \zeta)) + o(1). \end{aligned}$$

If $b_* = 0$, then

$$\int_{x_n}^1 f_n(x) \tilde{u}_n(x) dx = o(1).$$

If $b_* > 0$, then $x_n > 0$ and $w_n(x_n) = \alpha v_n(x_n)$. From

$$v_n(x_n) \rightarrow v_0 - \zeta = v_0 - \frac{\tau_*}{d_2 \beta} \tilde{C}(b_*)$$

and

$$w_n(x_n) = w_0 e^{-A_0 x_n - A \tau_n \int_0^{x_n} \tilde{u}_n dx} \rightarrow w_0 e^{-A_0 - A \tau_* C_0/2},$$

we obtain

$$(2.34) \quad \alpha \left(v_0 - \frac{\tau_*}{d_2 \beta} \tilde{C}(b_*) \right) = w_0 e^{-A_0 - A \tau_* C_0/2}.$$

Moreover,

$$\begin{aligned} \int_{x_n}^1 f_n(x)\tilde{u}_n(x)dx &= \int_{x_n}^1 f(w_n(x))\tilde{u}_n(x)dx \\ &= \int_{x_n}^1 f(w_0e^{-A_0x-A\tau_n\int_0^{x_n}\tilde{u}_n(s)ds-A\tau_n\int_{x_n}^x\tilde{u}_n(s)ds})\tilde{u}_n dx \\ &= \int_0^{\tilde{C}(b_*)-C_0/2} f(w_0e^{-A_0-A\tau_*C_0/2-A\tau_*y})dy + o(1). \end{aligned}$$

Therefore, whether $b_* = 0$ or $b_* > 0$, we always have

$$(2.35) \quad \int_0^1 f_n(x)\tilde{u}_n(x)dx \rightarrow \int_0^{\tilde{C}(b_*)} f(w_0e^{-A_0-A\tau_*\max\{C_0/2,y\}})dy.$$

We may now use

$$\int_0^1 [f_n(x) - m]\tilde{u}_n dx = 0$$

to obtain

$$m\tilde{C}(b_*) = \int_0^{\tilde{C}(b_*)} f(w_0e^{-A_0-A\tau_*\max\{C_0/2,y\}})dy,$$

which gives (2.27).

Note that if $b_* = 0$, then $x_n = 1$ is possible, and we have only $w_n(x_n) \geq \alpha v(x_n)$, so instead of (2.27), we should have

$$\alpha\left(v_0 - \frac{\tau_*}{d_2\beta}\tilde{C}(b_*)\right) \leq w_0e^{-A_0-A\tau_*C_0/2}.$$

Thus we have established (2.28) and (2.29). Clearly (2.31) follows from (2.33) and the fact that $v_n \rightarrow v_0 - \zeta$ uniformly in $[0, 1]$. The proof is complete. \square

3. Limiting profile of the positive solutions. We are now ready to state and prove our main results. We will show that the limiting equations obtained in the previous section uniquely determine x_* and τ_* , and the value of v_0 determines which set of limiting equations should be used for calculating x_* and τ_* . In this way, the asymptotic behavior of the positive solutions is completely determined.

Let us recall that m is fixed such that

$$(3.1) \quad 0 < m < f(\min\{\alpha v_0, w_0\}),$$

and $\sigma_n \rightarrow \infty$ is a sequence of positive numbers. Therefore by Theorems 1.1 and 1.2, problem (1.1) with $\sigma = \sigma_n$ has a positive solution (u_n, v_n) for all large n . Recall that $C_0 > 0$ is given in (2.6), which is completely determined by δ and d_1 . Due to (3.1) there exists a unique $\tau_0^* > 0$ such that

$$(3.2) \quad m = \int_0^1 f(w_0e^{-A\tau_0^*\max\{C_0/2,C_0y\}})dy.$$

Let us then define

$$(3.3) \quad v^* = v^*(m) := \frac{w_0}{\alpha}e^{-A\tau_0^*C_0/2} + \frac{\tau_0^*}{d_2}mC_0(1 + \beta^{-1}).$$

Let $\underline{v}(m) > 0$ be uniquely determined by

$$m = f(\alpha \underline{v}(m)).$$

By (3.1), we always have $v_0 > \underline{v}(m)$.

When $m < f(w_0 e^{-A_0})$, we can find a unique $\tau_1^* > 0$ such that

$$(3.4) \quad m = \int_0^1 f(w_0 e^{-A_0 - A\tau_1^* \max\{C_0/2, C_0 y\}}) dy.$$

We now define

$$(3.5) \quad v_* = v_*(m) := \begin{cases} \frac{w_0}{\alpha} e^{-A_0 - A\tau_1^* C_0/2} + \frac{\tau_1^*}{d_2} m C_0 \beta^{-1} & \text{if } m < f(w_0 e^{-A_0}), \\ \underline{v}(m) & \text{if } f(w_0 e^{-A_0}) \leq m < f(w_0). \end{cases}$$

It is easily seen that $v_*(m)$ is continuous in m .

As we will see below, to completely determine the asymptotic profile of (u_n, v_n) , it is necessary to distinguish the cases $v_0 \in [v_*(m), v^*(m)]$, $v_0 > v^*(m)$, and $v_0 < v_*(m)$.

THEOREM 3.1. *Suppose that $v_0 > \underline{v}(m)$ and*

$$(3.6) \quad v_*(m) \leq v_0 \leq v^*(m).$$

Then the system (2.8) and (2.9), namely,

$$\begin{cases} w_0 e^{-A_0 x_* - A\tau_*(C_0/2)} = \alpha \left[v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - x_*) \right], \\ m = \int_0^1 f(w_0 e^{-A_0 x_* - A\tau_* \max\{C_0/2, C_0 y\}}) dy, \end{cases}$$

has a unique solution pair (x_, τ_*) satisfying $x_* \in [0, 1]$ and $\tau_* > 0$. Moreover,*

$$u_n \rightarrow 0 \text{ in } C([0, 1] \setminus [x_* - \epsilon, x_* + \epsilon]) \quad \forall \epsilon > 0, \quad \int_0^1 u_n dx \rightarrow \tau_* C_0,$$

$$v_n(x) \rightarrow v_0 - \frac{\tau_*}{d_2} m C_0 (1 + \beta^{-1} - \max\{x, x_*\}) \quad \text{uniformly in } [0, 1].$$

Furthermore, $x_ = 0$ if $v_0 = v^*(m)$, $x_* \in (0, 1)$ if $v_*(m) < v_0 < v^*(m)$, and $x_* = 1$ if $v_0 = v_*(m)$.*

Proof. Using the notation of the previous section, by passing to a subsequence, $x_n \rightarrow x_* \in [0, 1]$. By possibly passing to a further subsequence, the behavior of (u_n, v_n) as $n \rightarrow \infty$ is then determined by Lemmas 2.2, 2.3 (if $x_* \in (0, 1)$), Lemma 2.4 (if $x_* = 0$ and subcases (a1) and (a2) occur), and Lemma 2.5 (if $x_* = 1$ and subcases (b1) and (b2) happen).

If we can show that x_* and τ_* are uniquely determined by the value of v_0 , then the corresponding results in the previous section would hold not only for a subsequence, but for the entire original sequence, and hence the behavior of (u_n, v_n) as $n \rightarrow \infty$ would be completely determined.

The rather long proof below is broken into several steps.

Step 1. Subcases (a2) and (b2) do not happen

First we observe that subcase (a2) does not happen. Indeed, if this case occurs, then since $C(a_*) < C_0/2$, we see (as explained below) from a careful comparison of (2.19) and (3.2) that

$$\tau_* > \tau_0^*, \quad \tau_* C(a_*) < \tau_0^* C_0/2, \quad \tau_* [C_0/2 + C(a_*)] > \tau_0^* C_0.$$

In the comparison, we can deduce these inequalities one at a time, in the above order, and the previous inequalities are used for obtaining the next inequality. For example, to deduce $\tau_* C(a_*) < \tau_0^* C_0/2$ from $\tau_* > \tau_0^*$, we observe that $\tau_* C(a_*) \geq \tau_0^* C_0/2$ and $\tau_* > \tau_0^*$ would imply

$$\begin{aligned} \tau_* \max\{C(a_*), [C_0/2 + C(a_*)]y\} &\geq \max\{\tau_0^* C_0/2, [\tau_* C_0/2 + \tau_0^* C_0/2]y\} \\ &\geq \tau_0^* \max\{C_0/2, C_0 y\} \end{aligned}$$

with strict inequality holding in the last step for $y \in [1/2, 1]$, which is impossible when one compares (2.19) with (3.2).

It then follows from (2.20) and (2.21) that $v_0 > v^*(m)$, contradicting (3.6).

Similarly, if subcase (b2) happens, then from (2.27) we deduce

$$\tau_* > \tau_1^* \quad \text{and} \quad \tau_* [C_0/2 + \tilde{C}(b_*)] < \tau_1^* C_0,$$

which imply, by (2.28) and (2.29), that $v_0 < v_*(m)$, again contradicting (3.6). Therefore subcase (b2) cannot happen.

Thus, by our discussion in the previous section, we have the cases where (2.8) and (2.9) hold. To show that (2.8) and (2.9) have a unique solution (x_*, τ_*) satisfying $x_* \in [0, 1]$ and $\tau_* > 0$, we establish a procedure to uniquely find x_* and τ_* . In the discussion below, we will treat $v_0 > 0$ as a varying parameter.

Step 2. A procedure to solve (2.8) and (2.9).

It is useful to use the new variable

$$\lambda = A_0 x_* + A \tau_* C_0/2.$$

Then

$$x_* = (\lambda - A \tau_* C_0/2)/A_0,$$

and (2.8) can be rewritten as

$$\frac{w_0}{\alpha} e^{-\lambda} = v_0 - \frac{\tau_*}{d_2} m C_0 \left(1 + \beta^{-1} - \frac{\lambda - A \tau_* C_0/2}{A_0} \right)$$

or

$$\frac{m C_0}{d_2 A_0} \tau_* [(1 + \beta^{-1}) A_0 - \lambda + A(C_0/2) \tau_*] = v_0 - \frac{w_0}{\alpha} e^{-\lambda}.$$

We now consider the quadratic equation of τ :

$$(3.7) \quad \frac{m C_0}{d_2 A_0} \tau [(1 + \beta^{-1}) A_0 - \lambda + A(C_0/2) \tau] = v_0 - \frac{w_0}{\alpha} e^{-\lambda}.$$

For each $v_0 > 0$, let $\lambda_0(v_0)$ denote the minimal nonnegative λ such that $v_0 - \frac{w_0}{\alpha} e^{-\lambda} \geq 0$. Clearly

$$(3.8) \quad \lambda_0(v_0) = 0 \text{ if } v_0 \geq w_0/\alpha, \quad \lambda_0(v_0) \text{ is decreasing in } (0, w_0/\alpha], \quad \lim_{v_0 \rightarrow 0} \lambda_0(v_0) = \infty.$$

For each $v_0 > 0$ and $\lambda \geq \lambda_0(v_0)$, the quadratic equation (3.7) has a maximal zero, which we denote by $\tau(\lambda, v_0)$. It is easily seen that $\tau(\lambda, v_0) \geq 0$ and

$$(3.9) \quad \text{when } v_0 \leq w_0/\alpha, \tau(\lambda_0(v_0), v_0) = \max\left\{0, \frac{\lambda_0(v_0) - A_0(1 + \beta^{-1})}{AC_0/2}\right\},$$

$$(3.10) \quad \tau(\lambda, v_0) \text{ is increasing in } \lambda \text{ and in } v_0, \quad \lim_{v_0 \rightarrow \infty} \tau(\lambda, v_0) = \infty \text{ for fixed } \lambda \geq 0.$$

Since $\lambda_0(w_0/\alpha) = 0$, by (3.9), $\tau(\lambda_0(w_0/\alpha), w_0/\alpha) = 0$. Let us consider the continuous function

$$M(v_0) = \int_0^1 f(w_0 e^{-\lambda_0(v_0) - A\tau(\lambda_0(v_0), v_0) \max\{0, C_0 y - C_0/2\}}) dy.$$

The above observation shows that $M(w_0/\alpha) = f(w_0) > m$. By (3.8), we have $M(v_0) \rightarrow 0$ as $v_0 \rightarrow 0$. By (3.10), we deduce $M(v_0) \rightarrow 0$ as $v_0 \rightarrow \infty$. Hence from the continuity of $M(v_0)$ we can find v_{min} and v_{max} such that

$$0 < v_{min} < w_0/\alpha < v_{max} < \infty,$$

$$M(v_0) > m \quad \forall v_0 \in (v_{min}, v_{max}), \quad M(v_{min}) = M(v_{max}) = m.$$

Now for each $v_0 \in (v_{min}, v_{max})$,

$$m < \int_0^1 f(w_0 e^{-\lambda_0(v_0) - A\tau(\lambda_0(v_0), v_0) \max\{0, C_0 y - C_0/2\}}) dy.$$

This and the monotonicity of $\tau(\lambda, v_0)$ in λ imply that for such v_0 we can find a unique $\lambda_* = \lambda_*(v_0) > \lambda_0(v_0)$ such that

$$m = \int_0^1 f(w_0 e^{-\lambda_* - A\tau(\lambda_*, v_0) \max\{0, C_0 y - C_0/2\}}) dy.$$

Clearly $v_0 \rightarrow \lambda_*(v_0)$ is continuous and

$$\lambda_*(v_{min} + 0) = \lambda_0(v_{min}), \quad \lambda_*(v_{max} - 0) = \lambda_0(v_{max}).$$

So we may define

$$\lambda_*(v_{min}) = \lambda_0(v_{min}), \quad \lambda_*(v_{max}) = \lambda_0(v_{max}).$$

We claim that the function $T(v_0) := \tau(\lambda_*(v_0), v_0)$ is increasing in $[v_{min}, v_{max}]$. Otherwise, we can find $v_{min} \leq s_1 < s_2 \leq v_{max}$ such that $T(s_1) \geq T(s_2)$. Since

$$\int_0^1 f(w_0 e^{-\lambda_*(s_1) - AT(s_1) \max\{0, C_0 y - C_0/2\}}) = \int_0^1 f(w_0 e^{-\lambda_*(s_2) - AT(s_2) \max\{0, C_0 y - C_0/2\}}),$$

$T(s_1) \geq T(s_2)$ implies that $\lambda_*(s_1) \leq \lambda_*(s_2)$, which implies, by the monotonicity of $\tau(\lambda, v_0)$,

$$T(s_1) = \tau(\lambda_*(s_1), s_1) < \tau(\lambda_*(s_2), s_2) = T(s_2).$$

This contradiction proves the claimed monotonicity of $T(v_0)$.

We show next that $T(v_{max}) > \tau_0^*$. Since $v_{max} > w_0/\alpha$, we have $\lambda_0(v_{max}) = 0$ and hence

$$m = M(v_{max}) = \int_0^1 f(w_0 e^{-A\tau(0, v_{max}) \max\{0, C_0 y - C_0/2\}}) dy.$$

By (3.2),

$$m = \int_0^1 f(w_0 e^{-A\tau_0^* C_0/2 - A\tau_0^* \max\{0, C_0 y - C_0/2\}}) dy.$$

Comparing the above two expressions, we obtain $\tau(0, v_{max}) > \tau_0^*$. Hence

$$T(v_{max}) = \tau(\lambda_*(v_{max}), v_{max}) = \tau(\lambda_0(v_{max}), v_{max}) = \tau(0, v_{max}) > \tau_0^*,$$

as we wanted.

We now consider $T(v_{min})$. We have two different cases: $m < f(w_0 e^{-A_0})$ and $m \geq f(w_0 e^{-A_0})$. First consider the case $m < f(w_0 e^{-A_0})$. We show that $T(v_{min}) < \tau_1^*$ in this case. Since $\lambda_*(v_{min}) = \lambda_0(v_{min})$ we have

$$T(v_{min}) = \tau(\lambda_0(v_{min}), v_{min}).$$

Hence, by (3.9),

$$T(v_{min}) = \max\left\{0, \frac{\lambda_0(v_{min}) - A_0(1 + \beta^{-1})}{AC_0/2}\right\}.$$

If $\frac{\lambda_0(v_{min}) - A_0(1 + \beta^{-1})}{AC_0/2} \leq 0$, then $T(v_{min}) = 0 < \tau_1^*$. If $\frac{\lambda_0(v_{min}) - A_0(1 + \beta^{-1})}{AC_0/2} > 0$, then

$$T(v_{min}) = \frac{\lambda_0(v_{min}) - A_0(1 + \beta^{-1})}{AC_0/2},$$

and hence

$$\begin{aligned} m &= \int_0^1 f(w_0 e^{-\lambda_0(v_{min}) - AT(v_{min}) \max\{0, C_0 y - C_0/2\}}) dy \\ &= \int_0^1 f(w_0 e^{-A_0(1 + \beta^{-1}) - AT(v_{min}) C_0/2 - AT(v_{min}) \max\{0, C_0 y - C_0/2\}}) dy \\ &= \int_0^1 f(w_0 e^{-A_0(1 + \beta^{-1}) - AT(v_{min}) \max\{C_0/2, C_0 y\}}) dy. \end{aligned}$$

Comparing this with (3.4), we find that $T(v_{min}) < \tau_1^*$.

With the above properties of $T(v_0)$, we can uniquely determine v_* and v^* with $v_{min} < v_* < v^* < v_{max}$ such that

$$T(v^*) = \tau_0^*, \quad T(v_*) = \tau_1^*.$$

We claim that $v^* = v^*(m)$ and $v_* = v_*(m)$. Indeed, from

$$m = \int_0^1 f(w_0 e^{-\lambda_*(v^*) - AT(v^*) \max\{0, C_0 y - C_0/2\}}) dy$$

and $T(v^*) = \tau_0^*$, we easily see by comparing with (3.2) that $\lambda_*(v^*) = \tau_0^*AC_0/2$. Hence

$$\tau_0^* = T(v^*) = \tau(\lambda_*(v^*), v^*) = \tau(\tau_0^*AC_0/2, v^*).$$

By the definition of $\tau(\lambda, v_0)$, the above identity means that $\tau = \tau_0^*$ solves (3.7) with $\lambda = \tau_0^*AC_0/2$ and $v_0 = v^*$. Therefore we may compare with (3.3) to deduce $v^* = v^*(m)$. Similarly, we can show that $v_* = v_*(m)$.

Since T is monotone, for each $v_0 \in [v_*(m), v^*(m)]$, $T(v_0) \in [\tau_1^*, \tau_0^*]$. Hence we can compare (3.2) and (3.4) with

$$m = \int_0^1 f(w_0e^{-\lambda_*(v_0)-AT(v_0)\max\{0, C_0y-C_0/2\}})dy$$

to find that, for such v_0 , we necessarily have

$$AT(v_0)C_0/2 \leq \lambda_*(v_0) \leq A_0 + AT(v_0)C_0/2;$$

otherwise we would arrive at contradictions to $T(v_0) \in [\tau_1^*, \tau_0^*]$. This implies that there exists a unique $x_* \in [0, 1]$ such that

$$\lambda_*(v_0) = A_0x_* + AT(v_0)C_0/2.$$

Let $\tau_* = T(v_0)$; we find that (x_*, τ_*) solves (2.8) and (2.9).

We next consider the case $m \geq f(w_0e^{-A_0})$. In this case, $v_*(m) = \underline{v}(m)$; moreover, we show that

$$T(v_{min}) = 0, \quad v_{min} = \underline{v}(m).$$

Indeed, from

$$T(v_{min}) = \tau(\lambda_*(v_{min}), v_{min}) = \tau(\lambda_0(v_{min}), v_{min}),$$

we obtain

$$\begin{aligned} m &= \int_0^1 f(w_0e^{-\lambda_0(v_{min})-A\tau(\lambda_0(v_{min}), v_{min})\max\{0, C_0y-C_0/2\}})dy \\ &\leq f(w_0e^{-\lambda_0(v_{min})}). \end{aligned}$$

It follows that $\lambda_0(v_{min}) \leq A_0 < A_0(1+\beta^{-1})$. By (3.9), we deduce $\tau(\lambda_0(v_{min}), v_{min}) = 0$, that is, $T(v_{min}) = 0$. This gives

$$m = \int_0^1 f(w_0e^{-\lambda_0(v_{min})})dy = f(w_0e^{\lambda_0(v_{min})}).$$

Hence

$$\alpha \underline{v}(m) = w_0e^{-\lambda_0(v_{min})}.$$

On the other hand, since $v_{min} < w_0/\alpha$, by the definition of the function λ_0 ,

$$v_{min} - \frac{w_0}{\alpha}e^{-\lambda_0(v_{min})} = 0.$$

Therefore we have $v_{min} = \underline{v}(m)$.

We can now conclude that there exists a unique $v^* \in (v_{min}, v_{max})$ such that $T(v^*) = \tau_0^*$. We may then prove $v^* = v^*(m)$ as before. Since T is monotone, for each $v_0 \in (v_*(m), v^*(m)] = (\underline{v}(m), v^*(m)]$, $T(v_0) \in (0, \tau_0^*]$. Hence we can compare (3.2) and $m \geq f(w_0 e^{-A_0})$ with (3.1) to deduce

$$AT(v_0)C_0/2 \leq \lambda_*(v_0) < A_0 + AT(v_0)C_0/2,$$

and there exists a unique $x_* \in [0, 1)$ such that

$$\lambda_*(v_0) = A_0 x_* + AT(v_0)C_0/2.$$

Let $\tau_* = T(v_0)$; we find that (x_*, τ_*) solves (2.8) and (2.9).

The above discussion shows that when (3.6) holds, (2.8) and (2.9) have at least one solution (x_*, τ_*) satisfying $x_* \in [0, 1]$ and $\tau_* > 0$, and such a solution can be found by following the above procedure.

Step 3. Uniqueness of (x_*, τ_*) and completion of the proof.

We next show that when (3.6) holds, (2.8) and (2.9) have a unique solution (x_*, τ_*) satisfying $x_* \in [0, 1]$ and $\tau_* > 0$. So let (x_*, τ_*) be an arbitrary solution of (2.8) and (2.9) with $v_0 \in [v_*(m), v^*(m)] \cap (\underline{v}(m), v^*(m)]$ and $x_* \in [0, 1]$, $\tau_* > 0$. Then τ_* must be the maximal zero of (3.7) with $\lambda = A_0 x_* + A\tau_* C_0/2 > 0$; this is the case because $v_0 - \frac{w_0}{\alpha} e^{-\lambda} > 0$, and thus the two zeros of (3.7) are of opposite sign. Therefore, using our earlier notation,

$$\tau_* = \tau(\lambda, v_0), \quad \lambda > \lambda_0(v_0).$$

Then (2.9) yields

$$\begin{aligned} m &= \int_0^1 f(w_0 e^{-A_0 x_* - A\tau_* \max\{C_0/2, C_0 y\}}) dy \\ &= \int_0^1 f(w_0 e^{-\lambda - A\tau(\lambda, v_0) \max\{0, C_0 y - C_0/2\}}) dy. \end{aligned}$$

Since $v_0 \in [v_*(m), v^*(m)] \cap (\underline{v}(m), v^*(m)] \subset (v_{min}, v_{max})$, in view of the above identity, our definition of $\lambda_*(v_0)$ implies that $\lambda = \lambda_*(v_0)$ and hence $\tau(\lambda, v_0) = T(v_0)$; that is, $\tau_* = \tau(\lambda, v_0) = T(v_0)$. This implies that the solution pair (x_*, τ_*) is the same as the one obtained through our procedure introduced above for solving (2.8) and (2.9). Hence there is a unique solution.

With τ_* and x_* uniquely determined now, it is easily seen that our conclusions for u_n and v_n follow from Lemmas 2.2, 2.3, 2.4, and 2.5.

Moreover, from the above procedure for finding (x_*, τ_*) , we easily see that $x_* = 0$ if $v_0 = v^*(m)$, $x_* \in (0, 1)$ if $v_*(m) < v_0 < v^*(m)$, and $x_* = 1$ if $v_0 = v_*(m)$.

The proof of the theorem is now complete. \square

Next we consider the case that $v_0 > v^*(m)$. Let $0 < \lambda_0 < \lambda^0$ be uniquely determined by

$$(3.11) \quad m = f(w_0 e^{-A\lambda_0}) = \int_0^1 f(w_0 e^{-A\lambda_0 y}) dy.$$

For each $\lambda \in [0, \lambda_0]$, we can find a unique $\Gamma = \Gamma(\lambda)$ such that

$$(3.12) \quad m = \int_0^1 f(w_0 e^{-A \max\{\lambda, \Gamma y\}}) dy.$$

Moreover, it is easily seen that $\lambda \mapsto \Gamma(\lambda)$ is a continuous decreasing function with

$$\Gamma(\lambda_0) = \lambda_0, \quad \Gamma(0) = \lambda^0.$$

Therefore we can find a unique $\lambda_*^0 \in (0, \lambda_0)$ such that

$$\Gamma(\lambda_*^0) = 2\lambda_*^0.$$

Comparing with (3.2), we find that actually

$$(3.13) \quad \lambda_*^0 = \tau_0^* C_0/2.$$

We define

$$\Lambda(\lambda) := \frac{w_0}{\alpha} e^{-A\lambda} + \frac{\Gamma(\lambda)}{d_2} m(1 + \beta^{-1}).$$

Clearly $\Lambda(\lambda)$ is a decreasing function on $[0, \lambda_0]$ with

$$\Lambda(0) = \frac{w_0}{\alpha} + \frac{\lambda^0}{d_2} m(1 + \beta^{-1}), \quad \Lambda(\lambda_*^0) = \frac{w_0}{\alpha} e^{-A\lambda_*^0} + \frac{2\lambda_*^0}{d_2} m(1 + \beta^{-1}).$$

Due to (3.13), we find that

$$\Lambda(\lambda_*^0) = v^*(m).$$

THEOREM 3.2. *Suppose that*

$$(3.14) \quad v_0 > v^*(m) = \Lambda(\lambda_*^0).$$

If $v_0 < \Lambda(0)$ and $\lambda^ \in (0, \lambda_*^0)$ is uniquely determined by $v_0 = \Lambda(\lambda^*)$, then*

$$u_n \rightarrow 0 \text{ in } C([\epsilon, 1]) \quad \forall \epsilon \in (0, 1), \quad \int_0^1 u_n dx \rightarrow \Gamma(\lambda^*),$$

$$v_n(x) \rightarrow v_0 - \frac{\Gamma(\lambda^*)}{d_2} m(1 + \beta^{-1} - x) \text{ uniformly in } [0, 1].$$

If $v_0 \geq \Lambda(0)$, then the above conclusions hold with $\lambda^ = 0$.*

Proof. We first show that case (a2) happens. Let us start by observing that the cases leading to (2.8) and (2.9) (namely, cases (i), (ii)(a1), and (iii)(b1)) cannot happen. Indeed, in these cases, (x_*, τ_*) solves (2.8) and (2.9) with $x_* \in [0, 1]$ and $\tau_* > 0$. As in Step 3 of the proof of Theorem 3.1, denoting $\lambda = A_0 x_* + A\tau_* C_0/2$, we must have $\tau_* = \tau(\lambda, v_0)$ and $\lambda > \lambda_0(v_0)$. Then (2.9) gives

$$m = \int_0^1 f(w_0 e^{-\lambda - A\tau(\lambda, v_0) \max\{0, C_0 y - C_0/2\}}) dy.$$

Since $v_0 > v^*(m)$, we have either

$$v_0 > v_{max} \quad \text{or} \quad v_0 \in (v^*(m), v_{max}].$$

If $v_0 \in (v^*(m), v_{max}] \subset (v_{min}, v_{max})$, then the above identity implies that $\lambda = \lambda_*(v_0)$ and hence $\tau(\lambda, v_0) = T(v_0)$. From $v_0 > v^*(m)$ we now deduce $\tau_* = T(v_0) > \tau_0^*$, and hence we can compare (2.9) with (3.2) to deduce $x_* < 0$, a contradiction.

If $v_0 > v_{max}$, then by the monotonicity of $\tau(\cdot, \cdot)$, we deduce

$$\tau(\lambda, v_0) > \tau(\lambda, v_{max}) > \tau(0, v_{max}).$$

Therefore, recalling $\lambda_*(v_{max}) = \lambda_0(v_{max}) = 0$, we obtain

$$\begin{aligned} m &= \int_0^1 f(w_0 e^{-\lambda - A\tau(\lambda, v_0) \max\{0, C_0 y - C_0/2\}}) dy \\ &< \int_0^1 f(w_0 e^{-A\tau(0, v_{max}) \max\{0, C_0 y - C_0/2\}}) dy \\ &= m, \end{aligned}$$

again a contradiction. Therefore none of the cases that lead to (2.8) and (2.9) can happen. This implies that either (a2) or (b2) happens.

Next we show that case (b2) cannot happen. Otherwise, by (2.27) we obtain

$$m < f(w_0 e^{-A_0}).$$

Hence $\tau_1^* > 0$ is defined. Moreover, comparing (2.27) with (3.4), we obtain

$$\tau_* > \tau_1^*, \quad \tau_* \tilde{C}(b_*) < \tau_1^* C_0,$$

which imply, by (2.28) and (2.29), that $v_0 < v_*(m) < v^*(m)$, contradicting (3.14).

Therefore we necessarily have case (a2). We now introduce the notation

$$\lambda = \tau_*[C(a_*) - C_0/2], \quad \Gamma = \tau_* C(a_*).$$

From (2.19), (2.20), and (2.21) we find that

$$(3.15) \quad m = \int_0^1 f(w_0 e^{-A \max\{\lambda, \Gamma y\}}) dy,$$

$$(3.16) \quad v_0 \geq \frac{w_0}{\alpha} e^{-A\lambda} + \frac{\Gamma}{d_2} m(1 + \beta^{-1}),$$

with equality holding if $a_* > 0$.

Suppose now that $v_0 \geq \Lambda(0)$. We claim that in this case we have $\lambda = 0$ and hence, by (3.15), $\Gamma = \Gamma(0) = \lambda^0$. Suppose for the sake of contradiction that $\lambda > 0$. From (3.15) and (3.11) we easily see that $\lambda \leq \lambda_0$. Now $C(a_*) - C_0/2 > 0$, and hence $a_* > 0$. Thus equality in (3.16) holds. By (3.15) we deduce $\Gamma = \Gamma(\lambda)$, and hence it follows from (3.16) that $v_0 = \Lambda(\lambda) < \Lambda(0)$, contradicting our assumption on v_0 above. Hence in this case, we have $\lambda = 0$ and thus

$$C(a_*) - C_0/2 = 0, \quad \tau_* = \Gamma(0)/(C_0/2).$$

Next we suppose that $v^*(m) < v_0 < \Lambda(0)$. From (3.15) we deduce $\Gamma = \Gamma(\lambda)$ for some $\lambda \in [0, \lambda_0]$. We must have $\lambda > 0$ for otherwise, from (3.15) and (3.16), we deduce $\Gamma = \Gamma(0)$ and $v_0 \geq \Lambda(0)$, contradicting our current assumption on v_0 . Therefore $\lambda > 0$ and hence $a_* > 0$, implying that equality in (3.16) holds. Recalling $\Gamma = \Gamma(\lambda)$, we thus obtain $v_0 = \Lambda(\lambda)$ and $\lambda = \lambda^*$. It follows that τ_* and a_* in Lemma 2.4 are uniquely determined by

$$\tau_*[C(a_*) - C_0/2] = \lambda^*, \quad \tau_* C(a_*) = \Gamma(\lambda^*),$$

namely,

$$\tau_* = \frac{\Gamma(\lambda^*) - \lambda^*}{C_0/2}, \quad a_* = C^{-1}(\lambda^*/\tau_* + C_0/2).$$

The rest of the proof now follows from Lemma 2.4. \square

We now consider the remaining case that $\underline{v}(m) < v_0 < v_*(m)$, which can happen only if $m < f(w_0e^{-A_0})$. Suppose that $\lambda_0, \lambda^0, \lambda_*^0$, and $\Gamma(\lambda)$ are as in Theorem 3.2 but with w_0 there replaced by $w_0e^{-A_0}$, and we denote them by $\tilde{\lambda}_0, \tilde{\lambda}^0, \tilde{\lambda}_*^0$, and $\tilde{\Gamma}(\lambda)$, respectively. Define

$$\Delta(\lambda) := \frac{w_0}{\alpha}e^{-A_0-A\lambda} + \frac{\tilde{\Gamma}(\lambda)}{d_2}m\beta^{-1}.$$

Then $\Delta(\lambda)$ is a decreasing function over $[0, \tilde{\lambda}_0]$ with

$$\Delta(0) = \frac{w_0}{\alpha}e^{-A_0} + \frac{\tilde{\lambda}_0^0}{d_2}m\beta^{-1}, \quad \Delta(\tilde{\lambda}_*^0) = \frac{w_0}{\alpha}e^{-A_0-A\tilde{\lambda}_*^0} + \frac{\tilde{\lambda}_*^0}{d_2}m\beta^{-1} = v_*(m).$$

THEOREM 3.3. *Suppose that $m < f(w_0e^{-A_0})$ and*

$$(3.17) \quad \underline{v}(m) < v_0 < v_*(m) = \Delta(\tilde{\lambda}_*^0).$$

If $v_0 > \Delta(0)$ and $\lambda_ \in (0, \tilde{\lambda}_*^0)$ is uniquely determined by $v_0 = \Delta(\lambda_*)$, then*

$$u_n \rightarrow 0 \text{ in } C([0, 1 - \epsilon]) \quad \forall \epsilon \in (0, 1), \quad \int_0^1 u_n dx \rightarrow \tilde{\Gamma}(\lambda_*),$$

$$v_n(x) \rightarrow v_0 - \frac{\tilde{\Gamma}(\lambda_*)}{d_2}m\beta^{-1} \text{ uniformly in } [0, 1].$$

If $\underline{v}(m) < v_0 \leq \Delta(0)$, then the above conclusions hold with $\lambda_ = 0$.*

Proof. This is similar to the proof of Theorem 3.2. Here we can show that case (b2) must happen, and then we use Lemma 2.5. We omit the details. \square

Remark 3.4. Our results in this section reveal an interesting property for the limiting total biomass $\lim_{n \rightarrow \infty} \int_0^1 u_n(x)dx$. First consider the case $v_*(m) \leq v_0 \leq v^*(m)$. By Theorem 3.1, in this case the above limit has value τ_*C_0 . By the proof of Theorem 3.1, we know that $\tau_* = T(v_0)$ with $T(v_0)$ a strictly increasing function of v_0 . Therefore the limit is strictly increasing with v_0 for $v_0 \in [v_*(m), v^*(m)]$.

If $v_0 \in (v^*(m), \Lambda(0))$, then by Theorem 3.2, $\lim_{n \rightarrow \infty} \int_0^1 u_n(x)dx = \Gamma(\lambda^*) = \Gamma \circ \Lambda^{-1}(v_0)$. Since $\Gamma(\cdot)$ and $\Lambda(\cdot)$ are both strictly decreasing functions, $\Gamma \circ \Lambda^{-1}(\cdot)$ is strictly increasing, and hence the limit is strictly increasing with v_0 for $v_0 \in (v^*(m), \Lambda(0))$. However, by Theorem 3.2, this limit takes the same value $\Gamma(0)$ for all $v_0 \geq \Lambda(0)$.

If $v_0 \in (\Delta(0), v_*(m))$, then by Theorem 3.3, the limit takes value $\tilde{\Gamma}(\lambda_*) = \tilde{\Gamma} \circ \Delta^{-1}(v_0)$, which is strictly increasing in v_0 , but it takes the same value $\tilde{\Gamma}(0)$ for all $v_0 \in (\underline{v}(m), \Delta(0)]$.

Therefore if we denote

$$v_{**} = v_{**}(m) := \Delta(0), \quad v^{**} = v^{**}(m) := \Lambda(0),$$

then $\lim_{n \rightarrow \infty} \int_0^1 u_n(x)dx$ is strictly increasing with v_0 as v_0 varies in the range $v_{**} \leq v_0 \leq v^{**}$, but is constant for $v_0 \geq v^{**}$ or for $v_0 \in (\underline{v}(m), v_{**}]$.

4. Biological interpretations. In this section we compare our results with the game theoretical model in [KL] and explain the predictions that our theoretical results offer for the phytoplankton problem being modelled.

Since $u(x)$ and $v(x)$ are, respectively, rescaled versions of the biomass function $b(x)$ and the nutrient function $R(x)$ used in the original model of [KL] (see Part I for details), we will interpret $u(x)$ and $v(x)$ as representing the (steady) distributions of the biomass and nutrient at depth x of a water column with surface at $x = 0$ and bottom at $x = 1$.

- (i) First we note that, if we replace $\max\{C_0/2, C_0y\}$ in (2.9) by C_0 , then the system of equations for (x_*, τ_*) in Theorem 3.1 reduces to the game theoretical model of [KL], namely, equations (4) and (5) in [KL] with $\hat{B} = \tau_*C_0$. Thus the game theoretical model of [KL] is a simplified version of the rigorous limiting equations here. It captures the main properties of the limiting equations but only in the case that $v_* \leq v_0 \leq v^*$.
- (ii) When $v_0 > v^* = v^*(m)$, from Theorem 3.2 and Step 1 of the proof of Theorem 3.1, we see that as $\sigma \rightarrow \infty$ the total biomass has limit

$$\Gamma(\lambda^*) = \tau_*C(a_*) > \tau_0^*C_0.$$

If we have simply used (2.8) and (2.9) with $x_* = 0$ to calculate the total biomass, we would have obtained the incorrect limit $\tau_0^*C_0$. Similarly, the limit of the total biomass in the case of Theorem 3.3 is less than the value one would have obtained by simply using (2.8) and (2.9) with $x_* = 1$.

- (iii) Remark 3.4 shows that as $\sigma \rightarrow \infty$, the limit of the total biomass is strictly increasing with v_0 in the range $v_{**} \leq v_0 \leq v^{**}$. It is interesting to notice that the layer position of the biomass (in the limit) already reaches the surface at $v_0 = v^*$ (i.e., $x_* = 0$ when $v_0 = v_*$), but as the nutrient level at the sediment v_0 increases past the critical value v^* , though the layer remains at the surface, the total biomass keeps increasing until v_0 reaches a second critical value v^{**} , where the total biomass reaches a maximum, and then remains at this value even if v_0 is further increased. On the other hand, if v_0 is decreased to v_* , the layer of the biomass lowers to the bottom ($x_* = 1$), but as v_0 decreases past v_* , though the layer remains at the bottom, the total biomass keeps decreasing until v_0 reaches the critical value v_{**} , where the total biomass reaches its minimal value, and then remains at this minimal value until v_0 is so low ($v_0 \leq \underline{v}(m)$) that the phytoplankton can no longer survive.
- (iv) Our results support the important predictions in [KL] that depth-regulating phytoplankton can form a thin layer in a poorly mixed water column and that the concentration of the limiting nutrient should be low and constant above the phytoplankton layer and linearly increasing with depth below the layer. The predictions in item (iii) above seem to provide new insights to this problem.
- (v) We could fix v_0 and use a different parameter in the model, say the surface light level w_0 , as a varying parameter to interpret the phenomena represented in items (i)–(iii) above.

Acknowledgment. Y. Du thanks NCTS for the hospitality during his visit, when a major part of this paper was written.

REFERENCES

- [DH] Y. DU AND S.-B. HSU, *Concentration phenomena in a nonlocal quasi-linear problem modelling phytoplankton I: Existence*, SIAM J. Math. Anal., 40 (2008), pp. 1419–1440.
- [GT] D. GILBARG AND N. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, Springer-Verlag, Berlin, 1983.
- [KL] C. A. KLAUSMEIER AND E. LITCHMAN, *Algal games: The vertical distribution of phytoplankton in poorly mixed water columns*, Limnol. Oceanogr., 46 (2001), pp. 1998–2007.

ON p -HARMONIC MAP HEAT FLOWS FOR $1 \leq p < \infty$ AND THEIR FINITE ELEMENT APPROXIMATIONS*

JOHN W. BARRETT[†], XIAOBING FENG[‡], AND ANDREAS PROHL[§]

Abstract. Motivated by emerging applications from imaging processing, this paper studies the heat flow of a generalized p -harmonic map into spheres for the whole spectrum, $1 \leq p < \infty$, in a unified framework. The existence of global weak solutions is established for the flow using the energy method together with a regularization and a penalization technique. In particular, a BV -solution concept is introduced and the existence of such a solution is proved for the 1-harmonic map heat flow. The main idea used to develop such a theory is to exploit the properties of measures of the forms $\mathcal{A} \cdot \nabla \mathbf{v}$ and $\mathcal{A} \wedge \nabla \mathbf{v}$, which pair a divergence- L^1 , or a divergence-measure, tensor field \mathcal{A} and a BV -vector field \mathbf{v} . Based on these analytical results, a practical fully discrete finite element method is then proposed for approximating weak solutions of the p -harmonic map heat flow, and the convergence of the proposed numerical method is also established.

Key words. p -harmonic maps, heat flow, penalization, energy method, color image denoising, finite element method

AMS subject classifications. 35K65, 58E20, 35Q80, 65M12, 65M60

DOI. 10.1137/070680825

1. Introduction. Let $\Omega \subset \mathbf{R}^m$ be a bounded domain with smooth boundary $\partial\Omega$, and let S^{n-1} denote the unit sphere in \mathbf{R}^n . A map $\mathbf{u} \in C^1(\Omega, S^{n-1})$ is called a p -harmonic map if it is a critical point of the following p -energy:

$$(1.1) \quad E_p(\mathbf{v}) := \frac{1}{p} \int_{\Omega} |\nabla \mathbf{v}|^p dx, \quad 1 \leq p < \infty.$$

It is well known [12, 21, 40] that the Euler–Lagrange equation of the p -energy is

$$(1.2) \quad -\Delta_p \mathbf{u} = |\nabla \mathbf{u}|^p \mathbf{u},$$

where

$$(1.3) \quad \Delta_p \mathbf{u} := \operatorname{div}(|\nabla \mathbf{u}|^{p-2} \nabla \mathbf{u}).$$

Note that Δ_p is often called the p -Laplacian. It is easy to see that (1.2) is a *degenerate elliptic* equation for $p > 2$ and a *singular elliptic* equation for $1 \leq p < 2$. These degeneracy and singular characteristics both disappear when $p = 2$.

We call a map $\mathbf{u} \in W^{1,p}(\Omega, S^{n-1})$ a *weakly p -harmonic map* if \mathbf{u} satisfies (1.2) in the distributional sense. Here $W^{1,p}(\Omega, S^{n-1})$ denotes the Sobolev space

$$W^{1,p}(\Omega, S^{n-1}) := \{ \mathbf{u} \in W^{1,p}(\Omega, \mathbf{R}^n); \mathbf{u}(x) \in S^{n-1} \text{ for a.e. } x \in \Omega \}.$$

*Received by the editors January 23, 2007; accepted for publication (in revised form) June 10, 2008; published electronically November 7, 2008.

<http://www.siam.org/journals/sima/40-4/68082.html>

[†]Department of Mathematics, Imperial College London, London, SW7 2AZ, UK (jwb@ic.ac.uk).

[‡]Department of Mathematics, The University of Tennessee, Knoxville, TN 37996 (xfeng@math.utk.edu). The work of this author was partially supported by NSF grant DMS-0410266.

[§]Mathematisches Institut, Universität Tübingen, Auf der Morgenstelle 10, D-72076 Tübingen, Germany (prohl@na.uni-tuebingen.de)

One well-known method for looking for a weakly p -harmonic map is the homotopy method (or the gradient descent method), which then leads to considering the following gradient flow (or heat flow) for the p -energy functional E_p complemented with some given boundary and initial conditions:

$$(1.4) \quad \mathbf{u}_t - \Delta_p \mathbf{u} = |\nabla \mathbf{u}|^p \mathbf{u} \quad \text{in } \Omega_T := \Omega \times (0, T),$$

$$(1.5) \quad |\mathbf{u}| = 1 \quad \text{in } \Omega_T.$$

Clearly, (1.4) is a *degenerate parabolic* equation for $p > 2$ and a *singular parabolic* equation for $1 \leq p < 2$. Again, these degeneracy and singular characteristics both disappear when $p = 2$.

We remark that p -harmonic maps and weakly p -harmonic maps between two Riemannian manifolds (M, g) and (N, h) and their heat flows can be defined in the same fashion (cf. [21, 43, 44, 46, 47]). In this paper, we shall only consider the case $M = \Omega$ and $N = S^{n-1}$, which is sufficient for the applications that we are interested in.

The p -harmonic map and its heat flow, in particular, the harmonic map and the harmonic map heat flow ($p = 2$), have been extensively studied in the past twenty years for $1 < p < \infty$ (cf. [9, 11, 12, 13, 16, 17, 18, 21, 22, 25, 30, 31, 32, 34, 38, 39, 40, 41, 46, 48] for $1 < p < \infty$; [26, 27] for $p = 1$). In the case when the target manifold is a sphere, the existence of a global weak solution for the harmonic flow was first proved by Chen in [11] using a penalization technique. The result and the penalization technique were extended to the p -harmonic flow for $p > 2$ by Chen, Hong, and Hungerbühler in [12]. The p -harmonic flow for $1 < p < 2$ was solved by Misawa in [41] using a time discretization technique (the method of Rothe) proposed in [31], and by Liu in [39] using a penalization technique similar to that of [12]. The p -harmonic flow ($1 < p < \infty$) from a unit ball in \mathbf{R}^m into $S^1 \subset \mathbf{R}^2$ was studied by Courilleau and Demengel in [16]. In the case of general target manifolds the p -harmonic flow ($1 < p < \infty$) with small initial data was treated by Fardoun and Regbaoui in [22], and the conformal case of the p -harmonic flow was considered by Hungerbühler in [33]. Nonuniqueness of the p -harmonic flow was addressed in [15, 32]. Recently, Giga, Kashima, and Yamazaki [28] proved existence of strong local solutions for 1-harmonic map heat flow using nonlinear semigroup theory. Besides the great amount of mathematical interests in the p -harmonic map and the p -harmonic flow, research on these problems has been strongly motivated by applications of the harmonic map and its heat flow in liquid crystals and micromagnetism; we refer to [2, 5, 7, 8, 19, 29, 36, 37, 51] and the references therein for detailed exposition in this direction.

Another reason, which is the main motivation of this paper, for studying the p -harmonic map and its heat flow, in particular, for $1 \leq p < 2$, arises from their emerging and intriguing applications to image processing for denoising color images (cf. [49, 50]). We recall that a color image is often expressed by the RGB color system, in which a vector $\mathbf{I}(x) = (r(x), g(x), b(x))$ for each pixel $x = (x_1, x_2)$ is used to represent the intensity of the three primary colors (red, green, and blue). The chromaticity and brightness of a color image are deduced from the RGB system by decomposing $\mathbf{I}(x)$ into two components,

$$\eta(x) := |\mathbf{I}(x)| \quad (\text{brightness}),$$

$$\mathbf{g} := \frac{\mathbf{I}(x)}{|\mathbf{I}(x)|} \quad (\text{chromaticity}),$$

where $|\mathbf{I}(x)|$ stands for the Euclidean norm of $\mathbf{I}(x)$. By definition, the chromaticity must lie on the unit sphere S^2 . One of the main benefits of the chromaticity and brightness decomposition is that it allows one to denoise η and \mathbf{g} separately using different methods. For example, one may denoise η by the well-known total variation (TV) model of Rudin, Osher, and Fatemi [42] (also see [24]), but denoise \mathbf{g} by another model. One such model is to define the recovered chromaticity \mathbf{u} as a (generalized) p -harmonic map [49, 50]

$$(1.6) \quad \mathbf{u} = \operatorname{argmin}_{\mathbf{v} \in W^{1,p}(\Omega; S^2)} J_{p,\lambda}(\mathbf{v}) \quad \text{for } p \geq 1,$$

where

$$(1.7) \quad J_{p,\lambda}(\mathbf{v}) := E_p(\mathbf{v}) + \frac{\lambda}{2} \int_{\Omega} |\mathbf{v} - \mathbf{g}|^2 dx \quad \text{for } \lambda > 0.$$

In particular, the cases $1 \leq p < 2$ are the most important and interesting since the recovered images keep geometric information such as edges and corners of the noisy color images. We shall call (1.6) the p -harmonic model for color image denoising. We also remark that the second term on the right-hand side of (1.7) is often called a fidelity term. As in the TV model [42, 24], the parameter λ controls the trade-off between goodness of fit-to-the-data and variability in \mathbf{u} .

Again, to find a solution for the p -harmonic map model (1.6), we consider the gradient flow (heat flow) for the energy functional $J_{p,\lambda}$, which is given by

$$(1.8) \quad \mathbf{u}_t - \Delta_p \mathbf{u} + \lambda(\mathbf{u} - \mathbf{g}) = \mu_{p,\lambda} \mathbf{u} \quad \text{in } \Omega_T,$$

$$(1.9) \quad |\mathbf{u}| = 1 \quad \text{in } \Omega_T,$$

$$(1.10) \quad \mathcal{B}_p \mathbf{n} = 0 \quad \text{on } \partial\Omega_T := \partial\Omega \times (0, T),$$

$$(1.11) \quad \mathbf{u} = \mathbf{u}_0 \quad \text{on } \Omega \times \{t = 0\},$$

where

$$(1.12) \quad \mathcal{B}_p := |\nabla \mathbf{u}|^{p-2} \nabla \mathbf{u}, \quad \mu_{p,\lambda} := |\nabla \mathbf{u}|^p + \lambda(1 - \mathbf{u} \cdot \mathbf{g}).$$

The goals of this paper are twofold. First, we shall present a general theory of weak solutions for the parabolic system (1.8)–(1.11) for the whole spectrum $1 \leq p < \infty$. To the best of our knowledge, there is no theory known in the literature for the 1-harmonic map and its heat flow, which on the other hand is the most important (and most difficult) case for the color image denoising application. So our theory and results fill this void. Furthermore, our theory handles the p -harmonic map heat flow (1.8)–(1.11) for all $1 \leq p < \infty$ in a unified fashion, rather than treating the system separately for different values of p and using different methods (cf. [11, 12, 13, 16, 18, 21, 22, 30, 31, 34, 41, 46, 48]). Second, based on the above theoretical work, we also develop and analyze a practical fully discrete finite element method for approximating the solutions of the p -harmonic map heat flow (1.8)–(1.11).

We now highlight the main ideas and key steps of our approach. Notice that there are two nonlinear terms in the p -harmonic flow: the p -Laplace term and the right-hand side due to the nonconvex constraint $|\mathbf{u}| = 1$, so the main difficulties are how to handle these two terms and how to pass to the limit in these two terms when a

compactness argument is employed. To handle the degeneracy of the p -Laplace term, we approximate the p -energy $E_p(\mathbf{v})$ by the following regularized energy

$$(1.13) \quad \begin{aligned} E_p^\varepsilon(\mathbf{v}) &:= \frac{b_p(\varepsilon)}{2} \int_\Omega |\nabla \mathbf{v}|^2 dx + \frac{1}{p} \int_\Omega |\nabla \mathbf{v}|_\varepsilon^p dx \\ &= \int_\Omega \left\{ \frac{b_p(\varepsilon)}{2} |\nabla \mathbf{v}|^2 + \frac{1}{p} [|\nabla \mathbf{v}|^2 + a_p(\varepsilon)^2]^{\frac{p}{2}} \right\} dx, \end{aligned}$$

where $\varepsilon > 0$ and

$$(1.14) \quad a_p(\varepsilon) := \begin{cases} 0 & \text{if } 2 \leq p < \infty, \\ \varepsilon & \text{if } 1 \leq p < 2, \end{cases} \quad b_p(\varepsilon) := \varepsilon^\alpha \quad \text{for } 1 \leq p < \infty$$

for some $\alpha > 0$. Here and in the rest of this paper we adopt the shorthand notation

$$(1.15) \quad |\nabla \mathbf{v}|_\varepsilon := \sqrt{|\nabla \mathbf{v}|^2 + a_p(\varepsilon)^2}.$$

To handle the nonconvex constraint $|\mathbf{u}| = 1$, we approximate it by the well-known Ginzburg–Landau penalization [7], that is, we abandon the exact constraint, but enforce it approximately by adding a penalization term to the regularized p -energy E_p^ε . To this end, we replace the energy E_p^ε by

$$(1.16) \quad E_p^{\varepsilon,\delta}(\mathbf{v}) := E_p^\varepsilon(\mathbf{v}) + L^\delta(\mathbf{v}) \quad \text{for } \varepsilon, \delta > 0,$$

where

$$(1.17) \quad L^\delta(\mathbf{v}) := \frac{1}{\delta} \int_\Omega F(\mathbf{v}) dx, \quad F(\mathbf{v}) := \frac{1}{4} (|\mathbf{v}|^2 - 1)^2 \quad \forall \delta > 0, \mathbf{v} \in \mathbf{R}^n.$$

So the idea is, as δ gets smaller and smaller, the energy functional $E_p^{\varepsilon,\delta}$ becomes more and more favorable for maps \mathbf{u} which take values close to the unit sphere S^{n-1} .

Consequently, the regularized model for the p -harmonic model (1.6) (with general m and n) reads

$$(1.18) \quad \mathbf{u}^{\varepsilon,\delta} = \operatorname{argmin}_{\mathbf{v} \in W^{1,p}(\Omega; \mathbf{R}^n)} J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{v}) \quad \text{for } p \geq 1,$$

where

$$(1.19) \quad J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{v}) := E_p^{\varepsilon,\delta}(\mathbf{v}) + \frac{\lambda}{2} \int_\Omega |\mathbf{v} - \mathbf{g}|^2 dx.$$

In addition, the gradient flow for the regularized energy functional $J_{p,\lambda}^{\varepsilon,\delta}$ is given by

$$(1.20) \quad \mathbf{u}_t^{\varepsilon,\delta} - \Delta_p^\varepsilon \mathbf{u}^{\varepsilon,\delta} + \frac{1}{\delta} (|\mathbf{u}^{\varepsilon,\delta}|^2 - 1) \mathbf{u}^{\varepsilon,\delta} + \lambda (\mathbf{u}^{\varepsilon,\delta} - \mathbf{g}) = 0 \quad \text{in } \Omega_T,$$

$$(1.21) \quad \mathcal{B}_p^{\varepsilon,\delta} \mathbf{n} = 0 \quad \text{on } \partial\Omega_T,$$

$$(1.22) \quad \mathbf{u}^{\varepsilon,\delta} = \mathbf{u}_0 \quad \text{on } \Omega \times \{t = 0\},$$

which is an approximation to the original flow (1.8)–(1.11), where

$$(1.23) \quad \mathcal{B}_p^{\varepsilon,\delta} := b_p(\varepsilon) \nabla \mathbf{u}^{\varepsilon,\delta} + |\nabla \mathbf{u}^{\varepsilon,\delta}|_\varepsilon^{p-2} \nabla \mathbf{u}^{\varepsilon,\delta}, \quad \Delta_p^\varepsilon \mathbf{u}^{\varepsilon,\delta} := \operatorname{div} \mathcal{B}_p^{\varepsilon,\delta}.$$

After having introduced the regularized flow (1.20)–(1.22), the next step is to analyze this regularized flow, in particular, to derive uniform (in ε and δ) a priori estimates. Finally, we pass to the limit in (1.20)–(1.22), first letting $\delta \rightarrow 0$ and then setting $\varepsilon \rightarrow 0$. As pointed out earlier, the main difficulty here is passing to the limit in two nonlinear terms on the left-hand side of (1.20). For $1 < p < \infty$, this will be done using a compactness technique and exploiting the symmetries of the unit sphere S^{n-1} , as done in [11, 12, 13, 21, 32, 41, 46]. However, for $p = 1$, since $L^1(\Omega)$ is not a reflexive Banach space, instead of working in the Sobolev space $W^{1,1}(\Omega)$, we are forced to work in $BV(\Omega)$, the space of functions of bounded variation, since solutions of the original 1-harmonic map heat flow (1.8)–(1.11) belong only to $L^\infty((0, T); [BV(\Omega)]^n)$ in general. This lack of regularity makes the analysis for $p = 1$ much more delicate than that for $1 < p < \infty$.

We note also that the regularized flow (1.20)–(1.22) not only plays an important role for proving existence of weak solutions for the flow (1.8)–(1.11), but also provides a practical and convenient formulation for approximating the solutions. The second goal of this paper is to develop a practical fully discrete finite element method for approximating solutions of the regularized flow (1.20)–(1.22); and hence, approximating solutions of the original p -harmonic flow (1.8)–(1.11) via the regularized flow. It is well known that explicitly enforcing the nonconvex constraint $|\mathbf{u}| = 1$ at the discrete level is hard to achieve. The penalization used in (1.20) allows one to get around this numerical difficulty at the expense of introducing an additional scale, δ . As expected, the numerical difficulty now is to control the dependence of the regularized solutions on δ and to establish scaling laws which relate the numerical scales (spatial and temporal mesh sizes) to the penalization scale δ for both stability and accuracy concerns. We refer to [6] and the references therein for more discussions in this direction and discussions on a related problem arising from liquid crystal applications. Borrowing a terminology from the phase transition of materials science, the regularized flow (1.20)–(1.22) may be regarded as a “diffuse interface” model for the original “sharp interface” model (1.8)–(1.11), and the “diffuse interface” is represented by the region $\{|\mathbf{u}^{\varepsilon, \delta}| < 1 - \delta\}$.

The remaining part of the paper is organized as follows. In section 2 we collect some known results and facts, which will be used in the later sections. In section 3 we present a complete analysis for the regularized flow (1.20)–(1.22), which includes proving its well-posedness, an energy law, a maximum principle, and uniform (in ε and δ) a priori estimates. As expected, these uniform a priori estimates serve as the basis for carrying out the energy method and the compactness arguments in the later sections. In section 4 we pass to the limit in (1.20)–(1.22) as $\delta \rightarrow 0$ for each fixed $\varepsilon > 0$. As in [13, 12, 32, 41], the main idea is to exploit the monotonicity of the operator Δ_p^ε and the symmetries of the unit sphere S^{n-1} . Sections 5 and 6 are devoted to passing to the limit as $\varepsilon \rightarrow 0$ in the ε -dependent limiting system obtained in section 4. For $1 < p < \infty$, this will be done by following the idea of section 4. However, for $p = 1$ the analysis becomes much more delicate because the nonreflexivity of $W^{1,1}(\Omega)$ forces us to work in the $BV(\Omega)$ space. The main idea used to develop a BV -solution concept is to exploit the properties of measures of the forms $\mathcal{A} \cdot \nabla \mathbf{v}$ and $\mathcal{A} \wedge \nabla \mathbf{v}$, which pair a divergence- L^1 , or a divergence-measure, tensor field \mathcal{A} and a BV -vector field \mathbf{v} . Finally, based on the theoretical results of sections 3–6, in section 7 we propose and analyze a practical fully discrete finite element method for approximating solutions of the p -harmonic flow (1.8)–(1.11) via the regularized flow (1.20)–(1.22). It is proved that the proposed numerical scheme

satisfies a discrete energy inequality, which mimics the differential energy inequality, and this leads to uniform (in ε and δ) a priori estimates and the convergence of the numerical approximations to the solutions of the flow (1.8)–(1.11) as the spatial and temporal mesh sizes, and the parameters ε and δ , all tend to zero.

2. Preliminaries. The standard notation for spaces is adopted in this paper (cf. [1, 3]). For example, $W^{k,p}(\Omega)$, $k \geq 0$ integer and $1 \leq p < \infty$, denotes the Sobolev spaces over the domain Ω and $\|\cdot\|_{W^{k,p}}$ denotes its norm. $W^{0,p}(\Omega) = L^p(\Omega)$ and $W^{k,2}(\Omega) = H^k(\Omega)$ are also used. $L^q((0, T); W^{k,p}(\Omega, \mathbf{R}^n))$ denotes the space of vector-valued functions (or maps), whose $W^{k,p}(\Omega)$ -norm is L^q -integrable as a function of t over the interval $(0, T)$, and $\|\cdot\|_{L^q(W^{k,p})} := (\int_0^T \|\cdot\|_{W^{k,p}}^q dt)^{\frac{1}{q}}$ for $q \in [1, \infty)$ denotes its norm, with the standard modification for $q = \infty$. We use $\langle \cdot, \cdot \rangle$ to denote a generic dual product between elements of a Banach space X and its dual space X' .

In addition, $\mathcal{M}(\Omega)$ (resp., $[\mathcal{M}(\Omega)]^n$) denotes the space of real-valued (resp., \mathbf{R}^n -valued) finite Radon measures on Ω . Recall that $\mathcal{M}(\Omega)$ is the dual space of $C_0(\Omega)$ (cf. [3]). For a positive non-Lebesgue measure $\mu \in \mathcal{M}(\Omega)$, $L^p(\Omega, \mu)$ is used to denote the space of L^p -integrable functions with respect to the measure μ , and for $f \in L^1(\Omega, \mu)$, the measure $f\mu$ is defined by

$$f\mu(A) := \int_A f d\mu \quad \text{for any Borel set } A \subset \Omega.$$

For any $\mu \in \mathcal{M}(\Omega)$, $\mu = \mu^a + \mu^s$ denotes its Radon–Nikodým decomposition, where μ^a and μ^s , respectively, denotes the absolute continuous part and the singular part of μ with respect to the Lebesgue measure \mathcal{L}^n .

Furthermore, $BV(\Omega)$ is used to denote the space of functions of bounded variation. Recall that a function $u \in L^1(\Omega)$ is called a function of *bounded variation* if all of its first order partial derivatives (in the distributional sense) are measures with finite total variations in Ω . Hence, the gradient of such a function u , denoted by Du , is a bounded vector-valued measure, with the finite total variation

$$(2.1) \quad |Du| \equiv |Du|(\Omega) := \sup \left\{ \int_{\Omega} -u \operatorname{div} \mathbf{v} dx; \mathbf{v} \in [C_0^1(\Omega)]^n, \|\mathbf{v}\|_{L^\infty} \leq 1 \right\}.$$

$BV(\Omega)$ is known to be a Banach space endowed with the norm

$$(2.2) \quad \|u\|_{BV} := \|u\|_{L^1} + |Du|.$$

For any $u \in BV(\Omega)$, $(Du)^a$ and $|Du|^a$ are used to denote, respectively, the absolute continuous part of Du and $|Du|$ with respect to the Lebesgue measure \mathcal{L}^n . We refer to [3] for detailed discussions about the space $BV(\Omega)$ and properties of BV functions.

For any vector $\mathbf{a} \in \mathbf{R}^n$ and any matrices $\mathcal{A}, \mathcal{B} \in \mathbf{R}^{n \times m}$ with j th column vectors $\mathcal{A}^{(j)}$ and $\mathcal{B}^{(j)}$, respectively, we define the following wedge products:

$$(2.3) \quad \mathcal{A} \wedge \mathbf{a} := [\mathcal{A}^{(1)} \wedge \mathbf{a}, \dots, \mathcal{A}^{(m)} \wedge \mathbf{a}], \quad \mathbf{a} \wedge \mathcal{A} := [\mathbf{a} \wedge \mathcal{A}^{(1)}, \dots, \mathbf{a} \wedge \mathcal{A}^{(m)}],$$

$$(2.4) \quad \mathcal{A} \wedge \mathcal{B} := \mathcal{A}^{(1)} \wedge \mathcal{B}^{(1)} + \dots + \mathcal{A}^{(m)} \wedge \mathcal{B}^{(m)}.$$

We point out that (2.4) defines $\mathcal{A} \wedge \mathcal{B}$ to be a vector instead of a matrix, which seems not to be natural. However, we shall see in section 6 that it turns out that this is a convenient and useful notation.

We conclude this section by citing some known results which will be used in the later sections. The first result is known as “the decisive monotonicity trick,” and its proof can be found in [52].

LEMMA 2.1. *Let X be a reflexive Banach space, and X' denote the dual space of X . Suppose that an operator $\mathcal{F} : X \rightarrow X'$ satisfies*

- (i) \mathcal{F} is monotone on X , that is, $\langle \mathcal{F}u - \mathcal{F}v, u - v \rangle \geq 0$ for all $u, v \in X$;
- (ii) \mathcal{F} is hemicontinuous, that is, the function $t \mapsto \langle \mathcal{F}(u + tv), w \rangle$ is continuous on $[0, 1]$ for all $u, v, w \in X$.

In addition, suppose that $u_k \rightarrow u$ and $\mathcal{F}u_k \rightarrow f$ weakly in X and X' , respectively, as $n \rightarrow \infty$, and

$$\overline{\lim}_{k \rightarrow \infty} \langle \mathcal{F}u_k, u_k \rangle \leq \langle f, u \rangle.$$

Then

$$\mathcal{F}u = f.$$

The second lemma is a compactness result, it can be proved following the proof of Theorem 2.1 of [12] (also see Theorem 3 of Chapter 4 of [20]) using the fact that the operator Δ_p^ε is uniformly elliptic.

LEMMA 2.2. *For $1 \leq p < \infty$, let $p^* = \max\{2, p\}$. For a fixed $\varepsilon > 0$, let $\{\mathbf{w}^{\varepsilon, \delta}\}_{\delta > 0}$ be bounded in $L^\infty((0, T); W^{1, p^*}(\Omega, \mathbf{R}^n))$, and let $\{\mathbf{f}^{\varepsilon, \delta}\}_{\delta > 0}$ be bounded in $L^1((0, T); L^1(\Omega, \mathbf{R}^n))$, both uniformly in δ . Moreover, suppose that $\{\frac{\partial \mathbf{w}^{\varepsilon, \delta}}{\partial t}\}_{\delta > 0}$ is bounded in $L^2((0, T); L^2(\Omega, \mathbf{R}^n))$ uniformly in δ and $\mathbf{w}^{\varepsilon, \delta}$ satisfies the equation*

$$\frac{\partial \mathbf{w}^{\varepsilon, \delta}}{\partial t} - \Delta_p^\varepsilon \mathbf{w}^{\varepsilon, \delta} = \mathbf{f}^{\varepsilon, \delta} \quad \text{in } \Omega_T, \quad \delta > 0,$$

in the distributional sense. Then $\{\mathbf{w}^{\varepsilon, \delta}\}_{\delta > 0}$ is precompact in $L^q((0, T); W^{1, q}(\Omega, \mathbf{R}^n))$ for all $1 \leq q < p^$.*

The last lemma is a variation of Lemma 2.2, and its proof can be found in [38] (also see Lemma 9 of [40]).

LEMMA 2.3. *For $1 < p < \infty$, let $\{\mathbf{w}^\varepsilon\}_{\varepsilon > 0}$ be bounded in $L^\infty((0, T); W^{1, p}(\Omega, \mathbf{R}^n))$, and $\{\mathbf{f}^\varepsilon\}$ be bounded in $L^1((0, T); L^1(\Omega, \mathbf{R}^n))$, both uniformly in ε . Moreover, suppose that $\{\frac{\partial \mathbf{w}^\varepsilon}{\partial t}\}$ is bounded uniformly in ε in $L^2((0, T); L^2(\Omega, \mathbf{R}^n))$ and \mathbf{w}^ε satisfies the equation*

$$\frac{\partial \mathbf{w}^\varepsilon}{\partial t} - \Delta_p^\varepsilon \mathbf{w}^\varepsilon = \mathbf{f}^\varepsilon \quad \text{in } \Omega_T, \quad \varepsilon > 0,$$

in the distributional sense. Then $\{\mathbf{w}^\varepsilon\}_{\varepsilon > 0}$ is precompact in $L^q((0, T); W^{1, q}(\Omega, \mathbf{R}^n))$ for all $1 \leq q < p$.

Remark 2.1. In Lemma 2.2, $\varepsilon > 0$ is a fixed parameter and the differential operator Δ_p^ε does not depend on the variable index δ . It can be shown that (cf. Theorem 2.1 of [12]) that the lemma still holds for $\varepsilon = 0$ when $p \geq 2$. On the other hand, ε is a variable index in Lemma 2.3, and the operator Δ_p^ε also depends on the variable index ε .

3. Well-posedness of the regularized flow (1.20)–(1.22). In this section we shall analyze the regularized flow (1.20)–(1.22) for each fixed pair of positive numbers (ε, δ) . We establish an energy law, a maximum principle, uniform (in both ε and δ) a priori estimates, and existence and uniqueness of weak and classical solutions. We begin with a couple of definitions for solutions to (1.20)–(1.22).

DEFINITION 3.1. For $1 \leq p < \infty$, a map $\mathbf{u}^{\varepsilon, \delta} : \Omega_T \rightarrow \mathbf{R}^n$ is called a global weak solution to (1.20)–(1.22) if

- (i) $\mathbf{u}^{\varepsilon, \delta} \in L^\infty((0, T); W^{1, p^*}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$, for $p^* = \max\{2, p\}$;
- (ii) $|\mathbf{u}^{\varepsilon, \delta}| \leq 1$ a.e. in Ω_T ;
- (iii) $\mathbf{u}^{\varepsilon, \delta}$ satisfies (1.20)–(1.22) in the distributional sense.

DEFINITION 3.2. A weak solution $\mathbf{u}^{\varepsilon, \delta}$ to (1.20)–(1.22) is called a strong solution if $\mathbf{u}^{\varepsilon, \delta} \in L^p((0, T); W^{2, p^*}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^{p^*}(\Omega, \mathbf{R}^n))$. It is called a regular solution if in addition $\mathbf{u}^{\varepsilon, \delta} \in H^1((0, T); W^{1, p^*}(\Omega, \mathbf{R}^n))$.

3.1. Energy law and a priori estimates. Since (1.20)–(1.22) is the gradient flow for the functional $J_{p, \lambda}^{\varepsilon, \delta}$, its regular solutions must satisfy a dissipative energy law. Indeed, we have the following lemma.

LEMMA 3.3. Let \mathbf{u}_0 and \mathbf{g} be sufficiently smooth, and suppose that $\mathbf{u}^{\varepsilon, \delta}$ is a regular solution to (1.20)–(1.22). Then $\mathbf{u}^{\varepsilon, \delta}$ satisfies the following energy law:

$$(3.1) \quad J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}^{\varepsilon, \delta}(s)) + \int_0^s \|\mathbf{u}_t^{\varepsilon, \delta}(t)\|_{L^2}^2 dt = J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}_0) \quad \text{for a.e. } s \in [0, T].$$

Proof. Testing (1.20) with $\mathbf{u}_t^{\varepsilon, \delta}$ we get

$$\begin{aligned} \|\mathbf{u}_t^{\varepsilon, \delta}(t)\|_{L^2}^2 + \frac{d}{dt} \int_{\Omega} \left\{ \frac{b_p(\varepsilon)}{2} |\nabla \mathbf{u}^{\varepsilon, \delta}(t)|^2 + \frac{1}{p} |\nabla \mathbf{u}^{\varepsilon, \delta}(t)|_{\varepsilon}^p \right. \\ \left. + \frac{1}{\delta} F(\mathbf{u}^{\varepsilon, \delta}(t)) + \frac{\lambda}{2} |\mathbf{u}^{\varepsilon, \delta}(t) - \mathbf{g}|^2 \right\} dx = 0. \end{aligned}$$

The desired identity (3.1) then follows from integrating the above equation in t over the interval $[0, s]$ and using the definition of $J_{p, \lambda}^{\varepsilon, \delta}$. \square

The above energy law immediately implies the following uniform (in ε and δ) a priori estimates.

COROLLARY 3.4. Suppose that \mathbf{u}_0 and \mathbf{g} satisfy

$$(3.2) \quad J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}_0) \leq c_0$$

for some positive constant c_0 independent of ε and δ . Then there exists another positive constant $C := C(p, \lambda, c_0)$ which is also independent of ε and δ such that

$$(3.3) \quad \|\mathbf{u}^{\varepsilon, \delta}\|_{L^\infty(W^{1, p})} + \|\mathbf{u}^{\varepsilon, \delta}\|_{H^1(L^2)} \leq C \quad \text{for } 1 \leq p < \infty,$$

$$(3.4) \quad \delta^{-\frac{1}{2}} \|\mathbf{u}^{\varepsilon, \delta}\|^2 - 1 \|_{L^\infty(L^2)} \leq C \quad \text{for } 1 \leq p < \infty,$$

$$(3.5) \quad \|\nabla \mathbf{u}^{\varepsilon, \delta}\|_{\varepsilon}^{p-2} \|\nabla \mathbf{u}^{\varepsilon, \delta}\|_{L^\infty(L^{p'})} \leq C \quad \text{for } p' = \frac{p}{p-1}, \quad 1 \leq p < \infty,$$

$$(3.6) \quad \sqrt{b_p(\varepsilon)} \|\nabla \mathbf{u}^{\varepsilon, \delta}\|_{L^\infty(L^2)} \leq C \quad \text{for } 1 \leq p < \infty.$$

Next, using a test function technique of [12], we show that the modulus $|\mathbf{u}^{\varepsilon, \delta}|$ of every weak solution $\mathbf{u}^{\varepsilon, \delta}$ to (1.20)–(1.22) satisfies a maximum principle.

LEMMA 3.5. Suppose that $|\mathbf{g}| \leq 1$ and $|\mathbf{u}_0| \leq 1$ a.e. in Ω . Then weak solutions to (1.20)–(1.22) satisfy $|\mathbf{u}^{\varepsilon, \delta}| \leq 1$ a.e. in Ω_T .

Proof. Define the function

$$(3.7) \quad \chi(z) := \frac{(z-1)_+}{z} = \begin{cases} 0 & \text{for } 0 \leq z \leq 1, \\ \frac{z-1}{z} & \text{for } z > 1. \end{cases}$$

It is easy to check that χ is a nonnegative monotone increasing function on the interval $[0, \infty)$.

Now testing (1.20) with $\mathbf{v} := \mathbf{u}^{\varepsilon, \delta} \chi(|\mathbf{u}^{\varepsilon, \delta}|)$ we get

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\{|\mathbf{u}^{\varepsilon, \delta}| > 1\}} (|\mathbf{u}^{\varepsilon, \delta}| - 1)^2 dx + \int_{\{|\mathbf{u}^{\varepsilon, \delta}| > 1\}} \left\{ b_p(\varepsilon) |\nabla \mathbf{u}^{\varepsilon, \delta}|^2 + |\nabla \mathbf{u}^{\varepsilon, \delta}|_\varepsilon^{p-2} |\nabla \mathbf{u}^{\varepsilon, \delta}|^2 \right. \\ & \quad \left. + \frac{1}{\delta} (|\mathbf{u}^{\varepsilon, \delta}|^2 - 1) |\mathbf{u}^{\varepsilon, \delta}|^2 + \lambda (|\mathbf{u}^{\varepsilon, \delta}|^2 - \mathbf{u}^{\varepsilon, \delta} \cdot \mathbf{g}) \right\} \chi(|\mathbf{u}^{\varepsilon, \delta}|) dx \\ & \quad + \frac{1}{4} \int_{\{|\mathbf{u}^{\varepsilon, \delta}| > 1\}} \frac{[b_p(\varepsilon) + |\nabla \mathbf{u}^{\varepsilon, \delta}|_\varepsilon^{p-2}] |\nabla |\mathbf{u}^{\varepsilon, \delta}|^2|^2}{|\mathbf{u}^{\varepsilon, \delta}|^3} dx = 0. \end{aligned}$$

Since $\mathbf{u}^{\varepsilon, \delta} \cdot \mathbf{g} \leq |\mathbf{u}^{\varepsilon, \delta}| \cdot |\mathbf{g}| \leq |\mathbf{u}^{\varepsilon, \delta}|$, the second integral is nonnegative. The assertion then follows from integrating the above inequality and using the assumption $|\mathbf{u}_0| \leq 1$ a.e. in Ω . \square

3.2. Existence of global weak and classical solutions. We now state and prove the existence of global weak and classical solutions to the regularized flow (1.20)–(1.22) for each fixed pair of positive numbers (ε, δ) . Since $-\Delta_p^\varepsilon$ is uniformly elliptic for all $1 \leq p < \infty$, the existence of classical solutions follows immediately from the classical theory of parabolic partial differential equations (cf. [35]).

THEOREM 3.6. *Let $\Omega \subset \mathbf{R}^m$ be a bounded domain with a smooth boundary. Suppose that \mathbf{u}_0 and \mathbf{g} are sufficiently smooth functions (say, $\mathbf{u}_0, \mathbf{g} \in [C^3(\bar{\Omega})]^n$) and satisfy (3.2). Then for each fixed pair of positive numbers (ε, δ) , the regularized flow (1.20)–(1.22) possesses a unique global classical solution $\mathbf{u}^{\varepsilon, \delta}$. Moreover, $\mathbf{u}^{\varepsilon, \delta}$ satisfies the following energy law:*

$$(3.8) \quad J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}^{\varepsilon, \delta}(s)) + \int_0^s \|\mathbf{u}_t^{\varepsilon, \delta}(t)\|_{L^2}^2 dt = J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}_0) \quad \text{for a.e. } s \in [0, T].$$

Proof. The existence and uniqueness follow from an application of the standard results for parabolic systems; see Chapter 5 of [35]. Equation (3.8) follows from Lemma 3.3. \square

For less regular functions \mathbf{u}_0 and \mathbf{g} , we have the following weaker result.

THEOREM 3.7. *Let $\Omega \subset \mathbf{R}^m$ be a bounded domain with smooth boundary. Suppose that $|\mathbf{u}_0| \leq 1$ and $|\mathbf{g}| \leq 1$ a.e. in Ω and $J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}_0) < \infty$. Then the regularized flow (1.20)–(1.22) has a unique global weak solution $\mathbf{u}^{\varepsilon, \delta}$ in the sense of Definition 3.1. Moreover, $\mathbf{u}^{\varepsilon, \delta}$ satisfies the following energy inequality:*

$$(3.9) \quad J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}^{\varepsilon, \delta}(s)) + \int_0^s \|\mathbf{u}_t^{\varepsilon, \delta}(t)\|_{L^2}^2 dt \leq J_{p, \lambda}^{\varepsilon, \delta}(\mathbf{u}_0) \quad \forall s \in [0, T].$$

Proof. Let η_ρ denote any well-known mollifier (cf. [35]), and let $\mathbf{u}_{0, \rho} := \eta_\rho * \mathbf{u}_0$ and $\mathbf{g}_\rho := \eta_\rho * \mathbf{g}$ denote the mollifications of \mathbf{u}_0 and \mathbf{g} , respectively. Let $\mathbf{u}_\rho^{\varepsilon, \delta}$ denote the classical solution to (1.20)–(1.22) corresponding to the smoothed datum functions $\mathbf{u}_{0, \rho}$ and \mathbf{g}_ρ . Hence, $\mathbf{u}_\rho^{\varepsilon, \delta}$ satisfies the energy law (3.8) with $\mathbf{u}_{0, \rho}$ and \mathbf{g}_ρ in the place of \mathbf{u}_0 and \mathbf{g} .

From Lemma 3.5 we know that $|\mathbf{u}_\rho^{\varepsilon, \delta}|$ satisfies the maximum principle, thus,

$$\max_{(x, t) \in \bar{\Omega}_T} |\mathbf{u}_\rho^{\varepsilon, \delta}(x, t)| \leq 1.$$

Next, since

$$\lim_{\rho \rightarrow 0} J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_{0,\rho}^\varepsilon) = J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_0) < \infty,$$

the energy law for $\mathbf{u}_\rho^{\varepsilon,\delta}$ immediately implies that $\mathbf{u}_\rho^{\varepsilon,\delta}$ satisfies estimates (3.3)–(3.6), uniformly in ρ and δ .

The remainder of the proof is to extract a convergent subsequence of $\{\mathbf{u}_\rho^{\varepsilon,\delta}\}_{\rho>0}$ and to pass to the limit as $\rho \rightarrow 0$ in the weak formulation of (1.20). This can be done easily in all terms of (1.20) except the nonlinear term in Δ_p^ε (see (1.23)). To overcome the difficulty, we appeal to Minty’s trick or the “decisive monotonicity trick” as described in Lemma 2.1. Since this part of proof is the same as that of Theorem 1.5 of [12], we omit it and refer to [12] for the details.

The uniqueness of weak solutions follows from a standard perturbation argument and using the fact that Δ_p^ε is a monotone operator. Finally, (3.9) follows from setting $\rho \rightarrow 0$ in the energy law (3.8) for $\mathbf{u}_\rho^{\varepsilon,\delta}$ and from using the lower semicontinuity of $J_{p,\lambda}^{\varepsilon,\delta}$ and the L^2 -norm with respect to L^2 weak convergence. \square

We conclude this section with some remarks.

Remark 3.1. (a) Since $W^{1,1}(\Omega)$ is not a reflexive Banach space, Minty’s trick would not apply in the case $p = 1$ without the help of the $b_p(\varepsilon)\Delta$ term in the operator Δ_p^ε (see (1.23)). On the other hand, in the presence of this term, Minty’s trick does apply since we now deal with the Sobolev space $H^1(\Omega)$ instead of $W^{1,1}(\Omega)$. In fact, if $b_p(\varepsilon) = 0$ in (1.20), BV solutions are what one can only expect in general for the regularized flow (1.20)–(1.22) in the case $p = 1$ (cf. [24]).

(b) We also point out that using nonzero parameter $b_p(\varepsilon)$ is not necessary in the case $1 < p < \infty$. For example, the conclusion of Theorem 3.7 still holds if we replace the above $b_p(\varepsilon)$ and p^* by

$$\widehat{b}_p(\varepsilon) = \begin{cases} 0 & \text{if } 1 < p < \infty, \\ \varepsilon^\alpha & \text{if } p = 1 \quad (\alpha > 0) \end{cases} \quad \text{and} \quad \widehat{p}^* = \begin{cases} p & \text{if } 1 < p < \infty, \\ 2 & \text{if } p = 1. \end{cases}$$

On the other hand, the conclusion of Theorem 3.6 may no longer be true after this modification.

(c) Theorem 3.7 also holds when Ω is a bounded Lipschitz domain. One way to prove this assertion is to use the Galerkin method as done in [12], or to use the finite element method to be introduced in section 7.

4. Passing to the limit as $\delta \rightarrow 0$. The goal of this section is to derive the limiting flow of (1.20)–(1.22) as $\delta \rightarrow 0$ for each fixed $\varepsilon > 0$. As in [11, 12, 13, 31, 41, 46], the key ideas for passing to the limit are to use the compactness result of Lemma 2.2 and to exploit the symmetries of the unit sphere S^{m-1} . Our main result of this section is the following existence theorem.

THEOREM 4.1. *Let $p^* := \max\{2, p\}$ and let $\Omega \subset \mathbf{R}^m$ be a bounded Lipschitz domain. For $1 \leq p < \infty$, suppose that $\mathbf{u}_0 \in W^{1,p^*}(\Omega, \mathbf{R}^n)$, $|\mathbf{u}_0| = 1$ and $|\mathbf{g}| \leq 1$ in Ω . Then there exists a map $\mathbf{u}^\varepsilon \in L^\infty((0, T); W^{1,p^*}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$ such that*

$$(4.1) \quad |\mathbf{u}^\varepsilon| = 1 \quad \text{a.e. in } \Omega_T,$$

and \mathbf{u}^ε is a weak solution (in the distributional sense; see (4.29)) to the problem

$$(4.2) \quad \mathbf{u}_t^\varepsilon - \Delta_p^\varepsilon \mathbf{u}^\varepsilon + \lambda(\mathbf{u}^\varepsilon - \mathbf{g}) = \mu_{p,\lambda}^\varepsilon \mathbf{u}^\varepsilon \quad \text{in } \Omega_T,$$

$$(4.3) \quad \mathcal{B}_p^\varepsilon \mathbf{n} = 0 \quad \text{on } \partial\Omega_T,$$

$$(4.4) \quad \mathbf{u}^\varepsilon = \mathbf{u}_0 \quad \text{on } \Omega \times \{t = 0\},$$

where

$$(4.5) \quad \mu_{p,\lambda}^\varepsilon := b_p(\varepsilon)|\nabla \mathbf{u}^\varepsilon|^2 + |\nabla \mathbf{u}^\varepsilon|_\varepsilon^{p-2} |\nabla \mathbf{u}^\varepsilon|^2 + \lambda(1 - \mathbf{u}^\varepsilon \cdot \mathbf{g}).$$

Moreover, \mathbf{u}^ε satisfies the energy inequality

$$(4.6) \quad J_{p,\lambda}^\varepsilon(\mathbf{u}^\varepsilon(s)) + \int_0^s \|\mathbf{u}_t^\varepsilon(t)\|_{L^2}^2 dt \leq J_{p,\lambda}^\varepsilon(\mathbf{u}_0) \leq J_{p,\lambda}^1(\mathbf{u}_0) \quad \text{for a.e. } s \in [0, T]$$

and the additional estimates

$$(4.7) \quad \|\nabla \mathbf{u}^\varepsilon|_\varepsilon^{p-2} \nabla \mathbf{u}^\varepsilon\|_{L^\infty(L^{p'})} \leq C \quad \text{for } p' = \frac{p}{p-1}, \quad 1 \leq p < \infty,$$

$$(4.8) \quad \|\operatorname{div} \mathcal{B}_p^\varepsilon\|_{L^2(L^1)} \leq C \quad \text{for } 1 \leq p < \infty,$$

$$(4.9) \quad \|\operatorname{div} (\mathcal{B}_p^\varepsilon \wedge \mathbf{u}^\varepsilon)\|_{L^2(L^2)} \leq C \quad \text{for } 1 \leq p < \infty,$$

$$(4.10) \quad \sqrt{b_p(\varepsilon)} \|\nabla \mathbf{u}^\varepsilon\|_{L^\infty(L^2)} \leq C$$

for some positive ε -independent constant C . Here $\mathcal{B}_p^\varepsilon$ is the $n \times m$ matrix

$$(4.11) \quad \mathcal{B}_p^\varepsilon := b_p(\varepsilon) \nabla \mathbf{u}^\varepsilon + |\nabla \mathbf{u}^\varepsilon|_\varepsilon^{p-2} \nabla \mathbf{u}^\varepsilon$$

and

$$(4.12) \quad J_{p,\lambda}^\varepsilon(\mathbf{v}) := E_p^\varepsilon(\mathbf{v}) + \frac{\lambda}{2} \int_\Omega |\mathbf{v} - \mathbf{g}|^2 dx.$$

Proof. In light of Remark 3.1(c) and the density argument, without loss of the generality we assume that Ω is a bounded smooth domain, and \mathbf{u}_0 and \mathbf{g} are smooth functions. On noting that $|\mathbf{u}_0| = 1$ implies that $L^\delta(\mathbf{u}_0) = 0$, then $E_p^{\varepsilon,\delta}(\mathbf{u}_0) = E_p^\varepsilon(\mathbf{u}_0)$ in (1.16). Since $E_p^\varepsilon(\mathbf{u}_0) \leq E_p^1(\mathbf{u}_0)$, hence the assumptions on \mathbf{u}_0 and \mathbf{g} ensure (3.2) holds. Since the proof is long, we divide it into three steps.

Step 1: Extracting a convergent subsequence. Let $\mathbf{u}^{\varepsilon,\delta}$ be the weak solution solution to (1.20)–(1.22) whose existence is established in Theorem 3.7. Since $E_p^{\varepsilon,\delta}(\mathbf{u}_0) \leq E_p^1(\mathbf{u}_0)$, then $J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_0)$ is uniformly bounded with respect to ε and δ . Hence, (3.9) implies that $\mathbf{u}^{\varepsilon,\delta}$ satisfies the uniform (in both ε and δ) estimates (3.3)–(3.6) and the maximum principle $|\mathbf{u}^{\varepsilon,\delta}| \leq 1$ on $\overline{\Omega}_T$.

By the weak compactness of $W^{1,p}(\Omega)$ and Sobolev embedding (cf. [1, 45]), there exists a subsequence of $\{\mathbf{u}^{\varepsilon,\delta}\}_{\delta>0}$ (still denoted by the same notation) and a map

$\mathbf{u}^\varepsilon \in L^\infty((0, T); W^{1,p^*}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$ such that as $\delta \rightarrow 0$

$$(4.13) \quad \mathbf{u}^{\varepsilon,\delta} \longrightarrow \mathbf{u}^\varepsilon \quad \text{weakly* in } L^\infty((0, T); W^{1,p^*}(\Omega, \mathbf{R}^n)),$$

$$(4.14) \quad \text{strongly in } L^2((0, T); L^2(\Omega, \mathbf{R}^n)),$$

$$(4.15) \quad \text{a.e. in } \Omega_T,$$

$$(4.16) \quad \mathbf{u}_t^{\varepsilon,\delta} \longrightarrow \mathbf{u}_t^\varepsilon \quad \text{weakly in } L^2((0, T); L^2(\Omega, \mathbf{R}^n)),$$

$$(4.17) \quad |\mathbf{u}^{\varepsilon,\delta}| \longrightarrow 1 \quad \text{strongly in } L^2((0, T); L^2(\Omega)).$$

It follows immediately from (4.15) and (4.17) that

$$(4.18) \quad |\mathbf{u}^\varepsilon| = 1 \quad \text{a.e. in } \Omega_T.$$

Step 2: Wedge product technique and passing to the limit. Since $|\mathbf{u}^{\varepsilon,\delta}| \leq 1$ in $\overline{\Omega}_T$, an application of the Lebesgue-dominated convergence theorem yields, on noting (4.15), that

$$(4.19) \quad \mathbf{u}^{\varepsilon,\delta} \xrightarrow{\delta \searrow 0} \mathbf{u}^\varepsilon \quad \text{strongly in } L^r((0, T); L^r(\Omega, \mathbf{R}^n)) \quad \forall r \in [1, \infty).$$

Let

$$\mathbf{f}^{\varepsilon,\delta} := \frac{1}{\delta}(1 - |\mathbf{u}^{\varepsilon,\delta}|^2)\mathbf{u}^{\varepsilon,\delta} - \lambda(\mathbf{u}^{\varepsilon,\delta} - \mathbf{g}).$$

It follows from the estimate $|\mathbf{u}^{\varepsilon,\delta}| \leq 1$ and the inequality $|\mathbf{u}^{\varepsilon,\delta}| \leq \frac{1}{2}(1 + |\mathbf{u}^{\varepsilon,\delta}|^2)$ that

$$\begin{aligned} \frac{1}{\delta} \int_0^T \int_\Omega (1 - |\mathbf{u}^{\varepsilon,\delta}|^2) |\mathbf{u}^{\varepsilon,\delta}| \, dxdt &\leq \frac{1}{\delta} \int_0^T \int_\Omega (1 - |\mathbf{u}^{\varepsilon,\delta}|^2)^2 \, dxdt \\ &\quad + \frac{1}{\delta} \int_0^T \int_\Omega (1 - |\mathbf{u}^{\varepsilon,\delta}|^2) |\mathbf{u}^{\varepsilon,\delta}|^2 \, dxdt. \end{aligned}$$

Inequality (3.4) immediately implies that the first term on the right-hand side is uniformly bounded in δ . Testing (1.20) by $\mathbf{u}^{\varepsilon,\delta}$ and using estimates (3.3), (3.5), and (3.6), we conclude that the second term on the right-hand side is also uniformly bounded in δ . Hence, $\mathbf{f}^{\varepsilon,\delta}$ is uniformly bounded with respect to δ in $L^1((0, T); L^1(\Omega, \mathbf{R}^n))$. By Lemma 2.2 we have that

$$(4.20) \quad \nabla \mathbf{u}^{\varepsilon,\delta} \xrightarrow{\delta \searrow 0} \nabla \mathbf{u}^\varepsilon \quad \text{strongly in } L^q((0, T); L^q(\Omega, \mathbf{R}^{n \times m})) \quad \forall q \in [1, p^*).$$

This and (3.5) imply that

$$(4.21) \quad |\nabla \mathbf{u}^{\varepsilon,\delta}|_e^{p-2} \nabla \mathbf{u}^{\varepsilon,\delta} \xrightarrow{\delta \searrow 0} |\nabla \mathbf{u}^\varepsilon|_e^{p-2} \nabla \mathbf{u}^\varepsilon$$

weakly* in $L^\infty((0, T); L^{p'}(\Omega, \mathbf{R}^{n \times m}))$ with $p' = \frac{p}{p-1}$ if $p \neq 1$ and weakly* in $L^\infty((0, T); L^\infty(\Omega, \mathbf{R}^{n \times m}))$ if $p = 1$.

Next, taking the wedge product of (1.20) with $\mathbf{u}^{\varepsilon,\delta}$ yields

$$(4.22) \quad \mathbf{u}_t^{\varepsilon,\delta} \wedge \mathbf{u}^{\varepsilon,\delta} - \operatorname{div}(\mathcal{B}_p^{\varepsilon,\delta} \wedge \mathbf{u}^{\varepsilon,\delta}) - \lambda \mathbf{g} \wedge \mathbf{u}^{\varepsilon,\delta} = 0,$$

where $\mathcal{B}_p^{\varepsilon, \delta}$ is defined in (1.23).

Testing (4.22) with any $\mathbf{w} \in L^\infty((0, T); W^{1, p^*}(\Omega, \mathbf{R}^n))$ we get

$$(4.23) \quad \int_0^T \int_\Omega \left\{ (\mathbf{u}_t^{\varepsilon, \delta} \wedge \mathbf{u}^{\varepsilon, \delta}) \cdot \mathbf{w} + (\mathcal{B}_p^{\varepsilon, \delta} \wedge \mathbf{u}^{\varepsilon, \delta}) \cdot \nabla \mathbf{w} - \lambda (\mathbf{g} \wedge \mathbf{u}^{\varepsilon, \delta}) \cdot \mathbf{w} \right\} dxdt = 0.$$

It follows from setting $\delta \rightarrow 0$ in (4.23) and using (4.16), (4.19), and (4.21) that

$$(4.24) \quad \int_0^T \int_\Omega \left\{ (\mathbf{u}_t^\varepsilon \wedge \mathbf{u}^\varepsilon) \cdot \mathbf{w} + (\mathcal{B}_p^\varepsilon \wedge \mathbf{u}^\varepsilon) \cdot \nabla \mathbf{w} - \lambda (\mathbf{g} \wedge \mathbf{u}^\varepsilon) \cdot \mathbf{w} \right\} dxdt = 0,$$

where $\mathcal{B}_p^\varepsilon$ is given by (4.11).

Note that (4.18) implies

$$(4.25) \quad \mathbf{u}_t^\varepsilon \cdot \mathbf{u}^\varepsilon = 0, \quad (\mathcal{B}_p^\varepsilon)^T \mathbf{u}^\varepsilon = 0 \quad \text{a.e. in } \Omega_T.$$

This in turn yields the following identity

$$(4.26) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t^\varepsilon \cdot \mathbf{u}^\varepsilon \varphi + \mathcal{B}_p^\varepsilon \cdot \nabla (\mathbf{u}^\varepsilon \varphi) + \lambda (\mathbf{u}^\varepsilon - \mathbf{g}) \cdot \mathbf{u}^\varepsilon \varphi \right\} dxdt = \int_0^T \int_\Omega \mu_{p, \lambda}^\varepsilon \varphi dxdt$$

for any $\varphi \in L^\infty(\Omega_T) \cap L^\infty((0, T); W^{1, p^*}(\Omega))$, where $\mu_{p, \lambda}^\varepsilon$ is defined by (4.5).

Finally, for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$, taking $\mathbf{w} = \mathbf{u}^\varepsilon \wedge \mathbf{v}$ in (4.24), $\varphi = \mathbf{u}^\varepsilon \cdot \mathbf{v}$ in (4.26), and using the formula $\mathbf{a} \cdot (\mathbf{b} \wedge \mathbf{c}) = (\mathbf{a} \wedge \mathbf{b}) \cdot \mathbf{c}$ yield

(4.27)

$$\int_0^T \int_\Omega \left\{ \mathbf{u}_t^\varepsilon \cdot (\mathbf{u}^\varepsilon \wedge (\mathbf{u}^\varepsilon \wedge \mathbf{v})) + \mathcal{B}_p^\varepsilon \cdot \nabla (\mathbf{u}^\varepsilon \wedge (\mathbf{u}^\varepsilon \wedge \mathbf{v})) - \lambda \mathbf{g} \cdot (\mathbf{u}^\varepsilon \wedge (\mathbf{u}^\varepsilon \wedge \mathbf{v})) \right\} dxdt = 0,$$

(4.28)

$$\int_0^T \int_\Omega \left\{ \mathbf{u}_t^\varepsilon \cdot \mathbf{u}^\varepsilon (\mathbf{u}^\varepsilon \cdot \mathbf{v}) + \mathcal{B}_p^\varepsilon \cdot \nabla (\mathbf{u}^\varepsilon (\mathbf{u}^\varepsilon \cdot \mathbf{v})) + \lambda (\mathbf{u}^\varepsilon - \mathbf{g}) \cdot \mathbf{u}^\varepsilon (\mathbf{u}^\varepsilon \cdot \mathbf{v}) \right\} dxdt = \int_0^T \int_\Omega \mu_{p, \lambda}^\varepsilon \mathbf{u}^\varepsilon \cdot \mathbf{v} dxdt.$$

Subtracting (4.27) from (4.28) and using the identity

$$\mathbf{v} = (\mathbf{u}^\varepsilon \cdot \mathbf{v}) \mathbf{u}^\varepsilon - \mathbf{u}^\varepsilon \wedge (\mathbf{u}^\varepsilon \wedge \mathbf{v}),$$

we obtain that

$$(4.29) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t^\varepsilon \cdot \mathbf{v} + \mathcal{B}_p^\varepsilon \cdot \nabla \mathbf{v} + \lambda (\mathbf{u}^\varepsilon - \mathbf{g}) \cdot \mathbf{v} \right\} dxdt = \int_0^T \int_\Omega \mu_{p, \lambda}^\varepsilon \mathbf{u}^\varepsilon \cdot \mathbf{v} dxdt$$

for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$. This is equivalent to saying that \mathbf{u}^ε is a weak solution (in the distributional sense) to (4.2)–(4.4).

Step 3: Wrapping up. We conclude the proof by showing the estimates (4.6)–(4.10). First, (4.6) follows immediately from letting $\delta \rightarrow 0$ in (3.9), appealing to

Fatou’s lemma and the lower semicontinuity of L^2 - and L^{p^*} -norms with respect to L^2 - and L^{p^*} -weak convergence. We emphasize again that this is possible in the case $p = 1$, because the uniform (in δ) estimate (3.6) implies that $u^\varepsilon \in L^\infty((0, T); H^1(\Omega, \mathbf{R}^n))$. Inequalities (4.7) and (4.10) are direct consequences of (4.6). Finally, the bounds (4.8) and (4.9) follow immediately from (4.29) and (4.24), respectively. Hence the proof is complete. \square

Remark 4.1. If $b_p(\varepsilon) = 0$ is used in the regularization, the solutions to (4.1)–(4.4) are only expected to belong to $L^\infty((0, T); [BV(\Omega)]^n)$ in general when $p = 1$.

5. Passing to the limit as $\varepsilon \rightarrow 0$: The case $1 < p < \infty$. In this section, we shall pass to the limit as $\varepsilon \rightarrow 0$ in (4.1)–(4.4) and show that the limit map is a weak solution to (1.8)–(1.11). Since the analysis and techniques for passing the limit for $1 < p < \infty$ and $p = 1$ are quite different, we shall first consider the case $1 < p < \infty$ in this section and leave the case $p = 1$ to the next section. We begin with a definition of weak solutions to (1.8)–(1.11) in the case $1 < p < \infty$.

DEFINITION 5.1. For $1 < p < \infty$, a map $\mathbf{u} : \Omega_T \rightarrow \mathbf{R}^n$ is called a global weak solution to (1.8)–(1.11) if

- (i) $\mathbf{u} \in L^\infty((0, T); W^{1,p}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$;
- (ii) $|\mathbf{u}| = 1$ a.e. on Ω_T ;
- (iii) \mathbf{u} satisfies (1.8)–(1.11) in the distributional sense.

Our main result of this section is the following existence theorem.

THEOREM 5.2. Let $1 < p < \infty$, and suppose that the assumptions on \mathbf{u}_0 and \mathbf{g} in Theorem 4.1 still hold. Then problem (1.8)–(1.11) has a weak solution \mathbf{u} in the sense of Definition 5.1. Moreover, \mathbf{u} satisfies the energy inequality

$$(5.1) \quad J_{p,\lambda}(\mathbf{u}(s)) + \int_0^s \|\mathbf{u}_t(t)\|_{L^2}^2 dt \leq J_{p,\lambda}(\mathbf{u}_0) \quad \text{for a.e. } s \in [0, T],$$

where $J_{p,\lambda}$ is defined by (1.7).

Proof. We divide the proof into three steps.

Step 1: Extracting a convergent subsequence. Let \mathbf{u}^ε denote the solution of (4.1)–(4.4) constructed in Theorem 4.1. From (4.1), (4.6)–(4.10), the weak compactness of $W^{1,p}(\Omega)$, and Sobolev embedding (cf. [1, 45]), there exists a subsequence of $\{\mathbf{u}^\varepsilon\}_{\varepsilon>0}$ (still denoted by the same notation) and a map $\mathbf{u} \in L^\infty((0, T); W^{1,p}(\Omega, \mathbf{R}^n)) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$ such that as $\varepsilon \rightarrow 0$

$$(5.2) \quad \mathbf{u}^\varepsilon \rightharpoonup \mathbf{u} \quad \text{weakly* in } L^\infty((0, T); W^{1,p}(\Omega, \mathbf{R}^n)),$$

$$(5.3) \quad \mathbf{u}^\varepsilon \rightarrow \mathbf{u} \quad \text{strongly in } L^2((0, T); L^p(\Omega, \mathbf{R}^n)),$$

$$(5.4) \quad |\mathbf{u}^\varepsilon| \rightarrow |\mathbf{u}| \quad \text{a.e. in } \Omega_T,$$

$$(5.5) \quad b_p(\varepsilon)\nabla \mathbf{u}^\varepsilon \rightharpoonup 0 \quad \text{weakly in } L^2((0, T); L^2(\Omega, \mathbf{R}^{n \times m})),$$

$$(5.6) \quad \mathbf{u}_t^\varepsilon \rightharpoonup \mathbf{u}_t \quad \text{weakly in } L^2((0, T); L^2(\Omega, \mathbf{R}^n)).$$

It follows immediately from (5.4) and (4.18) that

$$(5.7) \quad |\mathbf{u}| = 1 \quad \text{a.e. in } \Omega_T.$$

Step 2: Passing to the limit. First, by (5.4) and the Lebesgue-dominated convergence theorem we have that

$$(5.8) \quad \mathbf{u}^\varepsilon \xrightarrow{\varepsilon \searrow 0} \mathbf{u} \quad \text{strongly in } L^r((0, T); L^r(\Omega, \mathbf{R}^n)) \quad \forall r \in [1, \infty).$$

Next, let $\mathbf{f}^\varepsilon := \mu_{p,\lambda}^\varepsilon \mathbf{u}^\varepsilon - \lambda(\mathbf{u}^\varepsilon - \mathbf{g})$. Clearly, $\mathbf{f}^\varepsilon \in L^1((0, T); L^1(\Omega; \mathbf{R}^n))$ and is uniformly bounded, on noting (4.5)–(4.6). By Lemma 2.3 we get

$$(5.9) \quad \nabla \mathbf{u}^\varepsilon \xrightarrow{\varepsilon \searrow 0} \nabla \mathbf{u} \quad \text{strongly in } L^q((0, T); L^q(\Omega, \mathbf{R}^{n \times m})) \quad \forall q \in [1, p),$$

which, using (4.7) and (5.5), implies that

$$(5.10) \quad |\nabla \mathbf{u}^\varepsilon|^{p-2} \nabla \mathbf{u}^\varepsilon \xrightarrow{\varepsilon \searrow 0} |\nabla \mathbf{u}|^{p-2} \nabla \mathbf{u} \quad \text{weakly* in } L^\infty((0, T); L^{p'}(\Omega, \mathbf{R}^{n \times m})),$$

$$(5.11) \quad \mathcal{B}_p^\varepsilon \xrightarrow{\varepsilon \searrow 0} \mathcal{B}_p := |\nabla \mathbf{u}|^{p-2} \nabla \mathbf{u} \quad \text{weakly in } L^2((0, T); L^{p'_*}(\Omega, \mathbf{R}^{n \times m})),$$

where $p' = \frac{p}{p-1}$ and $p'_* := \min\{2, p'\}$.

It then follows from taking $\varepsilon \rightarrow 0$ in (4.24) and using (5.6), (5.8), and (5.11) that

$$(5.12) \quad \int_0^T \int_\Omega \left\{ (\mathbf{u}_t \wedge \mathbf{u}) \cdot \mathbf{w} + (\mathcal{B}_p \wedge \mathbf{u}) \cdot \nabla \mathbf{w} - \lambda(\mathbf{g} \wedge \mathbf{u}) \cdot \mathbf{w} \right\} dxdt = 0$$

for any $\mathbf{w} \in C^1(\overline{\Omega}_T)$. Since $\mathcal{B}_p \in L^\infty((0, T); L^{p'}(\Omega, \mathbf{R}^{n \times m}))$, by the standard density argument one can show that (5.12) also holds for all $\mathbf{w} \in L^\infty((0, T); W^{1,p}(\Omega, \mathbf{R}^n)) \cap L^\infty(\Omega_T)$.

Since $|\mathbf{u}| = 1$ a.e. in Ω_T , there holds the following identity, which is analogous to (4.26):

$$(5.13) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot \mathbf{u} \varphi + \mathcal{B}_p \cdot \nabla(\mathbf{u} \varphi) + \lambda(\mathbf{u} - \mathbf{g}) \cdot \mathbf{u} \varphi \right\} dxdt = \int_0^T \int_\Omega \mu_{p,\lambda} \varphi dxdt$$

for any $\varphi \in L^\infty(\Omega_T) \cap L^\infty((0, T); W^{1,p}(\Omega, \mathbf{R}^n))$, where $\mu_{p,\lambda}$ is defined by (1.12).

Finally, for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$, on choosing $\mathbf{w} = \mathbf{u} \wedge \mathbf{v}$ in (5.12) and $\varphi = \mathbf{u} \cdot \mathbf{v}$ in (5.13), subtracting the resulting equations, and using the identity

$$\mathbf{v} = (\mathbf{u} \cdot \mathbf{v}) \mathbf{u} - \mathbf{u} \wedge (\mathbf{u} \wedge \mathbf{v}),$$

we obtain that

$$(5.14) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot \mathbf{v} + \mathcal{B}_p \cdot \nabla \mathbf{v} + \lambda(\mathbf{u} - \mathbf{g}) \cdot \mathbf{v} \right\} dxdt = \int_0^T \int_\Omega \mu_{p,\lambda} \mathbf{u} \cdot \mathbf{v} dxdt$$

for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$. Hence, \mathbf{u} is a weak solution to (1.8)–(1.11).

Step 3: Wrapping up. We conclude the proof by showing the energy inequality (5.1). First, notice that (4.6) implies that

$$(5.15) \quad \int_\Omega \left\{ \frac{1}{p} |\nabla \mathbf{u}^\varepsilon(s)|^p + \frac{\lambda}{2} |\mathbf{u}^\varepsilon(s) - \mathbf{g}|^2 \right\} dx + \int_0^s \|\mathbf{u}_t^\varepsilon(t)\|_{L^2}^2 dt \leq J_{p,\lambda}^\varepsilon(\mathbf{u}_0) \quad \forall s \in [0, T].$$

Then (5.1) follows from taking $\varepsilon \rightarrow 0$ in (5.15) and using Fatou’s lemma and the lower semicontinuity of the L^p -norm with respect to L^p -weak convergence. The proof is complete. \square

Remark 5.1. (a) We remark that it was proved in [15, 32, 41] that weak solutions to (1.8)–(1.11) are not unique in general.

(b) Although the above proof is carried out for any $\alpha > 0$ in the definition of $b_p(\varepsilon)$ (cf. (1.15)), α should be chosen large enough so that the error due to the perturbation

term $b_p(\varepsilon)\Delta$ is much smaller than the error due to other regularization terms in numerical simulations.

(c) The existence result of Theorem 5.2 is established under the assumption $\mathbf{u}_0 \in W^{1,p^*}(\Omega, \mathbf{R}^n)$ with $p^* = \max\{p, 2\}$. This condition can be weakened to $\mathbf{u}_0 \in W^{1,p}(\Omega, \mathbf{R}^n)$ in the case $1 < p < 2$ by a smoothing technique.

6. Passing to the limit as $\varepsilon \rightarrow 0$: The case $p = 1$. In this section, we consider the case $p = 1$ and establish the existence of global weak solutions for the 1-harmonic map heat flow (1.8)–(1.11) by passing to the limit as $\varepsilon \rightarrow 0$ in (4.2)–(4.4). There are two main difficulties which prevent one from repeating the analysis and techniques of the previous section. First, the compactness result of Lemma 2.3 no longer holds when $p = 1$. Second, since the sequence $\{\mathbf{u}^\varepsilon\}_{\varepsilon>0}$ is uniformly bounded only in $L^\infty((0, T); [W^{1,1}(\Omega) \cap L^\infty(\Omega)]^n)$, and $W^{1,1}(\Omega, \mathbf{R}^n)$ is not a reflexive Banach space, hence the limiting map \mathbf{u} now belongs to $L^\infty((0, T); [BV(\Omega) \cap L^\infty(\Omega)]^n)$; i.e., $\mathbf{u}(t)$ is only a map of bounded variation. As expected, these two difficulties make the passage to the limit as $\varepsilon \rightarrow 0$ more difficult and delicate.

6.1. Technical tools and lemmas. In this subsection, we shall cite some technical tools and lemmas in order to develop a weak solution concept to be given in the next subsection for the 1-harmonic map heat flow. Specifically, we need the pairings $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ between a tensor field \mathcal{A} and a BV -vector field \mathbf{v} as a generalization of the pairing $\mathbf{b} \cdot Dv$ between a vector field \mathbf{b} and a BV -function v developed in [4, 10]. The proofs of the key results in achieving this can be found in [23], where an extension of the theory for the 1-harmonic case studied in this section is developed for a general linear growth functional on the gradient.

We recall from [23] the space of *divergence- L^q tensors*

$$(6.1) \quad \mathcal{Y}(\Omega)_q := \{\mathcal{A} \in L^\infty(\Omega; \mathbf{R}^{n \times m}); \operatorname{div} \mathcal{A} \in L^q(\Omega; \mathbf{R}^n)\} \quad \text{for } 1 \leq q < \infty,$$

and that $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ are defined as follows.

DEFINITION 6.1. For any $\mathcal{A} \in \mathcal{Y}(\Omega)_1$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$, we define $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ to be the functionals on $C_0^\infty(\Omega)$ and $[C_0^\infty(\Omega)]^n$, respectively, by

$$(6.2) \quad \langle \mathcal{A} \cdot D\mathbf{v}, \psi \rangle := - \int_{\Omega} (\mathcal{A}^T \mathbf{v}) \cdot \nabla \psi \, dx - \int_{\Omega} (\operatorname{div} \mathcal{A} \cdot \mathbf{v}) \psi \, dx \quad \forall \psi \in C_0^\infty(\Omega),$$

$$(6.3) \quad \langle \mathcal{A} \wedge D\mathbf{v}, \mathbf{w} \rangle := - \int_{\Omega} (\mathcal{A} \wedge \mathbf{v}) \cdot \nabla \mathbf{w} \, dx - \int_{\Omega} (\operatorname{div} \mathcal{A} \wedge \mathbf{v}) \cdot \mathbf{w} \, dx \quad \forall \mathbf{w} \in [C_0^\infty(\Omega)]^n,$$

where \mathcal{A}^T stands for the matrix transpose of \mathcal{A} and the notation (2.4) is used in (6.3).

We now list some properties of the pairings $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ and refer to section 2 of [23] for their proofs. The first lemma declares that $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ are Radon measures in Ω .

LEMMA 6.2. For any Borel set $E \subset \Omega$, there hold

$$|\langle \mathcal{A} \cdot D\mathbf{v}, \psi \rangle| \leq \max_{x \in E} |\psi(x)| \cdot \|\mathcal{A}\|_{L^\infty(E, \mathbf{R}^{n \times m})} \cdot |D\mathbf{v}|(E) \quad \forall \psi \in C_0(E),$$

$$|\langle \mathcal{A} \wedge D\mathbf{v}, \mathbf{w} \rangle| \leq \max_{x \in E} |\mathbf{w}(x)| \cdot \|\mathcal{A}\|_{L^\infty(E, \mathbf{R}^{n \times m})} \cdot |D\mathbf{v}|(E) \quad \forall \mathbf{w} \in [C_0(E)]^n.$$

Hence, it follows from the Riesz theorem (cf. Theorem 1.54 of [3]) that both functionals $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ are Radon measures in Ω .

COROLLARY 6.3. *The measures $\mathcal{A} \cdot D\mathbf{v}$, $|\mathcal{A} \cdot D\mathbf{v}|$, $\mathcal{A} \wedge D\mathbf{v}$, and $|\mathcal{A} \wedge D\mathbf{v}|$ all are absolutely continuous with respect to the measure $|D\mathbf{v}|$ in Ω . Moreover, there hold inequalities*

$$(6.4) \quad |(\mathcal{A} \cdot D\mathbf{v})(E)| \leq |\mathcal{A} \cdot D\mathbf{v}|(E) \leq \|\mathcal{A}\|_{L^\infty(E', \mathbf{R}^{n \times m})} \cdot |D\mathbf{v}|(E),$$

$$(6.5) \quad |(\mathcal{A} \wedge D\mathbf{v})(E)| \leq |\mathcal{A} \wedge D\mathbf{v}|(E) \leq \|\mathcal{A}\|_{L^\infty(E', \mathbf{R}^{n \times m})} \cdot |D\mathbf{v}|(E)$$

for all Borel sets E and for all open sets E' such that $E \subset E' \subset \Omega$.

Hence, by the Radon–Nikodým theorem (cf. Theorem 1.28 of [3]), there exist $|D\mathbf{v}|$ -measurable functions $\Theta := \Theta(\mathcal{A}, D\mathbf{v}, x) : \Omega \rightarrow \mathbf{R}$ and $\Lambda := \Lambda(\mathcal{A}, D\mathbf{v}, x) : \Omega \rightarrow \mathbf{R}^n$ such that

$$(6.6) \quad (\mathcal{A} \cdot D\mathbf{v})(E) = \int_E \Theta d|D\mathbf{v}| \quad \text{and} \quad \|\Theta\|_{L^\infty(\Omega, |D\mathbf{v}|)} \leq \|\mathcal{A}\|_{L^\infty(\Omega, \mathbf{R}^{n \times m})},$$

$$(6.7) \quad (\mathcal{A} \wedge D\mathbf{v})(E) = \int_E \Lambda d|D\mathbf{v}| \quad \text{and} \quad \|\Lambda\|_{L^\infty(\Omega, |D\mathbf{v}|)} \leq \|\mathcal{A}\|_{L^\infty(\Omega, \mathbf{R}^{n \times m})}$$

for all Borel sets $E \subset \Omega$.

The second lemma declares that every $\mathcal{A} \in \mathcal{Y}(\Omega)_1$ has a well-behaved traction $\mathcal{A}\mathbf{n}$ on the boundary of a Lipschitz domain Ω .

LEMMA 6.4. *Let Ω be a bounded domain with a Lipschitz continuous boundary $\partial\Omega$ in \mathbf{R}^m . Then there exists a linear operator $\beta : \mathcal{Y}(\Omega)_1 \rightarrow L^\infty(\partial\Omega; \mathbf{R}^n)$ such that*

$$(6.8) \quad \|\beta(\mathcal{A})\|_{L^\infty(\partial\Omega, \mathbf{R}^n)} \leq \|\mathcal{A}\|_{L^\infty(\Omega, \mathbf{R}^{n \times m})},$$

$$(6.9) \quad \langle \mathcal{A}, \mathbf{v} \rangle_{\partial\Omega} = \int_{\partial\Omega} \beta(\mathcal{A})(x) \mathbf{v}(x) d\mathcal{H}^{m-1} \quad \forall \mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n,$$

$$(6.10) \quad \beta(\mathcal{A})(x) = \mathcal{A}(x)\mathbf{n}(x) \quad \forall x \in \partial\Omega, \mathcal{A} \in C^1(\overline{\Omega}, \mathbf{R}^{n \times m}).$$

Remark 6.1. Since $\beta(\mathcal{A})$ is a weakly defined traction of \mathcal{A} on $\partial\Omega$, hence we shall use $\mathcal{A}\mathbf{n}$ to denote $\beta(\mathcal{A})$ in the rest of this section.

The third lemma declares that the following hold.

LEMMA 6.5. *Let Ω be a bounded domain with a Lipschitz continuous boundary $\partial\Omega$ in \mathbf{R}^m . Then for any $\mathcal{A} \in \mathcal{Y}(\Omega)_1$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$ there hold the identities*

$$(6.11) \quad \int_{\Omega} \operatorname{div} \mathcal{A} \cdot \mathbf{v} dx + (\mathcal{A} \cdot D\mathbf{v})(\Omega) = \int_{\partial\Omega} \mathcal{A}\mathbf{n} \cdot \mathbf{v} d\mathcal{H}^{m-1},$$

$$(6.12) \quad \int_{\Omega} \operatorname{div} \mathcal{A} \wedge \mathbf{v} dx + (\mathcal{A} \wedge D\mathbf{v})(\Omega) = \int_{\partial\Omega} \mathcal{A}\mathbf{n} \wedge \mathbf{v} d\mathcal{H}^{m-1}.$$

Lemmas 6.6 and 6.7 state continuity results for the measure $\mathcal{A} \cdot D\mathbf{v}$ and $\mathcal{A} \wedge D\mathbf{v}$ with respect to \mathcal{A} and \mathbf{v} , respectively.

LEMMA 6.6. *Let $\mathcal{A}_j, \mathcal{A} \in \mathcal{Y}(\Omega)_1$ and suppose that*

$$\mathcal{A}_j \rightharpoonup \mathcal{A} \quad \text{weakly* in } L^\infty(E),$$

$$\operatorname{div} \mathcal{A}_j \rightharpoonup \operatorname{div} \mathcal{A} \quad \text{weakly in } L^1(E)$$

for all open sets $E \subset\subset \Omega$. Then for all $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$ the following hold:

$$(6.13) \quad \mathcal{A}_j \cdot D\mathbf{v} \rightharpoonup \mathcal{A} \cdot D\mathbf{v} \quad \text{weakly* in } \mathcal{M}(\Omega),$$

$$(6.14) \quad \mathcal{A}_j \wedge D\mathbf{v} \rightharpoonup \mathcal{A} \wedge D\mathbf{v} \quad \text{weakly* in } [\mathcal{M}(\Omega)]^n$$

and

$$(6.15) \quad \Theta(\mathcal{A}_j, D\mathbf{v}, \cdot) \rightharpoonup \Theta(\mathcal{A}, D\mathbf{v}, \cdot) \quad \text{weakly* in } L^\infty(E) \forall E \subset\subset \Omega,$$

$$(6.16) \quad \Lambda(\mathcal{A}_j, D\mathbf{v}, \cdot) \rightharpoonup \Lambda(\mathcal{A}, D\mathbf{v}, \cdot) \quad \text{weakly* in } [L^\infty(E)]^n \forall E \subset\subset \Omega.$$

Here “ $E \subset\subset \Omega$ ” means that E is compactly contained in Ω ; that is, $E \subset \overline{E} \subset \Omega$ and \overline{E} is compact.

LEMMA 6.7. Let $\mathcal{A} \in \mathcal{Y}(\Omega)_1$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$. Suppose that $\{\mathbf{v}_j\} \subset [C^\infty(\Omega) \cap BV(\Omega)]^n$ strictly converges to \mathbf{v} (cf. Definition 3.14 of [3]). Then

$$(6.17) \quad \mathcal{A} \cdot D\mathbf{v}_j \rightharpoonup \mathcal{A} \cdot D\mathbf{v} \quad \text{weakly* in } \mathcal{M}(\Omega),$$

$$(6.18) \quad \mathcal{A} \wedge D\mathbf{v}_j \rightharpoonup \mathcal{A} \wedge D\mathbf{v} \quad \text{weakly* in } [\mathcal{M}(\Omega)]^n.$$

Moreover,

$$(6.19) \quad \int_{\Omega} \mathcal{A} \cdot D\mathbf{v}_j \, dx \longrightarrow \int_{\Omega} \mathcal{A} \cdot D\mathbf{v},$$

$$(6.20) \quad \int_{\Omega} \mathcal{A} \wedge D\mathbf{v}_j \, dx \longrightarrow \int_{\Omega} \mathcal{A} \wedge D\mathbf{v}.$$

Lemma 6.8 gives the precise representations for the density functions Θ and Λ defined in Corollary 6.3.

LEMMA 6.8. (i) If $\mathcal{A} \in \mathcal{Y}(\Omega)_1 \cap C(\Omega, \mathbf{R}^{n \times m})$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$, then there hold

$$(6.21) \quad \Theta(\mathcal{A}, D\mathbf{v}, x) = \mathcal{A}(x) \cdot \frac{D\mathbf{v}}{|D\mathbf{v}|}(x), \quad \Lambda(\mathcal{A}, D\mathbf{v}, x) = \mathcal{A}(x) \wedge \frac{D\mathbf{v}}{|D\mathbf{v}|}(x)$$

for $|D\mathbf{v}|$ a.e. in Ω .

(ii) If $\mathcal{A} \in \mathcal{Y}(\Omega)_1$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$, then there hold

$$(6.22) \quad \Theta(\mathcal{A}, D\mathbf{v}, x) = \mathcal{A}(x) \cdot \frac{D\mathbf{v}}{|D\mathbf{v}|}(x), \quad \Lambda(\mathcal{A}, D\mathbf{v}, x) = \mathcal{A}(x) \wedge \frac{D\mathbf{v}}{|D\mathbf{v}|}(x)$$

for $|D\mathbf{v}|^a$ a.e. in Ω , where $\frac{D\mathbf{v}}{|D\mathbf{v}|}$ denotes the density function of the measure $D\mathbf{v}$ with respect to the measure $|D\mathbf{v}|$, and $|D\mathbf{v}|^a$ denotes the absolute continuous part of the measure $|D\mathbf{v}|$ with respect to the Lebesgue measure \mathcal{L}^n .

Next, we recall from [23] the space of divergence-measure tensors

$$(6.23) \quad \mathcal{DT}(\Omega) := \{ \mathcal{A} \in L^\infty(\Omega, \mathbf{R}^{n \times m}); \operatorname{div} \mathcal{A} \in [\mathcal{M}(\Omega)]^n \}$$

and briefly discuss two of its important properties. First, as in the case of the space of the divergence- L^1 tensors $\mathcal{Y}(\Omega)_1$, Definition 6.1 is still valid for $\mathcal{A} \in \mathcal{DT}(\Omega)$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega) \cap C(\Omega)]^n$ (cf. [4]). Second, there is a well-behaved traction $\mathcal{A}\mathbf{n}$ for every $\mathcal{A} \in \mathcal{DT}(\Omega)$.

LEMMA 6.9. *Let Ω be a bounded domain with Lipschitz continuous boundary $\partial\Omega$ in \mathbf{R}^m . Then there exists a linear operator $\alpha : \mathcal{DT}(\Omega) \rightarrow L^\infty(\partial\Omega; \mathbf{R}^n)$ such that*

$$(6.24) \quad \|\alpha(\mathcal{A})\|_{L^\infty(\partial\Omega, \mathbf{R}^n)} \leq \|\mathcal{A}\|_{L^\infty(\Omega, \mathbf{R}^{n \times m})},$$

$$(6.25) \quad \langle \mathcal{A}, \mathbf{v} \rangle_{\partial\Omega} = \int_{\partial\Omega} \alpha(\mathcal{A})(x) \mathbf{v}(x) d\mathcal{H}^{m-1} \quad \forall \mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega) \cap C(\Omega)]^n,$$

$$(6.26) \quad \alpha(\mathcal{A})(x) = \mathcal{A}(x)\mathbf{n}(x) \quad \forall x \in \partial\Omega, \mathcal{A} \in C^1(\overline{\Omega}, \mathbf{R}^{n \times m}).$$

Moreover, for any $\mathcal{A} \in \mathcal{DT}(\Omega)$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega) \cap C(\Omega)]^n$, let $\mathbf{A}\mathbf{n} := \alpha(\mathcal{A})$ on $\partial\Omega$. Then there hold the following Green's formulas:

$$(6.27) \quad (\operatorname{div} \mathcal{A} \cdot \mathbf{v})(\Omega) + (\mathcal{A} \cdot D\mathbf{v})(\Omega) = \int_{\partial\Omega} \mathbf{A}\mathbf{n} \cdot \mathbf{v} d\mathcal{H}^{m-1},$$

$$(6.28) \quad (\operatorname{div} \mathcal{A} \wedge \mathbf{v})(\Omega) + (\mathcal{A} \wedge D\mathbf{v})(\Omega) = \int_{\partial\Omega} \mathbf{A}\mathbf{n} \wedge \mathbf{v} d\mathcal{H}^{m-1}.$$

Third, the following product rule holds.

LEMMA 6.10. *For any $\mathcal{A} \in \mathcal{DT}(\Omega)$ and $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega)]^n$, the identities*

$$(6.29) \quad \operatorname{div}(\mathcal{A}^T \mathbf{v}) = (\operatorname{div} \mathcal{A}) \cdot \overline{\mathbf{v}} + \overline{\mathcal{A} \cdot D\mathbf{v}},$$

$$(6.30) \quad \operatorname{div}(\mathcal{A} \wedge \mathbf{v}) = (\operatorname{div} \mathcal{A}) \wedge \overline{\mathbf{v}} + \overline{\mathcal{A} \wedge D\mathbf{v}}$$

hold in the sense of Radon measures in Ω , where $\overline{\mathbf{v}}$ denotes the limit of a mollified sequence for \mathbf{v} through a positive symmetric mollifier, and $\overline{\mathcal{A} \cdot D\mathbf{v}}$ (resp., $\overline{\mathcal{A} \wedge D\mathbf{v}}$) is a Radon measure which is absolutely continuous with respect to the measure $|D\mathbf{v}|$, and whose absolutely continuous part $\overline{(\mathcal{A} \cdot D\mathbf{v})}^a$ (resp., $\overline{(\mathcal{A} \wedge D\mathbf{v})}^a$) with respect to the Lebesgue measure \mathcal{L}^m in Ω coincides with $\mathcal{A} \cdot (\nabla \mathbf{v})^a$ (resp., $\mathcal{A} \wedge (\nabla \mathbf{v})^a$) a.e. in Ω , that is, $\overline{(\mathcal{A} \cdot D\mathbf{v})}^a = \mathcal{A} \cdot (\nabla \mathbf{v})^a$ (resp., $\overline{(\mathcal{A} \wedge D\mathbf{v})}^a = \mathcal{A} \wedge (\nabla \mathbf{v})^a$) for \mathcal{L}^m a.e. in Ω .

Remark 6.2. Equations (6.29) and (6.30) hold without all the overbars if either $\mathbf{v} \in [BV(\Omega) \cap L^\infty(\Omega) \cap C(\Omega)]^n$ or $\mathcal{A} \in \mathcal{Y}(\Omega)_1$.

Remark 6.3. It should be noted that the results of Lemmas 6.6–6.8 also hold for the tensor fields in $\mathcal{DT}(\Omega)$ and the vector fields in $[BV(\Omega) \cap L^\infty(\Omega) \cap C(\Omega)]^n$.

6.2. Existence of weak solutions of 1-harmonic map heat flow. Throughout the rest of this paper we let $B_1(\mathbf{R}^{n \times m})$ denote the unit ball in the Euclidean space $\mathbf{R}^{n \times m}$; that is,

$$B_1(\mathbf{R}^{n \times m}) = \left\{ \mathcal{A} \in \mathbf{R}^{n \times m}; |\mathcal{A}| := \left(\sum_{j=1}^m \sum_{i=1}^n \mathcal{A}_{ij}^2 \right)^{\frac{1}{2}} \leq 1 \right\}.$$

In addition, let

$$\sigma_\lambda^\varepsilon := \mu_{1,\lambda}^\varepsilon \mathbf{u}^\varepsilon, \quad \mathcal{B}^\varepsilon := \frac{\nabla \mathbf{u}^\varepsilon}{|\nabla \mathbf{u}^\varepsilon|_\varepsilon}.$$

Hence $\mathcal{B}_1^\varepsilon = b_1(\varepsilon) \nabla \mathbf{u}^\varepsilon + \mathcal{B}^\varepsilon$ (cf. (4.11)).

We now give a definition of weak solutions to (1.8)–(1.11) in the case $p = 1$.

DEFINITION 6.11. *For $p = 1$, a map $\mathbf{u} : \Omega_T \rightarrow \mathbf{R}^n$ is called a global weak solution to (1.8)–(1.11) if there exists a tensor (or matrix-valued function) \mathcal{B} such that*

- (i) $\mathbf{u} \in L^\infty((0, T); [BV(\Omega) \cap L^\infty(\Omega)]^n) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$;
- (ii) $|\mathbf{u}| = 1$ for \mathcal{L}^{m+1} a.e. in Ω_T ;
- (iii) $\mathcal{B} \in L^\infty((0, T); L^\infty(\Omega, B_1(\mathbf{R}^{n \times m}))) \cap L^2((0, T); \mathcal{DT}(\Omega))$;
- (iv) \mathbf{u} and \mathcal{B} satisfy $\mathcal{B} \wedge \mathbf{u} \in L^2((0, T); \mathcal{Y}(\Omega)_2)$, $\mathcal{B}^T \mathbf{u} = 0$, and $\mathcal{B} \cdot (D\mathbf{u})^a = |D\mathbf{u}|^a$ for \mathcal{L}^{m+1} a.e. in Ω_T ;
- (v) $\mathcal{B}\mathbf{n} = 0$ on $\partial\Omega_T$ in the sense of Lemma 6.9;
- (vi) there holds the identity

$$\int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot \mathbf{v} + \mathcal{B} \cdot \nabla \mathbf{v} + \lambda(\mathbf{u} - \mathbf{g}) \cdot \mathbf{v} \right\} dx dt = \int_0^T \left(\int_\Omega \mathbf{v} d\sigma_\lambda \right) dt$$

for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$, where σ_λ denotes the vector-valued Radon measure

$$(6.31) \quad \sigma_\lambda = (\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u} + \lambda(1 - \mathbf{g} \cdot \mathbf{u}) \mathbf{u}.$$

Moreover, the Radon measure $(\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u}$ is absolutely continuous with respect to the measure $|D\mathbf{u}|$, and for \mathcal{L}^1 a.e. $t \in (0, T)$ there exists a function $\Phi(\mathcal{B}, D\mathbf{u}, x, t) : \Omega \rightarrow \mathbf{R}$ such that

$$((\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u})(E) = \int_E \Phi(\mathcal{B} \wedge \mathbf{u}, D\mathbf{u}, x, t) d|D\mathbf{u}| \quad \text{for all Borel sets } E \subset \Omega,$$

$$\|\Phi(t)\|_{L^\infty(\Omega, |D\mathbf{u}|)} \leq \|\mathcal{B}(t)\|_{L^\infty(\Omega, \mathbf{R}^{n \times m})} \quad \text{for } \mathcal{L}^1 \text{ a.e. } t \in (0, T),$$

$$\Phi(\mathcal{B} \wedge \mathbf{u}, D\mathbf{u}, x, t) = (\mathcal{B}(t) \wedge \mathbf{u}) \wedge \frac{D\mathbf{u}}{|D\mathbf{u}|} \quad \text{for } |D\mathbf{u}|^a \text{ a.e. in } \Omega, \mathcal{L}^1 \text{ a.e. } t \in (0, T),$$

$$((\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u})^a = (\mathcal{B} \wedge \mathbf{u}) \wedge (D\mathbf{u})^a = |D\mathbf{u}|^a \mathbf{u} \quad \text{for } \mathcal{L}^{m+1} \text{ a.e. in } \Omega_T.$$

Remark 6.4. (a) The tensor field \mathcal{B} extends $\frac{D\mathbf{u}}{|D\mathbf{u}|}$ as a calibration across the gulfs.

(b) If $\mathbf{u}(t) \in W_{\text{loc}}^{1,1}(\Omega)$ and $\mathcal{B}(t) \in \mathcal{Y}(\Omega)_1$, using (ii) and (iv) and the identity $(\mathbf{a} \wedge \mathbf{b}) \wedge \mathbf{c} = (\mathbf{a} \cdot \mathbf{c})\mathbf{b} - (\mathbf{b} \cdot \mathbf{c})\mathbf{a}$ we get that $(\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u} = |\nabla \mathbf{u}| \mathbf{u}$. Hence, (6.31) can be rewritten as

$$(6.32) \quad \sigma_\lambda = |\nabla \mathbf{u}| \mathbf{u} + \lambda(1 - \mathbf{g} \cdot \mathbf{u}) \mathbf{u}.$$

Thus, (6.31) is a weak form of (6.32) since $|D\mathbf{u}| \mathbf{u}$ may not be defined for $\mathbf{u}(t) \in [BV(\Omega) \cap L^\infty(\Omega)]^n$ and $\mathcal{B}(t) \in \mathcal{DT}(\Omega)$.

Our main result of this section is the following existence theorem.

THEOREM 6.12. *Let $p = 1$, and suppose that $\mathbf{u}_0 \in H^1(\Omega, \mathbf{R}^n)$ and $\mathbf{g} \in L^2(\Omega, \mathbf{R}^n)$ with $|\mathbf{u}_0| = 1$ and $|\mathbf{g}| \leq 1$ a.e. in Ω . Then problem (1.8)–(1.11) has a global weak solution \mathbf{u} in the sense of Definition 6.11. Moreover, \mathbf{u} satisfies the energy inequality*

$$(6.33) \quad I_\lambda(\mathbf{u}(s)) + \int_0^s \|\mathbf{u}_t(t)\|_{L^2}^2 dt \leq I_\lambda(\mathbf{u}_0) \quad \text{for a.e. } s \in [0, T],$$

where

$$(6.34) \quad I_\lambda(\mathbf{u}) := |D\mathbf{u}|(\Omega) + \frac{\lambda}{2} \int_\Omega |\mathbf{u} - \mathbf{g}|^2 dx.$$

Proof. The proof is divided into four steps.

Step 1: Extracting a convergent subsequence and passing to the limit. Let \mathbf{u}^ε denote the solution of (4.1)–(4.4) constructed in Theorem 4.1. It follows from (4.1), (4.6), (4.7), and (4.8) that $\{\mathbf{u}^\varepsilon\}_{\varepsilon>0}$ is uniformly bounded in $L^\infty((0, T); [W^{1,1}(\Omega) \cap L^\infty(\Omega)]^n \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n)))$, $\{\mathcal{B}^\varepsilon\}_{\varepsilon>0}$ in $L^\infty((0, T); L^\infty(\Omega, B_1(\mathbf{R}^{n \times m})))$, and $\{\operatorname{div} \mathcal{B}_1^\varepsilon\}_{\varepsilon>0}$ in $L^2((0, T); L^1(\Omega; \mathbf{R}^n))$, and $\{\sigma_\lambda^\varepsilon\}_{\varepsilon>0}$ is uniformly bounded in $L^\infty((0, T); L^1(\Omega, \mathbf{R}^n))$. Since $L^1(\Omega) \subset \mathcal{M}(\Omega)$ and $W^{1,1}(\Omega) \subset BV(\Omega)$, by the weak compactness of $\mathcal{M}(\Omega)$ and $BV(\Omega)$ (cf. [3]) we have that there exist subsequences of $\{\mathbf{u}^\varepsilon\}_{\varepsilon>0}$, $\{\mathcal{B}^\varepsilon\}_{\varepsilon>0}$, $\{\mathcal{B}_1^\varepsilon\}_{\varepsilon>0}$, and $\{\sigma_\lambda^\varepsilon\}_{\varepsilon>0}$ (still denoted by the same notation), respectively, and maps $\mathbf{u} \in L^\infty((0, T); [BV(\Omega) \cap L^\infty(\Omega)]^n) \cap H^1((0, T); L^2(\Omega, \mathbf{R}^n))$, $\mathcal{B} \in L^\infty((0, T); L^\infty(\Omega, B_1(\mathbf{R}^{n \times m})))$, $\nu \in L^2((0, T); [\mathcal{M}(\Omega)]^n)$, and $\sigma_\lambda \in L^\infty((0, T); [\mathcal{M}(\Omega)]^n)$, respectively, such that as $\varepsilon \rightarrow 0$

$$(6.35) \quad \mathbf{u}^\varepsilon \rightharpoonup \mathbf{u} \quad \text{weakly* in } L^\infty((0, T); [BV(\Omega) \cap L^\infty(\Omega)]^n),$$

$$(6.36) \quad \text{strongly in } L^2((0, T); L^1(\Omega, \mathbf{R}^n))$$

$$(6.37) \quad \text{a.e. in } \Omega_T,$$

$$(6.38) \quad b_1(\varepsilon) \nabla \mathbf{u}^\varepsilon \rightharpoonup 0 \quad \text{weakly in } L^2((0, T); L^2(\Omega, \mathbf{R}^{n \times m})),$$

$$(6.39) \quad \mathbf{u}_t^\varepsilon \rightharpoonup \mathbf{u}_t \quad \text{weakly in } L^2((0, T); L^2(\Omega, \mathbf{R}^n)),$$

$$(6.40) \quad \operatorname{div} \mathcal{B}_1^\varepsilon \rightharpoonup \nu \quad \text{weakly* in } L^2((0, T); [\mathcal{M}(\Omega)]^n),$$

$$(6.41) \quad \mathcal{B}^\varepsilon \rightharpoonup \mathcal{B} \quad \text{weakly* in } L^\infty((0, T); L^\infty(\Omega, B_1(\mathbf{R}^{n \times m}))),$$

$$(6.42) \quad \sigma_\lambda^\varepsilon \rightharpoonup \sigma_\lambda \quad \text{weakly* in } L^\infty((0, T); [\mathcal{M}(\Omega)]^n).$$

It follows immediately from (6.37) and (4.18) that

$$(6.43) \quad |\mathbf{u}| = 1 \quad \text{a.e. in } \Omega_T,$$

and an application of the Lebesgue-dominated convergence theorem yields that

$$(6.44) \quad \mathbf{u}^\varepsilon \xrightarrow{\varepsilon \searrow 0} \mathbf{u} \quad \text{strongly in } L^r((0, T); L^r(\Omega, \mathbf{R}^n)) \quad \forall r \in [1, \infty).$$

Now, taking $\varepsilon \rightarrow 0$ in (4.24), and (4.29) (with $p = 1$) we have for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$ and $\mathbf{w} \in L^\infty((0, T); H^1(\Omega; \mathbf{R}^n))$

$$(6.45) \quad \int_0^T \int_\Omega \left\{ (\mathbf{u}_t \wedge \mathbf{u}) \cdot \mathbf{w} + (\mathcal{B} \wedge \mathbf{u}) \cdot \nabla \mathbf{w} - \lambda (\mathbf{g} \wedge \mathbf{u}) \cdot \mathbf{w} \right\} dxdt = 0,$$

$$(6.46) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot \mathbf{v} + \mathcal{B} \cdot \nabla \mathbf{v} + \lambda (\mathbf{u} - \mathbf{g}) \cdot \mathbf{v} \right\} dxdt = \int_0^T \left(\int_\Omega \mathbf{v} d\sigma_\lambda \right) dt.$$

Step 2: Identifying ν and σ_λ . First, it follows from the identity $(\mathcal{B}^\varepsilon)^T \mathbf{u}^\varepsilon = 0$ (cf. (4.25)), (6.41), and (6.44) that

$$(6.47) \quad \mathcal{B}^T \mathbf{u} = 0 \quad \text{a.e. in } \Omega_T.$$

Second, for any $\mathbf{v} \in [C^1(\overline{\Omega}_T)]^n$, it follows from (6.38), (6.40), and (6.41) that

$$\int_{\Omega_T} \mathbf{v} \cdot d\nu = \lim_{\varepsilon \rightarrow 0} \int_{\Omega_T} \mathbf{v} \cdot \operatorname{div} \mathcal{B}_1^\varepsilon dxdt = - \lim_{\varepsilon \rightarrow 0} \int_{\Omega_T} \nabla \mathbf{v} \cdot \mathcal{B}_1^\varepsilon dxdt = - \int_{\Omega_T} \nabla \mathbf{v} \cdot \mathcal{B} dxdt.$$

Hence, $\operatorname{div} \mathcal{B}$ exists and

$$(6.48) \quad \operatorname{div} \mathcal{B} = \nu, \quad \text{and therefore} \quad \mathcal{B} \in L^2((0, T); \mathcal{DT}(\Omega)).$$

It then follows from (6.11) that

$$\mathcal{B}\mathbf{n} = 0 \quad \text{on } \partial\Omega_T.$$

Third, notice that (6.45) immediately implies that

$$(6.49) \quad \operatorname{div}(\mathcal{B} \wedge \mathbf{u}) \in L^2(\Omega_T, \mathbf{R}^n), \quad \text{and hence} \quad \mathcal{B} \wedge \mathbf{u} \in L^2((0, T); \mathcal{Y}(\Omega)_2).$$

Let $\{\mathbf{u}_\rho\}$ denote the smooth approximation sequence of \mathbf{u} as constructed in Theorem 3.9 of [3]. For any $\mathbf{v} \in [C_0^1(\Omega_T)]^n$, setting $\mathbf{w} = \mathbf{u}_\rho \wedge \mathbf{v}$ in (6.45) we get

$$(6.50) \quad \int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot (\mathbf{u} \wedge (\mathbf{u}_\rho \wedge \mathbf{v})) + (\mathcal{B} \wedge \mathbf{u}) \cdot \nabla (\mathbf{u}_\rho \wedge \mathbf{v}) - \lambda \mathbf{g} \cdot (\mathbf{u} \wedge (\mathbf{u}_\rho \wedge \mathbf{v})) \right\} dx dt = 0.$$

It follows from (6.43), the convergence property of \mathbf{u}_ρ (cf. [3]), and the identity $\mathbf{u} \wedge (\mathbf{u} \wedge \mathbf{v}) = \mathbf{u}(\mathbf{u} \cdot \mathbf{v}) - \mathbf{v}$ that

$$\begin{aligned} \lim_{\rho \rightarrow 0} \int_\Omega \mathbf{u}_t \cdot (\mathbf{u} \wedge (\mathbf{u}_\rho \wedge \mathbf{v})) dx &= \int_\Omega \mathbf{u}_t \cdot (\mathbf{u} \wedge (\mathbf{u} \wedge \mathbf{v})) dx = \int_\Omega \mathbf{u}_t \cdot (\mathbf{u}(\mathbf{u} \cdot \mathbf{v}) - \mathbf{v}) dx \\ &= - \int_\Omega \mathbf{u}_t \cdot \mathbf{v} dx, \\ \lim_{\rho \rightarrow 0} \int_\Omega \mathbf{g} \cdot (\mathbf{u} \wedge (\mathbf{u}_\rho \wedge \mathbf{v})) dx &= \int_\Omega \mathbf{g} \cdot (\mathbf{u} \wedge (\mathbf{u} \wedge \mathbf{v})) dx \\ &= \int_\Omega \left\{ (\mathbf{u} - \mathbf{g}) \cdot \mathbf{v} - (1 - \mathbf{g} \cdot \mathbf{u}) \mathbf{u} \cdot \mathbf{v} \right\} dx. \end{aligned}$$

It follows from Theorem 3.9 of [3], Lemma 6.7, and the identity

$$\int_\Omega (\mathcal{B} \wedge \mathbf{u}) \cdot \nabla (\mathbf{u}_\rho \wedge \mathbf{v}) dx = \int_\Omega \left\{ ((\mathcal{B} \wedge \mathbf{u}) \wedge \mathbf{u}_\rho) \cdot \nabla \mathbf{v} + ((\mathcal{B} \wedge \mathbf{u}) \wedge \nabla \mathbf{u}_\rho) \cdot \mathbf{v} \right\} dx$$

that

$$\begin{aligned} \lim_{\rho \rightarrow 0} \int_\Omega (\mathcal{B} \wedge \mathbf{u}) \cdot \nabla (\mathbf{u}_\rho \wedge \mathbf{v}) dx &= \int_\Omega ((\mathcal{B} \wedge \mathbf{u}) \wedge \mathbf{u}) \cdot \nabla \mathbf{v} dx + \langle (\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u}, \mathbf{v} \rangle \\ &= - \int_\Omega \mathcal{B} \cdot \nabla \mathbf{v} dx + \langle (\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u}, \mathbf{v} \rangle. \end{aligned}$$

Here we have used the fact that $(\mathcal{B} \wedge \mathbf{u}) \wedge \mathbf{u} = -\mathcal{B}$ in light of (6.43) and (6.47), and the measure $(\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u}$ is defined by (6.3) with $\mathcal{A} = \mathcal{B} \wedge \mathbf{u}$ and $\mathbf{v} = \mathbf{u}$.

Finally, substituting the above three equations into (6.50) and multiplying the equation by (-1) we get

$$(6.51) \quad \begin{aligned} \int_0^T \int_\Omega \left\{ \mathbf{u}_t \cdot \mathbf{v} + \mathcal{B} \cdot \nabla \mathbf{v} + \lambda (\mathbf{u} - \mathbf{g}) \cdot \mathbf{v} \right\} dx dt \\ = \int_0^T \langle (\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u} + \lambda(1 - \mathbf{g} \cdot \mathbf{u})\mathbf{u}, \mathbf{v} \rangle dt \end{aligned}$$

for any $\mathbf{v} \in [C_0^1(\Omega_T)]^n$. This and (6.46) imply that

$$\sigma_\lambda = (\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u} + \lambda(1 - \mathbf{g} \cdot \mathbf{u})\mathbf{u}.$$

Step 3: Identifying \mathcal{B} . First, since $\{\mathcal{B}^\varepsilon \cdot \nabla \mathbf{u}^\varepsilon\}$ is uniformly bounded in $L^2((0, T); L^1(\Omega))$, then there exists a subsequence (still denoted by the same notation) and $\mu \in L^2((0, T); \mathcal{M}(\Omega))$ such that

$$(6.52) \quad \mathcal{B}^\varepsilon \cdot \nabla \mathbf{u}^\varepsilon \longrightarrow \mu \quad \text{weakly* in } L^2((0, T); \mathcal{M}(\Omega)).$$

Let $\mathbf{u}_\rho^\varepsilon$ and \mathbf{u}_ρ denote mollified sequences for \mathbf{u}^ε and \mathbf{u} , respectively, through a positive symmetric mollifier. For any open set $E \subset \Omega$ we have

$$\begin{aligned} (6.53) \quad \lim_{\varepsilon \rightarrow 0} \int_0^T \int_E \mathcal{B}^\varepsilon \cdot \nabla \mathbf{u}^\varepsilon \, dxdt &= \lim_{\varepsilon \rightarrow 0} \lim_{\rho \rightarrow 0} \int_0^T \int_E \mathcal{B}^\varepsilon \cdot D\mathbf{u}_\rho^\varepsilon \, dxdt \\ &= \lim_{\rho \rightarrow 0} \lim_{\varepsilon \rightarrow 0} \int_0^T \int_E \mathcal{B}^\varepsilon \cdot D\mathbf{u}_\rho^\varepsilon \, dxdt \\ &= \lim_{\rho \rightarrow 0} \int_0^T (\mathcal{B} \cdot D\mathbf{u}_\rho)(E) \, dt \quad (\text{by (6.41), (6.44)}) \\ &= \int_0^T \overline{\mathcal{B} \cdot D\mathbf{u}}(E) \, dt \quad (\text{by (6.29)}), \end{aligned}$$

where $\overline{\mathcal{B} \cdot D\mathbf{u}}$ is defined in Lemma 6.10.

Hence, it follows from (6.52), (6.53), and Lemma 6.10 that

$$(6.54) \quad \mu = \overline{\mathcal{B} \cdot D\mathbf{u}} \lll |D\mathbf{u}|.$$

We refer the reader to Definition 1.24 of [3] for the notation “ \lll .”

On the other hand, a direct calculation yields that

$$\mathcal{B}^\varepsilon \cdot \nabla \mathbf{u}^\varepsilon = \frac{|\nabla \mathbf{u}^\varepsilon|^2}{|\nabla \mathbf{u}^\varepsilon|_\varepsilon} \geq |\nabla \mathbf{u}^\varepsilon|_\varepsilon - \varepsilon.$$

Setting $\varepsilon \rightarrow 0$ and by the lower semicontinuity of the BV -norm (cf. [3]) we obtain

$$\mu \ggg |D\mathbf{u}|,$$

which together with (6.54) and Lemma 6.10 yield that

$$(6.55) \quad |D\mathbf{u}| = \mu = \overline{\mathcal{B} \cdot D\mathbf{u}},$$

and hence

$$(6.56) \quad |D\mathbf{u}|^a = (\overline{\mathcal{B} \cdot D\mathbf{u}})^a = \mathcal{B} \cdot (D\mathbf{u})^a.$$

Finally, we note that all other properties of \mathcal{B} and the measure $(\mathcal{B} \wedge \mathbf{u}) \wedge D\mathbf{u}$ listed in (vi) of Definition 6.11 are immediate consequences of (6.55) and Lemmas 6.2, 6.4–6.10, and Corollary 6.3.

Step 4: Finishing up. We now conclude the proof of the theorem by showing the energy inequality (6.33). First, notice that (4.6) implies that for a.e. $s \in [0, T]$

$$(6.57) \quad I_\lambda(\mathbf{u}^\varepsilon(s)) + \int_0^s \|\mathbf{u}_t^\varepsilon(t)\|_{L^2}^2 dt \leq J_{1,\lambda}^\varepsilon(\mathbf{u}_0) \leq \frac{b_1(\varepsilon)}{2} \|\nabla \mathbf{u}_0\|_{L^2}^2 + a_1(\varepsilon)|\Omega| + I_\lambda(\mathbf{u}_0).$$

Then (6.33) follows from taking $\varepsilon \rightarrow 0$ in (6.57), using the lower semicontinuity of the BV -seminorm and the L^2 -norm with respect to L^2 -weak convergence. The proof is complete. \square

Remark 6.5. (a) Since weak solutions to (1.8)–(1.11) are not unique in general for $1 < p < \infty$ (cf. [15, 32, 41]), we expect that this nonuniqueness also holds for the case $p = 1$.

(b) The existence result in Theorem 6.12 is proved under the assumption $\mathbf{u}_0 \in H^1(\Omega, \mathbf{R}^n)$. This assumption can be weakened to $\mathbf{u}_0 \in [BV(\Omega) \cap L^\infty(\Omega)]^n$ using a smoothing technique.

7. Fully discrete finite element approximations.

7.1. Formulation of fully discrete finite element methods. For ease of exposition, we assume Ω is a polytope in this section. Let \mathcal{T}_h be a quasi-uniform “triangulation” of the domain Ω of mesh size $0 < h < 1$ and $\bar{\Omega} = \bigcup_{K \in \mathcal{T}_h} \bar{K}$ ($K \in \mathcal{T}_h$ are tetrahedrons in the case $m = 3$). Let $J_\tau := \{t_k\}_{k=0}^L$ be a uniform partition of $[0, T]$ with mesh size $\tau := \frac{T}{L}$, and $\partial_t \mathbf{v}^k := (\mathbf{v}^k - \mathbf{v}^{k-1})/\tau$. For an integer $r \geq 1$, let $P_r(K)$ denote the space of polynomials of degree less than or equal to r on K . We introduce the finite element space

$$\mathbf{V}^h = \{ \mathbf{v}_h \in C(\bar{\Omega}, \mathbf{R}^n) \cap H^1(\Omega, \mathbf{R}^n); \mathbf{v}_h|_K \in [P_r(K)]^n \forall K \in \mathcal{T}_h \}.$$

Notice that the density function F defined in (1.17) is not a convex function. On the other hand, there exist two convex functions W_+ and W_- such that

$$(7.1) \quad F(\mathbf{v}) = W_+(\mathbf{v}) - W_-(\mathbf{v}).$$

One such an example is $W_+(\mathbf{v}) = \frac{|\mathbf{v}|^4}{4}$ and $W_-(\mathbf{v}) = \frac{|\mathbf{v}|^2}{2} - \frac{1}{4}$. Clearly, the above decomposition is not unique.

We are now ready to introduce our fully discrete finite element discretizations for the initial boundary value problem (1.20)–(1.22). Find $\mathbf{u}_h^k \in \mathbf{V}^h$ for $k = 1, 2, \dots, L$ such that

$$(7.2) \quad \int_\Omega \left\{ \partial_t \mathbf{u}_h^k \cdot \mathbf{v}_h + \mathcal{B}_h^k \cdot \nabla \mathbf{v}_h + \lambda (\mathbf{u}_h^k - \mathbf{g}) \cdot \mathbf{v}_h + \frac{1}{\delta} W'_+(\mathbf{u}_h^k) \cdot \mathbf{v}_h \right\} dx = \frac{1}{\delta} \int_\Omega W'_-(\mathbf{u}_h^{k-1}) \cdot \mathbf{v}_h dx \quad \forall \mathbf{v} \in \mathbf{V}^h,$$

where

$$(7.3) \quad \mathcal{B}_h^k = [b_p(\varepsilon) + |\nabla \mathbf{u}_h^k|_\varepsilon^{p-2}] \nabla \mathbf{u}_h^k,$$

with some starting value $\mathbf{u}_h^0 \in \mathbf{V}^h$ to be specified later. Note that for notational brevity we have omitted the indices ε, δ , and p on \mathbf{u}_h^k and \mathcal{B}_h^k

For each k , (7.2) is a nonlinear equation in \mathbf{u}_h^k . Hence, the above numerical method is an implicit scheme, and its well-posedness is ensured by the following theorem.

THEOREM 7.1. *For each fixed $k \geq 1$, suppose that $\mathbf{u}_h^{k-1} \in \mathbf{V}^h$ is known. Then there exists a unique solution $\mathbf{u}_h^k \in \mathbf{V}^h$ to (7.2)–(7.3). Moreover, $\{\mathbf{u}_h^k\}_{k=0}^L$ satisfies the following energy estimate:*

$$(7.4) \quad \frac{\tau}{2} \sum_{k=1}^{\ell} \|\partial_t \mathbf{u}_h^k\|_{L^2}^2 + J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_h^k) \leq J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_h^0) \quad \text{for } 1 \leq \ell \leq L.$$

Here $J_{p,\lambda}^{\varepsilon,\delta}$ is defined by (1.19).

Proof. For each fixed $k \geq 1$, it is easy to check that (7.2)–(7.3) is the Euler–Lagrange equation of the following functional over \mathbf{V}^h :

$$(7.5) \quad G_k(\mathbf{v}) := \int_{\Omega} \left\{ \frac{1}{2\tau} |\mathbf{v} - \mathbf{u}_h^{k-1}|^2 + \frac{b_p(\varepsilon)}{2} |\nabla \mathbf{v}|^2 + \frac{1}{p} |\nabla \mathbf{v}|_{\varepsilon}^p + \frac{\lambda}{2} |\mathbf{v} - \mathbf{g}|^2 + \frac{1}{\delta} W_+(\mathbf{v}) \right\} dx - \frac{1}{\delta} \int_{\Omega} W'_-(\mathbf{u}_h^{k-1}) \cdot \mathbf{v} dx.$$

Since G_k is a convex, coercive, and differentiable functional, then it has a unique minimizer $\mathbf{u}_h^k \in \mathbf{V}^h$ (cf. [47]), and hence (7.2)–(7.3) has a unique solution.

Since \mathbf{u}_h^k is the minimizer of G_k over \mathbf{V}^h , we have that

$$(7.6) \quad G_k(\mathbf{u}_h^k) \leq G_k(\mathbf{u}_h^{k-1}).$$

It follows from the convexity of W_- that

$$W'_-(\mathbf{u}_h^{k-1})(\mathbf{u}_h^k - \mathbf{u}_h^{k-1}) \leq W_-(\mathbf{u}_h^k) - W_-(\mathbf{u}_h^{k-1}).$$

This and (7.6) imply that

$$\frac{1}{2} \|\partial_t \mathbf{u}_h^k\|_{L^2}^2 + \frac{J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_h^k) - J_{p,\lambda}^{\varepsilon,\delta}(\mathbf{u}_h^{k-1})}{\tau} \leq 0.$$

The bound (7.4) then follows from applying the summation operator $\tau \sum_{k=1}^{\ell}$ ($1 \leq \ell \leq L$) to the last inequality. Hence the proof is complete. \square

7.2. Convergence analysis. The goal of this subsection is to show that the numerical solution of (7.2)–(7.3) converges to the unique weak solution of (1.20)–(1.22) as $h, \tau \rightarrow 0$. There are two approaches to reach this goal. The first approach assumes the existence of the solution of (1.20)–(1.22), which in fact has been proved in Theorem 3.7, and then proves that \mathbf{u}_h^k converges to that solution. The other approach shows the convergence *without* assuming the existence of the solution of (1.20)–(1.22). This can be done by applying the energy method and compactness argument used in the proof of Theorem 3.7 to the finite element solution $\{\mathbf{u}_h^k\}$. In the following, we shall go with the latter approach since this will also provide an alternative proof for Theorem 3.7 as alluded to in Remark 3.1(c).

For the fully discrete finite element solution $\{\mathbf{u}_h^k\}$, we define its linear interpolation in t as follows:

$$(7.7) \quad \mathbf{U}^{\varepsilon,\delta,h,\tau}(\cdot, t) := \frac{t - t_{k-1}}{\tau} \mathbf{u}_h^k(\cdot) + \frac{t_k - t}{\tau} \mathbf{u}_h^{k-1}(\cdot) \quad \forall t \in [t_{k-1}, t_k], \quad 1 \leq k \leq L.$$

Clearly, $\mathbf{U}^{\varepsilon,\delta,h,\tau}$ is continuous in both x and t .

The main result of this section is the following convergence theorem.

THEOREM 7.2. *For $1 \leq p < \infty$, suppose that $\mathbf{u}_0 \in W^{1,p^*}(\Omega, \mathbf{R}^n)$ with $p^* = \max\{2, p\}$, $|\mathbf{u}_0| = 1$, and $|\mathbf{g}| \leq 1$ in Ω . For each pair of positive numbers (ε, δ) , let $\mathbf{U}^{\varepsilon,\delta,h,\tau}$ be defined by (7.7). Then there exists $\mathbf{u}^{\varepsilon,\delta} \in L^\infty(\Omega_T)$ such that*

$$(7.8) \quad \lim_{h,\tau \rightarrow 0} \|\mathbf{u}^{\varepsilon,\delta} - \mathbf{U}^{\varepsilon,\delta,h,\tau}\|_{L^q(\Omega_T)} = 0 \quad \forall q \in [1, \infty),$$

provided that

$$\lim_{h \rightarrow 0} \|\mathbf{u}_0 - \mathbf{u}_h^0\|_{W^{1,p^*}(\Omega)} = 0.$$

Moreover, $\mathbf{u}^{\varepsilon,\delta}$ solves (1.20)–(1.22) in the sense of Definition 3.1.

Proof. The proof follows along the same lines as that of Theorem 1.5 of [12], where the convergence of a general Galerkin approximation was proved for the case $p \geq 2$. Since the finite element approximation is a special Galerkin approximation, the proof of Theorem 1.5 of [12] can easily be adapted to the finite element approximation $\mathbf{U}^{\varepsilon,\delta,h,\tau}$ for $p \geq 2$ thanks to the discrete energy estimate (7.4) and the facts that $\partial_t \mathbf{u}_h^k = \mathbf{U}_t^{\varepsilon,\delta,h,\tau}$ and $\tau \sum_{k=1}^L \|\partial_t \mathbf{u}_h^k\|_{L^2}^2 = \|\mathbf{U}_t^{\varepsilon,\delta,h,\tau}\|_{L^2(L^2)}$.

Since the operator $-\Delta_p^\varepsilon$ is uniformly elliptic, as a result, the compactness of Lemma 2.2 holds not only for $p \geq 2$ but also for $1 \leq p < 2$. In addition, note that $\mathbf{U}^{\varepsilon,\delta,h,\tau}$ is uniformly (in h and τ) bounded in $L^\infty((0, T); H^1(\Omega, \mathbf{R}^n))$ for $1 \leq p < 2$. Hence, the proof of Theorem 1.5 of [12] can be adapted with slight modifications to prove (7.8) for the case $1 \leq p < 2$. \square

Remark 7.1. Several practical choices of \mathbf{u}_h^0 are possible. For instance, both the L^2 -projection of \mathbf{u}_0 and the Clément finite element interpolation of \mathbf{u}_0 into \mathbf{V}^h (cf. [14]) are qualified candidates for \mathbf{u}_h^0 .

An immediate consequence of Theorems 4.1, 5.2, 6.12, 7.1, and 7.2 is the following convergence theorem.

THEOREM 7.3. *Let $1 \leq p < \infty$, let $\mathbf{U}^{\varepsilon,\delta,h,\tau}$ be defined by (7.7), and assume that the assumptions of Theorems 7.1, 7.2, 4.1, 5.2, and 6.12 hold. Then there exists a subsequence of $\{\mathbf{U}^{\varepsilon,\delta,h,\tau}\}$ (still denoted by the same notation) and a weak solution \mathbf{u} of (1.8)–(1.11) such that*

$$(7.9) \quad \lim_{\varepsilon,\delta \rightarrow 0} \lim_{h,\tau \rightarrow 0} \|\mathbf{u} - \mathbf{U}^{\varepsilon,\delta,h,\tau}\|_{L^q(\Omega_T)} = 0 \quad \forall q \in [1, \infty).$$

Acknowledgments. The authors would like to thank the Mathematisches Forschungsinstitut Oberwolfach for the kind hospitality and opportunity of its “Research in Pairs” program. The second author would like to thank Professor Fanghua Lin for a helpful discussion and the Institute of Mathematics and Its Applications (IMA) of the University of Minnesota for its support and hospitality during the author’s recent visit to the IMA. The authors would also like to thank two anonymous referees for carefully reading the paper and for their valuable suggestions.

REFERENCES

[1] R. A. ADAMS, *Sobolev Spaces*, Academic Press, New York, 1975.
 [2] F. ALOUGES, *A new algorithm for computing liquid crystal stable configurations: The harmonic mapping case*, SIAM J. Numer. Anal., 34 (1997), pp. 1708–1726.
 [3] L. AMBROSIO, N. FUSCO, AND D. PALLARA, *Functions of Bounded Variation and Free Discontinuity Problems*, Oxford University Press, New York, 2000.

- [4] G. ANZELLOTTI, *Pairing between measures and bounded functions and compensated compactness*, Ann. Mat. Pura Appl. (4), 135 (1983), pp. 293–318.
- [5] J. W. BARRETT, X. FENG, AND A. PROHL, *Convergence of a fully discrete finite element method for a degenerate parabolic system modeling nematic liquid crystals with variable degree of orientation*, M2AN Math. Model. Numer. Anal., 40 (2006), pp. 175–199.
- [6] R. BECKER, X. FENG, AND A. PROHL, *Finite element approximations of the Ericksen–Leslie model for nematic liquid crystal flow*, SIAM J. Numer. Anal., 46 (2008), pp. 1704–1731.
- [7] F. BETHUEL, H. BREZIS, AND F. HÉLEIN, *Ginzburg-Landau Vertices*, Birkhäuser, New York, 1994.
- [8] H. BREZIS, *Operateurs maximaux monotones et semi-groupes de contractions dans les espaces de Hilbert*, North-Holland, Amsterdam, The Netherlands, 1973.
- [9] K. C. CHANG, W. D. DING, AND R. YE, *Finite-time blow-up of the heat flow of harmonic maps from surfaces*, J. Differential Geom., 36 (1992), pp. 507–515.
- [10] G.-Q. CHEN AND H. FRID, *Divergence-measure fields and hyperbolic conservation laws*, Arch. Rational Mech. Anal., 147 (1999), pp. 89–118.
- [11] Y. CHEN, *The weak solutions to the evolution problem of harmonic maps*, Math. Z., 201 (1989), pp. 69–74.
- [12] Y. CHEN, M.-H. HONG, AND N. HUNGERBÜHLER, *Heat flow of p -harmonic maps with values into spheres*, Math. Z., 215 (1994), pp. 25–35.
- [13] Y. M. CHEN AND M. STRUWE, *Regularity for the heat flow for harmonic maps*, Math. Z., 201 (1989), pp. 83–103.
- [14] P. G. CIARLET, *The Finite Element Method for Elliptic Problems*, North-Holland, Amsterdam, The Netherlands, 1978.
- [15] J.-M. CORON, *Nonuniqueness for the heat flow of harmonic maps*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 7 (1990), pp. 335–344.
- [16] P. COURILLEAU AND F. DEMENGEL, *Heat flow for p -harmonic maps with values in the circle*, Nonlinear Anal., 41 (2000), pp. 689–700.
- [17] F. DUZAAR AND M. FUCHS, *Existence and regularity of functions which minimize certain energies in homotopy classes of mappings*, Asymptot. Anal., 5 (1991), pp. 129–144.
- [18] J. EELLS AND J. H. SAMPSON, *Harmonic mappings of Riemannian manifolds*, Amer. J. Math., 86 (1964), pp. 109–169.
- [19] J. ERICKSEN AND D. KINDERLEHRER, *Theory and applications of liquid crystals*, IMA Vol. Math. Appl. 5, Springer-Verlag, New York, 1997, pp. 99–122.
- [20] L. C. EVANS, *Weak Convergence Methods for Nonlinear Partial Differential Equations*, AMS, Providence, RI, 1990.
- [21] A. FARDOUN AND R. REGBAOUI, *Heat flow for p -harmonic maps between compact Riemannian manifolds*, Indiana Math. J., 40 (2002), pp. 1305–1320.
- [22] A. FARDOUN AND R. REGBAOUI, *Heat flow for p -harmonic maps with small initial data*, Calc. Var. Partial Differential Equations, 16 (2003), pp. 1–16.
- [23] X. FENG, *Divergence- L^q and divergence-measure tensors fields and gradient flow for linear growth functionals of maps into spheres*, Calc. Var. Partial Differential Equations, submitted.
- [24] X. FENG AND A. PROHL, *Analysis of total variation flow and its finite element approximations*, M2AN Math. Model. Numer. Anal., 37 (2003), pp. 533–556.
- [25] A. FREIRE, *Uniqueness for the harmonic map flow in two dimensions*, Calc. Var. Partial Differential Equations, 1 (1995), pp. 95–105.
- [26] M. GIAQUINTA, G. MODICA, AND J. SOUCEK, *Variational problems for maps of bounded variation with values in S^1* , Calc. Var. Partial Differential Equations, 1 (1993), pp. 87–121.
- [27] M. GIAQUINTA, G. MODICA, AND J. SOUCEK, *Cartesian Currents in the Calculus Variations, II: Variational Integrals*, Springer-Verlag, New York, 1998.
- [28] Y. GIGA, Y. KASHIMA, AND N. YAMAZAKI, *Local solvability of a constrained gradient system of total variation*, Abstr. Appl. Anal., 8 (2004), pp. 651–682.
- [29] B. GUO AND M. C. HONG, *The Landau-Lifshitz equation of the ferromagnetic spin chain and harmonic maps*, Calc. Var. Partial Differential Equations, 1 (1994), pp. 311–334.
- [30] R. HARDT AND F. H. LIN, *Mappings minimizing the L^p -norm of the gradient*, Comm. Pure Appl. Math., 15 (1987), pp. 555–588.
- [31] N. HUNGERBÜHLER, *Global weak solutions of the p -harmonic flow into homogeneous spaces*, Indiana Math. J., 45 (1996), pp. 275–288.
- [32] N. HUNGERBÜHLER, *Non-uniqueness for the p -harmonic flow*, Canad. Math. Bull., 40 (1997), pp. 793–798.
- [33] N. HUNGERBÜHLER, *m -harmonic flow*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 24 (1997), pp. 593–631.

- [34] N. HUNGERBÜHLER, *Heat flow into spheres for a class of energies*, in Variational Problems in Riemannian Geometry, Progr. Nonlinear Differential Equations Appl. 59, Birkhäuser, Basel, 2004, pp. 45–65.
- [35] O. A. LADYŽENSKAJA, V. A. SOLONNIKOV, AND N. N. UARLCEVA, *Linear and Quasilinear Equations of Parabolic Type*, Transl. Math. Monogr. 23, AMS, Providence, RI, 1967.
- [36] L. D. LANDAU AND E. M. LIFSHITZ, *Electrodynamics of Continuous Media*, Pergamon, Oxford, 1960.
- [37] F. H. LIN, *Static and moving defects in liquid crystals*, in Proceedings of the International Congress of Mathematicians (Kyoto, 1990), Math. Soc. Japan, Tokyo, 1991, pp. 1165–1171.
- [38] X.-G. LIU, *A note on heat flow of p -harmonic mappings*, Kexue Tongbao, 42 (1997), pp. 15–18 (in Chinese).
- [39] X.-G. LIU, *A remark on p -harmonic heat flows*, Chinese Sci. Bull., 42 (1997), pp. 441–444.
- [40] M. MISAWA, *Approximation of p -harmonic maps by the penalized equation* Nonlinear Anal., 47 (2001), pp. 1069–1080.
- [41] M. MISAWA, *On the p -harmonic flow into spheres in the singular case*, Nonlinear Anal., 50 (2002), pp. 485–494.
- [42] L. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Phys. D, 60 (1992), pp. 259–268.
- [43] R. SCHOEN AND K. UHLENBECK, *A regularity theorem for harmonic maps*, J. Differential Geom., 17 (1982), pp. 307–335.
- [44] R. SCHOEN AND K. UHLENBECK, *Regularity of minimizing harmonic maps into the sphere*, Invent. Math., 78 (1984), pp. 89–100.
- [45] J. SIMON, *Compact sets in the space $L^p(0, T; B)$* , Ann. Mat. Pura Appl. (4), 146 (1987), pp. 65–96.
- [46] M. STRUWE, *On the evolution of harmonic maps of Riemannian surfaces*, Comment. Math. Helv., 60 (1985), pp. 558–581.
- [47] M. STRUWE, *Variational Methods*, Springer-Verlag, New York, 1990.
- [48] M. STRUWE, *Geometric evolution problems*, in Nonlinear Partial Differential Equations in Differential Geometry (Park City, UT, 1992), IAS/Park City Math. Ser. 2, AMS, Providence, RI, 1996, pp. 257–339.
- [49] B. TANG, G. SAPIRO, AND V. CASELLES, *Diffusion of general data on non-flat manifolds via harmonic maps theory: The direction diffusion case*, Int. J. Comput. Vision, 36 (2000), pp. 149–161.
- [50] L. A. VESE AND S. J. OSHER, *Numerical methods for p -harmonic flows and applications to image processing*, SIAM J. Numer. Anal., 40 (2002), pp. 2085–2104.
- [51] E. VIRGA, *Variational Theories for Liquid Crystals*, Chapman and Hall, London, 1994.
- [52] E. ZEIDLER, *Nonlinear Functional Analysis and Its Applications*, Vol. II/B, Springer-Verlag, New York, 1990.

GLOBAL EXISTENCE OF SOME INFINITE ENERGY SOLUTIONS FOR A PERFECT INCOMPRESSIBLE FLUID*

RALPH SAXTON[†] AND FERIDE TIĞLAY[†]

Abstract. This paper provides results on local and global existence for a class of solutions to the Euler equations for an incompressible, inviscid fluid. By considering a class of solutions which exhibits a characteristic growth at infinity we obtain an initial value problem for a nonlocal equation. We establish local well-posedness in all dimensions and persistence in time of these solutions for three and higher dimensions. We also examine a weaker class of global solutions.

Key words. incompressible fluid, three dimensions, stagnation point, global existence, Cauchy problem for periodic initial data, infinite energy

AMS subject classifications. 35Q35, 76B03

DOI. 10.1137/080713768

1. Introduction. A fundamental question in the study of fluids concerns the possibility of finite time blow up of solutions to the Euler equations for a perfect, incompressible fluid. It is well known that blow up cannot take place in the two-dimensional case for solutions defined over a bounded domain subject to Dirichlet boundary conditions (see, for example, the work of Wolibner [17] and Ebin [7]), since smooth data in this case lead to solutions remaining smooth for all time. However, the question remains open in higher dimensions.

A separate class of solutions consists of those having “stagnation-point” form, which attracted early attention by Weyl [16] and Lin [13], and provides a set of equations which depend on only a single spatial variable and time. The resulting equations, once solved, provide exact solutions to the full Euler equations. However, the associated growth of the full solutions in certain directions means that the flows possess, at best, only locally finite kinetic energy. Nevertheless, one may discuss such questions as finite time blow up for this class, and it has been shown by Stuart in [15] that this can take place when the reduced equations are defined over the real line and those solutions decay at infinity. In this case the corresponding spatial domain for the full equations is \mathbb{R}^n for $n = 2$ and $n = 3$, with the full set of solutions growing linearly in the other direction(s).

The evolution of two-dimensional solutions, which can blow up on the unbounded domain \mathbb{R}^2 , therefore differs significantly from the class of globally defined solutions which exists for bounded subdomains of \mathbb{R}^2 . So the consequences of higher-dimensional stagnation-point solutions blowing up might be thought to result simply from their behavior at infinity rather than bearing on the question of singularity formation. This view is in some sense strengthened by the results of Childress et al. in [3] (see also a related result by Cox in [5]), where solutions defined over a two-dimensional, infinite strip were examined. Since blow up was still found over this, smaller, semi-infinite domain, stagnation-point solutions defined over such domains

*Received by the editors January 18, 2008; accepted for publication (in revised form) May 27, 2008; published electronically November 19, 2008.

<http://www.siam.org/journals/sima/40-4/71376.html>

[†]Department of Mathematics, University of New Orleans, New Orleans, LA 70148 (rsaxton@math.uno.edu, ftigley@uno.edu). The second author carried out part of this work at the Department of Mathematics at the University of Notre Dame.

would give the appearance of behaving, generally, much as those defined on the full space.

A similar approach is implemented by Constantin in [4] to reduce the three-dimensional Euler equations, periodic in two directions, to a nonlocal Riccati equation and prove the blow up in finite time by solving these equations on characteristics.

In this paper we consider a stagnation-point class of solutions defined over \mathbb{R}^n which is spatially periodic in one coordinate direction. In two dimensions, the equations reduce to those of [3] and, although we examine slightly different boundary conditions, the same blow up results essentially apply. In three and higher dimensions, however, we find a “regularizing” effect not present in solutions which decay to zero at infinity in the same coordinate direction, and this leads to the existence for all time of all such solutions, stemming from sufficiently smooth initial data.

Section 2 sets out the fundamental field equations, which in a basic sense date back to [13], [16]. Section 3 is devoted to local well-posedness (existence, uniqueness, and continuous dependence on initial data) of classical solutions to the initial value problem for the pseudodifferential equation derived in section 2. We establish this result by rewriting the problem on the topological group \mathcal{D} of C^1 class diffeomorphisms as an initial value problem for an ordinary differential equation. Section 4 provides a priori estimates for more regular classes of solutions, leading to global existence of such solutions in three and higher dimensions.

In section 5, we reconsider a class of piecewise affine solutions previously mentioned in [3]. These solutions are less regular than those arising in our existence results. It is found that they exist globally, independently of the underlying dimension.

2. The n -dimensional equation. Consider the n -dimensional Euler equations for an ideal, inviscid and incompressible fluid

$$(2.1a) \quad \partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p = 0,$$

$$(2.1b) \quad \nabla \cdot \mathbf{u} = 0,$$

where $\mathbf{x} = (x_1, \dots, x_n) \equiv (x_1, \mathbf{x}')$. Denoting x_1 by x , $\mathbf{u}(x, \mathbf{x}', t)$ represents the spatial velocity field of the fluid and $p(x, \mathbf{x}', t)$ its pressure. We impose the ansatz

$$(2.2) \quad \mathbf{u}(x, \mathbf{x}', t) = (u(x, t), -\partial_x u(x, t) \mathbf{v}(\mathbf{x}', t)),$$

where the $(n-1)$ -dimensional vector field, \mathbf{v} , will be chosen below. As a consequence of (2.2), (2.1a) may be written as

$$(2.3) \quad \partial_t u + u \partial_x u + \partial_x p = 0,$$

together with

$$(2.4) \quad \partial_t \partial_x u \mathbf{v} + \partial_x u \partial_t \mathbf{v} + u \partial_x^2 u \mathbf{v} - (\partial_x u)^2 \mathbf{v} \cdot \nabla' \mathbf{v} - \nabla' p = \mathbf{0},$$

where the primed operators refer to the variable \mathbf{x}' . Using (2.2) and (2.3), one sees that $\nabla' \partial_x p = 0$. Hence, differentiating (2.4) in x eliminates the pressure term to give

$$(2.5) \quad \partial_x (\partial_t \partial_x u + u \partial_x^2 u) \mathbf{v} + \partial_x^2 u \partial_t \mathbf{v} - \partial_x ((\partial_x u)^2) \mathbf{v} \cdot \nabla' \mathbf{v} = \mathbf{0}.$$

Applying the ∇' operator to (2.5) and using (2.1b) with (2.2) to find that $\nabla' \cdot \mathbf{v} = 1$ shows

$$(2.6) \quad \partial_x (\partial_t \partial_x u + u \partial_x^2 u) - \partial_x ((\partial_x u)^2) \nabla' \mathbf{v} : \nabla' \mathbf{v} = 0,$$

where $\nabla' \mathbf{v} : \nabla' \mathbf{v} = \text{tr}(\nabla' \mathbf{v})^2 = \partial_j v_k \partial_k v_j$ (summing over j, k from 2 to n). For compatibility, we must choose \mathbf{v} such that $\nabla' \mathbf{v} : \nabla' \mathbf{v}$ is independent of \mathbf{x}' . This can be done, for instance, by choosing $\mathbf{v} = \frac{1}{n-1} \mathbf{x}'$, in which case $\nabla' \mathbf{v} : \nabla' \mathbf{v} = \frac{1}{n-1}$ and (2.4) takes the form

$$(2.7) \quad \left(\partial_t \partial_x u + u \partial_x^2 u - \frac{1}{n-1} (\partial_x u)^2 \right) \mathbf{v} - \nabla' p = \mathbf{0}.$$

We note that this particular choice of \mathbf{v} corresponds to the stagnation-point form solution referred to in the introduction. As a result, the solution becomes unbounded in the \mathbf{x}' direction, and hence has infinite energy when considered over the entire n -dimensional domain. We next examine the periodic, initial-boundary value problem, with boundary conditions for $x \in \mathbb{T} \simeq \mathbb{R}/\mathbb{Z}, t \geq 0$, given by

$$(2.8) \quad \mathbf{u}(0, \mathbf{x}', t) = \mathbf{u}(1, \mathbf{x}', t)$$

and

$$(2.9) \quad p(0, \mathbf{x}', t) = p(1, \mathbf{x}', t).$$

Since (2.6) now becomes

$$(2.10) \quad \partial_x (\partial_t \partial_x u + u \partial_x^2 u) - \frac{1}{n-1} \partial_x ((\partial_x u)^2) = 0,$$

which we remark happens to be the x -derivative of a Calogero-class equation (see [2])

$$\partial_x \partial_t u + u \partial_x^2 u - \Phi(\partial_x u) = 0,$$

we obtain the equation

$$(2.11) \quad \partial_t \partial_x u + u \partial_x^2 u - \frac{1}{n-1} (\partial_x u)^2 = f$$

with f purely a function of time. This implies, by (2.7), that

$$(2.12) \quad \nabla' p = \frac{f}{n-1} \mathbf{x}',$$

while by (2.3), $-\partial_x^2 p = \frac{n}{n-1} (\partial_x u)^2 + f$ and so $\Delta p = -\frac{n}{n-1} (\partial_x u)^2$. Finally, for sufficiently smooth functions $u(x, t)$, using (2.8) and integrating (2.11) we have

$$(2.13) \quad f = -\frac{n}{n-1} \int_{\mathbb{T}} (\partial_x u)^2 dx,$$

while (2.3), (2.8), and (2.9) imply that

$$(2.14) \quad \frac{d}{dt} \int_{\mathbb{T}} u dx = 0.$$

Let us introduce the operator ∂_x^{-1} defined by

$$\partial_x^{-1} \phi(x, t) = \int_{x_0}^x \phi(y, t) dy - \int_{\mathbb{T}} \int_{x_0}^x \phi(y, t) dy dx.$$

We make the observation that ∂_x and ∂_x^{-1} generally do not commute since $[\partial_x, \partial_x^{-1}]\phi = \int_{\mathbb{T}} \phi \, dx$, where $[P, Q] = PQ - QP$. Consider (2.11), written in the form

$$(2.15) \quad \partial_x(\partial_t u + u\partial_x u) = \frac{n}{n-1} (\partial_x u)^2 + f(t)$$

with $n > 1$. As a result of (2.13) and the fact that $\partial_x^{-1}\partial_x\phi = \phi - \int_{\mathbb{T}} \phi \, dx$, we may write (2.15) in a nonlocal form as

$$(2.16) \quad \partial_x(\partial_t u + u\partial_x u) = \frac{n}{n-1} \partial_x^{-1}\partial_x((\partial_x u)^2)$$

and then, using the periodicity of u , we obtain

$$(2.17) \quad \partial_t u + u\partial_x u = \frac{n}{n-1} \partial_x^{-2}\partial_x((\partial_x u)^2).$$

Okamoto and Zhu [14] previously established local existence for (2.10) with $u \in H^2$, using a method introduced by Kato and Lai. Their approach requires showing uniqueness separately and then using uniqueness to prove continuous dependence on initial data. Here we instead derive a local well-posedness result in C^1 which follows from Picard iteration after rewriting the equation as an ordinary differential equation on an infinite-dimensional Banach space. This method can also be used to prove well-posedness in Sobolev spaces $H^s(\mathbb{T})$ for $s > 3/2$ (see [10] and [12], for example, for a similar result for the Camassa–Holm equation).

3. Local existence of classical solutions. In [1], Arnold observed that the initial value problem for the classical Euler equations of a perfect fluid can be stated as a geometric problem of finding geodesics on the group of volume preserving diffeomorphisms. Following this observation, Ebin and Marsden [8] developed the functional analytic tools to establish sharp local well-posedness results for the Euler equations. This method has since been used for other equations with similar geometric interpretations; for example, Misiolek [12] obtained local well-posedness in $C^1(\mathbb{T})$ for the Camassa–Holm equation, which is the equation for geodesics of the H^1 metric on the Virasoro group.

In this section we develop an appropriate analytic framework for (2.17), using a similar approach to prove the following theorem.

THEOREM 3.1. *Suppose that $n > 1$. Then there exists a unique solution*

$$u \in C^0([0, T], C^1(\mathbb{T})) \cap C^1([0, T], C^0(\mathbb{T}))$$

to the Cauchy problem for (2.17) with initial data $u_0 \in C^1(\mathbb{T})$ for some $T > 0$, and the solution depends continuously on the initial data.

Let γ be the flow generated by u , that is,

$$\frac{d\gamma}{dt}(x, t) = u(\gamma(x, t), t),$$

or $u = \dot{\gamma} \circ \gamma^{-1}$. Then we obtain the equation

$$(3.1) \quad \ddot{\gamma} = \frac{n}{n-1} \partial_x^{-2}\partial_x((\partial_x(\dot{\gamma} \circ \gamma^{-1}))^2) \circ \gamma$$

from (2.17). Therefore it is sufficient to prove that

$$F(X, \gamma) = \frac{n}{n-1} (\partial_x^{-2}\partial_x((\partial_x(X \circ \gamma^{-1}))^2)) \circ \gamma$$

defines a continuously differentiable vector field in a neighborhood of the identity on the topological group \mathcal{D} of C^1 class diffeomorphisms. Then Theorem 3.1 follows by Picard iteration over Banach spaces.

We remark that the smooth dependence on initial data for (3.1) implies only continuous dependence on initial data for (2.17). The map $\gamma \rightarrow \gamma^{-1}$ is continuous but not locally Lipschitz [8], and this prevents obtaining more regularity for the initial data to solution map by this method. The question of whether the regularity of the solution map $u_0 \rightarrow u(t)$ can be improved, or not, is open. It is known, for instance, that it is not possible to improve the regularity of this map for the Camassa–Holm equation in Sobolev spaces [11].

In the remainder of this section, C_γ will represent a generic constant depending only on the C^1 norms of γ and γ^{-1} .

Proof of Theorem 3.1. Let us denote by P_γ the operator given by conjugation

$$P_\gamma(\phi) := P(\phi \circ \gamma^{-1}) \circ \gamma$$

for any $\gamma \in \mathcal{D}$ and pseudodifferential operator P . Using this notation we write

$$F(X, \gamma) = \frac{n}{n-1} (\partial_x^{-2} \partial_x)_\gamma ((\partial_x)_\gamma X)^2.$$

Next we compute the directional derivative $\partial_\gamma F_{(X,\gamma)}$ and prove that it is a bounded linear map.

Note that $(\partial_x^{-2} \partial_x) f = \partial_x^{-1} \{f - \int_0^1 f \, dx\}$ is a bounded operator from $C^0(\mathbb{T})$ into $C^1(\mathbb{T})$. Furthermore we abuse the notation slightly and denote by \mathcal{D} the connected component of orientation preserving C^1 diffeomorphisms of \mathbb{T} .

Let $s \rightarrow \gamma_s$ be a smooth curve in \mathcal{D} such that $\gamma_0 = id$ and $\partial_s \gamma_s|_{s=0} = W$ for $W \in C^1(\mathbb{T})$. Then we have

$$(3.2) \quad \partial_\gamma F_{(X,\gamma)}(W) = \frac{n}{n-1} \{(\partial_\varepsilon G_\varepsilon)|_{\varepsilon=0} \circ \gamma + W(\partial_x G_\varepsilon)_{\varepsilon=0} \circ \gamma\},$$

where $G_\varepsilon = \partial_x^{-2} \partial_x \{(\partial_x(X \circ \gamma_\varepsilon^{-1}))^2\}$. We know that $\partial_x \partial_x^{-1}$ gives the identity; hence the second summand on the right in (3.2) can be written as

$$(3.3) \quad W \partial_x G_\varepsilon|_{\varepsilon=0} \circ \gamma = \left\{ ((\partial_x)_\gamma X)^2 - \int_0^1 (\partial_x(X \circ \gamma^{-1}))^2 dx \right\} W.$$

Moreover, the computation of the first summand on the right in (3.2) is reduced by

$$(3.4) \quad (\partial_\varepsilon G_\varepsilon)|_{\varepsilon=0} = \partial_x^{-2} \partial_x (\partial_\varepsilon H_\varepsilon|_{\varepsilon=0})$$

to determine $\partial_\varepsilon H_\varepsilon|_{\varepsilon=0}$, where $H_\varepsilon = (\partial_x(X \circ \gamma_\varepsilon^{-1}))^2$. Here a straightforward computation leads to

$$(3.5) \quad \partial_\varepsilon H_\varepsilon|_{\varepsilon=0} = -(\partial_x(X \circ \gamma^{-1}))^2 \partial_x(W \circ \gamma^{-1}) - \partial_x \{ (W \circ \gamma^{-1}) (\partial_x(X \circ \gamma^{-1}))^2 \}.$$

Then, after an integration by parts, we obtain

$$\begin{aligned} \partial_\varepsilon G_\varepsilon|_{\varepsilon=0} &= -(\partial_x(X \circ \gamma^{-1}))^2 (W \circ \gamma^{-1}) + \int_0^1 (\partial_x(X \circ \gamma^{-1}))^2 (W \circ \gamma^{-1}) dx \\ &\quad - \partial_x^{-2} \partial_x \{ (\partial_x(X \circ \gamma^{-1}))^2 \partial_x(W \circ \gamma^{-1}) \} \\ &\quad + \left(x - \frac{1}{2} \right) \{ (\partial_x(X \circ \gamma^{-1}))^2 (W \circ \gamma^{-1}) \}|_0^1. \end{aligned}$$

The last term on the left-hand side of the above inequality vanishes since X and W are periodic functions and γ is an orientation preserving diffeomorphism. Therefore we have

$$\begin{aligned}
 \partial_\gamma F_{(X,\gamma)}(W) &= \frac{n}{n-1} \left\{ -\partial_x^{-2} \partial_x \{ (\partial_x(X \circ \gamma^{-1}))^2 \partial_x(W \circ \gamma^{-1}) \} \circ \gamma \right. \\
 (3.6) \qquad \qquad \qquad &+ \int_0^1 (\partial_x(X \circ \gamma^{-1}))^2 W \circ \gamma^{-1} dx \\
 &\left. - W \int_0^1 (\partial_x(X \circ \gamma^{-1}))^2 dx \right\}.
 \end{aligned}$$

The linearity of the map $W \rightarrow \partial_\gamma F_{(X,\gamma)}(W)$ is clear. Thus we proceed to show that it is bounded. It is sufficient to estimate the C^1 norms of all three summands on the right in (3.6). The second and third terms are both bounded by $C_\gamma \|W\|_{C^0} \|X\|_{C^1}^2$. For the first term on the right in (3.6), we have

$$\begin{aligned}
 \|(\partial_x^{-2} \partial_x)_\gamma \{ ((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W \} \|_{C^1} &\leq \|(\partial_x^{-2} \partial_x)_\gamma \{ ((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W \} \|_{C^0} \\
 (3.7) \qquad \qquad \qquad &+ \|((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W \|_{C^0} \|\gamma\|_{C^1},
 \end{aligned}$$

which is bounded by $C_\gamma \|X\|_{C^1}^2 \|W\|_{C^1}$.

In the direction of X , the Gâteaux derivative of F is given by

$$\partial_X F_{(X,\gamma)}(W) = \frac{2n}{n-1} (\partial_x^{-2} \partial_x)_\gamma ((\partial_x)_\gamma X (\partial_x)_\gamma W),$$

and this is a bounded map since

$$\begin{aligned}
 \|\partial_X F_{(X,\gamma)}(W)\|_{C^1} &\leq C_n \|(\partial_x^{-2} \partial_x)_\gamma ((\partial_x)_\gamma X (\partial_x)_\gamma W)\|_{C^0} \\
 &+ \left\| (\partial_x)_\gamma X (\partial_x)_\gamma W - \int_{\mathbb{T}} (\partial_x)_\gamma X (\partial_x)_\gamma W dx \right\|_{C^0} \|\gamma\|_{C^1} \\
 &\leq C_{n,\gamma} \|X\|_{C^1} \|W\|_{C^1},
 \end{aligned}$$

where $C_{n,\gamma}$ depends only on n and C^1 norms of γ and γ^{-1} .

In order to complete the proof of Theorem 3.1 it is sufficient to show that F is Fréchet differentiable; i.e., both directional derivatives $\partial_\gamma F$ and $\partial_X F$ are continuous maps.

Continuity of $(X, \gamma) \rightarrow \partial_\gamma F_{(X,\gamma)}(W)$. The following inequality reduces the proof of continuity of $\partial_\gamma F_{(X,\gamma)}(W)$ to estimating the two summands on the right-hand side:

$$\begin{aligned}
 \|\partial_\gamma F_{(X_1,\gamma_1)}(W) - \partial_\gamma F_{(X_2,\gamma_2)}(W)\|_{C^1} &\leq \|\partial_\gamma F_{(X_1,\gamma_1)}(W) - \partial_\gamma F_{(X_2,\gamma_1)}(W)\|_{C^1} \\
 &+ \|\partial_\gamma F_{(X_2,\gamma_1)}(W) - \partial_\gamma F_{(X_2,\gamma_2)}(W)\|_{C^1}
 \end{aligned}$$

where the inequality holds up to a constant depending on n . We rewrite the C^1 norm that we wish to estimate to show continuity in X as

$$\begin{aligned}
 (3.8) \quad \|\partial_\gamma F_{(X_1, \gamma)} - \partial_\gamma F_{(X_2, \gamma)}\|_{C^1} &\leq \|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X_1)^2 - ((\partial_x)_\gamma X_2)^2\} (\partial_x)_\gamma W\|_{C^1} \\
 &+ \left| \int_0^1 (\partial_x(X_1 \circ \gamma - X_2 \circ \gamma))^2 W \circ \gamma^{-1} dx \right| \\
 &+ |W| \left| \int_0^1 (\partial_x(X_1 \circ \gamma - X_2 \circ \gamma))^2 dx \right|.
 \end{aligned}$$

The last two summands in (3.8) are bounded by $C_\gamma \|X_1 - X_2\|_{C^1}^2 \|W\|_{C^0}$. For the remaining term, we have

$$\begin{aligned}
 &\|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X_1)^2 - ((\partial_x)_\gamma X_2)^2\} (\partial_x)_\gamma W\|_{C^1} \\
 &\leq \|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X_1)^2 - ((\partial_x)_\gamma X_2)^2\} (\partial_x)_\gamma W\|_{C^0} \\
 &\quad + \|((\partial_x(X_1 \circ \gamma^{-1}))^2 - (\partial_x(X_2 \circ \gamma^{-1}))^2) \partial_x(W \circ \gamma^{-1})\|_{C^0} \|\partial_x \gamma\|_{C^0},
 \end{aligned}$$

which is bounded by $C_\gamma \|X_1 - X_2\|_{C^1} \|X_1 + X_2\|_{C^1} \|W\|_{C^1}$.

Our next estimate establishes continuity of $\gamma \rightarrow \partial_\gamma F_{(X, \gamma)}(W)$. Note that it is sufficient to consider

$$\begin{aligned}
 (3.9) \quad &\|\partial_\gamma F_{(X, \gamma)}(W) - \partial_\gamma F_{(X, id)}(W)\|_{C^1} \\
 &\leq \|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W\} - \partial_x^{-2} \partial_x ((\partial_x X)^2 \partial_x W)\|_{C^1}
 \end{aligned}$$

$$(3.10) \quad + \left| \int_0^1 \{(\partial_x(X \circ \gamma^{-1}))^2 (W \circ \gamma^{-1}) - (\partial_x X)^2 W\} dx \right|$$

$$(3.11) \quad + \|W\|_{C^0} \left| \int_0^1 \{(\partial_x(X \circ \gamma^{-1}))^2 - (\partial_x X)^2\} dx \right|,$$

where the inequality is up to a constant depending on n . After adding and subtracting the appropriate terms, (3.10) is bounded by

$$(3.12) \quad \|W\|_{C^0} \|\partial_x X\|_{C^0} \|(\partial_x X \circ \gamma^{-1}) \partial_x \gamma^{-1} - \partial_x X\|_{C^0} + \|\partial_x X\|_{C^0}^2 \|W\|_{C^1} \|\gamma - id\|_{C^0},$$

which is bounded (up to a constant C_γ) by $\|X\|_{C^1}^2 \|W\|_{C^1} \|\gamma - id\|_{C^1}$. The term in (3.11) is estimated similarly. Hence in order to establish continuity in γ it is sufficient to bound

$$\begin{aligned}
 (3.13) \quad &\|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W\} - \partial_x^{-2} \partial_x ((\partial_x X)^2 \partial_x W)\|_{C^1} \\
 &\leq \|(\partial_x^{-2} \partial_x)_\gamma \{((\partial_x)_\gamma X)^2 (\partial_x)_\gamma W\} - \partial_x^{-2} \partial_x ((\partial_x X)^2 \partial_x W)\|_{C^0}
 \end{aligned}$$

$$(3.14) \quad + \|\{(\partial_x(X \circ \gamma^{-1}))^2 \partial_x(W \circ \gamma^{-1})\} \circ \gamma \partial_x \gamma - (\partial_x X)^2 \partial_x W\|_{C^0}.$$

The norm in (3.14) is equal to

$$(3.15) \quad \|\partial_x W (\partial_x X)^2 \{(\partial_x \gamma^{-1})^2 \circ \gamma - 1\}\|_{C^0},$$

which is bounded by $C_\gamma \|\partial_x X\|_{C^0}^2 \|\partial_x W\|_{C^0} \|\gamma - id\|_{C^1}$. For (3.13) it is sufficient to estimate

$$(3.16) \quad \|(\partial_x^{-2} \partial_x)_\gamma S - \partial_x^{-2} \partial_x S\|_{C^0} + \|\partial_x^{-2} \partial_x S - \partial_x^{-2} \partial_x ((\partial_x X)^2 \partial_x W)\|_{C^0},$$

where

$$S = S(X, \gamma, W) = (\partial_x(X \circ \gamma^{-1}))^2 \circ \gamma \partial_x(W \circ \gamma^{-1}) \circ \gamma.$$

After a change of variables we regroup the terms in the first summand of (3.16) to obtain

(3.17)

$$(\partial_x^{-2} \partial_x)_\gamma S - \partial_x^{-2} \partial_x S = \int_0^1 \int_{\gamma^{-1}(x)}^x S(y) \partial_x \gamma(y) dy dx + \int_{x_0}^x S(y) (\partial_x \gamma(y) - 1) dy$$

(3.18)

$$- \int_0^1 \int_{x_0}^x S(y) (\partial_x \gamma(y) - 1) dy dx$$

(3.19)

$$- \left(\gamma(x) - \frac{1}{2} \right) \int_0^1 S \circ \gamma^{-1}(y) dy + \left(x - \frac{1}{2} \right) \int_0^1 S(y) dy.$$

The right-hand side of this equality can be simplified further as

(3.20)

$$(\partial_x^{-2} \partial_x)_\gamma S - \partial_x^{-2} \partial_x S = \int_0^1 \int_{\gamma^{-1}(x)}^x S(y) \partial_x \gamma(y) dy dx$$

(3.21)

$$+ \partial_x^{-2} \partial_x (S(\partial_x \gamma - 1)) - (\gamma(x) - x) \int_0^1 S \circ \gamma^{-1}(y) dy.$$

Clearly the sup norms of all three summands in (3.20)–(3.21) vanish in the limit as γ goes to id in C^1 . The second summand in (3.16) is equal to the C^0 norm of

(3.22)

$$\partial_x^{-2} \partial_x (\partial_x W (\partial_x X)^2 \{(\partial_x \gamma^{-1})^3 \circ \gamma - 1\}),$$

which is bounded by $C_\gamma \|W\|_{C^1} \|X\|_{C^1}^2 \|\gamma - id\|_{C^1}$. Hence the continuity of $\gamma \rightarrow \partial F_{(X,\gamma)}$ follows.

The continuity of $(X, \gamma) \rightarrow \partial_X F_{(X,\gamma)}$ can be shown analogously. Therefore $F(X, \gamma)$ defines a continuously differentiable map in a neighborhood of (id, u_0) . This completes the proof of Theorem 3.1. \square

4. Global existence for $n \geq 3$. In this section we investigate the persistence of solutions of the initial value problem for (2.17) and show that, unlike the two-dimensional case where solutions may blow up in finite time (see [3], [5]), they persist for $n \geq 3$ in the appropriate function spaces.

For the theorem below let us use the following notation:

$$X_n(\mathbb{T}) = \begin{cases} W^{2,\infty}(\mathbb{T}), & n = 3, \\ W^{2,\frac{n-1}{n-3}}(\mathbb{T}), & n > 3, \end{cases}$$

and

$$Y_n(\mathbb{T}) = \begin{cases} W^{1,\infty}(\mathbb{T}), & n = 3, \\ W^{1,\frac{n-1}{n-3}}(\mathbb{T}), & n > 3. \end{cases}$$

THEOREM 4.1. *Let $n \geq 3$ and assume that $u_0(x) \in X_n(\mathbb{T})$. Then the solution $u(x, t)$ from Theorem 3.1 has a unique extension to all $T < \infty$ such that $u(x, t) \in C^0([0, T], X_n(\mathbb{T})) \cap C^1([0, T], Y_n(\mathbb{T}))$.*

Proof. Consider (2.10), expressed in the form

$$(4.1) \quad \partial_t \partial_x^2 u + u \partial_x^3 u + \frac{n-3}{n-1} \partial_x u \partial_x^2 u = 0.$$

Using the flow γ of u ($\dot{\gamma} = u \circ \gamma$) in (4.1), we first solve for $\partial_x^2 u$:

$$(4.2) \quad \partial_x^2 u \circ \gamma(t) = u_0'' \exp\left(-\frac{n-3}{n-1} \int_0^t \partial_x u \circ \gamma(s) ds\right).$$

By the identity $\partial_x \dot{\gamma} = \partial_x u \circ \gamma \partial_x \gamma$, we also have

$$(4.3) \quad \partial_x \gamma(t) = \exp\left(\int_0^t \partial_x u \circ \gamma(s) ds\right).$$

Therefore

$$(4.4) \quad \partial_x^2 u \circ \gamma(t) (\partial_x \gamma(t))^{\frac{n-3}{n-1}} = u_0''.$$

Note that Theorem 3.1 implies that $\gamma \in C^1$ locally in time and it follows, given $\partial_x \gamma(0) = 1$, that there exists an interval, $t \in [0, \tau(\varepsilon))$, over which $0 < \varepsilon \leq \inf_{x \in \mathbb{T}} \partial_x \gamma(t) \leq \sup_{x \in \mathbb{T}} \partial_x \gamma(t) \leq \varepsilon^{-1}$. Equation (4.4) then implies that, locally, $u \in X_n(\mathbb{T})$, since γ maps \mathbb{T} diffeomorphically to itself. In turn, (2.16) shows that $\partial_t u \in Y_n(\mathbb{T})$ over the same time interval. (With additional assumptions on the data, further regularity can also be bootstrapped to higher derivatives.)

Assuming then that sufficient smoothness holds locally in time, we find on multiplying (4.1) by $|\partial_x^2 u|^{p-2} \partial_x^2 u$ that

$$(4.5) \quad \partial_t |\partial_x^2 u|^p + u \partial_x |\partial_x^2 u|^p + p \frac{n-3}{n-1} \partial_x u |\partial_x^2 u|^p = 0.$$

Since (2.2) and (2.8) imply that both u and $\partial_x u$ are periodic functions of x , the same is true of $\partial_x^2 u$, by (2.11). One therefore obtains, on integrating (4.5) over \mathbb{T} ,

$$(4.6) \quad \frac{d}{dt} \int_{\mathbb{T}} |\partial_x^2 u|^p dx + \left(p \frac{n-3}{n-1} - 1\right) \int_{\mathbb{T}} \partial_x u |\partial_x^2 u|^p dx = 0,$$

from which it follows that the $L^{\frac{n-1}{n-3}}(\mathbb{T})$ norm of $\partial_x^2 u$ is uniformly conserved in time for $n > 3$. The case $n = 3$ can either be considered as the limit $n \rightarrow 3$ with $p \rightarrow \infty$ in (4.6), or directly using (4.1) which shows that $\partial_x^2 u$ is constant along characteristics and hence its $L^\infty(\mathbb{T})$ norm is uniformly conserved.

Periodicity of $u(x, t)$ in x implies there exists a zero for $\partial_x u$, say at $x = x_0(t)$, and so for $x, x_0 \in \mathbb{T}$,

$$\partial_x u(x, t) = \int_{x_0}^x \partial_y^2 u(y, t) dy.$$

For $n > 3$, we therefore have the estimate

$$|\partial_x u(x, t)| \leq |x - x_0|^{\frac{2}{n-1}} \|u_0''\|_{\frac{n-1}{n-3}} \leq \|u_0''\|_{\frac{n-1}{n-3}}$$

using Hölder’s inequality, and so

$$\|\partial_x u\|_\infty \leq \|u_0''\|_{\frac{n-1}{n-3}}.$$

If $n = 3$, then

$$\|\partial_x^2 u\|_\infty = \|u_0''\|_\infty,$$

which means

$$|\partial_x u(x, t)| \leq |x - x_0| \|u_0''\|_\infty \leq \|u_0''\|_\infty$$

for all $x \in \mathbb{T}$, and so

$$\|\partial_x u\|_\infty \leq \|u_0''\|_\infty$$

for all $t > 0$.

Further, since $u(x, t) - u_0(x)$ has mean zero by (2.14), there exists $x = x_1(t)$, where $u(x_1, t) = u_0(x_1)$ and, for $x, x_1 \in \mathbb{T}$, we have

$$u(x, t) = u_0(x) + \int_{x_1}^x (\partial_y u(y, t) - u_0'(y)) dy.$$

It follows that

$$|u(x, t)| \leq \|u_0\|_\infty + |x - x_1| (\|\partial_x u\|_\infty + \|u_0'\|_\infty)$$

for all $x \in \mathbb{T}$, which gives the inequality

$$\|u\|_\infty \leq \|\partial_x u\|_\infty + \|u_0\|_{C^1}.$$

Combining the results of the previous two paragraphs shows that

$$(4.7) \quad \|u\|_{C^1} \leq \|u_0\|_{C^2} \quad \text{for } n = 3,$$

and

$$(4.8) \quad \|u\|_{C^1} \leq \|u_0\|_{C^1} + \|u_0''\|_{\frac{n-1}{n-3}} \quad \text{for } n > 3.$$

Finally, on using the properties of the operators $\partial_x^{-1} \partial_x$ and $\partial_x^{-2} \partial_x$ in (2.16) and (2.17) together with the above estimates, it is seen that $\|\partial_t u\|_\infty$ and $\|\partial_t \partial_x u\|_\infty$ are majorized by a function of $\|u_0\|_{C^2}$ for $n = 3$, while $\|\partial_t u\|_\infty$ and $\|\partial_t \partial_x u\|_{\frac{n-1}{n-3}}$ are majorized by a function of $\|u_0\|_{W^{2, \frac{n-1}{n-3}}}$ for $n > 3$.

In both of these cases it follows that the $C^1(\mathbb{T})$ norm of u and the $C^0(\mathbb{T})$ norm of $\partial_t u$ remain uniformly bounded in time over any interval of local existence and, by bootstrapping the arguments of Theorem 3.1, the solution can be continued, globally, in time. \square

As a remark, we note here how a blow up argument made in [3], which involves a nontrivial class of separable solutions to (2.10) for $n = 2$, fails to apply in the case $n > 3$. In particular, the possible appearance of a $(\tau - t)^{-1}$ factor, $\tau > 0$, in the two-dimensional case no longer exists in higher dimensions.

Given the solution form $u(x, t) = X(x)T(t)$, (2.10) reduces to

$$(4.9) \quad \lambda X''(x) + \frac{n-3}{n-1} X'(x)X''(x) + X(x)X'''(x) = 0, \quad x \in \mathbb{T},$$

where

$$(4.10) \quad \dot{T}(t) - \lambda T(t)^2 = 0, \quad t \geq 0,$$

and λ is a constant. Multiplying (4.9) by $|X''(x)|^{\frac{5-n}{n-3}} X''(x)$ now gives

$$(4.11) \quad \lambda |X''(x)|^{\frac{n-1}{n-3}} + \frac{n-3}{n-1} (X'(x) |X''(x)|^{\frac{n-1}{n-3}} + X(x) (|X''(x)|^{\frac{n-1}{n-3}})') = 0.$$

By using periodicity, an integration of (4.11) over \mathbb{T} for $n > 3$ therefore shows

$$(4.12) \quad \lambda \int_{\mathbb{T}} |X''(x)|^{\frac{n-1}{n-3}} dx = 0$$

and the result then follows. We note that there are in general one or more points of inflection in nontrivial, periodic solutions, which prevents this argument from holding in two dimensions.

5. Weak solutions. In this section, we construct a basic, piecewise differentiable class of weak solutions $u(x, t) \in C^0([0, T], PC^1(\mathbb{T})) \cap C^1([0, T], PC^0(\mathbb{T}))$ to (2.11), which are found to exist for all $T > 0$, regardless of the underlying dimension.

For every vector field $\Phi(\mathbf{x}, t) \in C_0^\infty([0, T] \times \mathbb{T} \times \mathbb{R}^{n-1}; \mathbb{R}^n)$ such that $\nabla \cdot \Phi = 0$, and for every scalar function $\theta(\mathbf{x}, t) \in C_0^\infty([0, T] \times \mathbb{T} \times \mathbb{R}^{n-1}; \mathbb{R})$, the velocity field $\mathbf{u}(\mathbf{x}, t)$ in (2.1b) satisfies

$$(5.1) \quad \int_Q \partial_t \Phi \cdot \mathbf{u} + (\nabla \Phi \mathbf{u}) \cdot \mathbf{u} \, d\mathbf{x} \, dt + \int_{\mathbb{T} \times \mathbb{R}^{n-1}} \Phi(\mathbf{x}, 0) \cdot \mathbf{u}(\mathbf{x}, 0) \, d\mathbf{x} = 0,$$

$$(5.2) \quad \int_{\mathbb{R}^n} \nabla \theta \cdot \mathbf{u} \, d\mathbf{x} = 0,$$

where $Q = [0, T] \times \mathbb{T} \times \mathbb{R}^{n-1}$.

In terms of (2.2), (5.1) and (5.2) reduce to

$$(5.3) \quad \begin{aligned} & \int_Q \partial_t \phi u - \partial_t \Phi' \cdot \mathbf{v} \partial_x u \, d\mathbf{x} \, dt \\ & + \int_Q \partial_x \phi u^2 - (\partial_x \Phi' + \nabla' \phi) \cdot \mathbf{v} u \partial_x u + (\nabla' \Phi' \mathbf{v}) \cdot \mathbf{v} (\partial_x u)^2 \, d\mathbf{x} \, dt \\ & + \int_{\mathbb{T} \times \mathbb{R}^{n-1}} \phi(\mathbf{x}, 0) u(x, 0) - \Phi'(\mathbf{x}, 0) \cdot \mathbf{v}(\mathbf{x}, 0) \partial_x u(x, 0) \, d\mathbf{x} = 0, \end{aligned}$$

$$(5.4) \quad \int_{\mathbb{T} \times \mathbb{R}^{n-1}} \partial_x \theta u - \nabla' \theta \cdot \mathbf{v} \partial_x u \, d\mathbf{x} = 0,$$

in which we have used the notation $\Phi = (\phi, \Phi')$ to distinguish the first component from the remaining $n - 1$ components of Φ . Denoting by $[\mathbf{u}] = \mathbf{u}_+ - \mathbf{u}_-$ the jump in \mathbf{u} across any smooth surface of discontinuity, S , and considering test functions whose support crosses S , (5.2) shows that

$$(5.5) \quad [\mathbf{u}] \cdot \mathbf{n} = 0,$$

where $\mathbf{n} = (\mathbf{n}, \mathbf{n}')$ is normal to S . Then, by (5.4), we have

$$(5.6) \quad [u]\mathbf{n} - [\partial_x u]\mathbf{v} \cdot \mathbf{n}' = 0, \text{ where } \mathbf{v} = \frac{1}{n-1}\mathbf{x}'.$$

In examining weak, frontlike, piecewise continuous solutions for which $[u] = 0$ and $[\partial_x u] \neq 0$ (see [6]), it follows that these discontinuities propagate so that $\mathbf{n}' \cdot \mathbf{x}' = 0$. A weak formulation specific to such discontinuities may be derived by means of appropriate choice of test functions from (5.3), or by observing that (2.11) may be written in conservation form as

$$(5.7) \quad \partial_t((\partial_x u)^{1-n}) + \partial_x(u(\partial_x u)^{1-n}) + (n-1)(\partial_x u)^{-n}f = 0, \quad x \in \mathbb{T}.$$

We will admit weak solutions, $u(x, t)$, which satisfy the relation

$$(5.8) \quad \int_{\Omega} \partial_t \varphi (\partial_x u)^{1-n} + \partial_x \varphi u (\partial_x u)^{1-n} - (n-1)f \varphi (\partial_x u)^{-n} dx dt$$

$$(5.9) \quad + \int_{\mathbb{T}} \varphi(x, 0) (\partial_x u)^{1-n}(x, 0) dx = 0$$

for all $\varphi(x, t) \in C_0^\infty(\Omega)$, where $\Omega = [0, T) \times \mathbb{T}$. Using standard Rankine–Hugoniot-type arguments [6], discontinuities in $\partial_x u$ that jump across a curve $x = \psi(t)$ are seen to satisfy

$$(5.10) \quad (-\dot{\psi} + u(\psi, t))[(\partial_x u)^{1-n}] = 0,$$

and such discontinuities therefore propagate with the flow of (2.11); i.e., $\psi(t)$ is a member of the characteristic family $\dot{\gamma} = u \circ \gamma$.

5.1. Piecewise affine solutions. We begin by commenting on the general case of periodic, N -phase, piecewise affine solutions. Given that both $\partial_x u$ and $\partial_t u$ may be discontinuous across the curves $x = \psi_i(t), 1 \leq i \leq N - 1$, in order for u to remain continuous there we must have $[u](\gamma(t), t) = 0$, and so $\frac{d}{dt}[u](\gamma(t), t) = 0$. As a result, $\partial_t[u] + u\partial_x[u] = 0$, and first derivative jumps are seen to satisfy the relations

$$(5.11) \quad [\partial_t u] + u[\partial_x u] = 0 \text{ and } [\partial_x p] = 0$$

from (2.3). Under these conditions on u , the expression for $f(t)$, which was obtained in (2.13) by integrating (2.11) for $u \in C^1(\mathbb{T})$, remains unchanged:

$$(5.12) \quad f = -\frac{n}{n-1} \int_{\mathbb{T}} (\partial_x u)^2 dx.$$

In the special case, $N = 2$, which we consider here, our form of solution is given by the periodic extension of the function

$$(5.13) \quad u(x, t) = \alpha + \begin{cases} xp, & x \in (0, \tilde{\phi}), \\ \tilde{\phi}p + (x - \tilde{\phi})q, & x \in (\tilde{\phi}, \tilde{\psi}), \\ \tilde{\phi}p + (\tilde{\psi} - \tilde{\phi})q + (x - \tilde{\psi})p, & x \in (\tilde{\psi}, 1), \end{cases}$$

where $\tilde{\phi} = \phi - [[\phi]]$, $\tilde{\psi} = \psi - [[\psi]]$, with $[[\cdot]]$ denoting the “integer part” of the argument. The functions $\tilde{\phi}(t)$ and $\tilde{\psi}(t)$ are the representatives in $[0, 1]$ of the phase

curves $\phi(t), \psi(t) \in (-\infty, \infty)$, which start out from $\phi(0), \psi(0) \in [0, 1]$ and separate those regions where $\partial_x u(x, t)$ periodically takes on values of $p(t)$ or $q(t)$.

Proceeding heuristically for the moment, periodicity of $u(x, t)$ requires that

$$(5.14) \quad \mathcal{N}(t) = \phi(t)p(t) + (\psi(t) - \phi(t))q(t) + (1 - \psi(t))p(t) = 0.$$

Given the spatial periodicity in pressure (see (2.9)), we recall that integration of (2.3) over one period showed the integral $\int_0^1 u(x, t) dx$ to be independent of time. This allows (5.13) to be used to give an expression for $\alpha(t)$. The result may be written, for instance, in terms of $p, \tilde{\phi}$, and $\tilde{\psi}$, as

$$(5.15) \quad \alpha + \frac{p}{2}(\tilde{\phi} + \tilde{\psi} - 1) = c,$$

where we have set $c = \int_0^1 u(x, 0) dx = \alpha(0) + \frac{p(0)-q(0)}{2}(\psi(0) - \phi(0))(\phi(0) + \psi(0) - 1)$. We will assume that $0 < \phi(0) < \psi(0) < 1$. Choosing, for convenience, $c = 0$, the ‘‘average characteristic’’ must propagate with speed zero and we will see, consequently, that both $\phi(t)$ and $\psi(t)$ remain in $[0, 1]$ for all time. The distinctions between $\tilde{\phi}$ and $\phi, \tilde{\psi}$, and ψ will therefore not be made further here.

Combining (2.11), (5.12), and (5.13) now gives

$$(5.16) \quad \dot{p} = \frac{1}{n-1}p^2 + f, \quad \dot{q} = \frac{1}{n-1}q^2 + f,$$

where

$$(5.17) \quad f = -\frac{n}{n-1}(\phi p^2 + (\psi - \phi)q^2 + (1 - \psi)p^2).$$

Also, by (5.10),

$$(5.18) \quad \dot{\phi} = \alpha + \phi p$$

and

$$(5.19) \quad \dot{\psi} = \alpha + \phi p + (\psi - \phi)q.$$

We next verify that (5.14) follows from the system of equations (5.16)–(5.19). Differentiation and some simplification gives

$$\begin{aligned} \dot{\mathcal{N}} &= -(\psi - \phi)(p - q)q + (\psi - \phi)q^2 + (\phi + (1 - \psi))\left(\frac{p^2}{n-1} + f\right) + (\psi - \phi)\left(\frac{q^2}{n-1} + f\right) \\ &= -(\psi - \phi)(p - q)q - ((\phi + (1 - \psi))p^2 + (\psi - \phi)q^2) \\ &= -p(q(\psi - \phi) + p(\phi + (1 - \psi))) \\ &= -p\mathcal{N}, \end{aligned}$$

and so $\mathcal{N}(t) = \mathcal{N}(0)e^{-\int_0^t p(s) ds}$. In particular, taking periodic data, $\mathcal{N}(0) = 0$, means that $\mathcal{N}(t) = 0, t > 0$. We may therefore use (5.14) to write (5.19) as

$$(5.20) \quad (1 - \psi)' = -\alpha + (1 - \psi)p,$$

and again employ (5.14) to express the following phase fractions as functions of p and q :

$$(5.21) \quad \psi - \phi = \frac{p}{p - q}, \quad \phi + (1 - \psi) = \frac{-q}{p - q}.$$

Using these relations in (5.17) leads to

$$(5.22) \quad f = \frac{n}{n - 1}pq$$

from which (5.16) reduces to the autonomous system

$$(5.23) \quad \dot{p} = \frac{1}{n - 1}(p^2 + n pq),$$

$$(5.24) \quad \dot{q} = \frac{1}{n - 1}(q^2 + n pq).$$

Subtracting (5.18) from (5.19) implies

$$(5.25) \quad \psi(t) - \phi(t) = (\psi(0) - \phi(0)) \exp\left(\int_0^t q(s) ds\right),$$

which means that the center phase fraction does not collapse as long as $\int_0^t q(s) ds$ remains bounded away from $-\infty$. Similarly, adding (5.18) and (5.20) gives

$$(5.26) \quad \phi(t) + (1 - \psi(t)) = (\phi(0) + (1 - \psi(0))) \exp\left(\int_0^t p(s) ds\right),$$

and the outer phase fraction exists as long as $\int_0^t p(s) ds > -\infty$. Comparing (5.25) and (5.26) shows also that periodicity imposes the following requirement on $p(t)$ and $q(t)$ in terms of their initial phase fractions:

$$(5.27) \quad (\phi(0) + (1 - \psi(0))) \exp\left(\int_0^t p(s) ds\right) + (\psi(0) - \phi(0)) \exp\left(\int_0^t q(s) ds\right) = 1.$$

Using (5.16) to compute $p - q$ next gives

$$(5.28) \quad p(t) - q(t) = (p(0) - q(0)) \exp\left(\frac{1}{n - 1} \int_0^t p(s) + q(s) ds\right),$$

which shows $p(t) - q(t)$ cannot change sign. By (5.21), $p(t)$ and $q(t)$ consequently keep their signs as long as neither phase fraction collapses. This can alternatively be seen by considering a sketch of u and observing that p and q have opposite signs and can vanish only simultaneously. Thus, without loss of generality, we subsequently assume $p(t) > 0 > q(t)$, for at least some $t \geq 0$.

In the following lemma, we solve the system (5.23), (5.24) implicitly in order to show that two-phase solutions of the type (5.13) exist for all time.

LEMMA 5.1. *Let $(p(t), q(t))$ satisfy (5.23), (5.24) with initial data $p(0) > 0 > q(0)$ (respectively, $p(0) < 0 < q(0)$). Then $(p(t), q(t))$ exists for all $t \in (-\infty, \infty)$ and satisfies $p(t) > 0 > q(t)$ (respectively, $p(t) < 0 < q(t)$). Further, $(p(t), q(t)) \rightarrow (0, 0)$ as $t \rightarrow \pm\infty$, and $\|u(\cdot, t)\|_{C^1} + \|\partial_t u(\cdot, t)\|_{C^0} \rightarrow 0$ as $t \rightarrow \pm\infty$.*

Proof. Writing (5.23) and (5.24) using polar variables $p(t) = r(t) \cos \theta(t)$, $q(t) = r(t) \sin \theta(t)$, $r(t) \geq 0$, leads to the following:

$$(5.29) \quad \dot{r}(t) = \frac{r^2}{n-1}(\cos^3(\theta(t)) + n(\cos \theta(t) + \sin \theta(t)) \cos \theta(t) \sin \theta(t) + \sin^3(\theta(t))),$$

and

$$(5.30) \quad \dot{\theta}(t) = r(\cos \theta(t) - \sin \theta(t)) \cos \theta(t) \sin \theta(t).$$

In the case $p(0) > 0 > q(0)$, $\theta(0) \in (-\pi/2, 0)$, so by (5.30) $\dot{\theta}(0) < 0$ and, as long as $\theta(t) \in (-\pi/2, 0)$, $\dot{\theta}(t) < 0$.

We will show that $\theta(t) \rightarrow -\pi/2$ as $t \rightarrow \infty$ ($\theta(t) \rightarrow 0$ as $t \rightarrow -\infty$) by using (5.29) and (5.30). Integrating the resulting expression for $\frac{d \ln r}{d\theta}$ gives

$$(5.31) \quad r(\theta) = c |\cos \theta \sin \theta|^{\frac{1}{n-1}} |\cos \theta - \sin \theta|^{-\frac{n+1}{n-1}},$$

where $c > 0$ denotes a generic constant. Inserting this expression for $r(\theta)$ in (5.30) results in

$$(5.32) \quad \begin{aligned} \dot{\theta}(t) &= c |\cos \theta \sin \theta|^{\frac{1}{n-1}} |\cos \theta - \sin \theta|^{-\frac{n+1}{n-1}} (\cos \theta - \sin \theta) \cos \theta \sin \theta \\ &= -c |\cos \theta \sin \theta|^{\frac{n}{n-1}} |\cos \theta - \sin \theta|^{-\frac{2}{n-1}} \end{aligned}$$

for $\theta(t) \in (-\pi/2, 0)$, and so

$$(5.33) \quad \int_{\theta(0)}^{\theta(t)} \frac{|\cos \theta - \sin \theta|^{\frac{2}{n-1}}}{|\cos \theta \sin \theta|^{\frac{n}{n-1}}} d\theta = -ct, \quad -\frac{\pi}{2} < \theta(0) < 0, \quad c > 0.$$

Since $\frac{n}{n-1} > 1$, the integral expression diverges, both as $\theta(t) \rightarrow -\pi/2$ ($t \rightarrow +\infty$) and as $\theta(t) \rightarrow 0$ ($t \rightarrow -\infty$). Thus $\theta(t)$ and, from (5.31), $r(t)$, are bounded, continuous functions of time, with $\theta(t) \in (-\pi/2, 0)$. By (5.31) then, $r(t) \rightarrow 0$ as $t \rightarrow \pm\infty$, and, noting that from (5.14), (5.15),

$$\alpha(t) = -\frac{p(t) - q(t)}{2}(\psi(t) - \phi(t))(\phi(0) + \psi(0) - 1),$$

it follows from (5.35) and (5.36) that

$$p(t), q(t), \alpha(t), \text{ and } u(x, t) \rightarrow 0 \text{ as } t \rightarrow \pm\infty.$$

The remaining conclusions are easily obtained. □

Finally, we show that for $c = 0$ the phases remain in the interval $[0, 1]$ for all time.

THEOREM 5.1. *Suppose $c = 0$ and $0 < \phi(0) < \psi(0) < 1$. Then the phases ϕ, ψ stay in $[0, 1]$. In particular*

$$\phi(t) \in [\phi(0), \frac{1}{2}(\phi(0) + \psi(0))] \text{ and } \psi(t) \in (\frac{1}{2}(\phi(0) + \psi(0)), \psi(0)]$$

for all $t > 0$. Further, the time asymptotic behavior satisfies

$$\lim_{t \rightarrow \infty} \phi(t) = \frac{1}{2}(\phi(0) + \psi(0)) = \lim_{t \rightarrow \infty} \psi(t).$$

Proof. With (5.15) giving α (for $c = 0$), substituting into (5.18) and (5.20) shows that, by (5.26),

(5.34)

$$\dot{\phi}(t) = (1 - \psi(t))\dot{\psi} = \frac{p}{2}(\phi + (1 - \psi)) = \frac{1}{2}(\phi(0) + (1 - \psi(0)))p(t) \exp\left(\int_0^t p(s)ds\right),$$

and the individual phases therefore satisfy

$$(5.35) \quad \phi(t) = \phi(0) + \frac{1}{2}(\phi(0) + (1 - \psi(0))) \left(\exp\left(\int_0^t p(s)ds\right) - 1 \right)$$

and

$$(5.36) \quad \psi(t) = \psi(0) - \frac{1}{2}(\phi(0) + (1 - \psi(0))) \left(\exp\left(\int_0^t p(s)ds\right) - 1 \right).$$

Now we examine (5.27). Assuming $q < 0 < p$, the first term is monotonically increasing in time and must converge, for $0 < \phi(0) < \psi(0) < 1$, to a positive limit. The second term is positive but monotonically decreasing, so it may converge, as $t \rightarrow \infty$, either to a positive limit or to zero. Formally, setting $0 < \int_0^\infty p(t)dt = L < \infty$ and $-\infty \leq \int_0^\infty q(t)dt = M < 0$, (5.27) and (5.28) imply

$$(5.37) \quad (\phi(0) + (1 - \psi(0)))e^L + (\psi(0) - \phi(0))e^M = 1$$

and

$$(5.38) \quad \lim_{t \rightarrow \infty} (p(t) - q(t)) = (p(0) - q(0))e^{\frac{L+M}{n+1}}.$$

If M is finite and $p(0) \neq q(0)$, then (5.38), together with the necessity for $p(t)$ to approach zero as $t \rightarrow \infty$, means that $\lim_{t \rightarrow \infty} q(t) \neq 0$. However, this implies that $M = -\infty$, a contradiction. On the other hand, if $M = -\infty$, then, since $L < \infty$, (5.38) requires that $\lim_{t \rightarrow \infty} q(t) = 0$, which is permitted. We conclude that, for smooth solutions $p(t), q(t)$ to exist with $p(0) > 0 > q(0)$, we require that $L = \int_0^\infty p(t)dt < \infty$ and $M = \int_0^\infty q(t)dt = -\infty$. Thus, by (5.37),

$$(5.39) \quad e^L = \frac{1}{\phi(0) + 1 - \psi(0)}.$$

Writing (5.35) in the form

$$(5.40) \quad \phi(t) = \frac{1}{2}(\phi(0) - (1 - \psi(0))) + \frac{1}{2}(\phi(0) + (1 - \psi(0))) \exp \int_0^t p(s)ds,$$

it follows that

$$(5.41) \quad \lim_{t \rightarrow \infty} \phi(t) = \frac{1}{2}(\phi(0) + \psi(0))$$

and, similarly,

$$(5.42) \quad \lim_{t \rightarrow \infty} \psi(t) = \frac{1}{2}(\phi(0) + \psi(0)).$$

By the monotonicity in time of $\int_0^t p(s)ds$, therefore $\phi(t) \in [\phi(0), \frac{1}{2}(\phi(0) + \psi(0))]$ and $\psi(t) \in (\frac{1}{2}(\phi(0) + \psi(0)), \psi(0)]$ for all $t > 0$.

Acknowledgment. The first author would like to thank Barbara Keyfitz, colleagues, and staff for their support at the Fields Institute, Toronto, where part of this work was done.

REFERENCES

- [1] V. ARNOLD, *Sur la géométrie différentielle des groupes de Lie de dimension infinie et ses applications à l'hydrodynamique des fluides parfaits*, Ann. Inst. Fourier (Grenoble), 16 (1966), pp. 319–361.
- [2] F. CALOGERO, *A solvable nonlinear wave equation*, Stud. Appl. Math., 70 (1984), pp. 189–199.
- [3] S. CHILDRRESS, G. R. IERLEY, E. A. SPIEGEL, AND W. R. YOUNG, *Blow-up of unsteady two-dimensional Euler and Navier-Stokes solutions having stagnation-point form*, J. Fluid Mech., 203 (1989), pp. 1–22.
- [4] P. CONSTANTIN, *The Euler equations and nonlocal conservative Riccati equations*, Int. Math. Res. Not., 9 (2000), pp. 455–465.
- [5] S. M. COX, *Two-dimensional flow of a viscous fluid in a channel with porous walls*, J. Fluid Mech., 227 (1991), pp. 1–33.
- [6] C. M. DAFERMOS, *Hyperbolic Conservation Laws in Continuum Physics*, Grundlehren Math. Wiss. 325, Springer-Verlag, Berlin, 2000.
- [7] D. G. EBIN, *A concise presentation of the Euler equations of hydrodynamics*, Comm. Partial Differential Equations, 9 (1984), pp. 539–559.
- [8] D. G. EBIN AND J. MARSDEN, *Groups of diffeomorphisms and the motion of an incompressible fluid*, Ann. of Math. (2), 92 (1970), pp. 102–163.
- [9] J. EELLS, *A setting for global analysis*, Bull. Amer. Math. Soc., 72 (1966), pp. 751–807.
- [10] A. A. HIMONAS AND G. MISIOLEK, *The Cauchy problem for an integrable shallow-water equation*, Differential Integral Equations, 14 (2001), pp. 821–831.
- [11] A. A. HIMONAS AND G. MISIOLEK, *High frequency smooth solutions and well-posedness of the Camassa-Holm equation*, Int. Math. Res. Not., 51 (2005), pp. 3135–3151.
- [12] G. MISIOLEK, *Classical solutions of the periodic Camassa-Holm equation*, Geom. Funct. Anal., 12 (2002), pp. 1080–1104.
- [13] C. C. LIN, *Note on a class of exact solutions in magnetohydrodynamics*, Arch. Rational Mech. Anal., 1 (1958), pp. 391–395.
- [14] H. OKAMOTO AND J. ZHU, *Some similarity solutions of the Navier-Stokes equations and related topics*, Taiwanese J. Math., 4 (2000), pp. 65–103.
- [15] J. T. STUART, *Nonlinear Euler partial differential equations: Singularities in their solution*, in Applied Mathematics, Fluid Mechanics, Astrophysics, World Scientific, Singapore, 1988, pp. 81–95.
- [16] H. WEYL, *On the differential equations of the simplest boundary-layer problems*, Ann. of Math. (2), 43 (1942), pp. 381–407.
- [17] W. WOLIBNER, *Un théorème sur l'existence du mouvement plan d'un fluide parfait, homogène, incompressible, pendant un temps infiniment long*, Math. Z., 37 (1933), pp. 698–726.

BREAKING OF SYMMETRICAL PERIODIC SOLUTIONS IN A SINGULARLY PERTURBED KDV MODEL*

ALEXANDER TOVBIS†

Abstract. There are several recent developments in the well-known problem of breaking of homoclinic orbits (splitting of separatrices) of a system that undergoes a singular perturbation. First, survival of a homoclinic orbit is an exceptional situation that can be linked to triviality of the Stokes phenomenon of the underlying “truncated” equation. Second, homoclinic connections to exponentially small periodic orbits survive the perturbation in the generic case. In this paper we consider a different problem: we study deformations of “genuine” periodic orbits of the second order equation $y'' = y + y^2$ that undergoes the singular perturbation $\varepsilon^2 y'''' + (1 - \varepsilon^2)y'' = y + y^2$, where $\varepsilon > 0$ is a small parameter. We prove that if the period and the constant of motion do not change too rapidly (in ε), a genuine (nontrivial) periodic solution does not survive the perturbation.

Key words. singular perturbations, periodic solutions, exponentially small phenomena

AMS subject classifications. 34E, 34M, 34D15, 37J

DOI. 10.1137/070694053

1. Introduction.

1.1. Breaking of homoclinic connections. Let $y = 0$ be a hyperbolic stationary point of an n -dimensional differential equation $y' = f(y)$, and let W_s, W_u be the stable and unstable manifolds at $y = 0$. It is said that the stationary point $y = 0$ has a homoclinic connection if there exists a phase trajectory, originating and ending at $y = 0$, that lies on $W_s \cap W_u$. For example, the (bounded) separatrix solution

$$(1) \quad y(x) = -\frac{3/2}{\cosh^2(x/2)}$$

to

$$(2) \quad y''(x) = y(x) + y^2(x)$$

corresponds to the homoclinic connection of the stationary point $(0, 0)$ in the phase plane of (2). There exists a general problem to describe how a singular perturbation of the original equation affects the homoclinic (or a heteroclinic) connection. A large number of particular singularly perturbed equations, originating from the corresponding physical or computational problems, have been discussed in the literature (see, for example, [STL] and the references therein). In particular, the singular perturbation

$$(3) \quad \varepsilon^2 y''''(x) + (1 - \varepsilon^2)y''(x) = y(x) + y^2(x)$$

of (2) is related to the traveling wave reduction of a fifth order KdV equation that models gravity water waves. Existence or nonexistence of homoclinic connections of (3) has been studied in a number of papers; see, for example, [HM], [PRG], [GJ]. A rigorous proof of nonexistence of such connections was obtained in [AM] (see also

*Received by the editors June 8, 2007; accepted for publication (in revised form) June 4, 2008; published electronically November 19, 2008. This work was supported by NSF grant DMS 0508779. <http://www.siam.org/journals/sima/40-4/69405.html>

†Department of Mathematics, University of Central Florida, Orlando, FL 32816 (atovbis@pegasus.cc.ucf.edu).

[Eck]) and later on in [To4], where the “asymptotic beyond all orders” approach, first suggested in [KS] for a simple crystal growth model, was put on a rigorous basis. This approach is, in a sense, a natural way to study exponentially small phenomena, such as the mismatch between the stable and the unstable manifolds of (3). The goal of the present paper is to show that a similar argument can be used to prove nonexistence of symmetrical periodic solutions to (3) (see below), subject to some additional requirements.

The approach of [KS] is to rescale (3), which is called the outer equation, to the inner equation

$$(4) \quad v''''(z) + (1 - \varepsilon^2)v''(z) = \varepsilon^2v(z) + v^2(z),$$

where $x = \varepsilon z$ and $v = \varepsilon^2y$. The main advantage of (4) versus (3) is that the (exponentially small) difference between the stable and unstable solutions is detectable even in the leading order part

$$(5) \quad v''''(z) + v''(z) = v^2(z),$$

of (4), which is called the truncated equation (4). This means that the difference between the stable and unstable solutions to (3) can be studied through the Stokes phenomenon of (5). In fact, in [TP] we considered the family of singular perturbations

$$(6) \quad \varepsilon^2y^{(iv)} + (1 - \varepsilon^2)y'' - y = y^2 + \varepsilon^2\gamma(2yy'' + (y')^2)$$

of (2), parametrized by $\gamma \in \mathbb{R}$, and found exact conditions for persistence of the homoclinic connection (under the singular perturbation) in terms of the Stokes constant for the corresponding truncated inner equation. The technique developed in [To4] can be modified to consider singularly perturbed problems in other settings. For example, it was used in [To5] to show that the discretized equation (3) does not have a homoclinic connection.

Notice that persistence of the homoclinic connection under singular perturbations (3), (6) is equivalent to existence of a symmetrical (even) stable solution to the corresponding equation. Since the family of stable solutions to (3) (or to (6)) is one-dimensional (translation), we can arrange for $y'(0) = 0$; then the symmetry of the stable solution $y(x)$ is equivalent to $y'''(0) = 0$.

1.2. Deformation of periodic solutions under singular perturbations.

The problem of deformations of periodic solutions under singular perturbations has not yet received as much attention as the problem of homoclinic connections. However, it is known that there exist symmetrical periodic solutions to (3) that are exponentially small in ε ; see [Lo] and later papers [IL1], [IL2]. We consider such periodic solutions as deformations of a constant solution $y \equiv 0$ of the unperturbed equation (2). Existence of even periodic solutions that are deformations of a trivial solution was also proved in [BPBA] for a certain class of reversible fourth order Hamiltonian systems with zero Hamiltonian. Discussion of even periodic solutions for some fourth order ODEs with nonzero Hamiltonian and cubic nonlinearity ((3) has a quadratic nonlinearity) can be found in [PT].

In the present paper we extended the technique of [To4] to study deformations of a genuine (nonconstant) periodic solution of (2). We limit our attention to periodic solutions inside the potential well of (2). Any such solution is given by the Weierstrass elliptic \wp -function with invariants g_2, g_3 ,

$$(7) \quad y(x) = 6\wp_{g_2, g_3}(x - x_0) - \frac{1}{2},$$

and has an integral of motion

$$(8) \quad y'^2 = y^2 + \frac{2}{3}y^3 + C,$$

where $x_0 \in \mathbb{C}$ is a translation and $g_2 = \frac{1}{12}$, $g_3 = -\frac{1}{36}(C + \frac{1}{6})$ (see section 1.3). The requirement $C \in (-\frac{1}{3}, 0)$ on the constant of motion (energy) C guarantees that solution $y(x)$ is inside the potential well. (The value $C = 0$ corresponds to the trivial solution and to the separatrix solution of (2), whereas $C = -\frac{1}{3}$ corresponds to the stationary solution $y(x) \equiv -1$.) Solution (7) is a periodic function with finite basic real and purely imaginary periods

$$(9) \quad 2\omega_1 = \int_{e_1}^{\infty} \frac{du}{\sqrt{4u^3 - \frac{1}{12}u - g_3}}, \quad 2\omega_3 = i \int_{-\infty}^{e_3} \frac{du}{\sqrt{4u^3 - \frac{1}{12}u - g_3}},$$

where e_1, e_3 are the largest and the smallest real roots of $4u^3 - \frac{1}{12}u - g_3 = 0$ (see below), respectively. It is also an even function provided that x_0 is an integer linear combination of ω_1, ω_3 . It will be convenient for us to choose

$$(10) \quad x_0 = \omega_3,$$

so that the unperturbed solution (7) is bounded, real-valued, symmetrical, and periodic with the period $2\omega_1$ on $x \in \mathbb{R}$. Equivalently, we can consider solution (7) with $x_0 = 0$ along the horizontal line $\omega_3 + \mathbb{R}$ in the complex x -plane.

The perturbed equation (3) also has an integral of motion

$$(11) \quad \varepsilon^2(2y'''y' - y''^2) + (1 - \varepsilon^2)y'^2 - y^2 - \frac{2}{3}y^3 = C(\varepsilon),$$

where $C(\varepsilon)$ is the constant of motion (energy). Since we are interested in deformations of periodic solutions satisfying (8), we require $C(0) = C \in (-\frac{1}{3}, 0)$.

Let $\alpha \geq 1$. We say that a solution $y(x, \varepsilon)$ to the perturbed equation (3) is a C^α -deformation of a periodic solution (7) (in the potential well) to the unperturbed equation (2) under the singular perturbation (3) on interval S if

$$(12) \quad Y(x, \varepsilon) = Y(x, 0) + \varepsilon^\alpha \tilde{Y}(x, \varepsilon),$$

where vector $\tilde{Y}(x, \varepsilon)$ is continuous in ε , uniformly on interval S of the x -axis. Here $Y(x, \varepsilon) = \text{Col} (y(x, \varepsilon), y'(x, \varepsilon), y''(x, \varepsilon), y'''(x, \varepsilon))$. Note that a C^α -deformation of a periodic solution (7) has constant of motion (11) satisfying

$$(13) \quad C(\varepsilon) = C + \varepsilon^\alpha \tilde{C}(\varepsilon),$$

where $\tilde{C}(\varepsilon)$ is a continuous function.

Our main result is the following theorem.

THEOREM 1.1. *Let $y(x, 0) = 6\wp_{\frac{1}{12}, g_3}(x - \omega_3) - \frac{1}{2}$, where $|g_3| < 6^{-3}$, be a periodic solution (inside the potential well) of the nonperturbed equation (2). Let $n \in \mathbb{N}$ and $y(x, \varepsilon)$ be a C^α -deformation, $\alpha \geq 1$, of $y(x, 0)$ under the perturbation (3) on some open interval S that contains the segment $[\omega_1, n\omega_1]$. Then the deformation $y(x, \varepsilon)$ does not contain a sequence of symmetric periodic solutions $y(x, \varepsilon_m)$ of (3) with periods $2k\omega_1(\varepsilon_m)$, where $k = 1, 2, \dots, n$, $\lim_{m \rightarrow \infty} \varepsilon_m = 0$, and $\omega_1(\varepsilon)$ is subject to*

$$(14) \quad \omega_1(\varepsilon) - \omega_1 = \varepsilon^{\frac{\alpha}{2}} \tau(\varepsilon),$$

where $\tau(\varepsilon)$ is a continuous function for small $\varepsilon \geq 0$.

Similarly to the homoclinic connection problem, the central idea of our proof here is uniform (in small $\varepsilon \geq 0$) control of an iterative solution of inner equation (4). In a certain sense, though, the homoclinic connection problem is simpler, because the manifolds of stable and unstable solutions to (3) are one-dimensional, whereas deformations of a periodic solution to (2) under the perturbation (3) form a manifold of the full dimension. We rescale the outer equation (3) to the inner equation (4) by $x = \varepsilon z$, $v = \varepsilon^2 y$ and approach the problem through the following sequence of steps: (a) we construct by iterations a two-parameter family \mathcal{F}_α of deformations of the rescaled periodic solution (7), (10), satisfying (4); (b) we show that solutions $v(z, \varepsilon)$ to the inner equation (4) that corresponds to a C^α -deformation $y(x, \varepsilon)$ of a periodic solution (7) with the period $2\omega_1$ belong to the family \mathcal{F}_α ; and (c) for any $k \in \mathbb{N}$ we prove that the family \mathcal{F}_α does not contain a sequence $\{v(z, \varepsilon_m)\}_1^\infty$ of symmetrical and periodic (with the period $2k \frac{\omega_1(\varepsilon_m)}{\varepsilon_m}$) solutions if $\omega_1(\varepsilon)$ is subject to (14).

The paper is organized in the following way. In the remaining part of this section we describe the Stokes phenomenon for the truncated inner equation (5) and some connections between deformations of homoclinic and periodic orbits of (3). Solution of (3) by iterations, which yields the two-parameter family of solutions \mathcal{F}_α , is obtained in sections 2–3. Finally, the proof of Theorem 1.1 is completed in section 4.

1.3. Periodic solutions of the unperturbed equation. It is easy to check that the values C from the interval $(-\frac{1}{3}, 0)$ define periodic solutions of the unperturbed equation (2) within the well of the potential $-y^2 - \frac{2}{3}y^3$. After the change of variables $y = 6X - \frac{1}{2}$, (8) is reduced to

$$(15) \quad X'^2 = 4X^3 - \frac{1}{12}X - g_3,$$

where $g_3 = -\frac{1}{36}(C + \frac{1}{6})$. The solution to this equation is given by the Weierstrass elliptic function $\wp(x) = \wp_{g_2, g_3}(x)$, where the invariant $g_2 = \frac{1}{12}$ and g_3 is defined above. Thus, periodic solutions to the unperturbed equation (2) are given by

$$(16) \quad y(x) = 6\wp_{g_2, g_3}(x) - \frac{1}{2}.$$

The number $\Delta = g_2^3 - 27g_3^2$ is called the discriminant of $\wp_{g_2, g_3}(x)$. The condition $\Delta > 0$ is equivalent to $|g_3| < 6^{-3}$, i.e., to $C \in (-\frac{1}{3}, 0)$. Under this condition the Weierstrass function \wp has real period $2\omega_1$ defined by (9) and purely imaginary period $2\omega_3$, where $\omega_1, -i\omega_3$ are positive numbers (the latter requires correct choice of the branch of the square root in (9)). Using the standard notation, we denote $\omega_2 = \omega_1 + \omega_3$ and $e_j = \wp(\omega_j)$, where $j = 1, 2, 3$. It is well known that (see, for example, [GR]) e_j are the roots of the cubic polynomial in (15),

$$(17) \quad \wp'(\omega_j) = 0, \quad j = 1, 2, 3, \quad e_1 + e_2 + e_3 = 0, \quad \text{and} \quad e_3 \leq e_2 \leq e_1.$$

When x varies along the real axis, the value of \wp varies between e_1 and $+\infty$ (as $\wp(x)$ is an even function that has a second order pole at the origin with the principle part $\frac{1}{x^2}$); see Figure 1.

This corresponds to the periodic solution of (2), defined by (16), with the range from \tilde{e}_1 to infinity outside the potential well, i.e., to the unbounded branch of the periodic orbit (the energy of the motion is fixed). Here $\tilde{e}_j = 6e_j - \frac{1}{2}$, $j = 1, 2, 3$; see Figure 2.

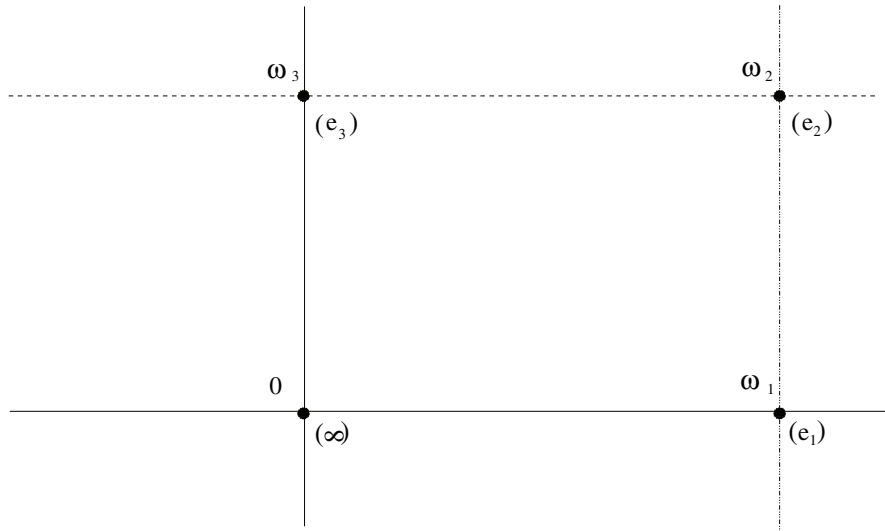


FIG. 1. Part of the parallelogram of periods $0\omega_1\omega_2\omega_3$ for the Weierstrass functions $\wp(x)$ in the complex x -plane together with the corresponding values of \wp shown in the brackets.

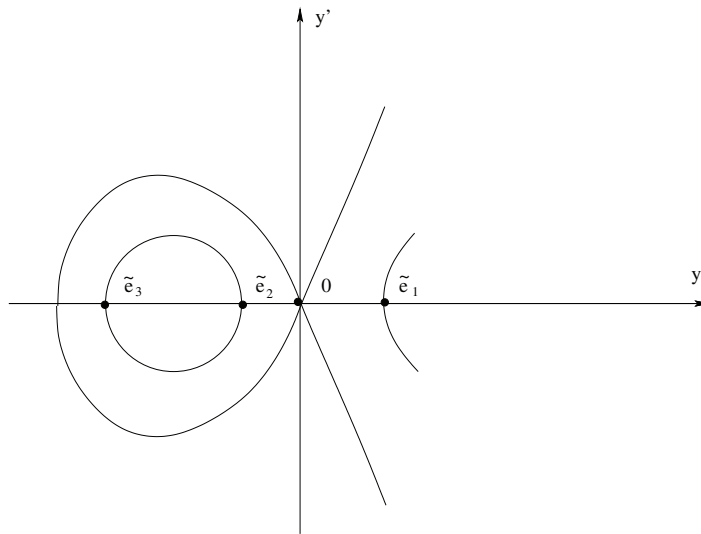


FIG. 2. Phase portrait of (2) showing periodic solutions ($\Delta > 0$) of the same period inside (bounded) and outside (unbounded) the potential well. Here $\tilde{e}_j = 6e_j - \frac{1}{2}$, $j = 1, 2, 3$.

When x varies along the horizontal line $\Im x = \omega_3$, the value of \wp varies between e_3 and e_2 . That corresponds to the periodic solution of (2), defined by (16), with the range from \tilde{e}_3 to \tilde{e}_2 inside the potential well, i.e., to the bounded branch of the periodic orbit. When x varies along the vertical line $\Re x = \omega_1$, the value of \wp varies between e_1 and e_2 . That corresponds to the tunneling between the unbounded and bounded branches of the periodic orbit defined by the solution (16); see Figures 1 and 2. For $\Delta < 0$ the periods ω_1 and ω_3 are complex conjugated so that ω_2 is real. In this case the value of \wp varies between e_2 and $+\infty$, which corresponds to infinite periodic motion outside the potential well; see Figure 3.

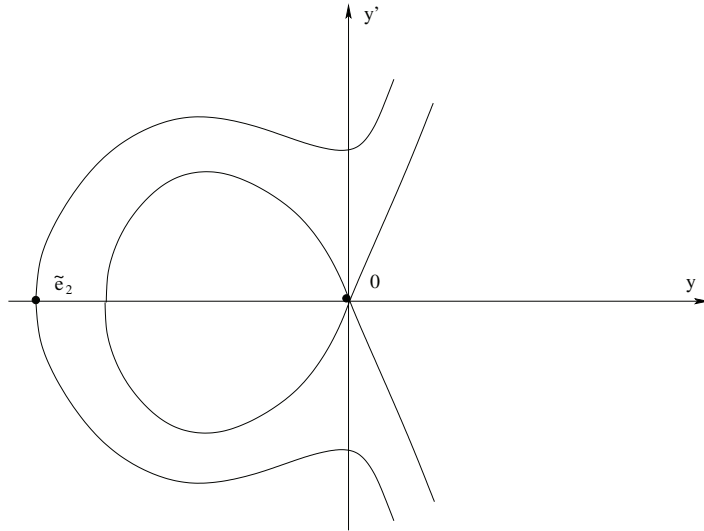


FIG. 3. Phase portrait of (2) showing an unbounded periodic solution ($\Delta < 0$) outside the potential well. Here $\tilde{e}_2 = 6e_2 - \frac{1}{2}$.

Here and henceforth it will be convenient for us to consider x_0 from (7) to be zero. Then we are interested in the behavior of $y(x)$ along the horizontal line $\omega_3(\varepsilon) + \mathbb{R}$.

1.4. Solutions to the perturbed equation. Considering perturbations of the separatrix solution (1) of the original equation, we focus our attention on the stable and unstable solutions of the perturbed equation (3) (those are solutions corresponding to the one-dimensional manifolds W_s and W_u in the phase space of (3)). We define the unstable and stable solutions to (3) as

$$(18) \quad y_{u,s}(x, \varepsilon) = \sum_{k=1}^{\infty} y_k(\varepsilon) e^{\pm kx},$$

respectively, where $y_1(\varepsilon)$ as an arbitrary continuous positive function and all $y_k, k > 1$, are uniquely defined through y_1 [To4]. The series (18) is convergent in properly chosen left- and right-half planes of the complex x -plane, respectively. It is easy to check that for $\varepsilon = 0$ the choice of $y_1(0) = 6$ yields the unbounded separatrix solution

$$(19) \quad y(x) = \frac{3/2}{\sinh^2(x/2)},$$

which is related to the bounded separatrix solution (1) through the shift $x \mapsto x + i\pi$. Considering x as a complex variable, now and henceforth we refer to (19) as the separatrix solution of (3).

Perturbations of periodic solutions to the original equation (2) that we are interested in are not as clearly identifiable as solutions (18) in the separatrix case. We start with considering the separatrix solution of (3) as the limit of periodic solutions (16) as $C \rightarrow 0^-$. In this limit the real half-period $\omega_1 \rightarrow \infty$, and the imaginary half-period $\omega_3 \rightarrow i\pi$. Thus, the point $\omega_2 \rightarrow i\pi + \infty$ as $C \rightarrow 0^-$.

Note that the stable and unstable solutions (18) are symmetrical with respect to $x = 0$ since $y_s(x, \varepsilon) = y_u(-x, \varepsilon)$. Alternatively, we can consider $y_u(x)$ as the

symmetrical continuation of $y_s(x)$ through $x = \infty$ on the compactification of \mathbb{R} . The same observation is correct for the horizontal line $\Im x = \pi$ and $x = i\pi + \infty$. The analogue of the point $x = i\pi + \infty$ in the periodic case is $x = \omega_2$. Therefore, by analogy with the separatrix case, we will consider perturbations of the periodic solution (16) that are symmetrical with respect to the point $\omega_2(\varepsilon)$. Equivalently, that means that the first and third derivatives of $y(x)$ turn zero at $x_0 = \omega_2$. The remaining two of the total of four conditions to define $y(x, \varepsilon)$ are $y'(\omega_3(\varepsilon), \varepsilon) = 0$ and the integral of motion (11). Thus, $y(x, \varepsilon)$ satisfies third order differential equation (11) with the prescribed energy $C(\varepsilon)$ together with three boundary conditions

$$(20) \quad y'(\omega_3(\varepsilon), \varepsilon) = 0, \quad y'(\omega_2(\varepsilon), \varepsilon) = y'''(\omega_2(\varepsilon), \varepsilon) = 0,$$

where $C(0) = C$ and $2\omega_1(\varepsilon), 2\omega_3(\varepsilon)$ are the fundamental periods of the Weierstrass function \wp in the solution (16) of (8).

It is easy to see that if $f(x)$ is an even (symmetrical at $x = 0$) periodic function with a period $2\omega_1$, then $f(x)$ is symmetrical at any point $x = k\omega_1, k \in \mathbb{N}$. Since a solution $y(x, \varepsilon)$ to the BVP (11), (20) is symmetrical at $x = \omega_2(\varepsilon)$ on the horizontal line $\Im x = \Im \omega_3(\varepsilon)$, $y(x, \varepsilon)$ is an even periodic solution if and only if

$$(21) \quad y'''(\omega_3(\varepsilon), \varepsilon) = 0.$$

1.5. The Stokes phenomenon for the inner equation. Periodic solutions (16) to (2) have second order poles at $x = 0$. The asymptotic beyond all orders approach of [KS] suggests blowing up this singularity by rescaling

$$(22) \quad x = \varepsilon z, \quad y(x, \varepsilon) = \varepsilon^{-2}v(z, \varepsilon),$$

where z is a new independent complex variable. The rescaled equation (inner equation) is given by (4), which we can rewrite as

$$(23) \quad (D_z^2 + 1)(D_z^2 - \varepsilon^2)v = v^2.$$

Let $v_+(z, \varepsilon)$ and $v_-(z, \varepsilon)$ denote, respectively, the rescaled solutions to the BVP (11), (20) and to its “symmetrical” problem, where conditions at $\omega_2(\varepsilon)$ are replaced by the same conditions at $-\bar{\omega}_2(\varepsilon) = -\omega_1(\varepsilon) + \omega_3(\varepsilon)$, i.e., $y'(-\bar{\omega}_2(\varepsilon), \varepsilon) = y'''(-\bar{\omega}_2(\varepsilon), \varepsilon) = 0$. Our observation (proved below) is that, similarly to the case of the separatrix solution, the leading order (in ε) of the difference between $v_+(z, \varepsilon)$ and $v_-(z, \varepsilon)$ can be derived from the truncated equation

$$(24) \quad (D_z^2 + 1)D_z^2v = v^2.$$

Note that $z = \infty$ is the irregular singular point of (24) and that this equation has a unique formal power series solution

$$(25) \quad \hat{v}(z) = \sum_{k=1}^{\infty} \frac{v_k}{z^{2k}},$$

where $v_1 = 6$.

The inverse Laplace transform \mathcal{L}^{-1} converts the truncated inner equation (24) into the convolution equation

$$(26) \quad (p^4 + p^2)V(p) = V(p) * V(p),$$

where p is a dual variable (called Borel variable), $V(p) = [\mathcal{L}^{-1}v](p)$, and $F(p) * G(p) = \int_0^p F(p-\tau)G(\tau)d\tau$. It is well known that the asymptotic expansion (25) of $v(z)$ yields the corresponding asymptotic expansion in powers of p of $V(p)$ at $p = 0$.

THEOREM 1.2. *Equation (26) admits a unique nontrivial power series solution in odd powers of p . This solution defines a function $V(p)$ that is analytic at the whole p -plane except for two vertical rays: from $p = i$ upward and from $p = -i$ downward. The function $V(p)$ is of exponential order 1 along any nonvertical ray in this cut p -plane.*

This theorem follows from a more general statement (Main Theorem) of [To2].

COROLLARY 1.3. *Let $v_{\pm}(z)$ be defined by*

$$(27) \quad v_{\pm}(z) = \int_0^{\pm\infty} e^{-z^p}V(p)dp.$$

These functions are the only analytic solutions of the truncated inner equation (24) that satisfy

$$(28) \quad v_{\pm}(z) \sim \hat{v}(z) \quad \text{as } z \rightarrow \infty, \quad z \in S_{\pm},$$

where S_{\pm} are sectors

$$(29) \quad S_+ = \{z : |\arg z| < \pi\}, \quad S_- = \{z : |\arg z - \pi| < \pi\}.$$

Moreover,

$$(30) \quad v_+(z) - v_-(z) = -2\pi i s e^{iz}(1 + o(1)) \quad \text{as } z \rightarrow \infty, \quad 0 < \arg z < \pi,$$

where the constant s is determined through

$$(31) \quad s = \lim_{p \rightarrow -i} (p + i)V(p).$$

Proof. The Taylor expansion of $V(p)$ at $p = 0$ can be obtained by applying the inverse Laplace transform (Borel transform) to the formal series $\hat{v}(z)$. The function $V(p)$ is analytic on a Riemann surface that has possible branch points of the logarithmic type at $p = ik$, where $k \in \mathbb{Z} \setminus \{0\}$. The statement on the behavior of $V(p)$ at singularities $p = \pm i$ follows from Theorem 2.2 in [To2]. The uniqueness of solutions $v_{\pm}(z)$ satisfying (28) follows from Theorem 1.2 and properties of the Laplace transform. \square

DEFINITION 1.4. *The constant s in (31) is called the Stokes constant for the truncated inner equation (24).*

PROPOSITION 1.5. *The following three conditions are equivalent: (i) the formal power series solution (25) has a positive radius of convergence; (ii) $s = 0$; and (iii) $v_+(z) \equiv v_-(z)$.*

Proof. The fact that (i) implies (iii) is obvious. The inverse statement follows from the fact that $v(z) = v_+(z) \equiv v_-(z)$ implies that the function $v(z)$ is single-valued near infinity and has asymptotic expansion $\hat{v}(z)$ in the full neighborhood of infinity (see [Wa]). It is also clear that (i) implies (ii), since in this case $\mathcal{L}^{-1}\hat{v}(z)$ is an entire function. Suppose now that $s = 0$. Since $V(p)$ is real-analytic on \mathbb{R} , we obtain that $V(p) = o(\frac{1}{p-(\pm i)})$ as $p \rightarrow \pm i$. That means, according to Corollary 2.2 in [To2], that $v_+(z)$ coincides with $v_-(z)$ in both the lower and upper z half-planes. Thus (ii) implies (iii). \square

Divergence (zero radius of convergence) of the formal power series (25) was proved in [To2, p. 247]. According to Proposition 1.5, this implies that $v_+(z) \not\equiv v_-(z)$. (In fact, nonzero numerical values of s for equations, equivalent to (24), were computed in [PRG], [GJ] by means of formal Borel summation.) This implies, at least at the formal level, that solution $v_+(z, \varepsilon)$ to (23) does not coincide with solution $v_-(z, \varepsilon)$ to (23). To make this statement rigorous, it is sufficient to show that (a)

$$(32) \quad v_{\pm}(z, \varepsilon) \rightarrow v_{\pm}(z)$$

as $\varepsilon \rightarrow 0$ in some regions R_{\pm} uniformly in $z \in R_{\pm}$, and (b) intersection $R_+ \cap R_-$ contains some segment of a positive length.

2. Inner equation. In the following theorem (Theorem 2.1) we construct a two-parameter family of solutions to (23) that are symmetrical at $z = \tilde{\omega}_2(\varepsilon)$ and that converge to $v_+(z)$, or to its translation $v_+(z - h)$, $h \in \mathbb{R}$, as $\varepsilon \rightarrow 0$ uniformly in the corresponding region (see below) of the complex z -plane. Here $\tilde{\omega}_j(\varepsilon) = \frac{\omega_j(\varepsilon)}{\varepsilon}$, $j = 1, 2, 3$.

This result constitutes an essential part of the method of [KS]; convergence of $v_+(z, \varepsilon)$ to $v_+(z)$ was mentioned only briefly in the original paper [KS] (with regard to a third order equation considered there), but explicit formulations and proofs were omitted. In other papers, connection between solutions of (23) and (24) was considered only on the formal level. The proof of convergence of $v_+(z, \varepsilon)$ to $v_+(z)$ in the separatrix case was given in [To4] (see also [To5] for the discretized equation (2), as well as [TP] for a more general family of singular perturbations of (2)). The major difficulty there was based on the fact that solutions of the full and of the truncated inner equations have different rates of convergence to 0 as $z \rightarrow \infty$ (in proper directions): $v_+(z, \varepsilon)$, $\varepsilon > 0$, approaches 0 exponentially fast (in z), while $v_+(z)$ has only power order convergence. One needs to find, however, some uniform majorization that will cover both cases. In order to find such a majorization, we suggested [To4], [To5] linearizing (23) around the unperturbed solution $6q^2$, where $q(z, \varepsilon) = \frac{\varepsilon}{\sinh \frac{z}{2}}$, and applying the contraction mapping principle to the obtained equation. A similar approach will be developed below with the natural choice

$$(33) \quad q(z, \varepsilon) = \varepsilon \sqrt{\wp_{\frac{1}{12}, g_3(\varepsilon)}(\varepsilon z) - \frac{1}{12}},$$

where $g_3 = g_3(\varepsilon)$ is determined through $\omega_1 = \omega_1(\varepsilon)$ by (9).

2.1. Linearized equation. The substitution

$$(34) \quad v(z, \varepsilon) = u(z, \varepsilon) + 6c(\varepsilon)q^2(z, \varepsilon),$$

where $q(z, \varepsilon)$ is given by (33), reduces (23) to

$$(35) \quad (D_z^2 + 1)(D_z^2 - \varepsilon^2)u = u^2 + 12c(\varepsilon)q^2u + f(q).$$

Here $c = c(\varepsilon)$ is a continuous function satisfying $c(0) = 1$. Calculations of f after some algebra yield

$$(36) \quad f(q) = c [36(c - 1 - 4\varepsilon^2)q^4 - 6!q^6 - 2\varepsilon^6 J(\varepsilon)],$$

where

$$(37) \quad J(\varepsilon) = -36g_3(\varepsilon) - \frac{1}{6}$$

is determined through $\omega_1 = \omega_1(\varepsilon)$. (To simplify notation, we omit the ε dependence in ω_j or $\tilde{\omega}_j$, $j = 1, 2, 3$, provided that such omission cannot cause any misunderstanding.) Note that $\lim_{\varepsilon \rightarrow 0} 6c(\varepsilon)q^2(z, \varepsilon) = \frac{6}{z^2}$. This is the leading term of the asymptotic expansion (25) of $v_+(z)$. So, substitution (34) linearizes (23) with respect to the leading term for both positive and zero values of ε .

Let

$$(38) \quad w = (D^2 + 1)u.$$

Then (35) becomes

$$(39) \quad (D^2 - \varepsilon^2)w = [(D^2 + 1)^{-1}w]^2 + 12cq^2(D^2 + 1)^{-1}w + f(q)$$

or

$$(40) \quad [D^2 - (\varepsilon^2 + 12q^2)]w = [(D^2 + 1)^{-1}w]^2 + 12cq^2[(D^2 + 1)^{-1} - 1]w - 12q^2(1 - c)w + f(q).$$

Using the identity

$$(41) \quad (D^2 + 1)^{-1} - 1 = -(D^2 + 1)^{-1}D^2$$

and solving (39) for D^2w , we finally get

$$(42) \quad [D^2 - (\varepsilon^2 + 12q^2)]w = [1 - 12cq^2(D^2 + 1)^{-1}] (f(q) + [(D^2 + 1)^{-1}w]^2) - 12q^2(1 - c)w - 12cq^2(D^2 + 1)^{-1} [\varepsilon^2 + 12cq^2(D^2 + 1)^{-1}] w.$$

This integrodifferential equation is equivalent to (35), (38). We solve (42) by the contraction mapping principle, considering it as a ‘‘perturbation’’ of the nonhomogeneous linear ODE

$$(43) \quad [D^2 - (\varepsilon^2 + 12q^2)]w = [1 - 12cq^2(D^2 + 1)^{-1}]f(q).$$

The corresponding homogeneous equation

$$(44) \quad [D^2 - (\varepsilon^2 + 12q^2)]w = 0$$

has two independent solutions $v_1(z)$ and $v_2(z)$, where $v_2(z)$ is symmetric (even) and $v_1(z)$ is antisymmetric (odd) at the origin. Indeed, $v = 6q^2$ is a solution to $(D^2 - \varepsilon^2)v = v^2$, which is the rescaled equation (2). Differentiating both sides of the latter equation, we get $(D^2 - \varepsilon^2)v' = 2vv'$. This equation coincides with (44), where $w = v'$. Thus, we obtain

$$(45) \quad v_1(z) = 6 \frac{d}{dz} q^2 = 6\varepsilon^3 \wp'(\varepsilon z).$$

The second linearly independent solution can be taken as

$$(46) \quad v_2(z) = v_1(z) \int_0^z \frac{d\xi}{v_1^2(\xi)}.$$

Note that $\wp(\varepsilon z)$ is symmetrical and, thus, v_1 is an antisymmetrical (odd) function with respect to any integer combination of half periods $\tilde{\omega}_1$ and $\tilde{\omega}_3$. It will be shown below (section 2.3) that solutions

$$(47) \quad v_{2j}(z) = v_2(z) - \frac{1}{6\varepsilon^7\Delta} \left[\frac{3g_3}{2}\omega_j - g_2\zeta(\omega_j) \right] v_1(z), \quad j = 1, 2, 3,$$

to (44), where

$$(48) \quad \zeta(x) = \frac{1}{x} - \int_0^x \left(\wp(y) - \frac{1}{y^2} \right) dy$$

is a “negative antiderivative” of \wp , are symmetrical at the points $z = \tilde{\omega}_j$, respectively.

2.2. Formulation of the theorem. Let $g(z)$ be a function, symmetrical (even) with respect to the point $\tilde{\omega}_2(\varepsilon)$, and let $\mathcal{I}_2 = (D^2 + 1)^{-1}$ denote the operator, inverse to $D^2 + 1$ and such that $\mathcal{I}_2 g(z)$ is symmetrical with respect to $\tilde{\omega}_2(\varepsilon)$. The construction of \mathcal{I}_2 will be given in section 3.1.

Equation (42) can be written in the operator form as $w = \mathcal{N}w$, where

$$(49) \quad \mathcal{N}w = \mathcal{I}_1 \circ [(1 - 12cq^2\mathcal{I}_2)[f + (\mathcal{I}_2 w)^2] - 12q^2(1 - c)w + 12cq^2\mathcal{I}_2(\varepsilon^2 + 12cq^2\mathcal{I}_2)w]$$

and

$$(50) \quad \mathcal{I}_1 g(z) = -v_1(z) \int_{\tilde{\omega}_2}^z v_2(\xi)g(\xi)d\xi + v_2(z) \int_{\tilde{\omega}_2}^z v_1(\xi)g(\xi)d\xi.$$

We define iterations

$$(51) \quad w_0 = 0, \quad w_1 = \mathcal{N}w_0 + \varepsilon^6 b(\varepsilon)v_{22}(z), \quad w_n = \mathcal{N}w_{n-1}, \quad n = 2, 3, \dots,$$

and $\Delta w_n = w_n - w_{n-1}$, $n = 1, 2, \dots$, where $b(\varepsilon)$ is a continuous function in a vicinity of $\varepsilon = 0$. The following theorem proves uniform convergence of the series $\sum_{n=1}^\infty \Delta w_n$ in a region $R_{z_0} \subset \mathbb{C}$, which is defined below. Thus, we obtain a family of solutions to (42), parametrized by $c(\varepsilon)$ and $b(\varepsilon)$, which are symmetrical at $z = \tilde{\omega}_2$.

Consider the parallelogram of periods of the elliptic function $q(z, \varepsilon)$ (the periods are $2\tilde{\omega}_1$ and $2\tilde{\omega}_3$) centered at the origin with the square cut that has vertices $\pm z_0(1+i)$ and $\pm z_0(1-i)$, where $z_0 > 0$ is a positive constant. We require $\varepsilon \in [0, \varepsilon_0]$, where

$$(52) \quad \varepsilon_0 z_0 < \min\{\omega_1, |\omega_3|\}.$$

By R_{z_0} we denote a quarter of this figure, i.e., of the parallelogram of periods with the square cut, that lies in the first quadrant; see Figure 4.

THEOREM 2.1. *Let E be a closed subinterval of $(-\frac{1}{3}, 0)$, and let $J(\varepsilon)$ from (37) be a function, defined on a neighborhood of $\varepsilon = 0$ with the range in E . Let continuous functions $b(\varepsilon), c(\varepsilon)$ satisfy*

$$(53) \quad |b(\varepsilon)| \leq l\varepsilon^\alpha \quad \text{and} \quad |c(\varepsilon) - 1| \leq l\varepsilon^\alpha$$

for some $l > 0$, $\alpha \geq 1$, and for all nonnegative ε from a vicinity of $\varepsilon = 0$. Then there exist $\varepsilon_0 > 0$, $z_0 > 0$ satisfying (52), which depend only on E and l , and such that the series

$$(54) \quad w(z, \varepsilon) = \sum_{n=1}^\infty \Delta w_n(z, \varepsilon)$$

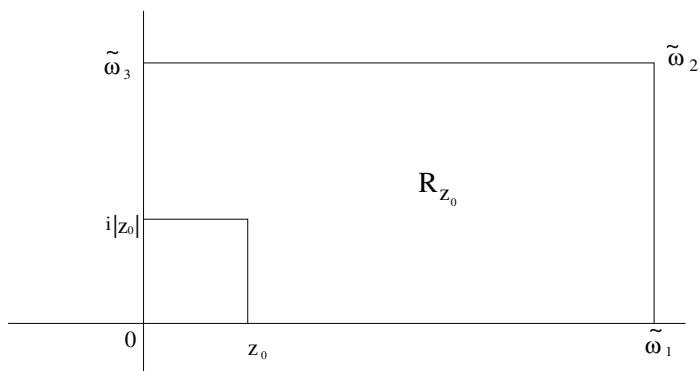


FIG. 4. Region R_{z_0} .

converges uniformly in $R_{z_0} \times [0, \varepsilon_0]$.

Remark 2.2. The rectangular “cut” in R_{z_0} can be replaced by a cut of any other shape so that the distance between R_{z_0} and the origin is $O(z_0)$. In fact, it will be replaced by a more convenient triangular “cut” in the course of the proof.

The proof of the theorem is given in section 3.4.

2.3. The second fundamental solution to the linearized equation. Note that, similarly to the separatrix case, the solution v_1 to (44) is an odd function, while the solution v_2 is an even function; see (45)–(46). Moreover, $\wp(\varepsilon z)$ is symmetrical, and, thus, v_1 is an antisymmetrical (odd) function with respect to any integer combination of half-periods $\tilde{\omega}_1$ and $\tilde{\omega}_3$.

Our computation of v_2 is based on the well-known fact (see, for example, [WW]) that the difference between two elliptic functions with the same periods and same principal parts at each singular point is a constant. The only zeros of the function $\wp'(\varepsilon z)$ within the parallelogram of periods are simple zeros at the points $\tilde{\omega}_j$, $j = 1, 2, 3$. Therefore

$$(55) \quad \frac{1}{\wp'^2(\varepsilon z)} = A_0 + \sum_{j=1}^3 A_j \wp(\varepsilon z - \omega_j)$$

for all z and

$$(56) \quad \int_0^z \frac{d\xi}{\wp'^2(\varepsilon \xi)} = A_0 z - \frac{1}{\varepsilon} \sum_{j=1}^3 A_j [\zeta(\varepsilon z - \omega_j) + \zeta(\omega_j)],$$

where the constants A_k , $k = 0, 1, 2, 3$, will be discussed below and $\zeta(x)$, defined by (48), is an odd function. The following arguments use some standard facts about the Weierstrass \wp -function that can be found, for example, in [WW], [Ha].

Using the “addition theorem” for ζ -function

$$(57) \quad \zeta(\varepsilon z - \omega_j) - \zeta(-\omega_j) = \zeta(\varepsilon z) + \frac{1}{2} \frac{\wp'(\varepsilon z)}{\wp(\varepsilon z) - e_j}, \quad j = 1, 2, 3,$$

we get

$$(58) \quad \int_0^z \frac{d\xi}{\wp'^2(\varepsilon \xi)} = \frac{1}{\varepsilon} \left[A_0 \varepsilon z - \zeta(\varepsilon z) \sum_{j=1}^3 A_j - \frac{1}{2} \wp'(\varepsilon z) \sum_{j=1}^3 \frac{A_j}{\wp(\varepsilon z) - e_j} \right].$$

Constants A_k , $k = 0, 1, 2, 3$, are calculated in Lemma 2.3 below. Differential equations

$$(59) \quad \wp'^2 = 4 \prod_{j=1}^3 (\wp - e_j) \quad \text{and} \quad \wp'' = 6\wp^2 - \frac{g_2}{2}$$

for the Weierstrass \wp -function and identities

$$(60) \quad e_1e_2 + e_1e_3 + e_2e_3 = -\frac{g_2}{4} \quad \text{and} \quad e_1e_2e_3 = \frac{g_3}{4}$$

are also utilized in this lemma.

LEMMA 2.3.

(61)

$$A_0 = \frac{9g_3}{\Delta}, \quad A_j = \left[\frac{1}{2(e_j - e_i)(e_j - e_k)} \right]^2, \quad j = 1, 2, 3, \quad \text{and} \quad A_+ = \sum_{j=1}^3 A_j = \frac{6g_2}{\Delta},$$

where the indices i, j, k are a permutation of $1, 2, 3$.

Proof. The Taylor expansion of \wp at ω_j , $j = 1, 2, 3$, yields

$$(62) \quad \wp(x) = e_j + a_j(x - \omega_j)^2 + O(x - \omega_j)^4.$$

Therefore the principal part of $(\wp')^{-2}$ at ω_j is $\frac{1}{4a_j^2(x - \omega_j)^2}$. According to (59), $a_j = \frac{1}{2}\wp''(\omega_j) = 3e_j^2 - \frac{g_2}{4}$. The principal part of $\wp(x)$ at the origin is $\frac{1}{x^2}$. Thus

$$(63) \quad A_j = \frac{1}{(6e_j^2 - \frac{g_2}{2})^2} = \left[\frac{1}{2(e_j - e_i)(e_j - e_k)} \right]^2,$$

where the indices i, j, k are a permutation of $1, 2, 3$. The latter expression follows from

$$(64) \quad \begin{aligned} 3e_j^2 - \frac{g_2}{4} &= 2e_j^2 + e_j(e_j + e_i) + e_je_k + e_ie_k = 2e_j^2 + e_ie_k = 2(e_i + e_k)^2 + e_ie_k \\ &= (2e_i + e_k)(2e_k + e_i) = (e_i - e_j)(e_k - e_j), \end{aligned}$$

where (17) and (60) were taken into account.

Direct computations, based on (17) and (60), show that

$$(65) \quad \sum_{j=1}^3 e_j^2 = \frac{g_2}{2}, \quad \sum_{i < k} e_i^2 e_k^2 = \frac{g_2^2}{16}, \quad \text{and} \quad \sum_{j=1}^3 e_j^3 = \frac{3}{4}g_3,$$

where summation in the second expression is taken over all possible pairs i, k such that $1 \leq i < k \leq 3$. Then

$$(66) \quad \sum_{j=1}^3 A_j = \frac{\sum_{i < k} (e_i - e_k)^2}{4 \prod_{i < k} (e_i - e_k)^2} = \frac{3g_2}{8 \prod_{i < k} (e_i - e_k)^2}.$$

Next we compute the denominator, using (60) and (63):

$$(67) \quad \begin{aligned} -8 \prod_{i < k} (e_i - e_k)^2 &= \prod_{j=1}^3 \left(6e_j^2 - \frac{g_2}{2} \right) = -\frac{g_2^3}{8} + 6\frac{g_2^2}{4} \sum_{j=1}^3 e_j^2 - 18g_2 \sum_{i < k} e_i^2 e_k^2 \\ &\quad + 6^3 \prod_{j=1}^3 e_j^2 = \frac{1}{2}(27g_3^2 - g_2^3) = -\frac{\Delta}{2}. \end{aligned}$$

Now the second statement of the lemma follows from (66).

In order to evaluate A_0 we put $z = 0$ at (55). Since \wp' has a pole at the origin, (55) yields

$$(68) \quad A_0 = - \sum_{j=1}^3 A_j e_j = \frac{\sum_{j=1}^3 e_j (e_i - e_k)^2}{-4 \prod_{i < k} (e_i - e_k)^2}.$$

Using (17) and (65), we calculate the numerator in the latter expression as $-\frac{9}{4}g_3$. Then the first statement of the lemma follows from (67)–(68). \square

Now, the combination of (45)–(46) and of (58)–(59) yields

$$(69) \quad v_2(z) = \frac{1}{6\varepsilon^4} \left[A_0 x - A_+ \zeta(x) \right] \wp'(x) - 2 \sum_{j=1}^3 A_j (\wp(x) - e_i)(\wp(x) - e_k),$$

where for convenience we use the notation $x = \varepsilon z$. The latter sum is a quadratic polynomial in \wp , which, according to (17) and (68), can be written as $-2A_+ \wp^2 + 2A_0 \wp + B$, where the constant $B = -2 \sum_{j=1}^3 A_j e_i e_k$. In order to determine B , note that (46) implies that v_2 is an even function, which has zero of order four at the origin. Using the Laurent expansions of \wp , \wp' , and ζ at the origin (see [GR]) and equating the constant term of (69) to be zero, we obtain

$$(70) \quad B = -2 \sum_{j=1}^3 A_j e_i e_k = A_+ \left(\frac{g_2}{10} + \frac{g_2}{30} \right) + 2A_+ \frac{g_2}{10} = \frac{2g_2^2}{\Delta}.$$

Combining the latter equation with (69) and Lemma 2.3 yields

$$(71) \quad v_2(z) = \frac{1}{\varepsilon^4 \Delta} \left(\left[\frac{3g_3}{2} \varepsilon z - g_2 \zeta(\varepsilon z) \right] \wp'(\varepsilon z) - 2g_2 \wp^2(\varepsilon z) + 3g_3 \wp(\varepsilon z) + \frac{g_2^2}{3} \right) \\ = \frac{1}{\varepsilon^4 \Delta} \left(\frac{3g_3}{2} [\varepsilon z \wp'(\varepsilon z) + 2\wp(\varepsilon z)] - \frac{g_2}{3} [3\zeta(\varepsilon z) \wp'(\varepsilon z) + 6\wp^2(\varepsilon z) - g_2] \right).$$

Note that $v_2(z)$ is an elliptic function if and only if $\frac{3g_3}{2}x - g_2\zeta(x)$ is doubly periodic with fundamental periods $2\omega_1, 2\omega_3$. Using $\zeta(x + 2\omega_j) = \zeta(x) + 2\zeta(\omega_j)$ for all x in the domain of ζ and $j = 1, 2, 3$ [GR], we can write the condition that v_2 is an elliptic function as

$$(72) \quad 3g_3\omega_j - 2g_2\zeta(\omega_j) = 0, \quad j = 1, 3.$$

Considering these equations as a linear system for unknowns $3g_3, -2g_2$, we come to the conclusion that it cannot have nontrivial solutions since its determinant $\omega_1\zeta(\omega_3) - \omega_3\zeta(\omega_1) \equiv -\frac{i\pi}{2}$ [GR]. Thus v_2 is not an elliptic function for any energy $C \in (-\frac{1}{3}, 0)$. Note that in the case of the separatrix solution v_2 is also not a hyperbolic function (see [To4]).

2.4. Estimates for $v_1(z), v_2(z)$. Let T denote trapezoid $ABCD$ in \mathbb{C} with the vertices $A(\tilde{\omega}_1 - \tilde{\omega}_3), B(\tilde{\omega}_1 + 2\tilde{\omega}_3), C(\frac{3}{2}\tilde{\omega}_3), D(-\frac{1}{2}\tilde{\omega}_3)$. We start with the estimate

$$(73) \quad |q^2(z, \varepsilon)| \leq \frac{Q}{|z|^2}$$

with some $Q > 0$, which is valid for any $\varepsilon \geq 0$ and for any z in the trapezoid T . The value of the constant Q can be chosen as $Q = 1 + |\omega_2|^2 \tilde{Q}$, where $\tilde{Q} = \frac{1}{12} + \max_{x \in \tilde{T}} |\wp(x) - \frac{1}{x^2}|$. Here \tilde{T} denotes the trapezoid in the x -plane, $x = \varepsilon z$, that corresponds to T . Note that ω_2 is defined by the energy $J(\varepsilon)$ of the periodic solution (33). However, we can choose Q such that (73) holds for all $J(\varepsilon) \in E$.

LEMMA 2.4. *There exist positive constants A and B such that the estimates*

$$(74) \quad |v_1(z)| \leq \frac{A}{|z|^3} \quad \text{and} \quad |v_2(z)| \leq B|z|^4$$

are valid in T for any $\varepsilon \geq 0$ and any $J(\varepsilon) \in E$.

Proof. The first estimate follows from (45) and the inequality

$$(75) \quad \left| \wp'(x) + \frac{2}{x^3} \right| \leq A^3$$

in \tilde{T} , similarly to (73).

The second inequality (74) follows from the fact that the expression in the square brackets in the right-hand side of (69) is bounded by $6|x|^4 B$ in \tilde{T} with some positive constant B that does not depend on ε and on $J(\varepsilon) \in E$. \square

3. Solution by iterations.

3.1. Operator $(D^2 + 1)^{-1}$. Let $\omega \in \mathbb{C}$, $a > 0$, and let an analytic function $g(z)$ be real-valued on $\Re z = \Re \omega$ and symmetrical with respect to ω . Note that under these assumptions $g(z)$ is also real-valued along $\Im z = \Im \omega$.

PROPOSITION 3.1. *The operator I_a^ω defined by*

$$(76) \quad I_a^\omega g(z) = \frac{1}{2i} \left[e^{iz} \int_{\omega-ia}^z e^{-i\xi} g(\xi) d\xi - e^{-iz} \int_{\omega+ia}^z e^{i\xi} g(\xi) d\xi \right]$$

is inverse to $D^2 + 1$, and $I_a^\omega g(z)$ is real-valued on $\Re z = \Re \omega$ and symmetrical with respect to ω .

Proof. The first statement of the proposition is obvious. In order to prove the second statement, we decompose

$$(77) \quad I_a^\omega g(z) = I^\omega g(z) + \frac{1}{2i} \left[e^{iz} \int_{\omega-ia}^\omega e^{-i\xi} g(\xi) d\xi - e^{-iz} \int_{\omega+ia}^\omega e^{i\xi} g(\xi) d\xi \right],$$

where $I^\omega = I_0^\omega$. Let $ib = \int_\omega^{\omega+ia} e^{i(\xi-\omega)} g(\xi) d\xi$. Since g is real-valued along the path of integration, $b \in \mathbb{R}$ and $-ib = \int_\omega^{\omega-ia} e^{-i(\xi-\omega)} g(\xi) d\xi$. Thus, the second term in (77) is symmetrical and real-valued as it is

$$(78) \quad \frac{1}{2i} [ibe^{i(z-\omega)} + ibe^{-i(z-\omega)}] = b \cos(z - \omega).$$

To complete the proof we observe that

$$(79) \quad \begin{aligned} I^\omega &= \frac{1}{2i} \left[e^{iz} \int_\omega^z e^{-i\xi} g(\xi) d\xi - e^{-iz} \int_\omega^z e^{i\xi} g(\xi) d\xi \right] = \int_\omega^z \sin(z - \xi) g(\xi) d\xi \\ &= \sin(z - \omega) \int_\omega^z \cos(\xi - \omega) g(\xi) d\xi - \cos(z - \omega) \int_\omega^z \sin(\xi - \omega) g(\xi) d\xi \end{aligned}$$

is also symmetrical with respect to ω and real-valued. \square

Using the reflection principle, it is easy to show that the function $g(z)$ of Proposition 3.1 is periodic with real or purely imaginary period $2P$ if and only if it is symmetrical with respect to $\omega + P$. Then Proposition 3.1 shows that there is a (real) one-parameter family of operators I_a^ω , preserving the symmetry of g at ω . The following proposition shows that symmetry can be preserved at both ω and $\omega + P$ if P is not a multiple of π . If P is a multiple of π , one can prove that a symmetrical and periodic $(D^2 + 1)^{-1}g(z)$ exists if and only if the Fourier expansion of $g(z)$ does not contain a $\cos z$ term.

PROPOSITION 3.2. *Let $g(z)$ be a function, symmetrical with respect to ω and $\omega + P$ and real-valued on $[\omega, \omega + P]$, where P is either real or purely imaginary. If $P \neq k\pi$, $k \in \mathbb{N}$, then there exists an operator I that is inverse to $D^2 + 1$ and such that Ig inherits the abovementioned properties of g .*

Proof. To satisfy the requirements at ω , the operator I has to be of the form $Ig = I^\omega g + d \cos(z - \omega)$, where $d \in \mathbb{R}$ (and depends on g). Then

(80)

$$Ig = I^{\omega+P}g + d \cos(z - \omega) + \frac{1}{2i} \left[e^{i(z-\omega-P)} \int_\omega^{\omega+P} e^{-i(\xi-\omega-P)}g(\xi)d\xi - e^{-i(z-\omega-P)} \int_\omega^{\omega+P} e^{i(\xi-\omega-P)}g(\xi)d\xi \right].$$

Let $A_{1,2}$ denote the (two) terms in the square brackets, respectively, and let $A = \frac{1}{2}(A_1 + A_2)$, $B = \frac{1}{2}(A_1 - A_2)$. Then the square bracket expression in (80) becomes $A \sin(z - \omega - P) - iB \cos(z - \omega - P)$. Therefore

$$(81) \quad Ig = I^{\omega+P}g + (d \cos P - iB) \cos(z - \omega - P) + (A - d \sin P) \sin(z - \omega - P).$$

Using (79) we see that the first two terms in the right-hand side of (81) are symmetrical at $\omega + P$. Thus Ig is symmetrical with respect to $\omega + P$ if and only if

$$(82) \quad d = \frac{2 \int_\omega^{\omega+P} \cos(\xi - \omega - P)g(\xi)d\xi}{\sin P}.$$

Equation (82) can be satisfied for any $P \neq k\pi$, $k \in \mathbb{Z}$. Noting that d and iB are real for both real and purely imaginary P concludes the proof. \square

3.2. Estimates for the operator \mathcal{I}_2 . Consider the trapezoid T introduced in section 2.4. Given a number $z_1 \in (0, \frac{1}{4}\tilde{\omega}_1)$, which is independent of ε , we cut a triangle near the origin with vertices z_1, z_2, z_3 , where z_2, z_3 are the points on the intersection of the imaginary axis and the lines through the points z_1, A and the points z_1, B , respectively. The obtained region is denoted by T_{z_1} . Clearly, for a given $z_0 > 0$ we have $R_{z_0} \subset T_{z_1}$ if z_0 is large enough (i.e., according to (52), if ε_0 is small enough). We define the operator

$$(83) \quad \mathcal{I}_2 = I_{-2i\tilde{\omega}_3}^{\tilde{\omega}_2}.$$

For a given $z \in T_{z_1}$, $\gamma_\pm(z)$ denote contours of integration $[z, B]$ and $[z, A]$, respectively; see Figure 5.

PROPOSITION 3.3. *The inequality*

$$(84) \quad \left| \frac{\xi}{z} \right| > \frac{\omega_1}{\sqrt{\omega_1^2 - \frac{25}{4}\omega_3^2}}$$

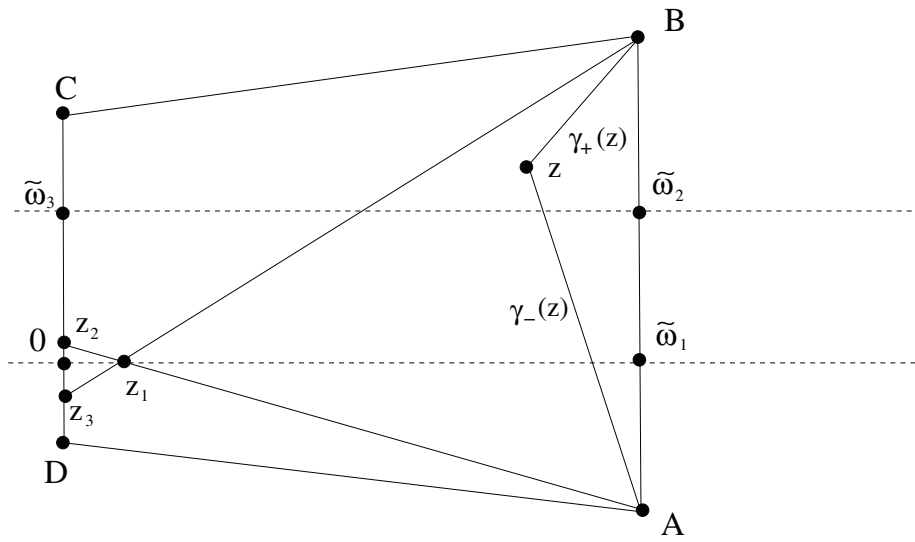


FIG. 5. Region T_{z_1} and contours $\gamma_{\pm}(z)$.

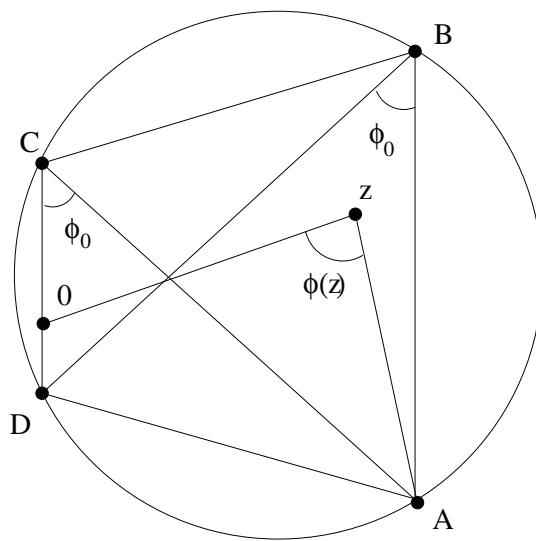


FIG. 6. Trapezoid $ABCD$ and angles $\phi(z)$, ϕ_0 .

holds for all $z \in T_{z_1}$ and all $\xi \in \gamma_+(z) \cap \gamma_-(z)$.

Proof. For an arbitrary $z \in T_{z_1}$, let $\phi(z)$ denote the angle between the segments $\gamma_-(z)$ and $[0, z]$ (see Figure 6), and let $\alpha(z) = \min_{\xi \in \gamma_-(z)} \frac{|\xi|}{|z|}$. We choose the acute angle for $\phi(\tilde{\omega}_1 + 2\tilde{\omega}_3)$ (at the vertex B) and then define $\phi(z)$ in T_{z_1} by continuity. Clearly $\alpha(z) = 1$ if $\phi(z) \geq \frac{\pi}{2}$ and $\alpha(z) = \sin \phi(z)$ if $\phi(z) < \frac{\pi}{2}$. The latter case can happen only if $\Im z \geq 0$; hence $\alpha(z) = 1$ if $\Im z < 0$. Let ϕ_0 be the acute angle between AB and DB . We show that $\phi(z) \geq \phi_0$ for every $z \in T_{z_1}$.

Indeed, let us inscribe the right trapezoid $ABCD$ into a circle. Then (at the vertex C) $\phi(\frac{3}{2}\tilde{\omega}_3) = \phi_0$. Simple geometrical considerations (see Figure 6) imply that

$\min_{z \in T_{z_1}} \phi(z)$ is attained at $z = \frac{3}{2}\tilde{\omega}_3$. Thus, $\alpha(z) \geq \sin \phi_0$. Therefore

$$(85) \quad \alpha(z) \geq \sin \phi_0 = \frac{\omega_1}{\sqrt{\omega_1^2 - \frac{25}{4}\omega_3^2}}$$

for any $z \in T_{z_1}$. The same estimate holds for the contour $\gamma_+(z)$. The proof is completed. \square

PROPOSITION 3.4. For any $z \in T_{z_1}$

$$(86) \quad |\mathcal{I}_2 f(z)| \leq \sqrt{1 - 4\frac{\omega_1^2}{\omega_3^2}} \max_{\xi \in \gamma_+(z) \cup \gamma_-(z)} |f(\xi)|,$$

where f is a continuous function on T_{z_1} .

Proof. According to the construction of γ_+, γ_- ,

$$(87) \quad \Re[i(\xi - z)] \leq -\frac{|\omega_3|}{\sqrt{4\omega_1^2 - \omega_3^2}}|z - \xi| \quad \text{and} \quad \Re[i(z - \xi)] \leq -\frac{|\omega_3|}{\sqrt{4\omega_1^2 - \omega_3^2}}|z - \xi|$$

for all $\xi \in \gamma_{\pm}(z)$, respectively. Therefore

$$(88) \quad \begin{aligned} |\mathcal{I}_2 f(z)| &\leq \frac{1}{2} \left| \int_{\gamma_-(z)} e^{i(z-\xi)} f(\xi) d\xi + \int_{\gamma_+(z)} e^{i(z-\xi)} f(\xi) d\xi \right| \\ &\leq \max_{\xi \in \gamma_+(z) \cup \gamma_-(z)} |f(\xi)| \int_0^\infty e^{-\left(|\omega_3|/\sqrt{4\omega_1^2 - \omega_3^2}\right)\lambda} d\lambda, \end{aligned}$$

which implies (86). \square

LEMMA 3.5. If

$$(89) \quad |f(z)| \leq \frac{1}{|z|^a}$$

in T_{z_1} with $a \geq 0$, then

$$(90) \quad |\mathcal{I}_2 f(z)| \leq \frac{\beta}{|z|^a},$$

where $\beta = \left(1 - \frac{25\omega_3^2}{4\omega_1^2}\right)^{\frac{1}{2}} \left(1 - 4\frac{\omega_1^2}{\omega_3^2}\right)^{\frac{a}{2}}$.

This lemma is a direct consequence of Propositions 3.3 and 3.4.

3.3. Estimates for the operator \mathcal{I}_1 . Solution $\mathcal{I}_1 g$ to linear differential equation (43), where g denotes the right-hand side of (43), that we consider in the paper is given by (50). In order to estimate \mathcal{I}_1 , we need the following simple statement.

PROPOSITION 3.6. Let a ray λ on the complex z -plane lie outside the disk $|z| < a$, $a > 0$. Then for any $\alpha \geq 2$

$$(91) \quad \left| \int_\lambda \frac{dz}{z^\alpha} \right| \leq \frac{\pi}{a^{\alpha-1}}.$$

Proof. Let m denote the line containing the ray λ . Suppose m is tangent to the disk $|z| < a$ at the point ω . If $z \in m$, then $|z|^2 = a^2 + x^2$, where $x = |z - \omega|$. Then

$$(92) \quad \left| \int_\lambda \frac{dz}{z^\alpha} \right| \leq \int_{-\infty}^\infty \frac{dx}{(a^2 + x^2)^{\frac{\alpha}{2}}} \leq \frac{2}{a^{\alpha-2}} \int_0^\infty \frac{dx}{(a^2 + x^2)} = \frac{\pi}{a^{\alpha-1}}.$$

This estimate will definitely hold if m does not intersect the closed disk $|z| \leq a$.

Suppose now that m intersects the disk $|z| < a$. Let ω denote the vertex of the ray λ , and let ζ denote one of the two points which lie on the distance a from both ω and m . Then $|\xi - \zeta| < |\xi|$ for any $\xi \in \lambda$, and we can apply the estimate (92) again, placing the origin at ζ . \square

Notice that estimates of Proposition 3.3 are applicable to the contours in (50). Then, using Lemma 2.4 and Proposition 3.6, we can easily establish the following estimates.

LEMMA 3.7. *If $|g(z)| \leq K|z|^{-n}$, $n \geq 6$, or if $|g(z)| \leq K|z|^{-4}$ in T_{z_1} , then*

$$(93) \quad |\mathcal{I}_1 g(z)| \leq \frac{M_1 K}{|z^{n-2}|} \quad \text{and} \quad |\mathcal{I}_1 g(z)| \leq \frac{M_2 K}{\varepsilon |z^{-3}|},$$

respectively, where $M_1 = (1 - \frac{25}{4} \frac{\omega_3^2}{\omega_1^2})^{\frac{n-5}{2}} \pi AB$ and $M_2 = \max\{2|\omega_3|, \sqrt{\omega_1^2 - \frac{9}{4}\omega_3^2}\} AB$.

3.4. Proof of Theorem 2.1. *Proof.* We start with the obvious observation that

$$(94) \quad \varepsilon \leq \frac{\sqrt{\omega_1^2 - 4\omega_3^2}}{|z|}$$

if $z \in T_{z_1}$. Using the first inequality in (53), we obtain that $f(q)$ from (36) satisfies

$$(95) \quad |f| \leq \tilde{K}_0 |z|^{-6} + \varepsilon \tilde{K}_1 |z|^{-4}$$

in T_{z_1} with some positive constants \tilde{K}_0, \tilde{K}_1 that do not depend on ε and $J(\varepsilon) \in E$. Thus, according to (49) and the estimates for $\mathcal{I}_{1,2}$,

$$(96) \quad |\mathcal{N}w_0| \leq \frac{K}{2} \left(\frac{\varepsilon^{\alpha-1}}{|z|^3} + \frac{1}{|z|^4} \right)$$

in T_{z_1} for some constant $K > 0$ that does not depend on $\varepsilon, J(\varepsilon)$, and z_1 .

Let us now prove by induction that by choosing a sufficiently large $z_1 > 0$ we can find $\delta > 0$, such that the estimate

$$(97) \quad \Delta w_n \leq (\delta M)^{n-1} \frac{K}{|z|^3}, \quad \text{where} \quad \delta M \leq \frac{1}{2},$$

holds in T_{z_1} for all $n = 1, 2, \dots$, all $J(\varepsilon) \in E$, and all sufficiently small ε . Here $M = \max\{M_1, M_2\}$, where M_1 is given in Lemma 3.7 with $n = 6$.

According to Lemma 2.4, the estimate for v_{22} for all $z \in T_{z_1}$ and $J(\varepsilon) \in E$ is

$$(98) \quad |v_{22}(z)| \leq B|z|^4 + \frac{L}{\varepsilon^7 |z|^3}$$

for some $L > 0$. Then, according to (94), (96), and (53), we can choose $K > 0$ to be so large that $w_1 = \mathcal{N}w_0 + \varepsilon^6 b(\varepsilon)v_{22}$ satisfies

$$(99) \quad |w_1| \leq \frac{K}{|z|^3}$$

in T_{z_1} for all $J(\varepsilon) \in E$ and all sufficiently small ε . Thus, for $n = 1$, (97) has been established. Assume that this estimate is true for $k = 1, 2, \dots, n$, and let us establish it for $k = n + 1$. First, we can represent

$$(100) \quad \Delta w_{n+1} = \mathcal{I}_1 \circ \left[(1 - 12cq^2 \mathcal{I}_2) [\mathcal{I}_2(w_n + w_{n-1}) \cdot \mathcal{I}_2 \Delta w_n] - 12q^2(1 - c) \Delta w_n \right. \\ \left. + 12cq^2 \mathcal{I}_2(\varepsilon^2 + 12cq^2 \mathcal{I}_2) \Delta w_n \right].$$

Taking into account (99) and Lemma 3.5, we obtain the estimates

(101)

$$\frac{12Ql\varepsilon^\alpha}{|z|} \left| \frac{\Delta w_n}{z} \right|, 12cQ\beta \left[\frac{\varepsilon}{|z|} \left| \frac{\varepsilon \Delta w_n}{z} \right| + \frac{12cQ\beta}{|z|} \left| \frac{\Delta w_n}{z^3} \right| \right], 4K\beta^2 \left(1 + \frac{12cQ\beta}{|z|^2} \right) \left| \frac{\Delta w_n}{z^3} \right|$$

for the second, third, and first (nonlinear) terms in the square brackets in (100), respectively. For the latter estimate, observe that both w_n and w_{n-1} are bounded by $\frac{2K}{|z|^3}$ according to (97). (Note that the constants β, K, M are independent of z_1 .) Applying operator \mathcal{I}_1 to these three terms and utilizing Lemma 3.7, we obtain the estimates

(102)

$$\frac{12Ql}{|z|} \varepsilon^{\alpha-1} M |\Delta w_n|, 12cQ\beta \left[\varepsilon + \frac{12cQ\beta}{|z|} \right] \frac{M}{|z|} |\Delta w_n|, 4K\beta^2 \left(1 + \frac{12cQ\beta}{|z|^2} \right) \left| \frac{M \Delta w_n}{z} \right|.$$

Factoring $M|\Delta w_n|$, we obtain $|\Delta w_{n+1}| \leq \delta M |\Delta w_n|$, where

$$(103) \quad \delta = \frac{12Ql\varepsilon^{\alpha-1} + 12cQ\beta\varepsilon}{|z|} + \frac{12^2c^2Q^2\beta^2}{|z|^2} + \frac{4K\beta^2}{|z|} \left(1 + \frac{12cQ\beta}{|z|^2} \right).$$

By choosing a sufficiently large $z_1 > 0$ and sufficiently small $\varepsilon_0 > 0$, we can guarantee the condition $\delta M \leq \frac{1}{2}$ in T_{z_1} for all $\varepsilon \leq \varepsilon_0$ and for all $J(\varepsilon) \in E$. Thus, (97) holds for $k = n + 1$. We can choose some $z_0 \geq z_1$ so that $R_{z_0} \subset T_{z_1}$ and $\varepsilon_0 z_0 \leq \max\{\omega_1, |\omega_3|\}$ (taking smaller ε_0 if necessary). Thus, inequalities (97) prove convergence of iterations to a solution of (42) that is uniform in $R_{z_0} \times (0, \varepsilon_0]$ for all $J(\varepsilon) \in E$.

In the case $\varepsilon = 0$ the region T_{z_1} becomes the right half-plane with the appropriate cut, $q^2 = \frac{6}{z^2}$, $f(q) = -\frac{6l}{z^6}$, and the fundamental solutions $v_{1,2}$ to the homogeneous equations (43) become $v_1(z) = \frac{-12}{z^3}$ and $v_2(z) = \frac{-z^4}{84}$. Thus, for the case $\varepsilon = 0$ Lemma 2.4 becomes trivial, and the proof of convergence of iterations $w = \sum_1^\infty \Delta w_n$ that solve (42) holds for the case $\varepsilon = 0$. (Note that in this case the contours of integration in the integral operators $\mathcal{I}_{1,2}$ are the rays, emanating from z , with the slope $\frac{\Im w}{\Re w}$, where w is the beginning of the corresponding contour in the case $\varepsilon > 0$; in other words, $\arg w$ does not depend on ε .) Thus, the proof of Theorem 2.1 is completed. \square

Remark 3.8. It is easy to verify that $\alpha > 1$ in Theorem 2.1 implies that there exist some $L_n > 0$ such that for all $n \in \mathbb{N}$

$$(104) \quad |\Delta w_n| \leq L_n |z|^{-n(1+\tilde{\alpha})-2}$$

in R_{z_0} , where $\tilde{\alpha} = \min\{1, \alpha - 1\}$. Moreover, for all $n \in \mathbb{N}$ there exist some $K_n > 0$ such that

$$(105) \quad \left| \sum_{j=n}^\infty \Delta w_j \right| \leq K_n |z|^{-n(1+\tilde{\alpha})-2}.$$

Let $P \subset T$ denote the triangle with vertices $\tilde{\omega}_3, \tilde{\omega}_1$, and $\tilde{\omega}_1 + 2\tilde{\omega}_3$. It is clear that $z \in P$ implies that all the contours of integration in operators $\mathcal{I}_{1,2}$ are in P . Notice that in P

$$(106) \quad \min\{\omega_1, |\omega_3|\} \leq \varepsilon |z| \leq \sqrt{\omega_1^2 - 4\omega_3^2}.$$

The latter condition allows us to prove convergence of iterations in P without the requirement (53) in Theorem 2.1.

COROLLARY 3.9. *If instead of (53) in Theorem 2.1 we require only*

$$(107) \quad \lim_{\varepsilon \rightarrow 0} b(\varepsilon) = \lim_{\varepsilon \rightarrow 0} [c(\varepsilon) - 1] = 0,$$

the statement of Theorem 2.1 is still valid in $P \times [0, \varepsilon_0]$; i.e., there exists some $\varepsilon_0 > 0$ such that the series (54) is uniformly convergent in $P \times [0, \varepsilon_0]$.

Proof. The function

$$(108) \quad h(\varepsilon) = \max\{|b(\varepsilon)|, |c(\varepsilon) - 1|\}$$

is continuous and nonnegative, and $\lim_{\varepsilon \rightarrow 0} h(\varepsilon) = 0$. Then inequalities (95) and (96) become

$$(109) \quad |f| \leq \tilde{K}_0 |z|^{-6} + h(\varepsilon) \tilde{K}_1 |z|^{-4} \quad \text{and} \quad |\mathcal{N}w_0| \leq \frac{K}{2} \left(\frac{h(\varepsilon)}{\varepsilon |z|^3} + \frac{1}{|z|^4} \right),$$

respectively. Then, as in Theorem 2.1, we can choose $K > 0$ so large that

$$(110) \quad |w_1| \leq \frac{\tilde{h}(\varepsilon)K}{|z|^2}$$

in P , where $\tilde{h}(\varepsilon) = \max\{h(\varepsilon), \varepsilon^2\}$.

As in Theorem 2.1, we will prove convergence of (54) by induction if we show that by choosing a sufficiently small $\varepsilon_0 > 0$ we can find $\delta > 0$, such that the estimate

$$(111) \quad \Delta w_n \leq (\delta M)^{n-1} \frac{\tilde{h}(\varepsilon)K}{|z|^2}, \quad \text{where} \quad \delta M \leq \frac{1}{2},$$

holds in P for all $n = 1, 2, \dots$ and all $J(\varepsilon) \in E$. Taking into account (111) and Lemma 3.5, we obtain the estimates

$$(112) \quad 12Ql\tilde{h}(\varepsilon) \left| \frac{\Delta w_n}{z^2} \right|, \\ 12cQ\beta \left[\left| \frac{\varepsilon^2 \Delta w_n}{z^2} \right| + \frac{12cQ\beta}{|z|^2} \left| \frac{\Delta w_n}{z^2} \right| \right], \quad 4\tilde{h}(\varepsilon)K\beta^2 \left(1 + \frac{12cQ\beta}{|z|^2} \right) \left| \frac{\Delta w_n}{z^2} \right|$$

for the second, third, and first (nonlinear) terms in the square brackets in (100), respectively. Making sure that $\tilde{h}(\varepsilon)$ and $|z|^{-2}$ are sufficiently small by choosing a sufficiently small ε_0 and using Lemma 3.7, we can repeat the arguments of Theorem 2.1 in order to complete the proof of convergence of iterations (54). \square

Remark 3.10. In fact, we can prove convergence of iterations (54) in larger domains than those stated in Theorem 2.1 and Corollary 3.9. Fix some $n \in \mathbb{N}$. For Theorem 2.1, we replace the trapezoid T with vertices $ABCD$ by the (concave) hexagon \mathcal{T} with vertices $AONBCD$, where $O = \tilde{\omega}_1$, $N = n\tilde{\omega}_1 + \tilde{\omega}_3$, and other vertices are the same as in T ; see Figure 7. For $z \in \mathcal{T}$, contours of integration in operators $\mathcal{I}_{1,2}$ are the same as before, except for contour $\gamma_-(z)$ that is now the union of segments $[z, O]$ and $[OA]$. For Corollary 3.9, we replace the triangle P with vertices $\tilde{\omega}_3, \tilde{\omega}_1, \tilde{\omega}_1 + 2\tilde{\omega}_3$ by a quadrilateral \mathcal{P} with vertices $\tilde{\omega}_1, n\tilde{\omega}_1 + \tilde{\omega}_3, \tilde{\omega}_1 + 2\tilde{\omega}_3, \tilde{\omega}_3$. For $z \in \mathcal{P}$, contours of integration in operators $\mathcal{I}_{1,2}$ are the same as before. It is easy to verify that the estimates for operators $\mathcal{I}_{1,2}$ from Lemmas 3.7 and 3.5 hold in domains \mathcal{T} and \mathcal{P} , possibly with larger constants that depend on n . Now proofs of Theorem 2.1 and Corollary 3.9 can be extended to the domains \mathcal{T} and \mathcal{P} , respectively.

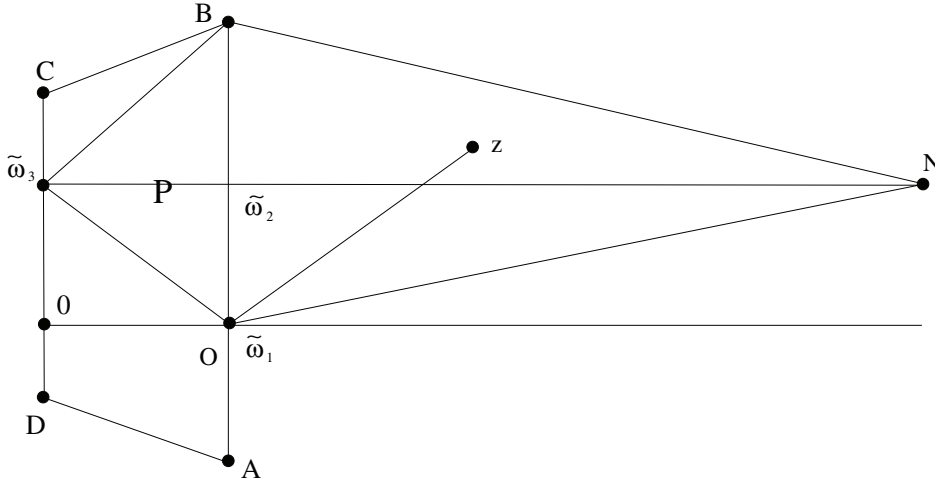


FIG. 7. Hexagon T with vertices $AONBCD$. The contour $\gamma_-(z)$ is the union of segments $[z, O]$ and $[OA]$.

3.5. Symmetry at $\tilde{\omega}_2$. Here we will show that any solution to inner equation (23), obtained through iterations (54), is symmetrical with respect to $z = \tilde{\omega}_2$. Consequently, corresponding solutions to the outer equation (1), connected with (23) via (22), are symmetric at $x = \omega_2$.

The operator \mathcal{I}_1 can be represented in the following form.

PROPOSITION 3.11. *If a function $g(z)$ is analytic and bounded in T_{z_1} , then*

$$(113) \quad \mathcal{I}_1 g = v_1(z) \int_{\tilde{\omega}_2}^z \frac{\int_{\tilde{\omega}_2}^t v_1(\tau)g(\tau)d\tau}{v_1^2(t)} dt.$$

Proof. Using (46) and integration by parts, one gets

$$(114) \quad \begin{aligned} \mathcal{I}_1 g &= -v_1(z) \int_{\tilde{\omega}_2}^z \left(v_1(t)g(t) \int_0^t \frac{d\tau}{v_1^2(\tau)} \right) dt + v_2(z) \int_{\tilde{\omega}_2}^z v_1(t)g(t)dt \\ &= v_1(z) \times \left[- \int_{\tilde{\omega}_2}^t v_1(\tau)g(\tau)d\tau \cdot \int_0^t \frac{d\tau}{v_1^2(\tau)} \Big|_{\tilde{\omega}_2}^z \right. \\ &\quad \left. + \int_{\tilde{\omega}_2}^z \frac{\int_{\tilde{\omega}_2}^t v_1(\tau)g(\tau)d\tau}{v_1^2(\tau)} dt + \int_0^z \frac{dt}{v_1^2(\tau)} \cdot \int_{\gamma(z)} v_1(t)g(t)dt \right] \\ &= v_1(z) \lim_{\zeta \rightarrow \tilde{\omega}_2} \left[\int_{\zeta}^z v_1(t)g(t)dt \cdot \int_0^{\zeta} \frac{dt}{v_1^2(\tau)} \right] + v_1(z) \int_{\tilde{\omega}_2}^z \frac{\int_{\tilde{\omega}_2}^t v_1(\tau)g(\tau)d\tau}{v_1^2(\tau)} dt. \end{aligned}$$

Note that the limit in the latter expression is zero because $\int_{\gamma(\zeta)} v_1(t)g(t)dt$ has a zero of order two at $\tilde{\omega}_2$, whereas $\int_0^{\zeta} \frac{dt}{v_1^2(\tau)}$ has a simple pole there. \square

Based on (45) and (73), solutions to the homogeneous equation (44) preserving the symmetry at the points $\tilde{\omega}_j, j = 1, 2, 3$, are given by

$$(115) \quad v_{2j}(z) = v_2(z) - \frac{1}{6\varepsilon^7 \Delta} \left[\frac{3g_3}{2} \omega_j - g_2 \zeta(\omega_j) \right] v_1(z), \quad j = 1, 2, 3,$$

respectively. Indeed, direct calculations show that $D_z v_{2j}(\tilde{\omega}_j) = 0$. Then, according to (44),

$$(116) \quad D_z^3 v_{2j} = (\varepsilon^2 + 12q^2)D_z v_{2j} + 12(q^2)'v_{2j} = 0$$

at $z = \tilde{\omega}_j$. By continuing this argument for higher derivatives, we prove symmetry of v_{2j} at $z = \tilde{\omega}_j$. Also, note that $v_{2j}(\tilde{\omega}_j) \in \mathbb{R}$. Then, according to the differential equation for v_{2j} , all even derivatives of v_{2j} at $z = \tilde{\omega}_j$ are also real.

The fact that operator \mathcal{I}_1 preserves the symmetry at the point $\tilde{\omega}_2$ and real-valuedness along the line $\Im z = \Im \tilde{\omega}_2$ is a direct consequence of Proposition 3.11. Thus, according to (36), Δw_1 is symmetric with respect to $\tilde{\omega}_2$ and real-valued on the line $\Im z = \Im \tilde{\omega}_2$. Then, according to (100) and to Proposition 3.1, the same is true for $w = \sum_1^\infty \Delta w_n$. Now, according to (34), the corresponding solution $v(z, \varepsilon)$ to (23) satisfies

$$(117) \quad v'(\tilde{\omega}_2, \varepsilon) = v'''(\tilde{\omega}_2, \varepsilon) = 0$$

and is real-valued on $\Im z = \Im \tilde{\omega}_2$.

Thus, in Theorem 2.1 and Corollary 3.9 we have constructed a two-parameter family of solutions to (23), where $b(\varepsilon)$ and $c(\varepsilon)$ are the parameters, that are real along $\Im z = \Im \tilde{\omega}_2$ and satisfy (117).

3.6. Calculation of $\mathcal{I}_1 f(q)$. In order to calculate $\mathcal{I}_1 f(q)$, we use (36) and (45) to evaluate the integral

$$(118) \quad \int_{\tilde{\omega}_2}^t v_1(\tau) f(q(\tau)) d\tau = \int f(q) d(6q^2) \\ = c \left[\frac{c-1-4\varepsilon^2}{3} (6q^2)^3 - 5/6(6q^2)^4 - 2\varepsilon^6 C(\varepsilon)(6q^2) \right] \Big|_{\tilde{\omega}_2}^z$$

in (113). Switching to the original variable $x = \varepsilon z$ and using (113), we obtain

$$(119) \quad \mathcal{I}_1 f = 6\varepsilon^3 c \wp'(x) \\ \times \int_{\omega_2}^x \frac{(6\wp(t) - \frac{1}{2}) \left[\frac{c-1-4m^2}{3\varepsilon} (6\wp(s) - \frac{1}{2})^2 - \frac{5\varepsilon}{6} (6\wp(s) - \frac{1}{2})^3 - 2\varepsilon C(\varepsilon) \right] \Big|_{s=\omega_2}^{s=t} dt}{144 \prod_1^3 (\wp(t) - e_j)}.$$

The integrand in (119) is nonsingular at $t = \omega_2$ since the numerator has zero at this point. Using this fact, we can cancel the factor $(\wp - e_2)$ in both the numerator and denominator. Then direct calculations show that the integrand in (119) becomes

$$(120) \quad \frac{c-1-4\varepsilon^2}{3\varepsilon} \cdot \frac{\wp^2 + (e_2 - \frac{1}{4})\wp + e_2^2 - \frac{1}{4}e_2 + \frac{1}{48}}{4(\wp - e_1)(\wp - e_3)} \\ - \frac{\frac{5\varepsilon}{6} (6\wp(t) - \frac{1}{2})^4}{144 \prod_1^3 (\wp(t) - e_j)} - \frac{2\varepsilon C(\varepsilon)(6\wp(t) - \frac{1}{2})}{144 \prod_1^3 (\wp(t) - e_j)}.$$

Let us focus on the first term of (120). Separating the integer part, we can rewrite it as

$$(121) \quad \frac{c-1-4\varepsilon^2}{3 \cdot 4\varepsilon} \left[1 - \frac{1}{4} \frac{\wp + e_2 - \frac{1}{6}}{(\wp - e_1)(\wp - e_3)} \right] \\ = \frac{c-1-4\varepsilon^2}{3 \cdot 4\varepsilon} \left[1 - \frac{B_0 + B_1\wp(t - \omega_1) + B_3\wp(t - \omega_3)}{4} \right],$$

where

$$(122) \quad B_1 = -\frac{e_3 + \frac{1}{6}}{(e_1 - e_3)^2(e_1 - e_2)}, \quad B_3 = \frac{e_1 + \frac{1}{6}}{(e_1 - e_3)^2(e_2 - e_3)}, \quad \text{and} \quad B_0 = -B_1e_1 - B_3e_3.$$

The expressions for B_1 , B_3 , and B_0 were obtained by evaluating the principal parts of the left-hand side of (122) at $t = \omega_1$, $t = \omega_3$, and $t = 0$, respectively.

Using similar calculations, we obtain expressions

$$(123) \quad -\frac{15\varepsilon}{2}\wp(t) + \frac{5\varepsilon}{2} - \frac{5\varepsilon}{24} [E_0 + E_1\wp(t - \omega_1) + E_3\wp(t - \omega_3)] \quad \text{and} \\ -\frac{\varepsilon C(\varepsilon)}{12} [D_0 + D_1\wp(t - \omega_1) + D_3\wp(t - \omega_3)]$$

for the second and the third terms of (120), respectively, where

$$(124) \quad D_1 = \frac{1}{(e_1 - e_3)^2(e_1 - e_2)}, \quad D_3 = -\frac{1}{(e_1 - e_3)^2(e_2 - e_3)}, \quad D_0 = -D_1e_1 - D_3e_3$$

and

$$(125) \quad E_3 = -\frac{\frac{9}{4}e_3 + \frac{3}{2}e_2 + 36e_2^3 - \frac{1}{3}}{(e_1 - e_3)^2(e_2 - e_3)}, \quad E_0 = -E_1e_1 - E_3e_3.$$

(The constant E_0 could be calculated similarly to E_3 , but we do not need its explicit value.)

Substituting (121)–(125) into (119), we obtain

$$(126) \quad \mathcal{I}_1 f = c\varepsilon^3 \wp'(x) \{ 45\varepsilon[\zeta(x) - \zeta(\omega_2)] + 6G_1[\zeta(x - \omega_1) - \zeta(\omega_1)] \\ + 6G_3[\zeta(x - \omega_3) - \zeta(\omega_3)] + [-15\varepsilon + (c - 1 - 4\varepsilon^2)(2\varepsilon)^{-1} - 6G_1e_1 - 6G_3e_3] (x - \omega_2) \},$$

where $x = \varepsilon z$ and

$$(127) \quad G_j = -\frac{B_j(c - 1 - 4\varepsilon^2)}{48\varepsilon} - \frac{5\varepsilon}{24} E_j - \frac{\varepsilon C(\varepsilon)}{12} D_j, \quad j = 1, 3.$$

3.7. The limit of iterations as $\varepsilon \rightarrow 0$. In this subsection we assume that conditions of Theorem 2.1 hold. In the $\varepsilon = 0$ case, direct calculation, based on (50) and (49), shows that

$$(128) \quad \Delta w_1(z, 0) = -\frac{90}{z^4} + \frac{12 \cdot 6!}{7} \left(-z^{-3} \int_{\infty}^z t^2 \mathcal{I}_2[t^{-6}] dt + z^4 \int_{\infty}^z t^{-5} \mathcal{I}_2[t^{-6}] dt \right).$$

It is easy to check that Δw_n , $n = 2, 3, \dots$, is a function, analytic in T_{z_1} and possessing an asymptotic expansion in powers of z^{-2} there, such that

$$(129) \quad \Delta w_n(z) \sim O(z^{-2(n+1)}), \quad z \rightarrow \infty, \quad |\arg z| < \frac{\pi}{2}.$$

Then, according to Corollary 1.3, the solution (34) to (23) in the case $\varepsilon = 0$ coincides with $v_+(z)$, i.e., $v(z, 0) = v_+(z)$.

Let $w(z, \varepsilon)$ denote the solution (54) from Theorem 2.1. In order to prove the continuity of $w(z, \varepsilon)$ in ε at $\varepsilon = 0$, we need the following statement.

PROPOSITION 3.12. *Let $g(z, \varepsilon)$ be analytic in z in T_{z_1} and continuous in ε for all sufficiently small $\varepsilon \geq 0$. Moreover, let $g(z, 0) \sim O(z^{-b})$ and $\Delta_\varepsilon g(z, \varepsilon) \sim O(z^{-a})$ in T_{z_1} , the latter uniformly in ε , where $a, b \geq 0$ and $\Delta_\varepsilon g(z, \varepsilon) = g(z, \varepsilon) - g(z, 0)$. Then*

$$(130) \quad \Delta_\varepsilon \mathcal{I}_2 g(z, \varepsilon) = O(z^{-a}) + O(\varepsilon^b),$$

and, if $b > 5$,

$$(131) \quad \Delta_\varepsilon \mathcal{I}_1 g(z, \varepsilon) = O(z^{2-a}) + O(\varepsilon^{b-5} z^{-3}) + O(z^{-b+4} \varepsilon^2),$$

where both (130) and (131) are uniform in T_{z_1} and in small $\varepsilon \geq 0$.

Proof. Since $g(z, \varepsilon)$ is analytic in T_{z_1} ,

$$(132) \quad \begin{aligned} \Delta_\varepsilon \mathcal{I}_2 g(z, \varepsilon) &= \frac{1}{2i} \left[\int_{\tilde{\omega}_1 - \tilde{\omega}_3}^z e^{i(z-\xi)} \Delta_\varepsilon g(\xi, \varepsilon) d\xi - \int_{\tilde{\omega}_1 + 2\tilde{\omega}_3}^z e^{-i(z-\xi)} \Delta_\varepsilon g(\xi, \varepsilon) d\xi \right] \\ &\quad + \frac{1}{2i} \left[\int_\infty^{\tilde{\omega}_1 - \tilde{\omega}_3} e^{i(z-\xi)} g(\xi, 0) d\xi - \int_\infty^{\tilde{\omega}_1 + 2\tilde{\omega}_3} e^{-i(z-\xi)} g(\xi, 0) d\xi \right]. \end{aligned}$$

According to Lemma 3.5, estimate (130) follows directly (132).

To prove (131), we first consider $\Delta_\varepsilon q^2(z, \varepsilon)$. Representing $\wp(x) = [\wp(x) - \frac{1}{x^2}] + \frac{1}{x^2} = \Psi(x) + \frac{1}{x^2}$, we note that $\Psi(x)$ is analytic in the trapezoidal region \tilde{T} (vertices at $\omega_1 - \omega_3, \omega_1 + 2\omega_3, \frac{3}{2}\omega_3, -\frac{1}{2}\omega_3$). (The trapezoid \tilde{T} is the image of the trapezoid T under the scaling $x = \varepsilon z$.) Since $q^2(z, 0) = \frac{1}{z^2}$ for any $z \in T$, we have

$$(133) \quad \Delta_\varepsilon q^2(z, \varepsilon) = \varepsilon^2 \left[\Psi(\varepsilon z) - \frac{1}{12} \right] = O(\varepsilon^2)$$

uniformly in T . Similarly, using (45), (46), and (71), we obtain

$$(134) \quad v_1(z, \varepsilon) - v_1(z, 0) = O(\varepsilon^3)$$

and

$$(135) \quad v_2(z, \varepsilon) - v_2(z, 0) = O(\varepsilon^2 z^6)$$

uniformly in T .

Now,

(136)

$$\begin{aligned} \Delta_\varepsilon \mathcal{I}_1 g(z, \varepsilon) = & - \left[v_1(z, \varepsilon) \int_{\tilde{\omega}_2}^z v_2(\xi, \varepsilon) g(\xi, \varepsilon) d\xi - v_1(z, 0) \int_{\tilde{\omega}_2}^z v_2(\xi, 0) g(\xi, 0) d\xi \right] \\ & + \left[v_2(z, \varepsilon) \int_{\tilde{\omega}_2}^z v_1(\xi, \varepsilon) g(\xi, \varepsilon) d\xi - v_2(z, 0) \int_{\tilde{\omega}_2}^z v_1(\xi, 0) g(\xi, 0) d\xi \right] \\ & + \left[-v_1(z, 0) \int_{\infty}^{\tilde{\omega}_2} v_2(\xi, 0) g(\xi, 0) d\xi + v_2(z, 0) \int_{\infty}^{\tilde{\omega}_2} v_1(\xi, 0) g(\xi, 0) d\xi \right]. \end{aligned}$$

The third term in the right-hand side of (136) is of the order $O(\varepsilon^{b-2})$ as $\varepsilon \rightarrow 0$, while the first term can be represented as

$$\begin{aligned} (137) \quad & - v_1(z, \varepsilon) \int_{\tilde{\omega}_2}^z [v_2(\xi, \varepsilon) \Delta_\varepsilon g(\xi, \varepsilon) + g(\xi, 0) \Delta_\varepsilon v_2(\xi, \varepsilon)] d\xi \\ & - \Delta_\varepsilon v_1(z, \varepsilon) \int_{\tilde{\omega}_2}^z v_2(\xi, 0) g(\xi, 0) d\xi. \end{aligned}$$

Using Lemma 2.4 together with (134), (135), we obtain estimates

$$(138) \quad O(z^{2-a}) + O(\varepsilon^2 z^{4-b}) + O(\varepsilon^3 z^{5-b})$$

for the three terms of (137) that are uniform in T and in small $\varepsilon \geq 0$. The second term in the right-hand side of (136) yields the same result. Now (131) follows from (138) and (94). \square

THEOREM 3.13. (1) *If in condition (53) of Theorem 2.1 we require $\alpha > 1$, then*

$$(139) \quad \lim_{\varepsilon \rightarrow 0} v(z, \varepsilon) = v_+(z)$$

uniformly on compact subsets of T_{z_1} .

(2) *If condition (53) of Theorem 2.1 is replaced by*

$$(140) \quad c(\varepsilon) = 1 + \varepsilon^\alpha \tilde{c}(\varepsilon) \quad \text{and} \quad b(\varepsilon) = \varepsilon^\alpha \tilde{b}(\varepsilon),$$

where $\alpha = 1$ and functions $\tilde{c}(\varepsilon), \tilde{b}(\varepsilon)$ are continuous for small $\varepsilon \geq 0$, then

$$(141) \quad \lim_{\varepsilon \rightarrow 0} v(z, \varepsilon) = v_+(z - z_*)$$

uniformly on compact subsets of T_{z_1} , where the translation z_ depends on $\tilde{b}(0), \tilde{c}(0)$.*

Proof. (1) As a consequence of (133), we obtain that (139) is equivalent to

$$(142) \quad \lim_{\varepsilon \rightarrow 0} w(z, \varepsilon) = w(z, 0)$$

as $\varepsilon \rightarrow 0$ uniformly on compact subsets of T_{z_1} . Since $w(z, \varepsilon) = \sum_{n=1}^\infty \Delta w_n(z, \varepsilon)$, where the series converges uniformly for all sufficiently small $\varepsilon \geq 0$ and all $z \in T_{z_1}$, it is sufficient to prove the continuity of each $\Delta w_n(z, \varepsilon)$ in $\varepsilon \geq 0$ in order to prove the continuity of $w(z, \varepsilon)$.

According to (36) and (94),

$$(143) \quad f(z, 0) = -6!z^{-6} \quad \text{and} \quad \Delta_\varepsilon f(z, \varepsilon) = O(\varepsilon^\alpha z^{-4}).$$

Then, according to Proposition 3.12,

$$(144) \quad \Delta_\varepsilon \mathcal{I}_1 f(z, \varepsilon) = O(\varepsilon^{\tilde{\alpha}} z^{-3}),$$

where $\tilde{\alpha} = \min\{\alpha - 1, 1\}$. Using

$$(145) \quad \Delta_\varepsilon [h(z, \varepsilon)g(z, \varepsilon)] = h(z, \varepsilon)\Delta_\varepsilon g(z, \varepsilon) + \Delta_\varepsilon h(z, \varepsilon)g(z, 0)$$

together with (36) and (130), we prove that $\Delta w_1(z, \varepsilon) = w_1(z, \varepsilon)$ also satisfies (144).

Now $\lim_{\varepsilon \rightarrow 0} \Delta w_n(z, \varepsilon) = \Delta w_n(z, 0)$, $n = 2, 3, \dots$, can be proved by induction. Assume that $\Delta_\varepsilon(\Delta w_k(z, \varepsilon))$ satisfies (144) for every $k \in \mathbb{N}$, $k \leq n$. We want to prove the statement for $k = n + 1$. According to (94), (129), and (130), $\Delta_\varepsilon g(z, \varepsilon) = O(\varepsilon^{\tilde{\alpha}} z^{-6})$, where $g(z, \varepsilon)$ denotes the argument of \mathcal{I}_1 in (100). It is easy to see that $g(z, 0) = O(z^{-2(n+3)})$. Then the fact that $\Delta_\varepsilon[\Delta w_{n+1}(z, \varepsilon)]$ satisfies (144) follows from (131). The proof of part (1) is completed.

(2) According to (126), for a given $z \in T_{z_1}$

$$(146) \quad \lim_{\varepsilon \rightarrow 0} \mathcal{I}_1 f = -\frac{90}{z^4} + \frac{\tilde{c}(0) \left\{ \sum_{j=1,3} B_j [2\zeta(\omega_j) - \omega_2 e_j] - 4\omega_2 \right\}}{4z^3},$$

where B_j are defined by (122). Similarly,

$$(147) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^7 \tilde{b}(\varepsilon) v_{22} = -\tilde{b}(0) \frac{\frac{3g_3}{2} \omega_2 - g_2 \zeta(\omega_2)}{3\Delta z^3}.$$

Thus,

$$(148) \quad \lim_{\varepsilon \rightarrow 0} w_1(z, \varepsilon) = -\frac{90}{z^4} + \frac{12z_*}{z^3},$$

where the constant z_* , which depend on $\tilde{b}(0), \tilde{c}(0)$, can be derived from (146), (147).

Now we can repeat the argument of part (1) to show that $\Delta w_n(z, \varepsilon)$ is continuous in ε . It is pretty straightforward to show that $\Delta w_n(z, 0) = O(z^{-n-2})$ as $z \rightarrow \infty$ and that $\Delta w_n(z, 0)$ has asymptotic expansion in z^{-1} . Due to uniform convergence of series (54) in Theorem 2.1, we proved that $\lim_{\varepsilon \rightarrow 0} v(z, \varepsilon)$ is a solution of the truncated equation (24) that has asymptotic expansion in z^{-1} as $z \rightarrow \infty$ with the leading terms $\frac{6}{z^2} + \frac{12z_*}{z^3} + \dots$. Then (141) follows from Corollary 1.3. \square

4. Nonexistence of symmetric periodic solutions to (3). Let S be an interval of \mathbb{R} that contains the point $\omega_1 = \omega_1(0)$, and let $y(x, \varepsilon)$ be a C^α -deformation of $y(x, 0) = 6\wp_{\frac{1}{12}, g_3}(x - \omega_3) - \frac{1}{2}$, where $|g_3| < 6^{-3}$, and where ω_1, ω_3 , and g_3 are related through (9). As mentioned earlier, it is more convenient for us to consider $y(x, \varepsilon)$ as a C^α -deformation of $y(x, 0) = 6\wp_{\frac{1}{12}, g_3}(x) - \frac{1}{2}$ on the interval $S + \omega_3$. Assume that $y(x, \varepsilon)$ contains a sequence $\{y(x, \varepsilon_m)\}_1^\infty$ of solutions that are symmetrical at $x = \omega_3(\varepsilon_m)$ (after a proper translation $\beta(\varepsilon)$) and periodic (along the horizontal line $x = \omega_3(\varepsilon_m) + \mathbb{R}$) with the period $2\omega_1(\varepsilon_m)$, where $\lim_{m \rightarrow \infty} \varepsilon_m = 0$ and $\omega_1(\varepsilon_m)$ satisfy (14). Here $\omega_1(\varepsilon), \omega_3(\varepsilon)$, and $g_3(\varepsilon)$ are related through (9).

Let us show that $\beta(\varepsilon_m) = O(\varepsilon_m^\alpha)$. Indeed, since $y(x, \varepsilon)$ is a C^α -deformation of $y(x, 0)$, we know that $y'(\omega_2(\varepsilon) + \beta(\varepsilon), \varepsilon) - y'(\omega_2(\varepsilon) + \beta(\varepsilon), 0) = O(\varepsilon^\alpha)$ and $\lim_{\varepsilon \rightarrow 0} \beta(\varepsilon) =$

0. But solutions $y(x + \beta(\varepsilon_m), \varepsilon_m)$ are symmetrical at $\omega_3(\varepsilon_m)$ and have period $2\omega_1(\varepsilon_m)$, $m \in \mathbb{N}$, so $y'(\omega_2(\varepsilon_m) + \beta(\varepsilon_m), \varepsilon_m) = 0$. Thus,

$$(149) \quad y'(\omega_2(\varepsilon_m) + \beta(\varepsilon_m), 0) = O(\varepsilon_m^\alpha).$$

Since $y'(\omega_2, 0) = 0$, then $y'(\omega_2(\varepsilon_m) + \beta(\varepsilon_m), 0) = \beta(\varepsilon_m)y''(\omega_2(\varepsilon_m) + \theta(\varepsilon_m)\beta(\varepsilon_m), 0)$ by the mean value theorem, where $\theta(\varepsilon_m) \in (0, 1)$. Combining this with (149) and the fact that φ'' is separated from zero in a vicinity of ω_2 (see (64)), we obtain $\beta(\varepsilon_m) = O(\varepsilon_m^\alpha)$. Since derivatives of $y(x, 0)$ are bounded on the line $\omega_3 + \mathbb{R}$, it is easy to show that $y(x + \beta(\varepsilon_m), \varepsilon_m)$ are also C^α -deformations of $y(x, 0)$. Thus, without any loss of generality, we will consider only such deformations $y(x, \varepsilon)$ of $y(x, 0)$ that are symmetrical at $\omega_2(\varepsilon)$. That means that solutions $v(z, \varepsilon)$ to the inner equation (4), that correspond to such $y(x, \varepsilon)$, are symmetrical at $\tilde{\omega}_2(\varepsilon)$.

Let us denote by \mathcal{F}_α , $\alpha \geq 1$, the two-parameter family of solutions constructed in Theorem 2.1 that satisfy an additional assumption: condition (53) in Theorem 2.1 is replaced by (140), where $\alpha \geq 1$. The proof of Theorem 1.1 is divided into the following two steps: (1) if $v(z, \varepsilon)$ is the inner solution, symmetrical at $z = \tilde{\omega}_2(\varepsilon)$, that corresponds to a C^α -deformation $y(x, \varepsilon)$ and if $\omega_1(\varepsilon)$ satisfies (14), then $v(z, \varepsilon) \in \mathcal{F}_\alpha$; (2) if $\omega_1(\varepsilon)$ satisfies (14), then there is no symmetrical and $2\omega_1(\varepsilon)$ periodic solution $v(z, \varepsilon) \in \mathcal{F}_\alpha$. Here and below we use $\tilde{\omega}_j$ to denote $\tilde{\omega}_j(\varepsilon)$, $j = 1, 2, 3$, wherever such notation is clear. Additionally, we prove that the BVP

$$(150) \quad 2v'''v' - v''^2 + (1 - \varepsilon^2)v'^2 - \varepsilon^2v^2 - \frac{2}{3}v^3 = \varepsilon^6C(\varepsilon),$$

$$v'(\tilde{\omega}_3, \varepsilon) = v'(\tilde{\omega}_2, \varepsilon) = v'''(\tilde{\omega}_2, \varepsilon) = 0$$

has a unique solution in \mathcal{F}_α .

If $y(x, \varepsilon)$ is a C^α -deformation of solution (16), then the constant of motion $C(\varepsilon)$ of $y(x, \varepsilon)$, given by (11), satisfies (13). Then the corresponding inner solution $v(z, \varepsilon)$ satisfies (150) and the two latter boundary conditions of (150). It is easy to see that $v(z, \varepsilon)$ is periodic with the period $2\omega_1(\varepsilon)$ and symmetrical if and only if, additionally, $v'(\tilde{\omega}_3, \varepsilon) = v'''(\tilde{\omega}_3, \varepsilon) = 0$. In other words, a solution $v(z, \varepsilon)$ to the BVP (150) is symmetric and periodic with the period $2\omega_1(\varepsilon)$ if and only if $v'''(\tilde{\omega}_3, \varepsilon) = 0$.

4.1. Solution of BVP (150). Here we prove existence and uniqueness of solution to BVP (150) within the family \mathcal{F}_α using the implicit function theorem.

If $v(z, \varepsilon)$ is a solution to the BVP (150), then the integral of motion (150) evaluated at the points $z = \tilde{\omega}_2, z = \tilde{\omega}_3$ becomes

$$(151) \quad \varepsilon^2v^2(\tilde{\omega}_j) + \frac{2}{3}v^3(\tilde{\omega}_j) + v''^2(\tilde{\omega}_j) + \varepsilon^6C(\varepsilon) = 0,$$

where $j = 2, 3$. We want to show that for sufficiently small ε there exists a unique $v \in \mathcal{F}_\alpha$ satisfying (151). To this end, we can restrict our attention on the triangle $P \subset T_{z_1}$, which was introduced in section 3.4.

According to (34), (49),

$$(152) \quad v = 6q^2 + \tilde{v} = 6q^2 + 6\varepsilon^\alpha \tilde{c}q^2 + \mathcal{I}_2 \left[\varepsilon^{6+\alpha} \tilde{b}v_{22} + \mathcal{I}_1 f(q) + \tilde{w} \right],$$

where

$$(153) \quad \tilde{w} = -12c\mathcal{I}_1q^2\mathcal{I}_2f(q) + \sum_{n=2}^{\infty} \Delta w_n.$$

It is easy to see that $v(z, \varepsilon) = O(\varepsilon^2)$ uniformly in P as $\varepsilon \rightarrow 0$. According to (106), (140), and Remark 3.8, we have $w(z, \varepsilon) = O(\varepsilon^{3+\tilde{\alpha}})$ uniformly in P , where $\tilde{\alpha} = \min\{1, \alpha - 1\}$ as $\varepsilon \rightarrow 0$. Then, according to (152) and Lemma 2.4,

$$(154) \quad \tilde{v} \sim O(\varepsilon^{3+\tilde{\alpha}}) \quad \text{and} \quad \tilde{w} \sim O(\varepsilon^{4+2\tilde{\alpha}})$$

uniformly in P as $\varepsilon \rightarrow 0$. Now combining (154) and (39), we see that

$$(155) \quad w'' = \varepsilon^2 w + [\mathcal{I}_2 w]^2 + 12cq^2 \mathcal{I}_2 w + f(q) \sim O(\varepsilon^{5+\tilde{\alpha}}) \quad \text{as} \quad \varepsilon \rightarrow 0,$$

uniformly on P . Then

$$(156) \quad w'(z, \varepsilon) = \int_{\tilde{\omega}_2}^z w''(t, \varepsilon) dt = O(\varepsilon^{4+\tilde{\alpha}})$$

uniformly in $z \in P$. To estimate v'' in (151) we need the following statement.

PROPOSITION 4.1. *Let $g(z, \varepsilon)$ be a differentiable in $z \in T_{z_1}$ function. If there exists a constant $M > 0$ such that for all small $\varepsilon > 0$ both $|g(z, \varepsilon)|$ and $|g'(z, \varepsilon)|$ are bounded by M in T_{z_1} , then there exists some $\delta_0 > 0$ such that*

$$(157) \quad \mathcal{I}_2 g(z, \varepsilon) = g(z, \varepsilon) - \mathcal{I}_2 g''(z, \varepsilon) + O\left(e^{-\frac{\delta_0}{\varepsilon} \rho(z)}\right)$$

as $\varepsilon \rightarrow 0$ uniformly in P , where $\rho(z) = \min\{|z - (\tilde{\omega}_1 + 3\tilde{\omega}_3)|, |z - (\tilde{\omega}_1 - \tilde{\omega}_3)|\}$.

Proof. Integrating by parts $\mathcal{I}_2 g(z, \varepsilon)$, expressed by (76) and (83), twice, we obtain

$$(158) \quad \begin{aligned} \mathcal{I}_2 g(z, \varepsilon) &= g(z, \varepsilon) - \mathcal{I}_2 g''(z, \varepsilon) - \frac{1}{2} e^{i(z-\xi)} [g(\xi, \varepsilon) - ig'(\xi, \varepsilon)] \Big|_{\xi=\tilde{\omega}_1-\tilde{\omega}_3} \\ &\quad - \frac{1}{2} e^{-i(z-\xi)} [g(\xi, \varepsilon) + ig'(\xi, \varepsilon)] \Big|_{\xi=\tilde{\omega}_1+3\tilde{\omega}_3}. \end{aligned}$$

The last two terms of (158) are exponentially small in ε^{-1} according to the construction of T_{z_1} . \square

Under the assumptions of Proposition 4.1,

$$(159) \quad D^2 \mathcal{I}_2 g = g - \mathcal{I}_2 g = \mathcal{I}_2 D^2 g + O\left(e^{-\frac{\delta_0}{\varepsilon} \rho(z)}\right).$$

According to (156), Proposition 4.1 is applicable for $g(z, \varepsilon) = w(z, \varepsilon)$. Thus,

$$(160) \quad v'' = 6c(\varepsilon)(q^2)'' + \mathcal{I}_2 w'' = O(\varepsilon^4) \quad \text{as} \quad \varepsilon \rightarrow 0,$$

uniformly on $[\tilde{\omega}_3, \tilde{\omega}_2]$.

Note that $y = \frac{6q^2}{\varepsilon^2}$ is the solution of the unperturbed equation (2) with the constant of motion $J(\varepsilon)$; see (37). Thus,

$$(161) \quad \left(\frac{6q^2}{\varepsilon^2}\right)^2 + \frac{2}{3} \left(\frac{6q^2}{\varepsilon^2}\right)^3 + J(\varepsilon) = 0$$

at $x = \varepsilon z = \omega_j$, $j = 1, 2, 3$. Notice that (9) and (14) imply that

$$(162) \quad J(\varepsilon) = C + \varepsilon^\alpha \tilde{J}(\varepsilon),$$

where $\tilde{J}(\varepsilon)$ is a continuous function. Now substituting $v = 6q^2 + \tilde{v}$ into (151), dividing both sides by $\varepsilon^{7+\tilde{\alpha}}$, and taking into account (161), we obtain

$$(163) \quad 2 \left(\frac{6q^2}{\varepsilon^2} \right) \check{v} + \varepsilon^{1+\tilde{\alpha}} \check{v}^2 + 2 \left(\frac{6q^2}{\varepsilon^2} \right)^2 \check{v} + 2\varepsilon^{1+\tilde{\alpha}} \left(\frac{6q^2}{\varepsilon^2} \right) \check{v}^2 + \frac{2}{3} \varepsilon^{2+2\tilde{\alpha}} \check{v}^3 + \frac{v'^2}{\varepsilon^{7+\tilde{\alpha}}} + \tilde{C}(\varepsilon) - \tilde{J}(\varepsilon) = 0$$

at $z = \tilde{\omega}_2, \tilde{\omega}_3$, where

$$(164) \quad \check{v}(z, \varepsilon) = \frac{\tilde{v}(z, \varepsilon)}{\varepsilon^{3+\tilde{\alpha}}} = \tilde{c}(\varepsilon) \left(\frac{6q^2}{\varepsilon^2} \right) + \varepsilon^4 \tilde{b}(\varepsilon) \mathcal{I}_2 v_{22} + \mathcal{I}_2 \frac{\mathcal{I}_1 f(q) + \tilde{w}}{\varepsilon^{3+\tilde{\alpha}}}.$$

Equation (163) evaluated at the points $z = \tilde{\omega}_2, \tilde{\omega}_3$ forms a system that we denote by

$$(165) \quad F(\varepsilon, \tilde{b}, \tilde{c}) = 0.$$

LEMMA 4.2. *There exist some $\varepsilon_1 \in (0, \varepsilon_0]$ and functions $\tilde{b}(\varepsilon), \tilde{c}(\varepsilon)$, continuous on $[0, \varepsilon_1]$, such that the system (165) holds identically on $[0, \varepsilon_1]$. Moreover, the solution $\tilde{b}(\varepsilon), \tilde{c}(\varepsilon)$ to (165) is unique.*

Proof. The proof is based on the implicit function theorem (see, for example, [MB, p. 122]). According to this theorem, we have to show that (a) $F(0, \tilde{b}(0), \tilde{c}(0)) = 0$ for some $\tilde{b}(0), \tilde{c}(0)$; (b) matrix

$$(166) \quad \text{Col} \left(\frac{\partial F}{\partial \tilde{b}}, \frac{\partial F}{\partial \tilde{c}} \right) \Big|_{\varepsilon=0}$$

is not singular; and (c) F and $\frac{\partial F}{\partial \varepsilon}, \frac{\partial F}{\partial \tilde{b}}$ are continuous in all variables in a vicinity of $(0, \tilde{b}(0), \tilde{c}(0))$.

According to (156), we can apply Proposition 4.1 to $\check{v}(z, \varepsilon) - \tilde{c}(\varepsilon) \frac{6q^2}{\varepsilon^2} = \frac{\mathcal{I}_2 w}{\varepsilon^{3+\tilde{\alpha}}}$. Then

$$(167) \quad \lim_{\varepsilon \rightarrow 0} \check{v}(z, \varepsilon) = \tilde{c}(0) \left(6\wp(\varepsilon z) - \frac{1}{2} \right) + \tilde{b}(0) \lim_{\varepsilon \rightarrow 0} \varepsilon^4 v_{22}(z, \varepsilon) + \lim_{\varepsilon \rightarrow 0} \frac{\mathcal{I}_1 f(q)}{\varepsilon^{3+\tilde{\alpha}}} + \lim_{\varepsilon \rightarrow 0} \frac{\tilde{w}}{\varepsilon^{3+\tilde{\alpha}}}$$

uniformly on $[\tilde{\omega}_3, \tilde{\omega}_2]$. Due to Remark 3.8, the last limit in (167) is zero.

Using (63), (69), (47), and the fact that $v_1(\tilde{\omega}_j) = 0$, we have

$$(168) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^4 v_{22}(\tilde{\omega}_j, \varepsilon) = -\frac{1}{12(e_j - e_i)(e_j - e_k)}.$$

Note that $\mathcal{I}_1 f(q)(\tilde{\omega}_2) = 0$. To calculate $\lim_{\varepsilon \rightarrow 0} \varepsilon^{-3-\tilde{\alpha}} \mathcal{I}_1 f(q)(\tilde{\omega}_3)$, we notice that

$$(169) \quad \wp'(x - \omega_j) = \left(6e_j^2 - \frac{g_2}{2} \right) (x - \omega_j) + O((x - \omega_j)^2) \\ = 2(e_i - e_j)(e_k - e_j)(x - \omega_j) + O((x - \omega_j)^2)$$

as $x \rightarrow \omega_j, j = 1, 2, 3$. It follows then from (126) that

$$(170) \quad \varepsilon^{-3} \mathcal{I}_1 f(q)(\omega_3, \varepsilon) = 12c(\varepsilon)(e_1 - e_3)(e_2 - e_3)G_3 \\ = \frac{c(\varepsilon)}{e_1 - e_3} \left\{ -\frac{c(\varepsilon) - 1 - 4\varepsilon^2}{4\varepsilon} \left(e_1 + \frac{1}{6} \right) - \frac{5\varepsilon}{2} \left(\frac{9}{4}e_3 + \frac{3}{2}e_2 + 36e_2^3 - \frac{1}{3} \right) + \varepsilon C(\varepsilon) \right\}.$$

Using (126), we obtain

$$(171) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-3-\tilde{\alpha}} \mathcal{I}_1 f(q)(\tilde{\omega}_3) = -\frac{\tilde{c}(\varepsilon)(e_1 + \frac{1}{6})}{4(e_1 - e_3)} - \delta_{\tilde{\alpha},1} \frac{e_1 + \frac{1}{6} + \frac{5}{2}(\frac{9}{4}e_3 + \frac{3}{2}e_2 + 36e_2^3 - \frac{1}{3}) + 36g_3 + \frac{1}{6}}{e_1 - e_3}$$

at $z = \tilde{\omega}_3$, where $\delta_{\tilde{\alpha},1} = 1$ if $\tilde{\alpha} = 1$ and $\delta_{\tilde{\alpha},1} = 0$ otherwise. Combining (167)–(171), we obtain

$$(172) \quad \lim_{\varepsilon \rightarrow 0} \check{v}(\tilde{\omega}_2) = \left(6e_2 - \frac{1}{2}\right) \tilde{c}(0) + \frac{\tilde{b}(0)}{12(e_1 - e_2)(e_2 - e_3)},$$

$$\lim_{\varepsilon \rightarrow 0} \check{v}(\tilde{\omega}_3) = \left[\left(6e_3 - \frac{1}{2}\right) - \frac{e_1 + \frac{1}{6}}{4(e_1 - e_3)}\right] \tilde{c}(0) - \frac{\tilde{b}(0)}{12(e_1 - e_3)(e_2 - e_3)} - \delta_{\tilde{\alpha},1} \tilde{D},$$

where \tilde{D} denotes the last fraction in (171).

According to (154) and (160), in the limit $\varepsilon \rightarrow 0$ equations (163) become

$$(173) \quad 2 \left[36\wp^2(\omega_j) - \frac{1}{4}\right] \lim_{\varepsilon \rightarrow 0} \check{v}(\tilde{\omega}_j) + \tilde{C}(0) - \tilde{J}(0) = 0,$$

where $j = 2, 3$. Clearly, (172)–(173) form a linear system of equations for the unknowns $\tilde{b}(0), \tilde{c}(0)$. Note that the first factor in (173) is different from zero since, otherwise, according to (64), we would have either $C = 0$ or $C = -\frac{1}{3}$, which contradicts our assumption $C \in (-\frac{1}{3}, 0)$. After some algebra, we calculate the determinant of the matrix of linear system (172)–(173) as

$$(174) \quad \begin{vmatrix} 6e_2 - \frac{1}{2} & \frac{1}{12(e_1 - e_2)(e_2 - e_3)} \\ 6e_3 - \frac{1}{2} - \frac{e_1 + \frac{1}{6}}{4(e_1 - e_3)} & -\frac{1}{12(e_1 - e_3)(e_2 - e_3)} \end{vmatrix} = \frac{1}{3\sqrt{\Delta}} \begin{vmatrix} (6e_2 - \frac{1}{2})(e_1 - e_2) & 1 \\ (6e_3 - \frac{1}{2})(e_1 - e_3) - \frac{e_1 + \frac{1}{6}}{4} & -1 \end{vmatrix} = \frac{7(e_1 + \frac{1}{6})}{12\sqrt{\Delta}},$$

where Δ , defined by (67), is different from zero. This determinant is different from zero since $e_1 > 0$. Thus, we have established parts (a) and (b) of the lemma; i.e., we have established that $F(0, \tilde{b}(0), \tilde{c}(0)) = 0$ at $z = \tilde{\omega}_2, \tilde{\omega}_3$ for some $\tilde{b}(0), \tilde{c}(0)$ and that matrix (166) is nonsingular. Continuity of F follows from Theorem 3.13 and the fact that iterations $\Delta w_n(z, \varepsilon)$ in the solution (54) are polynomials in $\tilde{b}(0), \tilde{c}(0)$ of degree not exceeding n . Estimates of Theorem 2.1 can be readily adjusted to prove convergence of the series $\sum_{n=1}^{\infty} \frac{\partial \Delta w_n(z, \varepsilon)}{\partial \tilde{c}}$ and $\sum_{n=1}^{\infty} \frac{\partial \Delta w_n(z, \varepsilon)}{\partial \tilde{b}}$ at $z = \tilde{\omega}_2, z = \tilde{\omega}_3$. Part (c) and the whole proof are completed. \square

THEOREM 4.3. *If $C(\varepsilon)$ and $\omega_1(\varepsilon)$ satisfy (13) and (14), respectively, where $\alpha \geq 1$, then the BVP (150) has a unique solution in \mathcal{F}_α .*

Proof. According to Lemma 4.2, it is sufficient to show that (151) at $z = \tilde{\omega}_3$ and the rescaled integral of motion (150) imply $v'(\tilde{\omega}_3) = 0$. Indeed, addition of these two equations yields

$$(175) \quad v'(\tilde{\omega}_3)[2v'''(\tilde{\omega}_3) + (1 - \varepsilon^2)v'(\tilde{\omega}_3)] = 0.$$

That means that $v'(\tilde{\omega}_3) = 0$ or $2v'''(\tilde{\omega}_3) + (1 - \varepsilon^2)v'(\tilde{\omega}_3) = 0$. According to (34), and since operators D and $D^2 + 1$ commute, this implies $u'(\tilde{\omega}_3) = 0$ or $2w'(\tilde{\omega}_3) = (1 + \varepsilon^2)u'(\tilde{\omega}_3)$. If $u'(\tilde{\omega}_3) = 0$, then $v'(\tilde{\omega}_3) = 0$, and the proof is completed. The assumption $u'(\tilde{\omega}_3) \neq 0$ leads to a contradiction. Indeed, on one hand, we have

$$(176) \quad u'(\tilde{\omega}_3) = \frac{2}{1 + \varepsilon^2}w'(\tilde{\omega}_3),$$

whereas, on the other hand, according to Proposition 4.1,

$$(177) \quad u'(\tilde{\omega}_3) = \mathcal{I}_2 w'|_{z=\tilde{\omega}_3} = w'(\tilde{\omega}_3) - \mathcal{I}_2 w'''|_{z=\tilde{\omega}_3} + O\left(e^{-\frac{\delta_0}{\varepsilon}}\right).$$

Proposition 4.1 is applicable to $\mathcal{I}_2 w'$ according to (39) and (156). Moreover, differentiating (39), we obtain $w''' = O(\varepsilon^{6+\alpha})$ uniformly on $[\tilde{\omega}_3, \tilde{\omega}_2]$. Thus, (177) contradicts (176). \square

4.2. Proof of Theorem 1.1.

Proof. As mentioned in the beginning of section 4, our proof consists of two steps. Step (1) is to show that the inner solution $v(z, \varepsilon) = \varepsilon^2 y(\varepsilon z, \varepsilon)$, corresponding to a C^α -deformation $y(x, \varepsilon)$ of the periodic solution $y(x, 0)$, is an \mathcal{F}_α solution. Let us consider first the case $n = 1$.

As discussed above, we can assume without any loss of generality that $v(z, \varepsilon_m)$ is symmetrical with respect to $z = \tilde{\omega}_3(\varepsilon_m)$ and $z = \tilde{\omega}_2(\varepsilon_m)$, where $g_3(\varepsilon_m)$ and $\omega_3(\varepsilon_m)$ are defined by $\omega_1(\varepsilon_m)$ through (9). Then $v'(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m) = v'''(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m) = 0$, so that $v(z, \varepsilon_m)$ is defined by initial conditions $v(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m)$ and $v''(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m)$. In order to show that there is a solution in \mathcal{F}_α with the abovementioned initial conditions at $z = \tilde{\omega}_2(\varepsilon_m)$, we first establish that

$$(178) \quad v(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m) - \varepsilon_m^2 \left(6e_2(0) - \frac{1}{2}\right) = \varepsilon_m^{2+\alpha} \tilde{v}(\varepsilon_m),$$

$$v''(\tilde{\omega}_2(\varepsilon_m), \varepsilon_m) - \varepsilon_m^4 \left(3e_2^2(0) - \frac{1}{48}\right) = \varepsilon_m^{4+\alpha} \tilde{v}''(\varepsilon_m),$$

where $\tilde{v}(\varepsilon)$, $\tilde{v}''(\varepsilon)$ are continuous functions. Indeed, according to our assumptions and taking into account (14),

$$(179) \quad y(\omega_2(\varepsilon_m), \varepsilon_m) - y(\omega_2, 0)$$

$$= [y(\omega_2(\varepsilon_m), \varepsilon_m) - y(\omega_2(\varepsilon_m), 0)] + [y(\omega_2(\varepsilon_m), 0) - y(\omega_2, 0)]$$

$$= \varepsilon_m^\alpha \tilde{y}(\omega_2(\varepsilon_m), \varepsilon_m) + \frac{1}{2} \varepsilon_m^\alpha \tau^2(\varepsilon_m) \wp''(\omega_2(0) + \theta(\omega_1(\varepsilon_m) - \omega_1)),$$

where $\theta \in (0, 1)$. Here we used the mean value theorem and the fact that $\wp'(\omega_2(0)) = 0$. A similar estimate holds for $y''(\omega_2(\varepsilon_m), \varepsilon_m) - y''(\omega_2, 0)$. Then, (178) follows from (179).

According to (164), equations (178) in the leading order can be written as

$$\begin{aligned}
 (180) \quad & \tilde{c}(\varepsilon) \left(\frac{6q^2}{\varepsilon^2} \right) + \varepsilon^4 \tilde{b}(\varepsilon) \mathcal{I}_2 v_{22} + \mathcal{I}_2 \frac{\mathcal{I}_1 f(q) + \tilde{w}}{\varepsilon^{3+\alpha}} \\
 & = \tilde{v}(\varepsilon) - \frac{1}{2} \tau^2(\varepsilon) \wp''(\omega_2(0) + \theta(\omega_1(\varepsilon) - \omega_1(0))), \\
 & \tilde{c}(\varepsilon) D_z^2 \left(\frac{6q^2}{\varepsilon^2} \right) + \varepsilon^4 \tilde{b}(\varepsilon) D_z^2 \mathcal{I}_2 v_{22} + D_z^2 \mathcal{I}_2 \frac{\mathcal{I}_1 f(q) + \tilde{w}}{\varepsilon^{3+\alpha}} \\
 & = \tilde{v}''(\varepsilon) - \frac{3}{2} (e_2(\varepsilon) + e_2(0)) \tau^2(\varepsilon) \wp''(\omega_2(0) + \theta(\omega_1(\varepsilon) - \omega_1(0))).
 \end{aligned}$$

According to (167) and Proposition 4.1, in the limit $\varepsilon \rightarrow 0$ equations (180) become

$$\begin{aligned}
 (181) \quad & \tilde{c}(0) \lim_{\varepsilon \rightarrow 0} \frac{6q^2}{\varepsilon^2} + \tilde{b}(0) \lim_{\varepsilon \rightarrow 0} \varepsilon^4 v_{22}(\tilde{\omega}_2(\varepsilon), \varepsilon) = \tilde{v}(0) - \frac{1}{2} \tau^2 \wp''(\omega_2(0)), \\
 & \tilde{c}(0) \lim_{\varepsilon \rightarrow 0} \left[\frac{6q^2}{\varepsilon^2} + \frac{36q^4}{\varepsilon^4} \right] + \tilde{b}(0) \lim_{\varepsilon \rightarrow 0} \left\{ \varepsilon^4 v_{22}(\tilde{\omega}_2(\varepsilon), \varepsilon) \left[1 + \frac{12q^2}{\varepsilon^2} \right] \right\} = \tilde{v}''(0) \\
 & \quad - 3e_2(0) \tau^2(0) \wp''(\omega_2(0)).
 \end{aligned}$$

Here we used the fact that $\frac{6q^2}{\varepsilon^2}$ and v_{22} satisfy differential equations (2) and (44), respectively. Calculation of the determinant of the latter system yields

$$\begin{aligned}
 (182) \quad & \lim_{\varepsilon \rightarrow 0} \begin{vmatrix} \frac{6q^2}{\varepsilon^2} & \varepsilon^4 v_{22} \\ \frac{6q^2}{\varepsilon^2} + \frac{36q^4}{\varepsilon^4} & \left[1 + \frac{12q^2}{\varepsilon^2} \right] \varepsilon^4 v_{22} \end{vmatrix} = \lim_{\varepsilon \rightarrow 0} \begin{vmatrix} \frac{6q^2}{\varepsilon^2} & \varepsilon^4 v_{22} \\ \frac{36q^4}{\varepsilon^4} & 12q^2 \varepsilon^2 v_{22} \end{vmatrix} \\
 & = 36 \lim_{\varepsilon \rightarrow 0} q^4 v_{22} = \frac{3(e_2(0) - \frac{1}{2})^2}{(e_1(0) - e_2(0))(e_2(0) - e_3(0))} \neq 0.
 \end{aligned}$$

Following the arguments of Lemma 4.2, we can now use the implicit function theorem to show that the system (180) has a unique solution $\tilde{c}(\varepsilon), \tilde{b}(\varepsilon)$, where $\tilde{c}(\varepsilon), \tilde{b}(\varepsilon)$ are continuous functions. Thus, $v(z, \varepsilon) \in \mathcal{F}_\alpha$. To prove step (1) for the case for general n we simply have to consider the region \mathcal{P} from Remark 3.10 instead of the triangle P and the point $n\tilde{\omega}_1(\varepsilon) + \tilde{\omega}_3(\varepsilon)$ instead of $\tilde{\omega}_2(\varepsilon)$.

To prove step (2), we assume that there is a family of solutions $v(z, \varepsilon_m) \subset \mathcal{F}_\alpha$ to (4), where $\omega_1(\varepsilon)$ satisfies (14), that is symmetrical at $\tilde{\omega}_3(\varepsilon_m)$. However, according to Theorem 3.13,

$$(183) \quad \lim_{m \rightarrow \infty} v(z, \varepsilon_m) = v_+(z)$$

uniformly in z on any closed segment of the imaginary axis that belongs to T_{z_1} . Due to the symmetry with respect to the imaginary axis, we also have

$$(184) \quad \lim_{m \rightarrow \infty} v(z, \varepsilon_m) = v_-(z)$$

uniformly in z on the same segment of the imaginary axis. But, according to Corollary 1.3, $v_+(z) \neq v_-(z)$. The obtained contradiction proves nonexistence of a family of symmetric periodic solutions $v(z, \varepsilon_m) \subset \mathcal{F}_\alpha$, $\alpha > 1$, where $\varepsilon_m \rightarrow 0$. The same result, according to Theorem 3.13, part (2), holds for the case $\alpha = 1$. \square

REFERENCES

- [AM] C. J. AMICK AND J. B. MCLEOD, *A singular perturbation problem in water waves*, Stability Appl. Anal. Contin. Media, 1 (1991), pp. 127–148.
- [BPBA] R. E. BEARDMORE, M. A. PELETIER, C. J. BUDD, AND M. AHMER WADEE, *Bifurcations of periodic solutions satisfying the zero-Hamiltonian constraint in reversible differential equations*, SIAM J. Math. Anal., 36 (2005), pp. 1461–1488.
- [Eck] W. ECKHAUS, *Singular perturbations of homoclinic orbits in \mathbb{R}^4* , SIAM J. Math. Anal., 23 (1992), pp. 1269–1290.
- [GR] I. S. GRADSTEYN AND I. W. RYZHIK, *Tables of Integrals, Series and Products*, Academic Press, New York, 1965.
- [GJ] R. GRIMSHAW AND N. JOSHI, *Weakly nonlocal solitary waves in a singularly perturbed Korteweg–de Vries equation*, SIAM J. Appl. Math., 55 (1995), pp. 124–135.
- [Ha] H. HANCOCK, *Theory of Elliptic Functions*, Dover, New York, 1958.
- [HM] J. M. HAMMERSLEY AND G. MAZZARINO, *Computational aspects of some autonomous differential equations*, Proc. Roy. Soc. London Ser. A, 424 (1989), pp. 19–37.
- [IL1] G. IOOSS AND E. LOMBARDI, *Normal forms with exponentially small remainder: Application to homoclinic connections for the reversible $0^{2+}\omega$ resonance*, C. R. Math. Acad. Sci. Paris, 339 (2004), pp. 831–838.
- [IL2] G. IOOSS AND E. LOMBARDI, *Normal forms with exponentially small remainder for analytic vector fields*, J. Differential Equations, 212 (2005), pp. 1–61.
- [KS] M. D. KRUSKAL AND H. SEGUR, *Asymptotics beyond all orders in a model of crystal growth*, Stud. Appl. Math., 85 (1991), pp. 129–181.
- [Lo] E. LOMBARDI, *Oscillatory Integrals and Phenomena Beyond All Algebraic Orders. With Applications to Homoclinic Orbits in Reversible Systems*, Lecture Notes in Math. 1741, Springer-Verlag, Berlin, 2000.
- [MB] E. J. MCSHANE AND T. A. BOTTS, *Real Analysis*, Van Nostrand Company, New York, 1959.
- [PT] L. A. PELETIER AND W. C. TROY, *Spatial Patterns: Higher Order Models in Physics and Mechanics*, Birkhäuser Boston, Boston, 2001.
- [PRG] Y. POMEAU, A. RAMANI, AND B. GRAMMATICOS, *Structural stability of the Korteweg–de Vries solutions under a singular perturbation*, Phys. D, 31 (1988), pp. 127–134.
- [STL] H. SEGUR, S. TANVEER, AND H. LEVINE, EDS., *Asymptotics beyond All Orders*, NATO ASI Ser. B. Phys. 284, Plenum, New York, 1991.
- [To2] A. TOVBIS, *Asymptotics beyond all orders and analytic properties of inverse Laplace transforms of solutions*, Comm. Math. Phys., 163 (1994), pp. 245–255.
- [To4] A. TOVBIS, *Breaking homoclinic connections for a singularly perturbed differential equation and the Stokes phenomenon*, Stud. Appl. Math., 104 (2000), pp. 353–386.
- [To5] A. TOVBIS, *On approximation of stable and unstable manifolds and the Stokes phenomenon*, in Nonlinear PDE’s, Dynamics and Continuum Physics, Contemp. Math. 255, AMS, Providence, RI, 2000, pp. 199–228.
- [TP] A. TOVBIS AND D. PELINOVSKY, *Exact conditions for existence of homoclinic orbits in the fifth-order KdV model*, Nonlinearity, 19 (2006), pp. 2277–2312.
- [WW] E. T. WHITTAKER AND G. N. WATSON, *A Course of Modern Analysis*, Cambridge University Press, Cambridge, UK, 1927.
- [Wa] W. WASOW, *Asymptotic Expansions for Ordinary Differential Equations*, Robert E. Krieger Publishing, Huntington, NY, 1976.

RISE OF CORRELATIONS OF TRANSFORMATION STRAINS IN RANDOM POLYCRYSTALS*

LEONID BERLYAND[†], OSCAR BRUNO[‡], AND ALEXEI NOVIKOV[§]

Abstract. We investigate the statistics of the transformation strains that arise in random martensitic polycrystals as boundary conditions cause its component crystallites to undergo martensitic phase transitions. In our laminated polycrystal model the orientation of the n grains (crystallites) is given by an uncorrelated random array of the orientation angles θ_i , $i = 1, \dots, n$. Under imposed boundary conditions the polycrystal grains may undergo a martensitic transformation. The associated transformation strains ε_i , $i = 1, \dots, n$ depend on the array of orientation angles, and they can be obtained as a solution to a nonlinear optimization problem. While the random variables θ_i , $i = 1, \dots, n$ are uncorrelated, the random variables ε_i , $i = 1, \dots, n$ may be correlated. This issue is central in our considerations. We investigate it in following three different scaling limits: (i) Infinitely long grains (laminated polycrystal of height $L = \infty$); (ii) Grains of finite but large height ($L \gg 1$); and (iii) Chain of short grains ($L = l_0/(2n)$, $l_0 \ll 1$). With references to de Finetti's theorem, Riesz' rearrangement inequality, and near neighbor approximations, our analyses establish that under the scaling limits (i), (ii), and (iii) the arrays of transformation strains arising from given boundary conditions exhibit no correlations, long-range correlations, and exponentially decaying short-range correlations, respectively.

Key words. polycrystals, misfit, phase transitions, correlations, De Finetti's theorem, Riesz' rearrangement inequality

AMS subject classifications. 35J20, 74N15, 82B44

DOI. 10.1137/070679685

1. Introduction. We investigate the statistics of the transformation strains (misfits) that arise in random martensitic polycrystals as boundary conditions cause its component crystallites to undergo solid-to-solid (martensitic) phase transitions. Martensitic transformations are shape-deforming phase transitions that can be induced in certain alloys as a result of changes in the imposed strains, stresses, or temperatures. These transitions occur when a crystalline solid transforms between its parent phase (austenite) and any of a number of variants of the product phase (martensite). We focus on a setting that, while sufficiently simple to allow for a complete analytical treatment, provides significant insights on the problem: We study laminated polycrystals that consist of sequences of n of grains of rectangular cross-section—of base $1/n$ and height $L = L(n)$, so that a complete polycrystal is an infinite parallelepiped with rectangular cross-section of base 1 and height L . The goal of this work is to provide a rigorous probabilistic theory for the misfit statistics in such polycrystals and, in particular, to provide a rationale for the approximations implicit in

*Received by the editors January 9, 2007; accepted for publication (in revised form) June 24, 2008; published electronically November 19, 2008.

<http://www.siam.org/journals/sima/40-4/67968.html>

[†]Department of Mathematics, Pennsylvania State University, University Park, PA 16802 (berlyand@math.psu.edu). This author was supported by NSF grants DMS-0204637 and DMS-0708324.

[‡]California Institute of Technology, Applied & Computational Mathematics, 1200 E. California Boulevard, MC 217-50, Pasadena, CA 91125 (bruno@acm.caltech.edu). This author acknowledges support from NSF grant DMS-0408040, AFOSR grant FA9550-05-1-0006, and NRC-JPL award 1263315.

[§]Department of Mathematics, Pennsylvania State University, University Park, PA 16802 (anovikov@math.psu.edu). This author was supported by NSF grant DMS-0604600.

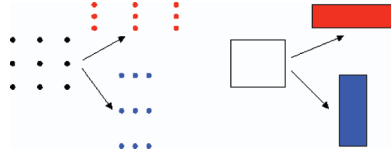


FIG. 1.1. A reference crystallite undergoes stress-free transformations: Atomic view (left) and macroscopic view (right).

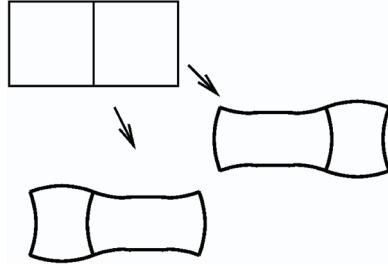


FIG. 1.2. One of two grains undergoes a stress-free transformation.

the numerical algorithms [6, 7, 8] for polycrystalline phase transitions in two- and three-dimensional space.

The microstructure in a laminated polycrystal is described by a sequence of the orientation angles θ_i , $i = 1, \dots, n$; θ_i represents the orientation of the two-dimensional lattice structure in the i th grain. We assume θ_i is a sequence of n independent identically distributed (i.i.d.) random variables. The transformation in the i th grain gives rise to a strain tensor, the transformation strain ε_i^T , ($i = 1, \dots, n$), which is constant, and it takes one of three *admissible* values: No deformation (the original square lattice remains square), or deformation into one of two rectangular crystalline lattices parallel to the original square lattice. The phase transition in the polycrystal gives rise to a sequence of transformation strains ε_i^T , $i = 1, \dots, n$ obtained by the minimization of the elastic *misfit energy* among all admissible configurations.

We briefly explain the concept of *misfit* using a simple example of a polycrystal with two grains. Assume a rectangular single-crystalline grain, considered separately, can undergo a stress-free (two-dimensional version of the) cubic-to-orthorhombic phase deformation into shapes depicted in Figure 1.1. A polycrystal with two square grains can undergo deformations as depicted in Figure 1.2. The elastic energy of the former transformation is zero, because it is stress-free. In contrast, the latter transformation requires some elastic *misfit energy* that arises because when two crystallites are combined in a polycrystal, their boundaries must remain coherent after the transformation. In general, minimization of misfit energy leads to interactions amongst all of the grains in a polycrystal. Our probabilistic setup allows us to provide a rigorous description of this phenomenology.

The main results of this paper characterize the probability distributions of the random variables ε_i^T that arise as minimizers of the overall elastic energy for a given i.i.d. distribution of the angle sequence θ_i . Such results are provided in three different cases according to whether the grains are (1) infinitely long ($L = \infty$), (2) of finite but large height ($L = L \gg 1$), and (3) short height ($L = l_0/(2n)$, $l_0 \ll 1$). In case (1) our treatment applies to arbitrary i.i.d. probability measures ρ defining the distribution of angles, and in cases (2) and (3), in turn, we restrict consideration to

i.i.d. distribution of angles with Bernoulli probability measures ρ . Our main results can be briefly described as follows:

1. *Infinitely long grains.* *Theorem 5.2.* For an arbitrary i.i.d. distribution of angles θ_i , $i = 1, 2, \dots, n$, under certain technical assumptions, in the limit $n \rightarrow \infty$, the transformation strains ε_i^T , $i = 1, 2, \dots, n$ are also i.i.d. with probability measure μ , where the measure μ is the minimizer of a certain functional ((5.8) below). In particular, in the case of infinitely long grains there are no correlations between transformation strains of any two grains.
2. *Long finite grains.* $L \gg 1$. *Theorem 6.4.* If θ_i , $i = 1, \dots, n$ are Bernoulli random variables (4.3), then in the limit $n \rightarrow \infty$, ε_i^T , $i = 1, 2, \dots, n$ have long-range but no short-range correlations.
3. *Short grains.* $L = l_0/(2n)$, $l_0 \ll 1$. *Theorems 7.4 and 7.5.* If θ_i , $i = 1, \dots, n$ are Bernoulli random variables (4.3), then in the limit $n \rightarrow \infty$, ε_i^T , $i = 1, 2, \dots, n$ have short-range but no long-range correlations.

Results 2 and 3 can be explained as follows. The cornerstone of our study is the maximization of an integral energy functional (see (3.4) below) of the form $\int K_L(x-t)f(x)f(t)dxdt$. Its integral kernel $K_L(x)$ decays on different length scales for long and short grains. For long grains it decays on the length scale of the composite (on $O(1)$ scale), while for short grains it decays on the length scale of a grain (on $O(1/n)$ scale). Maximization with respect to this integral kernel leads to long-range and short-range correlations for long and short grains, respectively. Formally, correlations arise because grains that undergo the stress-free transformation tend to “group together” on the scale of the decay of the integral kernel. We justify this heuristic idea in the case of long grains (see section 6) by applying a randomized version of the Riesz rearrangement inequality. In the case of short grains (see section 7) we show the transforming grains group together—by applying an isoperimetric inequality.

The paper is organized as follows: After describing in section 2 our model of the polycrystal, in section 3 we solve an auxiliary linear elasticity problem, and we obtain an explicit expression for the stored elastic energy for a fixed admissible array of transformation strains. In section 4 we describe our probabilistic model. Our main results are then established in the next three sections, where the nonlinear minimization problem for a random polycrystal is solved. The cases concerning infinitely long grains, finitely long grains, and short grains are studied in sections 5, 6, and 7, respectively.

2. Formulation.

Stress-free transformation. A two-dimensional polycrystal is a collection of grains. In our model, each grain is a single crystal (a crystallite) which can undergo a shape-deforming phase transition that results in a transformation strain. An untransformed grain with a horizontal-vertical square lattice (angle $\theta = 0$) may either elongate in the horizontal direction and remain unchanged in the vertical direction (the upper-right state in Figure 1.1); it may elongate in the vertical direction and remain unchanged in the horizontal direction (the lower right state in Figure 1.1); or, finally, it may not transform at all and thus have its size unchanged (the left state in Figure 1.1). These states correspond to the transformation strains:

$$(2.1) \quad \varepsilon_0^1 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad \varepsilon_0^2 = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}, \quad \varepsilon_0^0 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}.$$

The first and the second state correspond to a nontrivial transformation. The null strain ε_0^0 corresponds to absence of transformation.

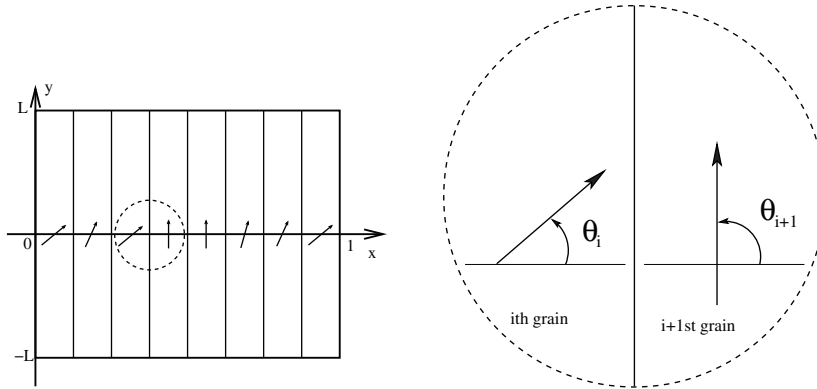


FIG. 2.1. *The laminated polycrystal.*

Mathematical model of a laminated polycrystal. Grains in a polycrystal have a varying orientation of the crystalline lattices. We consider a rectangular polycrystal $\Pi_L = [0, 1] \times [-L, L]$ (see Figure 2.1) partitioned into n vertical rectangular layers (the *grains*) of width $1/n$ and height $2L$,

$$\Pi_L = \cup_{i=1}^n \Pi_L^i, \quad \Pi_L^i = \left[\frac{i-1}{n}, \frac{i}{n} \right] \times [-L, L].$$

Each grain Π_L^i is occupied by a crystallite obtained by rotation by the *orientation angle*

$$\theta_i, \quad 0 \leq \theta_i \leq \pi/2,$$

of the reference crystallite (see Figure 2.1).

The array of crystallites' orientations is completely determined by the vector of the orientation angles

$$(2.2) \quad \boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n).$$

Using the matrix of rotation by an angle θ

$$R = R(\theta) = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix},$$

we see that the stress-free transformation strain for the grain Π_L^i must lie in the set

$$(2.3) \quad \mathcal{S}_{\theta_i} = \{ \boldsymbol{\varepsilon}^1(\theta_i), \boldsymbol{\varepsilon}^2(\theta_i), \boldsymbol{\varepsilon}^0(\theta_i) \}, \quad 0 \leq \theta_i \leq \pi/2,$$

where

$$(2.4) \quad \begin{aligned} \boldsymbol{\varepsilon}^1(\theta) &= R\boldsymbol{\varepsilon}^1(0)R^t = \begin{pmatrix} \cos^2(\theta) & \sin(\theta)\cos(\theta) \\ \sin(\theta)\cos(\theta) & \sin^2(\theta) \end{pmatrix}, \\ \boldsymbol{\varepsilon}^2(\theta) &= R\boldsymbol{\varepsilon}^2(0)R^t = \begin{pmatrix} \sin^2(\theta) & -\sin(\theta)\cos(\theta) \\ -\sin(\theta)\cos(\theta) & \cos^2(\theta) \end{pmatrix}, \\ \boldsymbol{\varepsilon}^0(\theta) &= R\boldsymbol{\varepsilon}^0(0)R^t = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}. \end{aligned}$$

The superscript t stands for the matrix transpose. The set of all sequences of strains that are admissible for some sequence of angles is denoted by $\tilde{\Omega}_n$:

$$(2.5) \quad \tilde{\Omega}_n = \{ \boldsymbol{\varepsilon}^T | \boldsymbol{\varepsilon}^T = (\boldsymbol{\varepsilon}_1^T, \boldsymbol{\varepsilon}_2^T, \boldsymbol{\varepsilon}_3^T, \dots, \boldsymbol{\varepsilon}_n^T), \boldsymbol{\varepsilon}_i^T \in \mathcal{S}_{\theta_i} \text{ for some } \theta_i \in [0, \pi/2] \}.$$

The set of all sequences of strains that are admissible for a given sequence $\boldsymbol{\theta}$ will be denoted by

$$(2.6) \quad \tilde{\Omega}_n(\boldsymbol{\theta}) = \left\{ \boldsymbol{\varepsilon}^T | \boldsymbol{\varepsilon}^T \in \tilde{\Omega}_n \text{ such that } \boldsymbol{\varepsilon}_i^T \in \mathcal{S}_{\theta_i} \right\}.$$

Linear elasticity equations for given transformation strains. For a given sequence of the orientation angles $\boldsymbol{\theta} = \{\theta_i, i = 1, \dots, n\}$ there are up to 3^n corresponding sequences $\boldsymbol{\varepsilon}^T = \{\boldsymbol{\varepsilon}_i^T, i = 1, \dots, n\}$ in the class $\tilde{\Omega}_n(\boldsymbol{\theta})$ defined in (2.6). Here we introduce the relevant elasticity PDEs on the domain Π_L for a *given* such $\boldsymbol{\varepsilon}^T$. We assume that each grain can be described by isotropic elasticity equations with elastic moduli given by

$$(2.7) \quad c_{ijkl} = \lambda \delta_{ij} \delta_{kl} + G(\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}),$$

where λ and G are the Lamé constants [19].

As an applied displacement is imposed, our polycrystal may acquire microscopic strains $\boldsymbol{\varepsilon}$ which contain combined contributions of elastic and stress-free transformations (see [9]):

$$(2.8) \quad \boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}^{\text{elastic}} + \boldsymbol{\varepsilon}^T.$$

Then Hooke’s law $\sigma_{ij} = c_{ijkl} \varepsilon_{kl}^{\text{elastic}}$ yields the stress-strain relation of *linear elasticity* under a given transformation strain $\boldsymbol{\varepsilon}^T$

$$(2.9) \quad \sigma_{ij} = c_{ijkl} (\varepsilon_{kl} - \varepsilon_{kl}^T).$$

Here the strain tensor ε_{kl} is determined by the displacement vector $\mathbf{u} = (u_1(x, y), u_2(x, y))$:

$$(2.10) \quad \varepsilon_{ij} = \frac{1}{2} (\partial_i u_j + \partial_j u_i) \text{ where } \partial_1 = \frac{\partial}{\partial x}, \partial_2 = \frac{\partial}{\partial y}.$$

The stress tensor satisfies the elasticity equations

$$(2.11) \quad \partial_j \sigma_{ij} = 0 \text{ for } i = 1, 2.$$

The above equation is to be understood in the distributional sense, and thus the traction must be continuous across the interfaces between grains:

$$(2.12) \quad [\sigma_{i1}](x, y) = 0, \text{ for } i = 1, 2 \text{ and } x = m/n, m = 1, 2, \dots, n - 1.$$

For such a *given* admissible configuration $\boldsymbol{\varepsilon}_i^T$, we assume a *given* imposed displacement that is chosen in the direction transversal to the laminates:¹

$$(2.13) \quad u_1(0, y) = 0, u_1(1, y) = U, u_2(0, 0) = 0,$$

together with the zero-traction boundary conditions

$$(2.14) \quad \sigma_{12}(0, y) = \sigma_{12}(1, y) = 0, \sigma_{i2}(x, \pm L) = 0.$$

¹Other boundary conditions (e.g., shears) could be treated similarly.

It is easy to check that, for a fixed admissible configuration $\boldsymbol{\varepsilon}^T$, (2.11), (2.13), and (2.14) are the Euler–Lagrange equations for the minimizer of the elastic energy

$$(2.15) \quad W(U, \boldsymbol{\varepsilon}^T) = \frac{1}{E_c} \min_{\mathbf{u}} \frac{1}{2L} \int_{\Pi_L} (\varepsilon_{ij} - \varepsilon_{ij}^T) c_{ijkl} (\varepsilon_{kl} - \varepsilon_{kl}^T) dx dy, \quad \mathbf{u} \text{ subject to (2.13),}$$

where

$$(2.16) \quad E_c = \frac{4G(\lambda + G)}{\lambda + 2G},$$

is the two-dimensional Young modulus. Thus it can be verified that the boundary value problem (2.11), (2.13), and (2.14) admits a unique solution.

Overall polycrystalline energy. As a displacement (2.13) is imposed on the polycrystal, each grain may undergo a stress-free transformation into one of the three possible stress-free states. The overall energy $\mathcal{W}_n(U, \boldsymbol{\theta})$ in the polycrystal is determined by global minimization of the *misfit energy* $W(U, \boldsymbol{\varepsilon}^T)$ of the polycrystal over all admissible configurations [9]

$$(2.17) \quad \mathcal{W}_n(U, \boldsymbol{\theta}) = \min_{\boldsymbol{\varepsilon}^T \in \Omega_n | \boldsymbol{\theta}} W(U, \boldsymbol{\varepsilon}^T).$$

The (possible nonunique) array(s) of transformations strains that arise in the polycrystal is (are) the minimizer(s) in (2.17).

Simplifying assumptions of our model. The idea that the energy minimization in composites and polycrystals can explain correlations has been long pursued in material science (see, e.g., [16, 22]). In this paper we use a quadratic form of the polycrystal's energy proposed in [9] and further developed analytically in [20, 6] and numerically in [6, 7, 8].

The probabilistic model introduced in this work captures many of the essential features of the general physical phenomenon of misfit and, at the same time, is amenable to rigorous analytical treatment.

Clearly, however, our model is too simple to reflect the rich phenomena that occur in actual three-dimensional polycrystals. For example, we consider isotropic elasticity, whereas typically, the crystalline lattice of each of the martensite variants has less symmetry than that of the austenite. Further, for sufficiently large grains, the lattices associated with the various martensite variants could be combined, giving rise to *twins* and/or higher-rank laminates of two or more different variants of martensite within each grain [21, 4, 15, 5, 17]—an effect that our model does not allow. We also note that, in general, a stress-free transformation is a time-dependent process that involves energy dissipation. Our study assumes that the final state of a polycrystal is determined by minimizers of a time-independent, dissipation-free misfit energy (see, e.g., [9, 6, 8] and references therein). Importantly, however, we do not assume that the grains in the polycrystal transform without elastic stresses (self-accommodation); see, e.g., [3, 2] and references therein.

Although not explicitly considered in this work, related phenomena, including electrical and magnetic polarizations in electro- and magneto-rheological materials and the combined elastic and magnetic-electric misfits arising from magnetostriction and electrostriction in composite materials, could be treated by similar methods.

3. Elasticity kernel. In this section we give a representation for the elastic energy $W(U, \boldsymbol{\varepsilon}^T)$ in terms of a certain integral kernel $K_L((x-t))$, and we then present

asymptotics of this kernel under two regimes that are relevant in our studies of the statistics of transformation strains in sections 5, 6, and 7. Denote spatial averages as

$$(3.1) \quad \langle g \rangle = \frac{1}{|\Pi_L|} \int_{\Pi_L} g(x, y) dx dy = \frac{1}{2L} \int_{\Pi_L} g(x, y) dx dy.$$

It turns out that the most convenient mathematical formulation of the elastic energy is in terms of

$$(3.2) \quad s(x) = \varepsilon_{22}^T(x),$$

and the volume fraction of grains that undergo a phase transition²

$$(3.3) \quad f = \langle I \rangle, \quad I = \varepsilon_{11}^T + \varepsilon_{22}^T.$$

Then the elastic energy

$$(3.4) \quad W(U, \varepsilon^T) = \langle (U - f + s)^2 \rangle - \int_{-1}^1 \int_{-1}^1 (s(x) - \langle s \rangle) K_L((x - t))(s(t) - \langle s \rangle) dx dt,$$

where K_L is an even, 2-periodic integral kernel whose cosine Fourier coefficients

$$(3.5) \quad \hat{K}_L(m) = \int_{-1}^1 K_L(x) \cos(\pi m x) dx$$

are explicitly given in Appendix A.1 by formulas (A.5) and (A.6).

The idea of the proof of (3.4) is to decompose the solution $\mathbf{u} = (u_1, u_2)$ of the boundary value problem (2.11), (2.13), and (2.14) in the form $\mathbf{u} = \tilde{\mathbf{u}} + \bar{\mathbf{u}}$, where $\tilde{\mathbf{u}}$ solves the elasticity equations (2.15) for infinitely long grains ($L = \infty$) and the remainder $\bar{\mathbf{u}}$ is the correction for finite L . It turns out that $\tilde{\mathbf{u}}$ is a piecewise linear function of the form

$$(3.6) \quad \begin{cases} \tilde{u}_1^i = a_i x + f_i, \\ \tilde{u}_2^i = c_i x + d y + g_i, \end{cases}$$

and $\bar{\mathbf{u}}$ satisfies the boundary value problem

$$(3.7) \quad \begin{cases} \partial_j \bar{\sigma}_{ij} = 0 \text{ for } i = 1, 2, \\ \bar{u}_1(0, y) = \bar{u}_1(1, y) = 0, \quad \partial_1 \bar{u}_2(0, y) = \partial_1 \bar{u}_2(1, y) = 0, \\ \bar{\sigma}_{12}(x, \pm L) = 0, \quad \bar{\sigma}_{22}(x, \pm L) = E_c(\varepsilon_{22}^T - \langle \varepsilon_{22}^T \rangle). \end{cases}$$

where

$$\begin{cases} \bar{\sigma}_{ij} = c_{ijkl} \bar{\varepsilon}_{kl}, \\ \bar{\varepsilon}_{11} = \partial_1 \bar{u}_1, \quad \bar{\varepsilon}_{12} = \bar{\varepsilon}_{21} = (\partial_2 \bar{u}_1 + \partial_1 \bar{u}_2)/2, \quad \bar{\varepsilon}_{22} = \partial_2 \bar{u}_2. \end{cases}$$

Both functions can be computed explicitly: $\tilde{\mathbf{u}}$ is obtained from direct computations and $\bar{\mathbf{u}}$ is found by the Airy function method. A detailed proof is provided in Appendix A.1. The asymptotics of the convolution kernel in two important limiting

²Since the transformation strains defined in (2.4) satisfy $\varepsilon_{11}^1 + \varepsilon_{22}^1 = 1$, $\varepsilon_{11}^2 + \varepsilon_{22}^2 = 1$, and $\varepsilon_{11}^0 + \varepsilon_{22}^0 = 0$, it follows that the quantity f equals the volume fraction of grains that undergo a phase transition.

cases, in turn, are summarized in the following two lemmas. The proofs of the lemmas are provided in Appendices A.2 and A.3, respectively.

LEMMA 3.1 (kernel asymptotics as $L \rightarrow \infty$). $K_L(x)$ can be represented in the form

$$(3.8) \quad K_L(x) = \frac{B}{L}K_\infty(x) + O(\exp(-L)), \text{ as } L \rightarrow \infty,$$

where

$$K_\infty = -\ln |\sin(\pi x/2)|,$$

and

$$B = \frac{5\lambda + 9G}{4\pi(\lambda + 2G)} > 0.$$

We now consider polycrystals for which the height is commensurate with the grain widths $L = l_0/(2n)$. The parameter l_0 is the height-to-width ratio. In particular, when $L = 1/(2n)$, such polycrystals can be viewed as chains of square grains.

LEMMA 3.2 (kernel asymptotics as $L \rightarrow 0$). Suppose $L = l_0/(2n)$, where $l_0 > 0$, and $n \rightarrow \infty$. Then for each fixed height-to-width ratio l_0 there exists a positive-definite function $\mathcal{K}_{l_0}(x)$, $x \in \mathbb{R}$ independent of n , such that

$$(3.9) \quad \|K_L(x) - n\mathcal{K}_{l_0}(nx)\|_{L^\infty([-1,1])} \rightarrow 0, \text{ as } n \rightarrow \infty,$$

where $\mathcal{K}_{l_0}(x)$ is even: $\mathcal{K}_{l_0}(-x) = \mathcal{K}_{l_0}(x)$, and

$$(3.10) \quad |\mathcal{K}_{l_0}(x)| \leq \frac{c}{l_0} \exp(-|x|/l_0).$$

The constant c in (3.10) does not depend on n and l_0 .

We will find it useful, especially for chains of rectangular grains, to identify sequences (vectors) (f_1, \dots, f_n) (of real number, matrices, etc.) with the corresponding piecewise constant (real valued, matrix valued, etc.) functions defined in the interval $[0, 1]$ that take the values $f(x) = f_k$ for $x \in \Pi_L^k$, $k = 1, \dots, n$. For example, the argument ε^T of W in (2.15), which is a matrix-valued function defined in the interval $[0, 1]$, will often be replaced by a sequence of n matrices $(\varepsilon_1^T, \varepsilon_2^T, \varepsilon_3^T, \dots, \varepsilon_n^T)$. As another example, note that the dependence on n of the quantity on the left-hand side of (2.17) arises merely from the fact that ε^T on the right-hand side of that formula is a piecewise constant function determined by a sequence of n matrices. For a function f defined by a sequence (f_1, \dots, f_n) the spatial average (3.1) is

$$\langle f \rangle = \frac{1}{n} \sum_{i=1}^n f_i.$$

Further, on the space of the piecewise constant functions the integral representation of the misfit energy (3.4) can be viewed as an algebraic nonnegative definite quadratic form:

$$(3.11) \quad W_n(U, \varepsilon^T) = \langle (U - f + s)^2 \rangle - \langle (s - \langle s \rangle)M(n, L)(s - \langle s \rangle) \rangle,$$

where $M = M(n, L)$ is a $n \times n$ symmetric Toeplitz matrix with entries $M_{ij}(n, L) = \lambda_{i-j}(n, L)$ defined by

$$(3.12) \quad (s - \langle s \rangle)M(s - \langle s \rangle) = \int_{-1}^1 \int_{-1}^1 (s(x) - \langle s \rangle)K_L(x - t)(s(t) - \langle s \rangle)dxdt$$

for piecewise constant functions $s = (s_1, \dots, s_n)$.

As an immediate consequence of Lemma 3.2 we can estimate the decay of the coefficients $\lambda_{i-j}(n, L)$ for a chain of rectangular grains ($L = l_0/(2n)$): As $n \rightarrow \infty$, $\lambda_k(n, L) \rightarrow \lambda_k(l_0)$ where

$$(3.13) \quad |\lambda_k(l_0)| \leq \frac{c}{n} \exp(-|k|/l_0).$$

By (3.13) for a chain of rectangular grains $\lambda_{i-j}(n, L)$ can be accurately approximated by a (truncated) Toeplitz matrix by setting $\lambda_k(n, L) \equiv 0$, for $|k| > k_0$, and the misfit energy is approximately

$$(3.14) \quad W_n^{k_0}(U, \boldsymbol{\varepsilon}^T) = \langle (U - f + s)^2 \rangle - \sum_{k=-k_0}^{k_0} \sum_{i=1}^n \lambda_k(n)(s_i - \langle s \rangle)(s_{i+k} - \langle s \rangle),$$

where we define s_{i+k} for $i + k > n$ by periodicity

$$s_{i+k} = \begin{cases} s_{i+k}, & \text{if } i + k \leq n, \\ s_{i+k-n}, & \text{if } i + k > n. \end{cases}$$

The approximation (3.14) provides a justification, in a one-dimensional context, of numerical schemes which are used in practical evaluation of the misfit energy [9, 6, 7, 8]. The approximation (3.14) takes into account only interaction with the nearest neighbors. Hence we call (3.14) k_0 -nearest neighbors energy. The next proposition shows that for any finite value of n and finite k_0 , this k_0 -nearest neighbors approximation has an exponential in k_0 error. Therefore the computational complexity of finding the misfit energy can be significantly reduced if (3.11) is replaced by (3.14). In [6, 7, 8] this truncation was implemented for general two- and three-dimensional polycrystals, and the convergence was verified numerically. The following proposition justifies this convergence analytically in the case of chains of rectangular grains, and provides an explicit exponential error estimate.

PROPOSITION 3.3. *For a given U and a given vector of orientation angles $\boldsymbol{\theta}$, suppose $\boldsymbol{\varepsilon}^T(k_0) \in \tilde{\Omega}_n(\boldsymbol{\theta})$ is a minimizer of the k_0 -nearest neighbors energy $W_n^{k_0}(U, \boldsymbol{\varepsilon}^T)$ given by (3.14). Then there is a universal constant c , independent of n , such that $W_n(U, \boldsymbol{\theta})$, the minimum of the misfit energy (3.11), satisfies*

$$(3.15) \quad |W_n(U, \boldsymbol{\theta}) - W_n^{k_0}(U, \boldsymbol{\varepsilon}^T(k_0))| \leq c \exp(-k_0).$$

A proof the theorem is given in Appendix A.4. Finally, applying Proposition 3.3 and Lemma 3.2, the misfit energy (3.4) of a chain of short grains becomes

$$(3.16) \quad W_n(U, \boldsymbol{\varepsilon}^T) = \langle (U - f + s)^2 \rangle - \lambda_0 \langle (s - \langle s \rangle)^2 \rangle - \frac{B e^{-1/l_0}}{n} \sum_{i=1}^n (s_i - \langle s \rangle)(s_{i+1} - \langle s \rangle) + O(e^{-2/l_0}), \text{ as } l_0 \rightarrow 0,$$

where $\lambda_0 > 0$, $B > 0$.³ Thus, when $l_0 \ll 1$, the misfit energy is approximated by the *nearest neighbor* energy.

³From numerical computations λ_k , $k \geq 2$ are negligible even for large $l_0 = 1$: $\sum_{k=2}^{\infty} |\lambda_k| \leq .1\lambda_1$, $\lambda_1 \leq .17\lambda_0$.

4. A probabilistic model. Our probabilistic model is set to describe energy minimizers within a random setting. In detail, we consider orientation angles θ_i , $i = 1, 2, \dots$ as a sequence of random variables⁴

$$\theta_i : \Omega \rightarrow \Omega^\theta = [0, \pi/4]$$

on a common probability space (Ω, \mathcal{F}, P) . We denote by ρ_n the usual induced probability measure of a sequence of the first n orientation angles on $(\Omega_n^\theta, \sigma_n^\theta)$, where $\Omega_n^\theta = (\Omega^\theta)^n = \Omega^\theta \times \Omega^\theta \times \dots \times \Omega^\theta$, and σ_n^θ is the Borel σ -algebra on Ω_n^θ .

We may define and work with the probability space $(\tilde{\Omega}_n, \tilde{\sigma}_n, \tilde{\mu}_n)$ of arrays of transformation strains $(\varepsilon_i^T)_{i=1}^n$, where we recall that $\tilde{\Omega}_n$ is the set of all sequences of strains that are admissible for some sequence of orientation angles (see (2.5)).

We will, however, work with a different, but equivalent probability space. Recall that the misfit energy (3.4) depends only on a sequence of pairs $\{(s_i, I_i)\}$, $i = 1, \dots, n$, where

$$I_i = \varepsilon_{11,i}^T + \varepsilon_{22,i}^T \in \{0, 1\}, \text{ and } s_i = \varepsilon_{22,i}^T \in [0, 1].$$

Therefore, we will study the probability space

$$(\Omega_n^T, \sigma_n^T, \mu_n), \quad \Omega_n^T = (\Omega^T)^n = \Omega^T \times \Omega^T \times \dots \times \Omega^T, \quad (s_i, I_i) \in \Omega^T = [0, 1] \times \{0, 1\},$$

where σ_n^T is the Borel σ -algebra on Ω_n^T and μ_n is a probability measure, which we will define next.

Probability measure in the space of transformation strains. Suppose the applied deformation U is given. For a fixed sequence of orientation angles $\boldsymbol{\theta}$, there are up to 3^n different admissible arrays of transformation strains $(\varepsilon_i^T)_{i=1}^n$, $T = 0, 1, 2$; see (2.6) with (2.3) and (2.4). They correspond to 3^n different arrays of pairs $\{(s_i, I_i)\}$, $i = 1, \dots, n$ using the rule that the matrices $\varepsilon^1(\theta_i)$, $\varepsilon^2(\theta_i)$, and $\varepsilon^0(\theta_i)$ in (2.4) correspond to the pairs $(s_i, I_i) = (\sin^2 \theta_i, 1)$, $(s_i, I_i) = (\cos^2 \theta_i, 1)$, and $(s_i, I_i) = (0, 0)$, respectively. Some of these arrays $\{(s_i, I_i)\}$, $i = 1, \dots, n$, say a number k of them, minimize the misfit energy (3.4) amongst all admissible arrays. We assume, as it may indeed be natural from a physics perspective, that each of these energy minimizing arrays occurs with equal probability. In other words, we will define the probability measure μ_n in such a way that the conditional probability measure $\mu_n((s_i, I_i)_{i=1}^n | \boldsymbol{\theta})$ satisfies

$$(4.1) \quad \mu_n((s_i, I_i)_{i=1}^n | \boldsymbol{\theta}) = \begin{cases} 1/k, & \text{if } \boldsymbol{\varepsilon}^T \text{ is a minimizer} \\ 0, & \text{if } \boldsymbol{\varepsilon}^T \text{ is not a minimizer,} \end{cases}$$

where k is the number of minimizers of the misfit energy (3.4) for the fixed $\boldsymbol{\theta}$. The probability measure μ_n on the sequences $(s_i, I_i)_{i=1}^n$ is thus defined by

$$(4.2) \quad \mu_n(A) = \int_{\Omega_n^\theta} \mu_n(A | \boldsymbol{\theta}) d\rho_n(\boldsymbol{\theta})$$

for any Borel set $A \in \sigma_n^T$.

⁴In principle $\theta \in [0, \pi/2]$. It follows from (2.4) that we are concerned only with $\sin^2 \theta$ and $\cos^2 \theta$. Since $\sin(\pi/2 - \theta) = \cos \theta$ and $\cos(\pi/2 - \theta) = \sin \theta$, we may and do assume that $\theta \in [0, \pi/4]$. The orientation of a square crystalline lattice can be described uniquely by a value $\theta \in [0, \pi/4]$.

The measure μ_n describes statistics of the transformation strains; the main objective of this paper is to describe it in detail for $n \gg 1$ and for various polycrystal configurations.

Distribution of angles. In the remainder of this paper we will assume that the angles θ_i are independent and identically distributed with the induced probability measure (distribution) ρ on $\Omega^\theta = [0, \pi/4]$ —although other types of angle distributions could be considered within the present context. In other words, the probability measure ρ_n on the sequences of the orientation angles will be taken to be a product measure of the form

$$\rho_n = \prod_{i=1}^n \rho.$$

In particular, to illustrate our theory we will consider two specific probability distributions ρ : (1) the *uniform distribution*, in which ρ is proportional to the Lebesgue measure, and (2) the *Bernoulli trials model* for which ρ is concentrated in two θ values

$$(4.3) \quad \theta_i = \begin{cases} \alpha, & \text{with probability } q, \\ \beta, & \text{with probability } 1 - q, \end{cases}$$

$$0 \leq \beta \leq \alpha \leq \pi/4.$$

5. Statistics of asymptotic energy minimizers 1: Infinitely long grains ($L = \infty$).

5.1. The main theorem. Suppose the grains are infinitely long ($L = \infty$). Then, by Lemma 3.1 the misfit energy for a given admissible sequence of transformation strains on the array of n grains is given by

$$(5.1) \quad W_n(U, \varepsilon^T) = \langle (U - f + \varepsilon_{22}^T)^2 \rangle = \langle (U - f + s)^2 \rangle.$$

The sequence of measures (4.2) contains convergent subsequences [18]; each such limit μ_{lim} is a measure on the set of infinite sequences of transformation strains; the limits along various subsequences may, in principle, not all coincide. In fact, in all cases we consider, however, all such limits do coincide, and the full sequences (4.2) are convergent. For the sake of simplicity, in the subsequent analysis we assume this is the case and we denote $\mu_{\text{lim}} = \lim_{n \rightarrow \infty} \mu_n$.

As we shall show the limits μ_{lim} are convex combinations of product measures. This is a consequence of the de Finetti’s representation theorem (see [13] for a general version of this theorem). In order to motivate the advantages of this observation in our context, we first consider one such limit μ_{lim} and we *assume* (this may or may not be true!) that (1) for each finite n the minimizers are unique, and (2) the measure μ_{lim} is given by a product of the form

$$(5.2) \quad \mu_{\text{lim}} = \tilde{\mu} := \prod_{i=1}^\infty \mu,$$

for a certain measure μ so that, according to μ_{lim} , the random variables $\{(s_i, I_i)\}_{i=1}^\infty$ are i.i.d.

Since, as we have seen above, for each i we must necessarily have $s_i = \sin^2(\theta_i)$ or $s_i = \cos^2(\theta_i)$ whenever $I_i = 1$, under the uniqueness assumption (1) above the measure μ in (5.2) must satisfy

$$\begin{cases} \mu(s = \sin^2(\theta)|\theta, I = 1) + \mu(s(\theta) = \cos^2(\theta)|\theta, I = 1) = 1 \\ \mu(s = \sin^2(\theta)|\theta, I = 1) = 0 \text{ or } 1 \quad \text{and} \quad \mu(s(\theta) = \cos^2(\theta)|\theta, I = 1) = 0 \text{ or } 1. \end{cases}$$

Hence, we can define a function $\kappa(\theta)$, $\kappa(\theta) = 0$, or $\kappa(\theta) = 1$ such that

$$(5.3) \quad \kappa(\theta) = \begin{cases} 1, & \text{if } \mu(s = \sin^2(\theta)|\theta, I = 1) = 1, \\ 0, & \text{if } \mu(s = \cos^2(\theta)|\theta, I = 1) = 1. \end{cases}$$

By the law of large numbers, as $n \rightarrow \infty$ we have for the misfit energy (5.1)

$$(5.4) \quad W_n(U, \epsilon^T) \rightarrow \int_0^{\pi/4} g(\theta)\chi(\theta)d\rho(\theta) + (U - f)^2(1 - f),$$

where $g(\theta) = (U - f + \sin^2 \theta)^2 \kappa(\theta) + (U - f + \cos^2 \theta)^2 (1 - \kappa(\theta))$ and where we have set

$$(5.5) \quad \chi(\theta) = \mu(I = 1|\theta), \quad (0 \leq \chi(\theta) \leq 1),$$

and $f = \int_0^{\pi/4} \chi(\theta)d\rho(\theta)$.

Clearly, in the present context the limiting values of the energy function $W_n(U, \epsilon^T)$ are determined uniquely by the functions $\kappa(\theta)$ and $\chi(\theta)$. Since μ_{lim} is the limit of probability measures $\{\mu_n\}$ given by (4.1) with (5.1), it follows that the measure μ in (5.2) must minimize (5.4). In other words, under the assumption (5.2), the overall minimization problem has been reduced to the following minimization problem for the functions $\kappa(\theta)$ and $\chi(\theta)$:

$$(5.6) \quad \begin{aligned} & \min_{\kappa(\theta), \chi(\theta)} \int_0^{\pi/4} g(\theta)\chi(\theta)d\rho(\theta) + (U - f)^2(1 - f) \\ & g(\theta) = (U - f + \sin^2 \theta)^2 \kappa(\theta) + (U - f + \cos^2 \theta)^2 (1 - \kappa(\theta)), \\ & f = \int_0^{\pi/4} \chi(\theta)d\rho(\theta), \quad 0 \leq \chi(\theta) \leq 1, \quad 0 \leq \kappa(\theta) \leq 1. \end{aligned}$$

One can anticipate that, generally, the assumption (5.2) does not hold. Indeed, even working under the assumption (5.2), we note that a solution μ (i.e., (κ, χ)) to the minimization problem (5.6) may not be unique. If there are two such solutions μ_1 and μ_2 to this problem, then the corresponding infinite products $\tilde{\mu}_1$ and $\tilde{\mu}_2$ could, conceivably, equal to the limit of a subsequence of μ_n . As shown in Theorem 5.2, however, in general μ_{lim} will equal a convex combination of such infinite products. The following definition will be useful in these regards.

DEFINITION 5.1. Consider the set \mathbb{S} of all measures on the set of pairs $(s, I) \in [0, 1] \times \{0, 1\}$. For each measure $\mu \in \mathbb{S}$ define the associated product measure

$$\tilde{\mu} = \prod_{i=1}^{\infty} \mu$$

on $\prod_{i=1}^{\infty} (s_i, I_i)$. A measure γ is called a convex combination of such product measures, if there exists a positive measure $\nu(\mu)$ on the set \mathbb{S} , such that

$$(5.7) \quad \gamma(A) = \int_{\mathbb{S}} \tilde{\mu}(A)d\nu(\mu), \quad \int_{\mathbb{S}} d\nu(\mu) = 1.$$

THEOREM 5.2. Consider infinitely long grains and an arbitrary i.i.d. angle distribution with the probability measure ρ , and assume the limit μ_{lim} of the sequence

μ_n defined by (4.2) and (4.1) exists. Then μ_{lim} is given by a convex combination of product measures arising from minimization problems. In detail, we have

$$\mu_{\text{lim}} = \int_{\mathbb{S}} \tilde{\mu} d\nu(\mu), \quad \int_{\mathbb{S}} d\nu(\mu) = 1,$$

where \mathbb{S} is the set of product measures: $\tilde{\mu} = \prod_{i=1}^{\infty} \mu$, and each μ is defined by

$$\mu((s, I)|\theta) = \begin{cases} \chi(\theta)\kappa(\theta), & \text{if } I = 1, \\ 1 - \chi(\theta), & \text{if } I = 0, \end{cases}$$

where $\kappa(\theta)$ and $\chi(\theta)$ are minimizers of

$$(5.8) \quad \min_{\kappa(\theta), \chi(\theta)} \int_0^{\pi/4} (U - f + s(\theta))^2 \chi(\theta) d\rho(\theta) + (U - f)^2(1 - f),$$

$$f = \int_0^{\pi/4} \chi(\theta) d\rho(\theta), \quad \chi : [0, \pi/4] \rightarrow [0, 1],$$

$$(5.9) \quad s(\theta) = \kappa(\theta) \sin^2(\theta) + (1 - \kappa(\theta)) \cos^2(\theta), \quad \text{where } \kappa : [0, \pi/4] \rightarrow \{0, 1\}.$$

If the minimizer of (5.8) is unique, then transformation strains in different grains are independent identically distributed; that is, μ_{lim} is a product measure $\mu_{\text{lim}} = \prod_{i=1}^{\infty} \mu$.

Proof. A key property of the energy W_n (5.1) in the case of infinitely long grains is that it is *invariant under permutations*; e.g., for a three-grain polycrystal, if $(s_1, I_1), (s_2, I_2), (s_3, I_3)$ is a minimizing sequence for the angles $(\theta_1, \theta_2, \theta_3)$, then $(s_2, I_2), (s_1, I_1), (s_3, I_3)$ is a minimizing sequence for the angles $(\theta_2, \theta_1, \theta_3)$ with the same probability. More generally, the form of the misfit energy and our assumption (4.1) imply that the probability measure μ_n (defined by (4.1), (4.2)) on the minimizers must be symmetric;⁵ that is, for any finite permutation $\tau \in S(n)$

$$\mu_n((s_1, I_1) \in A_1, \dots, (s_n, I_n) \in A_n) = \mu_n((s_{\tau(1)}, I_{\tau(1)}) \in A_1, \dots, (s_{\tau(n)}, I_{\tau(n)}) \in A_n).$$

As $n \rightarrow \infty$, the probability measure μ_n converges to a certain μ_{lim} . Clearly, μ_{lim} must be symmetric as well: For any n and $\tau \in S(n)$

$$\begin{aligned} & \mu_{\text{lim}}((s_1, I_1) \in A_1, \dots, (s_n, I_n) \in A_n, (s_{n+1}, I_{n+1}) \in A_{n+1}, \dots) \\ &= \mu_{\text{lim}}((s_{\tau(1)}, I_{\tau(1)}) \in A_1, \dots, (s_{\tau(n)}, I_{\tau(n)}) \in A_n, (s_{n+1}, I_{n+1}) \in A_{n+1}, \dots). \end{aligned}$$

Hence, we can apply the de Finetti's representation theorem [13], and μ_{lim} must be a convex combination of product measures $\tilde{\mu}$:

$$\mu_{\text{lim}} = \int_{\mathbb{S}} \tilde{\mu} d\nu(\mu), \quad \int_{\mathbb{S}} d\nu(\mu) = 1, \quad \nu(\mu) \geq 0;$$

see Definition 5.1, as claimed.

Let us now show that

$$(5.10) \quad \begin{aligned} & \lim_{n \rightarrow \infty} \int_{\Omega_n^T} W_n(U, \varepsilon^T) d\mu_n \\ &= \int_{\mathbb{S}} \left[\int_0^{\pi/4} (U - f + s(\theta))^2 \chi(\theta) d\rho(\theta) + (U - f)^2(1 - f) \right] d\nu(\mu). \end{aligned}$$

⁵Sometimes the term exchangeable is used instead of symmetric.

The key issue here is classical: Given that $\mu_n \rightarrow \mu_{\text{lim}}$ *weakly*, we cannot, in general, conclude convergence of $\int W_n d\mu_n$. In our case, however, we can, because the measures μ_n are symmetric. It follows that for an n -grain sample, the functions $\{(s_i, I_i)\}_{i=1}^n$ satisfy, for any $i, j, 1 \leq i, j \leq n$,

$$\int I_i I_j d\mu_n = \begin{cases} \int I_1 I_2 d\mu_n, & i \neq j, \\ \int I_1^2 d\mu_n, & i = j, \end{cases}$$

and similar equalities hold for $s_i I_j$, and $(U - \langle I \rangle + s_i)^2$. Therefore, W_n can be written as

$$\int W_n(U, \varepsilon^T) d\mu_n = \int F(U, s_1, I_1, I_2) d\mu_n + \frac{1}{n} \int G(s_1, I_1, I_2) d\mu_n,$$

where *both* functions $F(U, s_1, I_1, I_2)$ and $G(s_1, I_1, I_2)$ depend *continuously* (they are quadratic polynomials) on the values of s_i and I_i only in two grains $i = 1, 2$ (and, thus, do not depend on n), and they are explicitly given as

$$F = (U + s_1)^2 - 2UI_1 - 2s_1I_2 + I_1I_2, \quad G = -2s_1I_1 + 2s_1I_2 + I_1^2 - I_1I_2.$$

Therefore

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_{\Omega_n^T} W_n(U, \varepsilon^T) d\mu_n &= \lim_{n \rightarrow \infty} \int_{\Omega_n^T} F(U, s_1, I_1, I_2) d\mu_n + \lim_{n \rightarrow \infty} \frac{1}{n} \int_{\Omega_n^T} G(s_1, I_1, I_2) d\mu_n \\ &= \int_{\Omega_\infty^T} F(U, s_1, I_1, I_2) d\mu_{\text{lim}}. \end{aligned}$$

Hence it only remains to show that

$$\begin{aligned} &\int_{\Omega_\infty^T} F(U, s_1, I_1, I_2) d\mu_{\text{lim}} \\ (5.11) \quad &= \int_{\mathbb{S}} \left[\int_0^{\pi/4} (U - f + s(\theta))^2 \chi(\theta) d\rho(\theta) + (U - f)^2 (1 - f) \right] d\nu(\mu). \end{aligned}$$

The last equality is obtained by explicit computations provided in Appendix B.1. The proof of the identity (5.10) is now complete.

Further, up to a set of ν -measure zero, each μ must minimize (5.8). Otherwise, we can choose a $\delta > 0$ such that the set

$$A = \left\{ \mu : \int F(U, s_1, I_1, I_2) d\tilde{\mu} - \delta > \min \int F(U, s_1, I_1, I_2) d\tilde{\mu} \right\}$$

has positive measure: $\nu(A) > 0$. Then, if $\nu(\mathbb{S} \setminus A) \neq 0$ we define a new measure $\tilde{\nu}$ by

$$\tilde{\nu}(B) = \nu(B \setminus A) / \nu(\mathbb{S} \setminus A).$$

Clearly

$$\int_{\mathbb{S}} \left[\int F(U, s_1, I_1, I_2) d\tilde{\mu} \right] d\tilde{\nu}(\mu) < \int_{\mathbb{S}} \left[\int F(U, s_1, I_1, I_2) d\tilde{\mu} \right] d\nu(\mu),$$

which contradicts the assumption that ν yields a limit of minimizers of the misfit energy (5.8) as indicated in (5.11). If $\nu(\mathbb{S} \setminus A) \neq 0$, in turn, we can select a single minimizer and assign $\tilde{\nu}$ measure 1 to it, arriving again to a contradiction.

To establish (5.9), note that for a minimizer μ of (5.6), the μ probabilities conditional to a given angle θ and to $I = 1$, satisfy

$$\mu(s(\theta) = \cos^2(\theta)|\theta, I = 1) = 0, \text{ or } \mu(s(\theta) = \sin^2(\theta)|\theta, I = 1) = 0.$$

Hence $\kappa(\theta)$ takes only the values 0 and 1 and (5.9) holds.

Finally, suppose (5.8) admits a unique minimizer μ . Since, as established above, the limit μ_{lim} must be a convex combination of product measures that minimize (5.8); in the case of uniqueness of solution to (5.8), μ_{lim} must be the product measure $\mu_{\text{lim}} = \prod_{i=1}^{\infty} \mu$, as claimed. \square

A few remarks about Theorem 5.2 are in order. Firstly, we are aware of some examples when the minimization problem (5.8) has more than one solution. One of these examples is to consider a deterministic sequence $\theta_i = 0$ and $U = 1$. Then there are two solutions $\kappa_1 \equiv 0, \chi_1 \equiv 1$, and $\kappa_2 \equiv 0, \chi_2 \equiv 0$ to the minimization problem (5.8), which give rise to corresponding measures μ_1 and μ_2 , and, thus, product measures $\tilde{\mu}_1$ and $\tilde{\mu}_2$. For both $i = 1, 2$ we have

$$\int_0^{\pi/4} (U - f_i + s_i(\theta))^2 \chi_i(\theta) d\rho(\theta) + (U - f_i)^2 (1 - f_i) = 1.$$

In view of our symmetrization assumption 4.1, the limit of μ_n exists and it is equal to a convex combination of product measures as implied by Theorem (5.2); the convex combination is given by $\tilde{\mu}_1/2 + \tilde{\mu}_2/2$. We expect that generically the minimization problem (5.8) has a unique solution. We give two explicit examples when this measure is unique: Bernoulli trials, Lemma 5.3, and Uniform distribution, Lemma 5.6, in section 5.2.

Secondly, in principle, the (unique) solution to the minimization problem (5.8) may be such that $\chi(\theta)$ takes only two values 0 or 1, i.e., $\chi : [0, \pi/4] \rightarrow \{0, 1\}$. This, indeed, happens for uniform distribution (Lemma 5.6 in section 5.2). If we know that $\chi : [0, \pi/4] \rightarrow \{0, 1\}$, then the proof of Theorem 5.2 becomes straightforward. However, there are examples, when $0 < \chi(\theta) < 1$ for some θ , and one of them is the Bernoulli trials (Corollary 5.4 in section 5.2).

Finally, if we do not assume the uniform conditional probability (4.1), then μ_{lim} may not be unique even if the minimizer of (5.8) is unique. We discuss this issue for Bernoulli trials after the proof of Corollary 5.4 below.

Motivated by the above remarks, we next investigate in more detail how measure μ , a solution to the minimization problem (5.8), depends on the underlying probability measure ρ for two specific probability measures ρ : the Bernoulli trials model and the uniform distribution of θ .

5.2. Bernoulli trials and uniform distribution. For the Bernoulli trials model (4.3), the minimization problem (5.8) from Theorem 5.2 is

$$(5.12) \quad \min_{q_\alpha, q_\beta, s(\alpha), s(\beta)} W(q_\alpha, q_\beta, s(\alpha), s(\beta)),$$

$$W = (U - f + s(\alpha))^2 q_\alpha + (U - f + s(\beta))^2 q_\beta + (U - f)^2 (1 - f),$$

with $f = q_\alpha + q_\beta$,

$$(5.13) \quad 0 \leq q_\alpha \leq q, \quad 0 \leq q_\beta \leq 1 - q,$$

$0 \leq \beta \leq \alpha \leq \pi/4$, and $s(\alpha), s(\beta)$ are defined by $s(\theta) = \sin^2 \theta$ or $s(\theta) = \cos^2 \theta$. In particular, the minimization with respect to χ is reduced to determining the proportions q_α and q_β of grains with angles α and β that do not undergo a stress-free transformation.

LEMMA 5.3. For the Bernoulli trials model (4.3) with $0 < q < 1$ the minimizer (κ, χ) of (5.12) is unique. For the minimizer $\kappa(\theta) \equiv 1$, i.e.,

$$(5.14) \quad s(\theta) = \sin^2 \theta,$$

and χ depends on U and can be described as follows. For a given U the total proportion of grains that undergoes a stress-free transformation $f = f(U)$ is a (deterministic) nondecreasing function of U . For a given f we have several cases

- if $f < 1 - q$, then

$$(5.15) \quad q_\alpha = 0, \quad 0 \leq q_\beta \leq 1 - q,$$

i.e., $\chi(\alpha) = 0, \chi(\beta) = q_\beta/(1 - q)$,

- if $f > 1 - q$, then

$$(5.16) \quad 0 \leq q_\alpha \leq q, \quad q_\beta = 1 - q,$$

i.e. $\chi(\beta) = q_\alpha/q, \chi(\beta) = 1$.

A proof of the Lemma is in Appendix B.2.

COROLLARY 5.4. The probability distribution

$$\theta_i = \begin{cases} \pi/4, & \text{with probability } q, \quad q \geq 1/2, \\ 0, & \text{with probability } 1 - q \end{cases}$$

is an example, where the minimizer (κ, χ) of (5.12) is unique, but

$$\chi : [0, \pi/4] \rightarrow K \subset [0, 1], \quad K \neq \{0, 1\}.$$

Indeed, in this case

$$\begin{aligned} \chi(0) = U, \chi(\pi/4) = 0, & \text{ if } U \leq 1 - q, \\ \chi(0) = 1, \chi(\pi/4) = 0, & \text{ if } 1 - q \leq U \leq 5/4 - q, \\ \chi(0) = 1, \chi(\pi/4) = U - 5/4 + q, & \text{ if } 5/4 - q \leq U \leq 3/4, \text{ or} \\ \chi(0) = 1, \chi(\pi/4) = 1, & \text{ if } 3/4 \leq U. \end{aligned}$$

Let us now discuss our final remark that if we do not assume the uniform conditional probability (4.1), then μ_{lim} may not be unique even if the minimizer of (5.12) is unique. It depends on whether χ takes more than two values, that is, on whether $\chi : [0, \pi/4] \rightarrow K$, but $K \neq \{0, 1\}$. For example, consider the Bernoulli trials with $0 < q_\alpha < q$. Then one can choose the grains with $\theta_i = \alpha$, which do not undergo a stress-free transformation, arbitrarily, provided that their total proportion is q_α . Thus, if we remove our assumption of equal probability (4.1), in the case of infinitely long grains there are many minimizers of the misfit energy in addition to minimizers described in Theorem 5.2. Hence, it is possible to construct the limiting measure μ_{lim} on the infinite sequence of pairs $\{(s_i, I_i)\}_{i=1}^\infty$ so that it is not a product measure. Moreover, actual construction of the exact minimizers ϵ^T of the energy (5.1) (for a given sequence $\{\theta_i\}_{i=1}^n$) in practice [9, 6, 7, 8] is done numerically. Thus it typically results in finding an *almost* minimizer $\tilde{\epsilon}^T$, such that

$$(5.17) \quad |W_n(U, \epsilon^T) - W_n(U, \tilde{\epsilon}^T)| \leq \delta, \quad \delta > 0.$$

Thus, it is natural to ask which characteristic properties of exact minimizers are approximated by characteristic properties of almost minimizers. The property

that μ_{lim} is a product measure is *not* characteristic, but the proportion of grains that undergo a stress-free transformation is such property. For example, for Bernoulli trials, $q_\alpha(\boldsymbol{\theta})$ characterizes the proportion of grains with $\theta_i = \alpha$, $i = 1, \dots, n$, which undergo a stress-free transformation (grains for which $I = 1$), and we have the following immediate result.

LEMMA 5.5. *For every $\delta' > 0$ there exists $\delta > 0$ such that an almost minimizer $\tilde{\boldsymbol{\varepsilon}}^T$ in the sense (5.17) satisfies*

$$|q_\alpha(\boldsymbol{\theta}) - \tilde{q}_\alpha(\boldsymbol{\theta})| < \delta', \quad |q_\beta(\boldsymbol{\theta}) - \tilde{q}_\beta(\boldsymbol{\theta})| < \delta',$$

where $q_\alpha(\boldsymbol{\theta})$ and $\tilde{q}_\alpha(\boldsymbol{\theta})$ correspond to the exact and almost minimizers, respectively ($q_\beta(\boldsymbol{\theta})$ and $\tilde{q}_\beta(\boldsymbol{\theta})$ are defined analogously).

Analogous to the Bernoulli trials model, direct computations show the following result for the uniform distribution.

LEMMA 5.6. *For the uniform distribution of $\theta \in [0, \pi/4]$, the minimizer (κ, χ) of (5.12) is unique. For the minimizer $\kappa(\theta) \equiv 1$, hence*

$$(5.18) \quad s(\theta) = \sin^2 \theta,$$

and χ depends on U and can be described as follows. For a given U the total proportion of grains that undergoes a stress-free transformation $f(U)$ is a (deterministic) nondecreasing function of U given by

$$f(U) = \begin{cases} \frac{4\pi}{8+\pi}U + g(U), & \text{if } U \leq \frac{1}{4} + \frac{2}{\pi} \approx .88662, \\ 1 & \text{otherwise,} \end{cases}$$

where the small correction $g(U)$ is concave and it satisfies $g(0) = f(1/4 + 2/\pi) = 0$, $-0.055 < g(U) \leq 0$. For a given $f < 1$

$$\chi(\theta) = \begin{cases} 1, & \text{if } \sin^2 \theta \leq 2(f - U), \\ 0, & \text{otherwise.} \end{cases}$$

If $f = 1$, then $\chi \equiv 1$.

Lemmas 5.3 and 5.6 together with Theorem 5.2 imply the following.

COROLLARY 5.7. *For Bernoulli trials and uniform distribution, the unique minimizing sequence of transformation strains is i.i.d, and, in particular, it is uncorrelated.*

6. Statistics of asymptotic energy minimizers 2: Thin long grains (finite $L \gg 1$).

6.1. Basic definitions and formulation of the main theorem. In contrast to the case of infinitely long grains, if L is large but finite, then each grain may undergo a stress-free transformation which, as we show in this section, is correlated to stress-free transformations of other grains. In particular, for Bernoulli trials in case 1 ($L = \infty$, $n \rightarrow \infty$) the minimizers are shown to be i.i.d. (see Corollary 5.7 in the previous section), whereas in case 2 ($n \rightarrow \infty$, followed by $L \rightarrow \infty$) the minimizers are no longer i.i.d. (see Theorem 6.4 below).⁶

By Lemma 3.1 the misfit energy for $L \gg 1$ has the following asymptotic representation (up to higher order terms):

$$(6.1) \quad W_n(U, \boldsymbol{\varepsilon}^T) = \langle (U - f + s)^2 \rangle - \frac{B}{L} \bar{W}_n(U, \boldsymbol{\varepsilon}^T), \quad \text{as } L \rightarrow \infty,$$

⁶In this sense, we prove that the limits for large n and large L do not commute.

where

$$(6.2) \quad \bar{W}_n(U, \varepsilon^T) = \int_{-1}^1 \int_{-1}^1 (s(x) - \langle s \rangle)(s(t) - \langle s \rangle) K_\infty(x - t) dx dt,$$

$$(6.3) \quad K_\infty = -\ln |\sin(\pi x/2)| > 0.$$

For $L \gg 1$ the second term on the right-hand side of (6.1) amounts to a small correction to the misfit energy of infinitely long grains (5.1) (the first term in (6.1)), a situation that bears connections with the concept of almost minimizers (5.17) introduced in the previous section. In the present context we have:

LEMMA 6.1. *For any $\delta > 0$, there exists a sufficiently large $L_0 > 0$, so that for any $L \geq L_0$*

$$(6.4) \quad s(\theta) = \sin^2 \theta,$$

$$(6.5) \quad |q_\alpha(\theta) - q_\alpha^L(\theta)| < \delta, \quad |q_\beta(\theta) - q_\beta^L(\theta)| < \delta,$$

where $q_\alpha(\theta)$ and $q_\alpha^L(\theta)$ correspond to minimizers of (5.1) and (6.1), respectively ($q_\beta(\theta)$ and $q_\beta^L(\theta)$ are defined analogously).

A proof of Lemma 6.1 is in Appendix C.1. As we pointed out after Corollary 5.4, when $L = \infty$ and $0 < q_\alpha < q$ or $0 < q_\beta < 1 - q$ there is an ambiguity: The solution to the minimization problem (5.1) (the first term in (6.1)) is not unique. We now show that the second term (6.2) plays a role of regularization—it resolves this ambiguity by reducing the number of minimizers, and it gives rise to correlation of transformations in different grains of a laminated polycrystal. The correlations arise from maximization of \bar{W} (note the negative sign in front of the second term in (6.1)). We formalize this idea in the next definition.

DEFINITION 6.2. *For a given deformation U and a sequence of angles $\{\theta_i\}_{i=1}^n$, an asymptotic energy minimizer is a pair of piecewise constant functions*

$$(s(x), I(x)) = (s_i, I_i) \text{ if } x \in \left(\frac{i-1}{n}, \frac{i}{n} \right], \quad i = 1, 2, \dots, n, \quad s(x), I(x) \in \mathcal{H}_n,$$

which maximizes (6.2) (the second term in (6.1)) amongst all⁷ minimizers of the misfit energy of infinitely long grains (5.1) (the first term in (6.1)).

When the distribution of angles ρ_n is given, we denote by μ_n^a the corresponding probability measure of the distribution of asymptotic energy minimizers $(s_i, I_i)_{i=1}^n$. The behavior of the asymptotic energy minimizers for the distribution of angles ρ_n given by the Bernoulli trials (4.3) model will be described by the Riesz symmetrically rearranged minimizer which we define as follows.

DEFINITION 6.3. *Consider Bernoulli trials (4.3) model (θ is α or β with probabilities q or $1 - q$). For a given displacement U , let q_α and q_β be the proportion of grains for which $I = 1$ with $\theta = \alpha$ and $\theta = \beta$, respectively,*

$$(6.6) \quad 0 \leq q_\alpha \leq q, \quad 0 \leq q_\beta \leq 1 - q,$$

⁷Note that we consider here all possible minimizers of the misfit energy of infinitely long grains (5.1) (the first term in (6.1)). In other words, here minimizers of (5.1) may not satisfy our assumption of equal probability (4.1).

as given by Lemma 5.3. For a given sequence of angles $\{\theta_i\}_{i=1}^n$ a Riesz left-rearranged sequence of transformation strains is a pair of functions $(s_l(x), I_l(x)) \in \mathcal{H}_n$ given for each $x \in ((i-1)/n, i/n]$ as

$$(s_l(x), I_l(x)) = \left(s_l\left(\frac{i}{n}\right), I_l\left(\frac{i}{n}\right) \right) = \begin{cases} (\sin^2 \alpha, 1), & \text{if } \theta = \alpha, \quad 0 \leq i/n \leq q_\alpha/q, \\ (0, 0), & \text{if } \theta = \alpha, \quad q_\alpha/q \leq i/n \leq 1, \\ (\sin^2 \beta, 1), & \text{if } \theta = \beta, \quad 0 \leq i/n \leq q_\beta/(1-q), \\ (0, 0), & \text{if } \theta = \beta, \quad q_\beta/(1-q) \leq i/n \leq 1. \end{cases}$$

Similarly, a Riesz right-rearranged sequence $(s_r(x), I_r(x)) \in \mathcal{H}_n$ is determined by

$$(s_r(x), I_r(x)) = \left(s_r\left(\frac{i}{n}\right), I_r\left(\frac{i}{n}\right) \right) = \begin{cases} (\sin^2 \alpha, 1), & \text{if } \theta = \alpha, \quad 1 - q_\alpha/q \leq i/n \leq 1, \\ (0, 0), & \text{if } \theta = \alpha, \quad 0 \leq i/n \leq 1 - q_\alpha/q, \\ (\sin^2 \beta, 1), & \text{if } \theta = \beta, \quad 1 - q_\beta/(1-q) \leq i/n \leq 1, \\ (0, 0), & \text{if } \theta = \beta, \quad 0 \leq i/n \leq 1 - q_\beta/(1-q). \end{cases}$$

We denote by μ_n^l and μ_n^r the corresponding probability measures. A Riesz symmetrically rearranged probability measure μ_n^s is the average of right- and left-rearranged measures $\mu_n^s = 1/2\mu_n^l + 1/2\mu_n^r$. Finally, the Riesz symmetrically rearranged probability measure μ^s , the weak limit of probability measures: $\mu_n^s \Rightarrow \mu^s$, as $n \rightarrow \infty$.

The rearranged minimizers quantitatively describe the rise of correlations for asymptotic energy minimizers (Definition 6.2), because, as we prove in Lemma 6.7, the asymptotic and rearranged minimizers coincide in the limit as $n \rightarrow \infty$. In other words, Definitions 6.2 and 6.3 characterize the same measure as $n \rightarrow \infty$. Moreover, the following theorem shows that the minimizer probability measure of the full misfit energy (6.1) converges to the Riesz symmetrically rearranged probability measure when we let $n \rightarrow \infty$, and then let $L \rightarrow \infty$.

THEOREM 6.4. *Consider the Bernoulli trials (4.3). For a given U , let q_α and q_β be defined as in Lemma 5.3 by (5.15) or (5.16). Then as $n \rightarrow \infty$ and subsequently $L \rightarrow \infty$ the probability measure of the energy minimizer of the misfit energy (6.1) converges weakly to the Riesz symmetrically rearranged probability measure μ^s .*

6.2. Riesz rearrangement inequalities and proof of Theorem 6.4. The key idea of the proof comes from the classical Riesz rearrangement inequality (see e.g., [12], [14]). In particular, this inequality motivated the name for minimizers in Definition 6.3. The simplest form of this inequality, which is sufficient for our purposes is as follows.

LEMMA 6.5. *Riesz rearrangement inequality on a circle.*

Consider two classes of even, bounded, and positive functions on $[-1, 1]$:

$$\mathcal{A}_i = \{f(x) | f(x) = f(-x), 0 \leq f(x) \leq 1, \int_0^1 f(x)dx = p_i\}, \quad i = 1, 2.$$

Let $\chi_{p_i}^1(x) \in \mathcal{A}_i$ be the characteristic function of the set $[-p_i, p_i]$ and $\chi_{p_i}^2(x) \in \mathcal{A}_i$ be the characteristic function of the set $[-1, -1 + p_i] \cup [1 - p_i, 1]$. Suppose $K(x)$ is an even positive locally integrable 2-periodic function on \mathbb{R} that decreases on $[0, 1]$:

(6.7)

$$K(x) \geq 0, \quad K(x) = K(-x), \quad \int_0^1 K(x)dx < \infty, \quad K(x+2) = K(x), \quad \text{and } K(x_1) > K(x_2),$$

if $0 < x_1 < x_2 \leq 1$. Then for any $f(x) \in \mathcal{A}_1, g(x) \in \mathcal{A}_2$

$$(6.8) \quad \int_{-1}^1 \int_{-1}^1 (f(x) - p_1)(g(t) - p_2)K(x - t)dxdt \leq \int_{-1}^1 \int_{-1}^1 (\chi_{p_1}^1(x) - p_1)(\chi_{p_2}^1(t) - p_2)K(x - t)dxdt,$$

where the equality holds only in the following two cases

$$(6.9) \quad (a) f(x) = \chi_{p_1}^1(x), g(x) = \chi_{p_2}^1(x), (b) f(x) = \chi_{p_1}^2(x), g(x) = \chi_{p_2}^2(x).$$

Moreover, for any $\delta > 0$ there exists $\delta' > 0$ so that if $\min(e_1, e_2) \geq \delta'$,

$$(6.10) \quad e_1 = \int_0^1 (|f - \chi_{p_1}^1| + |g - \chi_{p_2}^1|)dx, \quad e_2 = \int_0^1 (|f - \chi_{p_1}^2| + |g - \chi_{p_2}^2|)dx,$$

then

$$(6.11) \quad \int_{-1}^1 \int_{-1}^1 (f(x) - p_1)(g(t) - p_2)K(x - t)dxdt \leq \int_{-1}^1 \int_{-1}^1 (\chi_{p_1}^1(x) - p_1)(\chi_{p_2}^1(t) - p_2)K(x - t)dxdt - \delta.$$

We assumed in this lemma that the functions are bounded from above by one. This assumption can be replaced by any positive number with obvious modifications of the results. The proof of Lemma 6.5 follows from considerations similar to those found in [1]. It basically says that among all possible functions $0 \leq f(x) \leq 1, 0 \leq g(x) \leq 1$ on a circle $[-1, 1]$ (where the endpoints $x = \pm 1$ are identified) the maximum of the integral

$$(6.12) \quad \int_{-1}^1 \int_{-1}^1 (f(x) - p_1)(g(t) - p_2)dxdt, \quad \text{with} \quad \int_0^1 f(x)dx = p_1, \int_0^1 g(x)dx = p_2$$

is achieved on characteristic functions of the intervals of length $2p_1$ and $2p_2$. The reason why the intervals centered at $x = 0$ and $x = 1$ is due to our assumption that $f(x)$ and $g(x)$ are even. In order to explain how Lemma 6.5 must be modified and applied for our case, we decompose

$$s(x) - \langle s \rangle = \sin^2 \alpha \varepsilon_\alpha + \sin^2 \beta \varepsilon_\beta, \quad \varepsilon_\alpha = (\chi_\alpha \tilde{\chi}_\alpha - q_\alpha), \quad \varepsilon_\beta = (\chi_\beta \tilde{\chi}_\beta - q_\beta),$$

where χ_α and $\chi_\beta, \chi_\alpha + \chi_\beta = 1$ are (random) characteristic functions of the angle distributions $\theta = \alpha$ and $\theta = \beta$, respectively; $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ are the characteristic function of the grains with $\theta = \alpha$ and $\theta = \beta$, respectively, for which $I = 1$.

The term (6.2) (the second term in (6.1)) equals

$$(6.13) \quad \begin{aligned} \bar{W}(U, \varepsilon^T) &= \sin^4 \alpha \int_{-1}^1 \int_{-1}^1 \varepsilon_\alpha(x) \varepsilon_\alpha(t) K_\infty(x - t) dx dt \\ &+ 2 \sin^2 \alpha \sin^2 \beta \int_{-1}^1 \int_{-1}^1 \varepsilon_\alpha(x) \varepsilon_\beta(t) K_\infty(x - t) dx dt \\ &+ \sin^4 \beta \int_{-1}^1 \int_{-1}^1 \varepsilon_\beta(x) \varepsilon_\beta(t) K_\infty(x - t) dx dt. \end{aligned}$$

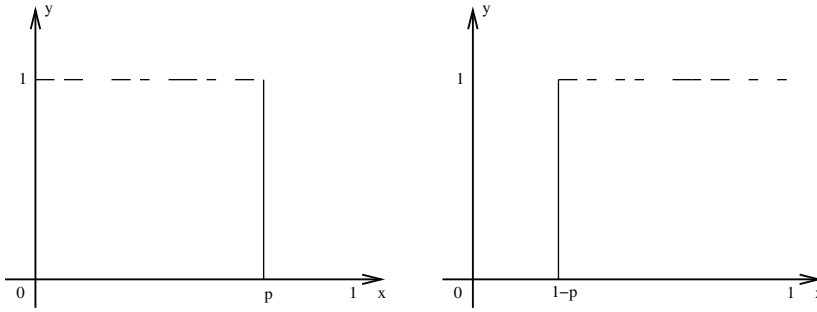


FIG. 6.1. A sample of two random intervals $\chi_\alpha(x)\chi_{q_\alpha}^1(x)$ and $\chi_\alpha(x)\chi_{q_\alpha}^2(x)$, where $p = q_\alpha/q$.

Each of the three integral terms in (6.13) has the form described in the previous Lemma 6.5, because by Lemma 3.1 the integral kernel $K_\infty = -\ln|\sin(\pi x/2)|$, and hence it satisfies all the conditions (6.7). As in Lemma 6.5 we need to maximize the integral (6.13) by varying the characteristic functions $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$. The only difference is the additional constraint that χ_α and χ_β are random characteristic functions. This additional constraint, loosely speaking, requires that the maximizers of (6.13) are “random intervals” still centered at $x = 0$ or $x = 1$. More precisely, note that the values of the characteristic functions $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ in (6.13) are important only where $\chi_\alpha = 1$ and $\chi_\beta = 1$, respectively. Hence for a sequence of Bernoulli random variables $\theta_i, i = 1, \dots, n$ we can define characteristic functions of random intervals of length $2q_\alpha$ on $[-1, 1]$ centered at $x = 0$ and $x = 1$ as a product of two characteristic functions $\chi_\alpha(x)\chi_{q_\alpha}^1(x)$ and $\chi_\alpha(x)\chi_{q_\alpha}^2(x)$, respectively, where

$$(6.14) \quad \chi_{q_\alpha}^1(x) = \begin{cases} 1, & -\frac{q_\alpha}{q} \leq x \leq \frac{q_\alpha}{q}, \\ 0, & \text{otherwise,} \end{cases} \quad \chi_{q_\alpha}^2(x) = \begin{cases} 0, & -1 + \frac{q_\alpha}{q} < x < 1 - \frac{q_\alpha}{q}, \\ 1, & \text{otherwise.} \end{cases}$$

Similarly, functions $\chi_\alpha(x)\chi_{q_\alpha}^1(x)$ and $\chi_\alpha(x)\chi_{q_\alpha}^2(x)$ are random intervals of length $2q_\beta$ centered at $x = 0$ and $x = 1$ where

$$(6.15) \quad \chi_{q_\beta}^1(x) = \begin{cases} 1, & -\frac{q_\beta}{1-q} \leq x \leq \frac{q_\beta}{1-q}, \\ 0, & \text{otherwise,} \end{cases} \quad \chi_{q_\beta}^2(x) = \begin{cases} 0, & -1 + \frac{q_\beta}{1-q} < x < 1 - \frac{q_\beta}{1-q}, \\ 1, & \text{otherwise.} \end{cases}$$

For an illustration see Figure 6.1. The above discussion is made rigorous by the following.

LEMMA 6.6 (randomized Riesz rearrangement inequality for asymptotic energy minimizers). *Consider Bernoulli trials (4.3). Suppose χ_α and χ_β , $\chi_\alpha + \chi_\beta = 1$ are (random) characteristic functions of the angle distributions $\theta = \alpha$ and $\theta = \beta$, respectively. Let $q_\alpha(\theta)$ and $q_\beta(\theta)$ be random variables of θ with values*

$$0 \leq q_\alpha(\theta) \leq q, \quad 0 \leq q_\beta(\theta) \leq 1 - q.$$

Suppose $K(x)$ satisfies (6.7). Then for every $\delta > 0$ there exists $\delta' > 0$ so that if

$$(6.16) \quad |q_\alpha(\theta) - q_\alpha| < \delta', \quad |q_\beta(\theta) - q_\beta| < \delta'$$

for some fixed q_α and q_β , then almost surely⁸ as $n \rightarrow \infty$ the maximizers of

(6.17)

$$\begin{aligned} & \max_{\tilde{\chi}_\alpha, \tilde{\chi}_\beta} \int_{-1}^1 \int_{-1}^1 (a\varepsilon_\alpha(x) + b\varepsilon_\beta(x))(a\varepsilon_\alpha(t) + b\varepsilon_\beta(t))K(x-t)dxdt, \\ \varepsilon_\alpha &= (\chi_\alpha \tilde{\chi}_\alpha - q_\alpha(\boldsymbol{\theta})), \varepsilon_\beta = (\chi_\beta \tilde{\chi}_\beta - q_\beta(\boldsymbol{\theta})), a > 0, b > 0 \\ \chi^\alpha(x) &= \chi^\alpha(-x), \chi^\beta(x) = \chi^\beta(-x), \int_0^1 \chi^\alpha \tilde{\chi}_\alpha dx = q_\alpha(\boldsymbol{\theta}), \int_0^1 \chi^\beta \tilde{\chi}_\beta dx = q_\beta(\boldsymbol{\theta}) \end{aligned}$$

satisfy

(6.18) $\min(e_1, e_2) < \delta,$

where

(6.19)

$$e_1 = \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^1| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^1| \chi_\beta) dx, e_2 = \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^2| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^2| \chi_\beta) dx,$$

$\chi_{q_\alpha}^i$ and $\chi_{q_\beta}^i, i = 1, 2$ are defined in (6.14) and (6.15), respectively.

The proof of Lemma 6.6 is by contradiction to the law of large numbers and it is provided in Appendix C.2. The next lemma shows the equivalence of Definitions 6.2 and 6.3 as $n \rightarrow \infty$.

LEMMA 6.7. Consider the Bernoulli trials (4.3). For a given U , let q_α and q_β be defined as in Lemma 5.3 by (5.15) or (5.16). Then the probability measure μ_n^a of the asymptotic energy minimizer of the misfit energy (6.1) (see Definition 6.2)) and the probability measure μ_n^s of the Riesz symmetrically rearranged energy minimizer (Definition 6.3) have the same weak limit μ^s . Moreover,

(6.20) $\lim_{n \rightarrow \infty} e_1 = 0, \text{ or } \lim_{n \rightarrow \infty} e_2 = 0 \text{ with equal probability } 1/2,$

where e_1 and e_2 are defined by (6.19) in Lemma 6.17 above.

Proof. By Lemmas 5.3 and 5.5, depending on U , the minimizer of the first term in (6.1) satisfies (5.15) or (5.16) for almost every $\boldsymbol{\theta}$ as $n \rightarrow \infty$. Since the condition (6.16) of the randomized Riesz rearrangement inequality is satisfied for any $\delta > 0$, by Lemma 6.17 for almost every $\boldsymbol{\theta}$,

$$\lim_{n \rightarrow \infty} \min(e_1, e_2) = 0.$$

Hence $\{(s_i, I_i)\}_{i=1}^n$ is either a left-rearranged or right-rearranged sequence almost surely as $n \rightarrow \infty$. For every n the measure μ_n^a must be symmetric with respect to the to reflection about the point $x = 1/2$, i.e. with equal probability either $e_1 \rightarrow 0$ or $e_2 \rightarrow 0$ as $n \rightarrow \infty$. This proves (6.20). Clearly (6.20) implies that μ_n^a and the Riesz symmetrically rearranged measure μ_n^s have the same weak limit as $n \rightarrow \infty$. \square

End of proof of Theorem 6.4. Again, by the symmetry of the problem with respect to reflection about the point $x = 1/2$, μ_n also must be similarly symmetric. By Lemma 6.1 the minimizer of (6.1) is an almost minimizer, i.e. for any $\delta > 0$ there is L_0 so that the condition (6.16) of the randomized Riesz rearrangement inequality holds.

⁸Here and in the sequel almost sure convergence is considered in the probability space (Ω, \mathcal{F}, P) set in the beginning of section 4.

Hence it implies that for sufficiently large L_0 the minimizing sequences $\{(s_i, I_i)\}_{i=1}^n$ for any $L > L_0$ are arbitrarily close to the Riesz symmetrically rearranged minimizing sequences, namely either

$$(6.21) \quad \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^1| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^1| \chi_\beta) dx < \delta, \text{ or } \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^2| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^2| \chi_\beta) dx < \delta,$$

with equal probability $1/2$ as $n \rightarrow \infty$. If $L_0 \rightarrow \infty$, then $\delta \rightarrow 0$, and this completes the proof.

It follows from (6.21) that long/short-range correlations in the minimizing sequences $\{(s_i, I_i)\}_{i=1}^n$ are determined by long/short-range correlations of the Riesz symmetrically rearranged minimizing sequences. There is no short-range correlations of the Riesz symmetrically rearranged minimizing sequences. Riesz rearranged measure is, however, correlated on the large-scale: For example, suppose $q_\alpha = q$ and $q_\beta < 1 - q$, then for the right-rearranged measure:

$$(6.22) \quad (s, I) = \begin{cases} (\sin^2 \alpha, 1), & \text{with probability } q, \\ (\sin^2 \beta, 1), & \text{with probability } 1 - q, \text{ if } 1 - q_\beta/(1 - q) \leq x \leq 1, \\ (0, 0), & \text{with probability } 1 - q, \text{ if } 0 \leq x < 1 - q_\beta/(1 - q), \end{cases}$$

and for the left-rearranged measure:

$$(6.23) \quad (s, I) = \begin{cases} (\sin^2 \alpha, 1), & \text{with probability } q, \\ (\sin^2 \beta, 1), & \text{with probability } 1 - q, \text{ if } 0 \leq x < q_\beta/(1 - q), \\ (0, 0), & \text{with probability } 1 - q, \text{ if } q_\beta/(1 - q) \leq x \leq 1. \end{cases}$$

The long-range correlation of transformation strains for the symmetrically rearranged minimizer probability measure can be read off the formulas (6.22) and (6.23).

7. Statistics of asymptotic energy minimizers 3: Chain of short grains.

7.1. Basic definitions and ideas. In this section we will show how exponentially decaying correlations arise when the scaling of the polycrystal is such that $L = l_0/(2n)$, $l_0 \ll 1$ (short grains). By estimate (3.16) the misfit energy for $l_0 \ll 1$ is given (up to higher order terms) by the nearest neighbor energy

$$(7.1) \quad W_n(U, \boldsymbol{\varepsilon}^T) = \langle (U - f + s)^2 \rangle - \lambda_0 \langle (s - \langle s \rangle)^2 \rangle - \frac{Bl_0}{n} \sum_{i=1}^n (s_i - \langle s \rangle)(s_{i+1} - \langle s \rangle),$$

where $B > 0$, $\lambda_0 > 0$. In this case, we show that for Bernoulli trials (4.3) exponentially decaying correlations arise when $n \rightarrow \infty$, followed by $l_0 \rightarrow 0$.

Qualitatively, the misfit energy $W_n^1(U, \boldsymbol{\varepsilon}^T)$ has three terms which are analogous to the case of thin long grains (6.1). The minimization of the first two terms,

$$(7.2) \quad W_n^0(U, \boldsymbol{\varepsilon}^T) = \langle (U - f + s)^2 \rangle - \lambda_0 \langle (s - \langle s \rangle)^2 \rangle$$

determines, as in Lemma 5.3, q_α and q_β , the *total* amount of the grains that undergo a stress-free transformation. The minimizers of (7.2) are, in general, not unique. The third term provides a small correction to (7.2), and, as in section 6, plays a role of regularization, that is it selects the unique minimizer of (7.2) that maximizes

$$(7.3) \quad \bar{W}_n(U, \boldsymbol{\varepsilon}^T) := \frac{1}{n} \sum_{i=1}^n (s_i - \langle s \rangle)(s_{i+1} - \langle s \rangle).$$

Analogous to Definition 6.2, the above considerations motivate the following definition:

DEFINITION 7.1. For a fixed sequence $\theta_i, i = 1, \dots, n$, the asymptotic energy minimizer of the nearest neighbor model (7.1) is the sequence $(s_i, I_i), i = 1, \dots, n$ such that it minimizes (7.2) and maximizes (7.3) among minimizers of (7.2).

LEMMA 7.2. For the Bernoulli trials model (4.3) the minimizing sequence $(s_i, I_i), i = 1, \dots, n$ of the misfit energy $W_n^0(U, \epsilon^T)$ given by (7.2) satisfies

$$(7.4) \quad s(\theta) = \sin^2 \theta,$$

and as $n \rightarrow \infty$,

$$q_\alpha(\theta) \rightarrow q_\alpha, q_\beta(\theta) \rightarrow q_\beta$$

almost surely. The values q_α and q_β are determined as follows. For a given U the total proportion of grains that undergoes a stress-free transformation $f = f(U)$ is a (deterministic) nondecreasing function of U . For a given f we have several cases:

- if $f < 1 - q$, then

$$(7.5) \quad q_\alpha = 0, 0 \leq q_\beta \leq 1 - q,$$

$$\chi(\alpha) = 0, \chi(\beta) = q_\beta / (1 - q),$$

- if $f > 1 - q$, then

$$(7.6) \quad 0 \leq q_\alpha \leq q, q_\beta = 1 - q,$$

$$\chi(\beta) = q_\alpha / q, \chi(\beta) = 1.$$

There are values of α, β , and U for which $0 < q_\beta < 1 - q$ or $0 < q_\alpha < q$.

The proof of this lemma is analogous to Lemma 5.3. Clearly, the function $f(U)$ in Lemma 7.2 is different from the one for the infinitely long grains in Lemma 5.3. However, the characteristic property of the measure that it is determined by q_α and q_β with either (7.5) or (7.6) still holds. One of the consequences of the previous lemma is that there are, again, some values q_α and q_β that determine the proportion of grains that undergo a stress-free transformation and they satisfy $q_\alpha = 0$ or $q_\beta = 1 - q$. This is exactly the characteristic property that we need to be able to prove exponential decay of correlations by applying the isoperimetric inequality (7.8) to the sequences described in Definition 7.1. Following the logic in section 6, we obtain that asymptotic energy minimizer of the nearest neighbor model (7.1) arises in the limit $n \rightarrow \infty$, followed by $l_0 \rightarrow \infty$. The proof of this statement is similar to the proof of the analogous statement in case 2; see end of the proof of Theorem 6.4 in section 6.2. Hence we only need to find a statistical characterization of maximizers of (7.3) for fixed q_α and q_β found from Lemma 7.2. This is given in Theorems 7.4 and 7.5 below.

7.2. Isoperimetric inequalities and characterization of maximizers of (7.3). Here, it is convenient to characterize any point in the composite $x \in [0, 1]$ as a point that belongs to a (maximal) uninterrupted string of identical values of θ .

DEFINITION 7.3. For a fixed $\theta = \theta_1, \theta_2, \dots, \theta_n$ we say that a string

$$\theta_\alpha^m = \theta_{i+1}, \theta_{i+2}, \dots, \theta_{i+m}, \theta_\alpha^m \subset \theta$$

is a (maximal) uninterrupted string of $\theta = \alpha$ of length m if all $\theta_{i+j} = \alpha, i = 1, \dots, m$ and $\theta_i = \theta_{i+1+m} = \beta$. We say that $x \in [0, 1]$ belongs to an uninterrupted string of values α of length m if $x \in \theta_\alpha^m$. The notion $x \in \theta_\beta^m$ is defined analogously.

Recall our notation

$$s(x) - \langle s \rangle = \sin^2 \alpha \varepsilon_\alpha + \sin^2 \beta \varepsilon_\beta,$$

$$\varepsilon_\alpha = (\chi_\alpha \tilde{\chi}_\alpha - q_\alpha), \quad \varepsilon_\beta = (\chi_\beta \tilde{\chi}_\beta - q_\beta),$$

where χ_α and χ_β , $\chi_\alpha + \chi_\beta = 1$ are (random) characteristic functions of the angle distributions $\theta = \alpha$ and $\theta = \beta$, respectively; $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ are the characteristic function of the grains with $\theta = \alpha$ and $\theta = \beta$, respectively, for which $I = 1$.

By Lemma 7.2, we have two cases: Either $q_\beta = 1 - q$, and then $\tilde{\chi}_\beta \equiv 1$, or $q_\alpha = 0$, and then $\tilde{\chi}_\alpha \equiv 0$. Let us study these two cases separately.

Suppose $q_\alpha = 0$. Let us look at maximization of

$$(7.7) \quad \frac{1}{n} \sum_{i=1}^n s_i s_{i+1}, \quad \langle s \rangle = q_\beta \sin^2 \beta$$

only. Each of s_i (up to the constant $\sin^2 \beta$) is either 1 or 0; therefore the maximization of the nearest neighbors term (7.7) can be understood as the minimization of the boundary of a set with constant area:

$$(7.8) \quad \min_{D \in \mathcal{A}} \partial D, \quad \mathcal{A} = \left\{ D \mid D = \{x \mid s_i(x) \neq 0\}, \int_D dx = q_\beta \right\}.$$

Then the usual isoperimetric inequality implies that the maximizer of (7.7) is such that the grains with $\theta_i = \beta$ undergo a stress-free transformation, if they belong to a “long” uninterrupted sequence θ_β^m of the grains with the same $\theta = \beta$. If $\theta_i = \beta$ belongs to a “short” uninterrupted sequence θ_β^m , then it does not undergo a stress-free transformation. Hence there should be short-range correlations. The notion of short and long sequences is relative to the value of the total number of grains that must undergo a stress-free transformation. The above ideas are formulated more precisely in the next theorem.

THEOREM 7.4. *Consider the Bernoulli trials (4.3). Denote by $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ the characteristic function of the grains with $\theta = \alpha$ and $\theta = \beta$, respectively, for which $I = 1$. Suppose U is such that the minimizer of the first term in (7.1) satisfies $q_\alpha = 0$. Then in the limit $n \rightarrow \infty$ the sequence (s_i, I_i) , $i = 1, 2, \dots, n, \dots$ is a stationary process with exponentially decaying short-range correlations, long-range correlations are zero, and $\tilde{\chi}_\alpha \equiv 1$. Moreover, almost surely as $n \rightarrow \infty$ the minimizer of the nearest neighbor model (7.1) satisfies*

$$\tilde{\chi}_\beta(x) = \begin{cases} 1, & \text{if } x \in \theta_\beta^m, \quad m > k, \\ 0, & \text{if } x \in \theta_\beta^m, \quad m < k, \\ 1, & \text{with probability } r, \quad \text{if } x \in \theta_\beta^k, \\ 0, & \text{with probability } 1 - r, \quad \text{if } x \in \theta_\beta^k, \end{cases}$$

where

$$(7.9) \quad k = \max(m) \text{ such that } q_\beta < (1 - q)^m,$$

and r is found from

$$(7.10) \quad q_\beta = r q (1 - q)^k - (1 - q)^{k+1}.$$

Proof. By the law of the large numbers, $\tilde{\chi}_\alpha = 0$ and $q_\beta(\boldsymbol{\theta}) \rightarrow q_\beta$ almost surely; therefore it is sufficient to study (7.7) or, equivalently, (7.8). By the isoperimetric inequality for every $\boldsymbol{\theta}$ the function $\tilde{\chi}_\beta(x)$ must be such that if $\theta_\beta^{m_1}, \theta_\beta^{m_2}, \dots, \theta_\beta^{m_t}$ are all the (maximal) uninterrupted sequences $\theta_\beta^{m_i} \in \boldsymbol{\theta}$ of values β ordered so that the indices are decreasing $m_1 \geq m_2 \geq \dots \geq m_t$, then there exists an $i: 1 \leq i \leq t$ so that $\tilde{\chi}(x) = 1$, if $x \in \theta_\beta^{m_j}, j \leq i$, and $\tilde{\chi}(x) = 0$, if $x \in \theta_\beta^{m_j}, j > i$ with the exception of at most one $j \geq i$.

By construction as $n \rightarrow \infty$, the process $(s_i, I_i), i = 1, \dots, n, \dots$ is stationary. Since all $\theta_\alpha^m, \theta_\beta^m$ are geometrically distributed independent random variables [11], it means explicitly that the probability of a string θ_β^m is given by

$$\rho_\infty(\dots \theta_\beta^m \dots) = q(1 - q)^{m-1}.$$

Hence if the total proportion of grains that undergoes a stress-free transformation is q_β , we must have, as $n \rightarrow \infty$

$$q_\beta = rq(1 - q)^k + \sum_{i=k+1} q(1 - q)^i = rq(1 - q)^k + (1 - q)^{k+1},$$

where, due to our assumption of equal probability (4.1), r is the probability that $\tilde{\chi}(x) = 1$, if $x \in \theta_\beta^k$. Therefore k is found so that $(1 - q)^{k+1} \leq q_\beta < (1 - q)^k$, i.e., (7.9), and r is found from (7.10).

Since $\theta_\alpha^m, \theta_\beta^m$ are independent random variables, the limiting process has exponentially decaying short-range correlations. It implies simultaneously two results: long-range correlations are zero, and short-range correlations decay exponentially with k . These correlations are not zero and can be computed explicitly. \square

Suppose $q_\beta = 1 - q$. This case is slightly more technically complicated, but the methods are the same as in the case $q_\alpha = 0$. The main new issue is that s_i may now take *three* values, and, therefore, we have to account for three possible different interfaces. Direct computations show that we have here three different situations, depending on the relative value of α and β . If $\sin^2 \alpha > 2 \sin^2 \beta$, then the maximizer of (7.7) is such that $\tilde{\chi}_\alpha(x) = 1$ if x belongs to the *longest* (maximal) uninterrupted strings θ_α^m . If, however, $\sin^2 \alpha < 2 \sin^2 \beta$, then $\tilde{\chi}_\alpha(x) = 1$ if x belongs to the *shortest* uninterrupted strings θ_α^m . If $\sin^2 \alpha = 2 \sin^2 \beta$, then there is no difference, and the only statement that is possible to make here is that $\tilde{\chi}_\alpha(x) = \tilde{\chi}_\alpha(y)$, if x and y belong to the same uninterrupted string θ_α^m . Due to our assumption of equal probability (4.1), it is possible to conclude that if $\sin^2 \alpha = 2 \sin^2 \beta$, then there is no correlation at all, therefore we will omit the discussion of this case. Combining these arguments with the arguments in the proof of Theorem 7.4 we have the following result.

THEOREM 7.5. *Consider the Bernoulli trials (4.3). Denote by $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ the characteristic function of the grains with $\theta = \alpha$ and $\theta = \beta$, respectively, for which $I = 1$. Suppose U is such that the minimizer of the first term in (7.1) satisfies $q_\beta = 1 - q$. Then in the limit $n \rightarrow \infty$ the sequence $(s_i, I_i), i = 1, 2, \dots, n, \dots$ is a stationary process with exponentially decaying short-range correlations; long-range correlations are zero, and $\tilde{\chi}_\beta \equiv 1$. Moreover, almost surely as $n \rightarrow \infty$, the minimizer of the nearest neighbor model (7.1) satisfies: If $\sin^2 \alpha > 2 \sin^2 \beta$, then*

$$\tilde{\chi}_\alpha(x) = \begin{cases} 1, & \text{if } x \in \theta_\alpha^m, m > k, \\ 0, & \text{if } x \in \theta_\alpha^m, m < k, \\ 1, & \text{with probability } r, \text{ if } x \in \theta_\alpha^k, \\ 0, & \text{with probability } 1 - r, \text{ if } x \in \theta_\alpha^k, \end{cases}$$

where $k = \max(m)$ such that $q_\alpha < q^m$ and r solves $q_\alpha = r(1 - q)q^k - q^{k+1}$; if $\sin^2 \alpha < 2 \sin^2 \beta$, then

$$\tilde{\chi}_\alpha(x) = \begin{cases} 1, & \text{if } x \in \theta_\alpha^m, \ m < k, \\ 0, & \text{if } x \in \theta_\alpha^m, \ m > k, \\ 1, & \text{with probability } r, \ \text{if } x \in \theta_\alpha^k, \\ 0, & \text{with probability } 1 - r, \ \text{if } x \in \theta_\alpha^k, \end{cases}$$

where $k = \max(m)$ such that $q - q_\alpha < q^m$ and r solves $q - q_\alpha = r(1 - q)q^k - q^{k+1}$.

Appendix A. Proofs for section 3.

A.1. Proof of formula (3.4). To obtain the representation (3.4) we begin by decomposing $\mathbf{u} = (u_1, u_2)$ of the boundary value problem (2.11), (2.13), and (2.14) in the form $\mathbf{u} = \tilde{\mathbf{u}} + \bar{\mathbf{u}}$, where $\tilde{\mathbf{u}}$ solves (3.6). The constants a_i, c_i, d, f_i , and g_i are chosen to satisfy the continuity of the displacement $\tilde{\mathbf{u}}$ and traction (condition (2.12)) and, denoting

$$\tilde{\sigma}_{ij} = c_{ijkl}(\tilde{\varepsilon}_{kl} - \varepsilon_{kl}^T), \quad \tilde{\varepsilon}_{11} = \partial_1 \tilde{u}_1, \quad \tilde{\varepsilon}_{12} = \tilde{\varepsilon}_{21} = (\partial_2 \tilde{u}_1 + \partial_1 \tilde{u}_2)/2, \quad \tilde{\varepsilon}_{22} = \partial_2 \tilde{u}_2,$$

the boundary conditions

$$\tilde{u}_1(0, y) = 0, \quad \tilde{u}_1(1, y) = U, \quad \tilde{u}_2(0, 0) = 0, \quad \tilde{\sigma}_{12}(x, y) = 0, \quad \text{for } x = 0, 1.$$

The stresses are

$$\begin{aligned} \tilde{\sigma}_{i,11} &= (\lambda + 2G)(a_i - \varepsilon_{i,11}^T) + \lambda(d - \varepsilon_{i,22}^T), \\ \tilde{\sigma}_{i,22} &= (\lambda + 2G)(d - \varepsilon_{i,22}^T) + \lambda(a_i - \varepsilon_{i,11}^T), \\ \tilde{\sigma}_{i,12} &= G(c_i - 2\varepsilon_{i,12}^T). \end{aligned}$$

Hence $[\tilde{\sigma}_{12}] = 0$, if $c_i = 2\varepsilon_{i,12}^T$. Similarly $[\tilde{\sigma}_{i,11}] = 0$, if

$$a_i = \varepsilon_{i,11}^T + \frac{\lambda}{\lambda + 2G} \varepsilon_{i,22}^T + \text{Const},$$

where the constant can be found from the condition that the displacement of the right boundary is U :

$$\text{Const} = U - \langle \varepsilon_{11}^T \rangle - \frac{\lambda}{\lambda + 2G} \langle \varepsilon_{22}^T \rangle.$$

Finally,

$$\begin{aligned} \langle \tilde{\sigma}_{22} \rangle &= \langle (\lambda + 2G)(d - \varepsilon_{22}^T) + \lambda \left(\frac{\lambda}{\lambda + 2G} \varepsilon_{22}^T + U - \langle \varepsilon_{11}^T \rangle - \frac{\lambda}{\lambda + 2G} \langle \varepsilon_{22}^T \rangle \right) \rangle \\ &= (\lambda + 2G)(d - \langle \varepsilon_{22}^T \rangle) + \lambda(U - \langle \varepsilon_{11}^T \rangle). \end{aligned}$$

Setting $\langle \tilde{\sigma}_{22} \rangle = 0$ we have

$$d = \langle \varepsilon_{22}^T \rangle - \frac{\lambda}{\lambda + 2G} (U - \langle \varepsilon_{11}^T \rangle).$$

The values of f_i and g_i in (3.6) are unimportant for our analysis, and we omit them. We have

$$(A.1) \quad \tilde{\sigma}_{11} = E_c(U - \langle \varepsilon_{11}^T \rangle), \quad \tilde{\sigma}_{i,22} = E_c(\langle \varepsilon_{22}^T \rangle - \varepsilon_{i,22}^T), \quad \tilde{\sigma}_{12} = 0,$$

where the Young’s modulus E_c is given by (2.16). The elastic misfit energy associated with $\tilde{\mathbf{u}}$ is

$$\frac{1}{E_c} \frac{1}{4L} \int_{-L}^L \int_{-1}^1 \left[\frac{1}{E_c} ((\tilde{\sigma}_{11})^2 + (\tilde{\sigma}_{22})^2) - 2 \frac{\lambda}{4G(\lambda + G)} \tilde{\sigma}_{11} \tilde{\sigma}_{22} + \frac{1}{2G} (\tilde{\sigma}_{12})^2 \right] dx dy.$$

Since $\tilde{\sigma}_{12} \equiv 0$, $\tilde{\sigma}_{11} = \text{const}$, and $\tilde{\sigma}_{22}$ is mean-zero, the above equation becomes

$$\frac{1}{E_c^2} \frac{1}{4L} \int_{-L}^L \int_{-1}^1 ((\tilde{\sigma}_{11})^2 + (\tilde{\sigma}_{22})^2) dx dy = \langle (U - f + \varepsilon_{22}^T)^2 \rangle = \langle (U - f + s)^2 \rangle,$$

which is the first term in (3.4).

Let us now find $\bar{\mathbf{u}}$. It solves (3.7) on a bounded domain Π_L . A useful periodic setting for (3.7) is obtained by assuming this equation is posed on an infinite strip $]-\infty, \infty[\times]-L, L[$ with data that is even and periodic in x :

$$(A.2) \quad \begin{aligned} (a) \quad & \theta(-x) = \theta(x), \quad \theta(x + 2) = \theta(x), \\ (b) \quad & \varepsilon^T(x, y) = \varepsilon^T(-x, y), \quad \varepsilon^T(x + 2, y) = \varepsilon^T(x, y). \end{aligned}$$

Thus, $\bar{\mathbf{u}}$ equals to the restriction of the solution of (3.7) on an infinite strip $x \in]-\infty, \infty[$, $y \in]-L, L[$ with periodicity conditions defined by (A.2). Since $\bar{\sigma}_{22}(x, \pm L)$ is a periodic, mean-zero, even function, it can be represented as a cosine Fourier series. The solution on the infinite strip with a sinusoidal symmetric stress $\cos(k\pi x)$ at $y = \pm L$ can be computed explicitly for any k by the Airy function method. Namely, since we are given that $\bar{\sigma}_{22}(x, \pm L)$ is a periodic, mean-zero, even function, it can be represented as

$$\bar{\sigma}_{22}(x, \pm L) = E_c \sum_{k=1}^{\infty} c_k \cos(k\pi x), \quad E_c = \frac{4G(\lambda + G)}{\lambda + 2G},$$

where c_k are the corresponding Fourier coefficients of $\varepsilon^T(x)$. The solution for the infinite strip with a sinusoidal symmetric stress $\cos(k\pi x)$ at $y = \pm L$ is given (see [19]) by the Airy function

$$\begin{aligned} & \Phi^k(x, y) \\ &= 2 \frac{\cos(k\pi x) k\pi y \sinh(k\pi L) \sinh(k\pi y) - [k\pi L \cosh(k\pi L) + \sinh(k\pi L)] \cosh(k\pi y)}{(k\pi)^2 (2k\pi L + \sinh(2k\pi L))}. \end{aligned}$$

This Airy function gives rise to the following stresses

$$\begin{aligned} \bar{\sigma}_{11}^k &= \frac{\partial^2 \Phi^k}{\partial y^2} = d_{11}^k \cos(k\pi x), \quad \bar{\sigma}_{12}^k = -\frac{\partial^2 \Phi^k}{\partial x \partial y} \\ &= d_{12}^k \sin(k\pi x), \quad \bar{\sigma}_{22}^k = \frac{\partial^2 \Phi^k}{\partial x^2} = d_{22}^k \cos(k\pi x), \end{aligned}$$

where

$$(A.3) \quad \begin{aligned} d_{11}^k &= 2 \frac{k\pi y \sinh(k\pi L) \sinh(k\pi y) - [k\pi L \cosh(k\pi L) - \sinh(k\pi L)] \cosh(k\pi y)}{2k\pi L + \sinh(2k\pi L)}, \\ d_{12}^k &= 2 \frac{k\pi y \sinh(k\pi L) \cosh(k\pi y) - k\pi L \cosh(k\pi L) \sinh(k\pi y)}{2k\pi L + \sinh(2k\pi L)}, \\ d_{22}^k &= 2 \frac{-k\pi y \sinh(k\pi L) \sinh(k\pi y) + [k\pi L \cosh(k\pi L) + \sinh(k\pi L)] \cosh(k\pi y)}{2k\pi L + \sinh(2k\pi L)}. \end{aligned}$$

Therefore the total stresses are

$$\begin{aligned} \sigma_{11} &= E_c(U - \langle \varepsilon_{11}^T \rangle) + E_c \sum_{k=1}^{\infty} c_k d_{11}^k \cos(k\pi x), \\ \sigma_{12} &= E_c \sum_{k=1}^{\infty} c_k d_{12}^k \sin(k\pi x), \quad \sigma_{22} = E_c \sum_{k=1}^{\infty} c_k (d_{22}^k - 1) \cos(k\pi x), \end{aligned}$$

where in the last equation we have $d_{22}^k - 1$ instead of d_{22}^k , because (see (A.1))

$$\tilde{\sigma}_{22} = -E_c \sum_{k=1}^{\infty} c_k \cos(k\pi x).$$

By definition

$$W = \frac{1}{E_c} \frac{1}{4L} \int_{-L}^L \int_{-1}^1 \left[\frac{1}{E_c} ((\sigma_{11})^2 + (\sigma_{22})^2) - 2 \frac{\lambda}{4G(\lambda + G)} \sigma_{11} \sigma_{22} + \frac{1}{2G} (\sigma_{12})^2 \right] dx dy.$$

Since $\int_{-1}^1 \cos(k\pi x) \cos(m\pi x) dx = \delta_{km}$

$$W(U, \varepsilon^T) = \langle (U - f + \varepsilon_{22}^T)^2 \rangle + \sum_{m=1}^{\infty} c_m^2 \hat{K}_L(m),$$

where

$$(A.4) \quad \hat{K}_L(m) = \frac{1}{4L} \int_{-L}^L \left[(d_{11}^m)^2 + d_{22}^m (d_{22}^m - 2) - \frac{2\lambda d_{11}^m (d_{22}^m - 1)}{\lambda + 2G} + \frac{2(\lambda + G)(d_{12}^m)^2}{\lambda + 2G} \right] dy.$$

Denoting $a = 2\pi mL$ and using Mathematica[®] we obtain an explicit form of (A.4):

$$(A.5) \quad \hat{K}_L(m) = \frac{5\lambda + 9G}{2(\lambda + 2G)} S_1(a) - \frac{\lambda + G}{\lambda + 2G} S_2(a), \quad a = 2\pi mL, \quad \text{where}$$

$$(A.6) \quad S_1(a) = \frac{\cosh(a) - 1}{a(a + \sinh(a))}, \quad S_2(a) = \frac{a^2(2 + \cosh(a))}{6(a + \sinh(a))^2}.$$

A.2. Proof of Lemma 3.1. Using (A.5) and (A.6) from Appendix A.1, direct computations show that the Fourier coefficients of $K_L(x)$ (3.5) are given by

$$(A.7) \quad \hat{K}_L(m) = \frac{B}{L} \frac{1}{m} + O(\exp(-L)),$$

so that, defining

$$(A.8) \quad K_{\infty}(x) = \sum_{m=1}^{\infty} \frac{1}{m} \cos(\pi m x),$$

(3.8) is satisfied. A closed form expression for this sum is known: $K_\infty(x) = -\ln|\sin(\pi x/2)|$ (see, e.g., [10]), and the Lemma follows.

A.3. Proof of Lemma 3.2. Set $L = l_0/(2n)$, then for a defined in (A.5) we obtain $a = 2\pi mL = l_0\pi m/n$. Substitute (A.6) into (A.5) and observe that the Fourier coefficients $\hat{K}_L(m)$ depend on the variable $a = l_0\pi m/n$. Thus we can introduce the notation

$$\hat{K}(a) := \hat{K}_L(m), \quad a = l_0\pi \frac{m}{n}.$$

In other words the Fourier series of $K_L(x)$ can be written in the form

$$K_L(x) = \sum_{m=-\infty}^{\infty} \hat{K}\left(l_0\pi \frac{m}{n}\right) \cos(\pi mx).$$

For $y = nx$, $y \in [-n, n]$ let

$$K^{(n)}(y) := \frac{1}{n}K_L(y/n) = \frac{1}{n} \sum_{m=-\infty}^{\infty} \hat{K}\left(l_0\pi \frac{m}{n}\right) \cos\left(\pi \frac{m}{n}y\right),$$

and $K^{(n)}(y) = 0$ for $|y| \geq n$. Set

$$\mathcal{K}_{l_0}(y) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{K}(l_0\zeta) \cos(\zeta y) d\zeta, \quad y \in \mathbb{R}.$$

Note that

$$\hat{K}(l_0\zeta) = \int_{-\infty}^{\infty} \mathcal{K}_{l_0}(y) \cos(\zeta y) dy,$$

provided $\mathcal{K}_{l_0}(y)$ is smooth and it decays sufficiently fast as $y \rightarrow \infty$. The function $\mathcal{K}_{l_0}(y)$ is exactly the limiting function mentioned in Lemma 3.2; that is, $K^{(n)}(y) = K_L(y/n)/n$ converges to $\mathcal{K}_{l_0}(y)$ as $n \rightarrow \infty$. Let us first verify properties of \mathcal{K}_{l_0} described in Lemma 3.2 and then establish convergence. Direct calculations using (A.5) show that $\hat{K}(a) > 0$ for all physical choices of the Lamé constants, so that, by Bochner’s theorem [18, Vol. 1], \mathcal{K}_{l_0} is positive-definite. We verify (3.10) by applying the Paley–Wiener type Theorems [18, Vol. 1]. Indeed, from (A.6) it follows that the Fourier transform of \mathcal{K}_{l_0} can be analytically extended into a finite strip $|\text{Im}(a)| \leq c_2$ around the real axis provided there is no solution of the equation $a + \sinh(a) = 0$, $a \neq 0$, or, equivalently

$$(A.9) \quad \text{Re}(a) = -\sinh(\text{Re}(a)) \cos(\text{Im}(a)), \quad \text{Im}(a) = -\cosh(\text{Re}(a)) \sin(\text{Im}(a)), \quad a \neq 0.$$

There is no solution of the last equation in (A.9) at least in the strip $|\text{Im}(a)| \leq \pi$. This implies the exponential decay (3.10) of $\mathcal{K}_{l_0}(y)$ by Paley–Wiener theorems. Finally, for every n consider the $2n$ -periodization of $\mathcal{K}_{l_0}(y)$:

$$\mathcal{K}_{l_0}^{(n)}(y) = \sum_{k=-\infty}^{+\infty} \mathcal{K}_{l_0}(y + 2nk).$$

Since $\mathcal{K}_{l_0}(y)$ decays exponentially as $y \rightarrow \infty$, we have that $\mathcal{K}_{l_0}^{(n)}(y)$ is a smooth $2n$ periodic function and its Fourier coefficients

$$\begin{aligned} \hat{\mathcal{K}}_{l_0}^{(n)}(m) &= \frac{1}{n} \int_{-n}^n \mathcal{K}_{l_0}^{(n)}(y) \cos\left(\pi \frac{m}{n} y\right) dy \\ &= \frac{1}{n} \int_{-n}^n \left(\sum_{k=-\infty}^{+\infty} \mathcal{K}_{l_0}(y + 2nk) \right) \cos\left(\pi \frac{m}{n} y\right) dy \\ &= \frac{1}{n} \int_{-\infty}^{\infty} \mathcal{K}_{l_0}(y) \cos\left(\pi \frac{m}{n} y\right) dy = \frac{1}{n} \hat{K}\left(l_0 \pi \frac{m}{n}\right) = \frac{1}{n} \hat{K}_L(m). \end{aligned}$$

Thus $\mathcal{K}_{l_0}^{(n)}(y) = K_L(y/n)/n = K^{(n)}(y)$ on $[-n, n]$. Therefore

$$\begin{aligned} \|K_L(x) - n\mathcal{K}_{l_0}(nx)\|_{L_\infty([-1,1])} &= n\|K^{(n)}(y) - \mathcal{K}_{l_0}(y)\|_{L_\infty([-n,n])} \\ &= n\|\mathcal{K}_{l_0}^{(n)}(y) - \mathcal{K}_{l_0}(y)\|_{L_\infty([-n,n])} \\ &= n \left\| \sum_{k \neq 0} \mathcal{K}_{l_0}(y + 2nk) \right\|_{L_\infty([-n,n])} \leq Cn \sum_{k=1}^{\infty} e^{-ckn} \\ &\leq Cne^{-cn} \rightarrow 0, \end{aligned}$$

as $n \rightarrow \infty$, and $C = C(l_0) > 0$ is independent of n .

A.4. Proof of Proposition 3.3. For any U and ε^T we can estimate the error of the truncation (3.14) as

$$\begin{aligned} |W_n(U, \varepsilon^T) - W_n^{k_0}(U, \varepsilon^T)| &\leq c \sum_{k > k_0} |\lambda_k(n)| \langle (s - \langle s \rangle)^2 \rangle \\ &\leq c \exp(-k_0) \langle (s - \langle s \rangle)^2 \rangle \leq c \exp(-k_0). \end{aligned}$$

Since

$$W_n^{k_0}(U, \varepsilon_{k_0}^T) = \min_{\tilde{\varepsilon}^T \in \tilde{\Omega}_n(\theta)} W_n^{k_0}(U, \tilde{\varepsilon}^T),$$

we have

$$W_n(U, \varepsilon^T) \leq \min_{\tilde{\varepsilon}^T \in \tilde{\Omega}_n(\theta)} W_n(U, \tilde{\varepsilon}^T) + c \exp(-k_0) = W_n(U, \theta) + c \exp(-k_0).$$

By definition of the minimizer $W_n(U, \theta) \leq W_n(U, \varepsilon^T)$. This implies (3.15).

Appendix B. Proofs for section 5.

B.1. Verification of (5.11). For a given product measure $\tilde{\mu}$ we have

$$\int I_1 I_2 d\tilde{\mu} = \left(\int I_1 d\tilde{\mu} \right)^2, \quad \int s_1 I_2 d\tilde{\mu} = \left(\int s_1 d\tilde{\mu} \right) \left(\int I_1 d\tilde{\mu} \right),$$

because the events in the first and the second grains are independent and identically distributed. Hence

$$\begin{aligned} \int I_1 I_2 d\tilde{\mu} &= \left(\int_0^{\pi/4} \chi(\theta) d\rho(\theta) \right)^2 \\ &= f^2, \quad \int (U + s_1) I_2 d\tilde{\mu} = f \left(\int_0^{\pi/4} (U + s(\theta)) \chi(\theta) d\rho(\theta) + (1 - f)U \right). \end{aligned}$$

It gives

$$\begin{aligned}
 & \int F(U, s_1, I_1, I_2) d\tilde{\mu} \\
 &= \int ((U + s_1)^2 - 2(U + s_1)I_2 + I_1I_2) d\tilde{\mu} \\
 &= \int_0^{\pi/4} (U + s(\theta))^2 \chi(\theta) d\rho(\theta) + (1 - f)U^2 \\
 &\quad - 2f \left(\int_0^{\pi/4} (U + s(\theta)) \chi(\theta) d\rho(\theta) + (1 - f)U \right) + f^2 \\
 &= \int_0^{\pi/4} ((U + s(\theta) - f)^2 - f^2) \chi(\theta) d\rho(\theta) + (1 - f)(U^2 - 2Uf) + f^2 \\
 &= \int_0^{\pi/4} (U + s(\theta) - f)^2 \chi(\theta) d\rho(\theta) + (U - f)^2(1 - f),
 \end{aligned}$$

where the last equality is obtained by noting the following identity:

$$f^2 = f^2(1 - f) + \int_0^{\pi/4} f^2 \chi(\theta) d\rho(\theta).$$

B.2. Proof of Lemma 5.3. The proof is the direct evaluation and comparison of all possible scenarios in (5.12). Using (5.12), we simply consider four functions of U , q , q_α , and q_β :

$$\begin{aligned}
 W_1 &= (U - q_\alpha - q_\beta + \cos^2 \alpha)^2 q_\alpha \\
 &+ (U - q_\alpha - q_\beta + \cos^2 \beta)^2 q_\beta + (U - q_\alpha - q_\beta)^2 (1 - q_\alpha - q_\beta), \\
 W_2 &= (U - q_\alpha - q_\beta + \cos^2 \alpha)^2 q_\alpha \\
 &+ (U - q_\alpha - q_\beta + \sin^2 \beta)^2 q_\beta + (U - q_\alpha - q_\beta)^2 (1 - q_\alpha - q_\beta), \\
 W_3 &= (U - q_\alpha - q_\beta + \sin^2 \alpha)^2 q_\alpha \\
 &+ (U - q_\alpha - q_\beta + \cos^2 \beta)^2 q_\beta + (U - q_\alpha + q_\beta)^2 (1 - q_\alpha - q_\beta), \\
 W_4 &= (U - q_\alpha - q_\beta + \sin^2 \alpha)^2 q_\alpha \\
 &+ (U - q_\alpha - q_\beta + \sin^2 \beta)^2 q_\beta + (U - q_\alpha - q_\beta)^2 (1 - q_\alpha - q_\beta),
 \end{aligned}$$

and compare their values for each fixed U and q , where q_α and q_β are in the range (5.13). It is easy to check that the minimum is always achieved for W_4 ; hence equation (5.14) is satisfied.

Appendix C. Proofs for section 6.

C.1. Proof of Lemma 6.4. Since the Fourier coefficients of a convolution equal the products of the Fourier coefficients, and since from (A.8) we know the Fourier coefficients of K_∞ are ≤ 1 in absolute value, and in view of Plancherel’s theorem, we have

$$\int_{-1}^1 \int_{-1}^1 (s(x) - \langle s \rangle)(s(t) - \langle s \rangle)K_\infty(x - t)dxdt \leq \int_{-1}^1 (s(x) - \langle s \rangle)^2 dx,$$

and $|s(x)| \leq 1$, for any $\delta > 0$, there exists $L_0 > 0$, so that for any $L \geq L_0$

$$\left| \frac{B}{L} \bar{W}_n(U, \varepsilon^T) \right| \leq \delta.$$

Hence for any $s = \varepsilon_{22}^T$

$$\langle (U - f + s)^2 \rangle < W_n(U, \varepsilon^T) \leq \langle (U - f + s)^2 \rangle + \delta.$$

It follows that

$$\min_{s,f} \langle (U - f + s)^2 \rangle < \min_{s,f} W_n(U, \varepsilon^T) \leq \min_{s,f} \langle (U - f + s)^2 \rangle + \delta$$

for any sequence of angles θ . Applying Lemma 5.5 we complete the proof of (6.5). Equality (6.4) follows from direct computations as in Lemma 5.3.

C.2. Proof of Lemma 6.6. Suppose (6.18) and (6.19) do not hold. It means that there exists $\delta > 0$ such that for every $\delta' > 0$ there is a sequence of sets $A_{n_k} \in \sigma_{n_k}^\theta$, $\{n_k\} \rightarrow \infty$ with probability $\rho_{n_k}(A_{n_k}) > 2C > 0$ such that for every fixed $\theta \in A_{n_k}$ (or, equivalently, $(\chi_\alpha, \chi_\beta) \in A_{n_k}$) there is a (at least one) maximizer $\tilde{\chi}_\alpha, \tilde{\chi}_\beta$ of (6.17) such that

$$(C.1) \quad \min(e_1, e_2) > \delta,$$

where

$$(C.2) \quad e_1 = \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^1| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^1| \chi_\beta) dx, \quad e_2 = \int_0^1 (|\tilde{\chi}_\alpha - \chi_{q_\alpha}^2| \chi_\alpha + |\tilde{\chi}_\beta - \chi_{q_\beta}^2| \chi_\beta) dx;$$

$\chi_{q_\alpha}^i$ and $\chi_{q_\beta}^i$ are defined in (6.14) and (6.15), respectively. Since $K(x) \in L^1[-1, 1]$ is fixed and $\tilde{\chi}_\alpha \chi_\alpha, \tilde{\chi}_\beta \chi_\beta$ are uniformly bounded in $L^\infty[-1, 1]$, for any $\delta_1 > 0$ there exists $h > 0$ so that

$$(C.3) \quad \begin{aligned} & \left| \int_{-1}^1 \int_{-1}^1 \varepsilon_\alpha(x) \varepsilon_\alpha(t) K(x - y) dx dt - \int_{-1}^1 \int_{-1}^1 \bar{\varepsilon}_\alpha(x) \bar{\varepsilon}_\alpha(t) K(x - y) dx dt \right| < \delta_1, \\ & \left| \int_{-1}^1 \int_{-1}^1 \varepsilon_\alpha(x) \varepsilon_\beta(t) K(x - y) dx dt - \int_{-1}^1 \int_{-1}^1 \bar{\varepsilon}_\alpha(x) \bar{\varepsilon}_\beta(t) K(x - y) dx dt \right| < \delta_1, \\ & \left| \int_{-1}^1 \int_{-1}^1 \varepsilon_\beta(x) \varepsilon_\beta(t) K(x - y) dx dt - \int_{-1}^1 \int_{-1}^1 \bar{\varepsilon}_\beta(x) \bar{\varepsilon}_\beta(t) K(x - y) dx dt \right| < \delta_1, \end{aligned}$$

where

$$\bar{\varepsilon}_\alpha(x) = \frac{1}{2h} \int_{-h}^h \varepsilon_\alpha(x + t) dt, \quad \bar{\varepsilon}_\beta(x) = \frac{1}{2h} \int_{-h}^h \varepsilon_\beta(x + t) dt.$$

By the law of the large numbers, for any $\delta_2 > 0$ there is $\bar{A}_{n_k} \in A_{n_k}$, $\rho_{n_k}(A_{n_k}) > C$ so that for every $\theta \in \bar{A}_{n_k}$ (or, equivalently, $(\chi_\alpha^{n_k}, \chi_\beta^{n_k}) \in \bar{A}_{n_k}$) and $n_k > N_0$, we have

$$\bar{\varepsilon}_\alpha(x) = q\chi_\alpha(x, \theta) - q_\alpha(\theta), \bar{\varepsilon}_\beta(x) = (1 - q)\chi_\beta(x, \theta) - q_\alpha(\theta),$$

where (random in $\theta \in \bar{A}_{n_k}$) functions $\chi_\alpha(x, \theta)$, $\chi_\beta(x, \theta)$ satisfy

$$0 \leq \chi_\alpha(x, \theta) \leq 1 + \delta_2, \quad 0 \leq \chi_\beta(x, \theta) \leq 1 + \delta_2,$$

and $\min(e_1, e_2) > \delta - \delta_2$ where

$$\begin{aligned} e_1 &= \int_0^1 (|\chi_\alpha(x, \theta) - \chi_{q_\alpha}^1| + |\chi_\beta(x, \theta) - \chi_{q_\beta}^1|) dx, \quad e_2 \\ &= \int_0^1 (|\chi_\alpha(x, \theta) - \chi_{q_\alpha}^2| + |\chi_\beta(x, \theta) - \chi_{q_\beta}^2|) dx. \end{aligned}$$

The classical Riesz rearrangement inequality (6.11) implies that when δ_1 and δ_2 are sufficiently small, there is $\delta' > 0$ so that for every $\theta \in \bar{A}_{n_k}$

$$\begin{aligned} &\int_{-1}^1 \int_{-1}^1 (a\varepsilon_\alpha(x) + b\varepsilon_\beta(x))(a\varepsilon_\alpha(t) + b\varepsilon_\beta(t))K(x - t) dx dt \\ &\leq \int_{-1}^1 \int_{-1}^1 (a\varepsilon_\alpha^1(x) + b\varepsilon_\beta^1(x))(a\varepsilon_\alpha^1(t) + b\varepsilon_\beta^1(t))K(x - t) dx dt - \delta', \end{aligned}$$

where

$$\varepsilon_\alpha^1(x) = \chi_{q_\alpha}^1 \chi_\alpha, \quad \varepsilon_\beta^1(x) = \chi_{q_\beta}^1 \chi_\beta.$$

Hence we have that $\tilde{\chi}_\alpha$ and $\tilde{\chi}_\beta$ with (C.1) and (C.2) cannot be maximizers of (6.17). This leads to contradiction with our assumption that (6.18) and (6.19) do not hold.

Acknowledgments. We thank Alexei Borodin and Omri Sarig for useful discussions. We are grateful to anonymous referees for their useful suggestions.

REFERENCES

- [1] A. BAERNSTEIN, II, *Convolution and rearrangement on the circle*, Complex Variables Theory Appl., 12 (1989), pp. 33–37.
- [2] K. BHATTACHARYA, B. LI, AND M. LUSKIN, *The simply laminated microstructure in martensitic crystals that undergo a cubic-to-orthorhombic phase transformation*, Arch. Ration. Mech. Anal., 149 (1999), pp. 123–154.
- [3] K. BHATTACHARYA AND R. V. KOHN, *Elastic energy minimization and the recoverable strains of polycrystalline shape-memory materials*, Arch. Ration. Mech. Anal., 139 (1997), pp. 99–180.
- [4] J. S. BOWLES AND J. K. MACKENZIE, *The crystallography of martensite transformations I*, Acta Metall., 2 (1954), pp. 129–137.
- [5] J. S. BOWLES AND J. K. MACKENZIE, *The crystallography of martensite transformations III. Face-centered cubic to body-centered tetragonal transformations*, Acta Metall., 2 (1954), pp. 224–234.
- [6] O. P. BRUNO AND G. H. GOLDSZTEIN, *A fast algorithm for the simulation of polycrystalline misfits: Martensitic transformations in two space dimensions*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 455 (1999), pp. 4245–4276.
- [7] O. P. BRUNO AND G. H. GOLDSZTEIN, *Numerical simulation of martensitic transformations in two- and three-dimensional polycrystals*, The J. R. Willis 60th anniversary volume, J. Mech. Phys. Solids, 48 (2000), pp. 1175–1201.

- [8] O. P. BRUNO AND G. H. GOLDSZTEIN, *A fast algorithm for the simulation of polycrystalline misfits. II. Martensitic transformations in three space dimensions*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 460 (2004), pp. 1613–1630.
- [9] O. P. BRUNO, F. REITICH, AND P. LEO, *The overall elastic energy of polycrystalline martensitic solids*, J. Mech. Phys. Solids, 44 (1996), pp. 1051–1101.
- [10] B. DEMIDOVICH, ED., *Problems in Mathematical Analysis*, Peace Publishers, Moscow, 1965.
- [11] W. FELLER, *An Introduction to Probability Theory and Its Applications*, 3rd ed., Vol. I, John Wiley & Sons, New York-London-Sydney, 1968.
- [12] G. H. HARDY, J. E. LITTLEWOOD, AND G. PÓLYA, *Inequalities*, 2nd ed., Cambridge University Press, Cambridge, 1952.
- [13] E. HEWITT AND L. J. SAVAGE, *Symmetric measures on Cartesian products*, Trans. Amer. Math. Soc., 80 (1955), pp. 470–501.
- [14] E. H. LIEB AND M. LOSS, *Analysis*, 2nd ed., Graduate Studies in Mathematics, 14, American Mathematical Society, Providence, RI, 2001.
- [15] J. K. MACKENZIE AND J. S. BOWLES, *The crystallography of martensite transformations II*, Acta Metall., 2 (1954), pp. 138–147.
- [16] G. W. MILTON, *The Theory of Composites*, Cambridge University Press, Cambridge, 2002.
- [17] R. L. PATTERSON AND C. M. WAYMAN, *Internal twinning in ferrous martensites*, Acta Met., 12 (1964), pp. 1306–11.
- [18] M. REED AND B. SIMON, *Methods of Modern Mathematical Physics*, 2nd ed., Academic Press, Inc., New York, 1980.
- [19] S. TIMOSHENKO AND J. N. GOODIER, *Theory of Elasticity*, 2nd ed., McGraw-Hill Book Company, Inc., New York, Toronto, London, 1951.
- [20] V. P. SMYSHLYAEV AND J. R. WILLIS, *A “non-local” variational approach to the elastic energy minimization of martensitic polycrystals*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci. 454 (1998), pp. 1573–1613.
- [21] M. S. WECHSLER, D. S. LIEBERMAN, AND T. A. READ, *On the theory of the formation of martensite*, AIME Trans. J. Metals, 197 (1953), pp. 1503–1515.
- [22] J. R. WILLIS, *Variational and related methods for the overall properties of composites*, Adv. Appl. Mech., 21 (1981), pp. 1–78.

THE PERIODIC UNFOLDING METHOD IN HOMOGENIZATION*

D. CIORANESCU[†], A. DAMLAMIAN[‡], AND G. GRISO[§]

Abstract. The periodic unfolding method was introduced in 2002 in [Cioranescu, Damlamian, and Griso, *C.R. Acad. Sci. Paris, Ser. 1*, 335 (2002), pp. 99–104] (with the basic proofs in [*Proceedings of the Narvik Conference 2004*, GAKUTO Internat. Ser. Math. Sci. Appl. 24, Gakkōtoshō, Tokyo, 2006, pp. 119–136]). In the present paper we go into all the details of the method and include complete proofs, as well as several new extensions and developments. This approach is based on two distinct ideas, each leading to a new ingredient. The first idea is the change of scale, which is embodied in the unfolding operator. At the expense of doubling the dimension, this allows one to use standard weak or strong convergence theorems in L^p spaces instead of more complicated tools (such as two-scale convergence, which is shown to be merely the weak convergence of the unfolding; cf. Remark 2.15). The second idea is the separation of scales, which is implemented as a macro-micro decomposition of functions and is especially suited for the weakly convergent sequences of Sobolev spaces. In the framework of this method, the proofs of most periodic homogenization results are elementary. The unfolding is particularly well-suited for multiscale problems (a simple backward iteration argument suffices) and for precise corrector results without extra regularity on the data. A list of the papers where these ideas appeared, at least in some preliminary form, is given with a discussion of their content. We also give a list of papers published since the publication [Cioranescu, Damlamian, and Griso, *C.R. Acad. Sci. Paris, Ser. 1*, 335 (2002), pp. 99–104], and where the unfolding method has been successfully applied.

Key words. homogenization, periodic unfolding, multiscale problems

AMS subject classifications. 49J45, 35B27, 74Q05

DOI. 10.1137/080713148

1. Introduction. The notion of two-scale convergence was introduced in 1989 by Nguetseng in [58], further developed by Allaire in [1] and by Lukkassen, Nguetseng, and Wall in [55] with applications to periodic homogenization. It was generalized to some multiscale problems by Ene and Saint Jean Paulin in [38], Allaire and Briane in [2], Lions et al. in [52] and Lukkassen, Nguetseng, and Wall in [55].

In 1990, Arbogast, Douglas, and Hornung defined a “dilation” operator in [5] to study homogenization for a periodic medium with double porosity. This technique was used again in [16], [3], [4], [48], [49], [50], [51], [54], [20], [21], [22], and [23].

In [24], we expanded on this idea and presented a general and quite simple approach for classical or multiscale periodic homogenization, under the name of “unfolding method.” Originally restricted to the case of domains consisting of a union of ε -cells, it was extended to general domains (see the survey of Damlamian [34]). In the present work, we give a complete presentation of this method, including all of the proofs, as well as several new extensions and developments. The relationship of the papers listed above with our work is discussed at the end of this introduction.

The periodic unfolding method is essentially based on two ingredients. The first one is the unfolding operator \mathcal{T}_ε (similar to the dilation operator), defined in section 2,

*Received by the editors January 1, 2008; accepted for publication (in revised form) May 5, 2008; published electronically November 26, 2008.

<http://www.siam.org/journals/sima/40-4/71314.html>

[†]Corresponding author. Laboratoire Jacques-Louis Lions, Université Pierre et Marie Curie, Boîte courrier 187, 4 Place Jussieu, 75252 Paris Cedex 05, France (cioran@ann.jussieu.fr).

[‡]Université Paris-Est, Laboratoire d'Analyse et de Mathématiques Appliquées, CNRS UMR 8050, Centre Multidisciplinaire de Créteil, 94010, Créteil, Cedex, France (damla@univ-paris12.fr).

[§]Corresponding author. Laboratoire Jacques-Louis Lions, Université Pierre et Marie Curie, Boîte courrier 187, 4 Place Jussieu, 75252 Paris Cedex 05, France (griso@ann.jussieu.fr).

where its properties are investigated. Let Ω be a bounded open set, and Y a reference cell in \mathbb{R}^n . By definition, the operator \mathcal{T}_ε associates to any function v in $L^p(\Omega)$, a function $\mathcal{T}_\varepsilon(v)$ in $L^p(\Omega \times Y)$. An immediate (and interesting) property of \mathcal{T}_ε is that it enables one to transform any integral over Ω in an integral over $\Omega \times Y$. Indeed, by Proposition 2.6 below

$$(1.1) \quad \int_{\Omega} w(x) \, dx \sim \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(w)(x, y) \, dx \, dy \quad \forall w \in L^1(\Omega).$$

Proposition 2.14 shows that the two-scale convergence in the $L^p(\Omega)$ -sense of a sequence of functions $\{v_\varepsilon\}$ is equivalent to the weak convergence of the sequence of unfolded functions $\{\mathcal{T}_\varepsilon(v_\varepsilon)\}$ in $L^p(\Omega \times Y)$. Thus, the two-scale convergence in Ω is reduced to a mere weak convergence in $L^p(\Omega \times Y)$, which conceptually simplifies proofs.

In section 2 are also introduced a local average operator \mathcal{M}_ε and an averaging operator \mathcal{U}_ε , the latter being, in some sense, the inverse of the unfolding operator \mathcal{T}_ε .

The second ingredient of the periodic unfolding method consists of separating the characteristic scales by decomposing every function φ belonging to $W^{1,p}(\Omega)$ in two parts. In section 3 it is achieved by using the local average. In section 4, the original proof of this scale-splitting, inspired by the finite element method (FEM), is given. The confrontation of the two methods of sections 3 and 4 is interesting in itself (Theorem 3.5 and Proposition 4.8). In both approaches, φ is written as $\varphi = \varphi_1^\varepsilon + \varepsilon\varphi_2^\varepsilon$, where φ_1^ε is a macroscopic part designed not to capture the oscillations of order ε (if there are any), while the microscopic part φ_2^ε is designed to do so. The main result states that, from any bounded sequence $\{w^\varepsilon\}$ in $W^{1,p}(\Omega)$, weakly convergent to some w , one can always extract a subsequence (still denoted $\{w^\varepsilon\}$) such that $w^\varepsilon = w_1^\varepsilon + \varepsilon w_2^\varepsilon$, with

$$(1.2) \quad \begin{aligned} \text{(i)} \quad & w_1^\varepsilon \rightharpoonup w \quad \text{weakly in } W^{1,p}(\Omega), \\ \text{(ii)} \quad & \mathcal{T}_\varepsilon(w^\varepsilon) \rightharpoonup w \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \text{(iii)} \quad & \mathcal{T}_\varepsilon(w_2^\varepsilon) \rightharpoonup \hat{w} \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \text{(iv)} \quad & \mathcal{T}_\varepsilon(\nabla w^\varepsilon) \rightharpoonup \nabla w + \nabla_y \hat{w} \quad \text{weakly in } L^p(\Omega \times Y), \end{aligned}$$

where \hat{w} belongs to $L^p(\Omega; W_{per}^{1,p}(Y))$.

In section 5 we apply the periodic unfolding method to a classical periodic homogenization problem. We point out that, in the framework of this method, the proof of the homogenization result is elementary. It relies essentially on formula (1.1), on the properties of \mathcal{T}_ε , and on convergences (1.2). It applies directly for both homogeneous Dirichlet or Neumann boundary conditions without hypothesis on the regularity of $\partial\Omega$. For nonhomogenous boundary conditions (or for Robin-type condition), some regularity of $\partial\Omega$ is required for the problem to make sense, in which case the method applies also directly (see Remark 5.12).

Section 6 is devoted to a corrector result, which holds without any additional regularity on the data (contrary to all previous proofs; see [11], [30], and [59]). This result follows from the use of the averaging operator \mathcal{U}_ε . The idea of using averages to improve corrector results first appeared in Dal Maso and Defranceschi [33]. We also give some error estimates and a new corrector result for the case of domains with

a smooth boundary (obtained by Griso in [42], [43], [44], and [45]). These results are explicitly connected to the unfolding method and improve on known classical ones (see [11] and [59]).

The periodic unfolding method is particularly well-suited for the case of multiscale problems. This is shown in section 7 by a simple backward iteration argument. This problem has a long history; one of the first papers on the subject is due to Bruggeman [19]. Its mathematical treatment by homogenization goes back to the book of Bensoussan, Lions, and Papanicolaou [11], where for this problem, the method of asymptotic expansions is used. For more recent references of multiscale homogenization and its applications, we refer to the books of Braides and Defranceschi [17], Milton [57], and the articles by Damlamian and Donato [35], Lukkassen and Milton [54], Lukkassen [53], Braides and Lukkassen [18], Babadjian and Baía [6], and Barchiesi [8].

The final section gives a list of papers where the method has been successfully applied since the publication of [24].

To conclude, let us turn back to the papers quoted at the beginning of this introduction and point out their relationships with our results. The dilation operation from Arbogast, Douglas, and Hornung [5] was defined in a domain which is an exact union of εY -cells. It consists in a change of variables, similar to that used in Definition 2.1 below. By this operation, any integral on Ω can be written as an integral over $\Omega \times Y$. Some elementary properties of the dilation operator in the space L^2 were also contained in Lemma 2 of [5].

The same dilation operator was used by Bourgeat, Luckhaus, and Mikelić in [16] under the name of “periodic modulation.” Proposition 4.6 of [16] showed that if a sequence two-scale converges and its periodic modulation converges weakly, they have the same limit.

In the context of two-scale convergence, Allaire and Conca [3] defined a pair of extension and projection operators (suited to Bloch decompositions) which are adjoint of each other. They are similar to our operators \mathcal{T}_ε and \mathcal{U}_ε and the equivalent of property (2.12) and Proposition 2.18(ii) below, are proved in Lemma 4.2 of [3]. These properties were exploited by Allaire, Conca, and Vanninathan in [4] for a general bounded domain by extending all functions by zero on its complement.

In [48], Lenczner used the dilation operator (here called “two-scale transformation”) in order to treat the homogenization of discrete electrical networks (by nature, the domain is a union of ε -cells). The convergence of the two-scale transform is called two-scale convergence (this would be confusing except that it was shown to be equivalent to the original two-scale convergence). As an aside, a convergence similar to (1.2)(iv) was also treated. In Lenczner and Mercier [49], Lenczner and Senouci-Bereksi [50], and Lenczner, Kader, and Perrier [51], this theory was applied to periodic electrical networks.

Finally, Casado Díaz and Luna-Laynez [21], Casado Díaz, Luna-Laynez, and Martín [22] and [23] used the dilation operator in the case of reticulated structures. In this framework, they obtained the equivalent of (3.7)(i) of Theorem 3.5 below.

2. Unfolding in L^p -spaces.

2.1. The unfolding operator \mathcal{T}_ε . In \mathbb{R}^n , let Ω be an open set and Y a reference cell (e.g., $]0, 1[^n$, or more generally, a set having the paving property, with respect to a basis (b_1, \dots, b_n) defining the periods).

By analogy with the notation in the one-dimensional case, for $z \in \mathbb{R}^n$, $[z]_Y$ denotes the unique integer combination $\sum_{j=1}^n k_j b_j$ of the periods such that $z - [z]_Y$ belongs

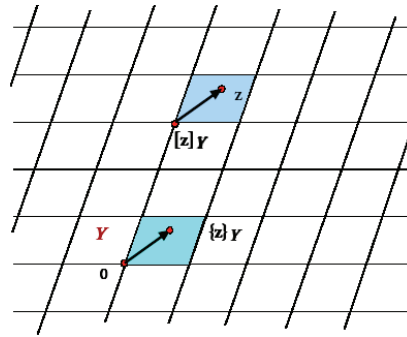


FIG. 1. Definition of $[z]_Y$ and $\{z\}_Y$.

to Y , and set

$$\{z\}_Y = z - [z]_Y \in Y \quad \text{a.e. for } z \in \mathbb{R}^n.$$

Then for each $x \in \mathbb{R}^n$, one has

$$x = \varepsilon \left(\left[\frac{x}{\varepsilon} \right]_Y + \left\{ \frac{x}{\varepsilon} \right\}_Y \right) \quad \text{a.e. for } x \in \mathbb{R}^n \text{ (See Figure 1).}$$

We use the following notations:

$$(2.1) \quad \begin{cases} \Xi_\varepsilon = \{ \xi \in \mathbb{Z}^N, \varepsilon(\xi + Y) \subset \Omega \}, \\ \widehat{\Omega}_\varepsilon = \text{interior} \left\{ \bigcup_{\xi \in \Xi_\varepsilon} \varepsilon(\xi + \overline{Y}) \right\}, \\ \Lambda_\varepsilon = \Omega \setminus \widehat{\Omega}_\varepsilon. \end{cases}$$

The set $\widehat{\Omega}_\varepsilon$ is the largest union of $\varepsilon(\xi + \overline{Y})$ cells ($\xi \in \mathbb{Z}^n$) included in Ω , while Λ_ε is the subset of Ω containing the parts from $\varepsilon(\xi + \overline{Y})$ cells intersecting the boundary $\partial\Omega$ (see Figure 2).

DEFINITION 2.1. For ϕ Lebesgue-measurable on Ω , the unfolding operator \mathcal{T}_ε is defined as follows:

$$\mathcal{T}_\varepsilon(\phi)(x, y) = \begin{cases} \phi \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon y \right) & \text{a.e. for } (x, y) \in \widehat{\Omega}_\varepsilon \times Y, \\ 0 & \text{a.e. for } (x, y) \in \Lambda_\varepsilon \times Y. \end{cases}$$

Observe that the function $\mathcal{T}_\varepsilon(\phi)$ is Lebesgue-measurable on $\Omega \times Y$ and vanishes for x outside of the set $\widehat{\Omega}_\varepsilon$.

As in classical periodic homogenization, two different scales appear in the definition of \mathcal{T}_ε : the “macroscopic” scale x gives the position of a point in the domain Ω , while the “microscopic” scale y ($= x/\varepsilon$) gives the position of a point in the cell Y . The unfolding operator doubles the dimension of the space and puts all of the oscillations in the second variable, in this way separating the two scales (see Figures 3, 4 and Figures 5, 6).

The following property of \mathcal{T}_ε is a simple consequence of Definition 2.1 for v and w Lebesgue-measurable; it will be used extensively:

$$(2.2) \quad \mathcal{T}_\varepsilon(vw) = \mathcal{T}_\varepsilon(v) \mathcal{T}_\varepsilon(w).$$

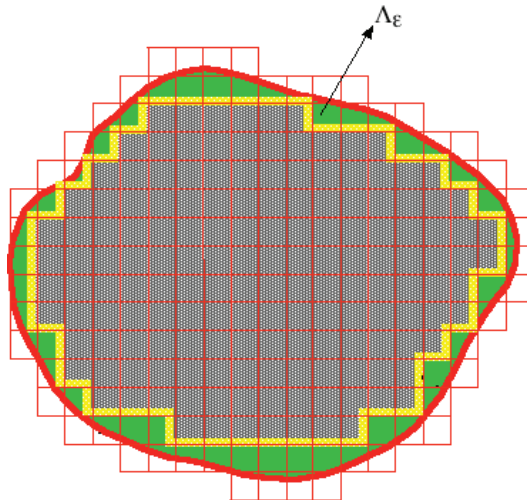


FIG. 2. The domains $\hat{\Omega}_\varepsilon$ and Λ_ε .

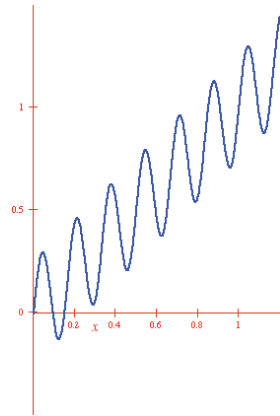


FIG. 3. $f_\varepsilon(x) = \frac{1}{4} \sin(2\pi \frac{x}{\varepsilon}) + x$; $\varepsilon = \frac{1}{6}$.

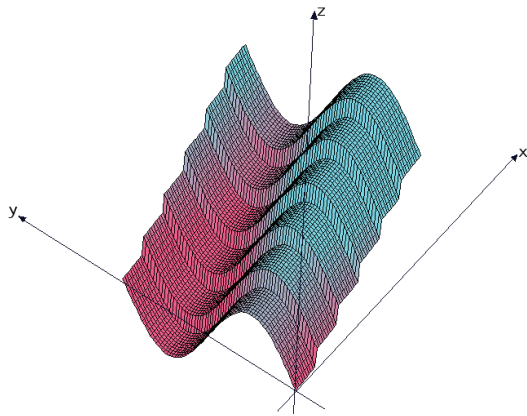


FIG. 4. $\mathcal{T}_\varepsilon(f_\varepsilon)$.

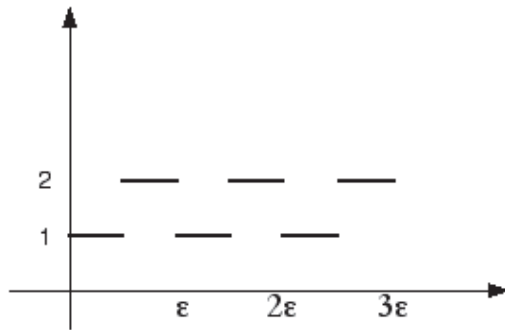


FIG. 5. $f_\varepsilon = f(\{\frac{x}{\varepsilon}\}_Y)$.

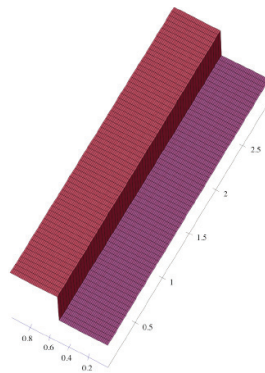


FIG. 6. $\mathcal{T}_\varepsilon(f_\varepsilon)$.

Another simple consequence of Definition 2.1 is the following result concerning highly oscillating functions.

PROPOSITION 2.2. *For f measurable on Y , extended by Y -periodicity to the whole of \mathbb{R}^n , define the sequence $\{f_\varepsilon\}$ by*

$$(2.3) \quad f_\varepsilon(x) = f\left(\frac{x}{\varepsilon}\right) \quad \text{a.e. for } x \in \mathbb{R}^n.$$

Then

$$\mathcal{T}_\varepsilon(f_\varepsilon|_\Omega)(x, y) = \begin{cases} f(y) & \text{a.e. for } (x, y) \in \widehat{\Omega}_\varepsilon \times Y, \\ 0 & \text{a.e. for } (x, y) \in \Lambda_\varepsilon \times Y. \end{cases}$$

If f belongs to $L^p(Y)$, $p \in [1, +\infty[$, and if Ω is bounded,

$$(2.4) \quad \mathcal{T}_\varepsilon(f_\varepsilon|_\Omega) \rightarrow f \quad \text{strongly in } L^p(\Omega \times Y).$$

Remark 2.3. An equivalent way to define f_ε in (2.3) is to take simply $f_\varepsilon(x) = f(\{\frac{x}{\varepsilon}\}_Y)$. For example, with

$$f(y) = \begin{cases} 1 & \text{for } y \in (0, 1/2), \\ 2 & \text{for } y \in (1/2, 1), \end{cases}$$

f_ε is the highly oscillating periodic function, with period ε from Figure 5.

Remark 2.4. Let f in $L^p(Y)$, $p \in [1, +\infty[$, and f_ε be defined by (2.3). It is well-known that $\{f_\varepsilon|_\Omega\}$ converges weakly in $L^p(\Omega)$ to the mean value of f on Y , and not strongly unless f is a constant (see Remark 2.11 below).

The next two results, essential in the study of the properties of the unfolding operator, are also straightforward from Definition 2.1.

PROPOSITION 2.5. *For $p \in [1, +\infty[$, the operator \mathcal{T}_ε is linear and continuous from $L^p(\Omega)$ to $L^p(\Omega \times Y)$. For every ϕ in $L^1(\Omega)$ and w in $L^p(\Omega)$,*

- (i) $\frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(\phi)(x, y) \, dx \, dy = \int_\Omega \phi(x) \, dx - \int_{\Lambda_\varepsilon} \phi(x) \, dx = \int_{\widehat{\Omega}_\varepsilon} \phi(x) \, dx,$
- (ii) $\frac{1}{|Y|} \int_{\Omega \times Y} |\mathcal{T}_\varepsilon(\phi)| \, dx \, dy \leq \int_\Omega |\phi| \, dx,$
- (iii) $\left| \int_\Omega \phi \, dx - \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(\phi) \, dx \, dy \right| \leq \int_{\Lambda_\varepsilon} |\phi| \, dx,$
- (iv) $\|\mathcal{T}_\varepsilon(w)\|_{L^p(\Omega \times Y)} = |Y|^{\frac{1}{p}} \|w 1_{\widehat{\Omega}_\varepsilon}\|_{L^p(\Omega)} \leq |Y|^{\frac{1}{p}} \|w\|_{L^p(\Omega)}.$

Proof. Recalling Definition 2.2 of $\widehat{\Omega}_\varepsilon$, one has

$$\begin{aligned} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(\phi)(x, y) \, dx \, dy &= \frac{1}{|Y|} \int_{\widehat{\Omega}_\varepsilon \times Y} \mathcal{T}_\varepsilon(\phi)(x, y) \, dx \, dy \\ &= \frac{1}{|Y|} \sum_{\xi \in \Xi_\varepsilon} \int_{(\varepsilon\xi + \varepsilon Y) \times Y} \mathcal{T}_\varepsilon(\phi)(x, y) \, dx \, dy. \end{aligned}$$

On each $(\varepsilon\xi + \varepsilon Y) \times Y$, by definition, $\mathcal{T}_\varepsilon(\phi)(x, y) = \phi(\varepsilon\xi + \varepsilon y)$ is constant in x . Hence, each integral in the sum on the right-hand side successively equals

$$\begin{aligned} \int_{(\varepsilon\xi + \varepsilon Y) \times Y} \mathcal{T}_\varepsilon(\phi)(x, y) \, dx \, dy &= |\varepsilon\xi + \varepsilon Y| \int_Y \phi(\varepsilon\xi + \varepsilon y) \, dy \\ &= \varepsilon^n |Y| \int_Y \phi(\varepsilon\xi + \varepsilon y) \, dy = |Y| \int_{(\varepsilon\xi + \varepsilon Y)} \phi(x) \, dx. \end{aligned}$$

By summing over Ξ_ε , the right-hand side becomes $\int_{\widehat{\Omega}_\varepsilon} \phi(x) \, dx$, which gives the result. \square

Property (iii) in Proposition 2.5 shows that any integral of a function on Ω is “almost equivalent” to the integral of its unfolded on $\Omega \times Y$; the “integration defect” arises only from the cells intersecting the boundary $\partial\Omega$ and is controlled by its integral over Λ_ε .

The next proposition, which we call **unfolding criterion for integrals** (u.c.i.), is a very useful tool when treating homogenization problems.

PROPOSITION 2.6 (u.c.i.). *If $\{\phi_\varepsilon\}$ is a sequence in $L^1(\Omega)$ satisfying*

$$\int_{\Lambda_\varepsilon} |\phi_\varepsilon| \, dx \rightarrow 0,$$

then

$$\int_\Omega \phi_\varepsilon \, dx - \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(\phi_\varepsilon) \, dx \, dy \rightarrow 0.$$

Based on this result, we introduce the following notation.

Notation. If $\{w_\varepsilon\}$ is a sequence satisfying u.c.i., we write

$$\int_{\Omega} w_\varepsilon dx \stackrel{\mathcal{T}_\varepsilon}{\simeq} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(w_\varepsilon) dx dy.$$

PROPOSITION 2.7. *Let $\{u_\varepsilon\}$ be a bounded sequence in $L^p(\Omega)$, with $p \in]1, +\infty[$ and $v \in L^{p'}(\Omega)$ ($1/p + 1/p' = 1$), then*

$$(2.5) \quad \int_{\Omega} u_\varepsilon v dx \stackrel{\mathcal{T}_\varepsilon}{\simeq} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(u_\varepsilon) \mathcal{T}_\varepsilon(v) dx dy.$$

Suppose $\partial\Omega$ is bounded. Let $\{u_\varepsilon\}$ be a bounded sequence in $L^p(\Omega)$ and $\{v_\varepsilon\}$ a bounded sequence in $L^q(\Omega)$, with $1/p + 1/q < 1$, then

$$(2.6) \quad \int_{\Omega} u_\varepsilon v_\varepsilon dx \stackrel{\mathcal{T}_\varepsilon}{\simeq} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(u_\varepsilon) \mathcal{T}_\varepsilon(v_\varepsilon) dx dy.$$

Proof. Observe that $1_{\Lambda_\varepsilon}(x) \rightarrow 0$ for all $x \in \Omega$. Consequently, by the Lebesgue dominated convergence theorem, one gets $\int_{\Lambda_\varepsilon} |v|^{p'} dx \rightarrow 0$, and then by the Hölder inequality, $\int_{\Lambda_\varepsilon} |u_\varepsilon v| dx \rightarrow 0$. This proves (2.5). If $\partial\Omega$ is bounded, then one immediately has $|\Lambda_\varepsilon| \rightarrow 0$ when $\varepsilon \rightarrow 0$, and this implies (2.6). \square

We now investigate the convergence properties related to the unfolding operator when $\varepsilon \rightarrow 0$. For ϕ uniformly continuous on Ω , with modulus of continuity m_ϕ , it is easy to see that

$$\sup_{x \in \widehat{\Omega}_\varepsilon, y \in Y} |\mathcal{T}_\varepsilon(\phi)(x, y) - \phi(x)| \leq m_\phi(\varepsilon).$$

So, as ε goes to zero, even though $\mathcal{T}_\varepsilon(\phi)$ is not continuous, it converges to ϕ uniformly on any open set strongly included in Ω . By density, and making use of Proposition 2.5, further convergence properties can be expressed using the mean value of a function defined on $\Omega \times Y$.

DEFINITION 2.8. *The mean value operator $\mathcal{M}_Y : L^p(\Omega \times Y) \mapsto L^p(\Omega)$ for $p \in [1, +\infty]$, is defined as follows:*

$$(2.7) \quad \mathcal{M}_Y(\Phi)(x) = \frac{1}{|Y|} \int_Y \Phi(x, y) dy \quad \text{a.e. for } x \in \Omega.$$

Observe that an immediate consequence of this definition is the estimate

$$\|\mathcal{M}_Y(\Phi)\|_{L^p(\Omega)} \leq |Y|^{-\frac{1}{p}} \|\Phi\|_{L^p(\Omega \times Y)} \quad \text{for every } \Phi \in L^p(\Omega \times Y).$$

PROPOSITION 2.9. *Let p belong to $[1, +\infty[$.*

(i) *For $w \in L^p(\Omega)$,*

$$\mathcal{T}_\varepsilon(w) \rightarrow w \quad \text{strongly in } L^p(\Omega \times Y).$$

(ii) *Let $\{w_\varepsilon\}$ be a sequence in $L^p(\Omega)$ such that*

$$w_\varepsilon \rightarrow w \quad \text{strongly in } L^p(\Omega).$$

Then

$$\mathcal{T}_\varepsilon(w_\varepsilon) \rightarrow w \quad \text{strongly in } L^p(\Omega \times Y).$$

(iii) For every relatively weakly compact sequence $\{w_\varepsilon\}$ in $L^p(\Omega)$, the corresponding $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$ is relatively weakly compact in $L^p(\Omega \times Y)$. Furthermore, if

$$\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup \widehat{w} \quad \text{weakly in } L^p(\Omega \times Y),$$

then

$$w_\varepsilon \rightharpoonup \mathcal{M}_Y(\widehat{w}) \quad \text{weakly in } L^p(\Omega).$$

(iv) If $\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup \widehat{w}$ weakly in $L^p(\Omega \times Y)$, then

$$(2.8) \quad \|\widehat{w}\|_{L^p(\Omega \times Y)} \leq \liminf_{\varepsilon \rightarrow 0} |Y|^{\frac{1}{p}} \|w_\varepsilon\|_{L^p(\Omega)}.$$

(v) Suppose $p > 1$, and let $\{w_\varepsilon\}$ be a bounded sequence in $L^p(\Omega)$. Then, the following assertions are equivalent:

- (a) $\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup \widehat{w}$ weakly in $L^p(\Omega \times Y)$ and $\limsup_{\varepsilon \rightarrow 0} |Y|^{\frac{1}{p}} \|w_\varepsilon\|_{L^p(\Omega)} \leq \|\widehat{w}\|_{L^p(\Omega \times Y)}$,
- (b) $\mathcal{T}_\varepsilon(w_\varepsilon) \rightarrow \widehat{w}$ strongly in $L^p(\Omega \times Y)$ and $\int_{\Lambda_\varepsilon} |w_\varepsilon|^p dx \rightarrow 0$.

Proof. (i) The result is obvious for any $w \in \mathcal{D}(\Omega)$. If $w \in L^p(\Omega)$, let $\phi \in \mathcal{D}(\Omega)$. Then, by using (iv) from Proposition 2.5,

$$\begin{aligned} \|\mathcal{T}_\varepsilon(w) - w\|_{L^p(\Omega \times Y)} &= \|\mathcal{T}_\varepsilon(w - \phi) + (\mathcal{T}_\varepsilon(\phi) - \phi) + (\phi - w)\|_{L^p(\Omega \times Y)} \\ &\leq 2|Y|^{\frac{1}{p}} \|w - \phi\|_{L^p(\Omega)} + \|\mathcal{T}_\varepsilon(\phi) - \phi\|_{L^p(\Omega \times Y)}, \end{aligned}$$

hence,

$$\limsup_{\varepsilon \rightarrow 0} \|\mathcal{T}_\varepsilon(w) - w\|_{L^p(\Omega \times Y)} \leq 2|Y|^{\frac{1}{p}} \|w - \phi\|_{L^p(\Omega)},$$

from which statement (i) follows by density.

(ii) The following estimate, a consequence of Proposition 2.5(iv), gives the result

$$\|\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{T}_\varepsilon(w)\|_{L^p(\Omega \times Y)} \leq |Y|^{\frac{1}{p}} \|w_\varepsilon - w\|_{L^p(\Omega)} \quad \forall w \in L^p(\Omega).$$

(iii) For $p \in]1, +\infty[$, by Proposition 2.5(iv), boundedness is preserved by \mathcal{T}_ε . Suppose that $\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup \widehat{w}$ weakly in $L^p(\Omega \times Y)$, and let $\psi \in L^{p'}(\Omega)$. From Proposition 2.7,

$$\int_\Omega w_\varepsilon(x) \psi(x) dx \stackrel{\mathcal{T}_\varepsilon}{\simeq} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(w_\varepsilon)(x, y) \mathcal{T}_\varepsilon(\psi)(x, y) dx dy.$$

In view of (i), one can pass to the limit in the right-hand side to obtain

$$\lim_{\varepsilon \rightarrow 0} \int_\Omega w_\varepsilon(x) \psi(x) dx = \int_\Omega \left\{ \frac{1}{|Y|} \int_Y \widehat{w}(x, y) dy \right\} \psi(x) dx.$$

For $p = 1$, one uses the extra property satisfied by weakly convergent sequences in $L^1(\Omega)$, in the form of the De La Vallée–Poussin criterion (which is equivalent to

relative weak compactness): there exists a continuous convex function $\Phi : \mathbb{R}^+ \mapsto \mathbb{R}^+$ such that

$$\lim_{t \rightarrow +\infty} \frac{\Phi(t)}{t} = +\infty, \quad \text{and the set } \left\{ \int_{\Omega} (\Phi \circ |w_{\varepsilon}|)(x) dx \right\} \text{ is bounded.}$$

Unfolding the last integral shows that

$$\left\{ \int_{\Omega \times Y} (\Phi \circ |\mathcal{T}_{\varepsilon}(w_{\varepsilon})|)(x, y) dx dy \right\} \text{ is bounded,}$$

which completes the proof of weak compactness of $\{\mathcal{T}_{\varepsilon}(w_{\varepsilon})\}$ in $L^1(\Omega \times Y)$ in the case of Ω with finite measure. For the case where the measure of Ω is not finite, a similar argument shows that the equiintegrability at infinity of the sequence $\{w_{\varepsilon}\}$ carries over to $\{\mathcal{T}_{\varepsilon}(w_{\varepsilon})\}$.

If $\mathcal{T}_{\varepsilon}(w_{\varepsilon}) \rightharpoonup \widehat{w}$ weakly in $L^1(\Omega \times Y)$, let ψ be in $\mathcal{D}(\Omega)$. For ε sufficiently small, one has

$$\int_{\Omega} w_{\varepsilon}(x) \psi(x) dx = \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_{\varepsilon}(w_{\varepsilon})(x, y) \mathcal{T}_{\varepsilon}(\psi)(x, y) dx dy.$$

In view of (i), one can pass to the limit in the right-hand side to obtain

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} w_{\varepsilon}(x) \psi(x) dx = \int_{\Omega} \left\{ \frac{1}{|Y|} \int_Y \widehat{w}(x, y) dy \right\} \psi(x) dx.$$

(iv) Inequality (2.8) is a simple consequence of Proposition 2.5(ii).

(v) Proposition 2.5(i) applied to the function $|w_{\varepsilon}|^p$ gives

$$\frac{1}{|Y|} \|\mathcal{T}_{\varepsilon}(w_{\varepsilon})\|_{L^p(\Omega \times Y)}^p + \int_{\Lambda_{\varepsilon}} |w_{\varepsilon}|^p dx = \|w_{\varepsilon}\|_{L^p(\Omega)}^p.$$

This identity implies the required equivalence. \square

COROLLARY 2.10. *Let f be in $L^p(Y)$, $p \in [1, +\infty[$, and $\{f_{\varepsilon}\}$ be the sequence defined by (2.3). Then*

$$(2.9) \quad f_{\varepsilon}|_{\Omega} \rightharpoonup \mathcal{M}_Y(f) \quad \text{weakly in } L^p(\Omega).$$

Proof. Proposition 2.2 gives the strong (hence weak) convergence of $\{\mathcal{T}_{\varepsilon}(f_{\varepsilon}|_{\Omega})\}$ to f in $L^p(\Omega \times Y)$. Convergence (2.9) follows from Proposition 2.9(iii).¹ \square

Remark 2.11. In general, in the case where Λ_{ε} is not null set (for every ε), the strong (resp. weak) convergence of the sequence $\{\mathcal{T}_{\varepsilon}(w_{\varepsilon})\}$ does not imply the corresponding convergence for the sequence $\{w_{\varepsilon}\}$, since it gives no control of the sequence $\{w_{\varepsilon}1_{\Lambda_{\varepsilon}}\}$. If $\{w_{\varepsilon}1_{\Lambda_{\varepsilon}}\}$ is bounded in $L^p(\Omega)$ and if $\{\mathcal{T}_{\varepsilon}(w_{\varepsilon})\}$ converges weakly, so does $\{w_{\varepsilon}\}$ by Proposition 2.9(iii). On the other hand, even if $\{w_{\varepsilon}1_{\Lambda_{\varepsilon}}\}$ converges strongly to 0 in $L^p(\Omega)$, the strong convergence of $\{\mathcal{T}_{\varepsilon}(w_{\varepsilon})\}$ does not imply that of $\{w_{\varepsilon}\}$ as it is shown by the sequence $\{f_{\varepsilon}|_{\Omega}\}$ in Corollary 2.10, unless f is a constant on Y .

COROLLARY 2.12. *Let p belong to $]1, +\infty[$, let $\{u_{\varepsilon}\}$ be a sequence in $L^p(\Omega)$ such that*

$$\mathcal{T}_{\varepsilon}(u_{\varepsilon}) \rightharpoonup u \quad \text{weakly in } L^p(\Omega \times Y),$$

¹Note that the proof of convergence (2.9) is really straightforward when using unfolding!

and let $\{v_\varepsilon\}$ be a sequence in $L^{p'}(\Omega)$ ($1/p + 1/p' = 1$), with

$$\mathcal{T}_\varepsilon(v_\varepsilon) \rightarrow v \quad \text{strongly in } L^{p'}(\Omega \times Y).$$

Then, for any φ in $C_c(\Omega)$, one has

$$\int_\Omega u_\varepsilon(x) v_\varepsilon(x) \varphi(x) dx \rightarrow \frac{1}{|Y|} \int_{\Omega \times Y} u(x, y) v(x, y) \varphi(x) dx dy.$$

Moreover, if

$$\int_{\Lambda_\varepsilon} |v_\varepsilon|^{p'} dx \rightarrow 0,$$

then, for any φ in $C(\bar{\Omega})$, one has

$$\int_\Omega u_\varepsilon(x) v_\varepsilon(x) \varphi(x) dx \rightarrow \frac{1}{|Y|} \int_{\Omega \times Y} u(x, y) v(x, y) \varphi(x) dx dy.$$

Proof. The result follows from the fact that, in both cases, the sequence $\{u_\varepsilon v_\varepsilon \phi\}$ satisfies the u.c.i. by the Hölder inequality. \square

Remark 2.13. A consequence of (iii) of Proposition 2.9, together with (iv) of Proposition 2.5, is the following. Suppose the sequence $\{w_\varepsilon\}$ converges weakly to w in $L^p(\Omega)$. Then the sequence $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$ is relatively weakly compact in $L^p(\Omega \times Y)$, and each of its weak-limit points \hat{w} satisfies $\mathcal{M}_Y(\hat{w}) = w$.

Now recall the following definition from Nguetseng [58] and Allaire [1].

Two-scale convergence. Let $p \in]1, +\infty[$. A bounded sequence $\{w_\varepsilon\}$ in $L^p(\Omega)$ two-scale converges to some w belonging to $L^p(\Omega \times Y)$, whenever, for every smooth function φ on $\Omega \times Y$, the following convergence holds:

$$\int_\Omega w_\varepsilon(x) \varphi\left(x, \frac{x}{\varepsilon}\right) dx \rightarrow \frac{1}{|Y|} \int \int_{\Omega \times Y} w(x, y) \varphi(x, y) dx dy.$$

The next result reduces two-scale convergence of a sequence to a mere weak $L^p(\Omega \times Y)$ -convergence of the unfolded sequence.

PROPOSITION 2.14. Let $\{w_\varepsilon\}$ be a bounded sequence in $L^p(\Omega)$, with $p \in]1, +\infty[$. The following assertions are equivalent:

- (i) $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$ converges weakly to w in $L^p(\Omega \times Y)$,
- (ii) $\{w_\varepsilon\}$ two-scale converges to w .

Proof. To prove this equivalence, it is enough to check that, for every φ in a set of admissible test functions for two-scale convergence (for instance, $\mathcal{D}(\Omega, L^q(Y))$), the sequence $\{\mathcal{T}_\varepsilon[\varphi(x, x/\varepsilon)]\}$ converges strongly to φ in $L^q(\Omega \times Y)$. This follows from the definition of \mathcal{T}_ε ; indeed

$$\mathcal{T}_\varepsilon\left[\varphi\left(x, \frac{x}{\varepsilon}\right)\right](x, y) = \varphi\left(\varepsilon \left[\frac{x}{\varepsilon}\right]_Y + \varepsilon y, y\right). \quad \square$$

Remark 2.15. Proposition 2.14 shows that the two-scale convergence of a sequence in $L^p(\Omega)$, $p \in]1, +\infty[$, is equivalent to the weak- $L^p(\Omega \times Y)$ convergence of the unfolded sequence. Notice that, by definition, to check the two-scale convergence, one has to use special test functions. To check a weak convergence in the space $L^p(\Omega \times Y)$, one simply makes the use of functions in the dual space $L^{p'}(\Omega \times Y)$. Moreover, due to density properties, it is sufficient to check this convergence only on smooth functions from $\mathcal{D}(\Omega \times Y)$.

2.2. The averaging operator \mathcal{U}_ε . In this section, we consider the adjoint \mathcal{U}_ε of \mathcal{T}_ε , which we call averaging operator. In order to do so, let v be in $L^p(\Omega \times Y)$, and let u be in $L^{p'}(\Omega)$. We have successively,

$$\begin{aligned} \frac{1}{|Y|} \int_{\Omega \times Y} \mathcal{T}_\varepsilon(u)(x, y) v(x, y) dx dy &= \frac{1}{|Y|} \int_{\widehat{\Omega}_\varepsilon \times Y} \mathcal{T}_\varepsilon(u)(x, y) v(x, y) dx dy \\ &= \frac{1}{|Y|} \sum_{\xi \in \Xi_\varepsilon} \int_{\varepsilon(\xi+Y) \times Y} u(\varepsilon\xi + \varepsilon y) v(x, y) dx dy \\ &= \sum_{\xi \in \Xi_\varepsilon} \frac{1}{|Y|} \int_{Y \times Y} u(\varepsilon\xi + \varepsilon y) v(\varepsilon\xi + \varepsilon z, y) \varepsilon^N dz dy \\ &= \sum_{\xi \in \Xi_\varepsilon} \frac{1}{|Y|} \int_Y dz \int_{\varepsilon(\xi+Y)} u(x) v\left(\varepsilon \left[\frac{x}{\varepsilon}\right]_Y + \varepsilon z, \left\{\frac{x}{\varepsilon}\right\}_Y\right) dx \\ &= \int_{\widehat{\Omega}_\varepsilon} u(x) \left(\frac{1}{|Y|} \int_Y v\left(\varepsilon \left[\frac{x}{\varepsilon}\right]_Y + \varepsilon z, \left\{\frac{x}{\varepsilon}\right\}_Y\right) dz \right) dx. \end{aligned}$$

This gives the formula for the averaging operator \mathcal{U}_ε .

DEFINITION 2.16. For p in $[1, +\infty]$, the averaging operator $\mathcal{U}_\varepsilon : L^p(\Omega \times Y) \mapsto L^p(\Omega)$ is defined as

$$\mathcal{U}_\varepsilon(\Phi)(x) = \begin{cases} \frac{1}{|Y|} \int_Y \Phi\left(\varepsilon \left[\frac{x}{\varepsilon}\right]_Y + \varepsilon z, \left\{\frac{x}{\varepsilon}\right\}_Y\right) dz & \text{a.e. for } x \in \widehat{\Omega}_\varepsilon, \\ 0 & \text{a.e. for } x \in \Lambda_\varepsilon. \end{cases}$$

Consequently, for $\psi \in L^p(\Omega)$ and $\Phi \in L^{p'}(\Omega \times Y)$, one has

$$\int_\Omega \mathcal{U}_\varepsilon(\Phi)(x) \psi(x) dx = \frac{1}{|Y|} \int_{\Omega \times Y} \Phi(x, y) \mathcal{T}_\varepsilon(\psi)(x, y) dx dy.$$

As a consequence of the duality (Hölder’s inequality) and of Proposition 2.5(iv), we get the following.

PROPOSITION 2.17. Let p belong to $[1, +\infty]$. The averaging operator is linear and continuous from $L^p(\Omega \times Y)$ to $L^p(\Omega)$ and

$$(2.10) \quad \|\mathcal{U}_\varepsilon(\Phi)\|_{L^p(\Omega)} \leq |Y|^{-\frac{1}{p}} \|\Phi\|_{L^p(\Omega \times Y)}.$$

The operator \mathcal{U}_ε maps $L^p(\Omega \times Y)$ into the space $L^p(\Omega)$. It allows one to replace the function $x \mapsto \Phi(x, \left\{\frac{x}{\varepsilon}\right\}_Y)$, which is meaningless, in general, by a function which always makes sense. Notice that this implies, in particular, that the largest set of test functions for two-scale convergence is actually the set $\mathcal{U}_\varepsilon(\Phi)$, with Φ in $L^{p'}(\Omega \times Y)$.

It is immediate from its definition that \mathcal{U}_ε is almost a left-inverse of \mathcal{T}_ε , since

$$(2.11) \quad \mathcal{U}_\varepsilon(\mathcal{T}_\varepsilon(\phi))(x) = \begin{cases} \phi(x) & \text{a.e. for } x \in \widehat{\Omega}_\varepsilon, \\ 0 & \text{a.e. for } x \in \Lambda_\varepsilon, \end{cases}$$

for every ϕ in $L^p(\Omega)$, while

$$(2.12) \quad \mathcal{T}_\varepsilon(\mathcal{U}_\varepsilon(\Phi))(x, y) = \begin{cases} \frac{1}{|Y|} \int_Y \Phi\left(\varepsilon \left[\frac{x}{\varepsilon}\right]_Y + \varepsilon z, y\right) dz & \text{a.e. for } (x, y) \in \widehat{\Omega}_\varepsilon \times Y, \\ 0 & \text{a.e. for } (x, y) \in \Lambda_\varepsilon \times Y, \end{cases}$$

for every Φ in $L^p(\Omega \times Y)$.

PROPOSITION 2.18 (properties of \mathcal{U}_ε). *Suppose that p is in $[1, +\infty[$.*

(i) *Let $\{\Phi_\varepsilon\}$ be a bounded sequence in $L^p(\Omega \times Y)$ such that $\Phi_\varepsilon \rightharpoonup \Phi$ weakly in $L^p(\Omega \times Y)$. Then*

$$\mathcal{U}_\varepsilon(\Phi_\varepsilon) \rightharpoonup \mathcal{M}_Y(\Phi) = \frac{1}{|Y|} \int_Y \Phi(\cdot, y) dy \quad \text{weakly in } L^p(\Omega).$$

In particular, for every $\Phi \in L^p(\Omega \times Y)$,

$$\mathcal{U}_\varepsilon(\Phi) \rightharpoonup \mathcal{M}_Y(\Phi) \quad \text{weakly in } L^p(\Omega),$$

but not strongly, unless Φ is independent of y .

(ii) *Let $\{\Phi_\varepsilon\}$ be a sequence such that $\Phi_\varepsilon \rightarrow \Phi$ strongly in $L^p(\Omega \times Y)$. Then*

$$\mathcal{T}_\varepsilon(\mathcal{U}_\varepsilon(\Phi_\varepsilon)) \rightarrow \Phi \quad \text{strongly in } L^p(\Omega \times Y).$$

(iii) *Suppose that $\{w_\varepsilon\}$ is a sequence in $L^p(\Omega)$. Then, the following assertions are equivalent:*

- (a) $\mathcal{T}_\varepsilon(w_\varepsilon) \rightarrow \widehat{w}$ strongly in $L^p(\Omega \times Y)$,
- (b) $w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{w}) \rightarrow 0$ strongly in $L^p(\Omega)$.

(iv) *Suppose that $\{w_\varepsilon\}$ is a sequence in $L^p(\Omega)$. Then, the following assertions are equivalent:*

- (c) $\mathcal{T}_\varepsilon(w_\varepsilon) \rightarrow \widehat{w}$ strongly in $L^p(\Omega \times Y)$ and $\int_{\Lambda_\varepsilon} |w_\varepsilon|^p dx \rightarrow 0$,
- (d) $w_\varepsilon - \mathcal{U}_\varepsilon(\widehat{w}) \rightarrow 0$ strongly in $L^p(\Omega)$.

Proof. (i) This follows from Proposition 2.9(ii) by duality for $p > 1$. It still holds for $p = 1$ in the same way as the proof of Proposition 2.9(ii). Indeed, if the De La Vallée–Poussin criterion is satisfied by the sequence $\{\Phi_\varepsilon\}$, it is also satisfied by the sequence $\{\mathcal{U}_\varepsilon(\Phi_\varepsilon)\}$, since for F convex and continuous, Jensen’s inequality implies that

$$F(\mathcal{U}_\varepsilon(\Phi_\varepsilon))(x) \leq \mathcal{U}_\varepsilon(F(\Phi_\varepsilon))(x).$$

(ii) The proof follows the same lines as that of (i)–(ii) of Proposition 2.9.

(iii) The implication (a) \Rightarrow (b) follows from (2.10) applied to $\Phi_\varepsilon = w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{w})$ and from (2.11).

As for the converse (b) \Rightarrow (a), Proposition 2.9(ii) implies that

$$\mathcal{T}_\varepsilon(w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{w})) \rightarrow 0 \quad \text{strongly in } L^p(\Omega \times Y).$$

Since $\mathcal{T}_\varepsilon(w_\varepsilon) = \mathcal{T}_\varepsilon(w_\varepsilon 1_{\widehat{\Omega}_\varepsilon})$, from (ii) above it converges to \widehat{w} strongly in $L^p(\Omega \times Y)$.

(iv) The implication (c) \Rightarrow (d) follows from (iii) and the second condition of (c).

Its converse (d) \Rightarrow (c) is a consequence of from (iii): since $\mathcal{U}_\varepsilon(\widehat{w}) 1_{\Lambda_\varepsilon} = 0$, (d) implies (b) and $w_\varepsilon 1_{\Lambda_\varepsilon} \rightarrow 0$ in $L^p(\Omega)$. \square

Remark 2.19. The statement of Proposition 2.18(iii) does not hold with weak convergences instead of strong ones, contrary to an erroneous statement made in [24]. In view of (2.11) and Proposition 2.18(i), if $\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup \widehat{w}$ weakly in $L^p(\Omega \times Y)$, then $w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{w})$ converges weakly to 0 in $L^p(\Omega)$.

But the converse of this last implication cannot hold. Indeed, choose \widehat{v} with $\mathcal{M}_Y(\widehat{v}) = 0$. By Proposition 2.18(i), $\mathcal{U}_\varepsilon(\widehat{v})$ converges weakly to $\mathcal{M}_Y(\widehat{v}) = 0$. Consequently, the weak limit of $w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{v})$ is also the weak limit of $w_\varepsilon 1_{\widehat{\Omega}_\varepsilon} - \mathcal{U}_\varepsilon(\widehat{w} + \widehat{v})$. If the converse were true, it would imply that $\mathcal{T}_\varepsilon(w_\varepsilon)$ converges weakly to both \widehat{w} and $\widehat{w} + \widehat{v}$. So $\widehat{v} = 0$. In other words, $\mathcal{M}_Y(\widehat{v}) = 0$ would imply $\widehat{v} = 0$.

Remark 2.20. Assertions (iii)(b) and (iv)(d) are corrector-type results.

Remark 2.21. The condition (iii)(a) is used by some authors to define the notion of “strong two-scale convergence.” From the above considerations, condition (c) of Proposition 2.18(iv) is a better candidate for this definition.

2.3. The local average operator \mathcal{M}_ε . In this section, we consider the classical average operator associated to the partition of Ω by ε -cells Y (setting it to be zero on the cells intersecting the boundary $\partial\Omega$).

DEFINITION 2.22. *The local average operator $\mathcal{M}_\varepsilon : L^p(\Omega) \mapsto L^p(\Omega)$, for $p \in [1, +\infty]$, is defined by*

$$(2.13) \quad \mathcal{M}_\varepsilon(\phi)(x) = \begin{cases} \frac{1}{\varepsilon^N |Y|} \int_\varepsilon \left[\frac{x}{\varepsilon} \right]_y \phi(\zeta) d\zeta & \text{if } x \in \widehat{\Omega}_\varepsilon, \\ 0 & \text{if } x \in \Lambda_\varepsilon. \end{cases}$$

Remark 2.23. It turns out that the local average \mathcal{M}_ε is connected to the unfolding operator \mathcal{T}_ε . Indeed, by the usual change of variable cell by cell,

$$\mathcal{M}_\varepsilon(\phi)(x) = \frac{1}{|Y|} \int_Y \mathcal{T}_\varepsilon(\phi)(x, y) dy = \mathcal{M}_Y(\mathcal{T}_\varepsilon(\phi))(x).$$

Remark 2.24. Note that, for any ϕ in $L^p(\Omega)$, one has $\mathcal{T}_\varepsilon(\mathcal{M}_\varepsilon(\phi)) = \mathcal{M}_\varepsilon(\phi)$ on the set $\Omega \times Y$. Moreover, one also has $\mathcal{U}_\varepsilon(\phi) = \mathcal{M}_\varepsilon(\phi)$.

PROPOSITION 2.25 (properties of \mathcal{M}_ε).

(i) *Suppose that p is in $[1, +\infty]$. For any ϕ in $L^p(\Omega)$,*

$$\|\mathcal{M}_\varepsilon(\phi)\|_{L^p(\Omega)} \leq \|\phi\|_{L^p(\Omega)}.$$

(ii) *Suppose that p is in $[1, +\infty]$. For $\phi \in L^p(\Omega)$ and $\psi \in L^{p'}(\Omega)$,*

$$(2.14) \quad \int_\Omega \mathcal{M}_\varepsilon(\phi) \psi dx = \int_\Omega \mathcal{M}_\varepsilon(\phi) \mathcal{M}_\varepsilon(\psi) dx = \int_\Omega \phi \mathcal{M}_\varepsilon(\psi) dx.$$

(iii) *Suppose that p is in $[1, +\infty[$. Let $\{v_\varepsilon\}$ be a sequence such that $v_\varepsilon \rightarrow v$ strongly in $L^p(\Omega)$. Then*

$$\mathcal{M}_\varepsilon(v_\varepsilon) \rightarrow v \quad \text{strongly in } L^p(\Omega).$$

In particular, for every $\phi \in L^p(\Omega)$,

$$(2.15) \quad \mathcal{M}_\varepsilon(\phi) \rightarrow \phi \quad \text{strongly in } L^p(\Omega).$$

(iv) *Suppose that p is in $[1, +\infty[$. Let $\{v_\varepsilon\}$ be a sequence such that $v_\varepsilon \rightharpoonup v$ weakly in $L^p(\Omega)$. Then*

$$\mathcal{M}_\varepsilon(v_\varepsilon) \rightharpoonup v \quad \text{weakly in } L^p(\Omega).$$

The same holds true for the weak- topology in $L^\infty(\Omega)$.*

Proof. The proofs of (i) and (ii) are straightforward. The proof of (iii) is a simple consequence of (ii) of Proposition 2.9. For the proof of (iv), let ϕ be in $L^{p'}(\Omega)$, with $p' \in [1, +\infty[$ ($p \neq 1$), and use (2.14) and (2.15) to obtain

$$\int_{\Omega} \phi \mathcal{M}_{\varepsilon}(v_{\varepsilon}) \, dx = \int_{\Omega} \mathcal{M}_{\varepsilon}(\phi) v_{\varepsilon} \, dx \rightarrow \int_{\Omega} \phi v \, dx.$$

For $p = 1$, in the same way as the proof of Proposition 2.9(ii) and Proposition 2.18(i), if the De La Vallée–Poussin criterion is satisfied by the sequence $\{v_{\varepsilon}\}$, it is also satisfied by the sequence $\{\mathcal{M}_{\varepsilon}(v_{\varepsilon})\}$, since for F convex and continuous, Jensen’s inequality implies that

$$F(\mathcal{M}_{\varepsilon}(v_{\varepsilon}))(x) \leq \mathcal{M}_{\varepsilon}(F(v_{\varepsilon}))(x),$$

which ends the proof. \square

COROLLARY 2.26. *Suppose that p is in $[1, +\infty[$. Let w be in $L^p(\Omega)$ and $\{w_{\varepsilon}\}$ be a sequence in $L^p(\Omega)$ satisfying $\mathcal{T}_{\varepsilon}(w_{\varepsilon}) \rightarrow w$ strongly in $L^p(\Omega \times Y)$. Then,*

$$w_{\varepsilon} 1_{\widehat{\Omega}_{\varepsilon}} \rightarrow w \quad \text{strongly in } L^p(\Omega).$$

Furthermore, if $\int_{\Lambda_{\varepsilon}} |w_{\varepsilon}|^p \rightarrow 0$, then, $w_{\varepsilon} \rightarrow w$ strongly in $L^p(\Omega)$.

Proof. Since w does not depend on y , one has $\mathcal{U}_{\varepsilon}(w) = \mathcal{M}_{\varepsilon}(w)$ which, by Proposition 2.25(iii), converges strongly to w . The conclusion follows from Proposition 2.18(iii), respectively, (iv). \square

3. Unfolding and gradients. This section is devoted to the properties of the restriction of the unfolding operator to the space $W^{1,p}(\Omega)$. Some results require no extra hypotheses, but many others are sensitive to the boundary conditions and the regularity of the boundary itself.

Observe that, for w in $W^{1,p}(\Omega)$, one has

$$(3.1) \quad \nabla_y(\mathcal{T}_{\varepsilon}(w)) = \varepsilon \mathcal{T}_{\varepsilon}(\nabla w), \quad \forall w \in W^{1,p}(\Omega) \quad \text{a.e. for } (x, y) \in \Omega \times Y.$$

Then, Proposition 2.5(iv) implies that $\mathcal{T}_{\varepsilon}$ maps $W^{1,p}(\Omega)$ into $L^p(\Omega; W^{1,p}(Y))$.

For simplicity, we assume that $Y =]0, 1[^n$. Nevertheless, the results we prove here hold true in the case of a general Y , with minor modifications.

PROPOSITION 3.1 (gradient in the direction of a period). *Let k in $[1, \dots, n]$ and $\{w_{\varepsilon}\}$ be a bounded sequence in $L^p(\Omega)$, with $p \in]1, +\infty]$, satisfying*

$$(3.2) \quad \varepsilon \left\| \frac{\partial w_{\varepsilon}}{\partial x_k} \right\|_{L^p(\Omega)} \leq C.$$

Then, there exist a subsequence (still denoted ε) and \widehat{w} in $L^p(\Omega \times Y)$, with $\frac{\partial \widehat{w}}{\partial y_k}$ in $L^p(\Omega \times Y)$ such that

$$(3.3) \quad \begin{aligned} \mathcal{T}_{\varepsilon}(w_{\varepsilon}) &\rightharpoonup \widehat{w} \quad \text{weakly in } L^p(\Omega \times Y), \\ \varepsilon \mathcal{T}_{\varepsilon} \left(\frac{\partial w_{\varepsilon}}{\partial x_k} \right) &= \frac{\partial \mathcal{T}_{\varepsilon}(w_{\varepsilon})}{\partial y_k} \rightharpoonup \frac{\partial \widehat{w}}{\partial y_k} \quad \text{weakly in } L^p(\Omega \times Y), \quad (\text{weakly-}^* \text{ for } p = +\infty). \end{aligned}$$

Moreover, the limit function \widehat{w} is 1-periodic, with respect to the y_k coordinate.

Proof. Convergences (3.3) are a simple consequence of (3.1) and (3.2). It remains to prove the periodicity of \widehat{w} . Without loss of generality, assume $k = n$ and write $y = (y', y_n)$, with y' in $Y' \doteq]0, 1[^{n-1}$ and $y_n \in]0, 1[$.

Let $\psi \in \mathcal{D}(\Omega \times Y')$. Convergences (3.3) imply that the sequence $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$ is bounded in $L^p(\Omega \times Y'; W^{1,p}(0, 1))$ so that $\{\mathcal{T}_\varepsilon(w_\varepsilon)|_{\{y_n=s\}}\}$ is bounded in $L^p(\Omega \times Y')$ for every $s \in [0, 1]$. The periodicity with respect to y_n results from the following computation with an obvious change of variable:

$$\begin{aligned} & \int_{\Omega \times Y'} [\mathcal{T}_\varepsilon(w_\varepsilon)(x, (y', 1)) - \mathcal{T}_\varepsilon(w_\varepsilon)(x, (y', 0))] \psi(x, y') \, dx \, dy' \\ &= \int_{\Omega \times Y'} \left\{ w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 1) \right) - w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 0) \right) \right\} \psi(x, y') \, dx \, dy' \\ &= \int_{\Omega \times Y'} w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 0) \right) [\psi(x - \varepsilon e_n, y') - \psi(x, y')] \, dx \, dy', \\ &= \int_{\Omega \times Y'} \mathcal{T}_\varepsilon(w_\varepsilon)(x, (y', 0)) [\psi(x - \varepsilon e_n, y') - \psi(x, y')] \, dx \, dy', \end{aligned}$$

which goes to zero. \square

COROLLARY 3.2. *Let $\{w_\varepsilon\}$ be in $W^{1,p}(\Omega)$, with $p \in]1, +\infty[$, and assume that $\{w_\varepsilon\}$ is a bounded sequence in $L^p(\Omega)$ satisfying*

$$\varepsilon \|\nabla w_\varepsilon\|_{L^p(\Omega)} \leq C.$$

Then, there exist a subsequence (still denoted ε) and $\widehat{w} \in L^p(\Omega; W^{1,p}(Y))$ such that

$$\begin{aligned} \mathcal{T}_\varepsilon(w_\varepsilon) &\rightharpoonup \widehat{w} \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \varepsilon \mathcal{T}_\varepsilon(\nabla w_\varepsilon) &\rightharpoonup \nabla_y \widehat{w} \quad \text{weakly in } L^p(\Omega \times Y). \end{aligned}$$

Moreover, the limit function \widehat{w} is Y -periodic, i.e., belongs to $L^p(\Omega; W^{1,p}_{per}(Y))$, where $W^{1,p}_{per}(Y)$ denotes the Banach space of Y -periodic functions in $W^{1,p}_{loc}(\mathbb{R}^n)$, with the $W^{1,p}(Y)$ norm.

COROLLARY 3.3. *Let p be in $]1, +\infty[$ and $\{w_\varepsilon\}$ be a sequence converging weakly in $W^{1,p}(\Omega)$ to w . Then,*

$$\mathcal{T}_\varepsilon(w_\varepsilon) \rightharpoonup w \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)).$$

Furthermore, if $\{w_\varepsilon\}$ converges strongly to w in $L^p(\Omega)$, the above convergence is strong (this is the case if, for example, $W^{1,p}(\Omega)$ is compactly embedded in $L^p(\Omega)$).

Proof. Using (3.1), since $\{w_\varepsilon\}$ weakly converges, one has the estimates

$$\begin{aligned} \|\mathcal{T}_\varepsilon(w_\varepsilon)\|_{L^p(\Omega \times Y)} &\leq C, \\ \|\nabla_y(\mathcal{T}_\varepsilon(w_\varepsilon))\|_{L^p(\Omega \times Y)} &\leq \varepsilon C, \end{aligned}$$

so that there exist a subsequence (still denoted ε) and \widehat{w} in $L^p(\Omega; W^{1,p}(Y))$ such that

$$(3.4) \quad \begin{aligned} \mathcal{T}_\varepsilon(w_\varepsilon) &\rightharpoonup \widehat{w} \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \nabla_y(\mathcal{T}_\varepsilon(w_\varepsilon)) &\rightarrow 0 \quad \text{strongly in } L^p(\Omega \times Y), \end{aligned}$$

and $\nabla_y \widehat{w} = 0$. Consequently, \widehat{w} does not depend on y , and Proposition 2.9(iii) immediately gives $w = \mathcal{M}_Y(\widehat{w}) = \widehat{w}$. Moreover, convergence (3.4) holds for the entire

sequence ε . Finally, if the sequence $\{w_\varepsilon\}$ converges strongly to w in $L^p(\Omega)$, so does the sequence $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$, thanks to Proposition 2.9(ii). \square

PROPOSITION 3.4. *Suppose that p is in $]1, +\infty[$. Let $\{w_\varepsilon\}$ be a sequence which converges strongly to some w in $W^{1,p}(\Omega)$. Then,*

- (i) $\mathcal{T}_\varepsilon(\nabla w_\varepsilon) \rightarrow \nabla w$ strongly in $L^p(\Omega \times Y)$,
- (ii) $\frac{1}{\varepsilon}(\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon)) \rightarrow y^c \cdot \nabla w$ strongly in $L^p(\Omega; W^{1,p}(Y))$,

where

$$y^c = \left(y_1 - \frac{1}{2}, \dots, y_n - \frac{1}{2} \right).$$

Proof. The first assertion follows from Proposition 2.9(i). To prove (ii), set

$$Z_\varepsilon \doteq \frac{1}{\varepsilon}(\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon)),$$

which has mean value zero in Y . Since

$$\nabla_y Z_\varepsilon = \frac{1}{\varepsilon} \nabla_y (\mathcal{T}_\varepsilon(w_\varepsilon)) = \mathcal{T}_\varepsilon(\nabla w_\varepsilon),$$

thanks to assertion (i),

$$\nabla_y Z_\varepsilon \rightarrow \nabla w \quad \text{strongly in } L^p(\Omega \times Y).$$

Then recall the Poincaré–Wirtinger inequality in Y :

$$(3.5) \quad \forall \psi \in W^{1,p}(Y), \quad \|\psi - \mathcal{M}_Y(\psi)\|_{L^p(Y)} \leq C \|\nabla \psi\|_{L^p(Y)}.$$

Applying it to the function $Z_\varepsilon - y^c \cdot \nabla w$ (which is of mean value zero) gives

$$(3.6) \quad \|Z_\varepsilon - y^c \cdot \nabla w\|_{L^p(\Omega \times Y)} \leq C \|\nabla_y Z_\varepsilon - \nabla w\|_{L^p(\Omega \times Y)},$$

and this concludes the proof. \square

THEOREM 3.5. *Suppose that p is in $]1, +\infty[$. Let $\{w_\varepsilon\}$ be a sequence converging weakly to some w in $W^{1,p}(\Omega)$. Up to a subsequence, there exists some \widehat{w} in $L^p(\Omega; W_{per}^{1,p}(Y))$ such that*

- (i) $\mathcal{T}_\varepsilon(\nabla w_\varepsilon) \rightharpoonup \nabla w + \nabla_y \widehat{w}$ weakly in $L^p(\Omega \times Y)$,
- (ii) $\frac{1}{\varepsilon}(\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon)) \rightharpoonup \widehat{w} + y^c \cdot \nabla w$ weakly in $L^p(\Omega; W^{1,p}(Y))$.

Moreover, $\mathcal{M}_Y(\widehat{w}) = 0$.

Proof. Following the same lines as in the previous proof, introduce

$$Z_\varepsilon = \frac{1}{\varepsilon}(\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon)),$$

which has mean value zero in Y . Since $\nabla_y Z_\varepsilon = \mathcal{T}_\varepsilon(\nabla w_\varepsilon)$, (ii) implies (i).

To prove (ii), note that the sequence $\{\nabla_y Z_\varepsilon\}$ is bounded in $L^p(\Omega \times Y)$. By (3.6),

$$\|Z_\varepsilon - y^c \cdot \nabla w\|_{L^p(\Omega \times Y)} \leq C$$

so that there exists \widehat{w} in $L^p(\Omega; W^{1,p}(Y))$ such that, up to a subsequence,

$$Z_\varepsilon - y^c \cdot \nabla w \rightharpoonup \widehat{w} \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)).$$

Since, by construction, $\mathcal{M}_Y(y^c)$ vanishes, so does $\mathcal{M}_Y(\widehat{w})$.

It remains to prove the Y -periodicity of \widehat{w} . This is obtained in the same way as in the proof of Proposition 3.1 by using a test function $\psi \in \mathcal{D}(\Omega \times Y')$. One has successively,

$$\begin{aligned} & \int_{\Omega \times Y'} [Z_\varepsilon(x, (y', 1)) - Z_\varepsilon(x, (y', 0))] \psi(x, y') \, dx \, dy' \\ &= \int_{\Omega \times Y'} \frac{1}{\varepsilon} \left\{ w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 1) \right) - w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 0) \right) \right\} \psi(x, y') \, dx \, dy' \\ &= \int_{\Omega \times Y'} w_\varepsilon \left(\varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon(y', 0) \right) \frac{1}{\varepsilon} [\psi(x - \varepsilon e_n, y') - \psi(x, y')] \, dx \, dy', \\ &= \int_{\Omega \times Y'} \mathcal{T}_\varepsilon(w_\varepsilon)(x, (y', 0)) \frac{1}{\varepsilon} [\psi(x - \varepsilon e_n, y') - \psi(x, y')] \, dx \, dy'. \end{aligned}$$

By Proposition 2.9(ii), $\{\mathcal{T}_\varepsilon(w_\varepsilon)\}$ converges strongly to w in $L^p(\Omega \times Y)$, and by (3.7) (i), it converges weakly to the same w in $L^p(\Omega; W^{1,p}(Y))$. By the trace theorem in $W^{1,p}(Y)$, the trace of $\mathcal{T}_\varepsilon(w_\varepsilon)$ on $\Omega \times Y'$ converges weakly to w in $L^p(\Omega \times Y')$. Hence, the last integral converges to

$$(3.8) \quad - \int_{\Omega \times Y'} w(x) \frac{\partial \psi}{\partial x_n}(x, y') \, dx \, dy'.$$

Similarly, since $(y^c \cdot \nabla w)(y', 1) - (y^c \cdot \nabla w)(y', 0) = \frac{\partial w}{\partial x_n}$, we obtain

$$\begin{aligned} & \int_{\Omega \times Y'} [(y^c \cdot \nabla w)(y', 1) - (y^c \cdot \nabla w)(y', 0)] \psi(x, y') \, dx \, dy' \\ &= \int_{\Omega \times Y'} \frac{\partial w}{\partial x_n} \psi(x, y') \, dx \, dy' = - \int_{\Omega \times Y'} w(x) \frac{\partial \psi}{\partial x_n}(x, y') \, dx \, dy'. \end{aligned}$$

This, together with (3.8) and convergence (3.7)(ii), shows that

$$\int_{\Omega \times Y'} [\widehat{w}(x, (y', 1)) - \widehat{w}(x, (y', 0))] \psi(x, y') \, dx \, dy' = 0,$$

so that \widehat{w} is y_n -periodic. The same holds in the directions of all of the other periods. \square

Theorem 3.5 can be generalized to the case of $W^{k,p}(\Omega)$ -spaces, with $k \geq 1$ and $p \in]1, +\infty[$. In order to do so, for $r = (r_1, \dots, r_n) \in \mathbb{N}^n$ with $|r| = r_1 + \dots + r_n \leq k$, introduce the notation D^r and D_y^r :

$$D^r = \frac{\partial^{|r|}}{\partial x_1^{r_1} \dots \partial x_n^{r_n}}, \quad D_y^r = \frac{\partial^{|r|}}{\partial y_1^{r_1} \dots \partial y_n^{r_n}}.$$

Then the following result holds.

THEOREM 3.6. *Let $\{w_\varepsilon\}$ be a sequence converging weakly in $W^{k,p}(\Omega)$ to w , $k \geq 1$, and $p \in]1, +\infty[$. There exist a subsequence (still denoted ε) and \widehat{w} in the space $L^p(\Omega; W_{per}^{k,p}(Y))$ such that*

$$(3.9) \quad \begin{cases} \mathcal{T}_\varepsilon(D^l w_\varepsilon) \rightharpoonup D^l w & \text{weakly in } L^p(\Omega; W^{k-l,p}(Y)), \quad |l| \leq k-1, \\ \mathcal{T}_\varepsilon(D^l w_\varepsilon) \rightharpoonup D^l w + D_y^l \widehat{w} & \text{weakly in } L^p(\Omega \times Y), \quad |l| = k. \end{cases}$$

Furthermore, if $\{w_\varepsilon\}$ converges strongly to w in $W^{k-1,p}(\Omega)$, the above convergences are strong in $L^p(\Omega; W^{k-l,p}(Y))$ for $|l| \leq k-1$.

Proof. We give a brief proof for $k = 2$. The same argument generalizes for $k > 2$. If $|l| = 1$, the first convergence in (3.9) follows directly from Corollary 3.3. Set

$$W_\varepsilon = \frac{1}{\varepsilon^2} [\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon) - y^c \cdot \mathcal{M}_\varepsilon(\nabla w_\varepsilon)].$$

The sequence $\{w_\varepsilon\}$ is bounded in $W^{2,p}(\Omega)$, hence proceeding as in the proof of Proposition 2.25(iii), one obtains

$$\|W_\varepsilon\|_{L^p(\Omega \times Y)} \leq C.$$

Moreover,

$$\nabla_y(W_\varepsilon) = \frac{1}{\varepsilon^2} (\mathcal{T}_\varepsilon(\nabla w_\varepsilon) - \mathcal{M}_\varepsilon(\nabla w_\varepsilon)),$$

and

$$D_y^l(W_\varepsilon) = \mathcal{T}_\varepsilon(D^l w_\varepsilon), \quad \text{with } |l| = 2.$$

This implies that the sequence $\{W_\varepsilon\}$ is bounded in $L^p(\Omega; W^{2,p}(Y))$. Therefore, there exist a subsequence (still denoted ε) and $\widetilde{w} \in L^p(\Omega; W^{2,p}(Y))$ such that

$$(3.10) \quad \begin{aligned} W_\varepsilon &\rightharpoonup \widetilde{w} \quad \text{weakly in } L^p(\Omega; W^{2,p}(Y)), \\ \frac{\partial W_\varepsilon}{\partial y_i} &= \frac{1}{\varepsilon^2} \left(\mathcal{T}_\varepsilon \left(\frac{\partial w_\varepsilon}{\partial x_i} \right) - \mathcal{M}_\varepsilon \left(\frac{\partial w_\varepsilon}{\partial x_i} \right) \right) \rightharpoonup \frac{\partial \widetilde{w}}{\partial y_i} \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)). \end{aligned}$$

Consequently,

$$(3.11) \quad D_y^l(W_\varepsilon) = \mathcal{T}_\varepsilon(D^l w_\varepsilon) \rightharpoonup D_y^l \widetilde{w} \quad \text{weakly in } L^p(\Omega \times Y), \quad |l| = 2.$$

Now, apply Theorem 3.5 to each of the derivatives $\frac{\partial w_\varepsilon}{\partial x_i}$, $i \in \{1, \dots, n\}$. There exist a subsequence (still denoted ε) and $\widehat{w}_i \in L^p(\Omega; W_{per}^{1,p}(Y))$ such that $\mathcal{M}_Y(\widehat{w}_i) \equiv 0$ and

$$\frac{1}{\varepsilon} \left(\mathcal{T}_\varepsilon \left(\frac{\partial w_\varepsilon}{\partial x_i} \right) - \mathcal{M}_\varepsilon \left(\frac{\partial w_\varepsilon}{\partial x_i} \right) \right) \rightharpoonup y^c \cdot \nabla \frac{\partial w}{\partial x_i} + \widehat{w}_i \quad \text{weakly in } L^p(\Omega \times Y).$$

Then (3.10) gives

$$(3.12) \quad \forall i \in \{1, \dots, n\}, \quad \frac{\partial \widetilde{w}}{\partial y_i} = y^c \cdot \nabla \frac{\partial w}{\partial x_i} + \widehat{w}_i.$$

Set

$$\widehat{w} = \widetilde{w} - \frac{1}{2} \sum_{i,j=1}^n (y_i^c y_j^c - \mathcal{M}_Y(y_i^c y_j^c)) \frac{\partial^2 w}{\partial x_i \partial x_j}.$$

By construction, the function \widehat{w} belongs to $L^p(\Omega; W^{2,p}(Y))$. Furthermore,

$$\mathcal{M}_Y(\widehat{w}) = 0, \quad \frac{\partial \widehat{w}}{\partial y_i} = \frac{\partial \widetilde{w}}{\partial y_i} - y^c \cdot \nabla \left(\frac{\partial w}{\partial x_i} \right) = \widehat{w}_i, \quad \text{and} \quad \mathcal{M}_Y(\nabla_y \widehat{w}) = 0.$$

The last equality implies that \widehat{w} belongs to $L^p(\Omega; W_{per}^{2,p}(Y))$. Finally from (3.12) one gets

$$D_y^l \widetilde{w} = D^l w + D_y^l \widehat{w}, \quad \text{with} \quad |l| = 2,$$

which together with (3.11), proves the last convergence of (3.9). \square

COROLLARY 3.7. *Let $\{w_\varepsilon\}$ be a sequence converging weakly in $W^{2,p}(\Omega)$ to w , and $p \in]1, +\infty[$. Then, there exist a subsequence (still denoted ε) and \widehat{w} in the space $L^p(\Omega; W_{per}^{2,p}(Y))$ such that*

$$\frac{1}{\varepsilon^2} \left[\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon) - y^c \cdot \mathcal{M}_\varepsilon(\nabla w_\varepsilon) \right] \rightharpoonup \frac{1}{2} \sum_{i,j=1}^n (y_i^c y_j^c - \mathcal{M}_Y(y_i^c y_j^c)) \frac{\partial^2 w}{\partial x_i \partial x_j} + \widehat{w}$$

weakly in $L^p(\Omega; W^{2,p}(Y))$, where \widehat{w} is such that $\mathcal{M}_Y(\widehat{w}) = 0$.

Remark 3.8. For the case $Y =]0, 1[^n$, y^c was defined in Proposition 3.4. For a general Y , all of the statements of this section hold true, with $y^c = y - \mathcal{M}_Y(y)$.

4. Macro-micro decomposition: The scale-splitting operators \mathcal{Q}_ε and \mathcal{R}_ε . In this section, we give a different proof of Theorem 3.5, which was the one given originally in [24]. It is based on a scale-separation decomposition which is useful in some specific situations, for example, in the statement of general corrector results (see section 6).

The procedure is based on a splitting of functions ϕ in $W^{1,p}(\Omega)$ (or in $W_0^{1,p}(\Omega)$) for $p \in [1, +\infty]$, in the form

$$\phi = \mathcal{Q}_\varepsilon(\phi) + \mathcal{R}_\varepsilon(\phi),$$

where $\mathcal{Q}_\varepsilon(\phi)$ is an approximation of ϕ having the same behavior as ϕ , while $\mathcal{R}_\varepsilon(\phi)$ is a remainder of order ε . Applied to the sequence $\{w_\varepsilon\}$ converging weakly to w in $W^{1,p}(\Omega)$, it shows that, while $\{\nabla w_\varepsilon\}$, $\{\nabla(\mathcal{Q}_\varepsilon(w_\varepsilon))\}$ and $\{\mathcal{T}_\varepsilon(\nabla \mathcal{Q}_\varepsilon(w_\varepsilon))\}$ have the same weak limit ∇w in $L^p(\Omega)$, respectively, in $L^p(\Omega \times Y)$, the sequence $\mathcal{T}_\varepsilon(\nabla(\mathcal{R}_\varepsilon(w_\varepsilon)))$ converges weakly in $L^p(\Omega \times Y)$ to $\nabla_y \widehat{w}'$ for some \widehat{w}' in $L^p(\Omega; W_{per}^{1,p}(Y))$.

We will distinguish between the case $W_0^{1,p}(\Omega)$ and the case $W^{1,p}(\Omega)$. For the former, any function ϕ in $W_0^{1,p}(\Omega)$ is extended by zero to the whole of \mathbb{R}^n , and this extension is denoted by $\widetilde{\phi}$. In the latter case, we suppose that $\partial\Omega$ is smooth enough so that there exists a continuous extension operator $\mathcal{P} : W^{1,p}(\Omega) \mapsto W^{1,p}(\mathbb{R}^n)$ satisfying

$$\|\mathcal{P}(\phi)\|_{W^{1,p}(\mathbb{R}^n)} \leq C \|\phi\|_{W^{1,p}(\Omega)}, \quad \forall \phi \in W^{1,p}(\Omega),$$

where C is a constant depending on p and $\partial\Omega$ only.

The construction of \mathcal{Q}_ε is based on the Q_1 interpolate of some discrete approximation, as is customary in FEM. The idea of using these types of interpolate was already present in Griso [40], [41] for the study of truss-like structures. For the purpose of this paper, it is enough to take the average on $\varepsilon\xi + \varepsilon Y$ to construct the discrete approximations, but the average on $\varepsilon\xi + \varepsilon Y'$, where Y' is any fixed open subset of Y ,

or any open subset of a manifold of codimension 1 in Y . The only property which is needed is the Poincaré–Wirtinger inequality, which holds in both of these cases.

DEFINITION 4.1. *The operator $\mathcal{Q}_\varepsilon : L^p(\mathbb{R}^n) \mapsto W^{1,\infty}(\mathbb{R}^n)$, for $p \in [1, +\infty]$, is defined as follows:*

$$(4.1) \quad \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi) = \mathcal{M}_\varepsilon(\phi)(\varepsilon\xi) \quad \text{for } \xi \in \varepsilon\mathbb{Z}^n,$$

and for any $x \in \mathbb{R}^n$, we set

$$(4.2) \quad \begin{aligned} &\mathcal{Q}_\varepsilon(\phi)(x) \text{ is the } Q_1 \text{ interpolate of the values of } \mathcal{Q}_\varepsilon(\phi) \text{ at the vertices} \\ &\text{of the cell } \varepsilon \left[\frac{x}{\varepsilon} \right]_Y + \varepsilon Y. \end{aligned}$$

In the case of the space $W_0^{1,p}(\Omega)$, the operator $\mathcal{Q}_\varepsilon : W_0^{1,p}(\Omega) \mapsto W^{1,\infty}(\Omega)$ is defined by

$$\mathcal{Q}_\varepsilon(\phi) = \mathcal{Q}_\varepsilon(\tilde{\phi})|_\Omega,$$

where $\mathcal{Q}_\varepsilon(\tilde{\phi})$ is given by (4.1).

In the case of the space $W^{1,p}(\Omega)$, the operator $\mathcal{Q}_\varepsilon : W^{1,p}(\Omega) \mapsto W^{1,\infty}(\Omega)$ is defined by

$$\mathcal{Q}_\varepsilon(\phi) = \mathcal{Q}_\varepsilon(\mathcal{P}(\phi))|_\Omega,$$

where $\mathcal{Q}_\varepsilon(\mathcal{P}(\phi))$ is given by (4.1).

We start with the following estimates.

PROPOSITION 4.2 (properties of \mathcal{Q}_ε on \mathbb{R}^n). *For ϕ in $L^p(\mathbb{R}^n)$, $p \in [1, +\infty]$, there exists a constant C depending on n and Y only, such that*

$$\begin{aligned} \text{(i)} \quad &\|\mathcal{Q}_\varepsilon(\phi)\|_{L^p(\mathbb{R}^n)} \leq C\|\phi\|_{L^p(\mathbb{R}^n)}, & \text{(ii)} \quad &\|\nabla \mathcal{Q}_\varepsilon(\phi)\|_{L^p(\mathbb{R}^n)} \leq \frac{C}{\varepsilon}\|\phi\|_{L^p(\mathbb{R}^n)}, \\ \text{(iii)} \quad &\|\mathcal{Q}_\varepsilon(\phi)\|_{L^\infty(\mathbb{R}^n)} \leq \frac{C}{\varepsilon^{n/p}}\|\phi\|_{L^p(\mathbb{R}^n)}, & \text{(iv)} \quad &\|\nabla \mathcal{Q}_\varepsilon(\phi)\|_{L^\infty(\mathbb{R}^n)} \leq \frac{C}{\varepsilon^{1+n/p}}\|\phi\|_{L^p(\mathbb{R}^n)}. \end{aligned}$$

Furthermore, for any ψ in $L^p(Y)$,

$$(4.3) \quad \left\| \mathcal{Q}_\varepsilon(\phi)\psi \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \right\|_{L^p(\mathbb{R}^n)} \leq C\|\phi\|_{L^p(\mathbb{R}^n)}\|\psi\|_{L^p(Y)};$$

if ψ is in $W_{per}^{1,p}(Y)$, then

$$(4.4) \quad \left\| \mathcal{Q}_\varepsilon(\phi)\psi \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \right\|_{W^{1,p}(\mathbb{R}^n)} \leq \frac{C}{\varepsilon}\|\phi\|_{L^p(\mathbb{R}^n)}\|\psi\|_{W^{1,p}(Y)}.$$

Proof. By definition, the Q_1 interpolate is Lipschitz-continuous and reaches its maximum at some $\varepsilon\xi$. So, to estimate the L^∞ norm of $\mathcal{Q}_\varepsilon(\phi)$, it suffices to estimate the $\mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi)$'s. By (4.1),

$$(4.5) \quad |\mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi)|^p \leq \frac{1}{|Y|} \int_Y |\phi(\varepsilon\xi + \varepsilon z)|^p dz = \frac{1}{\varepsilon^n|Y|} \int_{\varepsilon\xi + \varepsilon Y} |\phi(x)|^p dx.$$

Since

$$\frac{1}{\varepsilon^n|Y|} \int_{\varepsilon\xi + \varepsilon Y} |\phi(x)|^p dx \leq \frac{1}{\varepsilon^n|Y|} \|\phi\|_{L^p(\mathbb{R}^n)}^p,$$

estimate (iii) follows, with $C = \frac{1}{|Y|^{1/p}}$.

The space $Q_1(Y)$ is of dimension 2^n , hence all of the norms are equivalent. So, there are constants c_1, c_2 , and c_3 (depending only upon p and Y) such that, for every $\Phi \in Q_1(Y)$,

$$\begin{aligned} \|\nabla\Phi\|_{L^\infty(Y)} &\leq c_1 \sum_{\kappa \in \{0,1\}^n} \left| \Phi \left(\sum_{j=1}^n \kappa_j b_j \right) \right|, \\ \|\Phi\|_{L^p(Y)} &\leq c_2 \left(\sum_{\kappa \in \{0,1\}^n} \left| \Phi \left(\sum_{j=1}^n \kappa_j b_j \right) \right|^p \right)^{1/p}, \\ \|\nabla\Phi\|_{L^p(Y)} &\leq c_3 \left(\sum_{\kappa \in \{0,1\}^n} \left| \Phi \left(\sum_{j=1}^n \kappa_j b_j \right) \right|^p \right)^{1/p}. \end{aligned}$$

Rescaling these inequalities for $\Phi(y) \doteq \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon y)$, gives

$$\begin{aligned} \|\nabla\mathcal{Q}_\varepsilon(\phi)\|_{L^\infty(\varepsilon\xi+\varepsilon Y)} &\leq \frac{c_1}{\varepsilon} \sum_{\kappa \in \{0,1\}^n} \left| \mathcal{Q}_\varepsilon(\phi) \left(\varepsilon\xi + \varepsilon \sum_{j=1}^n \kappa_j b_j \right) \right|, \\ \|\mathcal{Q}_\varepsilon(\phi)\|_{L^p(\varepsilon\xi+\varepsilon Y)} &\leq c_2 \varepsilon^{n/p} \left(\sum_{\kappa \in \{0,1\}^n} \left| \mathcal{Q}_\varepsilon(\phi) \left(\varepsilon\xi + \varepsilon \sum_{j=1}^n \kappa_j b_j \right) \right|^p \right)^{1/p}, \\ \|\nabla\mathcal{Q}_\varepsilon(\phi)\|_{L^p(\varepsilon\xi+\varepsilon Y)} &\leq c_3 \varepsilon^{n/p-1} \left(\sum_{\kappa \in \{0,1\}^n} \left| \mathcal{Q}_\varepsilon(\phi) \left(\varepsilon\xi + \varepsilon \sum_{j=1}^n \kappa_j b_j \right) \right|^p \right)^{1/p}. \end{aligned}$$

Using (4.5), we have

$$\|\nabla\mathcal{Q}_\varepsilon(\phi)\|_{L^\infty(\mathbb{R}^n)} \leq \frac{2^n c_1}{\varepsilon^{1+n/p} |Y|^{1/p}} \|\phi\|_{L^p(\mathbb{R}^n)},$$

which gives (iv). Similarly,

$$\|\mathcal{Q}_\varepsilon(\phi)\|_{L^p(\varepsilon\xi+\varepsilon Y)}^p \leq \frac{c_2^p}{|Y|} \sum_{\kappa \in \{0,1\}^n} \int_{\varepsilon\xi+\varepsilon \sum_{j=1}^n \kappa_j b_j + \varepsilon Y} |\phi(x)|^p dx,$$

which, by summation on $\xi \in \Xi_\varepsilon$, gives (i), with $C = \frac{(2c_2)^{n/p}}{|Y|^{1/p}}$.

Estimate (ii), with $C = \frac{(2c_3)^{n/p}}{|Y|^{1/p}}$, follows by a similar computation.

To prove (4.3), observe first that the function $\mathcal{Q}_\varepsilon(\phi)\psi(\{\frac{\cdot}{\varepsilon}\}_Y)$ belongs to $L^p(\mathbb{R}^n)$, since $\mathcal{Q}_\varepsilon(\phi)$ is in $L^\infty(\mathbb{R}^n)$ and $\psi(\{\frac{\cdot}{\varepsilon}\}_Y)$ is in $L^p(\mathbb{R}^n)$. Moreover,

$$\left\| \psi \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \right\|_{L^p(\varepsilon\xi+\varepsilon Y)}^p = \varepsilon^n \|\psi\|_{L^p(Y)}^p,$$

while, by (4.5),

$$\|\mathcal{Q}_\varepsilon(\phi)\|_{L^\infty(\varepsilon\xi+\varepsilon Y)}^p \leq \sum_{\kappa \in \{0,1\}^n} \frac{1}{\varepsilon^n |Y|} \int_{\varepsilon\xi+\varepsilon Y + \varepsilon \sum_{j=1}^n \kappa_j b_j} |\phi(x)|^p dx.$$

Using these two estimates and summing on Ξ_ε gives (4.3), with $C = \frac{2^{n/p}}{|Y|^{1/p}}$.

Estimate (4.4) is obtained in a similar fashion, with $C = \frac{(2)^{n/p} + (2c_3)^{n/p}}{|Y|^{1/p}}$. \square

COROLLARY 4.3. For ϕ in $L^p(\mathbb{R}^n)$, $p \in [1, +\infty[$, the following convergences hold:

$$\begin{aligned} \mathcal{Q}_\varepsilon(\phi) &\rightarrow \phi \quad \text{strongly in } L^p(\mathbb{R}^n), \\ \varepsilon \nabla \mathcal{Q}_\varepsilon(\phi) &\rightarrow 0 \quad \text{strongly in } (L^p(\mathbb{R}^n))^n. \end{aligned}$$

DEFINITION 4.4. The remainder $\mathcal{R}_\varepsilon(\phi)$ is given by

$$\mathcal{R}_\varepsilon(\phi) = \phi - \mathcal{Q}_\varepsilon(\phi) \quad \text{for any } \phi \in W^{1,p}(\Omega).$$

The following proposition is well-known from the FEM.

PROPOSITION 4.5 (properties of \mathcal{Q}_ε and \mathcal{R}_ε). For the case $W_0^{1,p}(\Omega)$, one has

- (i) $\|\mathcal{Q}_\varepsilon(\phi)\|_{W^{1,p}(\Omega)} \leq C \|\phi\|_{W_0^{1,p}(\Omega)}$,
- (ii) $\|\mathcal{R}_\varepsilon(\phi)\|_{L^p(\Omega)} \leq \varepsilon C \|\phi\|_{W_0^{1,p}(\Omega)}$,
- (iii) $\|\nabla \mathcal{R}_\varepsilon(\phi)\|_{L^p(\Omega)} \leq C \|\nabla \phi\|_{L^p(\Omega)}$.

The constant C depends on Y (via its diameter and its Poincaré–Wirtinger constant) only, and depends neither on Ω nor on ε .

Similarly, for the case $W^{1,p}(\Omega)$, one has

- (iv) $\|\mathcal{Q}_\varepsilon(\phi)\|_{W^{1,p}(\Omega)} \leq C \|\mathcal{P}\| \|\phi\|_{W^{1,p}(\Omega)}$,
- (v) $\|\mathcal{R}_\varepsilon(\phi)\|_{L^p(\Omega)} \leq \varepsilon C \|\mathcal{P}\| \|\phi\|_{W^{1,p}(\Omega)}$,
- (vi) $\|\nabla \mathcal{R}_\varepsilon(\phi)\|_{L^p(\Omega)} \leq C \|\mathcal{P}\| \|\nabla \phi\|_{L^p(\Omega)}$.

Moreover, in both cases,

$$(4.6) \quad \left\| \frac{\partial^2 \mathcal{Q}_\varepsilon(\phi)}{\partial x_i \partial x_j} \right\|_{L^p(\Omega)} \leq \frac{C'}{\varepsilon} \|\nabla \phi\|_{L^p(\Omega)} \quad \text{for } i, j \in [1, \dots, n], \quad i \neq j,$$

where C' does not depend on ε .

Proof. We start with ϕ in $W^{1,p}(\mathbb{R}^n)$. From Proposition 2.5(i) and inequality (3.5), we get

$$(4.7) \quad \|\phi - \mathcal{M}_\varepsilon(\phi)\|_{L^p(\mathbb{R}^n)} = |Y|^{-\frac{1}{p}} \|\mathcal{T}_\varepsilon(\phi) - \mathcal{M}_\varepsilon(\phi)\|_{L^p(\mathbb{R}^n \times Y)} \leq \varepsilon C \|\nabla \phi\|_{L^p(\mathbb{R}^n)}.$$

On the other hand, for any $\psi \in W^{1,p}(\text{interior}(\overline{Y \cup (Y + e_i)}))$, $i \in \{1, \dots, n\}$, we have

$$\begin{aligned} |\mathcal{M}_{Y+e_i}(\psi) - \mathcal{M}_Y(\psi)| &= |\mathcal{M}_Y(\psi(\cdot + e_i) - \psi(\cdot))| \\ &\leq \|\psi(\cdot + e_i) - \psi(\cdot)\|_{L^p(Y)} \leq C \|\nabla \psi\|_{L^p(Y \cup (Y + e_i))}. \end{aligned}$$

By a scaling argument and using Definition 4.1, this gives

$$(4.8) \quad |\mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi) - \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon e_i)| \leq \varepsilon C \|\nabla \phi\|_{L^p(\varepsilon(\xi + Y) \cup \varepsilon(\xi + e_i + Y))}$$

for all $\xi \in \varepsilon\mathbb{Z}^n$.

Let $x \in \varepsilon(\xi + Y)$, and set for every $\kappa = (\kappa_1, \dots, \kappa_n) \in \{0, 1\}^n$,

$$\bar{x}_l^{(\kappa_l)} = \begin{cases} \frac{x_l - \xi_l}{\varepsilon} & \text{if } \kappa_l = 1, \\ 1 - \frac{x_l - \xi_l}{\varepsilon} & \text{if } \kappa_l = 0. \end{cases}$$

If $\xi \in \varepsilon\mathbb{Z}^n$, for every $\kappa \in \{0, 1\}^n$, by definition we have

$$(4.9) \quad \mathcal{Q}_\varepsilon(\phi)(x) = \sum_{\kappa \in \{0,1\}^n} \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon\kappa) \bar{x}_1^{(\kappa_1)} \dots \bar{x}_n^{(\kappa_n)},$$

and so, for example,

$$\begin{aligned} & \frac{\partial \mathcal{Q}_\varepsilon(\phi)}{\partial x_1}(x) \\ &= \sum_{\kappa_2, \dots, \kappa_n} \frac{\mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon(1, \kappa_2, \dots, \kappa_n)) - \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon(0, \kappa_2, \dots, \kappa_n))}{\varepsilon} \bar{x}_2^{(\kappa_2)} \dots \bar{x}_n^{(\kappa_n)}, \end{aligned}$$

and a same expression for the other derivatives. This last formula and (4.7)–(4.9) imply estimate (i) written in \mathbb{R}^n .

Now, from (4.9), we get

$$\phi(x) - \mathcal{Q}_\varepsilon(\phi)(x) = \sum_{\kappa \in \{0,1\}^n} \left(\phi(x) - \mathcal{Q}_\varepsilon(\phi)(\varepsilon\xi + \varepsilon\kappa) \right) \bar{x}_1^{(\kappa_1)} \dots \bar{x}_n^{(\kappa_n)},$$

and (ii) (in \mathbb{R}^n) follows by using estimate (4.7). Estimate (iii) (again in \mathbb{R}^n) is straightforward from the previous ones.

If ϕ is in $W_0^{1,p}(\Omega)$, let $\tilde{\phi}$ be its extension to the whole of \mathbb{R}^n . To derive (i)–(iii), it suffices to write down the estimates in \mathbb{R}^n obtained above. Similarly, applying them to $\mathcal{P}(\phi)$ for ϕ in $W^{1,p}(\Omega)$ gives (iv)–(vi).

To finish the proof, it remains to show estimate (4.6). To do so, it is enough to take the derivative with respect to any x_k , with $k \neq 1$ in the formula of $\frac{\partial \mathcal{Q}_\varepsilon(\phi)}{\partial x_1}$ above, and use estimate (4.8). \square

Remark 4.6. By construction (see explicit formula (4.9)), the function $\mathcal{Q}_\varepsilon(\phi)$ is separately piecewise linear on each cell. Observe also that, for any $k \in \{1, \dots, n\}$, $\frac{\partial \mathcal{Q}_\varepsilon(\phi)}{\partial x_k}$ is independent of x_k in each cell $\varepsilon(\xi + Y)$.

PROPOSITION 4.7. *Let $\{w_\varepsilon\}$ be a sequence converging weakly in $W_0^{1,p}(\Omega)$ (resp. $W^{1,p}(\Omega)$) to w . Then, the following convergences hold:*

- (i) $\mathcal{R}_\varepsilon(w_\varepsilon) \rightarrow 0$ strongly in $L^p(\Omega)$,
- (ii) $\mathcal{Q}_\varepsilon(w_\varepsilon) \rightharpoonup w$ weakly in $W^{1,p}(\Omega)$,
- (iii) $\mathcal{T}_\varepsilon(\nabla \mathcal{Q}_\varepsilon(w_\varepsilon)) \rightharpoonup \nabla w$ weakly in $L^p(\Omega \times Y)$.

Proof. Convergence (i) is a direct consequence of estimate (ii) (resp. (v)) of Proposition 4.5, and it implies convergence (ii). Together with (i), Proposition 2.9(ii) implies $\mathcal{T}_\varepsilon(\mathcal{Q}_\varepsilon(w_\varepsilon)) \rightharpoonup w$ weakly in $L^p(\Omega \times Y)$. From (4.5),

$$\left\| \frac{\partial}{\partial x_i} \left(\frac{\partial \mathcal{Q}_\varepsilon(w_\varepsilon)}{\partial x_j} \right) \right\|_{L^p(\Omega)} \leq \frac{C}{\varepsilon} \quad \text{for } i, j \in [1, \dots, n], \quad i \neq j.$$

Now, by Proposition 3.1, there exist a subsequence (still denoted ε) and $\hat{w}_j \in L^p(\Omega \times Y)$, with $\frac{\partial \hat{w}_j}{\partial y_i} \in L^p(\Omega \times Y)$ such that, for $i \neq j$,

$$\begin{aligned} \mathcal{T}_\varepsilon \left(\frac{\partial \mathcal{Q}_\varepsilon(w_\varepsilon)}{\partial x_j} \right) &\rightharpoonup \hat{w}_j \quad \text{weakly in } L^p(\Omega \times Y), \\ \varepsilon \mathcal{T}_\varepsilon \left(\frac{\partial^2 \mathcal{Q}_\varepsilon(w_\varepsilon)}{\partial x_i \partial x_j} \right) &\rightharpoonup \frac{\partial \hat{w}_j}{\partial y_i} \quad \text{weakly in } L^p(\Omega \times Y), \end{aligned}$$

where \widehat{w}_j is y_i -periodic for every $i \neq j$. Moreover, from Remark 4.6, the function \widehat{w}_j does not depend on y_j , hence it is Y -periodic. But, by Remark 4.6 again, \widehat{w}_j is also piecewise linear, with respect to any variable y_i . Consequently, \widehat{w}_j is independent of y . On the other hand, from (ii) above we have

$$\frac{\partial Q_\varepsilon(w_\varepsilon)}{\partial x_j} \rightharpoonup \frac{\partial w}{\partial x_j} \quad \text{weakly in } L^p(\Omega).$$

Now Proposition 2.9(iii) gives $\widehat{w}_j = \frac{\partial w}{\partial x_j}$, and convergence (iii) holds for the whole sequence ε . \square

PROPOSITION 4.8 (Theorem 3.5 revisited). *Let $\{w_\varepsilon\}$ be a sequence converging weakly in $W_0^{1,p}(\Omega)$ (resp. in $W^{1,p}(\Omega)$) to w . Then, up to a subsequence there exists some \widehat{w}' in the space $L^p(\Omega; W_{per}^{1,p}(Y))$ such that the following convergences hold:*

$$\begin{aligned} \frac{1}{\varepsilon} \mathcal{T}_\varepsilon(\mathcal{R}_\varepsilon(w_\varepsilon)) &\rightharpoonup \widehat{w}' \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \mathcal{T}_\varepsilon(\nabla \mathcal{R}_\varepsilon(w_\varepsilon)) &\rightharpoonup \nabla_y \widehat{w}' \quad \text{weakly in } L^p(\Omega \times Y), \\ \mathcal{T}_\varepsilon(\nabla w_\varepsilon) &\rightharpoonup \nabla w + \nabla_y \widehat{w}' \quad \text{weakly in } L^p(\Omega \times Y). \end{aligned}$$

Actually, the connection with the \widehat{w} of Theorem 3.5 is given by

$$\widehat{w} = \widehat{w}' - \mathcal{M}_Y(\widehat{w}').$$

Proof. Due to estimates of Proposition 4.5, up to a subsequence, there exists \widehat{w}' in $L^p(\Omega; W_{per}^{1,p}(Y))$ such that

$$\begin{aligned} \frac{1}{\varepsilon} \mathcal{T}_\varepsilon(\mathcal{R}_\varepsilon(w_\varepsilon)) &\rightharpoonup \widehat{w}' \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)), \\ \mathcal{T}_\varepsilon(\nabla \mathcal{R}_\varepsilon(w_\varepsilon)) &\rightharpoonup \nabla_y \widehat{w}' \quad \text{weakly in } L^p(\Omega \times Y). \end{aligned}$$

Combining with convergence (iii) of Proposition 4.7 shows that

$$\mathcal{T}_\varepsilon(\nabla w_\varepsilon) \rightharpoonup \nabla w + \nabla_y \widehat{w}' \quad \text{weakly in } L^p(\Omega \times Y).$$

So $\nabla_y \widehat{w} \equiv \nabla_y \widehat{w}'$ in $L^p(\Omega \times Y)$. Since $\mathcal{M}_Y(\widehat{w}) = 0$, it follows that $\widehat{w} = \widehat{w}' - \mathcal{M}_Y(\widehat{w}')$. \square

Remark 4.9. In the previous proposition, one can actually compute the average of \widehat{w}' . One can check that $\mathcal{M}_Y(\widehat{w}') = -\mathcal{M}_Y(y) \cdot \nabla w$, and consequently,

$$\frac{1}{\varepsilon} \left(\mathcal{T}_\varepsilon(w_\varepsilon) - \mathcal{M}_\varepsilon(w_\varepsilon) \right) \rightharpoonup y \cdot \nabla w + \widehat{w}' \quad \text{weakly in } L^p(\Omega; W^{1,p}(Y)).$$

5. Periodic unfolding and the standard homogenization problem.

DEFINITION 5.1. *Let $\alpha, \beta \in \mathbb{R}$, such that $0 < \alpha < \beta$ and \mathcal{O} be an open subset of \mathbb{R}^n . Denote by $M(\alpha, \beta, \mathcal{O})$ the set of the $n \times n$ matrices $A = (a_{ij})_{1 \leq i, j \leq n} \in (L^\infty(\mathcal{O}))^{n \times n}$ such that, for any $\lambda \in \mathbb{R}^n$ and a.e. on \mathcal{O} ,*

$$(A(x)\lambda, \lambda) \geq \alpha|\lambda|^2, \quad |A(x)\lambda| \leq \beta|\lambda|.$$

Let $A^\varepsilon = (a_{ij}^\varepsilon)_{1 \leq i, j \leq n}$ be a sequence of matrices in $M(\alpha, \beta, \Omega)$. For f given in $H^{-1}(\Omega)$, consider the Dirichlet problem

$$(5.1) \quad \begin{cases} -\operatorname{div} (A^\varepsilon \nabla u_\varepsilon) = f & \text{in } \Omega, \\ u^\varepsilon = 0 & \text{on } \partial\Omega. \end{cases}$$

By the Lax–Milgram theorem, there exists a unique $u^\varepsilon \in H_0^1(\Omega)$ satisfying

$$(5.2) \quad \int_{\Omega} A^\varepsilon \nabla u_\varepsilon \nabla v \, dx = \langle f, v \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}, \quad \forall v \in H_0^1(\Omega),$$

which is the variational formulation of (5.1). Moreover, one has the apriori estimate

$$(5.3) \quad \|u_\varepsilon\|_{H_0^1(\Omega)} \leq \frac{1}{\alpha} \|f\|_{H^{-1}(\Omega)}.$$

Consequently, there exist u_0 in $H_0^1(\Omega)$ and a subsequence, still denoted ε , such that

$$(5.4) \quad u_\varepsilon \rightharpoonup u_0 \quad \text{weakly in } H_0^1(\Omega).$$

We are now interested to give a limit problem, the “homogenized” problem, satisfied by u_0 . This is called standard homogenization, and the answer, for some classes of A^ε , can be found in many works, starting with the classical book by Bensoussan, Lions, and Papanicolaou [11] (see, for instance, Cioranescu and Donato [30] and the references herein). We now recall it.

THEOREM 5.2 (standard periodic homogenization). *Let $A = (a_{ij})_{1 \leq i, j \leq n}$ belong to $M(\alpha, \beta, Y)$, where $a_{ij} = a_{ij}(y)$ are Y -periodic. Set*

$$(5.5) \quad A^\varepsilon(x) = \left(a_{ij} \left(\frac{x}{\varepsilon} \right) \right)_{1 \leq i, j \leq n} \quad \text{a.e. on } \Omega.$$

Let u_ε be the solution of the corresponding problem (5.1), with f in $H^{-1}(\Omega)$. Then the whole sequence $\{u_\varepsilon\}$ converges to a limit u_0 , which is the unique solution of the homogenized problem

$$(5.6) \quad \begin{cases} -\operatorname{div} (A^0 \nabla u_0) = \sum_{i, j=1}^n a_{ij}^0 \frac{\partial^2 u_0}{\partial x_i \partial x_j} = f & \text{in } \Omega, \\ u_0 = 0 & \text{on } \partial\Omega, \end{cases}$$

where the constant matrix $A^0 = (a_{ij}^0)_{1 \leq i, j \leq n}$ is elliptic and given by

$$(5.7) \quad a_{ij}^0 = \mathcal{M}_Y \left(a_{ij} - \sum_{k=1}^n a_{ik} \frac{\partial \widehat{\chi}_j}{\partial y_k} \right) = \mathcal{M}_Y(a_{ij}) - \mathcal{M}_Y \left(\sum_{k=1}^n a_{ik} \frac{\partial \widehat{\chi}_j}{\partial y_k} \right).$$

In (5.7), the functions $\widehat{\chi}_j$ ($j = 1, \dots, n$), often referred to as correctors, are the solutions of the cell systems

$$(5.8) \quad \begin{cases} - \sum_{i, k=1}^n \frac{\partial}{\partial y_i} \left(a_{ik} \frac{\partial (\widehat{\chi}_j - y_j)}{\partial y_k} \right) = 0 & \text{in } Y, \\ \mathcal{M}_Y(\widehat{\chi}_j) = 0, \\ \widehat{\chi}_j \quad Y\text{-periodic.} \end{cases}$$

As will be seen below, using the periodic unfolding, the proof of this theorem is elementary! Actually, with the same proof, a more general result can be obtained, with matrices A^ε .

THEOREM 5.3 (periodic homogenization via unfolding). *Let u_ε be the solution of problem (5.1), with f in $H^{-1}(\Omega)$, and $A^\varepsilon = (a_{ij}^\varepsilon)_{1 \leq i, j \leq n}$ be a sequence of matrices in $M(\alpha, \beta, \Omega)$. Suppose that there exists a matrix B such that*

$$(5.9) \quad B^\varepsilon \doteq \mathcal{T}_\varepsilon(A^\varepsilon) \rightarrow B \quad \text{strongly in } [L^1(\Omega \times Y)]^{n \times n}.$$

Then there exists $u_0 \in H_0^1(\Omega)$ and $\hat{u} \in L^2(\Omega; H_{per}^1(Y))$ such that

$$(5.10) \quad \begin{aligned} u_\varepsilon &\rightharpoonup u_0 && \text{weakly in } H_0^1(\Omega), \\ \mathcal{T}_\varepsilon(u_\varepsilon) &\rightharpoonup u_0 && \text{weakly in } L^2(\Omega; H^1(Y)), \\ \mathcal{T}_\varepsilon(\nabla u_\varepsilon) &\rightharpoonup \nabla u_0 + \nabla_y \hat{u} && \text{weakly in } L^2(\Omega \times Y), \end{aligned}$$

and the pair (u_0, \hat{u}) is the unique solution of the problem

$$(5.11) \quad \begin{cases} \forall \Psi \in H_0^1(\Omega), \forall \Phi \in L^2(\Omega; H_{per}^1(Y)), \\ \frac{1}{|Y|} \int_{\Omega \times Y} B(x, y) [\nabla u_0(x) + \nabla_y \hat{u}(x, y)] [\nabla \Psi(x) + \nabla_y \Phi(x, y)] \, dx dy \\ = \langle f, \Psi \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}. \end{cases}$$

Remark 5.4. System (5.11) is the unfolded formulation of the homogenized limit problem. It is of standard variational form in the space

$$\mathcal{H} = H_0^1(\Omega) \times L^2(\Omega; H_{per}^1(Y)/\mathbb{R}).$$

Remark 5.5. Hypothesis (5.9) implies that $B \in M(\alpha, \beta, \Omega \times Y)$.

Remark 5.6. If A^ε is of the form (5.5), then $B(x, y) = A(y)$. In the case where $A^\varepsilon(x) = A_1(x)A_2(\frac{x}{\varepsilon})$, one has (5.9), with $B(x, y) = A_1(x)A_2(y)$.

Remark 5.7. Let us point out that every matrix $B \in M(\alpha, \beta, \Omega \times Y)$ can be approached by the sequence of matrices A^ε in $M(\alpha, \beta, \Omega)$, with A^ε defined as follows:

$$A^\varepsilon = \begin{cases} \mathcal{U}_\varepsilon(B) & \text{in } \hat{\Omega}_\varepsilon, \\ \alpha I_n & \text{in } \Lambda_\varepsilon. \end{cases}$$

Proof of Theorem 5.3. Convergences (5.10) follow from estimate (5.3), Corollary 3.3, and Proposition 4.7, respectively.

Let us choose $v = \Psi$, with $\Psi \in H_0^1(\Omega)$ as test function in (5.2). The integration formula (2.5) from Proposition 2.7 gives

$$(5.12) \quad \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \mathcal{T}_\varepsilon(\nabla \Psi) \, dx dy \stackrel{\mathcal{T}_\varepsilon}{\simeq} \langle f, \Psi \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}.$$

We are allowed to pass to the limit in (5.12), due to (5.9), (5.10), and Proposition 2.9(i), to get

$$(5.13) \quad \frac{1}{|Y|} \int_{\Omega \times Y} B(x, y) [\nabla u_0(x) + \nabla_y \hat{u}(x, y)] \nabla \Psi(x) \, dx dy = \langle f, \Psi \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}.$$

Now, taking in (5.2), as test function $v^\varepsilon(x) = \varepsilon\Psi(x)\psi(\frac{x}{\varepsilon})$, $\Psi \in \mathcal{D}(\Omega)$, $\psi \in H^1_{per}(Y)$, one has, due to (2.5) and Proposition 2.7,

$$\begin{aligned} & \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \varepsilon\psi(y)\mathcal{T}_\varepsilon(\nabla\Psi) \, dx dy \\ & + \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon)\nabla_y\psi(y)\mathcal{T}_\varepsilon(\Psi) \, dx dy \stackrel{\mathcal{T}_\varepsilon}{\simeq} \langle f, v_\varepsilon \rangle_{H^{-1}(\Omega), H^1_0(\Omega)}. \end{aligned}$$

Since $v^\varepsilon \rightharpoonup 0$ in $H^1_0(\Omega)$, we get at the limit

$$\frac{1}{|Y|} \int_{\Omega \times Y} B(x, y) [\nabla u_0(x) + \nabla_y \hat{u}(x, y)] \Psi(x) \nabla_y \psi(y) \, dx dy = 0.$$

By the density of the tensor product $\mathcal{D}(\Omega) \otimes H^1_{per}(Y)$ in $L^2(\Omega; H^1_{per}(Y))$, this holds for all Φ in $L^2(\Omega; H^1_{per}(Y))$. \square

Remark 5.8. As in the two-scale method, (5.11) gives \hat{u} in terms of ∇u_0 and yields the standard form of the homogenized equation, i.e., (5.6). In the simple case where $A(x, y) = A(y) = (a_{ij}(y))_{1 \leq i, j \leq n}$, it is easily seen that

$$(5.14) \quad \hat{u} = \sum_{i=1}^n \frac{\partial u_0}{\partial x_i} \hat{\chi}_i,$$

and that the limit B is precisely the matrix A^0 which was defined in Theorem 5.2 by (5.7) and (5.8).

PROPOSITION 5.9 (convergence of the energy). *Under the hypotheses of Theorem 5.3,*

$$(5.15) \quad \begin{cases} \lim_{\varepsilon \rightarrow 0} \int_{\Omega} A^\varepsilon \nabla u_\varepsilon \nabla u_\varepsilon \, dx = \frac{1}{|Y|} \int_{\Omega \times Y} B [\nabla u_0 + \nabla_y \hat{u}] [\nabla u_0 + \nabla_y \hat{u}] \, dx dy, \\ \lim_{\varepsilon \rightarrow 0} \int_{\Lambda_\varepsilon} |\nabla u_\varepsilon|^2 \, dx = 0. \end{cases}$$

Proof. By standard weak lower-semicontinuity, one successively obtains

$$\begin{aligned} & \frac{1}{|Y|} \int_{\Omega \times Y} B [\nabla u_0 + \nabla_y \hat{u}] [\nabla u_0 + \nabla_y \hat{u}] \, dx dy \\ & \leq \liminf_{\varepsilon \rightarrow 0} \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \, dx dy \\ & \leq \limsup_{\varepsilon \rightarrow 0} \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \, dx dy \\ & \leq \limsup_{\varepsilon \rightarrow 0} \int_{\Omega} A^\varepsilon \nabla u_\varepsilon \nabla u_\varepsilon \, dx = \limsup_{\varepsilon \rightarrow 0} \langle f, u_\varepsilon \rangle_{H^{-1}(\Omega), H^1_0(\Omega)} \\ & = \langle f, u_0 \rangle_{H^{-1}(\Omega), H^1_0(\Omega)} = \frac{1}{|Y|} \int_{\Omega \times Y} B [\nabla u_0 + \nabla_y \hat{u}] [\nabla u_0 + \nabla_y \hat{u}] \, dx dy, \end{aligned}$$

which gives the first convergence in (5.15), as well as

$$\limsup_{\varepsilon \rightarrow 0} \int_{\Lambda_\varepsilon} A^\varepsilon \nabla u_\varepsilon \nabla u_\varepsilon \, dx = 0,$$

which implies the second convergence in (5.15). \square

Remark 5.10. From the above proof, we also have

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \, dx \, dy \\ = \frac{1}{|Y|} \int_{\Omega \times Y} B[\nabla u_0 + \nabla_y \hat{u}] [\nabla u_0 + \nabla_y \hat{u}] \, dx \, dy. \end{aligned}$$

COROLLARY 5.11. *The following strong convergence holds:*

$$(5.16) \quad \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \rightarrow \nabla u_0 + \nabla_y \hat{u} \quad \text{strongly in } L^2(\Omega \times Y).$$

Proof. We have successively

$$\begin{aligned} \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon [\mathcal{T}_\varepsilon(\nabla u_\varepsilon) - \nabla u_0 - \nabla_y \hat{u}] [\mathcal{T}_\varepsilon(\nabla u_\varepsilon) - \nabla u_0 - \nabla_y \hat{u}] \, dx \, dy \\ = \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \, dx \, dy \\ - \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon [\nabla u_0 + \nabla_y \hat{u}] \mathcal{T}_\varepsilon(\nabla u_\varepsilon) \, dx \, dy \\ - \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon \mathcal{T}_\varepsilon(\nabla u_\varepsilon) [\nabla u_0 + \nabla_y \hat{u}] \, dx \, dy \\ + \frac{1}{|Y|} \int_{\Omega \times Y} B^\varepsilon [\nabla u_0 + \nabla_y \hat{u}] [\nabla u_0 + \nabla_y \hat{u}] \, dx \, dy. \end{aligned}$$

Each term in the right-hand side converges, the first one due to Remark 5.10, and the others due to (5.10) and hypothesis (5.9). So, the right-hand side term converges to zero. Then convergence (5.16) is a consequence of the ellipticity of B^ε . \square

Remark 5.12. One can consider problem (5.1) with a homogeneous Neumann boundary condition on $\partial\Omega$ provided a zero order term is added to the operator. This problem is variational on the space $H^1(\Omega)$ without any regularity condition on the boundary. The exact same method applies and gives the corresponding limit problem. In order for a nonhomogeneous Neumann boundary condition (or Robin condition) on $\partial\Omega$ to make sense, a well-behaved trace operator is needed from $H^1(\Omega)$ to $L^2(\Omega)$. In that case, the same method applies.

6. Some corrector results and error estimates. Under additional regularity assumptions on the homogenized solution u_0 and the cell-functions $\hat{\chi}_j$, the strong convergence for the gradient of u_0 with a corrector is known (cf. [11] Chapter 1, section 5, [30] Chapter 8, section 3 and references therein). More precisely, suppose that $\nabla_y \hat{\chi}_j \in (L^r(Y))^n$, $j = 1, \dots, n$ and $\nabla u^0 \in L^s(\Omega)$, with $1 \leq r, s < +\infty$ and such that $1/r + 1/s = 1/2$. Then

$$\nabla u_\varepsilon - \nabla u_0 - \sum_{j=1}^n \frac{\partial u_0}{\partial x_j} (\nabla_y \hat{\chi}_j) \left(\frac{\cdot}{\varepsilon} \right) \rightarrow 0 \quad \text{strongly in } L^2(\Omega).$$

Our next result gives a corrector result *without any additional regularity assumption* on $\hat{\chi}_j$, and its proof reduces to a few lines. We also include a new type of corrector.

THEOREM 6.1. *Under the hypotheses of Theorem 5.2, one has*

$$(6.1) \quad \nabla u_\varepsilon - \nabla u_0 - \mathcal{U}_\varepsilon(\nabla_y \hat{u}) \rightarrow 0 \quad \text{strongly in } L^2(\Omega).$$

In the case where $A^\varepsilon(x) = A(\{\frac{x}{\varepsilon}\}_Y)$, the function $u_0 + \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon(\frac{\partial u_0}{\partial x_i}) \chi_i(\{\frac{\cdot}{\varepsilon}\}_Y)$ belongs to $H^1(\Omega)$, and one has

$$(6.2) \quad u_\varepsilon - u_0 - \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \chi_i \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \rightarrow 0 \quad \text{strongly in } H^1(\Omega).$$

Proof. From (5.15), (5.16), and Proposition 2.18(iii), one immediately has

$$\nabla u_\varepsilon - \mathcal{U}_\varepsilon(\nabla u_0) - \mathcal{U}_\varepsilon(\nabla_y \hat{u}) \rightarrow 0 \quad \text{strongly in } L^2(\Omega).$$

But since ∇u_0 belongs to $L^2(\Omega)$, Corollary 2.26 implies that

$$\mathcal{U}_\varepsilon(\nabla u_0) \rightarrow \nabla u_0 \quad \text{strongly in } L^2(\Omega),$$

whence (6.1). From (4.4) in Proposition 4.2, the function $u_0 + \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon(\frac{\partial u_0}{\partial x_i}) \chi_i(\{\frac{\cdot}{\varepsilon}\}_Y)$ belongs to $H^1(\Omega)$. From (5.14), we obtain

$$\begin{aligned} \nabla u_0 + \mathcal{U}_\varepsilon(\nabla_y \hat{u}) - \nabla \left[u_0 + \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \chi_i \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \right] \\ = - \sum_{i=1}^n \left[\mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) - \mathcal{M}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \right] \nabla_y \chi_i \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right) \\ - \varepsilon \sum_{i=1}^n \nabla \left[\mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \right] \chi_i \left(\left\{ \frac{\cdot}{\varepsilon} \right\}_Y \right), \end{aligned}$$

and using estimate (4.2), Proposition 2.25(iii), and Corollary 4.3, one immediately gets the strong convergence in $L^2(\Omega)$ of the right-hand side in the above equality. Thanks to (6.1), one has (6.2). \square

We end this section by recalling the error estimates obtained by Griso in [42], [44], and [45] for problem (5.1), with $f \in L^2(\Omega)$.

THEOREM 6.2 (see [42], [44]). *Suppose that $\partial\Omega$ is of class $C^{1,1}$. The solution u_ε of (5.1) satisfies the following estimates:*

$$\begin{aligned} \left\| \nabla u_\varepsilon - \nabla u_0 - \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \nabla_y \hat{\chi}_i \left(\left\{ \frac{\cdot}{\varepsilon} \right\} \right) \right\|_{[L^2(\Omega)]^n} &\leq C\varepsilon^{1/2} \|f\|_{L^2(\Omega)}, \\ \|u_\varepsilon - u_0\|_{L^2(\Omega)} + \left\| \rho \left(\nabla u_\varepsilon - \nabla u_0 - \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \nabla_y \hat{\chi}_i \left(\frac{\cdot}{\varepsilon} \right) \right) \right\|_{[L^2(\Omega)]^n} &\leq C\varepsilon \|f\|_{L^2(\Omega)}, \end{aligned}$$

where $\hat{\chi}_i$ for $i = 1, \dots, n$ is defined by (5.8) and $\rho = \rho(x)$ is the distance between $x \in \Omega$ and the boundary $\partial\Omega$. The constant C depends on n, A , and $\partial\Omega$.

COROLLARY 6.3 (see [44]). *Let Ω' be an open set strongly included in Ω , then*

$$\left\| u_\varepsilon - u_0 - \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \hat{\chi}_i \left(\frac{\cdot}{\varepsilon} \right) \right\|_{H^1(\Omega')} \leq C\varepsilon \|f\|_{L^2(\Omega)}.$$

The constant depends on n, A, Ω' , and $\partial\Omega$.

In what follows in this paragraph, we suppose that the open set Ω is a bounded domain in \mathbb{R}^n , $n = 2$ or 3 , of polygonal ($n = 2$) or polyhedral ($n = 3$) boundary. We

assume that Ω is on one side only of its boundary, and that Γ_0 is the union of some edges ($n = 2$) or some faces ($n = 3$) of $\partial\Omega$. Recall that classical regularity results show that the solution u_0 of the homogenized problem (5.6) belongs to $H^{1+s}(\Omega)$ for s in $]1/2, 1[$ ($s = 1$ if the domain is convex) depending only on $\partial\Omega$, on A^0 , and satisfies the estimate

$$\|\nabla u_0\|_{H^{1+s}(\Omega)} \leq C\|f\|_{L^2(\Omega)}.$$

The error estimate for this case is given in the following result.

THEOREM 6.4 (see [45]). *The solution u_ε of problem (5.1) satisfies the estimate*

$$\begin{aligned} & \left\| \nabla u_\varepsilon - \nabla u_0 - \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \nabla_y \hat{\chi}_i \left(\frac{\cdot}{\varepsilon} \right) \right\|_{[L^2(\Omega)]^n} \leq C\varepsilon^{s/2} \|f\|_{L^2(\Omega)}, \\ \|u_\varepsilon - u_0\|_{L^2(\Omega)} + & \left\| \rho \left(\nabla u_\varepsilon - \nabla u_0 - \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \nabla_y \hat{\chi}_i \left(\frac{\cdot}{\varepsilon} \right) \right) \right\|_{[L^2(\Omega)]^n} \leq C\varepsilon^s \|f\|_{L^2(\Omega)}. \end{aligned}$$

The constants depend on n , A , and $\partial\Omega$.

COROLLARY 6.5 (see [45]). *Let Ω' be an open set strongly included in Ω , then*

$$\left\| u_\varepsilon - u_0 - \varepsilon \sum_{i=1}^n \mathcal{Q}_\varepsilon \left(\frac{\partial u_0}{\partial x_i} \right) \hat{\chi}_i \left(\frac{\cdot}{\varepsilon} \right) \right\|_{H^1(\Omega')} \leq C\varepsilon^s \|f\|_{L^2(\Omega)}.$$

The constant depends on n , A , Ω' , and $\partial\Omega$.

7. Periodic unfolding and multiscales. As we mentioned in the Introduction, the periodic unfolding method turns out to be particularly well-adapted to multiscales problems. As an example, we treat here a problem with two different small scales.

Consider two periodicity cells Y and Z , both having the properties introduced at the beginning of section 2 (each associated to its set of periods). Suppose that Y is “partitioned” in two nonempty disjoint open subsets Y_1 and Y_2 , i.e., such that $Y_1 \cap Y_2 = \emptyset$ and $\bar{Y} = \bar{Y}_1 \cup \bar{Y}_2$.

Let $A^{\varepsilon\delta}$ be a matrix field defined by

$$A^{\varepsilon\delta}(x) = \begin{cases} A_1 \left(\left\{ \frac{x}{\varepsilon} \right\}_Y \right) & \text{for } \left\{ \frac{x}{\varepsilon} \right\}_Y \in Y_1, \\ A_2 \left(\left\{ \frac{\left\{ \frac{x}{\varepsilon} \right\}_Y}{\delta} \right\}_Z \right) & \text{for } \left\{ \frac{x}{\varepsilon} \right\}_Y \in Y_2, \end{cases}$$

where A_1 is in $M(\alpha, \beta, Y_1)$ and A_2 in $M(\alpha, \beta, Z)$ (cf. Definition 5.1). Here we have two small scales, namely, ε and $\varepsilon\delta$, associated, respectively, to the cells Y and Z (see Figure 7).

Consider the problem

$$\int_\Omega A^{\varepsilon\delta} \nabla u_{\varepsilon\delta} \nabla w \, dx = \int_\Omega f w \, dx \quad \forall w \in H_0^1(\Omega),$$

with f in $L^2(\Omega)$. The Lax–Milgram theorem immediately gives the existence and uniqueness of $u_{\varepsilon\delta}$ in $H_0^1(\Omega)$ satisfying the estimate

$$\|u_{\varepsilon\delta}\|_{H_0^1(\Omega)} \leq \frac{1}{\alpha} \|f\|_{L^2(\Omega)}.$$

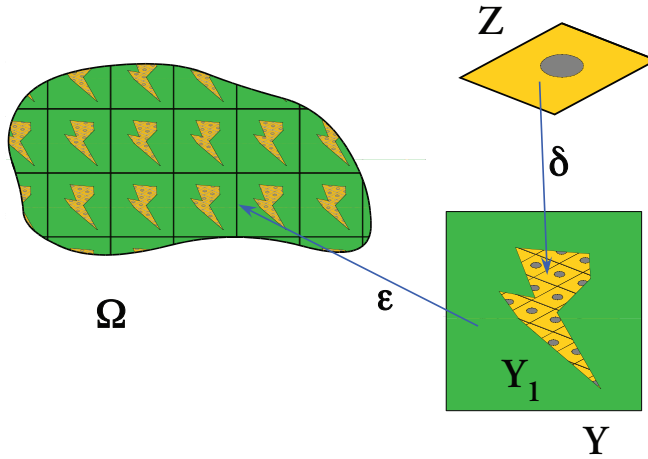


FIG. 7. A domain with periodic scales ε and $\varepsilon\delta$.

So, there is some u_0 such that, up to a subsequence,

$$u_{\varepsilon\delta} \rightharpoonup u_0 \quad \text{weakly in } H_0^1(\Omega).$$

Using the unfolding method for scale ε , as before we have

$$\begin{aligned} \mathcal{Q}_\varepsilon(u_{\varepsilon\delta}) &\rightharpoonup u_0 \quad \text{weakly in } H_0^1(\Omega), \\ \mathcal{T}_\varepsilon(u_{\varepsilon\delta}) &\rightharpoonup u_0 \quad \text{weakly in } L^2(\Omega; H^1(Y)), \\ \frac{1}{\varepsilon}\mathcal{T}_\varepsilon(\mathcal{R}_\varepsilon(u_{\varepsilon\delta})) &\rightharpoonup \hat{u} \quad \text{weakly in } L^2(\Omega; H^1(Y)), \\ \mathcal{T}_\varepsilon(\nabla u_{\varepsilon\delta}) &\rightharpoonup \nabla u_0 + \nabla_y \hat{u} \quad \text{in } L^2(\Omega \times Y). \end{aligned}$$

These convergences do not see the oscillations at the scale $\varepsilon\delta$. In order to capture them, one considers the restrictions to the set $\Omega \times Y_2$ defined by

$$v_{\varepsilon\delta}(x, y) \doteq \frac{1}{\varepsilon}\mathcal{T}_\varepsilon(\mathcal{R}_\varepsilon(u_{\varepsilon\delta}))|_{Y_2}.$$

Obviously,

$$v_{\varepsilon\delta} \rightharpoonup \hat{u}|_{Y_2} \quad \text{weakly in } L^2(\Omega; H^1(Y_2)).$$

Now, we apply to $v_{\varepsilon\delta}$, a similar unfolding operation, denoted \mathcal{T}_δ^y , for the variable y , thus adding a new variable $z \in Z$.

$$\mathcal{T}_\delta^y(v_{\varepsilon\delta})(x, y, z) = v_{\varepsilon\delta}\left(x, \delta\left[\frac{y}{\delta}\right]_Z + \delta z\right) \quad \text{for } x \in \Omega, \ y \in Y_2, \ \text{and } z \in Z.$$

It is essential to remark that all of the estimates and weak convergence properties which were shown for the original unfolding \mathcal{T}_ε still hold for \mathcal{T}_δ^y , with x being a mere parameter. For example, Proposition 4.7 and Theorem 3.5 adapted to this case imply that

$$\begin{aligned} \mathcal{T}_\delta^y(\nabla_y v_{\varepsilon\delta}) &\rightharpoonup \nabla_y \hat{u}|_{\Omega_2} + \nabla_z \tilde{u} \quad \text{weakly in } L^2(\Omega \times Y_2 \times Z), \\ \mathcal{T}_\delta^y(\mathcal{T}_\varepsilon(\nabla u_{\varepsilon\delta})) &\rightharpoonup \nabla u_0 + \nabla_y \hat{u} + \nabla_z \tilde{u} \quad \text{weakly in } L^2(\Omega \times Y_2 \times Z). \end{aligned}$$

Under these conditions, the limit functions u_0 , \hat{u} , and \tilde{u} are characterized by the following result.

THEOREM 7.1. *The functions*

$$u_0 \in H_0^1(\Omega), \quad \widehat{u} \in L^2(\Omega, H_{per}^1(Y)/\mathbb{R}), \quad \widetilde{u} \in L^2(\Omega \times \Omega_2, H_{per}^1(Z)/\mathbb{R})$$

are the unique solutions of the variational problem

$$\left\{ \begin{aligned} & \frac{1}{|Y||Z|} \int_{\Omega} \int_{Y_2} \int_Z A_2(z) \left\{ \nabla u_0 + \nabla_y \widehat{u} + \nabla_z \widetilde{u} \right\} \left\{ \nabla \Psi + \nabla_y \Phi + \nabla_z \Theta \right\} dx dy dz \\ & + \frac{1}{|Y|} \int_{\Omega} \int_{Y_1} A_1(y) \left\{ \nabla u_0 + \nabla_y \widehat{u} \right\} \left\{ \nabla \Psi + \nabla_y \Phi \right\} dx dy = \int_{\Omega} f \Psi dx, \\ & \forall \Psi \in H_0^1(\Omega), \quad \forall \Phi \in L^2(\Omega; H_{per}^1(Y)/\mathbb{R}), \forall \Theta \in L^2(\Omega \times \Omega_2, H_{per}^1(Z)/\mathbb{R}). \end{aligned} \right.$$

The proof uses test functions of the form

$$\Psi(x) + \varepsilon \Psi_1(x) \Phi_1\left(\frac{x}{\varepsilon}\right) + \varepsilon \delta \Psi_2(x) \Phi_2\left(\left\{\frac{x}{\varepsilon}\right\}_Y\right) \Theta_2\left(\frac{1}{\delta} \left\{\frac{x}{\varepsilon}\right\}_Y\right),$$

where Ψ, Ψ_1, Ψ_2 are in $\mathcal{D}(\Omega)$, Φ_1 in $H_{per}^1(Y)$, $\Phi_2 \in \mathcal{D}(Y_2)$, and $\Theta_2 \in H_{per}^1(Z)$.

Remark 7.2. The same theorem holds true for a general $A^{\varepsilon\delta}$ under the hypotheses

$$\begin{aligned} \mathcal{T}_{\varepsilon}(A^{\varepsilon\delta}) 1_{Y_1} &\rightarrow A_1 \quad \text{strongly in } [L^1(\Omega \times Y_1)]^{n \times n}, \\ \mathcal{T}_{\delta}^y(\mathcal{T}_{\varepsilon}(A^{\varepsilon\delta}) 1_{Y_2}) &\rightarrow A_2 \quad \text{strongly in } [L^1(\Omega \times Y_2 \times Z)]^{n \times n}. \end{aligned}$$

Proposition 5.9 (convergence of the energy) and Corollary 5.11 extend without any difficulty to the multiscale case.

PROPOSITION 7.3. *The convergence for the energy holds true:*

$$\begin{aligned} & \lim_{\varepsilon, \delta \rightarrow 0} \int_{\Omega} A^{\varepsilon\delta} \nabla u_{\varepsilon\delta} \nabla u_{\varepsilon\delta} dx \\ & = \frac{1}{|Y||Z|} \int_{\Omega} \int_{Y_2} \int_Z A_2(z) \left\{ \nabla u_0 + \nabla_y \widehat{u} + \nabla_z \widetilde{u} \right\} \left\{ \nabla u_0 + \nabla_y \widehat{u} + \nabla_z \widetilde{u} \right\} dx dy dz \\ & + \frac{1}{|Y|} \int_{\Omega} \int_{Y_1} A_1(y) \left\{ \nabla u_0 + \nabla_y \widehat{u} \right\} \left\{ \nabla u_0 + \nabla_y \widehat{u} \right\} dx dy. \end{aligned}$$

COROLLARY 7.4. *The following strong convergences hold true:*

$$\begin{aligned} \mathcal{T}_{\delta}^y(\nabla_y v_{\varepsilon\delta}) &\rightharpoonup \nabla_y \widehat{u}|_{\Omega_2} + \nabla_z \widetilde{u} \quad \text{strongly in } L^2(\Omega \times Y_2 \times Z), \\ \mathcal{T}_{\delta}^y(\mathcal{T}_{\varepsilon}(\nabla u_{\varepsilon\delta})) &\rightharpoonup \nabla u_0 + \nabla_y \widehat{u} + \nabla_z \widetilde{u} \quad \text{strongly in } L^2(\Omega \times Y_2 \times Z). \end{aligned}$$

Remark 7.5. A corrector result, similar to that of Theorem 6.1, can be obtained.

Remark 7.6. Theorem 7.1 can be extended to the case of any finite number of distinct scales by a simple reiteration.

8. Further developments. The unfolding method has some interesting properties which make it suitable for more general situations than that presented here. In problems which are set on a domain Ω_{ε} which depends on the parameter ε , it may be difficult to have a good notion of convergence for the sequence of solutions u_{ε} . The traditional way is to extend the solution by 0 outside Ω_{ε} ; however, this precludes the strong convergence of these extended functions in general. For the case of holes of the

size of order ε distributed ε -periodically, the unfolded sequence lives on a fixed domain. Similarly, for domains with ε -oscillating boundaries, a partial unfolding yields a function which is defined on a fixed domain.

We conclude by giving a list of publications making use of the unfolding method in several of these directions (both for linear and nonlinear problems).

- Reiterated homogenization: Meunier and Van Schaftingen [56].
- Electro-magnetism: Banks et al. [7], Bossavit, Griso, and Miara [15].
- Homogenization of thin piezoelectric shells: Ghergu et al. [39].
- Homogenization of diffusion deformation media: Griso and Rohan [46].
- Homogenization of the Stokes problem in porous media: Cioranescu, Damlamian, and Griso [25].
- Homogenization in perforated domains with Robin boundary conditions: Cioranescu, Donato and Zaki [31], [32].
- Homogenization in domains with oscillating boundaries: Damlamian and Pettersson [36].
- Homogenization of nonlinear integrals of the calculus of variations: Cioranescu, Damlamian, and De Arcangelis [27], [28], and [29].
- Homogenization of multivalued monotone operators of Leray–Lions type: Damlamian, Meunier, and Van Schaftingen [37].
- Thin junctions in linear elasticity: Blanchard, Gaudiello, and Griso [12], [13], Blanchard and Griso [14].
- Thin domains and free boundary problems arising in lubrication theory: Bayada, Martin, and Vazquez [9], [10].
- Elasticity problems in perforated domains: Griso and Sanchez-Rua [47].
- Neumann sieve and Dirichlet shield problems: Onofrei [60], Cioranescu et al. [26]. This last paper treats the case of domains with ε -periodically distributed “very small” holes (their size being a power of ε) on the boundary of which a homogeneous Dirichlet condition is prescribed. This requires the introduction of a rescaled unfolding operator (which originally appeared in the framework of the two-scale convergence in Casado-Díaz [20]).

Acknowledgments. We thank Petru Mironescu and Riccardo De Arcangelis for helpful comments and corrections.

REFERENCES

- [1] G. ALLAIRE, *Homogenization and two-scale convergence*, SIAM J. Math. Anal., 23 (1992), pp. 1482–1518.
- [2] G. ALLAIRE AND M. BRIANE, *Multiscale convergence and reiterated homogenization*, Proc. Roy. Soc. Edinburgh Sect. A, 126 (1996), pp. 297–342.
- [3] G. ALLAIRE AND C. CONCA, *Bloch wave homogenization and spectral asymptotic analysis*, J. Math. Pures Appl., 2 (1998), pp. 153–208.
- [4] G. ALLAIRE, C. CONCA, AND M. VANNINATHAN, *Spectral asymptotics of the Helmholtz model in fluid-solid structures*, Internat. J. Numer. Methods Engrg., 9 (1999), pp. 1463–1504.
- [5] T. ARBOGAST, J. DOUGLAS, AND U. HORNUNG, *Derivation of the double porosity model of single phase flow via homogenization theory*, SIAM J. Math. Anal., 21 (1990), pp. 823–836.
- [6] J. F. BABADJIAN AND M. BAÍA, *Multiscale nonconvex relaxation and application to thin films*, Asymptot. Anal., 48 (2006), pp. 173–218.
- [7] H. T. BANKS, V. A. BOKIL, D. CIORANESCU, N. L. GIBSON, G. GRISO, AND B. MIARA, *Homogenization of periodically varying coefficients in electromagnetic materials*, J. Sci. Comput., 28 (2006), pp. 191–221.
- [8] M. BARCHIESI, *Multiscale homogenization of convex functionals with discontinuous integrand*, J. Convex Anal., 14 (2007), pp. 205–226.

- [9] G. BAYADA, S. MARTIN, AND C. VAZQUEZ, *Two-scale homogenization of a hydrodynamic Elrod-Adams model*, *Asymptot. Anal.*, 44 (2005), pp. 75–110.
- [10] G. BAYADA, S. MARTIN, AND C. VAZQUEZ, *Homogenization of a nonlocal elastohydrodynamic lubrication problem: A new free boundary model*, *Math. Models Methods Appl. Sci.*, 15 (2005), pp. 1923–1956.
- [11] A. BENSOUSSAN, J. L. LIONS, AND G. PAPANICOLAOU, *Asymptotic Analysis for Periodic Structures*, *Stud. Math. Appl.* 5, North-Holland, Amsterdam, 1978.
- [12] D. BLANCHARD, A. GAUDIELLO, AND G. GRISO, *Junction of a periodic family of elastic rods with a 3d plate, Part I*, *J. Math. Pures Appl.*, 88 (2007), pp. 1–33.
- [13] D. BLANCHARD, A. GAUDIELLO, AND G. GRISO, *Junction of a periodic family of elastic rods with a thin plate, Part II*, *J. Math. Pures Appl.*, 88 (2007), pp. 149–190.
- [14] D. BLANCHARD AND G. GRISO, *Microscopic effects in the homogenization of the junction of rods and a thin plate*, *Asymptot. Anal.*, 56 (2008), pp. 1–36.
- [15] A. BOSSAVIT, G. GRISO, AND B. MIARA, *Modelling of periodic electromagnetic structures: Bianisotropic materials with memory effects*, *J. Math. Pures Appl.*, 84 (2005), pp. 819–850.
- [16] A. BOURGEAT, S. LUCKHAUS, AND A. MIKELIC, *Convergence of the homogenization process for a double-porosity model of immiscible two-phase flow*, *SIAM J. Math. Anal.*, 27 (1996), pp. 1520–1543.
- [17] A. BRAIDES AND A. DEFRANCESCHI, *Homogenization of Multiple Integrals*, *Oxford Lect. Ser. Math. Appl.* 12, Oxford University Press, Oxford, 1998.
- [18] A. BRAIDES AND D. LUKKASSEN, *Reiterated homogenization of integral functionals*, *Math. Models Methods Appl. Sci.*, 10 (2000), pp. 1–25.
- [19] D. A. G. BRUGGEMAN, *Berechnung verschiedener physikalischer konstanten von heterogenen substanzen*, *Ann. Physik.*, 24 (1935), p. 634.
- [20] J. CASADO-DÍAZ, *Two-scale convergence for nonlinear Dirichlet problems in perforated domains*, *Proc. Roy. Soc. Edinburgh Sect. A*, 130 (2000), pp. 249–276.
- [21] J. CASADO-DÍAZ AND M. LUNA-LAYNEZ, *A multiscale method to the homogenization of elastic thin reticulated structures*, in *Homogenization 2001*, GAKUTO Internat. Ser. Math. Sci. Appl. 18, Gakkōtoshō, Tokyo, 2003, pp. 155–168.
- [22] J. CASADO-DÍAZ, M. LUNA-LAYNEZ, AND J. D. MARTÍN, *An adaptation of the multi-scale methods for the analysis of very thin reticulated structures*, *C.R. Acad. Sci. Paris, Ser. 1*, 332 (2001), pp. 223–228.
- [23] J. CASADO-DÍAZ, M. LUNA-LAYNEZ, AND J. D. MARTÍN, *A new approach to the analysis of thin reticulated structures*, in *Homogenization 2001*, GAKUTO Internat. Ser. Math. Sci. Appl. 18, Gakkōtoshō, Tokyo, 2003, pp. 257–262.
- [24] D. CIORANESCU, A. DAMLAMIAN, AND G. GRISO, *Periodic unfolding and homogenization*, *C.R. Acad. Sci. Paris, Ser. 1*, 335 (2002), pp. 99–104.
- [25] D. CIORANESCU, A. DAMLAMIAN, AND G. GRISO, *The Stokes problem in perforated domains by the periodic unfolding method*, in *Proceedings of the Conference on New Trends in Continuum Mechanics*, M. Suliciu, ed., Theta, Bucarest, 2005, pp. 67–80.
- [26] D. CIORANESCU, A. DAMLAMIAN, G. GRISO, AND D. ONOFREI, *The periodic unfolding method for perforated domains and Neumann sieve models*, *J. Math. Pures Appl.*, 89 (2008), pp. 248–277.
- [27] D. CIORANESCU, A. DAMLAMIAN, AND R. DE ARCANGELIS, *Homogenization of nonlinear integrals via the periodic unfolding method*, *C.R. Acad. Sci. Paris, Ser. 1*, 339 (2005), pp. 77–82.
- [28] D. CIORANESCU, A. DAMLAMIAN, AND R. DE ARCANGELIS, *Homogenization of integrals with pointwise gradient constraints via the periodic unfolding method*, *Ricerche Mat.*, 55 (2006), pp. 31–54.
- [29] D. CIORANESCU, A. DAMLAMIAN, AND R. DE ARCANGELIS, *Homogenization of quasiconvex integrals via the periodic unfolding method*, *SIAM J. Math. Anal.*, 37 (2006), pp. 1435–1453.
- [30] D. CIORANESCU AND P. DONATO, *An Introduction to Homogenization*, *Oxford Lecture Ser. in Math. Appl.* 17, Oxford University Press, Oxford, 1999.
- [31] D. CIORANESCU, P. DONATO, AND R. ZAKI, *The periodic unfolding method in perforated domains*, *Port. Math.*, 63 (2006), pp. 467–496.
- [32] D. CIORANESCU, P. DONATO, AND R. ZAKI, *Asymptotic behavior of elliptic problems in perforated domains with nonlinear boundary conditions*, *Asymptot. Anal.*, 53 (2007), pp. 209–235.
- [33] G. DAL MASO AND A. DEFRANCESCHI, *Correctors for the homogenization of monotone operators*, *Differential Integral Equations*, 3 (1990), pp. 1151–1166.

- [34] A. DAMLAMIAN, *An elementary introduction to periodic unfolding*, in Proceedings of the Narvik Conference 2004, GAKUTO Internat. Ser. Math. Sci. Appl. 24, A. Damlamian, D. Lukkassen, A. Meidell, and A. Piatnitski, eds., Gakkōtoshō, Tokyo, 2006, pp. 119–136.
- [35] A. DAMLAMIAN AND P. DONATO, *H^0 -convergence and iterated homogenization*, Asymptot. Anal., 39 (2004), pp. 45–60.
- [36] A. DAMLAMIAN AND K. PETTERSSON, *Homogenization of oscillating boundaries*, Discrete Contin. Dyn. Syst., 23 (2009), pp. 197–219.
- [37] A. DAMLAMIAN, N. MEUNIER, AND J. VAN SCHAFTINGEN, *Periodic homogenization of monotone multivalued operators*, Nonlinear Anal., 67 (2007), pp. 3217–3239.
- [38] A. ENE AND J. SAINT JEAN PAULIN, *On a model of fractured porous media*, in Mathematical Modelling of Flow Through Porous Media, A. Bourgeat, C. Carasso, S. Luckhaus, and A. Mikelić, eds., World Scientific, River Edge, NJ, 1995, pp. 402–409.
- [39] M. GHERGU, G. GRISO, H. MECHKOUR, AND B. MIARA, *Homogenization of thin piezoelectric shells*, ESAIM: M2AN Math. Modeling Numer. Anal., 41 (2007), pp. 875–895.
- [40] G. GRISO, *Analyse asymptotique de structures réticulées*, Thèse Université Pierre et Marie Curie (Paris VI), Paris, 1996.
- [41] G. GRISO, *Thin reticulated structures*, in Progress in Partial Differential Equations, The Metz Surveys 3, M. Chipot, J. Saint Jean Paulin, and I. Shafirir, eds., Pitman, London, 1994, pp. 161–182.
- [42] G. GRISO, *Error estimate and unfolding for periodic homogenization*, Asymptot. Anal., 40 (2004), pp. 269–286.
- [43] G. GRISO, *Les méthodes d'éclatement en homogénéisation périodique et en élasticité linéarisée*, Thèse d'Habilitation, Université Pierre et Marie Curie, Paris, 2005.
- [44] G. GRISO, *Interior error estimates for periodic homogenization*, C.R. Acad. Sci. Paris, Ser. 1, 340 (2005), pp. 251–254.
- [45] G. GRISO, *Interior error estimates for periodic homogenization*, Anal. Appl., 4 (2006), pp. 61–79.
- [46] G. GRISO AND E. ROHAN, *On homogenization of diffusion-deformation problem in strongly heterogeneous media*, Ricerche Mat., 56 (2007), pp. 161–188.
- [47] G. GRISO AND T. SANCHEZ-RUA, *Homogenization of an elasticity problem for a catalyst support by using the unfolding method*, submitted.
- [48] M. LENCZNER, *Homogénéisation d'un circuit électrique*, C.R. Acad. Sci. Paris, Ser. 2, 324 (1997), pp. 537–542.
- [49] M. LENCZNER AND D. MERCIER, *Homogenization of periodic electrical networks including voltage to current amplifiers*, SIAM Multiscale Model. Simul., 2 (2004), pp. 359–397.
- [50] M. LENCZNER AND G. SENOUCI-BEREKSI, *Homogenization of electrical networks including voltage-to-voltage amplifiers*, Math. Models Methods Appl. Sci., 9 (1999), pp. 899–932.
- [51] M. LENCZNER, M. KADER, AND P. PERRIER, *Modèle à deux échelles de l'équation des ondes à coefficients oscillants*, C.R. Acad. Sci. Paris, Ser. 1, 328 (2000), pp. 335–340.
- [52] J. L. LIONS, D. LUKKASSEN, L. E. PERSSON, AND P. WALL, *Reiterated homogenization of monotone operators*, Chinese Ann. Math. Ser. B, 22 (2001), pp. 1–12.
- [53] D. LUKKASSEN, *Reiterated homogenization of non-standard Lagrangians*, C.R. Acad. Sci. Paris, Ser. 1, 332 (2001), pp. 999–1004.
- [54] D. LUKKASSEN AND G. W. MILTON, *On hierarchical structures and reiterated homogenization*, in Proceedings of the Conference on Function Spaces, Interpolation Theory and Related Topics in Honor of Jack Peetre on his 65th Birthday, M. Cwikl, M. Engliš, A. Kufner, L.-E. Persson, and G. Sparr, eds., Walter de Gruyter, Berlin, 2002, pp. 311–324.
- [55] D. LUKKASSEN, G. NGUETSENG, AND P. WALL, *Two-scale convergence*, Int. J. Pure Appl. Math., 2 (2002), pp. 35–86.
- [56] N. MEUNIER AND J. VAN SCHAFTINGEN, *Periodic reiterated homogenization for elliptic functions*, J. Math. Pures Appl., 9 (2005), pp. 1716–1743.
- [57] G. MILTON, *The Theory of Composites*, Cambridge University Press, Cambridge, 2002.
- [58] G. NGUETSENG, *A general convergence result for a functional related to the theory of homogenization*, SIAM J. Math. Anal., 20 (1989), pp. 608–629.
- [59] O. A. OLEINIK, A. S. SHAMAEV, AND G. A. YOSIFIAN, *Mathematical Problems in Elasticity and Homogenization*, North-Holland, Amsterdam, 1992.
- [60] D. ONOFREI, *The unfolding operator near a hyperplane and its applications to the Neumann sieve model*, Adv. Math. Sci. Appl., 16 (2006), pp. 239–258.

ON THE ASYMPTOTIC STABILITY OF SMALL NONLINEAR DIRAC STANDING WAVES IN A RESONANT CASE*

NABILE BOUSSAID†

Abstract. We study the behavior of perturbations of small nonlinear Dirac standing waves. We assume that the linear Dirac operator of reference $H = D_m + V$ has only two double eigenvalues and that degeneracies are due to a symmetry of H (theorem of Kramers). In this case, we can build a small four-dimensional manifold of stationary solutions tangent to the first eigenspace of H . Then we assume that a resonance condition holds, and we build a center manifold of real codimension 8 around each stationary solution. Inside this center manifold any H^s perturbation of stationary solutions, with $s > 2$, stabilizes towards a standing wave. We also build center-stable and center-unstable manifolds, each one of real codimension 4. Inside each of these manifolds, we obtain stabilization towards the center manifold in one direction of time, while in the other, we have instability. Eventually, outside all of these manifolds, we have instability in the two directions of time. For localized perturbations inside the center manifold, we obtain a nonlinear scattering result.

Key words. Dirac equation, nonlinear PDE, asymptotic stability, Strichartz estimates, smoothness estimates, center manifold

AMS subject classifications. 35Q40, 35Q55, 35Q75

DOI. 10.1137/070684641

Introduction. We study the asymptotic stability of stationary solutions of a time-dependent nonlinear Dirac equation.

A localized stationary solution of a given time-dependent equation represents a bound state of a particle. Like Rañada [39], we call it a *particle-like solution* (PLS). Many works have been devoted to the proof of the existence of such solutions for a wide variety of equations. Although their stability is a crucial problem (in particular, in a numerical computation or experiment), less attention has been devoted to this issue.

In this paper, we deal with the problem of stability of small PLSs of the following nonlinear Dirac equation:

$$(0.1) \quad i\partial_t\psi = (D_m + V)\psi + \nabla F(\psi),$$

where ∇F is the gradient of $F : \mathbb{C}^4 \mapsto \mathbb{R}$ for the standard scalar product of \mathbb{R}^8 . Here, D_m is the usual Dirac operator (see Thaller [48]) acting on $L^2(\mathbb{R}^3, \mathbb{C}^4)$:

$$D_m = \alpha \cdot (-i\nabla) + m\beta = -i \sum_{k=1}^3 \alpha_k \partial_k + m\beta,$$

where $m \in \mathbb{R}_+^*$, $\alpha = (\alpha_1, \alpha_2, \alpha_3)$, and β are \mathbb{C}^4 Hermitian matrices satisfying

$$\begin{cases} \alpha_i \alpha_k + \alpha_k \alpha_i = 2\delta_{ik} \mathbf{1}_{\mathbb{C}^4}, & i, k \in \{1, 2, 3\}, \\ \alpha_i \beta + \beta \alpha_i = \mathbf{0}_{\mathbb{C}^4}, & i \in \{1, 2, 3\}, \\ \beta^2 = \mathbf{1}_{\mathbb{C}^4}. \end{cases}$$

*Received by the editors March 8, 2007; accepted for publication (in revised form) June 4, 2008; published electronically November 26, 2008. This work was partially supported by EPSRC grant EP/D054621.

<http://www.siam.org/journals/sima/40-4/68464.html>

†Laboratoire de Mathématiques de Besançon, Université de Franche-Comté, 16 Route de Gray, F-25030 Besançon Cedex, France (nabile.boussaid@univ-fcomte.fr).

Here, we choose

$$\alpha_i = \begin{pmatrix} 0 & \sigma_i \\ \sigma_i & 0 \end{pmatrix} \quad \text{and} \quad \beta = \begin{pmatrix} I_{\mathbb{C}^2} & 0 \\ 0 & -I_{\mathbb{C}^2} \end{pmatrix},$$

where $\sigma_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, $\sigma_2 = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$, and $\sigma_3 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$.

In (0.1), V is the external potential field, and $F : \mathbb{C}^4 \mapsto \mathbb{R}$ is a nonlinearity with the following gauge invariance:

$$(0.2) \quad \forall(\theta, z) \in \mathbb{R} \times \mathbb{C}^4, \quad F(e^{i\theta}z) = F(z).$$

Some additional assumptions on F and V will be made in what follows. Nonlinearity with no potential arises in some Dirac models introduced by physicists to model either extended particles with self-interaction or particles in space-time with geometrical structure. In the latter case, physicists have shown that a relativistic theory sometimes imposes a fourth order nonlinear potential (i.e., a cubic nonlinearity) such as the square of a quadratic form on \mathbb{C}^4 ; see Rañada [39] and the references therein. We added a potential with special features to ensure the existence of small stationary solutions, whose existence and stability are easier to study.

Stationary solutions (PLSs) of (0.1) take the form $\psi(t, x) = e^{-iEt}\phi(x)$, where ϕ satisfies

$$(0.3) \quad E\phi = (D_m + V)\phi + \nabla F(\phi).$$

We show that there exists a manifold of small solutions to (0.3) tangent to the first eigenspace of $D_m + V$ (see Proposition 1.1 below).

Concerning the asymptotic stability in the Schrödinger equation, the question has been solved in several cases. For small stationary solutions in the simple eigenvalue case it has been studied by Soffer and Weinstein [43, 44], Pillet and Wayne [38], and Gustafson, Nakanishi, and Tsai [22]. For the two eigenvalue case under a resonance condition for an excited state, the problem has been studied by Tsai and Yau [51, 53, 54, 52, 50] and Soffer and Weinstein [45, 46]. Another problem has been studied by Cuccagna [16, 17, 18], who considered the case of a big PLS, when the linearized operator has only one eigenvalue, and obtained the asymptotic stability of the manifold of ground states. Schlag [41] proved that any ground state of the cubic nonlinear Schrödinger equation in dimension 3 is orbitally unstable but possesses a stable manifold of codimension 9.

We also would like to mention the works of Buslaev and Perel'mann [12, 10, 11, 9], Buslaev and Sulem [14, 13], Weder [55], and Krieger and Schlag [31] in the one-dimensional Schrödinger case. Krieger and Schlag [31] proved a result similar to [41] in the one-dimensional case. Concerning two-dimensional nonlinear Schrödinger equations, Mizumachi [35, 36] studied the properties of unstable solutions, while Kirr and Zarnescu [30] proved the existence of small asymptotically stable states in a semilinear case.

In [7], we prove that there are stable directions for the PLS manifold under a nonresonance assumption on the spectrum of $H := D_m + V$. This gives a stable manifold containing the PLS manifold. But we were not able to say anything about solutions starting outside the stable manifold.

The results we present here state the existence of a stable manifold and describe the behavior of solutions starting outside of it. In fact, we prove the instability of the perturbations starting outside the stable manifold. We also prove stabilization towards stationary solutions inside the stable manifold for H^s perturbation with $s > 2$. We have been able to obtain it since we impose a resonance condition (see Assumption 1.5 below), while in [7], we assumed that there is no resonance phenomena.

When the perturbations are localized, we are able to push the analysis further and obtain a nonlinear scattering.

This paper is organized as follows.

In section 1, we present our main results and the assumptions we need. Subsection 1.1 is devoted to the statement of the time decay estimates of the propagator associated with $H = D_m + V$ on the continuous subspace. One estimate is a kind of smoothness result, in the sense of Kato (see, e.g., [27]), and the other is a Strichartz-type result. We prove these estimates with the propagation and dispersive estimates proved in [7]. In subsection 1.2, we state the existence of small stationary states forming a manifold tangent to an eigenspace of H . The study of the dynamics around such states leads us to our main results; see subsection 1.3 and 1.4. In subsection 1.3, we split a neighborhood of a stationary state into different parts, each giving rise to stabilization or instability. In subsection 1.4, we state our scattering result.

To prove our theorems, we consider our nonlinear system as a small perturbation of a linear equation. More precisely, in subsection 2.2, we show that the spectral properties of the linearized operator around a stationary state, presented in section 2, permit us to obtain, as in the linear case, some properties of the dynamics around a stationary state. We obtain center, center-stable, and center-unstable manifolds. In section 3, we obtain, with our time decay estimates, a stabilization towards the PLS manifold for H^s perturbation with $s > 2$ in the center manifold. Section 4 deals with the dynamics outside the center manifold. Eventually in section 5, we conclude our study.

Our results are the analogue, in the Dirac case, of some results of Tsai and Yau [54], Soffer and Weinstein [43], Pillet and Wayne [38], and Gustafson, Nakanishi, and Tsai [22] about the semilinear Schrödinger equation. The point which we did not investigate is the long time behavior of perturbation in the unstable direction. In [51, 53, 54, 52, 50, 45, 46], the authors show that such perturbations relax towards the ground states. Their analyses are based on the asymptotic stability of ground states under a resonance condition. In our case this question is still open and should be investigated. But we believe that the techniques used in the previously cited works will allow us to obtain a similar result. Indeed the phenomenon that leads to this relaxation, namely nonlinear resonance, is not specific to the Schrödinger operator or nonnegative operators. Such a phenomenon comes from the interaction between discrete and continuous modes allowed by the nonlinearity. So a Fermi golden rule assumption should have the same consequences in both Schrödinger and Dirac equations.

We should also mention that four manifolds appear in this paper: the PLS manifold, the center manifold, the center-stable manifold, and the center-unstable manifold. The PLS manifold used to be called center manifold by previous papers on the same subject. The terminology we choose to adopt here comes from dynamical system theory and especially from the center manifold theorem. All these manifolds are invariant manifolds of our equation, at least locally in time. The PLS manifold is a set of stationary solutions, while the center manifold is a set of solutions converging

to the PLS manifold and containing the PLS manifold. The center manifold in turn is the intersection of the center-stable and center-unstable manifolds. The roles of the center-stable and center-unstable manifolds are exchanged when time is reversed.

This paper is devoted to the study of the asymptotic stability of some Dirac equations. So we are investigating convergence of orbits. As far as we know, the question of the orbital stability which requires the orbits to stay close and not necessarily to converge is still open for Dirac equations. In the Schrödinger case, orbital stability results (see, e.g., [15], [56, 57], or [42, 21]) give that any solution stays near the PLS manifold. Unfortunately, orbital stability criteria applied to Schrödinger equations use the fact that Schrödinger operators are bounded from below. Hence the question of orbital stability for Dirac standing waves cannot be solved by a straightforward application of the methods used in the Schrödinger case.

1. Assumptions and statements.

1.1. Time decay estimates. We generalize, to small nonlinear perturbations, stability results for linear systems. These results, as in [7], follow from linear decay estimates. Here we use smoothness-type and Strichartz-type estimates deduced from propagation and dispersive estimates of [7]. Hence, we work within the same assumptions for V and $D_m + V$.

Assumption 1.1. The potential $V : \mathbb{R}^3 \mapsto S_4(\mathbb{C})$ (*self-adjoint* 4×4 *matrices*) is a smooth function such that there exists $\rho > 5$ with

$$\forall \alpha \in \mathbb{N}^3, \exists C > 0, \quad \forall x \in \mathbb{R}^3, |\partial^\alpha V|(x) \leq \frac{C}{\langle x \rangle^{\rho+|\alpha|}}.$$

We notice that by the Kato–Rellich theorem, the operator

$$H := D_m + V$$

is essentially self-adjoint on $C_0^\infty(\mathbb{R}^3, \mathbb{C}^4)$ and self-adjoint on $H^1(\mathbb{R}^3, \mathbb{C}^4)$.

We also mention that Weyl's theorem gives that the essential spectrum of H is $(-\infty, -m] \cup [m, +\infty)$, and the work of Berthier and Georgescu [5, Theorems 6 and A] gives that there is no embedded eigenvalue. Hence the thresholds $\pm m$ are the only points of the continuous spectrum which can be associated with a wave of zero velocity. These waves perturb the spectral density and diminish the decay rate in the propagation and the dispersive estimates. We will work (as in [7]) within the following assumption.

Assumption 1.2. The operator H presents no resonance at thresholds and no eigenvalue at thresholds.

A resonance is a stationary solution in $H_{-\sigma}^{1/2} \setminus H^{1/2}$ for any $\sigma \in (1/2, \rho - 2)$, where H_σ^t is given by the following definition.

DEFINITION 1.1 (weighted Sobolev space). *The weighted Sobolev space is defined by*

$$H_\sigma^t(\mathbb{R}^3, \mathbb{C}^4) = \{f \in \mathcal{S}'(\mathbb{R}^3), \|\langle Q \rangle^\sigma \langle P \rangle^t f\|_2 < \infty\}$$

for $\sigma, t \in \mathbb{R}$. We endow it with the norm

$$\|f\|_{H_\sigma^t} = \|\langle Q \rangle^\sigma \langle P \rangle^t f\|_2.$$

If $t = 0$, we write L_σ^2 instead of H_σ^0 .

We have used the usual notation: $\langle u \rangle = \sqrt{1 + u^2}$, $P = -i\nabla$, and Q is the operator of multiplication by x in \mathbb{R}^3 .

Remark 1.1. Assumption 1.2 is generic in the following sense.

If V does not fulfill Assumption 1.2, then $(1 \pm \varepsilon)V$ does for nonzero small ε . As in [26, section 3] one can prove that resonances and eigenvalues at thresholds at $\pm m$ are in the kernel of $1 + (D_m \mp m)^{-1}V$ in $H_\sigma^{1/2}$ for any $\sigma \in (1/2, \rho - 2)$. Since $(D_m \mp m)^{-1}V$ is a compact operator in $H_\sigma^{1/2}$ (see proof of Lemma 1.1 below) for any $\sigma \in (1/2, \rho - 2)$, its spectrum is discrete, and hence if -1 is in its spectrum, then -1 is not in the spectrum of $\lambda(D_m \mp m)^{-1}V$ for $\lambda \neq 1$ but close to 1. This remark can be extended to any analytic operator valued functions $\lambda \rightarrow V(\lambda)$ satisfying Assumption 1.1.

Under the previous assumption one can prove the following lemma.

LEMMA 1.1. *The discrete spectrum of H is finite.*

Proof. To prove this lemma we show that there is no eigenvalue in a neighborhood of the thresholds $\pm m$. Since there is no other possible accumulation point for the discrete spectrum, this will be sufficient.

We will do the proof for m , since it is similar for $-m$. To do so, we show that $(H - z)^{-1}$ exists from $H_\sigma^{-1/2}$ to $H_{-\sigma}^{1/2}$ for any $\sigma \in (1/2, \rho - 2)$ and z close to m and $\Re z < m$. Since isolated eigenvectors belong to $H_\sigma^{-1/2}$ (see [23]), this will give the desired result.

We use the formula

$$H - z = (D_m - z) (1 + (D_m - z)^{-1}V)$$

for z close to m and $\Re z < m$. From [7, Proposition 2.4] (or [26, Lemma 2.1] and [48, section 1.E]), $z \mapsto (D_m - z)^{-1}$ is a continuous map from $\{z \in \mathbb{C}, \Im z > 0, \Re z > 0\}$ to $\mathcal{B}(H_\sigma^{-1/2}, H_{-\sigma}^{1/2})$ for any $\sigma \in (1/2, \rho - 2)$. The point is to prove the invertibility of $(1 + (D_m - z)^{-1}V)$ in $H_\sigma^{-1/2}$. Using continuity and Assumption 1.1, we see that it is enough to prove that $(1 + (D_m - m)^{-1}V)$ is invertible in $H_\sigma^{-1/2}$. This follows from Assumption 1.2 (see [7, Lemma 2.2 and Proof of Proposition 2.3]). \square

Now let

$$\mathbf{P}_c(H) = \mathbf{1}_{(-\infty, -m] \cup [m, +\infty)}(H)$$

be the projector associated with the continuous spectrum of H and let \mathcal{H}_c be its range. Using [7, Theorem 1.1], we obtain a limiting absorption principle which gives the H -smoothness of $\langle Q \rangle^{-1}$ in the sense of Kato.

THEOREM 1.1 (Kato smoothness estimates). *If Assumptions 1.1 and 1.2 hold, then for any $\sigma \geq 1$ and $s \in \mathbb{R}$ one has*

$$(1.1) \quad \left\| \langle Q \rangle^{-\sigma} e^{-itH} \mathbf{P}_c(H) \psi \right\|_{L_t^2(\mathbb{R}, H^s)} \leq C \|\psi\|_{H^s},$$

$$(1.2) \quad \left\| \int_{\mathbb{R}} e^{itH} \mathbf{P}_c(H) \langle Q \rangle^{-\sigma} F(t) dt \right\|_{H^s} \leq C \|F\|_{L_t^2(\mathbb{R}, H^s)},$$

$$(1.3) \quad \left\| \int_{s < t} \langle Q \rangle^{-\sigma} e^{-i(t-s)H} \mathbf{P}_c(H) \langle Q \rangle^{-\sigma} F(s) ds \right\|_{L_t^2(\mathbb{R}, H^s)} \leq C \|F\|_{L_t^2(\mathbb{R}, H^s)}.$$

Proof. We first prove (1.1). For $s = 0$, it is (see, e.g., [1, Proposition 7.11] or [40, Theorem XIII.25]) a consequence of the limiting absorption principle:

$$(1.4) \quad \sup_{\Im z \in (0,1)} \left\{ \left\| \langle Q \rangle^{-\sigma} (H - z)^{-1} P_c(H) \langle Q \rangle^{-\sigma} \right\|_2 \right\} < \infty,$$

which follows from [7, Theorem 1.1] (or Theorem 3.1 below) for $\sigma > 5/2$ using the fact that the Fourier transform in time of the propagator is the resolvent. Actually, the Fourier transform of

$$\langle Q \rangle^{-\sigma} e^{-it(H-i\varepsilon)} \mathbf{P}_c(H) \mathbf{1}_{\mathbb{R}_+^*}(t) \langle Q \rangle^{-\sigma} f$$

in time is

$$\langle Q \rangle^{-\sigma} (H - \lambda - i\varepsilon)^{-1} \mathbf{P}_c(H) \langle Q \rangle^{-\sigma} f$$

for $f \in L^2(\mathbb{R}^3, \mathbb{C}^4)$. Then we use the Born expansion

$$(H-z)^{-1} = (D_m-z)^{-1} - (D_m-z)^{-1}V(D_m-z)^{-1} + (D_m-z)^{-1}V(H-z)^{-1}V(D_m-z)^{-1},$$

the limiting absorption in [25, Theorem 2.1(i)] (the authors prove the identity (1.4) for $H = D_m$ when $\sigma = 1$), [7, Proposition 2.3], and the fact that $\|(1 - P_c(H))\|_{\mathcal{B}(H_\sigma^{1/2})}$ is bounded (since the discrete spectrum is finite and eigenvectors are exponentially decaying; see [23]) to obtain (1.4) for $\sigma = 1$. Hence we have concluded the proof for $s = 0$ and $\sigma \geq 1$. For $s \in 2\mathbb{Z}$ and $\sigma \geq 1$ it follows from the previous cases using the boundedness of $\langle H \rangle^s \langle D_m \rangle^{-s}$ and $\langle H \rangle^{-s} \langle D_m \rangle^s$ (which follow from the boundedness of V and its derivatives) and the boundedness of $\langle Q \rangle^{\mp\sigma} [\langle Q \rangle^{\pm\sigma}, H^s] \langle H \rangle^{-s}$ (which follow from multicommutator estimates; see [24, Appendix B]). The rest of the claim (1.1) follows by interpolation.

Estimates (1.1) and (1.2) are equivalent by duality.

To prove estimate (1.3) when $s = 0$ (the general case will follow in the same way as above), we notice that we have to prove that there exists $C > 0$ such that, for all $F, G \in L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))$, we have

$$\begin{aligned} & \left| \iint_{\mathbb{R}^2} \left\langle G(t), \langle Q \rangle^{-\sigma} e^{-i(t-s)H} \mathbf{P}_c(H) \mathbf{1}_{\mathbb{R}_+^*}(t-s) \langle Q \rangle^{-\sigma} F(s) \right\rangle ds dt \right| \\ & \leq C \|G\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))} \|F\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))}. \end{aligned}$$

We can suppose that F and G are smooth functions with compact support from $\mathbb{R} \times \mathbb{R}^3$ to \mathbb{C}^4 and we just need to prove that there exists $C > 0$ such that, for all $\varepsilon > 0$ and for all $F, G \in \mathcal{C}_0^\infty(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))$, we have

$$\begin{aligned} & \left| \iint_{\mathbb{R}^2} \left\langle G(t), \langle Q \rangle^{-\sigma} e^{-i(t-s)(H-i\varepsilon)} \mathbf{P}_c(H) \mathbf{1}_{\mathbb{R}_+^*}(t-s) \langle Q \rangle^{-\sigma} F(s) \right\rangle ds dt \right| \\ & \leq C \|G\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))} \|F\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))}. \end{aligned}$$

Then we take the limit as $\varepsilon \rightarrow 0$ and conclude using density arguments. Let us write $A_\varepsilon(t)$ for $\langle Q \rangle^{-\sigma} e^{-it(H-i\varepsilon)} \mathbf{P}_c(H) \mathbf{1}_{\mathbb{R}_+^*}(t) \langle Q \rangle^{-\sigma}$; then we have to prove

$$\left| \int_{\mathbb{R}} \langle G(t), (A_\varepsilon * F)(t) \rangle dt \right| \leq C \|G\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))} \|F\|_{L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))}.$$

Using Plancherel’s identity in $L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))$ and $A_\varepsilon * F \in L_t^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))$, we just need to prove

$$\left| \int_{\mathbb{R}} \langle \widehat{G}(\lambda), \widehat{A_\varepsilon * F}(\lambda) \rangle d\lambda \right| \leq C \left\| \widehat{G} \right\|_{L_\lambda^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))} \left\| \widehat{F} \right\|_{L_\lambda^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))}.$$

Since the Fourier transform in time of the propagator is the resolvent, F is smooth with compact support, and $\varepsilon > 0$, we obtain

$$\widehat{A_\varepsilon * F}(\lambda) = \langle Q \rangle^{-\sigma} (H - \lambda - i\varepsilon)^{-1} \mathbf{P}_c(H) \langle Q \rangle^{-\sigma} \widehat{F}(\lambda).$$

Hence we just have to prove

$$\begin{aligned} & \left| \int_{\mathbb{R}} \langle \widehat{G}(\lambda), \langle Q \rangle^{-\sigma} (H - \lambda - i\varepsilon)^{-1} \mathbf{P}_c(H) \langle Q \rangle^{-\sigma} \widehat{F}(\lambda) \rangle d\lambda \right| \\ & \leq C \left\| \widehat{G} \right\|_{L_\lambda^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))} \left\| \widehat{F} \right\|_{L_\lambda^2(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{C}^4))}. \end{aligned}$$

This in turn follows from the limiting absorption principle (1.4) just proved. \square

To state the next result, we need the following definition.

DEFINITION 1.2 (Besov space). *For $s \in \mathbb{R}$ and $1 \leq p, q \leq \infty$, the Besov space $B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^4)$ is the space of all $f \in \mathcal{S}'(\mathbb{R}^3, \mathbb{C}^4)$ (dual of the Schwartz space) such that*

$$\|f\|_{B_{p,q}^s} = \left(\sum_{j \in \mathbb{N}} 2^{jsq} \|\varphi_j * f\|_p^q \right)^{\frac{1}{q}} < +\infty$$

with $\widehat{\varphi} \in \mathcal{C}_0^\infty(\mathbb{R}^n \setminus \{0\})$ such that $\sum_{j \in \mathbb{Z}} \widehat{\varphi}(2^{-j}\xi) = 1$ for all $\xi \in \mathbb{R}^3 \setminus \{0\}$, $\widehat{\varphi}_j(\xi) = \widehat{\varphi}(2^{-j}\xi)$ for all $j \in \mathbb{N}^*$ and for all $\xi \in \mathbb{R}^3$, and $\widehat{\varphi}_0 = 1 - \sum_{j \in \mathbb{N}^*} \widehat{\varphi}_j$. It is endowed with the natural norm $f \in B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^4) \mapsto \|f\|_{B_{p,q}^s}$.

Using the dispersive estimates of [7, Theorem 1.2] and [29, Theorem 10.1], we obtain the following theorem.

THEOREM 1.2 (Strichartz estimates). *If Assumptions 1.1 and 1.2 hold, then for any $2 \leq p, q \leq \infty$, $\theta \in [0, 1]$, with $(1 - \frac{\theta}{2})(1 \pm \frac{\theta}{2}) = \frac{2}{p}$ and $(p, \theta) \neq (2, 0)$, and for any reals s, s' with $s' - s \geq \alpha(q)$, where $\alpha(q) = (1 + \frac{\theta}{2})(1 - \frac{\theta}{2})$, there exists a positive constant C such that*

$$(1.5) \quad \left\| e^{-itH} P_c(H) \psi \right\|_{L_t^p(\mathbb{R}, B_{q,2}^s(\mathbb{R}^3, \mathbb{C}^4))} \leq C \|\psi\|_{H^{s'}(\mathbb{R}^3, \mathbb{C}^4)},$$

$$(1.6) \quad \left\| \int e^{itH} P_c(H) F(t) dt \right\|_{H^s} \leq C \|F\|_{L_t^{p'}(\mathbb{R}, B_{q',2}^{s'}(\mathbb{R}^3, \mathbb{C}^4))},$$

$$(1.7) \quad \left\| \int_{s < t} e^{-i(t-s)H} P_c(H) F(s) ds \right\|_{L_t^p(\mathbb{R}, B_{q,2}^{-s}(\mathbb{R}^3, \mathbb{C}^4))} \leq C \|F\|_{L_t^{p'}(\mathbb{R}, B_{q',2}^{\tilde{s}}(\mathbb{R}^3, \mathbb{C}^4))},$$

for any $r \in [1, \infty]$, (\tilde{q}, \tilde{p}) chosen like (q, p) , and $s + \tilde{s} \geq \alpha(q) + \alpha(\tilde{q})$.

Proof. This is a consequence of [29, Theorem 10.1] applied to $U(t) = e^{-itH} P_c(H)$, using [7, Theorem 1.2] (or Theorem 3.2 below) and

$$B_{q,2}^{(1+\frac{\theta}{2})(1-\frac{\theta}{2})+s} \hookrightarrow (H^s, B_{1,2}^{1+\theta/2+s})_{2/((1\pm\theta/2)p),2}$$

continuously for $p \geq 2$ ($p \neq 2$ if $\theta = 0$) and $1/q = 1 - 1/((1 \pm \theta/2)p)$. For these embeddings, we refer to the proof of [4, Theorem 6.4.5] as well as the properties of the real interpolation (see [4] or [49]). More precisely, for $\theta = 0$ or 1 it is obvious. In the other cases, we work as in the proof of [4, Theorem 6.4.5(3)].

We use [4, Theorem 6.4.3] ($B_{p,2}^s$ is a retract of $l_2^s(L^p)$ for $s \in \mathbb{R}$ and $p, q \in [1, \infty]$) and [4, Theorem 5.6.2] (about the interpolation of $l_2^s(L^p)$ spaces) with [4, Theorem 5.2.1] (about the interpolation of L^p spaces). Then we conclude using the injection of L^p spaces into some Lorentz spaces [4, section 1.3 and Exercice 1.6.8].

In the case $\theta \neq 0$, the proof is actually simpler. We can prove it using the usual TT^* method and the Hölder inequality instead of the Hardy–Littlewood–Sobolev inequality. \square

1.2. The manifold of PLS. We study the nonlinear Dirac equation

$$(1.8) \quad \begin{cases} i\partial_t \psi = H\psi + \nabla F(\psi), \\ \psi(0, \cdot) = \psi_0 \end{cases}$$

with $\psi \in \mathcal{C}^1(I, H^1(\mathbb{R}^3, \mathbb{C}^4))$ for some real open interval I which contains 0 and $H = D_m + V$. The nonlinearity $F : \mathbb{C}^4 \mapsto \mathbb{R}$ is a differentiable map for the real structure of \mathbb{C}^4 , and hence the ∇ symbol has to be understood for the real structure of \mathbb{C}^4 : for the usual Hermitian product of \mathbb{C}^4 , one has

$$DF(v)h = \Re \langle \nabla F(v), h \rangle.$$

If F has a gauge invariance (see (0.2) or Assumption 1.4), this equation may have stationary solutions, i.e., solutions of the form $e^{-iEt}\phi_0$, where ϕ_0 satisfies the nonlinear stationary equation

$$E\phi_0 = H\phi_0 + \nabla F(\phi_0).$$

We notice that the Dirac operator D_m has an interesting invariance property due to its matrix structure. This invariance can be shared by some perturbed Dirac operators and gives a consequence of a theorem of Kramers; see [2, 37]. Indeed if we introduce K , the antilinear operator defined by

$$(1.9) \quad K \begin{pmatrix} \psi \\ \chi \end{pmatrix} = \begin{pmatrix} \sigma_2 \bar{\psi} \\ \sigma_2 \bar{\chi} \end{pmatrix} \quad \text{with } \sigma_2 = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix},$$

the operator D_m commutes with K . So if V also commutes with K , we obtain that the eigenspaces of H are always of even dimension. Here we work with the following assumption.

Assumption 1.3. The potential V commutes to K . The operator $H := D_m + V$ has only two double eigenvalues $\lambda_0 < \lambda_1$, with $\{\phi_0, K\phi_0\}$ and $\{\phi_1, K\phi_1\}$ as associated orthonormalized bases.

Remark 1.2. This assumption is fulfilled by $V = \lambda\phi\beta$, where ϕ is a smooth version of the characteristic function of the unit ball and λ is big enough to obtain a first set of double eigenvalues $\pm E$ with E close to m for H but not too much to avoid the appearance of a second set of eigenvalues. Computations have been sketched in [6], and similar ones have been done in detail for $V = \lambda\phi Id_{\mathbb{C}^4}$ in [47].

We also need the next assumption.

Assumption 1.4. The function $F : \mathbb{C}^4 \mapsto \mathbb{R}$ is in $C^\infty(\mathbb{R}^8, \mathbb{R})$ and satisfies $F(z) = O(|z|^4)$ as $z \rightarrow 0$. Moreover, it has the following invariance property:

$$\forall z \in \mathbb{C}^4, \forall u_1, u_2 \in \mathbb{C}^2 \text{ with } |u_1|^2 + |u_2|^2 = 1, F(u_1z + u_2Kz) = F(z).$$

Remark 1.3. The invariance property is actually equivalent to both the gauge invariance with respect to the semigroups generated by i and $e^{i\theta}K$ for all real θ .

We note that this assumption includes the cubic nonlinearity mentioned in the introduction. We were not able to go beyond the cubic order in the present work. This comes from the fact that here orbital stability does not hold or, equivalently, we have no a priori control in the L^2 norm of the perturbation. We bound it with $L_t^\infty H^s$ and $L_t^2 B_{\infty,2}^\beta$ ($s > \beta + 2 > 2$) using Duhamel’s formula and Theorem 1.2. Our method asks the nonlinearity to be at least cubic; see the proof of Lemma 3.8.

We also note that both the nonlinearity and the potential are smooth. This is not optimal. Concerning the nonlinearity, what we have in mind is a smooth cubic nonlinearity. For instance, a nonlinearity of class C^4 should be enough. It can be improved, but we are not looking for an optimal result in this direction. The advantage of choosing smooth potentials and nonlinearities here is that it allows stationary states to be smooth. It also avoids having to introduce a new parameter giving a bound for the regularity of our solutions.

We obtain the following proposition.

PROPOSITION 1.1 (PLS manifold). *If Assumptions 1.1–1.4 hold, then for any $\sigma \in \mathbb{R}^+$, there exist Ω , a neighborhood of 0 in \mathbb{C}^2 ; a smooth map*

$$h : \Omega \mapsto \{\phi_0, K\phi_0\}^\perp \cap H^2(\mathbb{R}^3, \mathbb{C}^4) \cap L_\sigma^2(\mathbb{R}^3, \mathbb{C}^4);$$

and a smooth map $E : \Omega \mapsto \mathbb{R}$ such that $S((u_1, u_2)) = u_1\phi_0 + u_2K\phi_0 + h((u_1, u_2))$ satisfies, for all $U \in \Omega$,

$$(1.10) \quad HS(U) + \nabla F(S(U)) = E(U)S(U),$$

with the following properties:

$$\begin{cases} h((u_1, u_2)) = \left(\frac{u_1}{|(u_1, u_2)|} Id_{\mathbb{C}^4} + \frac{u_2}{|(u_1, u_2)|} K \right) h(|(u_1, u_2)|, 0) \quad \forall U = (u_1, u_2) \in \Omega, \\ h(U) = O(|U|^2), \\ E(U) = E(|U|), \\ E(U) = \lambda_0 + O(|U|^2). \end{cases}$$

Proof. This result is adapted from [38, Proposition 2.2] after the reduction due to the invariance of the problem with respect to K . \square

Moreover, we have the following lemma.

LEMMA 1.2 (exponential decay). *For any $\beta \in \mathbb{N}^4$, $s \in \mathbb{R}^+$, and $p, q \in [1, \infty]$, there exist $\gamma > 0$, $\varepsilon > 0$, and $C > 0$ such that for all $U \in B_{\mathbb{C}^2}(0, \varepsilon)$ one has*

$$\|e^{\gamma(Q)} \partial_U^\beta S(U)\|_{B_{p,q}^s} \leq C \|S(U)\|_2,$$

where

$$\partial_{(u_1, u_2)}^\beta = \frac{\partial^{|\beta|}}{\partial^{\beta_1} \Re u_1 \partial^{\beta_2} \Im u_1 \partial^{\beta_3} \Re u_2 \partial^{\beta_4} \Im u_2}.$$

Proof. This is proved as in [7, Lemma 4.1], where we used ideas from [23]. \square

1.3. The unstable manifold and the stabilization. Each stationary solution previously introduced has, as in [7], a stable manifold. Under the following assumption, we can prove that a small perturbation of a stationary solution starting outside of this manifold leaves any neighborhood of this stationary solution. We work with the following assumption.

Assumption 1.5. The resonant condition

$$|\lambda_1 - \lambda_0| > \min\{|\lambda_0 + m|, |\lambda_0 - m|\}$$

holds. Moreover, we have the Fermi golden rule

$$(1.11) \quad \Gamma(\phi) > 0,$$

where, for any nonzero eigenvector ϕ associated with λ_0 , $\Gamma(\phi)$ is given by

$$\Gamma(\phi) = \lim_{\varepsilon \rightarrow 0, \varepsilon > 0} \left\langle d^2F(\phi)\phi_1, \Im((H - \lambda_0) + (\lambda_1 - \lambda_0) - i\varepsilon)^{-1} P_c(H)d^2F(\phi)\phi_1 \right\rangle.$$

In this assumption, the notation d^2F denotes the differential of ∇F with respect to the real structure of \mathbb{C}^4 . Let us introduce the linearized operator $JH(U)$ around a stationary state $S(U)$:

$$H(U) = H + d^2F(S(U)) - E(U).$$

We note that the operator $H(U)$ is not \mathbb{C} -linear but only \mathbb{R} -linear. Hence we work with the space $L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4)$ instead of $L^2(\mathbb{R}^3, \mathbb{C}^4)$ by writing

$$\begin{pmatrix} \Re\phi \\ \Im\phi \end{pmatrix}$$

instead of ϕ . The multiplication by $-i$ becomes the operator

$$J = \begin{pmatrix} 0 & I_{\mathbb{R}^4} \\ -I_{\mathbb{R}^4} & 0 \end{pmatrix}.$$

Now we mention some spectral properties of the *real operator* $JH(U)$ in $L^2(\mathbb{R}^3, \mathbb{C}^8)$ (the complexification of $L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4)$) which are needed to state and to understand our main theorem. These properties will be proved in subsection 2.1.

PROPOSITION 1.2 (spectrum of $JH(U)$). *The operator $JH(U)$ in $L^2(\mathbb{R}^3, \mathbb{C}^4 \times \mathbb{C}^4)$ has a four-dimensional algebraic kernel, and its spectrum is symmetric with respect to the imaginary and real axes.*

Outside the imaginary axis the spectrum is discrete. The sum of the algebraic eigenspaces associated with the right of the imaginary axis (i.e., associated with $\{z \in \mathbb{C}, \Re z > 0\}$) is stable under complex conjugation, and its dimension is 4. The same is true for the left of the imaginary axis.

The rest of the spectrum is the essential (or continuous) spectrum. We write $\mathcal{H}_c(U)$ for the space associated with the continuous spectrum. The space $J\mathcal{H}_c(U)$ is the orthogonal of the previous eigenspaces and the geometric kernel of $JH(U)$ and is invariant by the complex conjugation.

Proof. See subsection 2.1 below. \square

We will work on the real part of the previous spaces. Let $X_u(U) \subset L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4)$ be the real part of the sum of the spectral subspaces associated with $\{z \in \mathbb{C}, \Re z > 0\}$. We introduce a real basis $(\xi_i(U))_{i=1, \dots, 4}$ of $X_u(U)$.

We will also work in the real part of the sum of the spectral subspaces associated with $\{z \in \mathbb{C}, \Re z < 0\}$: $X_s(U) \subset L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4)$. We introduce a real basis $(\xi_i(U))_{i=5,\dots,8}$ of $X_s(U)$.

We define $\xi(U) = (\xi_i(U))_{i=1,\dots,8}$. We can state our main theorems which will be proved in sections 2–5.

THEOREM 1.3 (central manifold and asymptotic stability). *If Assumptions 1.1–1.5 hold, then, for $s > \beta + 2 > 2$ and $\sigma > 3/2$, there exist $\varepsilon > 0$, a continuous map $r : B_{\mathbb{C}^2}(0, \varepsilon) \mapsto \mathbb{R}$ with $r(U) = O(|U|^2)$, $C > 0$, \mathcal{V} a neighborhood of $(0, 0)$ in*

$$\mathcal{S} = \{(U, z); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H^s}(0, r(U))\}$$

endowed with the metric of $\mathbb{C}^2 \times H^s$, and a map $\Psi : \mathcal{V} \mapsto \mathbb{R}^8$, smooth on $\mathcal{V} \setminus (0, 0)$, satisfying for any nonzero $U \in B_{\mathbb{C}^2}(0, \varepsilon)$

$$\|\Psi(U, z)\| = O(\|z\|_{H^s}^2)$$

for all $z \in \mathcal{H}_c(U) \cap B_{H^s}(0, r(U))$ with $(U, z) \in \mathcal{V}$ such that the following is true.

For any initial condition of the form

$$\psi_0 = S(U_0) + z_0 + A \cdot \xi(U_0)$$

with $(U_0, z_0) \in \mathcal{V}$ and $A = \Psi(U_0, z_0)$, there exists a solution $\psi \in \cap_{k=0}^2 \mathcal{C}^k(\mathbb{R}, H^{s-k})$ of (1.8) with initial condition ψ_0 , and this solution is unique in $L^\infty((-T, T), H^s(\mathbb{R}^3, \mathbb{C}^4))$ for any $T > 0$.

Moreover, we have for all $t \in \mathbb{R}$

$$(1.12) \quad \psi(t) = e^{-i \int_0^t E(U(v)) \, dv} S(U(t)) + \varepsilon(t)$$

with $\|\dot{U}\|_{L^q(\mathbb{R})} \leq C\|z_0\|_{H^s}^2$ for all $q \in [1, \infty]$, $\lim_{t \rightarrow \pm\infty} U(t) = U_{\pm\infty}$, and

$$\max \left\{ \|\varepsilon\|_{L^\infty(\mathbb{R}^\pm, H^s)}, \|\varepsilon\|_{L^2(\mathbb{R}^\pm, H_{-\sigma}^s)}, \|\varepsilon\|_{L^2(\mathbb{R}^\pm, B_{\infty, 2}^\beta)} \right\} \leq C\|z_0\|_{H^s}.$$

Remark 1.4. We mention that to our knowledge the local well-posedness of the Cauchy problem associated with (1.8) has not been obtained for initial data in the energy space H^1 . Some results are available for slightly bigger spaces; see, for instance, the work of Escobedo and Vega [19] and Machihara, Nakamura, Nakanishi, and Ozawa [34, 33, 32].

In our case, we work with data in H^s with $s > 2$; this comes from the fact that we use inhomogeneous Strichartz estimates (1.7) with $q = \tilde{q} = \infty$ (see the proof of Lemma 3.8).

THEOREM 1.4 (center-stable and center-unstable manifolds). *With the same assumptions and notation as in Theorem 1.3, let \mathcal{CM} be the graph of $(U, z) \in \mathcal{V} \mapsto S(U) + z + \Psi(U, z) \cdot \xi(U)$. Then for the set*

$$\tilde{\mathcal{S}} = \{(U, z, p); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H^s}(0, r(U)), p \in B_{\mathbb{R}^4}(0, r(U))\}$$

endowed with the metric of $\mathbb{C}^2 \times H^s \times \mathbb{R}^4$, there exist $C > 0$, $\gamma > 0$, neighborhoods \mathcal{W}_\pm of $(0, 0, 0)$ in $\tilde{\mathcal{S}}$, and maps $\Phi_\pm : \mathcal{W}_\pm \mapsto \mathbb{R}^4$, smooth on $\mathcal{W}_\pm \setminus \{(0, 0, 0)\}$ with

$$\|\Phi_\pm(U, z, p)\| = O(\|z\|_{H^s}^2 + \|p\|^2)$$

for all $(U, z, p) \in \mathcal{W}_\pm$, such that for any initial condition of the form

$$\psi_0 = S(U_0) + z_0 + (p_+, p_-) \cdot \xi(U_0)$$

not in \mathcal{CM} , i.e.,

$$(p_+, p_-) \neq \Psi(U_0, z_0),$$

the following is true.

1. If $(U_0, z_0, p_+) \in \mathcal{W}_+$ and $p_- = \Phi_+(U_0, z_0, p_+)$ (resp., if $(U_0, z_0, p_-) \in \mathcal{W}_-$ and $p_+ = \Phi_-(U_0, z_0, p_-)$), then for any small neighborhood \mathcal{O} of $S(U_0)$ containing ψ_0 there exist $t_+ > 0$ (resp., $t_- > 0$) and a solution

$$\psi_+ \in \cap_{k=0}^2 \mathcal{C}^k([-t_+; +\infty), H^{s-k}),$$

respectively,

$$\psi_- \in \cap_{k=0}^2 \mathcal{C}^k((-\infty; t_-], H^{s-k}),$$

of (1.8) with initial condition ψ_0 , and in $L^\infty((-T', T), H^s(\mathbb{R}^3, \mathbb{C}^4))$ this solution is unique for any $T > 0$ (resp., $T \in (0, t_-)$) and any $T' \in (0, t_+)$ (resp., $T' < 0$).

Moreover, there exist $C > 0$, $\phi_\pm(t) \in \mathcal{CM}$, and $\rho_+(t) \in X_s(U_0)$ (resp., $\rho_-(t) \in X_u(U_0)$) for all $t > -t_+$ (resp., for all $t < t_-$) such that $\psi_\pm(t) = \phi_\pm(t) + \rho_\pm(t)$ with

$$\|\rho_\pm(t)\|_{H^s} \leq C \|\rho_\pm(0)\|_{H^s} e^{\mp \gamma t} \text{ as } t \rightarrow \pm\infty \quad \text{and} \quad \psi_\pm(\mp t_\pm) \notin \mathcal{O}.$$

We also have

$$\phi_\pm(t) = e^{-i \int_0^t E(U_\pm(v)) dv} S(U_\pm(t)) + \varepsilon_\pm(t) \quad \forall t > t_- \text{ (resp., } \forall t < -t_+)$$

with

$$\left\| \dot{U}_+ \right\|_{L^q((-t_+, +\infty))} \leq C (\|z_0\|_{H^s} + \|\rho_\pm(0)\|_{H^s})^2,$$

respectively,

$$\left\| \dot{U}_- \right\|_{L^q((-\infty, t_-))} \leq C (\|z_0\|_{H^s} + \|\rho_\pm(0)\|_{H^s})^2,$$

for all $q \in [1, \infty]$, $\lim_{t \rightarrow \pm\infty} U_\pm(t) = U_{\pm\infty}$, and

$$\begin{aligned} \max \left\{ \|\varepsilon_+\|_{L^\infty((-t_+, +\infty), H^s)}, \|\varepsilon_+\|_{L^2((-t_+, +\infty), H^s_\sigma)}, \|\varepsilon_+\|_{L^2((-t_+, +\infty), B^\beta_{\infty, 2})} \right\} \\ \leq C (\|z_0\|_{H^s} + \|\rho_\pm(0)\|_{H^s}), \end{aligned}$$

respectively,

$$\begin{aligned} \max \left\{ \|\varepsilon_-\|_{L^\infty((-\infty, t_-), H^s)}, \|\varepsilon_-\|_{L^2((-\infty, t_-), H^s_\sigma)}, \|\varepsilon_-\|_{L^2((-\infty, t_-), B^\beta_{\infty, 2})} \right\} \\ \leq C (\|z_0\|_{H^s} + \|\rho_\pm(0)\|_{H^s}). \end{aligned}$$

2. If $(U_0, z_0, p_+) \in \mathcal{W}_+$ and $(U_0, z_0, p_-) \in \mathcal{W}_-$ with $p_- \neq \Phi_+(U_0, z_0, p_+)$ and $p_+ \neq \Phi_-(U_0, z_0, p_-)$, then there exist $t_+(\psi_0) > 0$, $t_-(\psi_0) < 0$, and a unique solution ψ of (1.8) with initial condition ψ_0 such that, for any small neighborhood \mathcal{O} of $S(U_0)$ containing ψ_0 , $\phi \in \cap_{k=0}^2 \mathcal{C}^k([t_-; t_+], H^{s-k})$ with $\psi(t_+) \notin \mathcal{O}$ and $\psi(t_-) \notin \mathcal{O}$. This solution is unique in $L^\infty((T', T), H^s(\mathbb{R}^3, \mathbb{C}^4))$ for any $T \in (0, t_+)$ and any $T' \in (t_-, 0)$.

Remark 1.5. In Theorem 1.4, we note that when $U_0 = 0$, then $z_0 = 0$ and $p = 0$, and so the theorem does not say anything for this case. In fact, the charge conservation gives the orbital stability of 0. But we cannot extend the previous results to 0 since we can build a manifold of stationary states tangent to the eigenspace associated with λ_1 , similarly to Proposition 1.1.

One can also note that the previous result (and the following ones) are still valid when switching the roles of λ_0 and λ_1 as long as the resonance condition, Assumption 1.5, is fulfilled with λ_0 and λ_1 switched. If the resonance condition does not hold, the result in [7] can be applied if the strict nonresonance condition holds:

$$|\lambda_0 - \lambda_1| < \min\{|\lambda_1 + m|, |\lambda_1 - m|\},$$

and we obtain a stable manifold associated with λ_1 . As far as we know, the cases

$$|\lambda_0 - \lambda_1| = \min\{|\lambda_1 + m|, |\lambda_1 - m|\} \quad \text{and} \quad |\lambda_1 - \lambda_0| = \min\{|\lambda_0 + m|, |\lambda_0 - m|\}$$

are open.

1.4. The nonlinear scattering. If we choose a localized z_0 , we are able to further expand (1.12) as stated by the following theorems that are also proved in sections 2–5.

THEOREM 1.5. *With the assumptions and the notation of Theorem 1.3, for the set*

$$\mathcal{S}_\sigma = \{(U, z); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H_\sigma^s}(0, r(U))\}$$

endowed with the metric of $\mathbb{C}^2 \times H_\sigma^s$, there exists a neighborhood \mathcal{V}_σ of $(0, 0)$ in \mathcal{S}_σ such that the following is true. If $A = \Psi(U_0, z_0)$ with $(U_0, z_0) \in \mathcal{V}_\sigma$, there exist \mathcal{V}_σ^\pm open neighborhoods of $(0, 0)$ in \mathcal{S}_σ and $(V_{\pm\infty}; z_{\pm\infty}) \in \mathcal{V}_\sigma^\pm$, such that

$$\|V_{\pm\infty} - U_0\| \leq C \|z_0\|_{H_\sigma^s}^2, \quad \|z_{\pm\infty} - z_0\|_{H^s} \leq C \|z_0\|_{H_\sigma^s}^2,$$

and for all $t \in \mathbb{R}$

$$\psi(t) = e^{-itE(V_{\pm\infty})} S(V_\pm(t)) + e^{JtE(V_{\pm\infty})} e^{JtH(V_{\pm\infty})} z_{\pm\infty} + \varepsilon_\pm(t)$$

with

$$\left| \dot{V}_\pm(t) + i(E(V_\pm(t)) - E(V_{\pm\infty})) \right| \leq \frac{C}{\langle t \rangle^2} \|z_0\|_{H_\sigma^s}^2,$$

$$\|V_\pm(t) - V_{\pm\infty}\| \leq \frac{C}{\langle t \rangle} \|z_0\|_{H_\sigma^s},$$

$$\max \left\{ \|\varepsilon_\pm(t)\|_{H^s}, \|\varepsilon_\pm(t)\|_{H_{-\sigma}^s}, \|\varepsilon_\pm(t)\|_{B_{\infty,2}^\beta} \right\} \leq \frac{C}{\langle t \rangle^2} \|z_0\|_{H_\sigma^s}^2,$$

$$\text{and} \quad \left\| e^{-JtH(V_{\pm\infty})} e^{J \int_0^t (E(V_\pm(s)) - E(V_{\pm\infty})) ds} \varepsilon_\pm(t) \right\|_{H_{\frac{\sigma}{2}}^s} \leq \frac{C}{\langle t \rangle^{\frac{1}{2}}} \|z_0\|_{H_\sigma^s}^2$$

for all $t \in \mathbb{R}$.

Moreover, the maps

$$(U_0; z_0) \in \mathcal{V}_\sigma \mapsto (V_{\pm\infty}; z_{\pm\infty}) \in \mathcal{V}_\sigma^\pm$$

are bijective.

Remark 1.6. The fact that z_0 is localized gives us the convergence of

$$\int_0^t E(U(v)) \, dv - tE(U_{\pm\infty})$$

as $t \rightarrow \pm\infty$ and allows us to obtain an asymptotic profile for the dispersive part of the perturbed solution ϕ . Indeed, to obtain an asymptotic profile for z , we need an asymptotic profile for the phase and hence the existence of

$$(1.13) \quad \lim_{t \rightarrow \pm\infty} \int_0^t E(U(v)) \, dv - tE(U_{\pm\infty}) = \int_0^{\pm\infty} (E(U(v)) - E(U_{\pm\infty})) \, dv.$$

But in the present case, we control \dot{U} by the third power of some spatial norm of z . If z belongs just to some Lebesgue space in time, we are able to obtain just that

$$\lim_{t \rightarrow \pm\infty} (E(U(v)) - E(U_{\pm\infty})) = 0,$$

which gives nothing about the existence of (1.13). To obtain this existence, we need pointwise estimates in time for z (see subsection 3.2.4). This requires z_0 to be localized in space. Moreover, we believe that a result similar to [22, Theorem 1.9] showing slow pointwise decay for the perturbation in the nonlocalized case can be adapted here. Such phenomena prevent (1.13) from existing.

What we call the nonlinear scattering result is essentially the fact that the maps

$$(U_0; z_0) \in \mathcal{V}_\sigma \mapsto (U_{\pm\infty}; z_{\pm\infty}) \in \mathcal{V}_\sigma^\pm$$

are well defined and bijective.

Using wave operators for the pair $(JH(U), JD_m)$, we can obtain an expansion of the form $\psi(t) = e^{-i \int_0^t E(U(v)) \, dv} S(U_{\pm\infty}) + e^{-itD_m} z_\pm + \varepsilon_\pm(t)$ but will only have

$$\|z_{\pm\infty} - z_0\|_{H^s} \leq C \|z_0\|_{H^s}$$

and

$$\max \left\{ \|\varepsilon_\pm\|_{L^\infty(\mathbb{R}^\pm, H^s)}, \|\varepsilon_\pm\|_{L^2(\mathbb{R}^\pm, H_{-\sigma}^s)}, \|\varepsilon_\pm\|_{L^2(\mathbb{R}^\pm, B_{\infty,2}^\beta)} \right\} \leq C \|z_0\|_{H^s}.$$

Or, using wave operators for the pair $(JH(U), JH)$, we can obtain an expansion of the form $\psi(t) = e^{-i \int_0^t E(U(v)) \, dv} S(U_{\pm\infty}) + e^{-itH} z_\pm + \varepsilon_\pm(t)$ with

$$\|z_{\pm\infty} - z_0\|_{H^s} \leq C (|U_0| + \|z_0\|_{H^s}) \|z_0\|_{H^s}$$

and

$$\begin{aligned} \max \left\{ \|\varepsilon_\pm\|_{L^\infty(\mathbb{R}^\pm, H^s)}, \|\varepsilon_\pm\|_{L^2(\mathbb{R}^\pm, H_{-\sigma}^s)}, \|\varepsilon_\pm\|_{L^2(\mathbb{R}^\pm, B_{\infty,2}^\beta)} \right\} \\ \leq C (|U_0| + \|z_0\|_{H^s}) \|z_0\|_{H^s}. \end{aligned}$$

But in these cases, we cannot obtain a nice asymptotic completeness statement. We have (in the previous expansions)

$$\max \left\{ \sup_{t \in \mathbb{R}} \left(\| e^{JtH(U_{\pm\infty})} z_{\pm\infty} \|_{H^s} \right), \sup_{t \in \mathbb{R}} \left(\langle t \rangle^{3/2} \| e^{JtH(U_{\pm\infty})} z_{\pm\infty} \|_{H^s_{-\sigma}} \right), \right. \\ \left. \sup_{t \in \mathbb{R}} \left(\langle t \rangle^{3/2} \| e^{JtH(U_{\pm\infty})} z_{\pm\infty} \|_{B^{\beta}_{\infty,2}} \right) \right\} \leq C \| z_{\pm\infty} \|_{H^s}.$$

This follows from Lemmas 3.13 and 3.14.

Outside the center manifold, we can also have an expansion of the same type. But due to the presence of exponentially stable and unstable directions, one cannot expect a scattering result of the same type. Actually we cannot obtain the injectivity of the corresponding mappings. We have the following theorem.

THEOREM 1.6. *With the assumptions and notation of Theorem 1.4, for the sets*

$$\tilde{\mathcal{S}}_{\sigma} = \{ (U, z, p); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H^s_{\sigma}}(0, r(U)), p \in B_{\mathbb{R}^4}(0, r(U)) \}$$

endowed with the metric of $\mathbb{C}^2 \times H^s_{\sigma} \times \mathbb{R}^4$, there exist $C > 0, \gamma > 0$, and neighborhoods $\mathcal{W}_{\sigma}^{\pm}$ of $(0, 0, 0)$ in \mathcal{S}_{σ} such that the following is true.

If $\psi_0 \notin \mathcal{CM}$, $(U_0, z_0, p_+) \in \mathcal{W}_{\sigma}^+$ and $p_- = \Phi_+(U_0, z_0, p_+)$ (resp., $(U_0, z_0, p_-) \in \mathcal{W}_{\sigma}^-$ and $p_+ = \Phi_-(U_0, z_0, p_-)$), then there exist $C > 0, \phi_{\pm}(t) \in \mathcal{CM}$, and $\rho_{\pm}(t) \in X_s(U_0)$ for all $t > -t_+$ (resp., for all $t < t_-$) such that $\psi_{\pm}(t) = \phi_{\pm}(t) + \rho_{\pm}(t)$ with

$$\| \rho_{\pm}(t) \|_{H^s} \leq C \| \rho_{\pm}(0) \|_{H^s} e^{\mp \gamma t} \text{ as } t \rightarrow \pm\infty \quad \text{and} \quad \psi(\mp t_{\pm}) \notin \mathcal{O}.$$

Further, there exist $(V_{\pm\infty}; z_{\pm\infty}) \in \mathcal{S}$ such that

$$|V_{\pm\infty} - U_0| \leq C (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2,$$

$$\|z_{\pm\infty} - z_0\|_{H^s} \leq C (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2,$$

and for all $t > -t_+$ (resp., for all $t < t_-$)

$$\phi_{\pm}(t) = e^{-itE(V_{\pm\infty})} S(V_{\pm}(t)) + e^{JtE(V_{\pm\infty})} e^{JtH(V_{\pm\infty})} z_{\pm\infty} + \varepsilon_{\pm}(t)$$

with

$$\left| \dot{V}_{\pm}(t) + i(E(V_{\pm}(t)) - E(V_{\pm\infty})) \right| \leq \frac{C}{\langle t \rangle^2} (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2,$$

$$|V_{\pm}(t) - V_{\pm\infty}| \leq \frac{C}{\langle t \rangle} (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2,$$

$$\max \left\{ \|\varepsilon_{\pm}(t)\|_{H^s}, \|\varepsilon_{\pm}(t)\|_{H^s_{-\sigma}}, \|\varepsilon_{\pm}(t)\|_{B^{\beta}_{\infty,2}} \right\} \leq \frac{C}{\langle t \rangle^2} (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2,$$

$$\left\| e^{-JtH(V_{\pm\infty})} e^{J \int_0^t (E(V_{\pm}(s)) - E(V_{\pm\infty})) ds} \varepsilon_{\pm}(t) \right\|_{H^s_{\frac{3}{2}}} \leq \frac{C}{\langle t \rangle^{\frac{1}{2}}} (\|z_0\|_{H^s_{\sigma}} + \|\rho_{\pm}(0)\|_{H^s})^2$$

for all $t > -t_+$ (resp., for all $t < t_-$).

2. Linearized operator and exponentially stable and unstable manifolds. We study the dynamics associated with (1.8) around a stationary state. We will use spectral properties of the linearized operator around a stationary state.

2.1. The spectrum of the linearized operator. Here we study the spectrum of the linearized operator associated with (1.8) around a stationary state $S(U)$. Let us recall

$$H(U) = H + d^2F(S(U)) - E(U),$$

where d^2F is the differential of ∇F . The operator $H(U)$ is \mathbb{R} -linear but not \mathbb{C} -linear. Replacing $L^2(\mathbb{R}^3, \mathbb{C}^4)$ by $L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4)$ with the inner product obtained by taking the real part of the inner product of $L^2(\mathbb{R}^3, \mathbb{C}^4)$, we obtain a symmetric operator. We then complexify this real Hilbert space and obtain $L^2(\mathbb{R}^3, \mathbb{C}^4 \times \mathbb{C}^4)$ with its natural Hermitian product. This process transforms the operator $-i$ into

$$J = \begin{pmatrix} 0 & Id_{\mathbb{C}^4} \\ -Id_{\mathbb{C}^4} & 0 \end{pmatrix}.$$

For $\phi \in L^2(\mathbb{R}^3, \mathbb{R}^4 \times \mathbb{R}^4) \subset L^2(\mathbb{R}^3, \mathbb{C}^4 \times \mathbb{C}^4)$, we still write ϕ instead of

$$\begin{pmatrix} \Re\phi \\ \Im\phi \end{pmatrix}.$$

The extension of $H(U)$ to $L^2(\mathbb{R}^3, \mathbb{C}^4 \times \mathbb{C}^4)$ is also written $H(U)$ and is now a real operator. The extension of K (see (1.9)) is also written K .

The linearized operator associated with (1.8) around the stationary state $S(U)$ is given by $JH(U)$. We shall now study its spectrum.

Differentiating (1.10), we have that for $U = (u_1, u_2) \in \Omega$

$$\mathcal{H}_0(u_1, u_2) = \text{span} \left\{ \frac{\partial}{\partial \Re u_1} S(u_1, u_2), \frac{\partial}{\partial \Im u_1} S(u_1, u_2), \frac{\partial}{\partial \Re u_2} S(u_1, u_2), \frac{\partial}{\partial \Im u_2} S(u_1, u_2) \right\}$$

is invariant under the action of $JH(U)$. Differentiating the gauge invariance property for S , we notice that $JS(U) \in \mathcal{H}_0(U)$; differentiating the gauge invariance property for F , we also obtain

$$JH(U)JS(U) = 0;$$

and differentiating (1.10), we obtain for any $\beta \in \mathbb{N}^4$ with $|\beta| = 1$

$$JH(U)\partial_U^\beta S(U) = (\partial_U^\beta E)(U)JS(U).$$

The space $\mathcal{H}_0(U)$ is contained in the algebraic null space of $JH(U)$; in fact, it is exactly the algebraic null space as proved in what follows.

Now we state our results on the spectrum of $JH(U)$. The first deals with the excited states part. Using the function Γ introduced in Assumption 1.5, we have the following proposition.

PROPOSITION 2.1. *If Assumptions 1.1–1.5 hold, let*

$$\tilde{\Gamma}(|U|) = \inf_{V \in \mathbb{C}^2, |V|=|U|} \Gamma(S(V)).$$

For any sufficiently small U , we have

$$\tilde{\Gamma}(|U|) > 0.$$

There exist continuous maps $E_1^i : B_{\mathbb{C}^2}(0, \varepsilon) \mapsto \mathbb{R}$ smooth outside $U = 0$ with

$$\begin{cases} \Im E_1^i(U) = (\lambda_1 - E(U)) + O(|U|^4), \\ \Re E_1^i(U) = \Gamma^i(U) + O(|U|^6), \end{cases}$$

where

$$\Gamma^i(U) \geq \frac{1}{2} \tilde{\Gamma}(|U|)$$

for $i \in \{1, 2\}$ such that $E_1^i(U)$, $\overline{E_1^i(U)}$, $-E_1^i(U)$, and $-\overline{E_1^i(U)}$ are eigenvalues of $JH(U)$. These eigenvalues are simple if $E_1^1 \neq E_1^2$; otherwise the associated algebraic eigenspace is of dimension 2. We can define maps $\Phi_{\pm}^i(U) : B_{\mathbb{C}^2}(0, \varepsilon) \mapsto H^s$ such that

- $\{\Phi_+^1(U), \overline{\Phi_+^1(U)}\}$ is an orthonormal basis of the sum of the eigenspaces associated with $E_1^1(U)$ and $E_1^2(U)$,
- $\{\overline{\Phi_+^1(U)}, \Phi_+^2(U)\}$ is an orthonormal basis of the eigenspaces associated with $E_1^1(U)$ and $E_1^2(U)$,
- $\{\Phi_-^1(U), \overline{\Phi_-^1(U)}\}$ is an orthonormal basis of the sum of the eigenspaces associated with $-E_1^1(U)$ and $-E_1^2(U)$,
- $\{\overline{\Phi_-^1(U)}, \Phi_-^2(U)\}$ is an orthonormal basis of the sum of the eigenspaces associated with $-E_1^1(U)$ and $-E_1^2(U)$.

In the case of double eigenvalues $E_1^1(U) = E_1^2(U)$ one should consider algebraic eigenspaces instead of the sum of the eigenspaces. In any case, the associated projectors are continuous around $U = 0$ and smooth outside $U = 0$.

Moreover, for any $\beta \in \mathbb{N}^4$, $s \in \mathbb{R}^+$ and $p, q \in [1, \infty]$. There exist $\gamma > 0$, $\varepsilon > 0$, and $C > 0$ such that, for all $U \in B_{\mathbb{C}^2}(0, \varepsilon) \setminus \{0\}$ and for any $i \in \{1, 2\}$, one has

$$(2.1) \quad \|e^{\gamma \langle Q \rangle} \partial_U^\beta \Phi_{\pm}^i(U)\|_{B_{p,q}^s} \leq C \|\Phi_{\pm}^i(U)\|_2,$$

where

$$\partial_{(u_1, u_2)}^\beta = \frac{\partial^{|\beta|}}{\partial^{\beta_1} \Re u_1 \partial^{\beta_2} \Im u_1 \partial^{\beta_1} \Re u_2 \partial^{\beta_2} \Im u_2}.$$

Proof. Using Weyl's sequences, we prove that the essential spectrum of $JH(U)$, for small U , is the essential spectrum of $J(H - E(U))$, i.e., $i(\mathbb{R} \setminus (-c, c))$ with $c = \min\{m - E(U), m + E(U)\}$. So z with nonzero real part is in the spectrum of $JH(U)$ if and only if it is an isolated eigenvalue with finite multiplicity. The multiplicity of z here is the dimension of the algebraic kernel of $JH(U) - z$.

We now investigate properties of the eigenvalues. To do so, we use ideas from the proof of [54, Theorem 2.2]. The equation to solve for excited states is

$$(2.2) \quad (JH(U) - z)\phi = 0.$$

We consider solutions of the form $\phi = (v_1 S_1 + v_2 K S_1) + \eta$, where S_1 is the normalized eigenvector of JH ,

$$S_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} \phi_1 \\ -i\phi_1 \end{pmatrix},$$

with $(v_1, v_2) \in \mathbb{C}^2$ such that $|v_1|^2 + |v_2|^2 = 1$ and $\eta \in \{S_1, KS_1\}^\perp$. The orthogonal relation is taken, in fact, with respect to J (but since $JS_1 = iS_1$ and $JKS_1 = iKS_1$, we can take it in the usual way). For $z \in \mathbb{C} \setminus i\mathbb{R}$, we obtain the equation

$$(2.3) \quad \eta = (J(H - E(U)) - z)^{-1} P_1^\perp W(U) \{(v_1 S_1 + v_2 K S_1) + \eta\}$$

with P_1^\perp the orthogonal projector into $\{JS_1\}^\perp = \{S_1\}^\perp$ and $W(U) = JH(U) - J(H - E(U))$. We notice that $\{S_1\}^\perp$ is invariant under the action of $J(H - E(U))$. To solve this equation in η for a fixed U and z , we note that if

$$\Re z > 0 \quad \text{and} \quad |\Im z| \geq c,$$

where $c = \min\{m + E(U), m - E(U)\}$, the series

$$k(U, z)(v_1 S_1 + v_2 K S_1) = (J(H - E(U)) - z)^{-1} \times P_1^\perp \sum_{k \geq 0} \left(-W(U) (J(H - E(U)) - z)^{-1} P_1^\perp \right)^k W(U)(v_1 S_1 + v_2 K S_1)$$

is convergent in L^2 for sufficiently small $|U|$ and $|\Re z| = O(|U|^2)$ using the limiting absorption principle (1.4) and the bound of the resolvent $\|(H - z')^{-1}\| \leq |\Im z'|^{-1}$ in L^2 . Hence, we have a solution of (2.3).

Then we solve the equation in z . We obtain from (2.2) the equations

$$\langle (JH(U) - z)\phi, S_1 \rangle = 0 \quad \text{and} \quad \langle (JH(U) - z)\phi, K S_1 \rangle = 0$$

with $\phi = (v_1 S_1 + v_2 K S_1) + k(U, z)(v_1 S_1 + v_2 K S_1)$. Hence, z is an eigenvalue of the matrix $A(U, z)$ defined by

$$\begin{pmatrix} \langle JH(U)(S_1 + k(U, z)S_1), S_1 \rangle & \langle JH(U)(S_1 + k(U, z)S_1), K S_1 \rangle \\ \langle JH(U)(K S_1 + k(U, z)K S_1), S_1 \rangle & \langle JH(U)(K S_1 + k(U, z)K S_1), K S_1 \rangle \end{pmatrix}$$

in \mathbb{C}^2 . So we study the zeros of $\det(A(U, z) - z)$. We have that $d_\varepsilon : z \mapsto \det(A(U, z + \varepsilon) - z)$ is an analytic function for z in $\{z \in \mathbb{C}, \Re z > -\varepsilon\}$ for any $\varepsilon \geq 0$, $\det(A(0, z + \varepsilon) - z) = (z - i(\lambda_1 - \lambda_0))^2$ and $\det(A(U, z + \varepsilon) - z) - \det(A(0, z + \varepsilon) - z) \rightarrow 0$ as $U \rightarrow 0$ uniformly in z and ε in a bounded set. Since at $(U, \varepsilon) = (0, \varepsilon)$ ($\varepsilon > 0$) we have a double zero, for $(U, 0)$ close to $(0, 0)$, we have two zeros counted with multiplicity.

Hence if z is a zero of $d_0(U, \cdot)$ and hence an eigenvalue of $A(U, z)$, we have with an associated normalized eigenvector $(v_1, v_2) \in \mathbb{C}^2$ and $\psi = (v_1 S_1 + v_2 K S_1)$

$$\begin{aligned} z &= i(\lambda_1 - \lambda_0) + \langle W(U)\psi, \psi \rangle + \sum_{k \geq 0} \left\langle JH(U) (J(H - E(U)) - z)^{-1} \right. \\ &\quad \left. \times P_1^\perp \left(-W(U) (J(H - E(U)) - z)^{-1} P_1^\perp \right)^k W(U)\psi, \psi \right\rangle \\ &= i(\lambda_1 - \lambda_0) + \langle W(U)\psi, \psi \rangle \\ &\quad + \sum_{k \geq 0} \left\langle P_1^\perp \left(-W(U) (J(H - E(U)) - z)^{-1} P_1^\perp \right)^k W(U)\psi, \psi \right\rangle \end{aligned}$$

$$\begin{aligned}
 & + \sum_{k \geq 0} \left\langle (W(U) + z) (J(H - E(U)) - z)^{-1} \right. \\
 & \left. \times P_1^\perp \left(-W(U) (J(H - E(U)) - z)^{-1} P_1^\perp \right)^k W(U) \psi, \psi \right\rangle.
 \end{aligned}$$

Since $P_1^\perp(v_1 S_1 + v_2 K S_1) = 0$, we introduce the function

$$\begin{aligned}
 f(z) & = i(\lambda_1 - \lambda_0) + \left\langle W(U) \psi, \psi \right\rangle + \sum_{k \geq 0} \left\langle W(U) (J(H - E(U)) - z)^{-1} \right. \\
 & \left. \times P_1^\perp \left(-W(U) (J(H - E(U)) - z)^{-1} P_1^\perp \right)^k W(U) \psi, \psi \right\rangle.
 \end{aligned}$$

Since $J(v_1 S_1 + v_2 K S_1) = -i(v_1 S_1 + v_2 K S_1)$, we obtain that

$$\Re \langle W(U)(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle = 0.$$

Thus for $z \in \mathbb{C} \setminus i\mathbb{R}$, we have

$$\begin{aligned}
 & \Re f(z) \\
 & = \Re \left\langle W(U) (J(H - E(U)) - z)^{-1} P_1^\perp W(U)(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \right\rangle \\
 & \quad + O(|U|^6) \\
 & = \Im \langle d^2 F(S(U)) ((H - E(U)) + zJ)^{-1} \\
 & \quad \times P_1^\perp d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle + O(|U|^6).
 \end{aligned}$$

Then using (1.11) and

$$\begin{aligned}
 & ((H - E(U)) + zJ)^{-1} \\
 & = \frac{1}{2} \left(((H - E(U)) - iz)^{-1} (I_{C^2} + iJ) + ((H - E(U)) + iz)^{-1} (I_{C^2} - iJ) \right),
 \end{aligned}$$

we obtain

$$\begin{aligned}
 & \Im \langle d^2 F(S(U)) ((H - E(U)) + zJ)^{-1} P_1^\perp d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle \\
 & = \frac{1}{2} \left(\Im \langle d^2 F(S(U)) ((H - E(U)) - iz)^{-1} d^2 P_1^\perp F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle \right. \\
 & \left. + \Im \langle d^2 F(S(U)) ((H - E(U)) + iz)^{-1} P_1^\perp d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle \right) \\
 & \quad - \Im \langle d^2 F(S(U)) ((H - E(U))^2 + z^2)^{-1} z P_1^\perp d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle,
 \end{aligned}$$

and so, using regularity results of the resolvent of [20, Theorem 1.7], we obtain

$$\begin{aligned}
 & \Im \langle d^2 F(S(U)) ((H - E(U)) + (i(\lambda_1 - \lambda_0) + 0) J)^{-1} \\
 & \quad \times P_1^\perp d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle \\
 & = \frac{1}{2} \Im \langle d^2 F(S(U)) ((H - E(U)) + (\lambda_1 - \lambda_0) - i0)^{-1} \\
 & \quad \times P_c(H) d^2 F(S(U))(v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle.
 \end{aligned}$$

Using Assumption 1.5, the limiting absorption principle (1.4), and regularity results of [20, Theorem 1.7], we obtain

$$\Re f(z) = \frac{1}{2} \Im \langle d^2 F(S(U)) ((H - E(U)) + (\lambda_1 - \lambda_0) - i0)^{-1} \\ \times P_c(H) d^2 F(S(U)) (v_1 S_1 + v_2 K S_1), (v_1 S_1 + v_2 K S_1) \rangle + O(|U|^6)$$

for z in a ball of radius of order $|U|^2$ around $i(\lambda_1 - \lambda_0)$ and for small U . We also prove in the same way that

$$\Im f(z) = (\lambda_1 - \lambda_0) + O(|U|^4)$$

for z in a ball of radius of order $|U|^2$ around $i(\lambda_1 - \lambda_0)$ and for small U .

This proves that the zeros of $A(U, z)$ at the right of the imaginary axis are close to $i(\lambda_1 - \lambda_0)$. Hence we have obtained that, for small U , $JH(U)$ has at most two eigenvalues. Counted with algebraic multiplicity, we have that it is exactly two at the right side of the imaginary axis. Indeed, a complex number z is in the resolvent set of $JH(U)$ if and only if $JH(U) - z$ is invertible or if and only if $H(U) + Jz$ is invertible.

Let us consider for $\varepsilon > 0$

$$J_\varepsilon = \begin{pmatrix} \frac{\varepsilon}{2} & (1 + i\frac{\varepsilon}{2})Id_{\mathbb{C}^4} \\ (-1 - i\frac{\varepsilon}{2})Id_{\mathbb{C}^4} & \frac{\varepsilon}{2} \end{pmatrix}$$

and the set of z for which the operator

$$H(U) - J_\varepsilon z$$

is invertible. When $U = 0$, we have after multiplication by

$$P = \frac{1}{\sqrt{2}} \begin{pmatrix} i & 1 \\ 1 & i \end{pmatrix}$$

on the left and by P^{-1} on the right (that is to say, diagonalization of J_ε)

$$\begin{pmatrix} H - \lambda_0 + iz + i\varepsilon & 0 \\ 0 & H - \lambda_0 - iz \end{pmatrix},$$

which is invertible if

$$z \notin i \{ \mathbb{R} \setminus (-m - \lambda_0, m - \lambda_0) \cup \{0, \lambda_1 - \lambda_0\} \}$$

and

$$z \notin -i \{ \mathbb{R} \setminus (-m - \lambda_0, m - \lambda_0) \cup \{0, \lambda_1 - \lambda_0\} \} - \varepsilon.$$

Since we have closed operators, applying Theorems 3.16 and 5.33 of [28, Chapter IV] to $J_\varepsilon^{-1}H(U)$, we obtain that close to $i(\lambda_1 - \lambda_0)$ the operator $J_\varepsilon^{-1}H(U)$ has two simple eigenvalues or one double (algebraic) eigenvalue when U is small. Now consider ε as a parameter and U fixed, and applying again Theorems 3.16 and 5.33 of [28, Chapter IV], we obtain that, counted with multiplicity, $JH(U)$ has two eigenvalues since the discrete spectrum of $JH(U) = J_0^{-1}H(U)$ close to $i(\lambda_1 - \lambda_0)$ at the right

of the imaginary axis is outside the imaginary axis. Theorems 3.16 and 5.33 of [28, Chapter IV] give as well the continuity with respect to U of the associated projector.

The continuity of eigenvalues follows from the continuity of the zeros of d_ε introduced earlier.

Letting $P_1(U)$ be the spectral projector associated with $\{E_1^i(U), i \in \{1, 2\}\}$, the eigenvalue of $JH(U)$, we have the formula

$$P_1(U) = -\frac{1}{2i\pi} \int_\Gamma (JH(U) - z)^{-1} dz,$$

where Γ is a positively oriented circle at the right of the imaginary axis around $E_1^1(U')$ and $E_1^2(U')$ for some fixed U' and U sufficiently close to U' ; see [28, Formula I.1.16]. Since the resolvent is smooth outside $U = 0$, we have that $P_1(U)$ is smooth outside $U = 0$.

Using complex conjugation, we obtain the corresponding result for the spectrum at the right of the imaginary axis in a neighborhood of $i(\lambda_1 - \lambda_0)$.

Using Weyl's sequences, we prove that the essential spectrum of $(JH(U))^* = -H(U)J$, for small U , is the essential spectrum of $-(H - E(U))J = -J(H - E(U))$. So z with nonzero real part is in the spectrum of $(JH(U))^*$ if and only if it is an isolated eigenvalue; see [28, Theorem IV.5.33]. Then to obtain $-E_1^i(U)$ and $-\overline{E_1^i(U)}$, we notice that $E_1^i(U)$ and $\overline{E_1^i(U)}$ are eigenvalues of $(JH(U))^*$. Using the symmetry $J(JH(U)) = -(JH(U))^*J$, we show that any algebraic eigenvector ϕ of $(JH(U))^*$ associated with λ , $J\phi$ is an algebraic eigenvector of $JH(U)$ associated with $-\lambda$. Hence, repeating the previous proof for $(JH(U))^*$, we also obtain the multiplicity two. Except for the symmetry of the eigenvalues, this can also be obtained by adapting the previous proof to the left part of the complex plane.

The exponential decay works as in Lemma 1.2. □

Remark 2.1. If $F(z)$ is homogeneous of order p , then there exist $\varepsilon, \Gamma_1, \Gamma_2 > 0$ such that for all $U \in B_{\mathbb{C}^2}(0, \varepsilon)$

$$|U|^{p-2} \Gamma_1 \leq \tilde{\Gamma}(|U|) \leq |U|^{p-2} \Gamma_2.$$

We just write $S((u_1, u_2)) = u_1\phi_0 + u_2K\phi_0 + h((u_1, u_2))$, expand $\Gamma(U)$, and use Assumption 1.5 with the regularity results of the resolvent from [20, Theorem 1.7].

The following proposition deals with the essential spectrum of our linearized operator.

PROPOSITION 2.2. *If Assumptions 1.1–1.5 hold, for any sufficiently small nonzero $U \in \mathbb{C}^2$, let*

$$\mathcal{H}_1(U) = \text{span} \left\{ \Phi_+^1(U), \Phi_+^2(U), \overline{\Phi_+^1(U)}, \overline{\Phi_+^2(U)}, \Phi_-^1(U), \Phi_-^2(U), \overline{\Phi_-^1(U)}, \overline{\Phi_-^2(U)} \right\}.$$

The orthogonal space of $\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)$,

$$\mathcal{H}_c(U) = \{J\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\}^\perp,$$

is invariant under the action of $JH(U)$.

We also have that, for $\mathbf{P}_c(U)$, the orthogonal projector onto $J\mathcal{H}_c(U)$ is a bounded operator from $H_\sigma^s(\mathbb{R}^3, \mathbb{C}^8)$ or $B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^8)$ to itself for any reals $s, \sigma \in \mathbb{R}$ and $p, q \in [1, \infty]$, and for $U' \in B_{\mathbb{C}^2}(U, \varepsilon)$, with sufficiently small $\varepsilon > 0$, that

$$\mathbf{P}_c(U)|_{\mathcal{H}_c(U')} : \mathcal{H}_c(U') \mapsto \mathcal{H}_c(U)$$

is an isomorphism and its inverse $R(U', U)$ can be extended to a bounded operator from $H_\sigma^s(\mathbb{R}^3, \mathbb{C}^8)$ or $B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^8)$ to itself for any reals $s, \sigma \in \mathbb{R}$ and $p, q \in [1, \infty]$.

Moreover, there exists $C > 0$ such that we have

$$(2.4) \quad \begin{aligned} \|\psi\|_X &\leq C \|P_c \psi\|_X \\ \forall \psi \in \mathcal{H}_c(U) \cap X \text{ with } X &= H_\sigma^s(\mathbb{R}^3, \mathbb{C}^8) \text{ or } B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^8), \\ \forall s, \sigma \in \mathbb{R}, \quad \forall p, q \in [1, \infty]; \end{aligned}$$

$$(2.5) \quad \begin{aligned} \int_{\mathbb{R}} \|\langle Q \rangle^{-\sigma} e^{sJH(U)} \mathbf{P}_c(U) \psi\|_2^2 ds &\leq C \|\psi\|_2^2 \\ \forall \psi \in L^2, \quad \forall \sigma \geq 1; \end{aligned}$$

$$(2.6) \quad \begin{aligned} \left\| e^{tJH(U)} \mathbf{P}_c(U) \psi \right\|_2 &\leq C \|\psi\|_2 \\ \forall t \in \mathbb{R}, \quad \forall \psi \in L^2, \end{aligned}$$

and $\mathcal{H}_c(U)$ contains no eigenvector.

Remark 2.2. We use the same notation for $\mathcal{H}_c(U)$ and its real part (see definition below) that appears in our main theorems. We just note that $\mathcal{H}_c(U)$ appears when we discuss spectral properties in our proof. Then when we talk about dynamical properties, we deal with its real part. We remind the reader that the real part of $\mathcal{H}_c(U)$ is left invariant by $JH(U)$.

Proof. We prove that there is no other eigenvector by proving that smoothness estimate (2.5) takes place over

$$\mathcal{H}_c(U) = \{J\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\}^\perp.$$

First we prove that

$$\mathbf{P}_c((U))|_{\mathcal{H}_c(U')} : \mathcal{H}_c(U') \mapsto \mathcal{H}_c(U)$$

is an isomorphism. To prove it, we exhibit an inverse $R(U', U)$ which is the projector onto $\mathcal{H}_c(U')$ associated with the decomposition $\mathcal{H}_0(U) \oplus \mathcal{H}_1(U) \oplus \mathcal{H}_c(U')$ of $L^2(\mathbb{R}^3, \mathbb{C}^8)$. Indeed, we have $\{\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\} \cap \mathcal{H}_c(U') = \{0\}$ when U' and U are close to one another and $\text{codim} \mathcal{H}_c(U') = \dim\{\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\}$. We have a decomposition of $L^2(\mathbb{R}^3, \mathbb{C}^8)$ into closed subspaces; hence the associated projectors are continuous. So $R(U', U)$ should be of the form

$$R(U', U) = Id + \sum_i |J\xi_i(U)\rangle \langle \alpha_i(U', U)|,$$

where $\xi_i(U)$ is a basis of the eigenspaces of $JH(U)$ and $(\alpha_i(U', U))_i$ solve the equations

$$J\xi_j(U') + \sum_i \langle J\xi_i(U), J\xi_j(U') \rangle \alpha_i(U', U) = 0.$$

Such an α exists because the matrix $(\langle J\xi_i(U), J\xi_j(U') \rangle)_{i,j}$ is invertible. Indeed, if it is not invertible, there exists $\phi \in \mathcal{H}_1(U') = \text{span}\{\xi_j(U')\}$ orthogonal to $\mathcal{H}_1(U) =$

$\text{span}\{\xi_j(U)\}$, and hence $\phi \in \mathcal{H}_c(U)$ with $P_c(U')\phi = 0$. But for all $\psi \in \mathcal{H}_c(U)$, we have

$$\|\psi\| \leq \|P_c(U')\psi\| + \|(1 - P_c(U'))\psi\| \leq \|P_c(U')\psi\| + C|U - U'| \|\psi\|,$$

and hence for ϕ this gives a contradiction.

The boundedness of R in $\mathcal{B}(H_\sigma^s(\mathbb{R}^3, \mathbb{C}^8))$ or $\mathcal{B}(B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^8))$ follows from the exponential decay of eigenvectors and their derivatives; see Proposition 1.1, Lemma 1.2, and Proposition 2.1.

Let us now consider the orthogonal projector P_c associated with the continuous subspace of JH . Since the eigenvectors of JH are exponentially decaying, we can extend P_c to obtain an operator of $L_{\pm\sigma}^2$ into itself. The same is true for $P_c(U)$, and hence we can consider the extension of $\mathcal{H}_c(U)$ to $L_{\pm\sigma}^2$. We still call it $\mathcal{H}_c(U)$.

For all $\psi \in \mathcal{H}_c(U)$,

$$\|\psi\|_{L_{-\sigma}^2} \leq \|P_c\psi\|_{L_{-\sigma}^2} + \|(1 - P_c)\psi\|_{L_{-\sigma}^2}.$$

Since $1 - P_c$ is the projector into the eigenspaces of H and ψ is orthogonal to the eigenvectors of $JH(U)$, we obtain that

$$\|(1 - P_c)\psi\|_{L_{-\sigma}^2} = \|(1 - P_c(U))\psi\|_{L_{-\sigma}^2} + O(|U|)\|\psi\|_{L_{-\sigma}^2}.$$

Since the components of eigenvectors with respect to $\mathcal{H}_c(U)$ are small with respect to U (see the proof of Lemma 2.1), the projection of ψ in the eigenspaces of H is small. Moreover, for a sufficiently small nonzero U , we obtain estimate (2.4) for $X = L_{-\sigma}^2$ with $\sigma > 0$.

The rest of estimate (2.4) follows in the same way using the exponential decay of eigenvectors (estimate (2.1) and Lemma 1.2).

We infer that

$$\begin{aligned} & \| \langle Q \rangle^{-\sigma} e^{tJH(U)} \mathbf{P}_c(U)\psi \| \\ & \leq C \| \langle Q \rangle^{-\sigma} \mathbf{P}_c e^{tJH(U)} \mathbf{P}_c(U)\psi \| \\ & \leq C \| \langle Q \rangle^{-\sigma} \mathbf{P}_c e^{-it(H-E(U))} \mathbf{P}_c(U)\psi \| \\ & \quad + C \left\| \langle Q \rangle^{-\sigma} \int_0^t \mathbf{P}_c e^{-i(t-s)(H-E(U))} D\nabla F(S(U)) e^{sJH(U)} \mathbf{P}_c(U)\psi ds \right\|. \end{aligned}$$

Using estimates (1.1) and (1.3) of Theorem 1.1, we obtain estimate (2.5) for sufficiently small U :

$$\int_{\mathbb{R}} \| \langle Q \rangle^{-\sigma} e^{-sJH(U)} \mathbf{P}_c(U)\psi \|^2 ds \leq C \|\psi\|^2.$$

Hence there is no eigenvector in the range of $P_c(U)$. Using the inequalities (2.5); the conservation law for H ; and Duhamel’s formula,

$$e^{JtH(U)} = e^{-it(H-E(U))} + \int_0^t e^{-i(t-s)(H-E(U))} Jd^2\nabla F(S(U)) e^{JsH(U)} ds,$$

we prove estimate (2.6). □

Since $\mathcal{H}_c(U)$ is closed and $\text{codim}\mathcal{H}_c(U) = \dim\{\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\}$ and $\mathcal{H}_c(U) \cap \{\mathcal{H}_0(U) \oplus \mathcal{H}_1(U)\} = \{0\}$, we obtain $\mathcal{H}_0(U) \oplus \mathcal{H}_1(U) \oplus \mathcal{H}_c(U) = L^2(\mathbb{R}^3, \mathbb{C}^8)$ and the following proposition.

PROPOSITION 2.3. *Suppose that Assumptions 1.1–1.5 hold. Then the space $\mathcal{H}_0(U)$ is the geometric null space of $JH(U)$.*

2.2. Stable, unstable, and center manifolds. We can now obtain results similar to those of Bates and Jones [3]. We note that we will not prove that the Cauchy problem (1.8) is locally well-posed for an initial condition outside some manifolds (built below). However, this can be proved with the methods we present here or by generalizing to our case the results of Escobedo and Vega [19].

We have that $JH(U)$ as an operator in $L^2(\mathbb{R}^3, \mathbb{R}^8)$ is a closed densely defined operator that generates a continuous semigroup on $L^2(\mathbb{R}^3, \mathbb{R}^8)$. The spectrum of $JH(U)$ in $L^2(\mathbb{R}^3, \mathbb{R}^8)$ is the same as that of $JH(U)$ in $L^2(\mathbb{R}^3, \mathbb{C}^8)$, and so it splits into three parts,

$$\sigma_s(U) = \{\lambda \in \sigma(JH(u)), \Re\lambda < 0\} = \left\{-E_1^1(U), -E_1^2(U), -\overline{E_1^1(U)}, -\overline{E_1^2(U)}\right\},$$

$$\sigma_c(U) = \{\lambda \in \sigma(JH(u)), \Re\lambda = 0\} = \{0\} \cup i\{\mathbb{R} \setminus (-c, c)\},$$

$$\text{where } c = \min\{m + E(U), m - E(U)\},$$

$$\sigma_u(U) = \{\lambda \in \sigma(JH(u)), \Re\lambda > 0\} = \left\{E_1^1(U), E_1^2(U), \overline{E_1^1(U)}, \overline{E_1^2(U)}\right\},$$

each of which is associated with a spectral real subspace, respectively,

$$X_s(U) = \text{span}_{\mathbb{R}} \left\{\Re\Phi_-^1(U), \Im\Phi_-^1(U), \Re\Phi_-^2(U), \Im\Phi_-^2(U)\right\},$$

$$X_u(U) = \text{span}_{\mathbb{R}} \left\{\Re\Phi_+^1(U), \Im\Phi_+^1(U), \Re\Phi_+^2(U), \Im\Phi_+^2(U)\right\},$$

$$X_c(U) = \Re\mathcal{H}_0(U) \oplus \Re\mathcal{H}_c(U),$$

where we used the notation $\Re\Psi = (1/2)(\Psi + \overline{\Psi})$, $\Im\Psi = -(i/2)(\Psi - \overline{\Psi})$, and $\Re X = \{\Re\Psi, \Psi \in X\}$, the real part of the space X . The spaces $X_s(U)$ and $X_u(U)$ are finite dimensional. Let us write $\pi^c(U)$, $\pi^s(U)$, and $\pi^u(U)$ for the projectors associated with the decomposition $X_c(U) \oplus X_s(U) \oplus X_u(U)$. Since the eigenvectors also belong to L_σ^2 for any $\sigma \in \mathbb{R}$, the projectors $P_c(U)$ and $\pi^c(U)$ can be defined in L_σ^2 for any real σ . We can extend, in this way, the spaces $\mathcal{H}_c(U)$ and $X_c(U)$ to L_σ^2 for any $\sigma \in \mathbb{R}$. We have the following lemma.

LEMMA 2.1. *If Assumptions 1.1–1.5 hold, then, for any $\sigma \in \mathbb{R}$, there exist positive reals r, C_1, C_2 such that, for all $t \in \mathbb{R}$, we have*

$$(2.7) \quad C_1 \frac{1}{f_+(t)} \leq \left\| e^{tJH(U)} \pi^s(U) \right\|_{\mathcal{B}(L_\sigma^2)} \leq C_2 \frac{1}{f_-(t)},$$

$$(2.8) \quad C_1 f_-(t) \leq \left\| e^{tJH(U)} \pi^u(U) \right\|_{\mathcal{B}(L_\sigma^2)} \leq C_2 f_+(t),$$

$$(2.9) \quad \left\| e^{tJH(U)} \pi^c(U) \right\|_{\mathcal{B}(L_\sigma^2)} \leq C_2 \langle t \rangle^r,$$

where

$$f_\pm(t) = \begin{cases} e^{\gamma_\pm t} & \text{if } t \geq 0, \\ e^{\gamma_\mp t} & \text{if } t \leq 0, \end{cases}$$

with

$$\gamma_- < \Gamma_-(U) = \min\{\Re E_1^1(U), \Re E_1^2(U)\},$$

$$\gamma_+ > \Gamma_+(U) = \max\{\Re E_1^1(U), \Re E_1^2(U)\}.$$

Proof. The statements for the spaces $X^s(U)$ and $X^u(U)$ follow from (2.1).

The statement about $X^c(U)$ is a little more complicated. We note that we are not looking for an optimal r .

First, the result for $e^{-it(D_m+V)}$ in L^2_σ with $\sigma \in 2\mathbb{N}$ follows from [48, Theorem 8.5] (see also Proposition 3.1 below), which is based on the charge conservation. The case $\sigma \in \mathbb{R}$ follows by duality and interpolation.

Then for $e^{tJH(U)}\pi^c(U)$, we use Duhamel’s formula,

$$e^{JtH(U)}\pi^c(U) = e^{-it(H-E(U))}\pi^c(U) + \int_0^t e^{-i(t-s)(H-E(U))}Jd^2\nabla F(S(U))e^{JsH(U)}\pi^c(U) ds;$$

then the assertion for $e^{tJH(U)}\pi^c(U)$ follows from the assertion for $e^{-it(D_m+V)}$, the charge conservation of $e^{tJH(U)}P_c(U)$ (see (2.6)), the fact that $e^{tJH(U)}S(U) = S(U)$, $e^{tJH(U)}\partial_U^\beta S(U) = \partial_U^\beta S(U) + t\partial_U^\beta E(U)S(U)$, and Lemma 1.2. \square

By now we do not restrict our study to the space $L^2(\mathbb{R}^3, \mathbb{R}^8)$; we extend it to $L^2_\sigma(\mathbb{R}^3, \mathbb{R}^8)$ for any $\sigma \in \mathbb{R}$, but we still write $\mathcal{H}_c(U)$ and $X_c(U)$ for the extensions of these spaces to $L^2_\sigma(\mathbb{R}^3, \mathbb{R}^8)$ for any $\sigma \in \mathbb{R}$.

We now study the behavior of the solutions in L^2_σ of (1.8) centered around $S(U)$:

$$(2.10) \quad \partial_t\phi = JH(U)\phi + JN(U, \phi),$$

where $H(U) = H + d^2F(S(U)) - E(U)$, $N(U, \phi) = \nabla F(S(U) + \phi) - \nabla F(S(U)) - d^2F(S(U))\phi$, and d^2F is the differential of ∇F .

In this subsection, we study a modified equation which coincides with (2.10) as long as the solution stays in a neighborhood of a small $S(U)$:

$$(2.11) \quad \partial_t\phi = JH(U)\phi + JN_\varepsilon(U, \phi),$$

where $N_\varepsilon(U, \eta) = \rho(\varepsilon^{-1}\eta)N(U, \eta)$ and ρ is a smoothed version of the characteristic function of the unit ball of \mathbb{R}^8 .

We state the following proposition.

PROPOSITION 2.4 (center-stable manifold). *If Assumptions 1.1–1.5 hold, then for any sufficiently small nonzero U , there exists around $S(U)$ a unique invariant smooth center-stable manifold $W^{cs}(U)$ for (2.11) built as a graph with value in $X_u(U)$ and tangent to $S(U) + X_c(U) \oplus X_s(U)$ at $S(U)$.*

Any solution $\phi \in L^2_\sigma$ of (2.11) initially in the neighborhood of $S(U)$ tends as $t \rightarrow -\infty$ to $W^{cs}(U)$ with

$$dist_{L^2_\sigma}(\phi(t), W^{cs}(U)) = O(e^{\gamma t}) \text{ as } t \rightarrow -\infty$$

for any $\gamma \in (0, \Gamma_-(U))$ and any $s, \sigma \in \mathbb{R}$, and for any sufficiently small neighborhood V of $S(U)$ any solution in V not in $W^{cs}(U)$ leaves V in finite positive time.

Remark 2.3. For any $s \in \mathbb{R}^+$, due to the exponential decay of eigenvectors, if $\phi \notin H^s_\sigma$, there exists $\psi \in W^{cs}(U)$ such that $\phi - \psi \in H^s_\sigma$, and we have

$$dist_{H^s_\sigma}(\phi(t), W^{cs}(U)) = O(e^{\gamma t}) \text{ as } t \rightarrow -\infty,$$

as shown in the following proof.

If we consider only small solutions, we obtain a locally invariant manifold for (2.10); that is to say, for any initial condition in the manifold there exists a corresponding solution of (2.10) which stays in this manifold in a small interval of time

around 0. We notice that in the following proofs the size of this invariant manifold, which is given by ε , is a function of U and this function is $O(\Gamma_-(U))$. From now on, we call this function r .

Proof. Our proof is an adaptation of that of Bressan [8], and we refer to it for more details. We make the proof only for the case $\sigma = 0$; the proof in the general case is similar.

First we prove that there is a global solution of (2.11) which does not grow much as $t \rightarrow +\infty$. We look for the solution as a fixed point:

$$y(t) = \mathcal{G}_\varepsilon(y_0, y)(t)$$

for any $y_0 \in X_s(U) \oplus X_c(U)$, where for small positive ε

$$\begin{aligned} \mathcal{G}_\varepsilon(y_0, \eta)(t) &= e^{tJH(U)}y_0 + \int_0^t e^{(t-s)JH(U)}\pi^c(U)JN_\varepsilon(U, \eta(s)) \, ds \\ &+ \int_0^t e^{(t-s)JH(U)}\pi^s(U)JN_\varepsilon(U, \eta(s)) \, ds - \int_t^{+\infty} e^{(t-s)JH(U)}\pi^u(U)JN_\varepsilon(U, \eta(s)) \, ds, \end{aligned}$$

with $\pi^*(U)$ the projector into $X^*(U)$ with respect to the decomposition $\oplus_{* \in \{c, s, u\}} X^*(U)$.

Let us introduce, for any Γ smaller than $\Gamma_-(U)$, the space

$$Y_\Gamma = \left\{ y : \mathbb{R} \mapsto L^2(\mathbb{R}^3, \mathbb{C}^4), \exists C > 0, \|y(t)\|_2 \leq Ce^{\Gamma|t|} \forall t \in \mathbb{R} \right\}.$$

For sufficiently small $\varepsilon > 0$, the map $\mathcal{G}_\varepsilon(y_0, \cdot)$ leaves Y_Γ invariant and is continuous for the norm

$$N_\Gamma : y \mapsto \sup_{t \in \mathbb{R}} \left\{ \|y(t)\|_2 e^{-\Gamma|t|} \right\}.$$

Moreover, it is a strict contraction for sufficiently small U and $\varepsilon > 0$. Actually we choose ε as a function of Γ which is $O(\Gamma)$. In fact, since $Y_\Gamma \subset Y_{\Gamma'}$ for $\Gamma < \Gamma'$, we obtain that ε as a function of U is $O(\Gamma_-(U))$. This proves the existence of the fixed point y .

Then we fix $h_U^{cs}(y_0) = y(0) - y_0$. The invariance of the graph of h_U^{cs} by the flow of (2.11) is immediate.

Now we prove the smoothness property. We have that $N_\varepsilon(U, \eta)$ is l times differentiable in η from $Y_{\Gamma'}$ to Y_Γ if $(l + 1)\Gamma' \leq \Gamma$ and that \mathcal{G}_ε is l times differentiable from $X_c(U) \oplus X_s(U) \times Y_{\Gamma''}$ to Y_Γ if $2l\Gamma'' \leq \Gamma$ (see [8]). We introduce the family $(\eta_n)_{n \in \mathbb{N}}$ satisfying

$$\eta_0 = 0 \quad \text{and} \quad \eta_{n+1} = \mathcal{G}_\varepsilon(y_0, \eta_n).$$

This sequence converges to y (the fixed point) in Y_Γ . Moreover, as functions of y_0 , the convergence is uniform in Y_Γ (endowed with the norm N_Γ) on bounded sets of $X_s(U) \oplus X_c(U)$.

We want to prove that the sequence of their derivatives of order k with respect to η also converges in Y_Γ on bounded sets for any $\Gamma < \gamma(U)$. We prove it by induction in k . So suppose that $(\partial^j \eta_n)_{n \in \mathbb{N}}$ is converging in Y_Γ for all $j < k$ and any $\Gamma < \gamma(U)$. Then we have that for any $k \geq 2$ (see [8])

$$\begin{aligned} \partial \eta_n &= \partial \mathcal{G}_\varepsilon(y_0, \eta_{n-1}) = L + M(\partial_\eta N_\varepsilon(U, \eta_{n-1}) \partial \eta_{n-1}), \\ \partial^k \eta_n &= \partial^k \mathcal{G}_\varepsilon(y_0, \eta_{n-1}) = M(\partial_\eta N_\varepsilon(U, \eta_{n-1}) \partial^k \eta_{n-1} + \Psi_k(\eta_{n-1}, \dots, \partial^{k-1} \eta_{n-1})), \end{aligned}$$

with $L = e^{tJH(U)}$,

$$(M\eta)(t) = - \int_0^t e^{(t-s)JH(U)} \pi^{cs}(U) J\eta(s) ds + \int_t^{+\infty} e^{(t-s)JH(U)} \pi^u(U) J\eta(s) ds,$$

and Ψ_k a smooth function of k parameters. Hence since $M \circ \partial_\eta N_\varepsilon(U, y_{n-1})$ is a strict contraction in Y_Γ for sufficiently small ε and U (once more ε is $O(\Gamma_-(U))$), this proves the convergence of the sequence of k th derivatives in Y_Γ on bounded sets for any $\Gamma < \Gamma_-(U)$. Hence the sequences of derivatives of $(\eta_n)_{n \in \mathbb{N}}$ converge in Y_Γ on bounded sets for any $\Gamma < \Gamma_-(U)$. This gives the differentiability at any order of $y(0) = h(y_0)$. This also gives, since $N(U, \eta) = O(|\eta|^2)$ around zero, that $h(y_0) = O(|y_0|^2)$ around zero.

Now we want to prove that $W^{cs}(U)$ is attractive in negative time. In fact, $W^{cs}(U)$ is the graph of a smooth function $h : X^{cs} \mapsto X^u(U)$. Letting η be such that $S(U) + \eta$ is a solution of (1.8), we have

$$\partial_t \eta = JH(U)\eta + JN_\varepsilon(U, \eta),$$

where

$$\eta = y + r = y + h(y) + z$$

with $y = \pi^{cs}(U)\eta$, and we have the following equation for $z \in X^u(U)$:

$$\partial_t z = JH(U)z + M(U, y, z),$$

where

$$M(U, y, z) = \pi^u(U) \{JN_\varepsilon(U, \eta) - JN_\varepsilon(U, y + h(y))\} - Dh(y)\pi^{cs}(U) \{JN_\varepsilon(U, \eta) - JN_\varepsilon(U, y + h(y))\}.$$

Using Duhamel's formula, we obtain

$$z(t) = e^{tJH(U)} z(0) + \int_0^t e^{(t-s)JH(U)} M(U, y(s), z(s)) ds.$$

We obtain, since $z \in X^u(U)$,

$$\|z(t)\| \leq e^{\gamma(U)t} \|z(0)\| + C \int_0^t e^{(t-s)\gamma(U)} \|M(U, y(s), z(s))\| ds,$$

and so for $\gamma \in (0, \Gamma_-(U))$

$$e^{-\gamma t} \|z(t)\| \leq \|z(0)\| + C|U| \sup_{s \in [0, t]} \{e^{-\gamma s} \|z(s)\|\} + o\left(\sup_{s \in [0, t]} \{e^{-\gamma s} \|z(s)\|\}\right),$$

where C does not depend on U and z . Hence if $z(0)$ and U are small, we have that there exists $c > 0$ such that $\|z(t)\| \leq ce^{\gamma t}$ for all $t \leq 0$. We notice that, since $X^u(U) \subset H_\sigma^s$ for any $s \in \mathbb{R}^+$ and is finite dimensional (see Lemma 2.1), the time decay in L_σ^2 also gives a time decay in H_σ^s gives for any $s \in \mathbb{R}^+$.

Now choose V a sufficiently small neighborhood of 0 and ϕ a solution of (2.11) initially in V but not in $W^{cs}(U)$. Suppose ϕ stays in V in positive time. We obtain that $\phi \in Y_\Gamma$. We have

$$\begin{aligned} \phi(t) &= e^{tJH(U)} (\pi^s(U) + \pi^c(U)) \phi(0) \\ &+ \int_0^t e^{(t-s)JH(U)} \pi^s(U) JN_\varepsilon(U, \phi(s)) ds + \int_0^t e^{(t-s)JH(U)} \pi^c(U) JN_\varepsilon(U, \phi(s)) ds \\ &+ e^{tJH(U)} \pi^u(U) \left(\pi^u(U) \phi(0) + \int_0^\infty e^{-sJH(U)} \pi^u(U) JN_\varepsilon(U, \phi(s)) ds \right) \\ &- \int_t^\infty e^{(t-s)JH(U)} \pi^u(U) JN_\varepsilon(U, \phi(s)) ds \end{aligned}$$

with

$$\pi^u(U) \phi(0) + \int_0^\infty e^{-sJH(U)} \pi^u(U) JN_\varepsilon(U, \phi(s)) ds \neq 0.$$

Hence we obtain with (2.8) that $\phi(t)$ exponentially tends to infinity in norm. This is a contradiction, so ϕ leaves V in finite time. \square

Then reversing the time direction, that is to say, replacing H by $-H$ and F by $-F$, we obtain with this theorem a locally invariant center-unstable manifold with the following corresponding properties.

PROPOSITION 2.5 (center-unstable manifold). *If Assumptions 1.1–1.5 hold, then for any sufficiently small nonzero U , there exists around $S(U)$ a unique smooth invariant center-unstable manifold $W^{cu}(U)$ for (2.11), built as a graph with value in $X_S(U)$ and tangent to $S(U) + X_c(U) \oplus X_u(U)$ at $S(U)$.*

Any solution $\phi \in L^2_\sigma$ of (2.11) initially in the neighborhood of $S(U)$ tends as $t \rightarrow +\infty$ to $W^{cu}(U)$ with, for any $s \in \mathbb{R}^+$,

$$\text{dist}_{H^s_\sigma}(\phi(t), W^{cu}(U)) = O(e^{-\gamma t}) \text{ as } t \rightarrow +\infty,$$

and for $\gamma \in (0, \Gamma_+(U))$, any $s, \sigma \in \mathbb{R}$, and for any V a sufficiently small neighborhood of $S(U)$, any solution in V not in $W^{cu}(U)$ leaves V in finite negative time.

We can build in the same way a center manifold which is the intersection of the previous two manifolds.

PROPOSITION 2.6 (center manifold). *If Assumptions 1.1–1.5 hold, then for any sufficiently small nonzero U , there exists around $S(U)$ a unique smooth invariant center manifold $W^c(U)$ for (2.11), built as a graph with value in $X_s(U) \oplus X_u(U)$ and tangent to $S(U) + X_c(U)$ at $S(U)$.*

Moreover, we have that $W^c(U) = W^{cs}(U) \cap W^{cu}(U)$ and $W^c(U)$ contains the part of the PLS manifold which is in a small neighborhood of $S(U)$.

Proof. We build the center manifold with the same method as in the previous cases. We can also build a center-unstable manifold inside a center-stable manifold. More precisely, let $h^s_U : X_c(U) \oplus X_s(U) \mapsto X_u(U)$ be the map defining the center-stable manifold and let $h^u_U : X_c(U) \oplus X_u(U) \mapsto X_s(U)$ be the map defining the center-unstable manifold. A solution $y = S(U) + y_c + y_s + y_u$ with $y_* \in X_*(U)$ for $*$ $\in \{c, s, u\}$ is in the center-stable manifold if $y_u = h^s_U(y_c, y_s)$. Hence to obtain a center-unstable manifold inside a center-stable manifold one has to solve, for each y_c , the equation

$$y_s = h^u_U(y_c, h^s_U(y_c, y_s)),$$

which can be solved inside a small ball for small y_c and small U by means of the fixed point theorem, since $h_U^*(y_c, z)$ is $O(|y_c|^2 + |z|^2)$ around zero for $* \in \{s, u\}$.

In the same way, we can also build a center-stable manifold inside the center-unstable manifold.

Using the uniqueness of the center manifold, we obtain that these two manifolds are equal to the center manifold and $W^c(U) = W^{cs}(U) \cap W^{cu}(U)$.

Then any stationary states in a small neighborhood of $S(U)$ converge to $W^{cs}(U)$ and $W^{cu}(U)$ using the stabilization results of Propositions 2.4 and 2.5. Hence, we have that it belongs to $W^{cs}(U) \cap W^{cu}(U) = W^c(U)$. \square

In the following two sections, we study the dynamics inside and outside the center manifold, respectively.

3. The dynamics inside the center manifold. In this section, we prove that the dynamics inside the center manifold around $S(V_0)$, for small nonzero V_0 , relaxes towards the PLS manifold. To this end, we use Theorems 1.1 and 1.2 regarding the time decay of the propagator associated with H .

3.1. Decomposition of the system. As in [7], we decompose a solution $\phi \in W^c(V_0)$ of (1.8) with respect to the spectrum of $JH(U)$, with U specified in what follows, and we study the equations for these different parts of the decomposition. We introduce

$$\begin{aligned} \mathcal{H}_0^{\perp J}(u_1, u_2) &= \left\{ \eta \in L^2(\mathbb{R}^3, \mathbb{C}^8), \left\langle J\eta, \frac{\partial}{\partial \Re u_1} S(u_1, u_2) \right\rangle = \left\langle J\eta, \frac{\partial}{\partial \Im u_1} S(u_1, u_2) \right\rangle \right. \\ &= \left. \left\langle J\eta, \frac{\partial}{\partial \Re u_2} S(u_1, u_2) \right\rangle = \left\langle J\eta, \frac{\partial}{\partial \Im u_2} S(u_1, u_2) \right\rangle = 0 \right\}. \end{aligned}$$

In fact, we have

$$\mathcal{H}_0^{\perp J}(U) = \mathcal{H}_1(U) \oplus \mathcal{H}_c(U),$$

which is invariant under the action of $JH(U)$. We recall that $\mathcal{H}_1(U)$ is defined in Proposition 2.1 and $\mathcal{H}_c(U)$ in Proposition 2.2. We have the following lemma.

LEMMA 3.1. *If Assumptions 1.1–1.4 hold, let $s, \sigma \in \mathbb{R}$. Then there exist $\varepsilon, \varepsilon' > 0$ such that, for the manifold*

$$\Sigma = \left\{ (U, \eta), U \in B_{\mathbb{C}^2}(0, \varepsilon'), \eta \in \mathcal{H}_0^{\perp J}(U) \right\}$$

endowed with the metric of $\mathbb{C}^2 \times H_\sigma^s$ and any function $\phi \in B_{H_\sigma^s}(0, \varepsilon)$, there exists a unique $(U, \eta) \in \Sigma$ with

$$\phi = S(U) + \eta.$$

Moreover, there exists a neighborhood \mathcal{O} of $(0, 0) \in \Sigma$ such that the mapping $\phi \mapsto (U, \eta) \in \mathcal{O}$ is smooth.

Proof. To prove that Σ is a manifold, we use Proposition 2.2, which gives that it is locally isomorphic to some open subset of $\mathbb{C}^2 \times \mathcal{H}_c$ endowed with the metric of $\mathbb{C}^2 \times H_\sigma^s$. Then we work as in [22, Lemma 2.3]. Indeed we are looking, for all $\phi \in B_{H_\sigma^s}(0, \varepsilon)$, for a solution U in $\mathbb{C}^2 = \mathbb{R}^4$ of the equation

$$F(\phi, U) = (0, 0, 0, 0),$$

where

$$F(\phi, U) = \left(\left\langle J(\phi - S(U)), \frac{\partial}{\partial \Re u_1} S(u_1, u_2) \right\rangle, \left\langle J(\phi - S(U)), \frac{\partial}{\partial \Im u_1} S(u_1, u_2) \right\rangle, \right. \\ \left. \left\langle J(\phi - S(U)), \frac{\partial}{\partial \Re u_2} S(u_1, u_2) \right\rangle, \left\langle J(\phi - S(U)), \frac{\partial}{\partial \Im u_2} S(u_1, u_2) \right\rangle \right)$$

is smooth and satisfies $F(0, 0) = 0$ and

$$D_U F(0, 0) = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix},$$

which is invertible in \mathbb{R}^4 . Hence we apply the implicit function theorem. □

For any solution ϕ of (1.8) on an interval of time I containing 0, we write for $t \in I$

$$\phi(t) = e^{-i \int_0^t E(U(s)) ds} (S(U(t)) + \eta(t)),$$

where $\eta(t) \in \mathcal{H}_0^{\perp J}(U(t))$, and we want to solve the equation

$$(3.1) \quad \begin{aligned} i\partial_t \eta &= \{H - E(U)\} \eta + \{\nabla F(S(U) + \eta) - \nabla F(S(U))\} - idS(U)\dot{U} \\ &= \{H + d^2F(S(U)) - E(U)\} \eta + N(U, \eta) - idS(U)\dot{U} \end{aligned}$$

for $\eta(t) \in \mathcal{H}_0^{\perp J}(U(t))$. Here d^2F is the differential of ∇F and dS is the differential of S in \mathbb{R}^2 . To close the system, we need the equation for U . This follows from the condition

$$\langle \eta(t), JdS(U(t)) \rangle = 0.$$

After a time derivation (as in [7]), we obtain the equation

$$\dot{U}(t) = -A(U(t), \eta(t)) \langle N(U(t), \eta(t)), dS(U(t)) \rangle,$$

where

$$A(U, \eta) = [\langle JdS(U), dS(U) \rangle - \langle J\eta, d^2S(U) \rangle]^{-1};$$

indeed the matrix $[\langle JdS(U(t)), dS(U(t)) \rangle - \langle J\eta(t), d^2S(U(t)) \rangle]$ is invertible for small $|U(t)|$ and $\|\eta(t)\|_2$ since we have

$$[\langle JdS(U(t)), dS(U(t)) \rangle - \langle J\eta(t), d^2S(U(t)) \rangle] = \begin{pmatrix} J & 0_2 \\ 0_2 & J \end{pmatrix} + O(|U(t)| + \|\eta(t)\|_2).$$

We will need the following lemma.

LEMMA 3.2. *For any $s, s', \sigma \in \mathbb{R}$, any $p, q \in [1, \infty]$, and any $V_0 \in \mathbb{C}^2 \setminus \{0\}$ sufficiently small, there exist $\varepsilon, \varepsilon' > 0$ such that, for the manifold*

$$\mathcal{S}(V_0, \varepsilon) = \left\{ (U, z); U \in B_{\mathbb{C}^2}(V_0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H_\sigma^{s'}}(0, \varepsilon') \right\}$$

endowed with the metric of $\mathbb{C}^2 \times H_\sigma^{s'}$, there exists a unique map $g : S(V_0, \varepsilon) \mapsto B_{p,q}^s(\mathbb{R}^3, \mathbb{C}^4)$ which is smooth and satisfies $g(U, z) \in \mathcal{H}_1(U)$, $z + g(U, z) \in \mathcal{H}_0^{\perp J}(U)$, and $S(U) + z + g(U, z) \in W^c(V_0)$ for all $(U, z) \in S(V_0, \varepsilon)$. Moreover, we have $\|g(U, z)\|_{B_{p,q}^s} = O(\|z\|_{H^{s'}}^2)$.

Proof. The fact that $S(V_0, \varepsilon)$ is a manifold here is proved as in Lemma 3.1.

Let h^c be the function for which $W^c(V_0)$ is the graph. Any $\phi \in L^2(\mathbb{R}^3, \mathbb{R}^8)$ can be written in the form $S(V_0) + \tilde{U} \cdot DS(V_0) + \xi + \rho$ with $\rho \in \mathcal{H}_1(V_0)$ and $\xi \in \mathcal{H}_c(V_0)$. It can also be written in the form $S(U) + z + r$ with $r \in \mathcal{H}_1(U)$ and $z \in \mathcal{H}_c(U)$. These two decompositions in fact define two bijective smooth maps in sufficiently small sets (for the first we have a linear decomposition; for the second see Lemma 3.1). We write Ψ for the first and Φ for the second. Then $f = \Psi \circ \Phi^{-1}$ has three components following the decomposition $\mathcal{H}_0(V_0) \oplus \mathcal{H}_1(V_0) \oplus \mathcal{H}_c(V_0)$; we write them as (f_1, f_2, f_3) . Then g is the solution of the implicit equation in r ,

$$F(U, z, r) = f_2(U, z, r) - h^c(f_1(U, z, r), f_3(U, z, r)) = 0,$$

which can be solved by the implicit function theorem in $H_\sigma^{s'}$ since $\partial_r F(V_0, 0, 0)$ is invertible from $\mathcal{H}_1(V_0)$ to itself because $\partial_r f_2(V_0, 0, 0)$ ($f_2(V_0, r, 0) = r$) is invertible from $\mathcal{H}_1(V_0)$ to itself and $Dh_c(0, 0) = 0$.

The smoothness of g in the Besov spaces follows from the fact that $g(U, z) \in \mathcal{H}_1(U)$, and the exponential decay for excited states and their derivatives is given by (2.1).

Then we note that, for any U close to V_0 , the previous proof can be applied to $W^c(U)$. It shows that $W^c(U), W^c(V_0)$ are, in a neighborhood of $S(V_0)$, graphs of two functions on $S(V_0, \varepsilon)$, which are equal up to a translation in \mathbb{C}^2 of the first parameter. Hence their graphs are equal, so locally $W^c(U) = W^c(V_0)$. The last assertion then follows from the fact that at $S(U)$, $W^c(U)$ is tangent to $S(U) + X^c(U)$ and $X^c(U) \cap \mathcal{H}_1(U) = \{0\}$. \square

Hence decomposing η with respect to the spectrum of $JH(U)$, we write

$$\eta(t) = g(U(t), z(t)) + z(t)$$

with $z \in \mathcal{H}_c(U) \cap L^2(\mathbb{R}^3, \mathbb{R}^8)$. We obtain the system

$$\begin{cases} \dot{U} = -A(U, \eta)\langle N(U, \eta), dS(U) \rangle, \\ \partial_t z = JH(U)z + \mathbf{P}_c(U)JN(U, \eta) \\ \quad + \mathbf{P}_c(U(v))dS(U(v))A(U(v), \eta(v))\langle N(U(v), \eta(v)), dS(U(v)) \rangle \\ \quad + (dP_c(U))A(U, \eta)\langle N(U, \eta), dS(U) \rangle \eta \end{cases}$$

with

$$\eta(t) = z(t) + g(U(t), z(t)).$$

We note that this equation is defined only for z small with real values and U small. We now study this system.

3.2. The stabilization towards the PLS manifold. We now show that any solution of (1.8) which belongs to the center manifold $W^c(V_0)$, for a small nonzero V_0 , stabilizes as $t \rightarrow \pm\infty$ towards the manifold of the stationary states inside $W^c(V_0)$. To this end, we will use Theorems 1.1 and 1.2 to prove that z tends to zero in some sense.

Let us define for any $\varepsilon, \delta > 0$

$$\mathcal{U}(\varepsilon, \delta) = \left\{ U \in \mathcal{C}^1(\mathbb{R}, B_{\mathcal{C}^2}(V_0, \varepsilon)), \|\dot{U}\|_{L^1(\mathbb{R}) \cap L^\infty(\mathbb{R})} \leq \delta^2 \right\},$$

and for any $U \in \mathcal{U}(\varepsilon, \delta)$, let s, β be such that $s > \beta + 2 > 2$ and $\sigma > 3/2$,

$$\mathcal{Z}(U, \delta) = \left\{ z \in \mathcal{C}(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{R}^8)), z(t) \in \mathcal{H}_c(U(t)), \right. \\ \left. \max \left[\|z\|_{L^\infty(\mathbb{R}, H^s)}, \|z\|_{L^2(\mathbb{R}, H^s_\sigma)}, \|z\|_{L^2(\mathbb{R}, B^\beta_{\infty, 2})} \right] \leq \delta \right\},$$

and ε, δ are small enough to ensure that for $U \in \mathcal{U}(\varepsilon, \delta)$ and $z \in \mathcal{Z}(U, \delta)$

$$S(U) + z + g(U, z) \in W^c(V_0) \cap B_{H^s}(S(V_0), r(V_0)),$$

where g is defined in Lemma 3.2 and r is defined in Remark 2.3. It will appear later that δ is of the same order as $\|z_0\|_{H^s}$ (see Lemma 3.8).

3.2.1. Some useful lemmas. In the rest of our study, we will need some technical lemmas, which we collect here.

LEMMA 3.3. *If Assumptions 1.1–1.4 hold, let $\sigma, \sigma' \in \mathbb{R}, s > 1$, and $p, \tilde{p}_1, p_1, p_2, q$ in $[1, \infty]$ be such that*

$$\frac{1}{p} + \frac{s}{3} \geq \frac{1}{p_1} + \frac{1}{p_2} \geq \frac{1}{p}$$

and

$$\frac{1}{p} + \frac{s}{3} \geq \frac{1}{\tilde{p}_1}.$$

Then there exist $\varepsilon > 0$ and $C > 0$ such that, for all $U \in B_{\mathcal{C}}(0, \varepsilon)$ and $\eta \in B^s_{p_2, q}(\mathbb{R}^3, \mathbb{R}^8) \cap L^\infty(\mathbb{R}^3, \mathbb{R}^8)$ with $\langle Q \rangle^\sigma \eta \in B^s_{p_1, q}(\mathbb{R}^3, \mathbb{R}^8)$ and $\langle Q \rangle^{\sigma'} \eta \in B^s_{\tilde{p}_1, q}(\mathbb{R}^3, \mathbb{R}^8)$, we have

$$(3.2) \quad \|\langle Q \rangle^\sigma N(U, \eta)\|_{B^s_{p, q}} \leq C(s, F, |U| + \|\eta\|_{L^\infty}) |U| \|\eta\|_{L^\infty} \left\| \langle Q \rangle^{\sigma'} \eta \right\|_{B^s_{\tilde{p}_1, q}} \\ + C\left(s, F, |U| + \|\eta\|_{L^\infty \cap B^s_{p_2, q}}\right) \|\eta\|_{L^\infty}^2 \|\langle Q \rangle^\sigma \eta\|_{B^s_{p_1, q}}.$$

Proof. We recall the definition

$$N(U, \eta) = \nabla F(S(U) + \eta) - \nabla F(S(U)) - d^2 F(S(U))\eta.$$

We have

$$N(U, \eta) = \int_0^1 \int_0^1 d^3 F(S(U) + \theta' \theta \eta) \cdot \eta \cdot \theta \eta \, d\theta' \, d\theta$$

or

$$N(U, \eta) = \frac{1}{2} d^3 F(S(U)) \cdot \eta \cdot \eta + \int_0^1 \int_0^1 d^4 F(S(U) + \theta'' \theta' \theta \eta) \cdot \theta' \theta \eta \cdot \eta \cdot \theta \eta \, d\theta'' \, d\theta' \, d\theta.$$

Then we use, for $s \in \mathbb{R}^*_+$ and $p, p_1, p_2, q \in [1, \infty]$ such that $\frac{1}{p} + \frac{s}{3} \geq \frac{1}{p_1} + \frac{1}{p_2} \geq \frac{1}{p}$,

$$\|uv\|_{B^s_{p, q}} \leq C \|u\|_{B^s_{p_1, q}} \|v\|_{B^s_{p_2, q}},$$

and for $s > 1$, we use [19, Proposition 2.1]

$$\|d^k F(\psi)\|_{B_{p_2,q}^s} \leq C(s, F, \|\psi\|_{L^\infty}) \|\psi\|_{B_{p_2,q}^s}$$

for $k = 3$ or $k = 4$ and $d^4 F(z) = O(|z|)$; otherwise we decompose $d^4 F(z) = A + O(|z|)$, where A is a constant operator.

Eventually using Lemma 1.2 and

$$\left\| \langle Q \rangle^\sigma |\eta|^l \right\|_{B_{p_1,q}^s} \leq C \|\eta\|_{L^\infty}^{l-1} \|\langle Q \rangle^\sigma \eta\|_{B_{p_1,q}^s}$$

for $l \in \mathbb{N}$, we conclude the proof. \square

LEMMA 3.4. *If Assumptions 1.1–1.4 hold, let $\sigma \in \mathbb{R}$, $s > 1$, $p, p_1, p_2, q \in [1, \infty]$, and $\sigma_1, \sigma_2 \in \mathbb{R}$ such that*

$$\frac{1}{p} + \frac{s}{3} \geq \frac{1}{p_1} + \frac{1}{p_2} \geq \frac{1}{p}.$$

Then there exist $\varepsilon > 0$ and $C > 0$ such that, for all $U \in B_{\mathbb{C}}(0, \varepsilon)$ and $\eta \in B_{p,q}^s(\mathbb{R}^3, \mathbb{R}^8) \cap L^\infty(\mathbb{R}^3, \mathbb{R}^8)$ with $\langle Q \rangle^{\sigma_1} \eta \in B_{p_1,q}^s(\mathbb{R}^3, \mathbb{R}^8)$ and $\langle Q \rangle^{\sigma_2} \eta \in B_{p_2,q}^s(\mathbb{R}^3, \mathbb{R}^8)$, we have

$$\begin{aligned} & \|\langle Q \rangle^\sigma (\nabla F(S(U) + \eta) - \nabla F(S(U)) - \nabla F(\eta))\|_{B_{p,q}^s} \\ & \leq C(s, F, |U| + \|\eta\|_{L^\infty}) \left(|U| + \|\langle Q \rangle^{\sigma_1} \eta\|_{B_{p_1,q}^s} \right) |U| \|\langle Q \rangle^{\sigma_2} \eta\|_{B_{p_2,q}^s}. \end{aligned}$$

Proof. The proof is similar to that of Lemma 3.3. \square

LEMMA 3.5. *If Assumptions 1.1–1.4 hold, let $\sigma \in \mathbb{R}$, $s > 1$, and $p, q \in [1, \infty]$ such that $sp \geq 3$. Then there exist $\varepsilon > 0$ and $C > 0$ such that, for all $U, U' \in B_{\mathbb{C}^2}(0, \varepsilon)$ and $\eta, \eta' \in B_{p,q}^s(\mathbb{R}^3, \mathbb{R}^8)$, we have*

$$\begin{aligned} & \|\langle Q \rangle^\sigma \{N(U, \eta) - N(U', \eta')\}\|_{B_{p,q}^s} \leq C \left(s, F, |U| + |U'| + \|\eta\|_{B_{p,q}^s} + \|\eta'\|_{B_{p,q}^s} \right) \\ & \times \left\{ \left(\|\langle Q \rangle^{\sigma_1} \eta\|_{B_{p,q}^s} + \|\langle Q \rangle^{\sigma_1} \eta'\|_{B_{p,q}^s} \right)^2 \left(|U - U'| + \|\langle Q \rangle^{\sigma_2} (\eta - \eta')\|_{B_{p,q}^s} \right) \right. \\ & \quad \left. + \left(|U| + |U'| + \|\langle Q \rangle^{\sigma'_1} \eta\|_{B_{p,q}^s} + \|\langle Q \rangle^{\sigma'_1} \eta'\|_{B_{p,q}^s} \right) \right. \\ & \quad \left. \times \left(\|\langle Q \rangle^{\sigma'_2} \eta\|_{B_{p,q}^s} + \|\langle Q \rangle^{\sigma'_2} \eta'\|_{B_{p,q}^s} \right) \|\langle Q \rangle^{\sigma'_3} (\eta - \eta')\|_{B_{p,q}^s} \right\}, \end{aligned}$$

with $2\sigma_1 + \sigma_2 = \sigma'_1 + \sigma'_2 + \sigma'_3 = \sigma$ if $\langle Q \rangle^w \eta, \langle Q \rangle^w \eta' \in B_{p,q}^s(\mathbb{R}^3, \mathbb{R}^8)$ for $w \in \{\sigma_1, \sigma_2, \sigma'_1, \sigma'_2, \sigma'_3\}$.

Proof. Using the identity

$$N(u, \eta) = \int_0^1 \int_0^1 d^3 F(S(u) + \theta' \theta \eta) \cdot \eta \cdot \theta \eta \, d\theta' \, d\theta,$$

we can restrict our study to $d^3 F(\phi) - d^3 F(\phi')$. If $F = O(|z|^5)$, we have

$$\|\langle Q \rangle^\sigma (d^3 F(\phi) - d^3 F(\phi'))\|_{B_{p,q}^s} \leq \int_0^1 \|d^4 F(\phi + t(\phi - \phi'))\|_{B_{p,q}^s} \|\langle Q \rangle^\sigma (\phi - \phi')\|_{B_{p,q}^s} \, dt.$$

Then since $s > 1$ and $sp \geq 3$, we use

$$\|d^4 F(\psi)\|_{B_{p,q}^s} \leq C(s, F, \|\psi\|_{B_{p,q}^s}).$$

Using Lemma 1.2, we conclude the proof when $F = O(|z|^5)$.

Otherwise, if F is a homogeneous polynomial of order 4, the proof is easily adaptable since d^4F is a constant tensor.

The case $F = O(|z|^4)$ follows by summing the two previous cases since, as a function of $u \in \mathbb{R}^8$, $F(u) = Au^{\otimes 4} + O(|u|^5)$. \square

LEMMA 3.6. *If Assumptions 1.1–1.4 hold, let $\sigma \in \mathbb{R}$, $s \in \mathbb{R}$, and $p, q \in [1, \infty]$. Then there exist $\varepsilon > 0$, $M > 0$, and $C > 0$ such that, for all $U, U' \in B_{\mathbb{C}^2}(0, \varepsilon)$ and $\eta, \eta' \in B_{L^2(\mathbb{R}^3, \mathbb{R}^8)}(0, M)$ with $\langle Q \rangle^\sigma \{\eta - \eta'\} \in B_{p,q}^s(\mathbb{R}^3, \mathbb{R}^8)$, one has*

$$(3.3) \quad |A(U, \eta) - A(U', \eta')| \leq C \left\{ |U - U'| + \|\langle Q \rangle^\sigma \{\eta - \eta'\}\|_{B_{p,q}^s} \right\}.$$

Proof. We recall that

$$A(U, \eta) = [\langle JdS(U), dS(U) \rangle - \langle J\eta, d^2S(U) \rangle]^{-1}.$$

We have

$$\begin{aligned} A(U, \eta) - A(U', \eta') &= -[\langle JdS(U), dS(U) \rangle - \langle J\eta, d^2S(U) \rangle]^{-1} \\ &\times \{ \langle JdS(U), dS(U) \rangle - \langle J\eta, d^2S(U) \rangle - \langle JdS(U'), dS(U') \rangle + \langle J\eta', d^2S(U') \rangle \} \\ &\times [\langle JdS(U'), dS(U') \rangle - \langle J\eta', d^2S(U') \rangle]^{-1}. \end{aligned}$$

The lemma then follows from Lemma 1.2. \square

3.2.2. Global well-posedness for z . Let $U \in \mathcal{U}(\varepsilon, \delta)$ and $z_0 \in \mathcal{H}_c(U(0)) \cap H^s$. Let us write $U_\infty = \lim_{t \rightarrow +\infty} U(t)$; then we define $\mathcal{T}_{U, z_0}(z)$ by

$$\begin{aligned} \mathcal{T}_{U, z_0}(z)(t) &= e^{-itH+i \int_0^t E(U(r)) dr} z_0 \\ &+ \int_0^t e^{-i(t-v)H+i \int_v^t E(U(r)) dr} \mathbf{P}_c(U(v)) J \nabla F(\eta(v)) dv \\ &+ \int_0^t e^{-i(t-v)H+i \int_v^t E(U(r)) dr} \\ &\quad \times \mathbf{P}_c(U(v)) J \{ \nabla F(S(U(v)) + \eta(v)) - \nabla F(S(U(v)) - \nabla F(\eta(v)) \} dv \\ &+ \int_0^t e^{-i(t-v)H+i \int_v^t E(U(r)) dr} \\ &\quad \times \mathbf{P}_c(U(v)) dS(U(v)) A(U(v), \eta(v)) \langle N(U(v), \eta(v)), dS(U(v)) \rangle dv \\ &- \int_0^t e^{-i(t-v)H+i \int_v^t E(U(r)) dr} d\mathbf{P}_c(U(v)) \dot{U}(v) \eta(v) dv \end{aligned}$$

with

$$\eta(t) = z(t) + g(U(t), z(t)).$$

First, we have a local well-posedness result for z with the following lemma.

LEMMA 3.7. *If Assumptions 1.1–1.5 hold, then there exist $\varepsilon_0 > 0$ and $\delta_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U \in \mathcal{U}(\delta, \varepsilon)$, and $z_0 \in B_{H^s(0, \delta)} \cap \mathcal{H}_c(U(0))$ there are $T^\pm(z_0, U) > 0$ and a solution*

$$z \in \cap_{k=0}^2 C^k((-T^-(z_0, U); +T^+(z_0, U)), H^{s-k}(0, \delta))$$

of the equation

$$(3.4) \quad \begin{cases} \partial_t z = JH(U)z + \mathbf{P}_c(U)JN(u, \eta) - (d\mathbf{P}_c(U))\dot{U}\eta, \\ z(0) = z_0, \end{cases}$$

where $\eta(t) = z(t) + g(U(t), z(t))$.

Moreover, z is unique in $L^\infty((-T', T), H^s)$ for any $T \in (0, T^+(z_0, U))$ and $T' \in (0, T^-(z_0, U))$, and we have that if $T^+(z_0, U) < +\infty$, then

$$\lim_{t \rightarrow T^+(z_0, U)} \|z(t)\|_{H^s} \geq \delta,$$

and if $T^-(z_0, U) = +\infty$, then

$$\lim_{t \rightarrow -T^-(z_0, U)} \|z(t)\|_{H^s} \geq \delta.$$

Proof. The proof is a consequence of the fixed point theorem applied to \mathcal{T}_{U, z_0} .

Using Lemmas 3.3, 3.5, and 3.6 with the estimate (2.7)–(2.9) and the properties of g given by Lemma 3.2, we obtain that \mathcal{T}_{U, z_0} leaves a small ball in H^s invariant and is a contraction inside this ball.

Hence there exists a unique solution defined on the interval $[-T, T]$. Classical arguments permit us to extend the solution over a maximal interval $(-T^-(z_0, U), T^+(z_0, U))$ such that if $T^+(z_0, U) < \infty$, then necessarily the solution should leave a small ball in H^s at time $T^+(z_0, U)$. \square

We now have a global well-posedness result as stated in the following lemma.

LEMMA 3.8. *If Assumptions 1.1–1.5 hold, there exist $\varepsilon_0 > 0$, $\delta_0 > 0$, and $C > 0$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U \in \mathcal{U}(\varepsilon, \delta)$, and $z_0 \in B_{H^s}(0, \delta) \cap \mathcal{H}_c(U(0))$ we obtain, for the Cauchy problem (3.4), $T^+(U, z_0) = +\infty$, $T^-(U, z_0) = +\infty$, $z \in \mathcal{Z}(U, \delta)$, and*

$$\max \left[\|z\|_{L^\infty(\mathbb{R}, H^s)}, \|z\|_{L^2(\mathbb{R}, H^s_\sigma)}, \|z\|_{L^2(\mathbb{R}, B^\beta_{\infty, 2})} \right] \leq C \|z_0\|_{H^s}.$$

Proof. We have $(1 - P_c(U))z \equiv 0$ because $(1 - P_c(U(0)))z(0) = 0$ and its time derivative is zero.

Let us introduce for any $0 < T < T^+(U, z_0)$ the function

$$m(T) = \sup_{t \in (-T, T)} \left\{ \|z\|_{L^\infty((-T, T), H^s)}, \|z\|_{L^2((-T, T), H^s_\sigma)}, \|z\|_{L^2((-T, T), B^\beta_{\infty, 2})} \right\}.$$

First, we study the estimation of $L^2((-T, T), H^s_\sigma)$. We use estimate (2.4) and the estimates of Theorem 1.1:

$$\begin{aligned} & \|z\|_{L^2((-T, T), H^s_\sigma)} \\ & \leq C_0 \|z_0\|_{H^s} + C \left\| \mathbf{P}_c \int_0^t e^{-i(t-v)H + i \int_v^t E(U(r)) dr} \mathbf{P}_c(U(v)) J \nabla F(\eta(v)) dv \right\|_{L^2((-T, T), H^s_\sigma)} \\ & \quad + C \|\nabla F(S(U) + \eta) - \nabla F(S(U) - \nabla F(\eta))\|_{L^2((-T, T), H^s_\sigma)} \\ & \quad + C \|dS(U)A(U, \eta)\langle N(U, \eta), dS(U) \rangle\|_{L^2((-T, T), H^s_\sigma)} \\ & \quad + C \left\| (d\mathbf{P}_c(U)\dot{U}\eta) \right\|_{L^2((-T, T), H^s_\sigma)}. \end{aligned}$$

We now study the estimation of the third term of the right-hand side:

$$\begin{aligned} & \left\| \int_0^t e^{-i(t-v)H+i\int_v^t E(U(r)) dr} \mathbf{P}_c \mathbf{P}_c(U(v)) J \nabla F(\eta(v)) dv \right\|_{L^2_t((-T,T), H^s_{-\sigma})} \\ & \leq \int_{-T}^T \left\| e^{-i(t-v)H+i\int_v^t E(U(r)) dr} \mathbf{P}_c \mathbf{P}_c(U(v)) J \nabla F(\eta(v)) \right\|_{L^2_t((-T,T), H^s_{-\sigma})} dv \\ & \leq C(U) \|\nabla F(\eta)\|_{L^1((-T,T), H^s)} \\ & \leq C(U) \|\eta\|_{L^2((-T,T), L^\infty)}^2 \|\eta\|_{L^\infty((-T,T), H^s)}, \end{aligned}$$

where we used Theorem 1.1, estimate (1.2). Hence for the $L^2 H^s_{-\sigma}$ estimate, we obtain

$$\begin{aligned} \|z\|_{L^2((-T,T), H^s_{-\sigma})} & \leq C_0 \|z_0\|_{H^s} + C \|\eta\|_{L^2((-T,T), L^\infty)}^2 \|\eta\|_{L^\infty((-T,T), H^s)} \\ & \quad + C \left(\|U\|_\infty + \|\eta\|_{L^\infty((-T,T), H^s_{-\sigma})} \right) \|U\|_\infty \|\eta\|_{L^2((-T,T), H^s_{-\sigma})} \\ & \quad + C \|\eta\|_{L^2((-T,T), L^\infty)}^2 + C \|\dot{U}\|_{L^2} \|\eta\|_{L^\infty((-T,T), H^s)}. \end{aligned}$$

Using Lemma 3.2, we obtain

$$\|z\|_{L^2((-T,T), H^s_{-\sigma})} \leq C_0 \|z_0\|_{H^s} + Cm(T)^3 + Cm(T)^2 + C\epsilon m(T) + C\delta^2 m(T),$$

where C depends on $\|U\|_\infty$ and $\|\eta\|_{L^\infty((-T,T), H^s)}$.

Then, we estimate the H^s norm. Using estimate (2.4), we have

$$\begin{aligned} \|z(t)\|_{H^s} & \leq \|z_0\|_{H^s} + \int_{-T}^T \|\nabla F(\eta(v))\|_{H^s} dv \\ & \quad + \left\| \int_0^t e^{-i(t-v)H+i\int_v^t E(U(r)) dr} \mathbf{P}_c(U(v)) \right. \\ & \quad \left. \times J\{\nabla F(S(U(v)) + \eta(v)) - \nabla F(S(U(v))) - \nabla F(\eta(v))\} dv \right\|_{H^s} \\ & \quad + \int_{-T}^T \|dS(U(v))A(U(v), \eta(v))\langle N(U(v), \eta(v)), dS(U(v)) \rangle\|_{H^s} dv \\ & \quad + \int_{-T}^T \left\| d\mathbf{P}_c(U(v))\dot{U}(v)\eta(v) \right\|_{H^s} dv. \end{aligned}$$

To estimate the third term of the right-hand side, we use the H -smoothness estimates,

more precisely, Theorem 1.1, estimate (1.1), and then we use Lemma 3.4:

$$\begin{aligned} & \left\| \int_0^t e^{-i(t-v)H+i\int_v^t E(U(r)) dr} \right. \\ & \quad \times \mathbf{P}_c(U(v))J\{\nabla F(S(U(v)) + \eta(v)) - \nabla F(S(U(v))) - \nabla F(\eta(v))\} dv \left. \right\|_{H^s} \\ & \leq \left\| \int_0^t e^{ivH-i\int_0^v E(U(r)) dr} \right. \\ & \quad \times \mathbf{P}_c(U(v))J\{\nabla F(S(U(v)) + \eta(v)) - \nabla F(S(U(v))) - \nabla F(\eta(v))\} dv \left. \right\|_{H^s} \\ & \leq C \|\nabla F(S(U) + \eta) - \nabla F(S(U)) - \nabla F(\eta)\|_{L^2((-T,T),H^s_\sigma)} \\ & \leq C \left(\|U\|_\infty + \|\eta(v)\|_{L^\infty((-T,T),H^s)} \right) \|U\|_\infty \|\eta\|_{L^2((-T,T),H^s_\sigma)}. \end{aligned}$$

Hence for the $L^\infty H^s$ estimate, we obtain

$$\begin{aligned} \|z(t)\|_{H^s} & \leq \|z_0\|_{H^s} + C \|\eta\|_{L^\infty((-T,T),H^s)} \|\eta\|_{L^2((-T,T),L^\infty)}^2 \\ & + C \left(\|U\|_{L^\infty((-T,T))} + \|\eta(v)\|_{L^\infty((-T,T),H^s)} \right) \|U\|_{L^\infty((-T,T))} \|\eta\|_{L^2((-T,T),H^s_\sigma)} \\ & + C \|\eta\|_{L^2((-T,T),L^\infty)}^2 + \|\dot{U}\|_{L^1((-T,T))} \|\eta\|_{L^\infty((-T,T),H^s)}. \end{aligned}$$

Using Lemma 3.2, we obtain

$$\|z(t)\|_{H^s} \leq \|z_0\|_{H^s} + Cm(T)^3 + Cm(T)^2 + C\epsilon m(T) + C\delta^2 m(T),$$

where C depends on $\|U\|_\infty$ and $\|\eta\|_{L^\infty((-T,T),H^s)}$.

For the $L^2 B_{\infty,2}^\beta$ estimate, by Proposition 2.2 and Theorem 1.2, we have for any $\delta > 0$, any $p_\delta > 3/\delta$, and $\theta_\delta = \frac{4}{p_\delta - 2}$

$$\begin{aligned} \|z\|_{L^2((-T,T),B_{\infty,2}^\beta)} & \leq \|z\|_{L^2((-T,T),B_{p_\delta,2}^{\beta+\delta})} \\ & \leq C_0 \|z_0\|_{H^{\beta+1+\delta+\theta_\delta/2}} + C \|d^2 F(S(U)) \cdot \eta\|_{L^2((-T,T),B_{p_\delta,2}^{\beta+2+\delta+\theta_\delta})} \\ & \quad + C \|N(U, \eta)\|_{L^1((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} \\ & \quad + C \|dS(U)A(U, \eta)\langle N(U, \eta), dS(U) \rangle\|_{L^1((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} \\ & \quad + C \left\| (d\mathbf{P}_c(U))\dot{U}\eta \right\|_{L^1((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} dv. \end{aligned}$$

With Lemmas 3.3 and 3.4, we infer that

$$\begin{aligned} \|z\|_{L^2(\mathbb{R},B_{\infty,2}^\beta)} & \leq C_0 \|z_0\|_{H^{\beta+1+\delta+\theta_\delta/2}} + C|U|_\infty \|\eta\|_{L^2((-T,T),H_{-\sigma}^{\beta+2+\delta+\theta_\delta})} \\ & + C|U|_\infty \|z\|_{L^2((-T,T),L^\infty)} \|z\|_{L^2((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} \\ & + C(|U|_\infty + \|\eta\|_{L^\infty((-T,T),H^{\beta+1+\delta+\theta_\delta/2})}) \|\eta\|_{L^2((-T,T),L^\infty)}^2 \|z\|_{L^\infty((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} \end{aligned}$$

$$\begin{aligned}
 &+ C(\|U\|_\infty + \|\eta\|_{L^\infty((-T,T),H^{\beta+1+\delta+\theta_\delta/2})}) \|\eta\|_{L^2((-T,T),H_{-\sigma}^{\beta+1+\delta+\theta_\delta/2})} \\
 &\times \|\eta\|_{L^\infty((-T,T),H^{\beta+1+\delta+\theta_\delta/2})} \\
 &+ C\|\dot{U}\|_{L^1} \|\eta\|_{L^\infty((-T,T),H^{\beta+1+\delta+\theta_\delta/2})}.
 \end{aligned}$$

Since for small $\delta > 0$, $s \geq \beta + 2 + \delta + \theta_\delta$ and using Lemma 3.2, we infer that

$$\begin{aligned}
 \|z\|_{L^2((-T,T),B_{\infty,2}^\beta)} &\leq C_0\|z_0\|_{H^{\beta+1+\delta+\theta_\delta/2}} \\
 &+ Cm(T)^3 + Cm(T)^2 + C\varepsilon m(T) + C\delta^2 m(T).
 \end{aligned}$$

Hence we obtain

$$m(T) \leq C_0\|z_0\|_{H^{\beta+1+\delta+\theta_\delta/2}} + C\varepsilon m(T) + C\delta^2 m(T) + Cm(T)^3 + Cm(T)^2,$$

where C_0 does not depend on m and C is a nondecreasing function of $\|z\|_{L^\infty((-T,T),H^s)}$ and $\|U\|_\infty$ and hence can be bounded by a nondecreasing function of m .

If $\|z_0\|_{H^s}$ is small, then $m(0)$ is small and $m(T)$ stays small. Therefore we have that $z \in \mathcal{Z}(U, \delta)$ if $\|z_0\|_{H^s}, \delta$, and ε are small enough; moreover,

$$\max \left[\|z\|_{L^\infty(\mathbb{R},H^s)}, \|z\|_{L^2(\mathbb{R},H_{-\sigma}^s)}, \|z\|_{L^2(\mathbb{R},B_{\infty,2}^\beta)} \right] \leq f(\|z_0\|_{H^s}),$$

where f is such that there exists $C > 0$ with

$$f(\|z_0\|_{H^s}) \leq C \|z_0\|_{H^s}. \quad \square$$

The solution z just found is a function of z_0 and U ; writing it $z[z_0, U]$, we have the following important property given by the following lemma.

LEMMA 3.9. *If Assumptions 1.1–1.5 hold, then for any $T > 0$, there exist $\varepsilon_0 > 0$, $\delta_0 > 0$, $C > 0$, and $\kappa \in (0, 1)$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U, U' \in \mathcal{U}(\varepsilon, \delta)$, $z_0 \in \mathcal{H}_c(U(0))$, $z'_0 \in \mathcal{H}_c(U'(0))$, $z \in \mathcal{Z}(U, \delta)$, and $z' \in \mathcal{Z}(U', \delta)$, one has*

$$\begin{aligned}
 &\|z[z'_0, U'] - z[z_0, U]\|_{L^\infty((-T;T),H^s) \cap L^2((-T;T),L^\infty) \cap L^2((-T;T),H_{-\sigma}^s)} \\
 &\leq C \|z_0 - z'_0\|_{H^s} + \kappa \left\{ \|U - U'\|_{L^\infty((-T;T))} + \left\| \dot{U} - \dot{U}' \right\|_{L^\infty((-T;T))} \right\}.
 \end{aligned}$$

Proof. We use the techniques of the previous lemma. □

3.2.3. Global well-posedness for U and its stabilization. Here we want to solve the equation for U . We note that z has been built in the previous section and is a function of U and $z_0 \in \mathcal{H}_c(U(0))$. Let us introduce for any $U_0 \in B_{\mathbb{C}}(0, \varepsilon)$ the function on $\mathcal{U}(\varepsilon, \delta)$:

$$f_{U_0}(U)(t) = U_0 - \int_0^t A(U(v), \eta(v)) \langle N(U(v), \eta(v)), dS(U(v)) \rangle dv,$$

where $\eta = z(t) + g[U(t), z(t)]$. We have the following lemma.

LEMMA 3.10. *If Assumptions 1.1–1.5 hold, there exist $\varepsilon_0 > 0$ and $\delta_0 > 0$ such that, for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, the function f_{U_0} maps $\mathcal{U}(\varepsilon, \delta)$ into itself if U_0 and $z_0 \in H^s \cap \mathcal{H}_c(U_0)$ are small enough.*

Proof. By means of Lemma 3.3, we obtain

$$\|\partial_t f_{U_0}(U)\|_{L^1(\mathbb{R}) \cap L^\infty(\mathbb{R})} \leq C \|N(U(v), \eta(v))\|_{L^1(\mathbb{R}, H^s_\sigma) \cap L^\infty(\mathbb{R}, H^s)} \leq \delta^2$$

and

$$\|f_{U_0}(U)\|_{L^\infty(\mathbb{R})} \leq |U_0| + C \|N(U(v), \eta(v))\|_{L^1(\mathbb{R}, H^s)} \leq |U_0| + \delta^2;$$

hence for sufficiently small U_0 and δ , we obtain the lemma. \square

The function f_{U_0} also has a local Lipschitz property, as stated by the following lemma.

LEMMA 3.11. *If Assumptions 1.1–1.5 hold, for any $T > 0$, there exist $\varepsilon_0 > 0$, $\delta_0 > 0$, and $\kappa \in (0, 1)$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U, U' \in \mathcal{U}(\varepsilon, \delta)$, for any $z_0 \in \mathcal{H}_c(U(0)) \cap H^s$, for any $z'_0 \in \mathcal{H}_c(U'(0)) \cap H^s$ small enough, and for U_0, U'_0 small enough, one has*

$$\begin{aligned} & |f_{U_0}(U) - f_{U'_0}(U')|_{L^\infty((-T; T))} + |\partial_t f_{U_0}(U) - \partial_t f_{U'_0}(U')|_{L^1((-T; T))} \\ & \leq |U_0 - U'_0| + \kappa \left(\|U - U'\|_{L^\infty((-T; T))} + \|\dot{U} - \dot{U}'\|_{L^1((-T; T))} + \|z_0 - z'_0\|_{H^s} \right). \end{aligned}$$

Proof. The proof is a straightforward consequence of Lemmas 3.5, 3.6, and 3.9. \square

We now obtain the following lemma.

LEMMA 3.12. *If Assumptions 1.1–1.5 hold, there exist $\varepsilon > 0$ and $\delta > 0$ such that, for any $U_0 \in \mathbb{C}$ small and $z_0 \in \mathcal{H}_c(U_0) \cap H^s_\sigma$ small, the equation*

$$(3.5) \quad \begin{cases} \dot{U} &= -A(U, \eta) \langle N(U, \eta), dS(U) \rangle, \\ U(0) &= U_0, \end{cases}$$

where $\eta(t) = z(t) + g[U(t), z(t)]$, has a unique solution in $\mathcal{U}(\delta, \varepsilon)$. Moreover, there exists $C > 0$ such that

$$|U_{\pm\infty} - U_0| \leq C \|z_0\|_{H^s}^2.$$

Proof. This is also a fixed point result for f_{U_0} . Let us fix $T > 0$ and consider, for any $V \in \mathcal{U}(\delta, \varepsilon)$ with sufficiently small $\delta > 0$ and $\varepsilon > 0$, the sequence

$$\begin{cases} V_{n+1} = f_{U_0}(V_n) \quad \forall n \in \mathbb{N}, \\ V_0 = V \end{cases}$$

for any $n \in \mathbb{N}$, $V_n \in \mathcal{U}(\delta, \varepsilon)$. With Lemma 3.11, the fixed point theorem gives us the convergence for the norms of $L^\infty((-T, T))$ and $\dot{W}^{1,1}((-T, T))$ of $(V_n)_{n \in \mathbb{N}}$.

Then we notice that, for any $T' \in \mathbb{R}$, we have

$$V_{n+1}(t) = f_{f_{U_0}(V_n)(T')}(V_n)(t - T').$$

Since for $T' \in (-T; T)$, $(f_{U_0}(V_n)(T'))$ is a Cauchy sequence, Lemma 3.11 gives us the convergence of (V_n) for the norms of $L^\infty((T' - T; T' + T))$ and $\dot{W}^{1,1}((T' - T; T' + T))$.

Iterating this process, we obtain that the sequence converges uniformly locally in time and we prove the lemma since the other statements are classical. We note just that the last statement follows from the fact that there exists $C > 0$ such that

$$\int_{\mathbb{R}^\pm} |\dot{U}(v)| \, dv \leq \int_{\mathbb{R}^\pm} |A(U(v), \eta(v)) \langle N(U(v), \eta(v)), dS(U(v)) \rangle| \, dv \leq C \|z_0\|_{H^s}^2. \quad \square$$

3.2.4. The asymptotic profile of z . In this section, our aim is to specify the asymptotic profile of z when z_0 is localized. First we state the following proposition.

PROPOSITION 3.1. *There exists $\varepsilon > 0$ such that for all $U \in B_{\mathbb{C}^2}(0, \varepsilon)$ and $\alpha \in \mathbb{R}^+$ there exists $C > 0$ such that*

$$\left\| \langle Q \rangle^\alpha e^{JtH(U)} P_c(U) \psi \right\| \leq C_\alpha \sum_{\beta=0}^\alpha \langle t \rangle^\beta \left\| \langle Q \rangle^{\alpha-\beta} \psi \right\|$$

for any $\psi \in L^2(\mathbb{R}^3, \mathbb{C}^8)$.

Proof. From Proposition 2.2, we obtain the lemma for $\alpha = 0$; then we just need the estimate

$$\left\| Q^\alpha e^{JtH(U)} \psi \right\|^2 \leq C_\alpha^2 \sum_{0 \leq \beta \leq \alpha} |t|^{2|\beta|} \left\| Q^{\alpha-\beta} \psi \right\|^2$$

for any $\psi \in L^2(\mathbb{R}^3, \mathbb{C}^8)$, $\alpha \in \mathbb{N}^3$, and some $C > 0$ independent of ψ . The rest of the proposition will follow by interpolation.

For $U = 0$, this follows by an iterated proof from the identity

$$\frac{d}{dt} e^{itH} Q e^{-itH} = e^{itH} \alpha e^{-itH},$$

where α is the 3-vector of Dirac–Pauli matrices defined in the introduction. For $U \neq 0$, we use the same proof with the exponential decay of Proposition 1.1. \square

We can improve Lemma 3.13 if we use [7, Theorem 1.1] and [7, Theorem 1.2], which we repeat in what follows.

THEOREM 3.1 (Theorem 1.1 of [7]: propagation for perturbed Dirac dynamics). *Assume that Assumptions 1.1 and 1.2 hold, and let $\sigma > \frac{5}{2}$. Then one has*

$$\| e^{-itH} \mathbf{P}_c(H) \|_{B(L^2_\sigma, L^2_{-\sigma})} \leq C \langle t \rangle^{-\frac{3}{2}}.$$

We also have the following proposition.

PROPOSITION 3.2 (Proposition 2.2 of [7]: propagation far from thresholds). *Suppose that Assumption 1.1 holds. Then, for any $\chi \in C^\infty(\mathbb{R}^3, \mathbb{C}^4)$ bounded with support in $\mathbb{R} \setminus (-m; m)$ and for any $\sigma \geq 0$, there is $C > 0$ such that*

$$\| e^{-itH} \chi(H) \|_{B(L^2_\sigma, L^2_{-\sigma})} \leq C \langle t \rangle^{-\sigma}.$$

Using Duhamel’s formula as in Proposition 2.2 and interpolating with estimate (2.6), we obtain the following corollary.

COROLLARY 3.1. *Assume that Assumptions 1.1 and 1.2 hold, and let $\theta \geq 0$ and $\sigma > \frac{5}{2}\theta$. Then there exists $\varepsilon > 0$ such that for all $U \in B_{\mathbb{C}^2}(0, \varepsilon)$ one has*

$$\| e^{JtH(U)} \mathbf{P}_c(U) \|_{B(L^2_\sigma, L^2_{-\sigma})} \leq C \langle t \rangle^{-\frac{3\theta}{2}}.$$

THEOREM 3.2 (Theorem 1.2 of [7]: dispersion for perturbed Dirac dynamics). *Assume that Assumptions 1.1 and 1.2 hold. Then for $p \in [1, 2]$, $\theta \in [0, 1]$, $s - s' \geq (2 + \theta)(\frac{2}{p} - 1)$, and $q \in [1, \infty]$ there exists a constant $C > 0$ such that*

$$\| e^{-itH} \mathbf{P}_c(H) \|_{B_{p,q}^s, B_{p',q}^{s'}} \leq C (K(t))^{\frac{2}{p}-1}$$

with $\frac{1}{p} + \frac{1}{p'} = 1$, and

$$K(t) = \begin{cases} |t|^{-1+\theta/2} & \text{if } |t| \in (0, 1], \\ |t|^{-1-\theta/2} & \text{if } |t| \in [1, \infty). \end{cases}$$

Using Duhamel’s formula once more, Theorem 3.1, and Corollary 3.1, we obtain the following corollary.

COROLLARY 3.2. *Assume that Assumptions 1.1 and 1.2 hold, and let $p \in [1, 2]$, $\theta \in [0, 1]$, $s - s' \geq (2 + \theta)(\frac{2}{p} - 1)$, $q \in [1, \infty]$, and $\sigma > \max\{\frac{3}{2}, (\frac{2}{p} - 1)(1 + \frac{\theta}{2})\}$. Then there exists $\varepsilon > 0$ such that for all $U \in B_{C^2}(0, \varepsilon)$ one has*

$$\|e^{JtH(U)} \mathbf{P}_c(U)\|_{H^s_\sigma, B^{s'}_{p', q}} \leq C (K(t))^{\frac{2}{p}-1}$$

with $\frac{1}{p} + \frac{1}{p'} = 1$, and

$$K(t) = \begin{cases} |t|^{-1+\theta/2} & \text{if } |t| \in (0, 1], \\ |t|^{-1-\theta/2} & \text{if } |t| \in [1, \infty). \end{cases}$$

Proof. We first prove this for $U = 0$. We have to study the high and low energy parts in a different manner. For the low energy part, we iterate twice Duhamel’s formula with respect to D_m in order to use Theorems 3.1 and 3.2 for the free case.

In the high energy part, we use also Duhamel’s formula but with Theorem 3.2 for the free case and Proposition 3.2.

Then for $U \neq 0$, we work as for estimate (2.6). □

We obtain the following lemma.

LEMMA 3.13. *If Assumptions 1.1–1.5 hold, there exist $\varepsilon_0 > 0$, $\delta_0 > 0$, and $C > 0$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U_0 \in B_{C^2}(0, \varepsilon)$, and $z_0 \in B_{H^s_\sigma}(0, \delta) \cap \mathcal{H}_c(U_0)$ we obtain for the Cauchy problem (3.4) (with U the solution of (3.5)) a global solution z such that*

$$\max \left[\sup_{t \in \mathbb{R}} (\|z(t)\|_{H^s}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z(t)\|_{H^s_\sigma}), \right. \\ \left. \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z(t)\|_{B^\beta_{\infty, 2}}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{-3/2} \|z(t)\|_{H^s_{3/2}}) \right] \leq C \|z_0\|_{H^s_\sigma}.$$

Proof. The proof is similar to that of Lemma 3.8 with some adaptations involving the norm H^s_σ ; we also refer to the proof of [7, Lemma 5.5]. Indeed, we work in the spaces

$$\tilde{\mathcal{U}}(\varepsilon, \delta) = \left\{ U \in C^1(\mathbb{R}, B_{C^2}(V_0, \varepsilon)), \|\dot{V}(t)\| \leq \frac{\delta^2}{\langle t \rangle^3} \right\}$$

and

$$\mathcal{Z}(U, \delta) = \left\{ z \in \mathcal{C}(\mathbb{R}, L^2(\mathbb{R}^3, \mathbb{R}^8)), z(t) \in \mathcal{H}_c(U(t)), \right. \\ \max \left[\sup_{t \in \mathbb{R}} (\|z(t)\|_{H^s}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z(t)\|_{H^s_\sigma}), \right. \\ \left. \left. \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z(t)\|_{B^\beta_{\infty, 2}}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{-3/2} \|z(t)\|_{H^s_{3/2}}) \right] \leq \delta \right\}.$$

Let

$$t \mapsto \xi_{\pm}(t) = e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} z(t)$$

and

$$t \mapsto V_{\pm}(t) = e^{-i \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} U(t).$$

We have that $V_{\pm\infty} \lim_{\pm\infty} V_{\pm}(t)$ exist and we use exactly the same method as that of Lemma 3.8, applied to

$$\begin{aligned} \xi_{\pm}(t) &= e^{JtH(V_{\pm\infty})} z_0 \\ &+ \int_0^t e^{J(t-s)H(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) J (d^2 F(S(V_{\pm}(v))) - d^2 F(S(V_{\pm\infty}))) \xi_{\pm}(v) dv \\ &+ \int_0^t e^{J(t-s)H(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) JN(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) dv \\ &+ \int_0^t e^{J(t-s)H(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) dS(V(v)) A(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) \\ &\quad \times \langle N(V_{\pm}(v), \tilde{\eta}_{\pm}(v)), dS(V_{\pm}(v)) \rangle dv \\ &- \int_0^t e^{J(t-s)H(V_{\pm\infty})} (d\mathbf{P}_c(V_{\pm}(v))) A(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) \\ &\quad \times \langle N(V_{\pm}(v), \tilde{\eta}_{\pm}(v)), dS(V_{\pm}(v)) \rangle \tilde{\eta}_{\pm}(v) dv, \end{aligned}$$

with $\tilde{\eta}_{\pm}(t) = e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} (z(t) + g(U(t), z(t)))$, but using the previous time decay estimates.

There are two differences from the proof of Lemma 3.8:

One is in the estimate of the $H_{-\sigma}^s$ norm. In fact, before using the time decay estimates for $e^{-itH} P_c(H)$, we split the space associated with the continuous spectrum into two parts: one associated with energy close to the thresholds and one associated with the rest of the spectrum. In the first part, we use the fact that $\sigma > 3/2$ to estimate the $H_{-\sigma}^s$ by the $B_{\infty,2}^{\beta}$ norm since we work with bounded energies. In the second part, since we work far from thresholds, we use Proposition 3.2 after estimating the $H_{-\sigma}^s$ by the $H_{-3/2}^s$ norm.

The other difference is in the estimation of the $B_{\infty,2}^{\beta}$ norm. We use Corollary 3.2 for $e^{JtH(V_{\pm\infty})} z_0$ and Theorem 3.2 for the integrals. \square

We have that $\lim_{\pm\infty} U = U_{\pm\infty}$ exist. If $z_0 \in H_{\sigma}^s$, then the associated solution U satisfies

$$|\dot{U}| \leq \frac{C}{\langle t \rangle^3} \|z_0\|_{H_{\sigma}^s},$$

and we have

$$\int_0^t (E(U(v)) - E(U_{\pm\infty})) dv \rightarrow E_{\pm\infty} \quad \text{as } t \rightarrow \pm\infty$$

for some real $E_{\pm\infty}$. We introduce

$$V_{\pm}(t) = e^{-i \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} U(t),$$

which have a limit as $t \rightarrow \pm\infty$, respectively, as being

$$V_{\pm\infty} = e^{-iE_{\pm\infty}} U_{\pm\infty}.$$

Then we note that we can also obtain an asymptotic profile for $e^{itH+itE(U_{\infty})}z(t)$ if z_0 is localized. But we prefer to obtain a scattering result with respect to $e^{JtH(V_{\infty})}$. We have the following lemma.

LEMMA 3.14. *If Assumptions 1.1–1.5 hold, then there exist $\varepsilon_0 > 0$, $\delta_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U_0 \in B_{\mathbb{C}^2}(0, \varepsilon)$, and $z_0 \in B_{H^s_\sigma}(0, \delta) \cap \mathcal{H}_c(U_0)$ and for the solution z of (3.4) (with U the solution of (3.5)) given in Lemma 3.7, the limit*

$$z_{\pm\infty} = \lim_{t \rightarrow \pm\infty} e^{-JtH(V_{\pm\infty})} e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} z(t)$$

exists in H^s . Moreover, we have $z_{\pm\infty} \in \mathcal{H}_c(V_{\pm\infty}) \cap H^s_\sigma$, and there exists $C > 0$ such that

$$\begin{aligned} & \max \left\{ \|e^{-J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} e^{JtH(V_{\pm\infty})} z_{\pm\infty} - z(t)\|_{H^s}, \right. \\ & \|e^{-J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} e^{JtH(V_{\pm\infty})} z_{\pm\infty} - z(t)\|_{H^s_{-\sigma}}, \\ & \left. \|e^{-J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} e^{JtH(V_{\pm\infty})} z_{\pm\infty} - z(t)\|_{B^\beta_{\infty,2}} \right\} \leq \frac{C}{\langle t \rangle^2} \|z_0\|_{H^s_\sigma}^2 \end{aligned}$$

and

$$\|z_{\pm\infty} - e^{-JtH(V_{\pm\infty})} e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} z(t)\|_{H^{s}_{3/2}} \leq \frac{C}{\langle t \rangle^{\frac{1}{2}}} \|z_0\|_{H^s_\sigma}^2.$$

Proof. Let

$$t \mapsto \xi_{\pm}(t) = e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} z(t)$$

and

$$t \mapsto V_{\pm}(t) = e^{-i \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} U(t).$$

Using exactly the same method as that of Lemma 3.8, applied to

$$\begin{aligned} & e^{-JtH(V_{\pm\infty})} \xi_{\pm}(t) = z_0 \\ & + \int_0^t e^{-JsH(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) J (d^2F(S(V_{\pm}(v))) - d^2F(S(V_{\pm\infty}))) \xi_{\pm}(v) dv \\ & + \int_0^t e^{-JsH(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) JN(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) dv \\ & + \int_0^t e^{-JsH(V_{\pm\infty})} \mathbf{P}_c(V_{\pm}(v)) dS(V(v)) A(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) \\ & \quad \times \langle N(V_{\pm}(v), \tilde{\eta}_{\pm}(v)), dS(V_{\pm}(v)) \rangle dv \\ & - \int_0^t e^{-JsH(V_{\pm\infty})} (d\mathbf{P}_c(V_{\pm}(v))) A(V_{\pm}(v), \tilde{\eta}_{\pm}(v)) \\ & \quad \times \langle N(V_{\pm}(v), \tilde{\eta}_{\pm}(v)), dS(V_{\pm}(v)) \rangle \tilde{\eta}_{\pm}(v) dv, \end{aligned}$$

with $\tilde{\eta}_{\pm}(t) = e^{J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} (z(t) + g(U(t), z(t)))$, we prove that the limits

$$\lim_{t \rightarrow \pm\infty} e^{-JtH(V_{\pm\infty})} \xi_{\pm}(t) = z_{\pm\infty}$$

exist. If we use the method of Lemma 3.13, we obtain the estimates on the convergence of $e^{JtH(V_{\pm\infty})} z_{\pm\infty} - \xi_{\pm}(t)$. Then multiplying by $e^{-J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv}$, we obtain the estimates and the convergence of

$$e^{-J \int_0^t (E(U(v)) - E(U_{\pm\infty})) dv} e^{-JtH(V_{\pm\infty})} z_{\pm\infty} - z(t).$$

Then since $(1 - P_c(U(t))) z(t) = 0$, we have $(1 - P_c(V_{\pm\infty})) z_{\pm\infty} = 0$, and hence $z_{\pm\infty}$ belongs to $\mathcal{H}_c(V_{\pm\infty})$. \square

4. The dynamics outside the center manifold. We can make the same study in the center-stable manifold and the center-unstable manifold but only in one direction of time. Let us explain it for the center-stable manifold in positive time since it is similar for the center-unstable manifold. Actually it is equivalent if we revert the time direction.

We give just a sketch of the proof since it is similar to the previous study. Using the idea of the proof of exponential stabilization for Proposition 2.4, we write any solution ψ in the form $\phi + \rho + f(\phi, \rho)$ with ϕ in the center manifold, $\rho \in X_s(V_0)$, and f a function to be specified and which ensures that ϕ is in the center stable manifold.

Indeed, $W^c(V_0)$ is the graph of a smooth function $h^c : X_c(V_0) \mapsto X_s(V_0) \oplus X_u(V_0)$, and $W^{cs}(V_0)$ is the graph of a smooth function $h^u : X_c(V_0) \oplus X_s(V_0) \mapsto X_u(V_0)$. Let ν be such that $\psi = S(V_0) + \nu$ satisfies (1.8); then we have

$$\begin{aligned} \partial_t \nu &= JH(V_0)\nu + JN(V_0, \nu), \\ \nu &= y + h_c(y) + \rho + h^u(y, h^c(y) + \rho) \\ &= \phi(y) - S(V_0) + (\rho - \pi^s(V_0)h^c(y)) + (h^u(y, h^c(y) + \rho) - \pi^u(V_0)h^c(y)) \\ &= \phi(y) - S(V_0) + \rho + f(y, \rho) \end{aligned}$$

with $y = \pi^c(V_0)\nu = \pi^c(V_0)(\psi - S(V_0))$, $\phi(y) = S(V_0) + y + h^c(y)$ in the center manifold, and $\rho \in X_s(V_0)$. We have the following equation for ρ :

$$(4.1) \quad \partial_t \rho = JH(V_0)\rho + M(V_0, y, \rho),$$

where

$$\begin{aligned} M(V_0, y, \rho) &= \pi^s(V_0) \{JN(V_0, y + h^c(y) + \rho + f(y, \rho)) - JN(V_0, y + h^c(y))\} \\ &\quad - \pi^s(V_0) Dh^c(y) \pi^c(V_0) \{JN(V_0, y + h^c(y) + \rho + f(y, \rho)) - JN(V_0, y + h^c(y))\}. \end{aligned}$$

Then we obtain for ϕ the equation

$$\begin{aligned} \partial_t \phi &= JH\phi + J\nabla F(\phi) + R(\phi, \rho), \\ R(\phi, \rho) &= J\nabla F(\phi + \rho + f(y, \rho)) - J\nabla F(\phi) \\ &\quad - Jd^2 F(S(V_0))\rho - M(V_0, \pi^c(V_0)(\phi - S(V_0)), \rho), \end{aligned}$$

noting that $|R(\phi, \rho)| \leq C(\|\phi\|_{H^s}, \|\rho\|_{L^\infty})|\rho|$.

Working as in section 3, we write $\phi = S(U) + \eta$ with $\eta = z + g(U, z)$ and we have the following equations for U and z :

$$\begin{cases} \dot{U} = -A(U, \eta)\langle N(U, \eta) - JR(U, \eta, \rho), dS(U) \rangle, \\ \partial_t z = JH(U)z + \mathbf{P}_c(U)JN(U, \eta) \\ + \mathbf{P}_c(U(v))dS(U(v))A(U(v), \eta(v))\langle N(U(v), \eta(v)) - JR(U, \eta, \rho), dS(U(v)) \rangle \\ + (dP_c(U))A(U, \eta)\langle N(U, \eta) - JR(U, \eta, \rho), dS(U) \rangle\eta + \mathbf{P}_c(U)R(U, \eta, \rho) \end{cases}$$

with

$$\eta(t) = z(t) + g(U(t), z(t)),$$

where g is defined by Lemma 3.2 and

$$R(U, \eta, \rho) = R(S(U) + \eta, \rho).$$

These equations are similar to those we have studied but with an extra term coming from R which is exponentially decaying in positive time. Indeed, let us introduce for any $T_0 < 0$, $\gamma \in (0, \gamma(V_0))$, and $\delta > 0$ the set

$$\mathcal{R}_{T_0, \gamma}(\delta) = \{ \rho \in \mathcal{C}((T_0, +\infty), X_s(V_0)), |\rho(t)|_{H^s} \leq \delta e^{-\gamma t} \forall t > T_0 \}.$$

We study (4.1) in $\mathcal{R}_{T_0, \gamma}(\delta)$ with small initial condition ρ_0 . We also define for any $\varepsilon > 0$

$$\mathcal{U}_{T_0}(\varepsilon, \delta) = \left\{ U \in \mathcal{C}^1((T_0, +\infty), B_{\mathcal{C}^2}(V_0, \varepsilon)), \|\dot{U}\|_{L^1((T_0, +\infty)) \cap L^\infty((T_0, +\infty))} \leq \delta^2 \right\},$$

and for any $U \in \mathcal{U}_{T_0}(\varepsilon)$, let s, β be such that $s > \beta + 2 > 2$ and $\sigma > 3/2$,

$$\mathcal{Z}_{T_0}(U, \delta) = \{ z \in \mathcal{C}((T_0, +\infty), L^2(\mathbb{R}^3, \mathbb{R}^8)), z(t) \in \mathcal{H}_c(U(t)),$$

$$\max [\|z\|_{L^\infty((T_0, +\infty), H^s)}, \|z\|_{L^2((T_0, +\infty), H^s_{-\sigma})}, \|z\|_{L^2((T_0, +\infty), B^\beta_{\infty, 2})}] \leq \delta \},$$

and ε, δ are small enough to ensure that for $U \in \mathcal{U}(\varepsilon, \delta)$ and $z \in \mathcal{Z}(U, \delta)$

$$S(U) + z + g(U, z) \in W^c(V_0) \cap B_{H^s}(S(V_0), r(V_0)).$$

For a sufficiently small T_0 , we solve the equation for z first, then the one for ρ , and eventually the one for U using the method of section 3. This gives us the desired exponential decay for ρ as well as similar results for U and z .

We note that instead of Lemma 3.9, we obtain the following lemma.

LEMMA 4.1. *If Assumptions 1.1–1.5 hold, then for any $T > 0$, there exist $T_0 > 0$, $\varepsilon_0 > 0$, $\delta_0 > 0$, $C > 0$, and $\kappa \in (0, 1)$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U, U' \in \mathcal{U}_{T_0}(\varepsilon, \delta)$, $\rho, \rho' \in \mathcal{R}_{T_0, \gamma}$, $z_0 \in \mathcal{H}_c(U(0))$, $z'_0 \in \mathcal{H}_c(U'(0))$, $z \in \mathcal{Z}_{T_0}(U, \delta)$, and $z' \in \mathcal{Z}_{T_0}(U', \delta)$, one has*

$$\begin{aligned} & \|z[z'_0, U', \rho'] - z[z_0, U, \rho]\|_{L^\infty((T_0, T), H^s) \cap L^2((T_0, T), L^\infty) \cap L^2((T_0, T), H^s_{-\sigma})} \\ & \leq C \|z_0 - z'_0\|_{H^s} + \kappa \{ \|U - U'\|_{L^\infty((T_0, T), \mathbb{C}^2)} + \|\dot{U} - \dot{U}'\|_{L^\infty((T_0, T), \mathbb{C}^2)} \\ & \quad + \|e^{\gamma t}(\rho - \rho')(t)\|_{L^\infty_t((T_0, T), X^s(V_0))} \}. \end{aligned}$$

Then for ρ as a function of U , z_0 , and ρ_0 (the initial condition for ρ), we obtain the following lemma.

LEMMA 4.2. *If Assumptions 1.1–1.5 hold, then for any $T > 0$ there exist $T_0 > 0$, $\varepsilon_0 > 0$, $\delta_0 > 0$, $C > 0$, and $\kappa \in (0, 1)$ such that for any $\varepsilon \in (0, \varepsilon_0)$, $\delta \in (0, \delta_0)$, $U, U' \in \mathcal{U}_{T_0}(\varepsilon, \delta)$, $r_0, r'_0 \in X_s(V_0)$, $z_0 \in \mathcal{H}_c(U(0))$, $z'_0 \in \mathcal{H}_c(U'(0))$, $z \in \mathcal{Z}_{T_0}(U, \delta)$, and $z' \in \mathcal{Z}_{T_0}(U', \delta)$ one has*

$$\begin{aligned} & \|e^{\gamma t} (\rho[z'_0, U', \rho'_0] - \rho[z_0, U, \rho_0])\|_{L^\infty((T_0, T), X^s(V_0))} \\ & \leq C \|z_0 - z'_0\|_{H^s} + \kappa \{ \|U - U'\|_{L^\infty((T_0, T), \mathbb{C}^2)} + \|\dot{U} - \dot{U}'\|_{L^\infty((T_0, T), \mathbb{C}^2)} + \|\rho_0 - \rho'_0\|_{L^2} \}. \end{aligned}$$

We also note that the proof gives the well-posedness of (4.1) in $\mathcal{R}_{T_0, \gamma}(\delta)$ with small initial condition ρ_0 and that there exists $C > 0$ such that the solution ρ satisfies

$$\|\rho(t)\|_{H^s} \leq C \|\rho_\pm(0)\| e^{-\gamma t} \quad \forall t > T_0.$$

The asymptotic behaviors of U and z are obtained as in the previous section when z_0 is localized.

5. End of the proof of main theorems. We note that the small *locally invariant* center manifold built in section 2.2 for (2.10) is now a small *invariant* (globally in time) center manifold. Indeed, we have just proved the stabilization towards the PLS manifold; this ensures that a solution in the center manifold will stay inside this manifold in the two directions of time.

Now let us consider \mathcal{CM} as being the union of all these small globally invariant center manifolds and 0. Using the uniqueness of the center manifold and Lemma 3.2, we prove that $\mathcal{CM} \setminus \{0\}$ is a manifold. Now we generalize Lemma 3.2 by the following lemma.

LEMMA 5.1. *For any $s, s', \sigma \in \mathbb{R}$ and $p, q \in [1, \infty]$, there exist $\varepsilon > 0$, a continuous map $r : B_{\mathbb{C}}^2(0, \varepsilon) \mapsto \mathbb{R}^+$ with $r(U) = O(\Gamma(U))$, and a continuous map $\Psi : S \mapsto \mathcal{CM}$, where*

$$S_\sigma = \{(U, z); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H_\sigma^{s'}}(0, r(U))\}$$

is endowed with the metric of $\mathbb{C}^2 \times H_\sigma^{s'}$.

Moreover, Ψ is bijective from S to an open neighborhood of $(0, 0)$ in \mathcal{CM} and smooth on $S \setminus \{(0, 0)\}$. For all $U \in B_{\mathbb{C}^2}(0, \varepsilon)$, there exists $C > 0$ such that, for all $z \in \mathcal{H}_c(U) \cap B_{H_\sigma^{s'}}(0, r(U))$, $\Psi(U, z) \in \mathcal{H}_1(U)$, $z + \Psi(U, z) \in \mathcal{H}_0(U)^\perp$, and $S(U) + z + \Psi(U, z) \in \mathcal{CM}$. For sufficiently small nonzero U , we have $\|\Psi(U, z)\|_{B_{p,q}^s} = O(\|z\|_{H_\sigma^{s'}}^2)$ for $z \in H_\sigma^{s'}$ such that $(U, z) \in S$.

Proof. The proof works like that of Lemma 3.2. The statements for r follow from Remark 2.3. \square

The scattering result follows from a one-to-one correspondence of the initial profile with the asymptotic profile as stated in the following proposition.

PROPOSITION 5.1. *If Assumptions 1.1–1.5 hold, there exist $\varepsilon > 0$ and a continuous map $r : B_{\mathbb{C}}^2(0, \varepsilon) \mapsto \mathbb{R}^+$ with $r(U) = O(\Gamma(U))$ and $\mathcal{V}_\sigma, \mathcal{V}_\pm$ neighborhoods of $(0, 0)$ in*

$$S_\sigma = \{(U, z); U \in \mathbb{C}^2, z \in \mathcal{H}_c(U) \cap B_{H_\sigma^s}(0, r(U))\}$$

endowed with the norm of $\mathbb{C}^2 \times H_\sigma^s$ such that the maps

$$\mathcal{P}_\pm : \begin{pmatrix} U_0 \\ z_0 \end{pmatrix} \in \mathcal{V}_\sigma \mapsto \begin{pmatrix} V_{\pm\infty} \\ z_{\pm\infty} \end{pmatrix} \in \mathcal{V}_\pm$$

are bijections and are smooth on $\mathcal{V}_0 \setminus \{(0, 0)\}$.

Proof. We choose, for example,

$$\mathcal{V}_\sigma = \{(U, z); U \in B_{\mathbb{C}^2}(0, \varepsilon), z \in \mathcal{H}_c(U) \cap B_{H_\sigma^s}(0, r(U))\}$$

for some positive ε , and we work on the manifold $\mathcal{V}_\sigma \setminus \{(0, 0)\}$ which is locally isomorphic to an open set of $\mathbb{C}^2 \times \mathcal{H}_c(U) \cap H_\sigma^s$. We write

$$\mathcal{P}_\pm^{U_0}(U, z) = (U, z) + \mathcal{R}_\pm^{U_0}(U, z).$$

Since

$$\|(U_\infty, z_\infty) - (U_0, z_0)\|_{H_\sigma^s} = O(|U_0|^2 + \|z_0\|_{H_\sigma^s}^2),$$

we need only prove the statement locally. Hence we prove that, in a neighborhood of $(U_0, 0)$, the maps $\mathcal{P}_\pm^{U_0}(U, z) \mapsto (Id_{\mathbb{C}^2}, P_c(U_0))\mathcal{P}_\pm(U, R(U, U_0)z)$ are bijective (P_c and R are defined in Proposition 2.2).

To prove that $\mathcal{P}_\pm^{U_0}$ is bijective (i.e., the scattering exists), let us prove it for $\mathcal{P}_+^{U_0}$ (it is similar for $\mathcal{P}_-^{U_0}$). It is enough to prove that the following system has a unique solution in an open neighborhood of $(0, 0)$ in \mathcal{S}_σ :

$$\begin{aligned} V_\pm(t) &= V_{\pm\infty} \\ &+ \int_t^\infty A(V_\pm(v), e^{JsH(V_{\pm\infty})}\tilde{\eta}_\pm(v)) \langle N(U(v), e^{JsH(V_{\pm\infty})}\tilde{\eta}_\pm(v)), dS(V_\pm(v)) \rangle dv \end{aligned}$$

and

$$\begin{aligned} \tilde{\xi}_+(t) &= z_\infty - \int_t^\infty e^{-JsH(V_{+\infty})} \\ &\quad \times \mathbf{P}_c(V_+(v))J(d^2F(S(V_+(v))) - d^2F(S(V_{+\infty}))) e^{JsH(V_{+\infty})}\tilde{\xi}_+(v) dv \\ &- \int_t^\infty e^{-JsH(V_{+\infty})} \\ &\quad \times \mathbf{P}_c(V_+(v))JN(V_+(v), e^{JsH(V_{+\infty})}\tilde{\eta}_+(v)) dv \\ &- \int_t^\infty e^{-JsH(V_{+\infty})}\mathbf{P}_c(V_+(v))dS(V(v))A(V_+(v), e^{JsH(V_{+\infty})}\tilde{\eta}_+(v)) \\ &\quad \times \langle N(V_+(v), e^{JsH(V_{+\infty})}\tilde{\eta}_+(v)), dS(V_+(v)) \rangle dv \\ &+ \int_t^\infty e^{-JsH(V_{+\infty})}(d\mathbf{P}_c(V_+(v)))A(V_+(v), e^{JsH(V_{+\infty})}\tilde{\eta}_+(v)) \\ &\quad \times \langle N(V_+(v), e^{JsH(V_{+\infty})}\tilde{\eta}_+(v)), dS(V_+(v)) \rangle e^{JsH(V_{+\infty})}\tilde{\eta}_+(v) dv \end{aligned}$$

with $\tilde{\eta}_+(t) = \tilde{\xi}_+(t) + e^{-JsH(V_{+\infty})}g(V_+(t), e^{JsH(V_{+\infty})}\tilde{\xi}_+(t))$.

This system can be solved by a fixed point argument in the set of functions such that

$$\begin{aligned} &\max \left[\sup_{t \in \mathbb{R}} (\|z_{+\infty} - \tilde{\xi}_+(t)\|_{H^s}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z_{+\infty} - \tilde{\xi}_+(t)\|_{H_{-\sigma}^s}), \right. \\ &\left. \sup_{t \in \mathbb{R}} (\langle t \rangle^{3/2} \|z_{+\infty} - \tilde{\xi}_+(t)\|_{B_{\infty,2}^\beta}), \sup_{t \in \mathbb{R}} (\langle t \rangle^{-3/2} \|z_{+\infty} - \tilde{\xi}_+(t)\|_{H_{3/2}^s}) \right] \end{aligned}$$

and

$$\langle t \rangle^2 |V_+(t) - V_{+\infty}|$$

are small with the method we used in Lemma 3.14. \square

For the same reasons, the small locally invariant center-stable manifold built in section 2.2 is invariant in positive time. We can also consider the union of these manifolds, and we can obtain a map Φ_+ similar to the map Ψ built in Lemma 5.1.

The instability in negative time is, in fact, a consequence of Proposition 2.5.

The corresponding conclusion holds for the center-unstable manifold.

The statements on the instability outside these manifolds follow from Propositions 2.4 and 2.5.

Acknowledgments. I would like to thank Éric Séré for fruitful discussions and advice during the preparation of this work. I am indebted to Galina Perel'man for her careful reading and her comments and suggestions, to Michael Levitin for his advice, to Sylvain Golénia for having checked some proofs, to Matthieu Brassart for his answers to my questions, and to Mathieu Lewin for all his help. I am grateful to the referees for their careful reading, corrections, and remarks.

REFERENCES

- [1] W. O. AMREIN, A. BOUTET DE MONVEL, AND V. GEORGESCU, *C₀-Groups, Commutator Methods and Spectral Theory of N-Body Hamiltonians*, Progr. Math. 135, Birkhäuser Verlag, Basel, 1996.
- [2] E. BALSLEV AND B. HELFFER, *Limiting absorption principle and resonances for the Dirac operator*, Adv. in Appl. Math., 13 (1992), pp. 186–215.
- [3] P. W. BATES AND C. K. R. T. JONES, *Invariant manifolds for semilinear partial differential equations*, in Dynamics Reported, Vol. 2, Dynam. Report. Ser. Dynam. Systems Appl. 2, Wiley, Chichester, UK, 1989, pp. 1–38.
- [4] J. BERGH AND J. LÖFSTRÖM, *Interpolation Spaces. An Introduction*, Grundlehren Math. Wiss. 223, Springer-Verlag, Berlin, 1976.
- [5] A. BERTHIER AND V. GEORGESCU, *On the point spectrum of Dirac operators*, J. Funct. Anal., 71 (1987), pp. 309–338.
- [6] N. BOUSSAID, *Étude de la stabilité des petites solution stationnaires pour une classe d'équations de Dirac non linéaires*, Ph.D. thesis, Université Paris-Dauphine, Paris, 2006.
- [7] N. BOUSSAID, *Stable directions for small nonlinear Dirac standing waves*, Comm. Math. Phys., 268 (2006), pp. 757–817.
- [8] A. BRESSAN, *A tutorial on the center manifold theorem*, in Proceedings of the C.I.M.E. Course, Cetraro, Italy, 2003, Lecture Notes in Math. 1911, P. Marcati, ed., Springer-Verlag, New York, 2007, pp. 327–344.
- [9] V. S. BUSLAEV AND G. S. PEREL'MAN, *Nonlinear scattering: States that are close to a soliton*, Zap. Nauchn. Sem. S.-Peterburg. Otdel. Mat. Inst. Steklov. (POMI), 200 (1992), pp. 38–50, 70, 187.
- [10] V. S. BUSLAEV AND G. S. PEREL'MAN, *On nonlinear scattering of states which are close to a soliton*, Astérisque, no. 210 (1992), pp. 6, 49–63.
- [11] V. S. BUSLAEV AND G. S. PEREL'MAN, *Scattering for the nonlinear Schrödinger equation: States that are close to a soliton*, Algebra i Analiz, 4 (1992), pp. 63–102.
- [12] V. S. BUSLAEV AND G. S. PEREL'MAN, *On the stability of solitary waves for nonlinear Schrödinger equations*, in Nonlinear Evolution Equations, Amer. Math. Soc. Transl. Ser. 2 164, AMS, Providence, RI, 1995, pp. 75–98.
- [13] V. S. BUSLAEV AND C. SULEM, *Asymptotic stability of solitary waves for nonlinear Schrödinger equations*, in The Legacy of the Inverse Scattering Transform in Applied Mathematics (South Hadley, MA, 2001), Contemp. Math. 301, AMS, Providence, RI, 2002, pp. 163–181.
- [14] V. S. BUSLAEV AND C. SULEM, *On asymptotic stability of solitary waves for nonlinear Schrödinger equations*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 20 (2003), pp. 419–475.
- [15] T. CAZENAVE AND P.-L. LIONS, *Orbital stability of standing waves for some nonlinear Schrödinger equations*, Comm. Math. Phys., 85 (1982), pp. 549–561.

- [16] S. CUCCAGNA, *Stabilization of solutions to nonlinear Schrödinger equations*, Comm. Pure Appl. Math., 54 (2001), pp. 1110–1145.
- [17] S. CUCCAGNA, *On asymptotic stability of ground states of NLS*, Rev. Math. Phys., 15 (2003), pp. 877–903.
- [18] S. CUCCAGNA, *Erratum: “Stabilization of solutions to nonlinear Schrödinger equations”* [Comm. Pure Appl. Math., 54 (2001), pp. 1110–1145, MR1835384], Comm. Pure Appl. Math., 58 (2005), p. 147.
- [19] M. ESCOBEDO AND L. VEGA, *A semilinear Dirac equation in $H^s(\mathbf{R}^3)$ for $s > 1$* , SIAM J. Math. Anal., 28 (1997), pp. 338–362.
- [20] V. GEORGESCU AND M. MĂNTOIU, *On the spectral theory of singular Dirac type Hamiltonians*, J. Operator Theory, 46 (2001), pp. 289–321.
- [21] M. GRILLAKIS, J. SHATAH, AND W. STRAUSS, *Stability theory of solitary waves in the presence of symmetry. I*, J. Funct. Anal., 74 (1987), pp. 160–197.
- [22] S. GUSTAFSON, K. NAKANISHI, AND T.-P. TSAI, *Asymptotic stability and completeness in the energy space for nonlinear Schrödinger equations with small solitary waves*, Int. Math. Res. Not., no. 66 (2004), pp. 3559–3584.
- [23] P. D. HISLOP, *Exponential decay of two-body eigenfunctions: A review*, in Proceedings of the Symposium on Mathematical Physics and Quantum Field Theory (Berkeley, CA, 1999), Electron. J. Differ. Equ. Conf. 4, Southwest Texas State University, San Marcos, TX, 2000, pp. 265–288.
- [24] W. HUNZIKER AND I. M. SIGAL, *Time-dependent scattering theory of N -body quantum systems*, Rev. Math. Phys., 12 (2000), pp. 1033–1084.
- [25] A. IFTIMOVICI AND M. MĂNTOIU, *Limiting absorption principle at critical values for the Dirac operator*, Lett. Math. Phys., 49 (1999), pp. 235–243.
- [26] A. JENSEN AND T. KATO, *Spectral properties of Schrödinger operators and time-decay of the wave functions*, Duke Math. J., 46 (1979), pp. 583–611.
- [27] T. KATO, *Wave operators and similarity for some non-selfadjoint operators*, Math. Ann., 162 (1965/1966), pp. 258–279.
- [28] T. KATO, *Perturbation Theory for Linear Operators*, Classics Math., Springer-Verlag, Berlin, 1995 (reprint of the 1980 edition).
- [29] M. KEEL AND T. TAO, *Endpoint Strichartz estimates*, Amer. J. Math., 120 (1998), pp. 955–980.
- [30] E. KIRR AND A. ZARNESCU, *On the asymptotic stability of bound states in 2D cubic Schrödinger equation*, Comm. Math. Phys., 272 (2007), pp. 443–468.
- [31] J. KRIEGER AND W. SCHLAG, *On the focusing critical semi-linear wave equation*, Amer. J. Math., 129 (2007), pp. 843–913.
- [32] S. MACHIHARA, M. NAKAMURA, K. NAKANISHI, AND T. OZAWA, *Endpoint Strichartz estimates and global solutions for the nonlinear Dirac equation*, J. Funct. Anal., 219 (2005), pp. 1–20.
- [33] S. MACHIHARA, M. NAKAMURA, AND T. OZAWA, *Small global solutions for nonlinear Dirac equations*, Differential Integral Equations, 17 (2004), pp. 623–636.
- [34] S. MACHIHARA, K. NAKANISHI, AND T. OZAWA, *Small global solutions and the nonrelativistic limit for the nonlinear Dirac equation*, Rev. Mat. Iberoamericana, 19 (2003), pp. 179–194.
- [35] T. MIZUMACHI, *Instability of bound states for 2D nonlinear Schrödinger equations*, Discrete Contin. Dyn. Syst., 13 (2005), pp. 413–428.
- [36] T. MIZUMACHI, *A remark on linearly unstable standing wave solutions to NLS*, Nonlinear Anal., 64 (2006), pp. 657–676.
- [37] B. PARISSÉ, *Résonances paires pour l’opérateur de Dirac*, C. R. Acad. Sci. Paris Sér. I Math., 310 (1990), pp. 265–268.
- [38] C.-A. PILLET AND C. E. WAYNE, *Invariant manifolds for a class of dispersive, Hamiltonian, partial differential equations*, J. Differential Equations, 141 (1997), pp. 310–326.
- [39] A. F. RAÑADA, *Classical nonlinear Dirac field models of extended particles*, in Quantum Theory, Groups, Fields and Particles, Vol. 198, A. O. Barut, ed., Reidel, Amsterdam, pp. 271–291.
- [40] M. REED AND B. SIMON, *Methods of Modern Mathematical Physics. IV. Analysis of Operators*, Academic Press [Harcourt Brace Jovanovich], New York, 1978.
- [41] W. SCHLAG, *Stable Manifolds for an Orbitally Unstable NLS*, preprint, University of Chicago, Chicago, IL, 2004.
- [42] J. SHATAH AND W. STRAUSS, *Instability of nonlinear bound states*, Comm. Math. Phys., 100 (1985), pp. 173–190.
- [43] A. SOFFER AND M. I. WEINSTEIN, *Multichannel nonlinear scattering for nonintegrable equations*, Comm. Math. Phys., 133 (1990), pp. 119–146.
- [44] A. SOFFER AND M. I. WEINSTEIN, *Multichannel nonlinear scattering for nonintegrable equations. II. The case of anisotropic potentials and data*, J. Differential Equations, 98 (1992), pp. 376–390.

- [45] A. SOFFER AND M. I. WEINSTEIN, *Selection of the ground state for nonlinear Schrödinger equations*, Rev. Math. Phys., 16 (2004), pp. 977–1071.
- [46] A. SOFFER AND M. I. WEINSTEIN, *Theory of nonlinear dispersive waves and selection of the ground state*, Phys. Rev. Lett., 95 (2005), 213905.
- [47] N. SZPAK, *Spontaneous Particle Creation in Time-Dependent Overcritical Fields of QED*, Ph.D. thesis, J. W. Goethe University, Frankfurt am Main, Germany, 2005.
- [48] B. THALLER, *The Dirac Equation*, Texts Monogr. Phys., Springer-Verlag, Berlin, 1992.
- [49] H. TRIEBEL, *Interpolation Theory, Function Spaces, Differential Operators*, North-Holland Math. Library 18, North-Holland, Amsterdam, 1978.
- [50] T.-P. TSAI, *Asymptotic dynamics of nonlinear Schrödinger equations with many bound states*, J. Differential Equations, 192 (2003), pp. 225–282.
- [51] T.-P. TSAI AND H.-T. YAU, *Asymptotic dynamics of nonlinear Schrödinger equations: Resonance-dominated and dispersion-dominated solutions*, Comm. Pure Appl. Math., 55 (2002), pp. 153–216.
- [52] T.-P. TSAI AND H.-T. YAU, *Classification of asymptotic profiles for nonlinear Schrödinger equations with small initial data*, Adv. Theor. Math. Phys., 6 (2002), pp. 107–139.
- [53] T.-P. TSAI AND H.-T. YAU, *Relaxation of excited states in nonlinear Schrödinger equations*, Int. Math. Res. Not., no. 31 (2002), pp. 1629–1673.
- [54] T.-P. TSAI AND H.-T. YAU, *Stable directions for excited states of nonlinear Schrödinger equations*, Comm. Partial Differential Equations, 27 (2002), pp. 2363–2402.
- [55] R. WEDER, *Center manifold for nonintegrable nonlinear Schrödinger equations on the line*, Comm. Math. Phys., 215 (2000), pp. 343–356.
- [56] M. I. WEINSTEIN, *Modulational stability of ground states of nonlinear Schrödinger equations*, SIAM J. Math. Anal., 16 (1985), pp. 472–491.
- [57] M. I. WEINSTEIN, *Lyapunov stability of ground states of nonlinear dispersive evolution equations*, Comm. Pure Appl. Math., 39 (1986), pp. 51–67.

ANALYSIS OF MODEL EQUATIONS FOR STRESS-ENHANCED DIFFUSION IN COAL LAYERS. PART I: EXISTENCE OF A WEAK SOLUTION*

ANDRO MIKELIĆ[†] AND HANS BRUINING[‡]

Abstract. This paper is motivated by the study of the sorption processes in the coal. They are modeled by a nonlinear degenerate pseudoparabolic equation for stress-enhanced diffusion of carbon dioxide (CO₂) in coal, $\partial_t \phi = \partial_x \left\{ D(\phi) \partial_x \phi + \frac{D(\phi)\phi}{B} \partial_x (e^{-m\phi} \partial_t \phi) \right\}$, where B, m are positive constants and the diffusion coefficient $D(\phi)$ has a small value when the CO₂ volume fraction ϕ is $0 \leq \phi < \phi_c$, representative of coal in the glass state and orders of magnitude higher value for $\phi > \phi_c$, when coal is in the rubber-like state. These types of equations arise in a number of cases when nonequilibrium thermodynamics or extended nonequilibrium thermodynamics is used to compute the flux. For this equation, existence of the travelling wave-type solutions was extensively studied. Nevertheless, the existence seems to be known only for a sufficiently short time. We use the corresponding entropy functional in order to get existence, for any time interval, of an appropriate weak solution with square integrable first derivatives and satisfying uniform L^∞ -bounds. Due to the degeneracy, we obtain square integrability of the mixed second order derivative only in the region where the concentration ϕ is strictly positive. In obtaining the existence result it was crucial to have the regularized entropy as unknown for the approximate problem and not the original unknown (the concentration).

Key words. degenerate pseudoparabolic equation, entropy methods, stress-enhanced diffusion

AMS subject classifications. 35K70, 35K65, 76R50, 80A17

DOI. 10.1137/070710172

1. Introduction. One of the promising methods for reducing the discharge of the “greenhouse gas” carbon dioxide (CO₂) into the atmosphere is its sequestration in unminable coal seams. A typical procedure is the injection of CO₂ via deviated wells drilled inside the coal seams. CO₂ displaces the methane adsorbed on the internal surface of the coal. A production well gathers the methane as free gas. This process, known as CO₂-enhanced coal bed methane production (CO₂-ECBM), is a producer of energy and at the same time reduces greenhouse concentrations as about two CO₂ molecules displace one molecule of methane. Worldwide application of ECBM can reduce greenhouse gas emissions by a few percent. Coal has an extensive fracturing system called the cleat system. In fact, it is possible to discern a number of cleat

*Received by the editors December 5, 2007; accepted for publication (in revised form) July 8, 2008; published electronically November 26, 2008.

<http://www.siam.org/journals/sima/40-4/71017.html>

[†]Université de Lyon, Lyon, F-69003, France, and Université Lyon 1, Institut Camille Jordan, UFR Mathématiques, Site de Gerland, Bât. A, 50, avenue Tony Garnier, 69367 Lyon Cedex 07, France (mikelic@univ-lyon1.fr). This author’s research was partially supported by the Groupement MOMAS (Modélisation mathématique et simulations numériques liées aux problèmes de gestion des déchets nucléaires) (PACEN/CNRS, ANDRA, BRGM, CEA, EDF, IRSN) as a part of the project “Modèles de dispersion efficace pour des problèmes de Chimie-Transport: Changement d’échelle dans la modélisation du transport réactif en milieux poreux, en présence des nombres caractéristiques dominants.” It was initiated during this author’s sabbatical visit to the TU Eindhoven in 2006, supported by Visitors grant B-61-602 of the Netherlands Organisation for Scientific Research (NWO).

[‡]Dietz Laboratory, Geo-Environmental Engineering, Stevinweg 1, 2628 CN Delft, The Netherlands (j.bruining@tudelft.nl). This author’s research was supported by Shell International Exploration and Production B.V., NWO-Novem, and the Greenhouse-gas Removal Apprenticeship and Student Program (GRASP) funded by the EC. The idea of stress-enhanced diffusion [29], [30] for coal, worked out in the doctoral thesis of Mazumder [23], was the inspiration for this work.

systems at different scales. In the end, the matrix blocks between the smallest cleat systems typically have diameters of a few tens of microns [13].

The matrix blocks have a polymeric structure (dehydrated cellulose [32]), which provides the adsorption sites for the gases. At low temperatures or low sorption concentration the coal structure behaves like a rigid glassy polymer, in which movement is difficult. At high temperatures or high sorption concentrations, the glassy structure is converted to the less rigid and open rubber-like (swollen) structure [29], [30]. As coal is less dense in the rubber-like state, a conversion from the glassy state to the rubber-like state exhibits swelling. Therefore modeling of diffusion is not only relevant for modeling transport into the matrix blocks, but also for the modeling of swelling, which affects the permeability of the coal seam.

Ritger and Peppas [29], [30] distinguish between transport by Fickian diffusion and a process that occurs on the interface between the glass state and the rubber-like state. Ritger and Peppas state that the conversion process from the glass to the rubber state is controlled by a rate-limiting relaxation phenomenon (see also [2]). Thomas and Windle [31] (see also [16], [17], [19]), however, suggested in their classic paper that the diffusion transport was enhanced by stress gradients that resulted from the accommodation of large molecules in the small cavities providing the adsorption sites. For this, Alfrey, Gurnee, and Lloyd [1] coined the term superdiffusion or case II diffusion. At a critical concentration of the penetrants the glassy polymer is transformed to a rubber state, where the diffusion coefficient is of the order of a factor 1000 larger than in the glassy state.

This paper is the first of a series in which the model equations for case II diffusion [31], [16], [17], [19] will be analyzed. Our longtime interest is to investigate the one-dimensional sorption rate behavior, i.e., whether the equations indeed lead to a rate faster than the square root of time. In this paper, we establish existence of a weak solution for all times.

Nonlinear diffusion equations with a pseudoparabolic regularizing term being the Laplacian of the time derivative are considered in [25] and [26]. Global existence of a strong solution is proved by writing the problem as a linear elliptic operator, acting on the time derivative, equal to the nonlinear diffusion term. Then the linear elliptic operator, acting on the time derivative, is inverted, and the standard geometric theory of nonlinear parabolic equations (see, e.g., [15]) is applicable.

In our situation the physical model leads to a degenerate nonlinear second order elliptic operator, acting on the time derivative, in place of the Laplacian. The invertibility of this nonlinear elliptic operator is not clear anymore and depends on the solution itself. The same type of equation can occur in models that use classical irreversible thermodynamics (CIT) or extended irreversible thermodynamics (EIT). An important example is the model of the two-phase flow through porous media introduced in [14], where the capillary pressure relation is extended with a dynamic term, which contains the time derivative of the saturation. We also refer the reader to [5] for the modeling. This application to multiphase and unsaturated flows through porous media motivated a number of recent papers. In [18] one finds a detailed study of possible travelling wave solutions and in particular of the behavior of such travelling waves near fronts where the concentration is zero. Further studies of the travelling waves are in [8] and [7]. The small- and waiting-time behavior of the equations is studied in [20]. Study of the viscosity limit for the linear relaxation model of the dynamic term is in [10]. Nevertheless, the study of existence of a solution to the nonlinear model from [14] is undertaken only in [4] and [5], where the nondegeneracy is supposed and existence is local in time. Another existence result, also local in

time, is in the paper [9] by Düll, where a related pseudoparabolic equation modeling solvent uptake in polymeric solids is studied. Düll proved the short-time existence of a solution for the problem in \mathbb{R} , supposing nonnegative compactly supported initial datum. Contrary to our approach, the problem was written as a system containing a linear elliptic equation and an evolution equation. With such an approach, we did not manage to get estimates as good as those with the entropy approach undertaken in this paper. For studies of travelling waves and sharp fronts in case II diffusion models, we refer the reader to [17] and [33].

We consider the evolution problem

$$\begin{aligned}
 (1) \quad & \partial_t \phi = \partial_x \left\{ D(\phi) \partial_x \phi + \frac{D(\phi)\phi}{B} \partial_x (e^{-m\phi} \partial_t \phi) \right\} \quad \text{in } (0, L) \times (0, T), \\
 (2) \quad & D(\phi) \partial_x \phi + \frac{D(\phi)\phi}{B} \partial_x (e^{-m\phi} \partial_t \phi) = 0 \quad \text{on } \{x = L\} \times (0, T), \\
 (3) \quad & \phi(0, t) = \phi_g(t) \quad \text{on } (0, T), \quad \phi = 0 \quad \text{on } (0, L) \times \{0\}.
 \end{aligned}$$

Our goal is to obtain a global existence of a weak solution, for any time interval, as was obtained in [3] for a degenerate pseudoparabolic regularization of a nonlinear forward-backward heat equation. Our PDE allows a natural generalization of the classic Kullback entropy, and its integrand is given by

$$(4) \quad \mathcal{E}(\varphi) = \int_0^\varphi \frac{\varphi - \xi}{\xi} \left(e^{-m\xi} \frac{1}{D(\xi)} - \frac{1}{D(0)} \right) d\xi + \frac{1}{D(0)} (\varphi \log \varphi - \varphi).$$

As in [24], we will use $\mathcal{E}'(\varphi)$ as a test function, with the hope of obtaining a convenient a priori estimate. Formal calculation gives the equality

$$\begin{aligned}
 (5) \quad & \partial_t \int_0^L \left\{ \mathcal{E}(\phi) - \varphi \mathcal{E}'(\varphi_g) + \frac{1}{2B} (e^{-m\phi} \partial_x \phi)^2 \right\} dx \\
 & + \int_0^L \left(\frac{1}{\phi} e^{-m\phi} (\partial_x \phi)^2 + \phi \partial_t \mathcal{E}'(\varphi_g) \right) dx = 0.
 \end{aligned}$$

The presence of the initial and boundary conditions leads to unbounded nonintegrable \mathcal{E}' . The equality (5) cannot be used directly, and we do not get the entropy estimates as in [12]. We had to obtain an additional estimate for the time derivative, and our calculations are more complicated than in the literature.

Existence is proved by showing that the “energy” of the system remains bounded during the time evolution of the system. The “energy” equation is derived from the differential equation by multiplying with an appropriate test function and integrating over the domain. The choice of the test function depends strongly on the choice of the nonlinearities. With an appropriate approximation, this can also be the basis of a numerical scheme that leads to an implicit first order nonlinear system of ODEs. The implicit dependence on the time derivative makes its solvability nontrivial. Solvability of our system of ODEs depends strongly on the initial conditions. The fact that the “energy” is bounded means that the numerical scheme is stable. If convergence can be proved, it shows that at least one solution exists.

As already stated, in this case an appropriate test function is $\Phi(\phi)$, where

$$\Phi'(\xi) = \frac{e^{-m\xi}}{\xi D(\xi)},$$

is, however, singular for $\xi = 0$. Another problem of the test function is that for large values of ξ , Φ' is exponentially small. In order to prove existence we need Φ that is bijective from \mathbb{R} to \mathbb{R} . Concerning the diffusion coefficient D , we extend it by setting $D(\xi) = D(-\xi)$ for $\xi < 0$.

We introduce Φ_δ by

$$(6) \quad \Phi'_\delta := \frac{e^{-m \min\{|\xi|, 1/\delta\}}}{(|\xi| + \delta) D(\xi)}, \quad \delta > 0, \quad \xi \in \mathbb{R},$$

and

$$(7) \quad \Phi_\delta(\phi) := \begin{cases} \int_0^\phi \frac{e^{-m \min\{\xi, 1/\delta\}}}{(\xi + \delta) D(\xi)} d\xi, & \phi > 0, \\ -\int_\phi^0 \frac{e^{-m \min\{-\xi, 1/\delta\}}}{(-\xi + \delta) D(-\xi)} d\xi, & \phi < 0. \end{cases}$$

Obviously, Φ_δ is odd and strictly increasing on \mathbb{R}_+ .

In order to obtain an existence result for problem (1)–(3) we study the following regularized problem in $Q_T = (0, L) \times (0, T)$:

$$(8) \quad \partial_t \phi = \partial_x \left\{ D(\phi) \partial_x \phi + \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_t \phi \right) \right\}$$

with boundary condition at $x = L$,

$$(9) \quad D(\phi) \partial_x \phi + \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_t \phi \right) \Big|_{x=L} = 0,$$

and boundary and initial conditions (3).

We start by introducing a variational solution for the problem (8), (9), and (3).

DEFINITION 1. *Let*

$$(10) \quad \mathcal{V} := \{z \in C^\infty[0, L], z|_{x=0} = 0\} \quad \text{and} \quad \mathcal{H} := \{C^\infty[0, T], h(T) = 0\}.$$

Then the variational formulation corresponding to the problem (3), (8), and (9) is

$$(11) \quad -\int_0^T \int_0^L \phi(x, t) g(x) \partial_t h(t) dx dt + \int_0^T \int_0^L D(\phi) \partial_x \phi(x, t) \partial_x g(x) h(t) dx dt + \int_0^T \int_0^L \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x g(x) h(t) \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_t \phi \right) dx dt = 0$$

for all $g \in \mathcal{V}$ and for all $h \in \mathcal{H}$, and at the boundary $x = 0$ we have

$$(12) \quad \phi - \phi_g = 0.$$

Our goal is to prove existence for (11)–(12). In order to have the entropy estimate, we should formulate the approximate problem in terms of it. Otherwise it *would not be possible to use it as a test function* for the approximate problem, which is finite dimensional. Getting a priori estimates without this approach is not clear.

Let $z := \Phi_\delta(\phi)$, $\phi = \Phi_\delta^{-1}(z)$, $z|_{x=0} = \Phi_\delta(\phi_g(t))$. We reformulate the problem (3), (8), and (9) in terms of z :

$$(13) \quad \frac{1}{\Phi'_\delta(\Phi_\delta^{-1}(z))} \partial_t z = \partial_x \left\{ \frac{D(\Phi_\delta^{-1}(z))}{\Phi'_\delta(\Phi_\delta^{-1}(z))} \partial_x z + \frac{D(\Phi_\delta^{-1}(z)) (|\Phi_\delta^{-1}(z)| + \delta)}{B} \partial_x (D(\Phi_\delta^{-1}(z)) (|\Phi_\delta^{-1}(z)| + \delta) \partial_t z) \right\} \quad \text{in } Q_T.$$

Moreover we can express the boundary and initial conditions in z as

$$(14) \quad z(0, t) = \Phi_\delta(\phi_g(t)) \quad \text{on } (0, T), \quad z(x, t=0) = \Phi_\delta(0) = 0 \quad \text{on } (0, L),$$

$$(15) \quad \frac{1}{\Phi'_\delta(\Phi_\delta^{-1}(z))} \partial_x z + \frac{(|\Phi_\delta^{-1}(z)| + \delta)}{B} \partial_x (D(\Phi_\delta^{-1}(z)) (|\Phi_\delta^{-1}(z)| + \delta) \partial_t z) = 0 \quad \text{at } x = L.$$

Our paper is organized as follows: section 2 describes the physical model, first proposed in [31]. We repeat the derivations from [16], [17], [19] for reasons of easy reference and unified notation.

In section 3 we introduce the discretization of the problem (13)–(15). We get the Cauchy problem for an implicit first order system of ODEs. Next the solvability of the discretized problem is proved, and the uniform L^2 a priori estimates for the first derivatives and the mixed second derivative are obtained for a small time interval $(0, T_0)$. They imply the short-time existence for the regularized problem.

We continue with section 4 where we use the entropy to establish the existence of a solution for the regularized problem for all times. Next, we establish L^∞ -bounds independent of the regularization parameter.

The last section 5 concerns the existence for the original problem. Using the entropy, the estimates for the time derivative, and the L^∞ -bounds again, we are able to pass to the limit when the regularization parameter tends to zero and prove the existence of at least one solution for the original problem.

2. Model equation for stress-induced diffusion. Consider a coal particle between the fractured cleat systems in coal. The matrix block can be considered as a small (30 μm diameter) cubical particle consisting of glassy coal. The coal face is exposed to the penetrant, in our case CO_2 . The coal face of the particle and the mechanism of the sorption process is shown schematically in Figure 1. The coal originates from a cellulose-like polymeric structure [32], with the chemical formula $\text{C}_n(\text{H}_2\text{O})_m$, from which part of the hydrogen and oxygen have disappeared during the coalification process, which took millions of years. The remaining structure behaves like a glassy polymer, which contains holes (sites) that can accommodate CO_2 , CH_4 , etc. In other words, sorption of gases by coal is more a dissolution process than is adsorption of gases at a coal surface. The holes receiving the CO_2 are originally too small to accommodate the molecule and need to expand. Consequently, the expanded hole exerts a stress on the neighboring molecules constituting the polymeric coal. Therefore the penetration of CO_2 will lead to both a stress gradient and a concentration gradient. The concentration will be expressed as a volume fraction ϕ , i.e., $\phi = c/\Omega$, where c is the molecular concentration and Ω is the molecular volume. As the CO_2 likes to move toward a region of smaller stress, the transport of the molecule will be caused by both a concentration gradient and a stress gradient. When the stresses become too high, a deformation occurs, in which the glassy polymeric structure is converted to a rubber-like (swollen) structure, which is much more open. Consequently the diffusion coefficient in the rubber-like structure is much higher (more

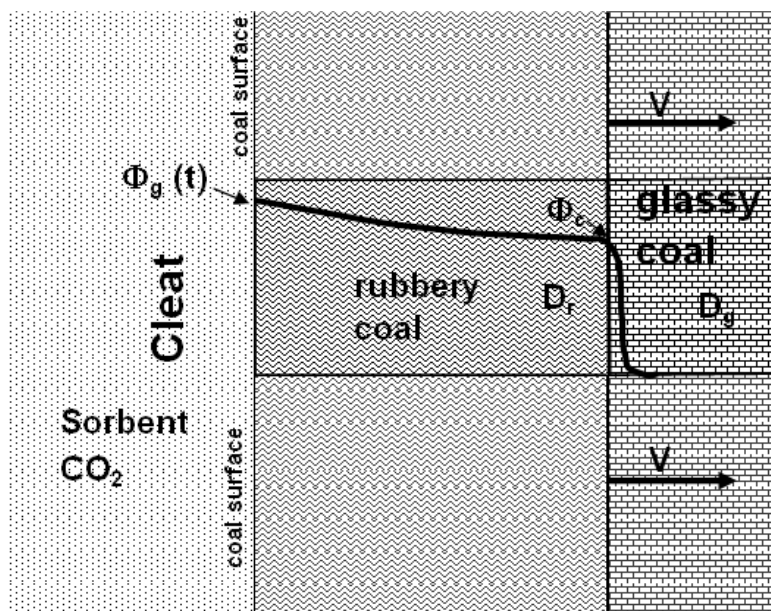


FIG. 1. A coal face exposed to a sorbent (CO_2). To the far right is the virgin coal, which behaves as a glassy polymer. As the sorbent penetrates in the coal, a reorientation of the polymeric coal structure occurs, and the coal becomes rubber-like. The diffusion coefficient in the rubber-like structure is much higher ($> 1000 \times$) than in the glassy structure. The rubber-like structure has also a lower density leading to swelling.

than 1000 times) than the diffusion coefficient in the glassy structure. The stresses are considered to depend on the CO_2 concentration in the coal, and conversion to the rubber-like structure occurs instantaneously when a certain critical concentration is exceeded.

These ideas were formulated for the first time by Thomas and Windle [31], and the derivation of the model equations will be explained below.

2.1. Derivation of model equations. The salient features of the Thomas and Windle model [31] are well summarized by Hui et al. [16], [17]. We summarize the derivation here with the help of the article by Hui et al. and the book of Landau and Lifshitz [21]; i.e., the molar (diffusive) flux J is not only driven by the volume fraction (ϕ) (concentration) gradient, but also by the stress (P_{xx}) gradient, i.e.,

$$(16) \quad J = -D \left(\frac{\partial \phi}{\partial x} + \frac{\Omega \phi}{kT} \frac{\partial P_{xx}}{\partial x} \right),$$

where k is the Boltzmann constant. As opposed to the equation in [21], which contains a scalar pressure gradient, the idea here is extended in [19] with the use of the stress gradient $\partial_x P_{xx}$. Hui (see [16], [17]) interprets P as the osmotic pressure. Note that J is the flux of a volume fraction and behaves as a velocity. The diffusion coefficient depends on the concentration. Below a critical volume fraction ϕ_c , a diffusion coefficient $D_g > 0$ characteristic of a glassy state is used, and above ϕ_c the diffusion coefficient $D_r > 0$ characteristic of the rubber (swollen) state is used. It can be expected that $D_r/D_g \gg 1$. In the model an abrupt change of the diffusion coefficients at ϕ_c is used, but D_r and D_g are considered constant for $\phi > \phi_c$ and $\phi < \phi_c$, respectively:

$$(17) \quad D(\xi) := \begin{cases} D_g, & 0 \leq \xi < \phi_c - \kappa, \\ D_g + (D_r - D_g)(\xi - \phi_c + \kappa)/(2\kappa), & \phi_c - \kappa \leq \xi \leq \phi_c + \kappa, \\ D_r, & \phi_c + \kappa < \xi < +\infty, \end{cases}$$

where $\kappa > 0$ is a small parameter. Extended nonequilibrium thermodynamics [19] suggests that vice versa also the stress (P_{xx}) is related to the volumetric flux gradient as

$$(18) \quad P_{xx} = -\eta_l \frac{\partial J}{\partial x} = \eta_l \frac{\partial \phi}{\partial t},$$

where the second equality follows from a mass conservation law that assumes incompressible flow,

$$(19) \quad \frac{\partial \phi}{\partial t} + \frac{\partial J}{\partial x} = 0.$$

With η_l we denote the elongational velocity [6], i.e., the resistance of movement due to a velocity gradient $\frac{\partial J}{\partial x}$ in the direction of flow. Elongational viscosity is caused by a resistance force of a fluid to accelerate. Hence, the force is proportional to the component of the gradient of the velocity in the flow direction. Elongational viscosity η_l is always larger than the shear viscosity η_s , e.g., in Newtonian fluids $\eta_l = 3\eta_s$. In this case the “fluid CO₂” is moving in the coal medium. The resistance to flow is largely determined by the coal-CO₂ interaction and not, as in the usual definition of viscosity, as CO₂-CO₂ interaction. Hence here we deal with an apparent or pseudoviscosity. With increasing CO₂ concentration the coal becomes more rubber-like, i.e., it acquires a more open structure, and the apparent viscosity decreases with increasing concentration (see (20)). The elongational viscosity η_l is supposed to depend on the volume fraction of the penetrant as

$$(20) \quad \eta_l = \eta_o \exp(-m\phi),$$

where m is a material constant and η_o is the volumetric viscosity of the unswollen coal sample.

Substituting expression (16) for the flux into the mass balance equation (19), where we also use (18), we obtain

$$(21) \quad \partial_t \phi = \partial_x \left\{ D(\phi) \partial_x \phi + \frac{D(\phi) \phi}{B} \partial_x (e^{-m\phi} \partial_t \phi) \right\},$$

where the constant $B = k_B T / (\eta_o \Omega)$. This equation is defined in $Q_T = (0, L) \times (0, T)$.

As initial condition we have that the concentration is

$$(22) \quad \phi(x, t = 0) = 0 \quad \text{on } (0, T).$$

The boundary condition at $x = 0$ must be derived from thermodynamic arguments. The final equilibrium concentration is reached when the coal has swollen to make the stress $P_{xx} = P_{xx}^0$ equal to zero. In this case the volume fraction of CO₂ in the coal is in equilibrium with the CO₂ in the fluid phase outside the coal. Also the CO₂ in the stressed coal is in equilibrium with the CO₂ in the fluid phase. The change in chemical potential is $d\mu = \Omega dP_{xx} + k_B T d \ln \phi$. Equating the chemical potential in the unstressed and stressed state leads to

$$(23) \quad \Omega P_{xx} + k_B T \ln \phi = \Omega P_{xx}^0 + k_B T \ln \phi_o,$$

where ϕ_o is the volume fraction at the coal boundary that would be in equilibrium with the CO₂ in the gas phase if the coal has relaxed to the rubber state with $P_{xx}^0 = 0$.

Substituting (18) and (20) into (23) leads to

$$(24) \quad t = -\phi_o \frac{\eta_0 \Omega}{k_B T} \int_0^{\phi/\phi_o} \frac{\exp(-m\phi_o y)}{\ln y} dy,$$

where we use the initial condition that $\phi = 0$ at $t = 0$. Singularity of the integrand at $y = 1$ guarantees that ϕ remains bounded by ϕ_o for all times.

At $x = L$ we have the boundary condition on $(0, T)$,

$$(25) \quad D(\phi) \left(\partial_x \phi + \frac{1}{B} \phi \partial_x (\exp(-m\phi) \partial_t \phi) \right)_{x=L} = 0.$$

In summary, we have one initial condition equation (22), one boundary condition at $x = L$, viz. (25), and the implicit boundary condition equation (24), which specifies $\phi(x = 0, t)$ as

$$(26) \quad \phi(0, t) = \phi_g(t).$$

ϕ_g satisfies the conditions

$$(27) \quad 0 \leq \phi_g \leq A_0, \quad \phi_g(0) = 0.$$

Remark 2. Equations like (21) can occur in many transport problems in which the flux is calculated using CIT or EIT. A well-known example for CIT in porous media flow is that the deviation of the capillary pressure P_c from its equilibrium value at a given oil saturation S_o , i.e., $P_c^o = P_c^o(S_o)$, is a driving force leading to a rate of change of the saturation (scalar flux). This leads [14], [27], [28], [22] to $\partial_t S_o = L(P_c - P_c^o)$ and to the transport equation for counter current imbibition:

$$\begin{aligned} \varphi \partial_t S_o &= \partial_x (\Lambda(S_o) \partial_x P_c) \\ &= \partial_x (\Lambda(S_o) \partial_x P_c^o(S_o)) + \partial_x \left(\Lambda(S_o) \partial_x \frac{1}{L(S_o)} \partial_t S_o \right). \end{aligned}$$

EIT [19] differs from CIT as it characterizes a system not only by its local thermodynamic variables (pressure, temperature, and concentration) but also by its gradients. The explanation in [19] is difficult to follow by nonspecialists as many thermodynamic relations are considered to be known by the reader. In isothermal systems and in the absence of other applied fields, e.g., electric fields, the volumetric flux J is, according to EIT, given by the following system of equations:

$$(28) \quad \tau_1 \partial_t J + J = -D \left(\frac{\partial \phi}{\partial x} + \frac{\Omega \phi}{kT} \frac{\partial P_{xx}}{\partial x} \right),$$

$$(29) \quad \tau_2 \partial_t P_{xx} + P_{xx} = -\eta_l \frac{\partial J}{\partial x}.$$

Reference [19] uses a mass flux instead of a volumetric flux and therefore uses a factor v_1 , being the partial volume per unit mass. Here τ_1, τ_2 are time constants, which are small with respect to L^2/D . The first terms on the left-hand sides of (28) and (29) appear only in EIT and not in CIT. The first terms on the left are of interest for

short-time behavior and are omitted from the model discussed here. Another example from EIT is the Taylor dispersion problem (see equation 10.34 in [19]), where there is an “xxt” derivative in the concentration, apart from many other terms. Hence, EIT or CIT can lead to transport equations of the form of (21).

3. Short-time existence for the regularized problem. In this section, we first introduce an approximate problem corresponding to (13)–(15). It is a first order system of ODEs for expansion coefficients, with implicit dependence on the time derivatives. First we prove the solvability on some interval $(0, T_N)$, where N is the parameter describing discretization in space. Then, we use the entropy to prove the solvability on interval $(0, T_0)$, where T_0 does not depend on N . Finally, we pass to the limit $N \rightarrow +\infty$ and prove that the problem (13)–(15) itself has a solution on $(0, T_0)$.

Let $V := \{g \in H^1(0, L) \mid g(0) = 0\}$ be the closure of \mathcal{V} in $H^1(0, L)$, and let $\{\alpha_j\}_{j \in \mathbb{N}}$ be a C^∞ -basis for V . We set $V_N := \text{span}\{\alpha_1, \dots, \alpha_N\}$ and introduce the following coefficients:

$$(30) \quad \begin{aligned} d_1(z) &:= \frac{1}{(\Phi'_\delta(\Phi_\delta^{-1}(z)))}, & d_2(z) &:= \frac{D(\Phi_\delta^{-1}(z))}{(\Phi'_\delta(\Phi_\delta^{-1}(z)))}, & \text{and} \\ d(z) &:= D(\Phi_\delta^{-1}(z)) (|\Phi_\delta^{-1}(z)| + \delta). \end{aligned}$$

The coefficients d_1 and d_2 are continuous, nonnegative, and bounded functions of z . d is a continuous function of z , bounded away from zero.

We start study of the initial boundary problem (13)–(15) by constructing an approximate solution for every N . It is defined as follows.

APPROXIMATE PROBLEM 3. For $q \in (2, +\infty)$, find $z_N = \sum_{j=1}^N c_j(t) \alpha_j(x) + \Phi_\delta(\phi_g(t)) \in W^{1,q}([0, T]; V_N)$ such that

$$(31) \quad \begin{aligned} &\int_0^L \partial_t z_N d_1(z_N) \alpha_k dx + \int_0^L d_2(z_N) \partial_x z_N \partial_x \alpha_k dx \\ &+ \int_0^L \frac{1}{B} d(z_N) \partial_x (d(z_N) \partial_t z_N) \partial_x \alpha_k dx = 0 \text{ for } k = 1, \dots, N, \quad \text{and} \end{aligned}$$

$$(32) \quad z_N|_{t=0} = P_N(z|_{t=0} - \Phi_\delta(\phi_g(0))) = 0,$$

where $P : V \rightarrow V_N$ is the projector $P_N(f)(x) := \sum_{j=1}^N \alpha_j(x) (f, \alpha_j)_V$.

Let the vector valued function \mathbf{F} be given by $F_k(t, \mathbf{c}, \partial_t \mathbf{c}) =$ left-hand side of (31) and let \mathbf{c} be the column vector consisting of elements $(c_1(t) \dots c_N(t))$; then (31), (32) are equivalent to the following Cauchy problem in \mathbb{R}^N :

$$(33) \quad \begin{cases} \mathbf{F}(t, \mathbf{c}, \partial_t \mathbf{c}) = 0, \\ \mathbf{c}|_{t=0} = 0. \end{cases}$$

The Cauchy problem (33) is difficult to solve, since the dependence of \mathbf{F} on $\partial_t \mathbf{c}$ is implicit. It is crucial to reduce it to an ordinary Cauchy problem of the form $\partial_t \mathbf{c} = \varrho(t, \mathbf{c})$.

We note that

$$(34) \quad \begin{aligned} F_k &:= \sum_{j=1}^N \left\{ \int_0^L d_1(z_N) \alpha_k \alpha_j dx + \int_0^L \frac{1}{B} d(z_N) \partial_x (d(z_N) \alpha_j) \partial_x \alpha_k dx \right\} \frac{dc_j}{dt} \\ &+ \sum_{j=1}^N \left\{ \int_0^L d_2(z_N) \partial_x \alpha_j \partial_x \alpha_k dx \right\} c_j + \int_0^L d_1(z_N) \alpha_k \partial_t \Phi_\delta(\phi_g(t)) dx. \end{aligned}$$

Then, after introducing the matrices $\mathcal{A}(\mathbf{c})$ and $\mathcal{B}(\mathbf{c})$ and the vector $\mathbf{f}(\mathbf{c})$ by

$$(35) \quad \mathcal{A}_{kj}(\mathbf{c}) := \int_0^L d_1(z_N) \alpha_k \alpha_j \, dx + \int_0^L \frac{1}{B} d(z_N) \partial_x (d(z_N) \alpha_j) \partial_x \alpha_k \, dx,$$

(36)

$$\mathcal{B}_{kj}(\mathbf{c}) := \int_0^L d_2(z_N) \partial_x \alpha_j \partial_x \alpha_k \, dx, \quad \text{and} \quad f_k(\mathbf{c}) = \int_0^L d_1(z_N) \alpha_k \partial_t \Phi_\delta(\phi_g(t)) \, dx,$$

$1 \leq k, j \leq N$, we see that the problem (31)–(32) is equivalent to the following Cauchy problem:

Find $\mathbf{c} \in W^{1,q}(0, T)^N$ such that

$$(37) \quad \mathcal{A}(\mathbf{c}) \frac{d\mathbf{c}}{dt} = -\mathcal{B}(\mathbf{c})\mathbf{c} - \mathbf{f}(\mathbf{c}) \quad \text{a.e. in } (0, T); \quad \mathbf{c}|_{t=0} = 0.$$

PROPOSITION 4. *There is a $T_N > 0$ such that the problem (31)–(32) has a unique solution $z_N \in W^{1,q}(0, T_N; V_N)$ for all $q < +\infty$.*

Proof. It is enough to prove that the Cauchy problem (37) has a solution.

Obviously, \mathcal{A} , \mathcal{B} , and \mathbf{f} are smooth functions of \mathbf{c} . Because of the singularity of $\partial_t \varphi_g$ at $t = 0$, $\mathbf{f}(\mathbf{c}) \in L^q(0, T)$ for all $q < +\infty$, but it is not bounded. Hence, the only property to check is the invertibility of the matrix \mathcal{A} . Let \mathbf{b} be an arbitrary vector from \mathbb{R}^N and let $b_\alpha(x) = \mathbf{b} \cdot \alpha(x) = \sum_{j=1}^N b_j \alpha_j(x)$. Then we have

$$\begin{aligned} (\mathcal{A}\mathbf{b}) \cdot \mathbf{b} &= \sum_{k,j=1}^N \mathcal{A}_{k,j} b_k b_j = \int_0^L d_1(z_N) (b_\alpha)^2 \, dx + \frac{1}{B} \int_0^L d(z_N) \partial_x b_\alpha \partial_x (d(z_N) b_\alpha) \, dx \\ &= \int_0^L d_1(z_N) (b_\alpha)^2 \, dx + \frac{1}{B} \int_0^L (d(z_N) \partial_x b_\alpha)^2 \, dx \\ &\quad + \frac{1}{B} \int_0^L d(z_N) \partial_x b_\alpha b_\alpha d'(z_N) \partial_x z_N \, dx \\ (38) \quad &\geq \int_0^L \left\{ d_1(z_N) - \frac{1}{4B} (d'(z_N))^2 (\partial_x z_N)^2 \right\} (b_\alpha)^2 \, dx. \end{aligned}$$

Since $\partial_x z_N(x, 0) = 0$ and functions $\{\alpha_j\}_{j \in \mathbb{N}}$ are linearly independent, the matrix \mathcal{A} is by (38) invertible in a neighborhood of $t = 0$. Then by the classical theory, the problem (37) has a unique solution on some interval $(0, T_N)$. \square

Next, we want to prove that the existence interval does not depend on N .

PROPOSITION 5. *There is a constant C , independent of N , such that*

$$(39) \quad \|\partial_x z_N\|_{L^\infty(0, T_N; L^2(0, L))} \leq C.$$

Consequently, the vector valued function \mathbf{c} remains bounded at $t = T_N$.

Proof. In (31) we can replace α_k by $z_N - \Phi_\delta(\phi_g)$. Then after using that $\partial_x (d(z_N) \partial_t z_N) = \partial_t (d(z_N) \partial_x z_N)$, we get

$$(40) \quad \begin{aligned} &\int_0^L d_1(z_N) z_N \partial_t z_N \, dx + \int_0^L d_2(z_N) (\partial_x z_N)^2 \, dx \\ &\quad + \int_0^L \frac{1}{B} \partial_t (d(z_N) \partial_x z_N) d(z_N) \partial_x z_N \, dx \end{aligned}$$

$$\begin{aligned}
 &= \int_0^L d_1(z_N) \Phi_\delta(\phi_g) \partial_t z_N \, dx = \partial_t \int_0^L \Phi_\delta(\phi_g)(t) \int_0^{z_N} d_1(\xi) \, d\xi \, dx \\
 &\quad - \partial_t \Phi_\delta(\phi_g)(t) \int_0^L \int_0^{z_N} d_1(\xi) \, d\xi \, dx.
 \end{aligned}$$

Integrating over t leads to

$$\begin{aligned}
 (41) \quad &\int_0^L \left(\int_0^{z_N(x,t)} d_1(\xi) \, \xi \, d\xi \right) dx + \int_0^t \int_0^L d_2(z_N) (\partial_x z_N)^2 \, dx \, d\tau \\
 &\quad + \frac{1}{2B} \int_0^L d(z_N)^2 (\partial_x z_N)^2 \, dx \\
 &= \int_0^L \left(\int_0^{z_N(x,t)} d_1(\xi) \, d\xi \right) dx \Phi_\delta(\phi_g)(t) - \int_0^t \partial_\tau \Phi_\delta(\phi_g)(\tau) \left(\int_0^L \int_0^{z_N} d_1(\xi) \, d\xi \, dx \right) d\tau.
 \end{aligned}$$

We easily find out that

$$(42) \quad \int_0^z d_1(\xi) \, \xi \, d\xi = \int_0^{\Phi_\delta^{-1}(z)} \Phi_\delta(\eta) \, d\eta \quad \text{and} \quad \int_0^z d_1(\xi) \, d\xi = \Phi_\delta^{-1}(z).$$

The growth of the terms in (42) indicates that it will be possible to control the two terms on the right-hand side of (41) by the first term on the left-hand side of (41).

Let $M_\phi := \max_{0 \leq t \leq T} |\Phi_\delta(\phi_g(t))|$. Then by the definition of $\Phi_\delta(\varphi)$, we have $C_0(\delta) \log(1 + \varphi/\delta) \leq \Phi_\delta(\varphi)$ for all $\varphi \geq 0$. Hence $\int_0^z d_1(\xi) \, \xi \, d\xi \geq C_0(\delta)(|\Phi_\delta^{-1}(z)| + \delta) \log(1 + |\Phi_\delta^{-1}(z)|/\delta) - |\Phi_\delta^{-1}(z)|$, and there is a large enough constant $C_\varphi = C_\varphi(M_\phi, \delta)$ such that $g(z) = C_0(\delta)(|\Phi_\delta^{-1}(z)| + \delta) \log(1 + |\Phi_\delta^{-1}(z)|/\delta) - |\Phi_\delta^{-1}(z)| - M_\phi |\Phi_\delta^{-1}(z)| + C_\varphi > |\Phi_\delta^{-1}(z)|$ for all z . The equality (41) now implies

$$\begin{aligned}
 (43) \quad &\int_0^L g(z_N(x,t)) \, dx + \int_0^t \int_0^L d_2(z_N) (\partial_x z_N)^2 \, dx \, d\tau + \frac{1}{2B} \int_0^L d(z_N)^2 (\partial_x z_N(t))^2 \, dx \\
 &\leq C_\varphi(M_\phi, \delta) L + \int_0^t |\partial_\tau \Phi_\delta(\phi_g)(\tau)| \left(\int_0^L g(z_N(x,\tau)) \, dx \right) d\tau.
 \end{aligned}$$

Since $\partial_\tau \Phi_\delta(\phi_g) \in L^1(0, T)$, we apply Gronwall's inequality, and estimate (39) follows. Hence \mathbf{c} remains bounded at $t = T_N$. \square

Nevertheless, since the matrix \mathcal{A} could degenerate, some components of $\frac{\partial \mathbf{c}}{\partial t}$ could blow up at $t = T_N$. In order to exclude this possibility and to prove that the maximal solution for (33) exists on $[0, T]$, we need an estimate for the time derivatives. Furthermore, if we want to pass to the limit $N \rightarrow +\infty$ in (31), estimate (39) is not sufficient. Our strategy is to obtain an estimate, uniform with respect to N , for $\partial_{xt} z_N$ in $L^2(Q_T)$.

THEOREM 6. *There exists $T_0 > 0$, independent of N , such that*

$$(44) \quad \|\partial_x z_N\|_{L^\infty(0, T_0; L^2(0, L))} \leq C,$$

$$(45) \quad \|\partial_t z_N\|_{L^2(0, T_0; L^2(0, L))} \leq C,$$

$$(46) \quad \|\partial_{xt} z_N\|_{L^2(0, T_0; L^2(0, L))} \leq C,$$

$$(47) \quad \left\| \partial_{xt} \int_0^{z_N} d(\xi) \, d\xi \right\|_{L^2((0, T_0) \times (0, L))} \leq C,$$

with constants independent of N . Consequently, the maximal solution for (33) exists on $[0, T_0]$.

Proof. We replace α_k in (31) by $\partial_t z_N - \partial_t \Phi_\delta(\phi_g)$. This yields

$$(48) \quad \int_0^L d_1(z_N) (\partial_t z_N)^2 dx + \int_0^L d_2(z_N) \partial_x z_N \partial_{xt} z_N dx + \frac{1}{B} \int_0^L d(z_N) \partial_t (d(z_N) \partial_x z_N) \partial_{xt} z_N dx = \int_0^L d_1(z_N) \partial_t z_N \partial_t \Phi_\delta(\phi_g) dx.$$

In the estimates which follow we will use the fact that integrability of higher order derivatives implies continuity and boundedness in x or in t . We recall that for one-dimensional Sobolev embeddings Morrey’s theorem applies and $H^1(0, t)$ (respectively, $H^1(0, L)$) is continuously embedded into the Hölder space $C^{0,1/2}[0, t]$ (respectively, into $C^{0,1/2}[0, L]$). See, e.g., [11] for more details. In our particular situation, we use the explicit dependence of the embedding constant on the length of the time interval and we prefer to derive the estimates directly.

First, as $\partial_x z_N \in L^2(0, L; H^1(0, t))$ and $\partial_x z_N|_{\tau=0} = 0$, we have for a.e. $x \in (0, L)$ and for every $\tau \in (0, t)$

$$(49) \quad |\partial_x z_N(x, \tau)| = \left| \int_0^\tau \partial_\xi \partial_x z_N(x, \xi) d\xi \right| \leq \sqrt{\tau} \sqrt{\int_0^\tau |\partial_\xi \partial_x z_N(x, \xi)|^2 d\xi}.$$

Next, as $\partial_\tau z_N \in L^2(0, t; H^1(0, L))$ and $\partial_\tau z_N|_{\tau=0} = \partial_\tau \Phi(\phi_g)$, we have for a.e. $\tau \in (0, t)$ and for every $x \in (0, L)$

$$(50) \quad \begin{aligned} |\partial_\tau z_N(x, \tau)| &\leq |\partial_\tau \Phi(\phi_g(\tau))| + \left| \int_0^x \partial_\xi \partial_\tau z_N(\xi, \tau) d\xi \right| \\ &\leq |\partial_\tau \Phi(\phi_g(\tau))| + \sqrt{L} \sqrt{\int_0^L |\partial_\xi \partial_\tau z_N(\xi, \tau)|^2 d\xi}. \end{aligned}$$

Estimates (49)–(50) imply

$$(51) \quad \begin{aligned} &\int_0^L \int_0^t |\partial_\tau z_N(x, \tau)|^2 |\partial_x z_N(x, \tau)|^2 dx d\tau \\ &\leq 2 \int_0^L \int_0^t \tau \left(\int_0^t |\partial_\xi \partial_x z_N(x, \xi)|^2 d\xi \right) \left(|\partial_\tau \Phi(\phi_g(\tau))|^2 + L \int_0^L |\partial_\xi \partial_\tau z_N(\xi, \tau)|^2 d\xi \right) d\tau dx \\ &\leq 2Lt \|\partial_{x\tau} z_N\|_{L^2((0,t) \times (0,L))}^4 + 2 \|\partial_{x\tau} z_N\|_{L^2((0,t) \times (0,L))}^2 \int_0^t \tau |\Phi'_\delta(\phi_g)|^2 |\partial_\tau \phi_g|^2 d\tau \\ &\leq \left(2\sqrt{L}\sqrt{t} \|\partial_{x\tau} z_N\|_{L^2((0,t) \times (0,L))}^2 + \frac{1}{\sqrt{2L}} \int_0^t \tau^{1/2} |\Phi'_\delta(\phi_g)|^2 |\partial_\tau \phi_g|^2 d\tau \right)^2. \end{aligned}$$

Now we integrate (48) with respect to time, over $(0, t)$, and estimate the obtained terms. The second term is estimated as follows:

$$(52) \quad \begin{aligned} &\left| \int_0^t \int_0^L d_2(z_N) \partial_x z_N \partial_{x\tau} z_N dx d\tau \right| \leq C \int_0^t \|\partial_{x\tau} z_N(\tau)\|_{L^2(0,L)} \|\partial_x z_N(\tau)\|_{L^2(0,L)} d\tau \\ &\leq C \sqrt{\int_0^t \|\partial_{x\tau} z_N(\tau)\|_{L^2(0,L)}^2 d\tau} \sqrt{\int_0^t \int_0^L (\partial_x z_N)^2 dx d\tau} \leq C \|\partial_{x\tau} z_N\|_{L^2((0,t) \times (0,L))}, \end{aligned}$$

where we have used the estimate (39). We rewrite the third term of (48), omitting the $1/B$ factor, as

$$(53) \quad \int_0^t \int_0^L d(z_N) \partial_\tau (d(z_N) \partial_x z_N) \partial_{x\tau} z_N \, dx \, d\tau = \int_0^t \int_0^L d(z_N) (\partial_{x\tau} z_N)^2 \, dx \, d\tau + \int_0^t \int_0^L d(z_N) d'(z_N) \partial_\tau z_N \partial_x z_N \partial_{x\tau} z_N \, dx \, d\tau.$$

The last term in (53) is cubic in derivatives of z_N . Our idea is to use the estimate (51), showing that for small times it enters with a small coefficient and then controlling it using other terms. Using the estimate (51), we find out that it satisfies the following inequality:

$$(54) \quad \left| \int_0^t \int_0^L d(z_N) \partial_{x\tau} z_N d'(z_N) \partial_\tau z_N \partial_x z_N \, dx \, d\tau \right| \leq C \|\partial_{x\tau} z_N\|_{L^2(Q_t)} \|\partial_\tau z_N \partial_x z_N\|_{L^2(Q_t)} \leq C\sqrt{t} \left(\|\partial_{x\tau} z_N\|_{L^2(Q_t)}^3 + \left\| \tau^{1/6} \Phi'_\delta(\phi_g) \partial_\tau \phi_g \right\|_{L^3(0,t)}^3 \right),$$

where $Q_t = (0, t) \times (0, L)$. Let $X^2(t) := \int_0^t \int_0^L |\partial_{x\tau} z_N|^2 \, dx \, d\tau$. Since $\partial_\tau \Phi_g(\phi_g) \in L^2(0, t)$, estimates (39), (52), (53), and (54) imply

$$(55) \quad \left\| \sqrt{d_1(z_N)} \partial_\tau z_N \right\|_{L^2((0,t) \times (0,L))}^2 + X^2(t) - C_1 \sqrt{t} X^3(t) \leq C_o,$$

where C_o depends on $\|\partial_\tau \Phi_\delta(\phi_g)\|_{L^2(0,t)}$ and on the constant from estimate (39). We note that the last term on the left-hand side corresponds to the lower bound for the cubic term, corresponding to the stress gradient part of the diffusive flux. Inequality (55) is satisfied for $t = 0$. The function $\varrho(X) = X^2 - C_1 \sqrt{t} X^3$ has its maximum on $(0, +\infty)$ in the point $X_o = 3/(2C_1 \sqrt{t})$. If $C_o < \varrho(X_o)$, then inequality (55) gives an estimate for $X(t)$. We note that $C_o < \varrho(X_o)$ if $t < \frac{4}{27C_1^2 C_o}$. Hence for $T \leq \frac{4}{27C_1^2 C_o} = T_0$ we have estimates (44)–(46).

From (44)–(46) it follows that $\partial_x z_N \partial_t z_N \in L^2(0, T_0; L^2(0, L)) \leq C$, and we have (47) as well. \square

The estimates (44)–(47) allow us to pass to the limit $N \rightarrow +\infty$. Using classical compactness and weak compactness arguments and due to the a priori estimates (44)–(47), we can extract a subsequence of z_N , denoted by the same subscripts, which converges to an element $z \in H^1((0, T_0) \times (0, L))$, $\partial_{xt} z \in L^2((0, T_0) \times (0, L))$, in the following sense:

$$(56) \quad z_N \rightarrow z \text{ strongly in } L^2((0, T_0) \times (0, L)) \text{ and a.e. on } (0, T_0) \times (0, L),$$

$$(57) \quad \partial_x z_N \rightharpoonup \partial_x z \text{ weakly in } L^2((0, T_0) \times (0, L)),$$

$$(58) \quad \partial_t z_N \rightharpoonup \partial_t z \text{ weakly in } L^2((0, T_0) \times (0, L)),$$

$$(59) \quad \partial_{xt} z_N \rightharpoonup \partial_{xt} z \text{ weakly in } L^2((0, T_0) \times (0, L)),$$

$$(60) \quad \partial_{xt} \int_0^{z_N} d(\xi) \, d\xi \rightharpoonup \partial_{xt} \int_0^z d(\xi) \, d\xi \text{ weakly in } L^2((0, T_0) \times (0, L)).$$

Now passing to the limit $N \rightarrow \infty$ in (31) does not pose problems, and we conclude that z satisfies (13)–(15).

We summarize the results in the following theorem.

THEOREM 7. *Let $\phi_g \in H^1(0, T)$. Then there exists $T_0 > 0$ such that problem (13)–(15) has at least one variational solution $z \in H^1((0, T_0) \times (0, L))$, $\partial_{xt}z \in L^2((0, T_0) \times (0, L))$.*

COROLLARY 8. *Let $\phi_g \in H^1(0, T)$. Then there exists $T_0 > 0$ such that the variational formulation (11)–(12) of the problem (8), (9), and (3) has at least one solution $\phi = \Phi_\delta^{-1}(z) \in H^1((0, T_0) \times (0, L))$, $\partial_{xt}\phi \in L^2((0, T_0) \times (0, L))$.*

4. Existence of the regularized problem. In this section, we first use the regularized entropy function to prove the global existence for the problem (8), (9), and (3) (i.e., for the regularized problem). Then we establish the L^∞ -bounds for the solution, independent of the regularization parameter.

Let us prove that any solution ϕ for the problem (8), (9), and (3), constructed in Corollary 8, could be extended from $(0, T_0)$ to arbitrary time interval $(0, T)$. First we test (11) by $\Phi_\delta(\phi) - \Phi_\delta(\phi_g(t))$. We have

$$\begin{aligned} & \int_0^t \int_0^L \partial_\tau \phi \Phi_\delta(\phi) \, dx \, d\tau + \int_0^t \int_0^L D(\phi) \partial_x \phi \partial_x \Phi_\delta(\phi) \, dx \, d\tau \\ & + \int_0^t \int_0^L \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_\tau \phi \right) \partial_x \Phi_\delta(\phi) \, dx \, d\tau \\ & = \int_0^t \int_0^L \partial_\tau \phi \Phi_\delta(\phi_g) \, dx \, d\tau, \end{aligned}$$

and it follows that

$$\begin{aligned} & \int_0^L \left(\int_0^{\phi(t)} \Phi_\delta(\xi) \, d\xi \right) dx + \int_0^t \int_0^L \frac{1}{2B} \partial_\tau \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_x \phi \right)^2 \, dx \, d\tau \\ & + \int_0^t \int_0^L D(\phi) \Phi'_\delta(\phi) (\partial_x \phi)^2 \, dx \, d\tau = \int_0^L \phi(t) \Phi_\delta(\phi_g(t)) \, dx \\ & - \int_0^t \int_0^L \phi \partial_\tau \Phi_\delta(\phi_g) \, dx \, d\tau \end{aligned}$$

and we get as in the proof of Proposition 5

$$(61) \quad \|\partial_x \phi\|_{L^\infty(0,t;L^2(0,L))} \leq C.$$

This estimate implies the boundedness of ϕ . We note that C does not depend on the smoothing of D at $\phi = \phi_c$.

Next we test (11) by

$$e^{-m \min\{|\phi|, 1/\delta\}} \partial_t \phi - e^{-m \min\{|\phi_g|, 1/\delta\}} \partial_t \phi_g$$

and get

$$\begin{aligned} & \int_0^t \int_0^L (\partial_\tau \phi)^2 e^{-m \min\{|\phi|, 1/\delta\}} \, dx \, d\tau \\ & + \int_0^t \int_0^L D(\phi) \partial_x \phi \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_\tau \phi \right) \, dx \, d\tau \\ & + \int_0^t \int_0^L \frac{D(\phi)(|\phi| + \delta)}{2B} \left(\partial_{x\tau} \int_0^\phi e^{-m \min\{|\xi|, 1/\delta\}} \, d\xi \right)^2 \, dx \, d\tau \end{aligned}$$

$$= \int_0^t \int_0^L \partial_\tau \phi e^{-m \min\{|\phi_g|, 1/\delta\}} \partial_\tau \phi_g \, dx \, d\tau.$$

Then, as in the proof of Proposition 5, by estimating the second and the fourth terms and after using (61), we conclude that

$$(62) \quad \|\partial_\tau \phi\|_{L^2((0,t) \times (0,L))} \leq C,$$

$$(63) \quad \left\| \partial_{x\tau} \int_0^\phi e^{-m \min\{|\xi|, 1/\delta\}} \, d\xi \right\|_{L^2((0,t) \times (0,L))} \leq C,$$

and from this it follows that

$$(64) \quad \|\partial_{x\tau} \phi\|_{L^2((0,t) \times (0,L))} \leq C.$$

Therefore, we arrive at the following theorem.

THEOREM 9. *Let $\phi_g \in H^1(0, T)$. Then for all $T > 0$ there exists a weak solution $\phi \in H^1((0, T) \times (0, L))$, $\partial_{xt}\phi \in L^2((0, T) \times (0, L))$ for the variational formulation (11)–(12) of the problem (8), (9), and (3).*

We conclude this section by establishing uniform L^∞ -bounds for ϕ . We have the following proposition.

PROPOSITION 10. *Let $\phi_g \in H^1(0, T)$ and $\phi_g \geq 0$. Then any weak solution ϕ of the problem (8), (9), and (3), obtained in Theorem 9, satisfies $\phi(x, t) \geq 0$ a.e. on Q_T .*

Proof. Let $a_- = -\inf\{a, 0\}$ and $a_+ = \sup\{a, 0\}$. Then $a = a_+ - a_-$ and $\Phi_\delta((\phi_g)_-) = \Phi_\delta(0) = 0$. We test (11) by $\Phi_\delta(\phi_-)$. Note that $\Phi_\delta(\phi_-)|_{x=0} = 0$ and $\Phi_\delta(\phi_-) \geq 0$. Then we have

$$\begin{aligned} & \int_0^t \int_0^L (\partial_\tau \phi) \Phi_\delta(\phi_-) \, dx \, d\tau + \int_0^t \int_0^L D(\phi) \partial_x \phi \partial_x \Phi_\delta(\phi_-) \, dx \, d\tau \\ & + \int_0^t \int_0^L \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_\tau \phi \right) \Phi'_\delta(\phi_-) \partial_x \phi_- \, dx \, d\tau = 0. \end{aligned}$$

Since $\phi_-|_{t=0} = 0$, $\phi_+ \phi_- = 0$, and $|\phi| \phi_- = \phi_-^2$, we get

$$\begin{aligned} & \int_0^L \left(\int_0^{\phi_-(x,t)} \Phi_\delta(\xi) \, d\xi \right) dx + \int_0^t \int_0^L D(\phi_-) \Phi'_\delta(\phi_-) (\partial_x \phi_-)^2 \, dx \, d\tau \\ & + \int_0^L \frac{D(\phi)(|\phi| + \delta)}{2B} \left(e^{-m \min\{|\phi|, 1/\delta\}} \partial_x \phi_- \right)^2 (t) \, dx = 0. \end{aligned}$$

It follows that $\partial_x \phi_- = 0$ and $\phi_-|_{x=0} = 0$. Therefore $\phi_- = 0$, and consequently $\phi = \phi_+ \geq 0$. \square

In the uniform bounds which follow, we use, for given positive constants m and δ , the function

$$(65) \quad G(z) := \int_0^z \exp\{-m \min\{\xi, 1/\delta\}\} \, d\xi, \quad z \geq 0.$$

Then we have the following bounds.

PROPOSITION 11. *Let $\phi_g \in H^1(0, T)$, $\phi_g \geq 0$ and $\partial_t \phi_g \geq 0$ a.e. on $(0, T)$. Then any weak solution ϕ of the problem (8), (9), and (3), obtained in Theorem 9, satisfies $\phi_g(t) \geq \phi(x, t)$ a.e. on Q_T .*

Proof. Let G be given by (65). We test (11) by $(G(\phi) - G(\phi_g))_+$. Note that $(G(\phi) - G(\phi_g))_+|_{x=0} = 0$. Then we have

$$(66) \quad \int_0^t \int_0^L \partial_\tau \phi (G(\phi) - G(\phi_g))_+ dx d\tau + \int_0^t \int_0^L D(\phi) \partial_x \phi \partial_x (G(\phi) - G(\phi_g))_+ dx d\tau + \int_0^t \int_0^L \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x (e^{-m \min\{|\phi|, 1/\delta\}} \partial_\tau \phi) \partial_x (G(\phi) - G(\phi_g))_+ dx d\tau = 0.$$

Note that

$$(67) \quad \partial_\tau \phi (G(\phi) - G(\phi_g))_+ = \partial_\tau \left(\int_0^\phi (G(\xi) - G(\phi_g))_+ d\xi \right) + G'(\phi_g) \partial_\tau \phi_g (\phi - \phi_g)_+$$

and

$$(68) \quad \frac{D(\phi)(|\phi| + \delta)}{B} \partial_\tau \partial_x G(\phi) \partial_x (G(\phi) - G(\phi_g))_+ = \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} (\partial_x (G(\phi) - G(\phi_g))_+)^2 \right) - (\partial_x (G(\phi) - G(\phi_g))_+)^2 \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right).$$

Then using the monotonicity of ϕ_g and G we obtain from (66), (67), and (68) the following inequality:

$$\begin{aligned} & \int_0^L \left(\int_0^{\phi(x,t)} (G(\xi) - G(\phi_g))_+ d\xi \right) dx + \int_0^t \int_0^L \frac{D(\phi)}{G'(\phi)} \left(\partial_x (G(\phi) - G(\phi_g))_+ \right)^2 dx d\tau \\ & \quad + \int_0^L \frac{D(\phi)(|\phi| + \delta)}{2B} (\partial_x (G(\phi) - G(\phi_g))_+)^2 dx \\ & \leq \int_0^t \int_0^L (\partial_x (G(\phi) - G(\phi_g))_+)^2 \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right) dx d\tau. \end{aligned}$$

Since $\partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right) \in L^2(0, T; L^\infty(0, L))$, we apply Gronwall's lemma and conclude that $(G(\phi) - G(\phi_g))_+ = 0$, from which it follows that $G(\phi) \leq G(\phi_g)$. Inversion of this equation leads to $\phi(x, t) \leq \phi_g(t)$ a.e. on Q_T . \square

PROPOSITION 12. *Let $\phi_g \in H^1(0, T)$, and let us suppose in addition that there are constants $A_0 > 0$, $\alpha > 0$, and $C_0 > 0$ such that*

$$(69) \quad A_0 \geq \phi_g(t) \geq C_0 t^\alpha \quad \forall t \in [0, T].$$

Then any weak solution ϕ of the problem (8), (9), and (3), obtained in Theorem 9, satisfies $A_0 \geq \phi(x, t) \geq C_0 t^\alpha$ a.e. on Q_T .

Proof. The proof follows the lines of Proposition 11. It is enough to prove the lower bound. We test (11) by $(G(C_0 t^\alpha) - G(\phi))_-$. Note that $(G(C_0 t^\alpha) - G(\phi))_-|_{x=0} = 0$. Then as in the proof of Proposition 11 we have

$$(70) \quad \int_0^t \int_0^L \partial_\tau \phi (G(C_0 \tau^\alpha) - G(\phi))_- dx d\tau + \int_0^t \int_0^L D(\phi) \partial_x \phi \partial_x (G(C_0 \tau^\alpha) - G(\phi))_- dx d\tau + \int_0^t \int_0^L \frac{D(\phi)(|\phi| + \delta)}{B} \partial_x (e^{-m \min\{|\phi|, 1/\delta\}} \partial_\tau \phi) \partial_x (G(C_0 \tau^\alpha) - G(\phi))_- dx d\tau = 0.$$

Note that

$$(71) \quad G(\phi) = G(C_0 t^\alpha) - (G(C_0 t^\alpha) - G(\phi))_+ + (G(C_0 t^\alpha) - G(\phi))_-,$$

$$\partial_\tau \phi (G(C_0 \tau^\alpha) - G(\phi))_- = \frac{\partial_\tau G(C_0 \tau^\alpha)}{G'(\phi)} (G(C_0 \tau^\alpha) - G(\phi))_-$$

$$(72) \quad + \frac{1}{2G'(\phi)} \partial_\tau (G(C_0 \tau^\alpha) - G(\phi))_-^2 \geq \frac{1}{2G'(\phi)} \partial_\tau (G(C_0 \tau^\alpha) - G(\phi))_-^2,$$

and

$$(73) \quad \begin{aligned} & \frac{D(\phi)(|\phi| + \delta)}{B} \partial_\tau \partial_x G(\phi) \partial_x (G(C_0 \tau^\alpha) - G(\phi))_- \\ &= \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} (\partial_x (G(C_0 \tau^\alpha) - G(\phi))_-)^2 \right) \\ & - (\partial_x (G(C_0 \tau^\alpha) - G(\phi))_-)^2 \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right). \end{aligned}$$

Then using the monotonicity of G we obtain from (70), (72), and (73) the following inequality:

$$\begin{aligned} & \int_0^L \frac{(G(C_0 t^\alpha) - G(\phi))_-^2}{2G'(\phi)} dx + \int_0^t \int_0^L \frac{D(\phi)}{G'(\phi)} (\partial_x (G(C_0 \tau^\alpha) - G(\phi))_-)^2 dx d\tau \\ & \quad + \int_0^L \frac{D(\phi)(|\phi| + \delta)}{2B} (\partial_x (G(C_0 t^\alpha) - G(\phi))_-)^2 dx \\ & \leq \int_0^t \int_0^L (\partial_x (G(C_0 \tau^\alpha) - G(\phi))_-)^2 \partial_\tau \frac{1}{2G'(\phi)} dx d\tau \\ & \quad + \int_0^t \int_0^L (G(C_0 \tau^\alpha) - G(\phi))_-^2 \partial_\tau \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right) dx d\tau. \end{aligned}$$

Since $\partial_t \left(\frac{D(\phi)(|\phi| + \delta)}{2B} \right)$ and $\partial_t \frac{1}{2G'(\phi)}$ are elements of $L^2(0, T; L^\infty(0, L))$, we apply Gronwall's lemma to the function

$$\int_0^t \|\partial_\tau \phi(\tau)\|_{L^\infty(0, L)} \|(G(C_0 \tau^\alpha) - G(\phi(\tau)))_-\|_{H^1(0, L)}^2 d\tau$$

and conclude that $(G(C_0 t^\alpha) - G(\phi))_- = 0$, from which it follows that $G(\phi) \geq G(C_0 t^\alpha)$. Inversion of G leads to $\phi(x, t) \geq C_0 t^\alpha$ a.e. on Q_T . \square

THEOREM 13. *Let $\phi_g \in H^1(0, T)$, $A_0 = \max_{0 \leq t \leq T} \phi_g(t)$, $A_0 \geq \phi_g \geq C_0 t^\alpha$, and $\alpha > 1$. Then there exists a weak solution ϕ , $C_0 t^\alpha \leq \phi(x, t) \leq A_0$, $\partial_{xt} \phi \in L^2((0, T) \times (0, L))$, $\phi \in H^1((0, T) \times (0, L))$, for the problem (8), (9), and (3).*

Remark 14.

- By choosing $\delta < 1/A_0$, we can replace $e^{-m \min\{|\phi|, 1/\delta\}}$ by $e^{-m\phi}$ and $|\phi| + \delta$ by $\phi + \delta$.
- In addition to the assumptions of Theorem 13 let us suppose that $\partial_t \phi_g \geq 0$. Then there exists a weak solution ϕ , $C_0 t^\alpha \leq \phi(x, t) \leq \phi_g(t)$, $\partial_{xt} \phi \in L^2((0, T) \times (0, L))$, $\phi \in H^1((0, T) \times (0, L))$, for the problem (8), (9), and (3).

5. Existence for the original problem. It remains to pass to the limit $\delta \rightarrow 0$. This limit will give us the solvability of the starting problem (1)–(3).

After Theorem 13, we are free to replace the nonlinearity $\exp\{-m\phi\}$ by $h(\xi) = e^{-m} \min\{\xi, A_0\}$, $\xi \geq 0$. We have existence for the system (11)–(12); i.e., for every $g \in L^2(0, T; V)$, $V = \{g \in H^1(0, L) \mid g(0) = 0\}$, we have

$$(74) \quad \int_0^T \int_0^L \partial_t \phi_\delta g \, dx \, dt + \int_0^T \int_0^L D(\phi_\delta) \left\{ \partial_x \phi_\delta + \frac{(\phi_\delta + \delta)}{B} \partial_x (h(\phi_\delta) \partial_t \phi_\delta) \right\} \partial_x g \, dx \, dt = 0,$$

$$(75) \quad \phi_\delta|_{x=0} = \phi_g(t) \quad \text{and} \quad \phi_\delta|_{t=0} = 0,$$

and we want to pass to the limit $\delta \rightarrow 0$.

Let

$$(76) \quad \Psi'_\delta(\xi) := \frac{h(\xi)}{D(\xi)(\xi + \delta)}, \quad \xi \geq 0,$$

and

$$(77) \quad \Psi_\delta(\phi) := \int_0^\phi \frac{1}{\xi + \delta} \left(\frac{h(\xi)}{D(\xi)} - \frac{h(0)}{D(0)} \right) d\xi + \frac{h(0)}{D(0)} \log(\phi + \delta) \quad \text{for } \phi \geq 0.$$

It should be noted that $\Psi_\delta(0) = \frac{h(0)}{D(0)} \log \delta < 0$, which would cause some complications.

THEOREM 15. *Let $\alpha > 0$, C_0 , and A_0 be positive constants and let*

$$(78) \quad \phi_g \in H^1(0, T), \quad C_0 t^\alpha \leq \phi_g \leq A_0 \quad \text{and} \quad \log \phi_g \in L^2(0, T).$$

Then problem (1)–(3) has at least one weak solution $\phi \in H^1((0, T) \times (0, L))$ such that $\sqrt{\phi} \partial_x (e^{-m\phi} \partial_t \phi) \in L^2((0, T) \times (0, L))$ and $C_0 t^\alpha \leq \phi \leq A_0$.

Proof.

Step 1 (a priori estimates uniform in δ). We test (74) by $\Psi_\delta(\phi_\delta) - \Psi_\delta(\phi_g)$ and get

$$\begin{aligned} & \int_0^t \int_0^L \partial_t \phi_\delta \Psi_\delta(\phi_\delta) \, dx \, d\tau + \int_0^t \int_0^L \frac{h(\phi_\delta)}{\phi_\delta + \delta} (\partial_x \phi_\delta)^2 \, dx \, d\tau \\ & + \frac{1}{B} \int_0^t \int_0^L D(\phi_\delta) (\phi_\delta + \delta) \partial_t (h(\phi_\delta) \partial_x \phi_\delta) \frac{h(\phi_\delta) \partial_x \phi_\delta}{D(\phi_\delta) (\phi_\delta + \delta)} \, dx \, d\tau \\ & = \int_0^t \int_0^L \partial_t \phi_\delta \Psi_\delta(\phi_g) \, dx \, d\tau. \end{aligned}$$

This yields

$$(79) \quad \int_0^L \left(\int_0^{\phi_\delta(t)} \Psi_\delta(\xi) \, d\xi + \frac{1}{2B} (h(\phi_\delta) \partial_x \phi_\delta)^2 \right) dx + \int_0^t \int_0^L \frac{h(\phi_\delta)}{\phi_\delta + \delta} (\partial_x \phi_\delta)^2 \, dx \, d\tau$$

$$= \int_0^t \int_0^L \partial_t \phi_\delta \Psi_\delta(\phi_g) \, dx \, d\tau.$$

In order to get a useful estimate we should find a bound for the first term on the left-hand side of (79). First we note that $\int_0^{\phi_\delta} \int_0^\xi \frac{1}{\eta + \delta} \left(\frac{h(\eta)}{D(\eta)} - \frac{h(0)}{D(0)} \right) d\eta \, d\xi$ defines a continuous function of ϕ_δ . Since ϕ_δ takes values between 0 and A_0 , it is bounded independently of δ . Hence

$$(80) \quad \left| \int_0^L \int_0^{\phi_\delta(t)} \Psi_\delta(\xi) d\xi dx \right| \leq \int_0^L \frac{h(0)}{D(0)} |\{\phi_\delta + \delta\} \log\{\phi_\delta + \delta\} - \phi_\delta - \delta \log \delta| dx + C.$$

Next $(\phi_\delta(t) + \delta) \log(\phi_\delta(t) + \delta) - \phi(t) - \delta \log \delta$ takes value zero at $t = 0$. It is a continuous function of ϕ_δ . Obviously $|(\phi(t) + \delta) \log(\phi(t) + \delta) - \phi(t) - \delta \log \delta| \leq \max\{1 - \delta + \delta \log \delta, (A_0 + \delta) \log(A_0 + \delta) - A_0 - \delta \log \delta\}$, and it is uniformly bounded with respect to δ .

With (80), (79) leads to

$$(81) \quad \int_0^t \int_0^L \frac{h(\phi_\delta)}{\phi_\delta + \delta} (\partial_x \phi_\delta)^2 dx d\tau \leq C + \left| \int_0^t \int_0^L \partial_t \phi_\delta \Psi_\delta(\phi_g) dx d\tau \right|.$$

Next we test (74) by $h(\phi_\delta) \partial_t \phi_\delta - h(\phi_g) \partial_t \phi_g$ and get

$$\begin{aligned} & \int_0^t \int_0^L h(\phi_\delta) (\partial_\tau \phi_\delta)^2 dx d\tau + \int_0^t \int_0^L D(\phi_\delta) \partial_x \phi_\delta \partial_x (h(\phi_\delta) \partial_\tau \phi_\delta) dx d\tau \\ & + \frac{1}{B} \int_0^t \int_0^L D(\phi_\delta) (\phi_\delta + \delta) (\partial_x (h(\phi_\delta) \partial_\tau \phi_\delta))^2 dx d\tau = \int_0^t \int_0^L \partial_t \phi_\delta h(\phi_g) \partial_\tau \phi_g dx d\tau, \end{aligned}$$

and from this

$$(82) \quad \begin{aligned} & \int_0^t \int_0^L h(\phi_\delta) (\partial_\tau \phi_\delta)^2 dx d\tau + \frac{1}{B} \int_0^t \int_0^L D(\phi_\delta) (\phi_\delta + \delta) (\partial_x (h(\phi_\delta) \partial_\tau \phi_\delta))^2 dx d\tau \\ & \leq B \int_0^t \int_0^L \frac{D(\phi_\delta)}{\phi_\delta + \delta} (\partial_x \phi_\delta)^2 dx d\tau + \int_0^t \int_0^L \frac{h^2(\phi_g)}{h(\phi_\delta)} (\partial_\tau \phi_g)^2 dx d\tau. \end{aligned}$$

Let $h_{\min} = e^{-mA_0}$. Then inserting (81) into (82) yields

$$\begin{aligned} & \int_0^t \int_0^L h(\phi_\delta) (\partial_\tau \phi_\delta)^2 dx d\tau + \frac{1}{B} \int_0^t \int_0^L D(\phi_\delta) (\phi_\delta + \delta) (\partial_x (h(\phi_\delta) \partial_\tau \phi_\delta))^2 dx d\tau \\ & \leq C + \frac{BD_r}{h_{\min}} \left| \int_0^t \int_0^L \partial_\tau \phi_\delta \Psi_\delta(\phi_g) dx d\tau \right| + \int_0^t \int_0^L \frac{h^2(\phi_g)}{h(\phi_\delta)} (\partial_\tau \phi_g)^2 dx d\tau \leq C \\ & \quad + \frac{1}{2} \int_0^t \int_0^L h(\phi_\delta) (\partial_\tau \phi_\delta)^2 dx d\tau + \frac{B^2(D_r)^2}{2h_{\min}^3} \|\Psi_\delta(\phi_g)\|_{L^2((0,t) \times (0,L))}^2 \\ & \quad + \frac{1}{h_{\min}} \int_0^t \int_0^L (\partial_\tau \phi_g)^2 dx d\tau. \end{aligned}$$

Step 2 (weak and strong compactness). From the above a priori estimate and assumptions (78) on ϕ_g , we conclude that

$$(83) \quad \|\partial_t \phi_\delta\|_{L^2((0,T) \times (0,L))} + \left\| \frac{1}{\sqrt{\phi_\delta + \delta}} \partial_x \phi_\delta \right\|_{L^2((0,T) \times (0,L))} \leq C,$$

$$(84) \quad \left\| \sqrt{\phi_\delta + \delta} \partial_x (h(\phi_\delta) \partial_t \phi_\delta) \right\|_{L^2((0,T) \times (0,L))} \leq C.$$

Hence there are a $\phi \in H^1((0, T) \times (0, L))$ and a subsequence $\{\phi_\delta\}$, denoted by the same subscripts, such that

$$(85) \quad \phi_\delta \rightarrow \phi \quad \text{strongly in } L^2((0, T) \times (0, L)) \text{ and a.e. on } (0, T) \times (0, L),$$

$$(86) \quad \partial_t \phi_\delta \rightharpoonup \partial_t \phi \quad \text{weakly in } L^2((0, T) \times (0, L)),$$

$$(87) \quad \partial_x \phi_\delta \rightharpoonup \partial_x \phi \quad \text{weakly in } L^2((0, T) \times (0, L)).$$

With the part of the flux containing the second order operator, the situation is more complicated. Obviously, there is $F \in L^2((0, T) \times (0, L))$ such that

$$(88) \quad \sqrt{\phi_\delta + \delta} \partial_{xt} \int_0^{\phi_\delta} h(\xi) d\xi \rightharpoonup F \quad \text{weakly in } L^2((0, T) \times (0, L)).$$

Using the lower bound $\phi_\delta \geq C_0 t^\alpha$, we get from the estimate (84) and convergence (85) that

$$(89) \quad \partial_{xt} \int_0^{\phi_\delta} h(\xi) d\xi \rightharpoonup \partial_{xt} \int_0^\phi h(\xi) d\xi \quad \text{weakly in } L^2((0, T) \times (0, L)).$$

The convergences (85) and (89) imply that F in (88) is given by $F = \sqrt{\phi} \partial_{xt} \int_0^\phi h(\xi) d\xi$.
Step 3 (passing to the limit). Consequently for every $g \in L^2(0, T; V)$ we have

$$(90) \quad \int_0^T \int_0^L \partial_t \phi_\delta g \, dx \, dt \rightarrow \int_0^T \int_0^L \partial_t \phi g \, dx \, dt \quad \text{for } \delta \rightarrow 0,$$

$$(91) \quad \int_0^T \int_0^L D(\phi_\delta) \partial_x \phi_\delta \partial_x g \, dx \, dt \rightarrow \int_0^T \int_0^L D(\phi) \partial_x \phi \partial_x g \, dx \, dt \quad \text{for } \delta \rightarrow 0,$$

$$(92) \quad \int_0^T \int_0^L \frac{D(\phi_\delta)}{B} (\phi_\delta + \delta) \partial_x (h(\phi_\delta) \partial_t \phi_\delta) \partial_x g \, dx \, dt \\ \rightarrow \int_0^T \int_0^L \frac{D(\phi)}{B} \phi \partial_x (h(\phi) \partial_t \phi) \partial_x g \, dx \, dt \quad \text{for } \delta \rightarrow 0.$$

Furthermore, for every $\zeta \in C^\infty([0, L] \times [0, T])$, such that $\zeta(L, t) = 0$ on $[0, T]$, we have

$$\int_0^T \phi|_{x=0} \zeta|_{x=0} \, dt = \int_0^T \phi_g(t) \zeta|_{x=0} \, dt - \int_0^T \int_0^L \partial_x ((\phi - \phi_\delta) \zeta) \, dx \, dt,$$

and, using the convergences (85) and (87), we obtain that the trace of ϕ at $x = 0$ satisfies the boundary condition (26). Hence, we conclude that ϕ satisfies the system (1)–(3). \square

Acknowledgment. The authors are grateful to the (anonymous) referees for their careful reading of the manuscript and helpful remarks.

REFERENCES

[1] T. ALFREY, E. F. GURNEE, AND W. G. LLOYD, *Diffusion in glassy polymers*, J. Polym. Sci. Part C, 12 (1966), pp. 249–261.
 [2] A. S. ARGON, R. E. COHEN, AND A. C. PATEL, *A mechanistic model of case-II diffusion of a diluent into a glassy polymer*, Polymer, 40 (1999), pp. 6991–7012.
 [3] G. I. BARENBLATT, M. BERTSCH, R. DAL PASSO, AND M. UGHI, *A degenerate pseudoparabolic regularization of a nonlinear forward-backward heat equation arising in the theory of heat and mass exchange in stably stratified turbulent shear flow*, SIAM J. Math. Anal., 24 (1993), pp. 1414–1439.

- [4] A. BELIAEV, *Homogenization of two-phase flows in porous media with hysteresis in the capillary relation*, European J. Appl. Math., 14 (2003), pp. 61–84.
- [5] A. YU. BELIAEV AND S. M. HASSANIZADEH, *A theoretical model of hysteresis and dynamic effects in the capillary relation for two-phase flow in porous media*, Transp. Porous Media, 43 (2001), pp. 487–510.
- [6] R. B. BIRD, R. C. AMSTRONG, AND O. HASSAGER, *Dynamics of Polymeric Liquids, Volume 1: Fluid Mechanics*, 2nd ed., John Wiley, New York, 1987.
- [7] C. CUESTA AND J. HULSHOF, *A model problem for groundwater flow with dynamic capillary pressure: Stability of travelling waves*, Nonlinear Anal., 52 (2003), pp. 1199–1218.
- [8] C. CUESTA, C. J. VAN DUJIN, AND J. HULSHOF, *Infiltration in porous media with dynamic capillary pressure: Travelling waves*, European J. Appl. Math., 11 (2000), pp. 381–397.
- [9] W. P. DÜLL, *Some qualitative properties of solutions to a pseudoparabolic equation modeling solvent uptake in polymeric solids*, Comm. Partial Differential Equations, 31 (2006), pp. 1117–1138.
- [10] C. J. VAN DUJIN, L. A. PELETIER, AND I. S. POP, *A new class of entropy solutions of the Buckley–Leverett equation*, SIAM J. Math. Anal., 39 (2007), pp. 507–536.
- [11] L. C. EVANS, *Partial Differential Equations*, Grad. Stud. Math. 19, AMS, Providence, RI, 1998.
- [12] L. C. EVANS, *A survey of entropy methods for partial differential equations*, Bull. Amer. Math. Soc. (N.S.), 41 (2004), pp. 409–438.
- [13] P. D. GAMSON, B. B. BEAMISH, AND D. P. JOHNSON, *Coal microstructure and microporosity and their effects on natural-gas recovery*, Fuel, 72 (1993), pp. 87–99.
- [14] S. M. HASSANIZADEH AND W. G. GRAY, *Thermodynamic basis of capillary pressure in porous media*, Water Resources Res., 29 (1993), pp. 3389–3405.
- [15] D. HENRY, *Geometric Theory of Semilinear Parabolic Equations*, Lecture Notes in Math. 840, Springer-Verlag, Berlin, 1993.
- [16] C. Y. HUI, K. C. WU, R. C. LASKY, AND E. J. KRAMER, *Case II diffusion in polymers I: Transient swelling*, J. Appl. Phys., 61 (1987), pp. 5129–5136.
- [17] C. Y. HUI, K. C. WU, R. C. LASKY, AND E. J. KRAMER, *Case II diffusion in polymers II: Steady-state front motion*, J. Appl. Phys., 61 (1987), pp. 5137–5149.
- [18] J. HULSHOF AND J. R. KING, *Analysis of a Darcy flow model with a dynamic pressure saturation relation*, SIAM J. Appl. Math., 59 (1998), pp. 318–346.
- [19] D. JOU, J. CASAS-VASQUEZ, AND G. LEBON, *Extended Irreversible Thermodynamics*, Springer-Verlag, Heidelberg, 2001.
- [20] J. R. KING AND C. M. CUESTA, *Small- and waiting-time behavior of a Darcy flow model with a dynamic pressure saturation relation*, SIAM J. Appl. Math., 66 (2006), pp. 1482–1511.
- [21] L. D. LANDAU AND E. M. LIFSHITZ, *Fluid Mechanics*, Pergamon Press, Oxford, London, Paris, Frankfurt, 1959.
- [22] C. M. MARLE, *On macroscopic equations governing multiphase flow with diffusion and chemical reactions in porous media*, Internat. J. Engrg. Sci., 20 (1982), pp. 643–662.
- [23] S. MAZUMDER, *Dynamics of CO₂ in Coal as a Reservoir*, Doctoral thesis, TU-Delft Civil Engineering and Geosciences, Delft, The Netherlands, 2007.
- [24] A. MIKELIĆ AND R. ROBERT, *On the equations describing a relaxation toward a statistical equilibrium state in the two-dimensional perfect fluid dynamics*, SIAM J. Math. Anal., 29 (1998), pp. 1238–1255.
- [25] A. NOVICK-COHEN AND R. L. PEGO, *Stable patterns in a viscous diffusion equation*, Trans. Amer. Math. Soc., 324 (1991), pp. 331–351.
- [26] V. PADRON, *Sobolev regularization of a nonlinear ill-posed parabolic problem as a model for aggregating populations*, Comm. Partial Differential Equations, 23 (1998), pp. 457–486.
- [27] D. R. PAVONE, *Macroscopic equations derived from space averaging for immiscible two-phase flow in porous media*, Rev. Inst. Français Pétrole, 44 (1989), pp. 29–41.
- [28] D. R. PAVONE, *A Darcy’s law extension and a new capillary pressure equation for two-phase flow in porous media*, in Proceedings of the 65th Annual Technical Conference of SPE, paper SPE 20474, 1990.
- [29] L. P. RITGER AND N. A. PEPPAS, *Transport of penetrants in the macromolecular structure of coals IV: Models for analysis of dynamic penetrant transport*, Fuel, 66 (1987), pp. 815–826.
- [30] L. P. RITGER AND N. A. PEPPAS, *Transport of penetrants in the macromolecular structure of coals VII: Transport in thin coal sections*, Fuel, 66 (1987), pp. 1379–1388.
- [31] N. L. THOMAS AND A. H. WINDLE, *A theory of case II diffusion*, Polymer, 23 (1982), pp. 529–542.
- [32] D. W. VAN KREVELEN, *Coal: Typology-Physics-Chemistry-Constitution*, 3rd ed., Elsevier, Amsterdam, 1993.
- [33] T. P. WITELSKI, *Traveling wave solutions for case II diffusion in polymers*, J. Polym. Sci. B: Polymer Physics, 34 (1996), pp. 141–150.

ASYMPTOTIC STABILITY OF THE STATIONARY SOLUTION FOR A HYPERBOLIC FREE BOUNDARY PROBLEM MODELING TUMOR GROWTH*

SHANGBIN CUI†

Abstract. In this paper we study the asymptotic behavior of solutions for a free boundary problem modeling the growth of tumors containing two species of cells: proliferating cells and quiescent cells. This tumor model was proposed by Pettet, Please, and McElwain [*Bull. Math. Biol.*, 63 (2001), pp. 231–257]. By using a functional approach and the C_0 semigroup theory, we prove that the unique stationary solution of this model ensured by the work of Cui and Friedman [*Trans. Amer. Math. Soc.*, 355 (2003), pp. 3537–3590] is locally asymptotically stable in certain function spaces. Key techniques used in the proof include an improvement of the linear estimate obtained by the work of Chen, Cui, and Friedman [*Trans. Amer. Math. Soc.*, 357 (2005), pp. 4771–4804] and a similarity transformation.

Key words. free boundary problem, hyperbolic equations, tumor growth, stationary solution, asymptotic stability

AMS subject classifications. 35C10, 35Q80, 92C15

DOI. 10.1137/080717778

1. Introduction. During the past thirty years, an increasing number of free boundary problems of partial differential equations have been proposed by groups of researchers to model the growth of various in vivo and in vitro tumors; cf. the reviewing articles [1], [15], [16], and [22] and the references cited therein. Such free boundary problems usually contain one or more reaction-diffusion equations describing the distribution of nutrient and inhibitory materials, and several first-order nonlinear partial differential equations or nonlinear conservation laws with source terms describing the evolution and movement of various tumor cells (proliferating cells, quiescent cells, and dead cells). Rigorous analysis of such tumor models is evidently a significant topic of research and has drawn great attention during the past few years. The main concern of this topic is the dynamics or the long-term behavior of the solutions of such free boundary problems.

Based on the applications of the well-established theories of elliptic and parabolic partial differential equations, parabolic differential equations in Banach spaces (i.e., differential equations in Banach spaces that are treatable with the analytic semigroup theory), and the bifurcation theory, rigorous analysis of models for the growth of tumors containing only one species of tumor cells has achieved great success; cf. [2], [3], [6], [7], [9], [10], [11], [17], [18], [19], [23], [24], and the references cited therein. As far as models for tumors containing more than one species of tumor cells are concerned, however, the results are much less. This is caused by the fact that such tumor models are much more difficult to analyze because they contain nonlinear conservation laws whose dynamical behavior is very hard to grasp.

In this paper we study the following free boundary problem modeling the growth of an in vitro tumor containing two species of cells—proliferating cells and quiescent

*Received by the editors March 5, 2008; accepted for publication (in revised form) July 23, 2008; published electronically November 26, 2008. This work is supported by the China National Natural Science Foundation under grant number 10471157 and funds in Sun Yat-Sen University.

<http://www.siam.org/journals/sima/40-4/71777.html>

†Institute of Mathematics, Sun Yat-Sen University, Guangzhou, Guangdong 510275, People's Republic of China (cuisb3@yahoo.com.cn).

cells:

$$(1.1) \quad \nabla^2 C = F(C) \quad \text{for } x \in \Omega(t), \quad t \geq 0,$$

$$(1.2) \quad C = C_0 \quad \text{for } x \in \partial\Omega(t), \quad t \geq 0,$$

$$(1.3) \quad \frac{\partial P}{\partial t} + \nabla \cdot (\vec{u}P) = [K_B(C) - K_Q(C)]P + K_P(C)Q \quad \text{for } x \in \Omega(t), \quad t \geq 0,$$

$$(1.4) \quad \frac{\partial Q}{\partial t} + \nabla \cdot (\vec{u}Q) = K_Q(C)P - [K_D(C) + K_P(C)]Q \quad \text{for } x \in \Omega(t), \quad t \geq 0,$$

$$(1.5) \quad P + Q = N \quad \text{for } x \in \Omega(t), \quad t \geq 0,$$

$$(1.6) \quad \frac{dR}{dt} = \vec{u} \cdot \vec{\nu} \quad \text{for } x \in \partial\Omega(t), \quad t \geq 0.$$

Here C denotes the concentration of nutrient (with all nutrient materials regarded as one species), P and Q denote the densities of proliferating cells and quiescent cells, respectively, whose mixture makes up the tumor tissue and has a constant density N , \vec{u} denotes the velocity of the cell movement, R denotes the radius of the tumor, $\Omega(t) = \{x \in \mathbb{R}^3 : r = |x| < R(t)\}$ is the domain occupied by the tumor at time t , and $\vec{\nu}$ is the unit outward normal of $\partial\Omega(t)$. Besides, C_0 is a positive constant reflecting the constant nutrient supply that the tumor receives from its surface, $F(C)$ is the nutrient consumption rate function, and $K_B(C)$, $K_D(C)$, $K_P(C)$, and $K_Q(C)$ are the birth rate of proliferating cells, the death rate of quiescent cells, the transferring rate of proliferating cells to quiescent cells, and the transferring rate of quiescent cells to proliferating cells, respectively. We shall only consider radially symmetric solutions of the above problem, so that C , P , and Q are functions of the radial space variable $r = |x|$ and the time variable t , and $\vec{u} = u(r, t)r^{-1}x$, where u is a scalar function. Note that in the above problem, the unknowns are C , P , Q , u , and R , and all of the rest of the constants and functions are given.

The above tumor model was proposed by Pettet et al. in the literature [21]. Its global well-posedness has been established by Cui and Friedman in [12]. A challenging task concerning this free boundary problem is the study of the asymptotic behavior of its solutions as time goes to infinity. For the corresponding model of the growth of tumors with one species of cells, it is known that there exists a unique stationary solution and that all time-dependent solutions converge to it as time goes to infinity, or, in other words, this unique stationary solution is globally asymptotically stable; cf. [7] and [18]. Since the above problem is a natural extension of such one species tumor model to the two species case, we are naturally led to the conjecture that a similar result holds for it. The advancement of the study toward this goal is as follows. In [13], Cui and Friedman proved that the problem (1.1)–(1.6) has a unique stationary solution. In [5], Chen, Cui, and Friedman further proved that this stationary solution is linearly asymptotically stable, namely, the trivial solution of the linearization of (1.1)–(1.6) at the stationary solution is asymptotically stable. However, this last-mentioned result does not imply that the stationary solution of (1.1)–(1.6) is asymptotically stable. In fact, the asymptotic stability of this stationary solution of (1.1)–(1.6), which is one of a number of interesting open problems raised in [15] and [16], is a very difficult problem due to a number of reasons. Firstly, since the problem (1.1)–(1.6) is of the hyperbolic type, not of the parabolic type, there is not a well-developed geometric theory like that for parabolic problems used in [6], [7], [9], [10], [11], [23], [24] to study this problem. Secondly, the stationary solution of this problem does not have an explicit expression, and, in particular, it possibly possesses a singularity at the point $r = 0$ (see Lemma 3.1 below for this point). It follows that

the method of linearized stability as in [3] and [17] (tempted by [5]) does not work. Thirdly, the commonly used characteristic method for hyperbolic conservation laws also does not apply to this problem, because all of its characteristic curves approach to the same point $r = 0$ as $t \rightarrow \infty$, and, consequently, analysis along characteristic curves does not give us much information of the asymptotic behavior of the solution. Finally, it is our experience that the other commonly used methods for investigating the asymptotic behavior of the solutions of partial differential equations, such as upper and lower solutions and so on (as either explicitly or implicitly discussed in [5]), also do not work. In order to solve this problem, we have to use some new ideas and develop some new techniques.

In this paper we shall prove that the unique stationary solution of (1.1)–(1.6) ensured by [13] is locally asymptotically stable. Recall that conditions given in [13], which ensure that (1.1)–(1.6) has a unique stationary solution are as follows:

$$(1.7) \quad F(C), K_B(C), K_D(C), K_P(C) \text{ and } K_Q(C) \text{ are analytic for } 0 \leq C \leq C_0;$$

$$(1.8) \quad F(0) = 0, \quad F'(C) > 0 \quad \text{for } 0 \leq C \leq C_0;$$

$$(1.9) \quad \begin{cases} K'_B(C) > 0 \text{ and } K'_D(C) < 0 \text{ for } 0 \leq C \leq C_0, \quad K_B(0) = 0 \text{ and } K_D(C_0) = 0; \\ K_P(C) \text{ and } K_Q(C) \text{ satisfy the same conditions as } K_B(C) \text{ and } K_D(C), \\ \text{respectively;} \\ K'_B(C) + K'_D(C) > 0 \text{ for } 0 \leq C \leq C_0. \end{cases}$$

The main result of this paper is the following.

THEOREM 1.1. *Assume that the conditions (1.7)–(1.9) are satisfied. Let $(C_*, P_*, Q_*, \vec{u}_*, R_*)$ be the unique stationary solution of the problem (1.1)–(1.6), and let (C, P, Q, \vec{u}, R) be a time-dependent solution of it such that $P|_{t=0} = P_0, Q|_{t=0} = Q_0$, and $R|_{t=0} = R_0$, where P_0, Q_0 , and R_0 are given initial data satisfying $0 \leq P_0 \leq N, 0 \leq Q_0 \leq N$, and $P_0 + Q_0 = N$. Then there exist positive constants μ, ε , and K such that if P_0, Q_0 , and R_0 satisfy*

$$\max_{0 \leq r \leq 1} \{|P_0(rR_0) - P_*(rR_*)| + |Q_0(rR_0) - Q_*(rR_*)|\} < \varepsilon,$$

$$\sup_{0 < r < 1} r(1-r) \left\{ \left| \frac{d}{dr} (P_0(rR_0) - P_*(rR_*)) \right| + \left| \frac{d}{dr} (Q_0(rR_0) - Q_*(rR_*)) \right| \right\} < \varepsilon,$$

and $|R_0 - R_*| < \varepsilon$; then, for all $t \geq 0$, we have

$$\max_{0 \leq r \leq 1} \{|P(rR(t), t) - P_*(rR_*)| + |Q(rR(t), t) - Q_*(rR_*)|\} < K\varepsilon e^{-\mu t},$$

$$\sup_{0 < r < 1} r(1-r) \left\{ \left| \frac{\partial}{\partial r} (P(rR(t), t) - P_*(rR_*)) \right| + \left| \frac{\partial}{\partial r} (Q(rR(t), t) - Q_*(rR_*)) \right| \right\} < K\varepsilon e^{-\mu t},$$

and $|R(t) - R_*| < K\varepsilon e^{-\mu t}$.

We shall use a functional approach to prove the above theorem. More precisely, we shall first reduce the problem (1.1)–(1.6) into a differential equation for the unknown $U = (p, z)$ in the Banach space $X = C[0, 1] \times \mathbb{R}$, where $p = p(r, t) = P(rR(t), t)$ and $z = z(t) = \log R(t)$. The reduced equation is of the hyperbolic type in the sense of Pazy [20] and is quasi-linear. We next use the Banach fixed point theorem to prove that, for any $U_0 = (p_0, z_0)$ sufficiently closed to the stationary point $U_* = (p_*, z_*)$, where $p_* = p_*(r) = P_*(rR_*)$ and $z_* = \log R_*$, this differential equation imposed, with the initial condition $U|_{t=0} = U_0$ and the decay estimate $\sup_{t \geq 0} e^{\mu t} \|U(t) - U_*\|_{X_0} < \infty$, where X_0 is a subspace of X , has a unique solution in the space $C([0, \infty), X_0)$

(endowed with the norm $\|U\| = \sup_{t \geq 0} e^{\mu t} \|U(t)\|_{X_0}$). To attain this goal we shall use some abstract results for hyperbolic differential equations in Banach spaces presented in [20]. In particular, a family of evolution systems for the linear equations related to the semilinearization of the reduced equation are obtained and applied to convert the semilinearized equations into integral equations. The main difficult and key step in the proof of Theorem 1.1 is the establishment of a uniform decay estimate for the family of evolution systems. In order to obtain such an estimate, we first use a localization technique to get an improvement of the linear estimate established in [5], removing the singularities at $r = 0$ contained in that estimate, and next use a special technique—similarity transformation—to extend this improved linear estimate to the family of evolution systems mentioned above. The similarity transformation technique is the core of this paper. It enables us to transform a transport equation of the form $\partial u / \partial t + w(r, t) \partial u / \partial r = f(r, t)$, where $0 < r < 1$ and $w(0, t) = w(1, t) = 0$, into a similar equation of the form $\partial \tilde{u} / \partial t + w_*(r) \partial \tilde{u} / \partial r = \tilde{f}(r, t)$, where $w_*(0) = w_*(1) = 0$. See section 5 for details of this transformation.

An interesting question raised by one of the referees is: can the above result be extended to the radially nonsymmetric version of the problem (1.1)–(1.6)? Up to now we do not have an answer to this question. A basic difficulty encountered in answering this question is that P_*, Q_*, \vec{u}_* are possibly singular at the center point $r = 0$. More precisely, we know only that P_* is continuous at $r = 0$, but we do not know if it is smooth at $r = 0$ (in fact, we even do not know if it is differentiable at $r = 0$). This possibly existing singularity of the stationary solution induces many obstacles in the analysis of its asymptotic stability for nonradial perturbations. Our conjecture is that P_*, Q_*, \vec{u}_* are smooth for all $r \leq R_*$, and a similar result as that in [11] proved for the one-species model also holds for the radially nonsymmetric version of the problem (1.1)–(1.6).

The structure of the rest part is as follows. In the following section we reduce the problem (1.1)–(1.6) into a differential equation in the Banach space $X = C[0, 1] \times \mathbb{R}$. In section 3 we summarize some basic properties of the stationary solution. In section 4 we study a linear equation obtained from semilinearizing the reduced equation and prove that its solution operator is an evolution system so that the semilinearized equation can be converted into an integral equation. Section 5 aims at developing the similarity transformation technique mentioned above. In section 6 we first derive an improvement of the linear estimate established in [5] and next use the similarity transformation technique to extend this estimate to the evolution systems obtained in section 4. After these preparations, in the last section we use the Banach fixed point theorem to prove Theorem 1.1.

Throughout this paper the notation “'” denotes both the ordinary derivatives of functions in \mathbb{R} and the Fréchet derivatives of mappings between Banach spaces.

Finally, to end this introduction, we would like to refer the reader to see [4], [8], and [14] for other work related with the problem (1.1)–(1.6).

2. Reduction of the problem. In this section we reduce the system of equations (1.1)–(1.6) into a differential equation in the Banach space $X = C[0, 1] \times \mathbb{R}$.

We first note that, by summing up (1.3), (1.4), and using (1.5), we get the following equation:

$$(2.1) \quad \nabla \cdot \vec{u} = \frac{1}{N} [K_B(C)P - K_D(C)Q].$$

Conversely, from (1.3), (1.5), and (2.1) we immediately obtain (1.4). Hence, the two groups of equations (1.3), (1.4), (1.5), and (1.3), (1.5), (2.1) are equivalent. We recall that $\vec{u} = u(r, t)r^{-1}x$ for some scalar function u .

By rescaling the space and time variables, setting

$$\begin{aligned}\bar{c}(\bar{r}, t) &= \frac{C(\bar{r}e^{z(t)}, t)}{C_0}, & \bar{p}(\bar{r}, t) &= \frac{P(\bar{r}e^{z(t)}, t)}{N}, \\ \bar{u}(\bar{r}, t) &= u(\bar{r}e^{z(t)}, t)e^{-z(t)}, & R(t) &= e^{z(t)},\end{aligned}$$

and using the equivalence of (1.3), (1.4), and (1.5) with (1.3), (1.5), and (2.1), we see that the problem (1.1)–(1.6) can be reformulated into the following problem (for the simplicity of the notation we omit all bars):

$$(2.2) \quad \frac{\partial^2 c}{\partial r^2} + \frac{2}{r} \frac{\partial c}{\partial r} = e^{2z} F(c) \quad \text{for } 0 < r \leq 1, \quad t \geq 0,$$

$$(2.3) \quad \left. \frac{\partial c}{\partial r} \right|_{r=0} = 0, \quad c|_{r=1} = 1 \quad \text{for } t \geq 0,$$

$$(2.4) \quad \frac{\partial p}{\partial t} + [u(r, t) - ru(1, t)] \frac{\partial p}{\partial r} = K_P(c) + [K_M(c) - K_N(c)]p - K_M(c)p^2$$

for $0 \leq r \leq 1, \quad t \geq 0,$

$$(2.5) \quad \frac{\partial u}{\partial r} + \frac{2}{r}u = -K_D(c) + K_M(c)p \quad \text{for } 0 < r \leq 1 \text{ and } u|_{r=0} = 0, \quad t \geq 0,$$

$$(2.6) \quad \frac{dz}{dt} = u(1, t) \quad \text{for } t \geq 0,$$

where $K_M(c) = K_B(c) + K_D(c)$, $K_N(c) = K_P(c) + K_Q(c)$, and $F(c)$, $K_B(c)$, $K_D(c)$, $K_P(c)$, and $K_Q(c)$ are the rescaled forms of the corresponding functions appearing in (1.1)–(1.6).

Clearly, (2.2) and (2.3) can be solved to get c as a function of z . Thus, instead of $c(r, t)$, later on we use the notation $c(r, z(t))$, or simply $c(r, z)$, to denote the solution of (2.2) and (2.3). Similarly, (2.5) can be solved to get u as a functional of p and z . Thus, later on we use the notation $u_{p,z}$ to redenote u . By a simple computation, we have

$$(2.7) \quad u_{p,z}(r, t) = \frac{1}{r^2} \int_0^r [-K_D(c(\rho, z(t))) + K_M(c(\rho, z(t)))p(\rho, t)]\rho^2 d\rho$$

for $0 < r \leq 1, t \geq 0$, and $u_{p,z}(0, t) = 0$ for $t \geq 0$. We also denote

$$(2.8) \quad w_{p,z}(r, t) = u_{p,z}(r, t) - ru_{p,z}(1, t).$$

Then (2.2)–(2.6) reduces into the following system of equations:

$$(2.9) \quad \begin{cases} \frac{\partial p}{\partial t} + w_{p,z}(r, t) \frac{\partial p}{\partial r} = f(r, p, z) & \text{for } 0 \leq r \leq 1, \quad t > 0, \\ \frac{dz}{dt} = u_{p,z}(1, t) & \text{for } t > 0, \end{cases}$$

where

$$f(r, p, z) = K_P(c(r, z)) + [K_M(c(r, z)) - K_N(c(r, z))]p - K_M(c(r, z))p^2.$$

In what follows we rewrite (2.9) as a differential equation in the Banach space $X = C[0, 1] \times \mathbb{R}$. To this end we denote $X_0 = C_V^1[0, 1] \times \mathbb{R}$, where $C_V^1[0, 1]$ is the function space

$$C_V^1[0, 1] = \{p \in C[0, 1] \cap C^1(0, 1) : r(1-r)p'(r) \in C[0, 1]\}$$

endowed with the norm

$$\|p\|_{C^1_V[0,1]} = \max_{0 \leq r \leq 1} |p(r)| + \sup_{0 < r < 1} |r(1-r)p'(r)| \quad \text{for } p \in C^1_V[0,1].$$

Clearly, endowed with the product norm, X_0 is a Banach space densely and continuously embedded into X . We introduce a mapping $\mathbb{F} : X_0 \rightarrow X$ as follows. First, for given $p \in C[0,1]$ and $z \in \mathbb{R}$, let $\mathcal{A}_0(p, z) : C^1_V[0,1] \rightarrow C(0,1)$ be the following linear operator: For any $q \in C^1_V[0,1]$,

$$\mathcal{A}_0(p, z)q(r) = -w_{p,z}(r)q'(r) \quad \text{for } 0 < r < 1.$$

Here and hereafter, $w_{p,z}(r)$ represents the function defined by similar formulations as in (2.7) and (2.8), with $p(r, t)$ and $z(t)$ there replaced by $p(r)$ and z , respectively. Later on we shall use the convention that, for a function $f \in C(0,1)$, if both limits $\lim_{r \rightarrow 0^+} f(r)$ and $\lim_{r \rightarrow 1^-} f(r)$ exist and are finite, then we write $f \in C[0,1]$. Furthermore, when we are concerned with the values of f at $r = 0$ and $r = 1$, we mean that $f(0) = \lim_{r \rightarrow 0^+} f(r)$ and $f(1) = \lim_{r \rightarrow 1^-} f(r)$. Using this convention, we see easily that, for any $p \in C[0,1]$ and $z \in \mathbb{R}$, we have $w_{p,z}(r)/r(1-r) \in C[0,1]$. It follows that, for any $q \in C^1_V[0,1]$, both limits $\lim_{r \rightarrow 0^+} w_{p,z}(r)q'(r)$ and $\lim_{r \rightarrow 1^-} w_{p,z}(r)q'(r)$ exist so that $\mathcal{A}_0(p, z)q \in C[0,1]$. It can also be easily seen that $\mathcal{A}_0(p, z) \in L(C^1_V[0,1], C[0,1])$, and

$$\|\mathcal{A}_0(p, z)\|_{L(C^1_V[0,1], C[0,1])} \leq \sup_{0 < r < 1} \left| \frac{w_{p,z}(r)}{r(1-r)} \right|.$$

Next we denote by $\mathcal{F} : C[0,1] \times \mathbb{R} \rightarrow C[0,1]$ and $\mathcal{G} : C[0,1] \times \mathbb{R} \rightarrow \mathbb{R}$, respectively, the following nonlinear operators:

$$\begin{aligned} \mathcal{F}(p, z)(r) &= f(r, p(r), c(r, z)), \\ \mathcal{G}(p, z) &= \int_0^1 [-K_D(c(r, z)) + K_M(c(r, z))p(r)]r^2 dr. \end{aligned}$$

We now define $\mathbb{F} : X_0 \rightarrow X$ to be the following nonlinear operator:

$$\mathbb{F}(U) = (\mathcal{A}_0(p, z)p + \mathcal{F}(p, z), \mathcal{G}(p, z)) \quad \text{for } U = (p, z) \in X_0.$$

It is obvious that $\mathbb{F} \in C^\infty(X_0, X)$. Later on we shall also regard \mathbb{F} as an unbounded nonlinear operator in X with domain X_0 .

With the above notations, we see that (2.9) can be rewritten as the following differential equation in the Banach space X :

$$(2.10) \quad \frac{dU}{dt} = \mathbb{F}(U).$$

Here $U = U(t)$ represents an X_0 -valued unknown function for $t \geq 0$, and the left-hand side denotes the Fréchet derivative of $U = U(t)$ regarded as a mapping from $[0, \infty)$ to the X space.

The equation (2.10) has a quasi-linear structure. To see this we define $\mathbb{A}_0 : X \rightarrow L(X_0, X)$ and $\mathbb{F}_0 : X \rightarrow X$ as follows: For $U = (p, z) \in X$ and $V = (q, y) \in X_0$, we let

$$\mathbb{A}_0(U)V = (\mathcal{A}_0(p, z)q, 0) \quad \text{and} \quad \mathbb{F}_0(U) = (\mathcal{F}(p, z), \mathcal{G}(p, z)).$$

Then we have

$$\mathbb{F}(U) = \mathbb{A}_0(U)U + \mathbb{F}_0(U) \quad \text{for } U \in X_0.$$

Clearly, $\mathbb{A}_0 \in C^\infty(X, L(X_0, X))$ and $\mathbb{F}_0 \in C^\infty(X, X) \cap C^\infty(X_0, X_0)$. Hence, the desired assertion follows. Later on, for given $U \in X$, we shall also regard $\mathbb{A}_0(U)$ as an unbounded linear operator in X with domain X_0 .

From [13] we know that under the assumptions (1.7)–(1.9), the problem (2.2)–(2.6) has a unique stationary solution which we denote as (c_*, p_*, u_*, z_*) . By definition, $(c_*, p_*, u_*, z_*) = (c_*(r), p_*(r), u_*(r), z_*)$ ($0 \leq r \leq 1$) is the solution of the following problem:

$$(2.11) \quad c_*'' + \frac{2}{r}c_*' = e^{2z_*}F(c_*) \quad \text{for } 0 < r \leq 1,$$

$$(2.12) \quad c_*'(0) = 0, \quad c_*(1) = 1,$$

$$(2.13) \quad u_*p_*' = f(r, p_*, z_*) \quad \text{for } 0 \leq r \leq 1,$$

$$(2.14) \quad u_*' + \frac{2}{r}u_* = -K_D(c_*) + K_M(c_*)p_* \quad \text{for } 0 < r \leq 1,$$

$$(2.15) \quad u_*(0) = 0, \quad u_*(1) = 0.$$

Let $U_* = (p_*, z_*)$. Then $U_* \in X_0$ (see Lemma 2.1 below), and, by the equivalence of the system (2.2)–(2.6) with the equation (2.10), it follows that U_* is the unique equilibrium of (2.10), i.e., $\mathbb{F}(U_*) = 0$ or $\mathbb{A}_0(U_*)U_* + \mathbb{F}_0(U_*) = 0$. Hence, to study the asymptotic stability of the stationary solution (c_*, p_*, u_*, z_*) of the problem (2.2)–(2.6), we need only to study the asymptotic stability of the stationary solution U_* of (2.10). For this purpose, we rewrite (2.10) into an equivalent equation for the difference $V = U - U_*$. Let $\mathbb{A} : X \rightarrow L(X_0, X)$ and $\mathbb{G} : X \rightarrow X$ be the following operators:

$$\begin{aligned} \mathbb{A}(V)W &= \mathbb{A}_0(U_* + V)W + [\mathbb{A}'_0(U_*)W]U_* + \mathbb{F}'_0(U_*)W \quad \text{for } V \in X, \quad W \in X_0, \\ \mathbb{G}(V) &= [\mathbb{A}_0(U_* + V) - \mathbb{A}_0(U_*) - \mathbb{A}'_0(U_*)V]U_* \\ &\quad + [\mathbb{F}_0(U_* + V) - \mathbb{F}_0(U_*) - \mathbb{F}'_0(U_*)V] \quad \text{for } V \in X. \end{aligned}$$

Then clearly (2.10) is equivalent to the following equation for $V = U - U_*$:

$$(2.16) \quad \frac{dV}{dt} = \mathbb{A}(V)V + \mathbb{G}(V).$$

Later on we shall concentrate our attention on (2.16).

The above deduction leads to the following preliminary result.

LEMMA 2.1. *The system (2.2)–(2.6) is equivalent to (2.16), and the asymptotic stability of the stationary solution (c_*, p_*, u_*, z_*) of (2.2)–(2.6) is equivalent to the asymptotic stability of the trivial solution of (2.16). Moreover, the assertion of Theorem 1.1 is equivalent to the following assertion for (2.16): There exist positive constants K, μ , and ε such that if the solution $V = V(t)$ of (2.16) satisfies $\|V(0)\|_{X_0} \leq \varepsilon$, then it also satisfies*

$$(2.17) \quad \|V(t)\|_{X_0} \leq K\varepsilon e^{-\mu t} \quad \text{for all } t \geq 0. \quad \square$$

Since it has been proved (see [12]) that the initial value problem of the system (2.2)–(2.6) is globally well-posed for any C^1 initial data, by the equivalence of (2.2)–(2.6) with (2.16), it follows that, for any $V_0 \in X_0$, (2.16) has a unique solution

$V = V(t)$ for all $t \geq 0$ such that $V(0) = V_0$. Our analysis later on is to prove that (2.17) holds provided $\|V_0\|_{X_0} \leq \varepsilon$.

We note that $\mathbb{A} \in C^\infty(X, L(X_0, X))$, $\mathbb{G} \in C^\infty(X, X)$, and, by using the Taylor expansions up to second-order for Fréchet derivatives of \mathbb{A}_0 and \mathbb{F}_0 , we have

$$(2.18) \quad \|\mathbb{G}(V)\|_X = O(\|V\|_X^2) \quad \text{as } \|V\|_X \rightarrow 0.$$

We also note that, by introducing an operator $\mathbb{B} : X \rightarrow X$ by

$$\mathbb{B}W = [\mathbb{A}'_0(U_*)W]U_* + \mathbb{F}'_0(U_*)W \quad \text{for } W \in X,$$

we have

$$\mathbb{A}(0) = \mathbb{F}'(U_*) = \mathbb{A}_0(U_*) + \mathbb{B} \quad \text{and} \quad \mathbb{A}(V) = \mathbb{A}_0(U_* + V) + \mathbb{B}.$$

Note that as an immediate consequence of the facts that $\mathbb{A}_0 \in C^\infty(X, L(X_0, X))$ and $\mathbb{F}_0 \in C^\infty(X, X)$, we have $\mathbb{B} \in L(X)$. Finally, we note that $[V \rightarrow \mathbb{A}_0(U_* + V)] \in C^\infty(X, L(X_0, X))$. These facts will play an important role in later analysis.

A simple computation shows that if we denote

$$(2.19) \quad a(r) = K_M(c_*(r)) - K_N(c_*(r)) - 2K_M(c_*(r))p_*(r),$$

$$b(r) = \{K'_P(c_*(r)) + [K'_M(c_*(r)) - K'_N(c_*(r))]p_*(r) - K'_M(c_*(r))p_*^2(r)\} c_z(r)$$

$$(2.20) \quad + rp'_*(r) \left[\int_0^1 g_c(\rho)c_z(\rho)\rho^2 d\rho - \frac{1}{r^3} \int_0^r g_c(\rho)c_z(\rho)\rho^2 d\rho \right],$$

$$(2.21) \quad \mathcal{B}q(r) = rp'_*(r) \left[\int_0^1 g_p(\rho)q(\rho)\rho^2 d\rho - \frac{1}{r^3} \int_0^r g_p(\rho)q(\rho)\rho^2 d\rho \right],$$

$$(2.22) \quad \mathcal{F}(q) = \int_0^1 g_p(\rho)q(\rho)\rho^2 d\rho,$$

$$(2.23) \quad \kappa = \int_0^1 g_c(\rho)c_z(\rho)\rho^2 d\rho,$$

where

$$g_p(r) = K_M(c_*(r)), \quad g_c(r) = -K'_D(c_*(r)) + K'_M(c_*(r))p_*(r), \quad c_z(r) = \frac{\partial c}{\partial z}(r, z_*),$$

then we have

$$(2.24) \quad \mathbb{B} = \begin{pmatrix} a(r) + \mathcal{B} & b(r) \\ \mathcal{F} & \kappa \end{pmatrix}.$$

Here and hereafter, when we write $\mathbb{M} = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix}$ for bounded linear operators $M_{11} \in L(X_1, Y_1)$, $M_{12} \in L(X_2, Y_1)$, $M_{21} \in L(X_1, Y_2)$, and $M_{22} \in L(X_2, Y_2)$, where X_1, X_2, Y_1 , and Y_2 are Banach spaces, we mean that \mathbb{M} is the bounded linear operator from $X_1 \times X_2$ to $Y_1 \times Y_2$ defined by $\mathbb{M}(x_1, x_2) = (M_{11}x_1 + M_{12}x_2, M_{21}x_1 + M_{22}x_2)$ for $(x_1, x_2) \in X_1 \times X_2$. Using this notation we see that

$$(2.25) \quad \mathbb{A}_0(U_*) = \begin{pmatrix} \mathcal{L}_0 & 0 \\ 0 & 0 \end{pmatrix} \quad \text{and} \quad \mathbb{A}_0(U_* + V) = \begin{pmatrix} \mathcal{L}_V & 0 \\ 0 & 0 \end{pmatrix},$$

where $\mathcal{L}_0 = \mathcal{A}_0(p_*, z_*)$ and $\mathcal{L}_V = \mathcal{A}_0(p_* + \varphi, z_* + \zeta)$ for $V = (\varphi, \zeta) \in X$. Finally, we recall that $a(r) < 0$ for all $0 \leq r \leq 1$; see (2.7) in section 2 of [5].

3. Some basic facts. In this section we summarize some basic properties of the functions $c_*(r) = c(r, z_*)$, $c_z(r) = \frac{\partial c}{\partial z}(r, z_*)$, $p_*(r)$, and $u_*(r)$. These properties will play an important role in later discussions.

LEMMA 3.1. *We have the following assertions:*

(1) $c_*, c_z \in C^\infty[0, 1]$, and

(3.1)

$$0 < c_*(0) \leq c_*(r) \leq 1 \quad \text{for } 0 \leq r \leq 1, \quad c'_*(r) > 0 \quad \text{for } 0 < r \leq 1, \quad c'_*(0) = 0.$$

(2) $p_* \in C[0, 1] \cap C^\infty(0, 1]$,

(3.2)

$$0 < p_*(0) \leq p_*(r) \leq 1 \quad \text{for } 0 \leq r \leq 1, \quad p'_*(r) > 0 \quad \text{for } 0 < r \leq 1,$$

and either $p_* \in C^1[0, 1]$ or there exists $0 < \gamma < 1$ such that $\lim_{r \rightarrow 0^+} r^\gamma p'_*(r)$ exists and is finite, so that $r^\gamma p'_*(r) \in C[0, 1]$. Moreover,

(3.3)

$$\lim_{r \rightarrow 0^+} r p'_*(r) = 0, \quad \lim_{r \rightarrow 0^+} r^2 p''_*(r) = 0, \quad \lim_{r \rightarrow 0^+} r^3 p'''_*(r) = 0.$$

(3) $u_* \in C^1[0, 1] \cap C^\infty(0, 1]$, and there exist positive constants C_1, C_2 such that

(3.4)

$$-C_1 r(1-r) \leq u_*(r) \leq -C_2 r(1-r) \quad \text{for } 0 \leq r \leq 1.$$

Besides, either $u_* \in C^2[0, 1]$ or there exists $0 < \gamma < 1$ such that $\lim_{r \rightarrow 0^+} r^\gamma u''_*(r)$ exists and is finite, so that $r^\gamma u''_*(r) \in C[0, 1]$. Moreover,

(3.5)

$$\lim_{r \rightarrow 0^+} r u''_*(r) = 0, \quad \lim_{r \rightarrow 0^+} r^2 u'''_*(r) = 0.$$

Proof. The assertions that $c_*, c_z \in C^\infty[0, 1]$ and relations in (3.1) are immediate. The assertions that $p_* \in C[0, 1] \cap C^\infty(0, 1]$, $u_* \in C^1[0, 1] \cap C^\infty(0, 1]$, and relations in (3.2) follow from Theorem 2.1 of [13], by which we also know that $u_*(r) < 0$ for $0 < r < 1$. The last assertion combined with the facts that $u'_*(0) < 0$ and $u'_*(1) > 0$ (see Theorem 7.1 of [13]) immediately yields (3.4). To prove (3.3) we compute

$$\lim_{r \rightarrow 0} u_*(r) p'_*(r) = K_P(c_*(0)) + [K_M(c_*(0)) - K_N(c_*(0))] p_*(0) - K_M(c_*(0)) p_*^2(0) = 0$$

(see (8.4) in section 8 of [13]), so that

$$\lim_{r \rightarrow 0} r p'_*(r) = \lim_{r \rightarrow 0} \frac{r}{u_*(r)} \cdot \lim_{r \rightarrow 0} u_*(r) p'_*(r) = 0,$$

and

$$\begin{aligned} \lim_{r \rightarrow 0} u_*(r) r p''_*(r) &= \lim_{r \rightarrow 0} r \{ K'_P(c_*(r)) + [K'_M(c_*(r)) - K'_N(c_*(r))] p_*(r) \\ &\quad - K'_M(c_*(r)) p_*^2(r) \} c'_*(r) + \lim_{r \rightarrow 0} \{ [K_M(c_*(r)) - K_N(c_*(r))] \\ &\quad - 2K_M(c_*(r)) p_*(r) \} r p'_*(r) - \lim_{r \rightarrow 0} u'_*(r) r p'_*(r) = 0, \end{aligned}$$

so that

$$\lim_{r \rightarrow 0} r^2 p''_*(r) = \lim_{r \rightarrow 0} \frac{r}{u_*(r)} \cdot \lim_{r \rightarrow 0} u_*(r) r p''_*(r) = 0.$$

This proves the first two relations in (3.3). The proof of the third relation is similar and is omitted. Next, from (2.14) we can easily deduce that

$$\begin{aligned}
 u_*''(r) &= [-K_D'(c_*(r)) + K_M'(c_*(r))p_*(r)]c_*'(r) \\
 &\quad + K_M(c_*(r))p_*'(r) + \frac{2}{r}[K_D(c_*(r)) - K_M(c_*(r))p_*(r)] \\
 &\quad + \frac{6}{r^4} \int_0^r [-K_D(c_*(\rho)) + K_M(c_*(\rho))p_*(\rho)]\rho^2 d\rho \\
 &= [-K_D'(c_*(r)) + K_M'(c_*(r))p_*(r)]c_*'(r) + \frac{6}{r^4} \int_0^r [K_D(c_*(r)) - K_D(c_*(\rho))]\rho^2 d\rho \\
 (3.6) \quad &\quad + K_M(c_*(r))p_*'(r) - \frac{6}{r^4} \int_0^r [K_M(c_*(r))p_*(r) - K_M(c_*(\rho))p_*(\rho)]\rho^2 d\rho.
 \end{aligned}$$

From this expression and the first relation in (3.3) we readily obtain the first relation in (3.5). The proof of the second relation in (3.5) is similar and is omitted. Finally, by Theorems 5.3 and 5.4 of [13], we know that either $p_* \in C^1[0, 1]$ or there exist constants $-1 < \alpha < 0$ and C such that¹

$$(3.7) \quad p_*'(r) = Cr^\alpha + O(1) \quad \text{for } r \rightarrow 0.$$

Suppose that it is the second case. Then, by letting $\gamma = |\alpha|$, we see that $0 < \gamma < 1$ and $r^\gamma p_*'(r) \in C[0, 1]$. Finally, from (3.7) we see that

$$(3.8) \quad p_*(r) = p_*(0) + C(1 + \alpha)^{-1}r^{1+\alpha} + O(r) \quad \text{for } r \rightarrow 0.$$

Substituting (3.7) and (3.8) into (3.6), we get $r^\gamma u_*''(r) \in C[0, 1]$. This completes the proof of Lemma 3.1. \square

COROLLARY 3.2. *Let $a, b, \mathcal{B}, \mathcal{F}$, and \mathbb{B} be as in (2.19)–(2.22) and (2.24). Then we have $a, b \in C_V^1[0, 1], \mathcal{B} \in L(C[0, 1], C_V^1[0, 1]) \subseteq L(C[0, 1]) \cap L(C_V^1[0, 1]), \mathcal{F} \in L(C[0, 1], \mathbb{R})$, and $\mathbb{B} \in L(X) \cap L(X_0)$. Moreover, we also have $r^2(1-r)^2 a''(r), r^2(1-r)^2 b''(r) \in C[0, 1]$. \square*

COROLLARY 3.3. $\mathbb{G} \in C^\infty(X, X) \cap C^\infty(X_0, X_0)$ and in addition to (2.18), we also have

$$(3.9) \quad \|\mathbb{G}(V)\|_{X_0} = O(\|V\|_{X_0}^2) \quad \text{as } \|V\|_{X_0} \rightarrow 0.$$

Proof. We have $\mathbb{G}(V) = \mathbb{G}_1(V) + \mathbb{G}_2(V)$, where $\mathbb{G}_1(V) = [\mathbb{A}_0(U_* + V) - \mathbb{A}_0(U_*) - \mathbb{A}'_0(U_*)V]U_*$ and $\mathbb{G}_2(V) = \mathbb{F}_0(U_* + V) - \mathbb{F}_0(U_*) - \mathbb{F}'_0(U_*)V$. Since $\mathbb{F}_0 \in C^\infty(X_0, X_0)$, it is evident that $\mathbb{G}_2 \in C^\infty(X_0, X_0)$ and $\|\mathbb{G}_2(V)\|_{X_0} = O(\|V\|_{X_0}^2)$ as $\|V\|_{X_0} \rightarrow 0$. Next, let $V = (\varphi, \zeta)$ and $(p, z) = (p_* + \varphi, z_* + \zeta)$. Then, by (2.25), we have $\mathbb{A}_0(U_* + V)U_* = (-w_{p,z}(r)p_*'(r), 0)$. Using this expression and the first two relations in (3.3), we can easily show that, for every $V \in X$, we have $\mathbb{A}_0(U_* + V)U_* \in X_0$, and the mapping $V \rightarrow \mathbb{A}_0(U_* + V)U_*$ belongs to $C^\infty(X, X_0)$. Hence we have $\mathbb{G}_1 \in C^\infty(X, X_0) \subseteq C^\infty(X_0, X_0)$ and $\|\mathbb{G}_1(V)\|_{X_0} = O(\|V\|_X^2) = O(\|V\|_{X_0}^2)$ as $\|V\|_{X_0} \rightarrow 0$. Combining these assertions together, we see that the desired assertion follows. \square

4. Evolution system. In this section we study the following initial value problem:

$$(4.1) \quad \frac{dU}{dt} = \mathbb{A}(V(t))U + F(t) \quad \text{for } t > 0, \quad U(0) = U_0,$$

¹In the notation of Theorem 5.4 of [13], we have $\alpha = \alpha(\lambda)$ and $C = (1 + \alpha(\lambda))\omega$.

where $V \in C([0, \infty), X)$, $F \in C([0, \infty), X_0)$, and $U_0 \in X_0$ are given. For a small $\varepsilon > 0$, we denote

$$S_\varepsilon = \{V = (\varphi, \zeta) \in X = C[0, 1] \times \mathbb{R} : \|\varphi\|_\infty \leq \varepsilon, |\zeta| \leq \varepsilon\}.$$

We want to prove that given $V \in C([0, \infty), X)$ such that $V(t) \in S_\varepsilon$ for all $t \geq 0$, the problem (4.1) has a unique solution $U \in C([0, \infty), X_0) \cap C^1([0, \infty), X)$, and there exists a family of bounded linear operators $\{\mathbb{U}(t, s, V) : t \geq s \geq 0\}$ on X_0 , the so-called *evolution system* determined by the family of operators $\{\mathbb{A}(V(t)) : t \geq 0\}$, such that the solution of this problem is given by

$$(4.2) \quad U(t) = \mathbb{U}(t, 0, V)U_0 + \int_0^t \mathbb{U}(t, s, V)F(s)ds \quad \text{for } t \geq 0.$$

For this purpose we shall prove that the family of operators $\{\mathbb{A}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on X , and its part in X_0 is a stable family of the infinitesimal generators of C_0 semigroups on X_0 . When these assertions are proved, the desired assertion then follows from the results of sections 5.2–5.5 and 6.4 of [20].

In what follows, for $V \in S_\varepsilon$, we denote by $\tilde{\mathbb{A}}(V)$ the part of $\mathbb{A}(V)$ in X_0 . Recall that $\text{Dom}(\tilde{\mathbb{A}}(V)) = \{U \in X_0 : \mathbb{A}(V)U \in X_0\}$, and $\tilde{\mathbb{A}}(V)U = \mathbb{A}(V)U$ for $U \in \text{Dom}(\tilde{\mathbb{A}}(V))$.

Let $w \in C^1[0, 1]$ and assume that it satisfies the following condition: There exist positive constants C_1 and C_2 such that

$$(4.3) \quad -C_1r(1-r) \leq w(r) \leq -C_2r(1-r) \quad \text{for } 0 \leq r \leq 1.$$

Note that this assumption particularly implies that $w(0) = w(1) = 0$, $w'(0) < 0$ and $w'(1) > 0$. For a such w , we denote by \mathcal{L}_0 the bounded linear operator from $C_V^1[0, 1]$ to $C[0, 1]$ defined by

$$\mathcal{L}_0q(r) = -w(r)q'(r) \quad \text{for } 0 < r < 1, \quad \text{for } q \in C_V^1[0, 1].$$

Later on we shall also regard \mathcal{L}_0 as an unbounded linear operator in $C[0, 1]$ with domain $C_V^1[0, 1]$. Note that if $w = u_*$, then $\mathcal{L}_0 = \mathcal{A}_0(p_*, z_*)$.

LEMMA 4.1. *Let the notation and assumption be as above. Then \mathcal{L}_0 generates a C_0 semigroup of contractions $e^{t\mathcal{L}_0}$ on $C[0, 1]$, i.e.,*

$$(4.4) \quad \|e^{t\mathcal{L}_0}\|_{L(C[0,1])} \leq 1 \quad \text{for } t \geq 0.$$

Moreover, $C_V^1[0, 1]$ is \mathcal{L}_0 -admissible,² and the restriction of $e^{t\mathcal{L}_0}$ on $C_V^1[0, 1]$ is a uniformly bounded C_0 semigroup on $C_V^1[0, 1]$, i.e., there exists constant $C > 0$ depending only on the constants C_1 and C_2 in (4.3) such that

$$(4.5) \quad \|e^{t\mathcal{L}_0}\|_{L(C_V^1[0,1])} \leq C \quad \text{for } t \geq 0.$$

²Recall that for a C_0 semigroup $T(t)$ ($t \geq 0$) on a Banach space X generated by an unbounded linear operator A in X , a linear subspace Y of X is called *A-admissible* if it is an invariant subspace of $T(t)$ for all $t \geq 0$ and the restriction of $T(t)$ ($t \geq 0$) to Y is a C_0 semigroup in Y . A necessary and sufficient condition for Y to be *A-admissible* is that (1) Y is an invariant subspace of $R(\lambda, A)$ for all $\lambda > \omega$, and (2) the part \tilde{A} of A in Y is an infinitesimal generator of a C_0 semigroup on Y . In this case we have $e^{t\tilde{A}} = e^{tA}|_Y$. See Theorem 5.5 in Chapter 4 of [20].

Proof. Let $\lambda \in \mathbb{C}$ be such that $\operatorname{Re}\lambda > 0$, and let $f \in C[0, 1]$. Consider the equation

$$(4.6) \quad -w(r)q'(r) - \lambda q(r) = f(r) \quad \text{for } 0 < r < 1.$$

By using some similar arguments as in section 2 of [8], we can easily show that this equation has a unique solution defined and bounded for $0 < r < 1$, given by

$$(4.7) \quad q(r) = e^{-\lambda \int_{r_0}^r \frac{d\rho}{w(\rho)}} \int_r^1 \frac{f(\eta)}{w(\eta)} e^{\lambda \int_{r_0}^\eta \frac{d\rho}{w(\rho)}} d\eta, \quad 0 < r < 1,$$

where r_0 is an arbitrarily chosen number in $(0, 1)$. Clearly $q \in C^1(0, 1)$. By using the L'Hospital's rule, we can easily verify that $q(r)$ has finite limits as $r \rightarrow 0^+$ and $r \rightarrow 1^-$, so that $q \in C[0, 1] \cap C^1(0, 1)$. Furthermore, since $r(1-r)/w(r) \in C[0, 1]$, from (4.6) we see that also $r(1-r)q'(r) \in C[0, 1]$, so that $q \in C_V^1[0, 1]$, and, moreover,

$$(4.8) \quad \begin{aligned} |q(r)| &\leq \|f\|_\infty e^{(\operatorname{Re}\lambda) \int_{r_0}^r \frac{d\rho}{|w(\rho)|}} \int_r^1 \frac{1}{|w(\eta)|} e^{-(\operatorname{Re}\lambda) \int_{r_0}^\eta \frac{d\rho}{|w(\rho)|}} d\eta \\ &= \frac{\|f\|_\infty}{\operatorname{Re}\lambda} \left[1 - e^{-(\operatorname{Re}\lambda) \int_r^1 \frac{d\rho}{|w(\rho)|}} \right] \leq \frac{\|f\|_\infty}{\operatorname{Re}\lambda} \quad \text{for } 0 < r < 1. \end{aligned}$$

This proves that, for any $\lambda \in \mathbb{C}$ with $\operatorname{Re}\lambda > 0$, we have $\lambda \in \rho(\mathcal{L}_0)$ and $\|R(\lambda, \mathcal{L}_0)\|_{L(C[0,1])} \leq 1/\operatorname{Re}\lambda$. Thus, by the Hille–Yosida theorem, we see that \mathcal{L}_0 generates a strongly continuous semigroup $e^{t\mathcal{L}_0}$ on $C[0, 1]$, which satisfies the estimate (4.4).

Next we prove that $C_V^1[0, 1]$ is \mathcal{L}_0 -admissible. Clearly, $C_V^1[0, 1]$ is an invariant subspace of $R(\lambda, \mathcal{L}_0)$ for all $\lambda \in \mathbb{C}$, with $\operatorname{Re}\lambda > 0$. Let $\tilde{\mathcal{L}}_0$ be the part of \mathcal{L}_0 in $C_V^1[0, 1]$. Since $r(1-r)/w(r)$, $w(r)/r(1-r) \in C[0, 1]$, we see that, for $q \in C[0, 1] \cap C^1(0, 1)$, $r(1-r)q'(r) \in C[0, 1]$ if and only if $w(r)q'(r) \in C[0, 1]$, and there exist positive constants C_1 and C_2 such that

$$C_1 \sup_{0 < r < 1} |w(r)q'(r)| \leq \sup_{0 < r < 1} |r(1-r)q'(r)| \leq C_2 \sup_{0 < r < 1} |w(r)q'(r)|.$$

By the above assertion, it follows easily that

$$(4.9) \quad \operatorname{Dom}(\tilde{\mathcal{L}}_0) = \{q \in C[0, 1] \cap C^2(0, 1) : w(r)q'(r) \in C[0, 1], w^2(r)q''(r) \in C[0, 1]\},$$

and $\|q\|'_{C_V^1[0,1]} = \|q\|_\infty + \|wq'\|_\infty$ is an equivalent norm in $C_V^1[0, 1]$. For $q \in \operatorname{Dom}(\tilde{\mathcal{L}}_0)$, we have, by definition, $\tilde{\mathcal{L}}_0 q = \mathcal{L}_0 q$. Now let $f \in C_V^1[0, 1]$ and $\operatorname{Re}\lambda > 0$. Let q be the solution of (4.6). Using (4.9) we can easily verify that $q \in \operatorname{Dom}(\tilde{\mathcal{L}}_0)$, so that it is the solution of the equation $\tilde{\mathcal{L}}_0 q - \lambda q = f$. Moreover, a simple computation shows that $wq' = R(\mathcal{L}_0, \lambda)(wf')$. Thus, by (4.8), we have

$$\begin{aligned} \|q\|'_{C_V^1[0,1]} &= \|q\|_\infty + \|wq'\|_\infty = \|R(\mathcal{L}, \lambda)f\|_\infty + \|R(\mathcal{L}, \lambda)(wf')\|_\infty \\ &\leq \frac{1}{\operatorname{Re}\lambda} \|f\|_\infty + \frac{1}{\operatorname{Re}\lambda} \|wf'\|_\infty = \frac{1}{\operatorname{Re}\lambda} \|f\|'_{C_V^1[0,1]}. \end{aligned}$$

Hence, $\lambda \in \rho(\tilde{\mathcal{L}}_0)$ and $\|R(\lambda, \tilde{\mathcal{L}}_0)\|'_{C_V^1[0,1]} \leq (\operatorname{Re}\lambda)^{-1}$. The desired assertion then follows from the Hille–Yosida theorem and footnote 2. \square

Given $V = (\varphi, \zeta) \in S_\varepsilon$, we set $p(r) = p_*(r) + \varphi(r)$, $z = z_* + \zeta$, and, as before, we denote

$$u_{p,z}(r) = \frac{1}{r^2} \int_0^r [-K_D(c(\rho, z)) + K_M(c(\rho, z))p(\rho)] \rho^2 d\rho, \quad w_{p,z}(r) = u_{p,z}(r) - ru_{p,z}(1).$$

Since $\|\varphi\|_\infty \leq \varepsilon$ and $|\zeta| \leq \varepsilon$, by a simple computation, we see that

$$(4.10) \quad -C\varepsilon r(1-r) \leq w_{p,z}(r) - u_*(r) \leq C\varepsilon r(1-r) \quad \text{for } 0 \leq r \leq 1.$$

Since $-C_1 r(1-r) \leq u_*(r) \leq -C_2 r(1-r)$, it follows that, for ε sufficiently small, we have

$$(4.11) \quad -C_1 r(1-r) \leq w_{p,z}(r) \leq -C_2 r(1-r) \quad \text{for } 0 \leq r \leq 1.$$

Later on we shall also use the notation $w_V(r)$ to redenote $w_{p,z}(r)$. We note that all constants C , C_1 , and C_2 that appear in (4.10)–(4.11) are independent of V and ε .

LEMMA 4.2. $\{\mathbb{A}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on $X = C[0, 1] \times \mathbb{R}$, and $\{\tilde{\mathbb{A}}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on X_0 .

Proof. Let $\mathcal{L}_V q(r) = -w_V(r)q'(r)$. Then, by Lemma 4.1, we see that, for any $V \in S_\varepsilon$, \mathcal{L}_V is an infinitesimal generator of a C_0 semigroup of contractions on $C[0, 1]$. From this assertion and the second expression in (2.25), we easily infer that, for any $V \in S_\varepsilon$, $\mathbb{A}_0(U_* + V)$ is an infinitesimal generator of a C_0 semigroup of contractions on $X = C[0, 1] \times \mathbb{R}$. Hence, $\{\mathbb{A}_0(U_* + V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on X , with stability constants $(M, \omega) = (1, 0)$. Since $\mathbb{A}(V) = \mathbb{A}_0(U_* + V) + \mathbb{B}$ and \mathbb{B} is a bounded linear operator on X independent of V , by a standard perturbation result (see, e.g., Theorem 2.3 in section 5.2 of [20]), we see immediately that $\{\mathbb{A}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on $X = C[0, 1] \times \mathbb{R}$, with stability constants $(M, \omega) = (1, \|\mathbb{B}\|)$.

In order to prove that $\{\tilde{\mathbb{A}}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of a C_0 semigroup on $X_0 = C_V^1[0, 1] \times \mathbb{R}$, we first establish an estimate for the semigroup $e^{t\tilde{\mathcal{L}}_V}$ on $C_V^1[0, 1]$ different from (4.5), where $\tilde{\mathcal{L}}_V$ represents the part of \mathcal{L}_V in $C_V^1[0, 1]$. Let $q_0 \in C_V^1[0, 1]$ and $q = e^{t\tilde{\mathcal{L}}_V} q_0 = e^{t\mathcal{L}_V} q_0$. Then q is the solution of the problem

$$\frac{\partial q}{\partial t} + w_V(r) \frac{\partial q}{\partial r} = 0 \quad \text{for } 0 \leq r \leq 1 \text{ and } t > 0, \quad q|_{t=0} = q_0.$$

Let $l(r, t) = r(1-r) \frac{\partial q(r,t)}{\partial r}$ and $l_0(r) = r(1-r)q'_0(r)$. Formally differentiating the above equation in r and multiplying it with $r(1-r)$, we get

$$\frac{\partial l}{\partial t} + w_V(r) \frac{\partial l}{\partial r} = a_V(r)l \quad \text{for } 0 \leq r \leq 1 \text{ and } t > 0, \quad l|_{t=0} = l_0,$$

where $a_V(r) = (1-2r) \frac{w_V(r)}{r(1-r)} - w'_V(r)$. Since $w_V \in C^1[0, 1]$, $w_V(0) = w_V(1) = 0$, and $a_V \in C[0, 1]$, by a standard characteristics argument, we can easily show that this problem has a unique solution $l \in C([0, 1] \times [0, \infty))$ such that $[\frac{\partial}{\partial t} + w_V(r) \frac{\partial}{\partial r}]l \in C([0, 1] \times [0, \infty))$. Thus, the above formal computation makes sense. Clearly, there exists a nonnegative constant c_0 independent of V such that $a_V(r) \leq c_0$ for $0 < r < 1$, for all $V \in S_\varepsilon$. Using this fact and the characteristics argument, we can easily obtain $\|l(\cdot, t)\|_\infty \leq \|l_0\|_\infty e^{c_0 t}$ for $t \geq 0$. Combining this estimate with $\|q(\cdot, t)\|_\infty \leq \|q_0\|_\infty$ ensured by (4.4), we get $\|q(\cdot, t)\|_{C_V^1[0,1]} \leq \|q_0\|_{C_V^1[0,1]} e^{c_0 t}$ for $t \geq 0$, or

$$\|e^{t\tilde{\mathcal{L}}_V}\|_{L(C_V^1[0,1])} \leq e^{c_0 t} \quad \text{for } t \geq 0, \quad \text{for all } V \in S_\varepsilon.$$

Hence, $\{\tilde{\mathcal{L}}_V : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on $C_V^1[0, 1]$, with stability constants $(M, \omega) = (1, c_0)$. Using this assertion and

the second expression in (2.25), we see easily that $\{\tilde{\mathbb{A}}_0(U_* + V) : V \in S_\varepsilon\}$, the part of $\{\mathbb{A}_0(U_* + V) : V \in S_\varepsilon\}$ on $X_0 = C[0, 1] \times \mathbb{R}$, is a stable family of the infinitesimal generators of C_0 semigroups on X_0 , with stability constants $(M, \omega) = (1, c_0)$. Since $\hat{\mathbb{A}}(V) = \tilde{\mathbb{A}}_0(U_* + V) + \mathbb{B}$, and, by Corollary 3.2, \mathbb{B} is a bounded linear operator on X_0 independent of V , we conclude, as before, that $\{\mathbb{A}(V) : V \in S_\varepsilon\}$ is a stable family of the infinitesimal generators of C_0 semigroups on $X_0 = C_V^1[0, 1] \times \mathbb{R}$, with stability constants $(M, \omega) = (1, c_0 + \|\mathbb{B}\|_{L(C_V^1[0,1])})$. This completes the proof of Lemma 4.2. \square

Since $\mathbb{A} \in C^\infty(X, L(X_0, X))$, by Lemma 4.2, we see that, for any $V \in C([0, \infty), X)$ such that $V(t) \in S_\varepsilon$ for all $t \geq 0$, $\{\mathbb{A}(V(t)) : t \geq 0\}$ satisfies the conditions (H_1) – (H_3) in section 5.3 of [20]. It follows by Theorem 3.1 in section 5.3 of [20] that given such a function $V = V(t)$, there exists an evolution system determined by $\{\mathbb{A}(V(t)) : t \geq 0\}$, which we denote as $\mathbb{U}(t, s, V)$. By definition, this means that

- (1) for any $t \geq s \geq 0$, $\mathbb{U}(t, s, V)$ is a bounded linear operator on X ;
- (2) $\mathbb{U}(s, s, V) = \text{id}$ for all $s \geq 0$, $\mathbb{U}(t, s, V)\mathbb{U}(s, r, V) = \mathbb{U}(t, r, V)$ for all $t \geq s \geq r$;
and
- (3) the mapping $(t, s) \rightarrow \mathbb{U}(t, s, V)$ is strongly continuous for $t \geq s \geq 0$.

However, the theory developed in [20] does not ensure that $U = \mathbb{U}(t, s, V)U_0$ is a solution of the problem

$$(4.12) \quad \frac{dU}{dt} = \mathbb{A}(V(t))U \quad \text{for } t > s, \quad U|_{t=s} = U_0,$$

even if $U_0 \in X_0$, unless some other conditions are satisfied by $\mathbb{U}(t, s, V)$. These conditions are as follows (see conditions (E_4) and (E_5) in Theorem 4.3 in section 5.4 of [20]):

- (4) $\mathbb{U}(t, s, V)X_0 \subseteq X_0$ for any $t \geq s \geq 0$; and
- (5) For any $U_0 \in X_0$, the mapping $(t, s) \rightarrow \mathbb{U}(t, s, V)U_0$ is continuous in X_0 for $t \geq s \geq 0$.

In the following lemma we shall directly prove that, for any $U_0 \in X_0$, the problem (4.12) has a unique solution $U = U_s(t) \in C([s, \infty), X_0) \cap C^1([s, \infty), X)$. By Theorem 4.2 in section 5.4 of [20], it then follows that $U_s(t) = \mathbb{U}(t, s, V)U_0$, and consequently, the conditions (4) and (5) above are satisfied.

LEMMA 4.3. *Given $V \in C([0, \infty), X)$ such that $V(t) \in S_\varepsilon$ for all $t \geq 0$, for any $s \geq 0$ and any $U_0 \in X_0$, the problem (4.12) has a unique solution $U = U_s(t) \in C([s, \infty), X_0) \cap C^1([s, \infty), X)$.*

Proof. Let $U = (q, y)$ and $U_0 = (q_0, y_0)$. Then (4.12) can be rewritten as follows:

$$(4.13) \quad \begin{cases} \frac{\partial q}{\partial t} + w_V(r, t) \frac{\partial q}{\partial r} = a(r)q + \mathcal{B}(q) + b(r)y & \text{for } 0 \leq r \leq 1, \quad t > s, \\ \frac{dy}{dt} = \mathcal{F}(q) + \kappa y & \text{for } t > s, \\ q|_{t=s} = q_0(r) & \text{for } 0 \leq r \leq 1, \quad \text{and } y|_{t=s} = y_0. \end{cases}$$

Using the characteristic method and the Banach fixed point theorem, we can easily show that this problem has a unique local solution (q, y) , with $q \in C([0, 1] \times [0, \delta])$ and $y \in C^1[0, \delta]$ for some $\delta > 0$. Since $w_V(0, t) = w_V(1, t) = 0$ for all $t \geq 0$, we see that the two lines $r = 0$ and $r = 1$ are characteristic curves. It follows that all characteristic curves starting from the open interval $(0, 1)$ always lie in it, so that the solution of the above problem exists for all $t \geq s$. It remains to prove that $q \in C([0, \infty), C_V^1[0, 1])$. To this end we formally differentiate the first equation in (4.13) in r and multiply it

with $r(1 - r)$, which gives, by letting $l(r, t) = r(1 - r)\frac{\partial q(r, t)}{\partial r}$, that

$$(4.14) \quad \frac{\partial l}{\partial t} + w_V(r, t)\frac{\partial l}{\partial r} = a_1(r, t)u + f_1(r, t) \quad \text{for } 0 \leq r \leq 1, \quad t > 0,$$

where $a_1(r, t) = a(r) + (1 - 2r)\frac{w_V(r, t)}{r(1 - r)} - \frac{\partial w_V(r, t)}{\partial r}$, and

$$f_1(r, t) = r(1 - r)a'(r)q(r, t) + r(1 - r)\frac{\partial \mathcal{B}q(r, t)}{\partial r} + r(1 - r)b'(r)y(t).$$

Clearly, $a_1 \in C([0, 1] \times [0, \infty))$. By Corollary 3.2, we see that also $f_1 \in C([0, 1] \times [0, \infty))$. Thus, by using the characteristic method, we can easily prove that (4.14) imposed with the initial condition $l(r, 0) = r(1 - r)q'_0(r)$ has a unique solution $l \in C([0, 1] \times [0, \infty))$ such that $[\frac{\partial}{\partial t} + w_V(r, t)\frac{\partial}{\partial r}]l \in C([0, 1] \times [0, \infty))$. Thus, the above formal computation makes sense, and consequently, $q \in C([0, \infty), C^1_V[0, 1]) \cap C^1([0, \infty), C[0, 1])$. The desired assertion now becomes immediate. \square

By the above results and Theorems 4.2 and 5.2 in sections 5.4 and 5.5 of [20], we get the following.

COROLLARY 4.4. *Let $V \in C([0, \infty), X)$ be as in Lemma 4.3, and let $F \in C([0, \infty), X_0)$. Then, for any $U_0 \in X_0$, the initial value problem (4.1) has a unique solution $U \in C([0, \infty), X_0) \cap C^1([0, \infty), X)$, and it is given by (4.2). \square*

5. Similarity transformation. In the previous section we proved that given $V \in C([0, \infty), X)$ such that $V(t) \in S_\varepsilon$ for all $t \geq 0$, the family of operators $\{\mathbb{A}(V(t)) : t \geq 0\}$ determines an evolution system $\mathbb{U}(t, s, V)$ ($t \geq s \geq 0$), so that the solution of (4.1) can be expressed as (4.2). However, the deduction in the previous section does not provide us with an estimate of the form $\|\mathbb{U}(t, s, V)\| \leq Ce^{-\mu(t-s)}$, which, however, is essential in order to establish the estimate (2.17). To establish such an estimate, we shall follow a different approach as follows: For $V \equiv 0$, the desired estimate will be obtained by improving the linear estimate established in [5]. For a general $V = V(t)$, we shall use a special transformation, which we call *similarity transformation*, to transform the problem into an equivalent problem which can be estimated by using the estimate for the case $V \equiv 0$. The similarity transformation is a family of C^1 diffeomorphisms $\bar{r} = T(r, t, s)$ of the unit interval $0 \leq r \leq 1$ to itself, where $t \geq s \geq 0$ are parameters. The aim of this section is to establish this transformation.

Let $w \in C([0, \infty), C^1[0, 1])$. We assume that w satisfies the following condition: For some small parameter $\varepsilon > 0$,

$$(5.1) \quad -C\varepsilon r(1 - r)e^{-\mu t} \leq w(r, t) - u_*(r) \leq C\varepsilon r(1 - r)e^{-\mu t} \quad \text{for } 0 \leq r \leq 1, \quad t \geq 0,$$

where C is a positive constant independent of ε and w . Since $-C_1r(1 - r) \leq u_*(r) \leq -C_2r(1 - r)$, we see that

$$(5.2) \quad \sup_{0 < r < 1} \left| \frac{w(r, t)}{u_*(r)} - 1 \right| \leq C\varepsilon e^{-\mu t} \quad \text{for } 0 \leq r \leq 1, \quad t \geq 0,$$

and, for ε sufficiently small, $w(r, t)$ satisfies (4.3) for all $t \geq 0$.

Let $0 \leq \xi \leq 1$ and $s \geq 0$. Consider the following initial value problem:

$$(5.3) \quad \frac{dr}{dt} = u_*(r) \quad \text{for } t > s, \quad r|_{t=s} = \xi.$$

Since $u_* \in C^1[0, 1]$, $u_*(r) < 0$ for $0 < r < 1$, and, in particular, $u_*(0) = u_*(1) = 0$, it can be easily shown that this problem has a unique solution $r = \Phi_*(\xi, t, s)$ for all $t \geq s \geq 0$, satisfying the following properties (A):

$$\begin{aligned} &\Phi_*(\xi, t, s) \text{ is continuously differentiable in } (\xi, t, s), \\ &\Phi_*(0, t, s) = 0, \quad \Phi_*(1, t, s) = 1 \quad \text{for } t \geq s \geq 0, \\ &0 < \Phi_*(\xi, t, s) < 1 \quad \text{for } 0 < \xi < 1, \quad t \geq s \geq 0, \\ &\frac{\partial \Phi_*(\xi, t, s)}{\partial \xi} > 0, \quad \frac{\partial \Phi_*(\xi, t, s)}{\partial t} < 0 \quad \text{for } 0 < \xi < 1, \quad t \geq s \geq 0, \\ &\Phi_*(\xi, s, s) = \xi \quad \text{for } 0 \leq \xi \leq 1, \quad s \geq 0. \end{aligned}$$

From these properties we see that, for any $s \geq 0$ and $t \geq s$, the mapping $\xi \rightarrow r = \Phi_*(\xi, t, s)$ is a C^1 diffeomorphism of $[0, 1]$ to itself. Let $\xi = \Psi_*(r, t, s)$ be the inverse of this mapping. Clearly, Ψ_* also satisfies the properties (A) (with ξ replaced by r). Furthermore, from the relation $\Psi_*(\Phi_*(\xi, t, s), t, s) = \xi$ (for $0 \leq \xi \leq 1$ and $t \geq s \geq 0$), it can be easily shown that $\xi = \Psi_*(r, t, s)$ is the unique solution of the following initial value problem:

$$(5.4) \quad \frac{\partial \xi}{\partial t} + u_*(r) \frac{\partial \xi}{\partial r} = 0 \quad \text{for } t > s, \quad \xi|_{t=s} = r.$$

Next, let $r = \Phi(\xi, t, s)$ ($0 \leq \xi \leq 1, t \geq s \geq 0$) be the solution of the following problem:

$$(5.5) \quad \frac{dr}{dt} = w(r, t) \quad \text{for } t > s, \quad r|_{t=s} = \xi.$$

Similarly as before, $\Phi(\xi, t, s)$ is well-defined for all $0 \leq \xi \leq 1$ and $t \geq s \geq 0$, and it also satisfies the properties (A). It follows that, for any $s \geq 0$ and $t \geq s$, the mapping $\xi \rightarrow r = \Phi(\xi, t, s)$ is a C^1 diffeomorphism of $[0, 1]$ to itself. Let $\xi = \Psi(r, t, s)$ be the inverse of this mapping. Clearly, Ψ also satisfies the properties (A) (with ξ replaced by r). Furthermore, from the relation $\Psi(\Phi(\xi, t, s), t, s) = \xi$ (for $0 \leq \xi \leq 1$ and $t \geq s \geq 0$), it can be easily shown that $\xi = \Psi(r, t, s)$ is the unique solution of the following initial value problem:

$$(5.6) \quad \frac{\partial \xi}{\partial t} + w(r, t) \frac{\partial \xi}{\partial r} = 0 \quad \text{for } t > s, \quad \xi|_{t=s} = r.$$

In what follows we consider the following initial value problem:

$$(5.7) \quad \begin{cases} \frac{\partial \bar{r}}{\partial t} + w(r, t) \frac{\partial \bar{r}}{\partial r} = u_*(\bar{r}) & \text{for } 0 \leq r \leq 1, \quad t > s, \\ \bar{r}|_{t=s} = r & \text{for } 0 \leq r \leq 1. \end{cases}$$

LEMMA 5.1. *For any $0 \leq r \leq 1$ and $s \geq 0$, the problem (5.7) has a unique solution $\bar{r} = T(r, t, s)$ for all $t \geq s$, and the following relation holds:*

$$(5.8) \quad T(r, t, s) = \Phi_*(\Psi(r, t, s), t, s) \quad \text{for } 0 \leq r \leq 1, \quad t \geq s \geq 0.$$

Proof. Using (5.3) and (5.6) we can easily verify that $\bar{r} = \Phi_*(\Psi(r, t, s), t, s)$ is a solution of the problem (5.7). Thus, (5.8) follows by the uniqueness of the solution. \square

By (5.8), it is evident that, for any $s \geq 0$ and $t \geq s$, the mapping $r \rightarrow \bar{r} = T(r, t, s)$ is a C^1 diffeomorphism of $[0, 1]$ to itself, satisfying $T(0, t, s) = 0$, $T(1, t, s) = 1$ for $t \geq s \geq 0$, and $\partial T(r, t, s)/\partial r > 0$ for $0 < r < 1$, $t \geq s$. We denote by $r = S(\bar{r}, t, s)$ the inverse of this mapping. By (5.8), it is clear that

$$(5.9) \quad S(\bar{r}, t, s) = \Phi(\Psi_*(\bar{r}, t, s), t, s) \quad \text{for } 0 \leq \bar{r} \leq 1, \quad t \geq s \geq 0.$$

It is also clear that $S(\bar{r}, t, s)$ satisfies $S(0, t, s) = 0$, $S(1, t, s) = 1$ for $t \geq s \geq 0$ and that $\partial S(\bar{r}, t, s)/\partial \bar{r} > 0$ for $0 < \bar{r} < 1$, $t \geq s \geq 0$.

T and S can be expressed in more explicit formulations. To show this we introduce a function F_* as follows:

$$(5.10) \quad F_*(r) = - \int_{\frac{1}{2}}^r \frac{d\eta}{u_*(\eta)} = \int_{\frac{1}{2}}^r \frac{d\eta}{|u_*(\eta)|} \quad \text{for } 0 < r < 1.$$

Clearly, $F_* \in C^1(0, 1)$, $F'_*(r) > 0$ for all $0 < r < 1$, and $\lim_{r \rightarrow 0^+} F_*(r) = -\infty$, $\lim_{r \rightarrow 1^-} F_*(r) = \infty$. Hence, $\bar{r} = F_*(r)$ is a C^1 diffeomorphism of the open unit interval $(0, 1)$ onto the real line $(-\infty, \infty)$. From (5.3) we easily obtain $F_*(\Phi_*(\xi, t, s)) - F_*(\xi) = -t + s$. Thus, we have

$$(5.11) \quad \Phi_*(\xi, t, s) = F_*^{-1}(F_*(\xi) - t + s),$$

and consequently,

$$(5.12) \quad \Psi_*(r, t, s) = F_*^{-1}(F_*(r) + t - s).$$

Next, let

$$g(\xi, t, s) = G(\Phi(\xi, t, s), t), \quad \text{where } G(r, t) = \frac{w(r, t)}{u_*(r)} - 1.$$

Since $w(r, t) = [1 + G(r, t)]u_*(r)$, from (5.5) we see that $r = \Phi(\xi, t, s)$ is a solution of the following problem:

$$(5.13) \quad \frac{dr}{dt} = [1 + g(\xi, t, s)]u_*(r) \quad \text{for } t > s, \quad r|_{t=s} = \xi.$$

Thus, similarly, as before, we have $F_*(\Phi(\xi, t, s)) - F_*(\xi) = -t + s - \int_s^t g(\xi, \tau, s)d\tau$, so that

$$(5.14) \quad \Phi(\xi, t, s) = F_*^{-1} \left(F_*(\xi) - t + s - \int_s^t g(\xi, \tau, s)d\tau \right),$$

$$(5.15) \quad \Psi(r, t, s) = F_*^{-1} \left(F_*(r) + t - s + \int_s^t g(\Psi(r, t, s), \tau, s)d\tau \right).$$

Combining (5.8), (5.9), (5.11), (5.12), (5.14), and (5.15), we see that

$$(5.16) \quad T(r, t, s) = F_*^{-1} \left(F_*(r) + \int_s^t g(\Psi(r, t, s), \tau, s)d\tau \right),$$

$$(5.17) \quad S(\bar{r}, t, s) = F_*^{-1} \left(F_*(\bar{r}) - \int_s^t g(\Psi_*(\bar{r}, t, s), \tau, s)d\tau \right).$$

LEMMA 5.2. Assume that $|\zeta| \leq C$. Then there exist positive constants C_1 and C_2 depending only on C such that, for any $0 < r < 1$, we have

$$(5.18) \quad C_1 r(1-r) \leq F_*^{-1}(F_*(r) + \zeta) [1 - F_*^{-1}(F_*(r) + \zeta)] \leq C_2 r(1-r).$$

Proof. Since $-C \leq \zeta \leq C$, by the monotonicity of F_* , we have

$$F_*^{-1}(F_*(r) - C) \leq F_*^{-1}(F_*(r) + \zeta) \leq F_*^{-1}(F_*(r) + C).$$

Thus,

$$\frac{F_*^{-1}(F_*(r) + \zeta)}{r} \leq \frac{F_*^{-1}(F_*(r) + C)}{r} \quad \text{and} \\ \frac{1 - F_*^{-1}(F_*(r) + \zeta)}{1-r} \leq \frac{1 - F_*^{-1}(F_*(r) - C)}{1-r}.$$

Moreover, using the facts that $u_*(r) = u'_*(0)r[1 + O(r^\beta)]$ as $r \rightarrow 0^+$ for some $0 < \beta \leq 1$ (see Assertion (3) of Lemma 3.1) and that $u_*(r) = -u'_*(1)(1-r)[1 + O(1-r)]$ as $r \rightarrow 1^-$, we can easily show that

$$\lim_{r \rightarrow 0^+} \frac{F_*^{-1}(F_*(r) + C)}{r} = e^{C|u'_*(0)|} \quad \text{and} \quad \lim_{r \rightarrow 1^-} \frac{1 - F_*^{-1}(F_*(r) - C)}{1-r} = e^{Cu'_*(1)}.$$

From these relations we immediately obtain the second inequality in (5.18). The proof for the first inequality in (5.18) is similar. \square

COROLLARY 5.3. For ε sufficiently small, we have

$$(5.19) \quad C_1 r(1-r) \leq T(r, t, s)[1 - T(r, t, s)] \leq C_2 r(1-r),$$

$$(5.20) \quad C_1 \bar{r}(1-\bar{r}) \leq S(\bar{r}, t, s)[1 - S(\bar{r}, t, s)] \leq C_2 \bar{r}(1-\bar{r}).$$

Proof. Let $\zeta = \int_s^t g(\Psi(r, t, s), \tau, s) d\tau$. By (5.2), we have

$$|\zeta| \leq \int_s^t |g(\Psi(r, t, s), \tau, s)| d\tau \leq C\varepsilon \int_s^t e^{-\mu\tau} d\tau \leq C\varepsilon \int_0^\infty e^{-\mu\tau} d\tau \leq C\varepsilon \leq C.$$

Hence, (5.19) follows from (5.16) and (5.18). Similarly, (5.20) follows from (5.17) and (5.18). \square

As an immediate consequence of Corollary 5.3, we see that there exists a constant $C > 1$ such that, for ε sufficiently small, we have $C^{-1} \leq T(r, t, s)/r \leq C$ and $C^{-1} \leq S(\bar{r}, t, s)/\bar{r} \leq C$.

COROLLARY 5.4. For ε sufficiently small, we have the following inequalities:

$$(5.21)$$

$$C_1 \Psi_*(r, t, s)[1 - \Psi_*(r, t, s)] \leq \Psi(r, t, s)[1 - \Psi(r, t, s)] \leq C_2 \Psi_*(r, t, s)[1 - \Psi_*(r, t, s)],$$

$$(5.22)$$

$$C_1 \Phi_*(\bar{r}, t, s)[1 - \Phi_*(\bar{r}, t, s)] \leq \Phi(\bar{r}, t, s)[1 - \Phi(\bar{r}, t, s)] \leq C_2 \Phi_*(\bar{r}, t, s)[1 - \Phi_*(\bar{r}, t, s)].$$

Proof. Let $\bar{r} = \Psi_*(r, t, s)$ and $\zeta = \int_s^t g(\Psi(r, t, s), \tau, s) d\tau$. Then, by (5.12) and (5.15), we have $\Psi(r, t, s) = F_*^{-1}(F_*(\bar{r}) + \zeta)$. By this expression and (5.18), we immediately obtain (5.21). The proof of (5.22) is similar. \square

As an immediate consequence of Corollary 5.4, we see that there exists a constant $C > 1$ such that, for ε sufficiently small, we have $C^{-1} \leq \Psi(r, t, s)/\Psi_*(r, t, s) \leq C$ and $C^{-1} \leq \Phi(\bar{r}, t, s)/\Phi_*(\bar{r}, t, s) \leq C$.

LEMMA 5.5. *We have the following inequalities:*

$$(5.23) \quad |T(r, t, s) - r| \leq C\varepsilon(e^{-\mu s} - e^{-\mu t})r(1 - r),$$

$$(5.24) \quad |S(\bar{r}, t, s) - \bar{r}| \leq C\varepsilon(e^{-\mu s} - e^{-\mu t})\bar{r}(1 - \bar{r}).$$

Proof. Similarly as in the proof of Corollary 5.3, we have

$$\int_s^t |g(\Psi(r, t, s), \tau, s)| d\tau \leq C\varepsilon \int_s^t e^{-\mu\tau} \leq C\varepsilon(e^{-\mu s} - e^{-\mu t}).$$

Thus, by noticing that $\frac{dF_*^{-1}(\eta)}{d\eta} = \frac{1}{F_*'(F_*^{-1}(\eta))} = |u_*(F_*^{-1}(\eta))|$, we see that

$$\begin{aligned} |T(r, t, s) - r| &= \left| F_*^{-1} \left(F_*(r) + \int_s^t g(\Psi(r, t, s), \tau, s) d\tau \right) - F_*^{-1}(F_*(r)) \right| \\ &\leq \int_0^1 |u_*(F_*^{-1}(F_*(r) + \zeta_\theta))| d\theta \cdot \int_s^t |g(\Psi(r, t, s), \tau, s)| d\tau \\ &\leq C\varepsilon(e^{-\mu s} - e^{-\mu t}) \cdot \int_0^1 |u_*(F_*^{-1}(F_*(r) + \zeta_\theta))| d\theta, \end{aligned}$$

where $\zeta_\theta = \theta \int_s^t g(\Psi(r, t, s), \tau, s) d\tau$. Since $|\zeta_\theta| \leq C$ and $|u_*(\eta)| \leq C\eta(1 - \eta)$, by Lemma 5.2, we have

$$|u_*(F_*^{-1}(F_*(r) + \zeta_\theta))| \leq CF_*^{-1}(F_*(r) + \zeta_\theta) [1 - F_*^{-1}(F_*(r) + \zeta_\theta)] \leq Cr(1 - r).$$

Substituting this estimate into the above inequality, we see that (5.23) follows. The proof of (5.24) is similar. \square

COROLLARY 5.6. *We have the following inequalities:*

$$(5.25) \quad |\Phi(\xi, t, s) - \Phi_*(\xi, t, s)| \leq C\varepsilon(e^{-\mu s} - e^{-\mu t})\Phi_*(\xi, t, s)[1 - \Phi_*(\xi, t, s)],$$

$$(5.26) \quad |\Psi(r, t, s) - \Psi_*(r, t, s)| \leq C\varepsilon(e^{-\mu s} - e^{-\mu t})\Psi_*(r, t, s)[1 - \Psi_*(r, t, s)].$$

Proof. Let $r = \Phi(\xi, t, s)$. Then $\xi = \Psi(r, t, s)$ and $\Phi_*(\xi, t, s) = \Phi_*(\Psi(r, t, s), t, s) = T(r, t, s)$. Thus, by (5.23), we have

$$|\Phi(\xi, t, s) - \Phi_*(\xi, t, s)| = |r - T(r, t, s)| \leq C\varepsilon(e^{-\mu s} - e^{-\mu t})r(1 - r).$$

Substituting $r = \Phi(\xi, t, s)$ into the right-hand side of the last inequality and using (5.22), we see that (5.25) follows. The proof of (5.26) is similar. \square

LEMMA 5.7. *Assume that, in addition to (5.2), there also holds*

$$(5.27) \quad \max_{0 \leq r \leq 1} \left| \frac{\partial w(r, t)}{\partial r} - u'_*(r) \right| \leq C\varepsilon e^{-\mu t}.$$

Then we have the following estimates:

$$(5.28) \quad e^{-C\varepsilon(e^{-\mu s} - e^{-\mu t})} \leq \frac{\partial T(r, t, s)}{\partial r} \leq e^{C\varepsilon(e^{-\mu s} - e^{-\mu t})},$$

$$(5.29) \quad e^{-C\varepsilon(e^{-\mu s} - e^{-\mu t})} \leq \frac{\partial S(\bar{r}, t, s)}{\partial \bar{r}} \leq e^{C\varepsilon(e^{-\mu s} - e^{-\mu t})}.$$

Proof. Recalling that $T(r, t, s) = \Phi_*(\Psi(r, t, s), t, s)$, we have

$$\begin{aligned}
 \frac{\partial T(r, t, s)}{\partial r} &= \frac{\partial \Phi_*}{\partial \xi}(\Psi(r, t, s), t, s) \frac{\partial \Psi(r, t, s)}{\partial r} \\
 &= \frac{\partial \Phi_*}{\partial \xi}(\Psi(r, t, s), t, s) \left[\frac{\partial \Phi}{\partial \xi}(\Psi(r, t, s), t, s) \right]^{-1} \\
 (5.30) \quad &= \exp \left(\int_s^t \left[u'_*(\Phi_*(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) \right] d\tau \right) \Bigg|_{\xi=\Psi(r, t, s)}.
 \end{aligned}$$

We write

$$\begin{aligned}
 u'_*(\Phi_*(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) &= [u'_*(\Phi_*(\xi, \tau, s)) - u'_*(\Phi(\xi, \tau, s))] \\
 &\quad + \left[u'_*(\Phi(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) \right].
 \end{aligned}$$

From assumption (5.27) we see that

$$\sup_{\xi \in \mathbb{R}} \left| u'_*(\Phi(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) \right| \leq \sup_{0 < r < 1} \left| u'_*(r) - \frac{\partial w}{\partial r}(r, \tau) \right| \leq C\varepsilon e^{-\mu\tau}.$$

Next, by assertion (3) of Lemma 3.1, we know that there exists $0 \leq \gamma < 1$ such that $r^\gamma u''_*(r) \in C[0, 1]$. With this fact in mind, we use the mean value theorem to compute

$$\begin{aligned}
 &|u'_*(\Phi_*(\xi, \tau, s)) - u'_*(\Phi(\xi, \tau, s))| \Big|_{\xi=\Psi(r, t, s)} \\
 &= |u''_*(\zeta)| |\Phi_*(\xi, \tau, s) - \Phi(\xi, \tau, s)| \Big|_{\xi=\Psi(r, t, s)} \\
 &= |\zeta^\gamma u''_*(\zeta)| \cdot \left[\left(\frac{\Phi(\xi, \tau, s)}{\zeta} \right)^\gamma \cdot (\Phi(\xi, \tau, s))^{-\gamma} |\Phi_*(\xi, \tau, s) - \Phi(\xi, \tau, s)| \right] \Bigg|_{\xi=\Psi(r, t, s)},
 \end{aligned}$$

where $\zeta = \theta\Phi_*(\xi, \tau, s) + (1-\theta)\Phi(\xi, \tau, s)$ for some $0 < \theta < 1$ (depending on ξ, τ, s). Since there exists a constant $0 < c < 1$ such that $\frac{\Phi_*(\xi, \tau, s)}{\Phi(\xi, \tau, s)} \geq c$ for ε sufficiently small, we have $\zeta \geq c\Phi(\xi, \tau, s)$. Thus,

$$\begin{aligned}
 &\left| u'_*(\Phi_*(\xi, \tau, s)) - u'_*(\Phi(\xi, \tau, s)) \right| \Big|_{\xi=\Psi(r, t, s)} \\
 &\leq C(\Phi(\xi, \tau, s))^{-\gamma} |\Phi_*(\xi, \tau, s) - \Phi(\xi, \tau, s)| \Big|_{\xi=\Psi(r, t, s)} \\
 &\leq C\varepsilon(e^{-\mu s} - e^{-\mu\tau})(\Phi(\xi, \tau, s))^{-\gamma} \Phi_*(\xi, \tau, s)[1 - \Phi_*(\xi, \tau, s)] \Big|_{\xi=\Psi(r, t, s)} \\
 &\leq C\varepsilon(e^{-\mu s} - e^{-\mu\tau})(\Phi(\xi, \tau, s))^{1-\gamma} [1 - \Phi(\xi, \tau, s)] \Big|_{\xi=\Psi(r, t, s)} \\
 &= C\varepsilon(e^{-\mu s} - e^{-\mu\tau})(\Psi(r, t, \tau))^{1-\gamma} [1 - \Psi(r, t, \tau)].
 \end{aligned}$$

In getting the last equality, we used the following relation:

$$(5.31) \quad \Phi(\Psi(r, t, s), \tau, s) = \Psi(r, t, \tau) \quad \text{for } 0 \leq r \leq 1, \quad s \leq \tau \leq t.$$

The proof of this relation is as follows: From (5.6) we know that $\rho = \Psi(r, t, \tau)$ is a solution of the following problem:

$$\frac{\partial \rho}{\partial t} + w(r, t) \frac{\partial \rho}{\partial r} = 0 \quad \text{for } 0 \leq r \leq 1, \quad t > \tau, \quad \rho|_{t=\tau} = r.$$

But it is easy to verify that $\rho = \Phi(\Psi(r, t, s), \tau, s)$ is also a solution of this problem. Hence, by uniqueness, we have (5.31). Hence, using (5.21) and (5.12), we get

$$\begin{aligned} & \left| u'_*(\Phi_*(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) \right|_{\xi=\Psi(r,t,s)} \\ & \leq C\varepsilon (e^{-\mu s} - e^{-\mu \tau}) (\Psi_*(r, t, \tau))^{1-\gamma} [1 - \Psi_*(r, t, \tau)] + C\varepsilon e^{-\mu \tau} \\ & = C\varepsilon (e^{-\mu s} - e^{-\mu \tau}) [F_*^{-1}(F_*(r) + t - \tau)]^{1-\gamma} [1 - F_*^{-1}(F_*(r) + t - \tau)] + C\varepsilon e^{-\mu \tau}. \end{aligned}$$

It follows that

$$\begin{aligned} & \int_s^t \left| u'_*(\Phi_*(\xi, \tau, s)) - \frac{\partial w}{\partial r}(\Phi(\xi, \tau, s), \tau) \right|_{\xi=\Psi(r,t,s)} d\tau \\ & \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}) \int_s^t [F_*^{-1}(F_*(r) + t - \tau)]^{1-\gamma} [1 - F_*^{-1}(F_*(r) + t - \tau)] d\tau \\ & \quad + C\varepsilon \int_s^t e^{-\mu \tau} d\tau \\ & = C\varepsilon (e^{-\mu s} - e^{-\mu t}) \int_{F_*(r)}^{F_*(r)+t-s} (F_*^{-1}(\xi))^{1-\gamma} [1 - F_*^{-1}(\xi)] d\xi + C\varepsilon (e^{-\mu s} - e^{-\mu t}) \\ & \quad \left(\xi = F_*(\eta), \quad d\xi = F'_*(\eta) d\eta = \frac{d\eta}{|u_*(\eta)|} \right) \\ & = C\varepsilon (e^{-\mu s} - e^{-\mu t}) \int_r^{F_*^{-1}(F_*(r)+t-s)} \frac{\eta^{1-\gamma}(1-\eta)}{|u_*(\eta)|} d\eta + C\varepsilon (e^{-\mu s} - e^{-\mu t}) \\ & \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}) \int_0^1 \eta^{-\gamma} d\eta + C\varepsilon (e^{-\mu s} - e^{-\mu t}) = C\varepsilon (e^{-\mu s} - e^{-\mu t}). \end{aligned}$$

Combining this result with (5.30), we see that (5.28) follows. Finally, (5.29) is an immediate consequence of (5.28). \square

COROLLARY 5.8. *Under the assumption of Lemma 5.7, for ε sufficiently small, we have*

$$\left| \frac{\partial T(r, t, s)}{\partial r} - 1 \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}), \quad \left| \frac{\partial S(\bar{r}, t, s)}{\partial \bar{r}} - 1 \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}),$$

and

$$C^{-1} \leq \frac{\partial T(r, t, s)}{\partial r} \leq C, \quad C^{-1} \leq \frac{\partial S(\bar{r}, t, s)}{\partial \bar{r}} \leq C. \quad \square$$

LEMMA 5.9. *Assume that $a \in C^1_V[0, 1]$, and the assumption of Lemma 5.7 holds. Then*

$$(5.32) \quad \|a(S(\cdot, t, s)) - a\|_\infty \leq C \|a\|_1 \varepsilon (e^{-\mu s} - e^{-\mu t}),$$

where $\|a\|_1 = \max_{0 \leq r \leq 1} r(1-r)|a'(r)|$. If further $r^2(1-r)^2 a''(r) \in C[0, 1]$, then we also have

$$(5.33) \quad \|a(S(\cdot, t, s)) - a\|_{C^1_V[0,1]} \leq C \|a\|_2 \varepsilon (e^{-\mu s} - e^{-\mu t}),$$

where $\|a\|_2 = \|a\|_1 + \max_{0 \leq r \leq 1} r^2(1-r)^2 |a''(r)|$.

Proof. We have

$$\begin{aligned} |a(S(r, t, s)) - a(r)| &= |a'(\eta)||S(r, t, s) - r| \\ &\leq C\varepsilon\eta(1 - \eta)|a'(\eta)| \cdot \frac{r(1 - r)}{\eta(1 - \eta)} (e^{-\mu s} - e^{-\mu t}) \\ &\leq C\|a\|_1\varepsilon (e^{-\mu s} - e^{-\mu t}), \end{aligned}$$

where $\eta = (1 - \theta)r + \theta S(r, t, s)$ for some $0 < \theta < 1$ (depending on $r, t,$ and s). In getting the last inequality, we used the inequality $\eta(1 - \eta) \geq Cr(1 - r)$ for $0 \leq r \leq 1$, which follows from (5.20) and the following identity:

$$\eta(1 - \eta) = (1 - \theta)r(1 - r) + \theta S(r, t, s)[1 - S(r, t, s)] + \theta(1 - \theta)[r - S(r, t, s)]^2.$$

Hence, (5.32) is proved. Next, we compute

$$\begin{aligned} r(1 - r) \left| \frac{\partial a(S(r, t, s))}{\partial r} - a'(r) \right| &= r(1 - r) \left| a'(S(r, t, s)) \frac{\partial S(r, t, s)}{\partial r} - a'(r) \right| \\ &\leq r(1 - r) |a'(S(r, t, s))| \left| \frac{\partial S(r, t, s)}{\partial r} - 1 \right| \\ &\quad + r(1 - r) |a'(S(r, t, s)) - a'(r)|. \end{aligned}$$

By Corollaries 5.3 and 5.8, we see that the first term on the right-hand side is bounded by $C\|a\|_1 \cdot C\varepsilon(e^{-\mu s} - e^{-\mu t})$, and, by a similar argument as before, we can easily show that the second term on the right-hand side is bounded by $C(\max_{0 \leq r \leq 1} r^2(1 - r)^2 |a''(r)|)\varepsilon(e^{-\mu s} - e^{-\mu t})$. Hence, (5.33) follows. \square

LEMMA 5.10. *Given $a \in C[0, 1]$, we define a bounded linear operator L in $C[0, 1]$ by*

$$L(q)(r) = \frac{1}{r^3} \int_0^r a(\rho)q(\rho)\rho^2 d\rho \quad \text{for } q \in C[0, 1], \quad 0 < r \leq 1,$$

and $L(q)(0) = \lim_{r \rightarrow 0^+} L(q)(r) = \frac{1}{3}a(0)q(0)$. Let $\bar{r} = T(r, t, s)$ and $r = S(\bar{r}, t, s)$ be as before, and let \tilde{L} be the following bounded linear operator in $C[0, 1]$:

$$\tilde{L}(q)(\bar{r}) = \frac{1}{r^3} \int_0^r a(\rho)q(T(\rho, t, s))\rho^2 d\rho \Big|_{r=S(\bar{r}, t, s)} \quad \text{for } q \in C[0, 1], \quad 0 < \bar{r} \leq 1,$$

and $\tilde{L}(q)(0) = \lim_{\bar{r} \rightarrow 0^+} \tilde{L}(q)(\bar{r}) = \frac{1}{3}a(0)q(0)$. Assume that $a \in C_V^1[0, 1]$, and the assumption of Lemma 5.7 holds. Then both L and \tilde{L} are bounded linear operators from $C[0, 1]$ to $C_V^1[0, 1]$, and we have

$$(5.34) \quad \left\| \tilde{L} - L \right\|_{L(C[0,1], C_V^1[0,1])} \leq C\|a\|_{C_V^1[0,1]}\varepsilon (e^{-\mu s} - e^{-\mu t}).$$

Proof. We give only the proof of (5.34), because the proof of the assertion that both L and \tilde{L} are bounded linear operators from $C[0, 1]$ to $C_V^1[0, 1]$ follows by a similar argument.

We first note that, for $q \in C[0, 1]$ and $0 < \bar{r} \leq 1$, $\tilde{L}(q)(\bar{r})$ can be rewritten as follows:

$$\tilde{L}(q)(\bar{r}) = \frac{1}{[S(\bar{r}, t, s)]^3} \int_0^{\bar{r}} a(S(\rho, t, s))q(\rho) \left[\frac{S(\rho, t, s)}{\rho} \right]^2 \frac{\partial S(\rho, t, s)}{\partial \rho} \rho^2 d\rho.$$

Thus,

$$\begin{aligned} \tilde{L}(q)(\bar{r}) - L(q)(\bar{r}) &= \left[\frac{\bar{r}}{S(\bar{r}, t, s)} \right]^3 \cdot \frac{1}{\bar{r}^3} \int_0^{\bar{r}} a(S(\rho, t, s))q(\rho) \\ &\quad \left[\frac{S(\rho, t, s)}{\rho} \right]^2 \left[\frac{\partial S(\rho, t, s)}{\partial \rho} - 1 \right] \rho^2 d\rho \\ &\quad + \left[\frac{\bar{r}}{S(\bar{r}, t, s)} \right]^3 \cdot \frac{1}{\bar{r}^3} \int_0^{\bar{r}} a(S(\rho, t, s))q(\rho) \left\{ \left[\frac{S(\rho, t, s)}{\rho} \right]^2 - 1 \right\} \rho^2 d\rho \\ &\quad + \left[\frac{\bar{r}}{S(\bar{r}, t, s)} \right]^3 \cdot \frac{1}{\bar{r}^3} \int_0^{\bar{r}} [a(S(\rho, t, s)) - a(\rho)]q(\rho)\rho^2 d\rho \\ &\quad + \left\{ \left[\frac{\bar{r}}{S(\bar{r}, t, s)} \right]^3 - 1 \right\} \cdot \frac{1}{\bar{r}^3} \int_0^{\bar{r}} a(\rho)q(\rho)\rho^2 d\rho. \end{aligned}$$

From Corollary 5.3, Lemma 5.5, Corollary 5.8, and Lemma 5.9, we know that

$$\begin{aligned} \left| \frac{\bar{r}}{S(\bar{r}, t, s)} \right| &\leq C, \quad \left| \frac{\bar{r}}{S(\bar{r}, t, s)} - 1 \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}), \\ \left| \frac{S(\rho, t, s)}{\rho} \right| &\leq C, \quad \left| \frac{\partial S(\rho, t, s)}{\partial \rho} - 1 \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}), \\ |a(S(\rho, t, s)) - a(\rho)| &\leq C\|a\|_1 \varepsilon (e^{-\mu s} - e^{-\mu t}). \end{aligned}$$

Using the above estimates, we see easily that

$$(5.35) \quad \max_{0 \leq \bar{r} \leq 1} \left| \tilde{L}(q)(\bar{r}) - L(q)(\bar{r}) \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}) \|a\|_{C^1_V[0,1]} \|q\|_\infty.$$

Next, by a simple computation, we have

$$\begin{aligned} \bar{r}(1 - \bar{r})L(q)'(\bar{r}) &= (1 - \bar{r})a(\bar{r})q(\bar{r}) - \frac{3(1 - \bar{r})}{\bar{r}^3} \int_0^{\bar{r}} a(\rho)q(\rho)\rho^2 d\rho, \\ \bar{r}(1 - \bar{r})\tilde{L}(q)'(\bar{r}) &= \frac{\bar{r}(1 - \bar{r})}{S(\bar{r}, t, s)} a(S(\bar{r}, t, s))q(\bar{r}) \frac{\partial S(\bar{r}, t, s)}{\partial \bar{r}} \\ &\quad - \frac{3\bar{r}(1 - \bar{r})}{[S(\bar{r}, t, s)]^4} \frac{\partial S(\bar{r}, t, s)}{\partial \bar{r}} \int_0^{\bar{r}} a(S(\rho, t, s))q(\rho) \\ &\quad \left[\frac{S(\rho, t, s)}{\rho} \right]^2 \frac{\partial S(\rho, t, s)}{\partial \rho} \rho^2 d\rho. \end{aligned}$$

Using these expressions and a similar argument as before, we have

$$(5.36) \quad \sup_{0 < \bar{r} < 1} \bar{r}(1 - \bar{r}) \left| \tilde{L}(q)'(\bar{r}) - L(q)'(\bar{r}) \right| \leq C\varepsilon (e^{-\mu s} - e^{-\mu t}) \|a\|_{C^1_V[0,1]} \|q\|_\infty.$$

To save space, we omit the details here. By (5.35) and (5.36), we see that (5.34) follows. \square

6. Decay estimates. In this section we establish a decay estimate for the evolution system $\{U(t, s, V) : t \geq s \geq 0\}$ obtained in section 4, where $V = V(t) \in C([0, \infty), S_\varepsilon)$, under an additional assumption that $V(t)$ is exponentially decaying as $t \rightarrow \infty$.

We first consider the special case where $V = 0$. In this case we have $\mathbb{U}(t, s, V) = e^{(t-s)\mathbb{A}(0)}$. The main result, Theorem 5.1 of [5], gives a decay estimate for $e^{t\mathbb{A}(0)}$ (see (6.8) below). But that estimate contains some singularity at $r = 0$, so that it does not meet our requirement. In what follows we use the fact that $\lim_{r \rightarrow 0^+} r p'_*(r) = 0$ ensured by Lemma 3.1 to establish an improved estimate. To this end we need a preliminary lemma, which gives an estimate for the semigroup generated by the following operator $\mathcal{L} = \mathcal{L}_0 + a$:

$$\mathcal{L}q(r) = -w(r)q'(r) + a(r)q(r) \quad \text{for } 0 \leq r \leq 1,$$

where w and a are given functions.

LEMMA 6.1. *Assume that $w \in C^1[0, 1]$ and satisfies (4.3), and $a \in C^1_V[0, 1]$. Then \mathcal{L} generates a C_0 semigroup $e^{t\mathcal{L}}$ on $C[0, 1]$ satisfying the following estimate:*

$$(6.1) \quad \|e^{t\mathcal{L}}\|_{L(C[0,1])} \leq e^{\omega_0 t} \quad \text{for } t \geq 0,$$

where $\omega_0 = \max_{0 \leq r \leq 1} a(r)$. Moreover, $C^1_V[0, 1]$ is \mathcal{L} -admissible, and, for any $\omega > \omega_0$, we have

$$(6.2) \quad \|e^{t\mathcal{L}}\|_{L(C^1_V[0,1])} \leq C_\omega e^{\omega t} \quad \text{for } t \geq 0.$$

Here, C_ω is independent of the function w .

The proof is similar to that of Lemma 4.1, so that is omitted. \square

LEMMA 6.2. *There exists a constant $\mu^* > 0$ such that, for any $0 < \mu < \mu^*$, the semigroup $e^{t\mathbb{A}(0)}$ ($t \geq 0$) generated by $\mathbb{A}(0)$ satisfies the following estimate:*

$$(6.3) \quad \|e^{t\mathbb{A}(0)}\|_{L(C[0,1])} \leq C e^{-\mu t} \quad \text{for } t \geq 0.$$

Proof. Given $U_0 = (\phi_0, \zeta_0) \in X$, let $U(t) = e^{t\mathbb{A}(0)}U_0 = (\phi(r, t), \zeta(t))$. Then (ϕ, ζ) is the unique solution of the following initial value problem:

$$(6.4) \quad \partial_t \phi + u_*(r) \partial_r \phi = a(r)\phi + \mathcal{B}(\phi) + b(r)\zeta \quad \text{for } 0 \leq r \leq 1, \quad t > 0,$$

$$(6.5) \quad \frac{d\zeta}{dt} = \mathcal{F}(\phi) + \kappa \zeta \quad \text{for } t > 0,$$

$$(6.6) \quad \phi(r, 0) = \phi_0(r) \quad \text{for } 0 \leq r \leq 1, \quad \zeta(0) = \zeta_0,$$

where $a(r), b(r), \mathcal{B}(\phi), \mathcal{F}(\phi)$, and κ are given in (2.19)–(2.23). By Theorem 5.1 of [5] and the remark in the end of section 8 of [5], we know that there exists a constant $\sigma^* > 0$ and a function $\hat{\phi} \in C^1(0, 1]$ satisfying

$$(6.7) \quad \hat{\phi}(r) > 0 \quad \text{for } 0 < r \leq 1, \quad \hat{\phi}(r) \sim Cr^{-\theta} \quad \text{for } r \rightarrow 0$$

for some constants $1 \leq \theta < 3$ and $C > 0$, such that the solution of the above problem satisfies the following estimate:

$$(6.8) \quad |\zeta(t)| + \sup_{0 < r \leq 1} \left| \frac{\phi(r, t)}{\hat{\phi}(r)} \right| \leq C \left(|\zeta_0| + \sup_{0 < r \leq 1} \left| \frac{\phi_0(r)}{\hat{\phi}(r)} \right| \right) (1+t)^2 e^{-\sigma^* t} \quad \text{for } t \geq 0.$$

This particularly implies that, for any $0 < \sigma < \sigma^*$ and $\delta \in (0, 1)$, we have

$$(6.9) \quad |\zeta(t)| + \sup_{\delta \leq r \leq 1} |\phi(r, t)| \leq C \left(|\zeta_0| + \sup_{0 \leq r \leq 1} |\phi_0(r)| \right) e^{-\sigma t} \quad \text{for } t \geq 0,$$

because $1/\hat{\phi}(r)$ has a positive lower bound for $\delta \leq r \leq 1$ and a finite upper bound for $0 \leq r \leq 1$. In what follows we prove that, for δ sufficiently small, there also holds

$$(6.10) \quad \sup_{0 \leq r \leq \delta} |\phi(r, t)| \leq C e^{-\mu t} \quad \text{for } t \geq 0$$

for some $\mu > 0$.

Take a nonnegative cut-off function $\varphi \in C[0, 1]$ such that

$$\varphi(r) \leq 1 \quad \text{for } 0 \leq r \leq 1, \quad \varphi(r) = 1 \quad \text{for } 0 \leq r \leq \delta, \quad \varphi(r) = 0 \quad \text{for } 2\delta \leq r \leq 1.$$

We split \mathcal{B} into a sum of two operators as follows:

$$(6.11) \quad \mathcal{B}(q) = \mathcal{B}_1(q) + \mathcal{B}_2(q) \quad \text{for } q \in C[0, 1],$$

where

$$\begin{aligned} \mathcal{B}_1(q) &= -r p'_*(r) \varphi(r) \cdot \frac{1}{r^3} \int_0^{\min\{r, \delta\}} g_p(\rho) q(\rho) \rho^2 d\rho, \\ \mathcal{B}_2(q) &= r p'_*(r) \int_0^1 g_p(\rho) q(\rho) \rho^2 d\rho - r^{-2} p'_*(r) [1 - \varphi(r)] \int_0^{\min\{r, \delta\}} g_p(\rho) q(\rho) \rho^2 d\rho \\ &\quad - r^{-2} p'_*(r) \int_{\min\{r, \delta\}}^r g_p(\rho) q(\rho) \rho^2 d\rho \end{aligned}$$

and introduce $f(r, t) = \mathcal{B}_2(\phi(\cdot, t))(r) + b(r)\zeta(t)$. By (6.4) and (6.11), we see that ϕ is the solution of the equation

$$(6.12) \quad \partial_t \phi + u_*(r) \partial_r \phi = a(r) \phi + \mathcal{B}_1(\phi) + f(r, t) \quad \text{for } 0 \leq r \leq 1, \quad t > 0$$

subject to the initial condition $\phi(r, 0) = \phi_0(r)$. Introducing operators $\mathcal{L}(q) = -u_* q' + aq$, we see that (6.12) can be rewritten as the following differential equation in $C[0, 1]$:

$$(6.13) \quad \frac{d\phi(\cdot, t)}{dt} = (\mathcal{L} + \mathcal{B}_1)(\phi(\cdot, t)) + f(\cdot, t).$$

Using Lemma 6.1 we see that the operator \mathcal{L} generates a strongly continuous semigroup $e^{t\mathcal{L}}$ on $C[0, 1]$, and $\|e^{t\mathcal{L}}\| \leq e^{-\omega t}$ for all $t \geq 0$ where $\omega = \min_{0 \leq r \leq 1} |a(r)| > 0$ and using the facts that $\lim_{r \rightarrow 0} r p'_*(r) = 0$ and $\varphi(r) = 0$ for $r \geq 2\delta$, we can easily deduce that, for any given $\varepsilon > 0$, there exists corresponding $\delta > 0$ such that $\|\mathcal{B}_1(q)\|_\infty \leq \varepsilon \|q\|_\infty$, or, in other words, \mathcal{B}_1 is a bounded linear operator on $C[0, 1]$, with norm dominated by ε . Thus, the operator $\mathcal{L} + \mathcal{B}_1$ also generates a strongly continuous semigroup $e^{t(\mathcal{L} + \mathcal{B}_1)}$ on $C[0, 1]$, and it satisfies

$$(6.14) \quad \|e^{t(\mathcal{L} + \mathcal{B}_1)}\| \leq e^{-(\omega - \varepsilon)t} \quad \text{for all } t \geq 0.$$

In what follows we assume that ε is sufficiently small such that $\omega - \varepsilon > 0$. By (6.13), we have

$$(6.15) \quad \phi(\cdot, t) = e^{t(\mathcal{L} + \mathcal{B}_1)} \phi_0 + \int_0^t e^{(t-\tau)(\mathcal{L} + \mathcal{B}_1)} f(\cdot, \tau) d\tau.$$

Using (6.7), (6.8), and (6.9), we can easily show that $\|f(\cdot, t)\|_\infty \leq C_\delta \|U_0\| e^{-\sigma t}$ for $t \geq 0$. Hence, from (6.15) and (6.14) we see that, for any $0 < \mu < \min\{\sigma, \omega - \varepsilon\}$, there holds

$$\|\phi(\cdot, t)\|_\infty \leq \|\phi_0\|_\infty e^{-(\omega - \varepsilon)t} + C \|U_0\| e^{-\mu t} \quad \text{for } t \geq 0,$$

from which (6.10) immediately follows.

Combining (6.9) and (6.10), we get (6.3). This completes the proof. \square

LEMMA 6.3. *Let μ^* be as in Lemma 6.2. Then, for any $0 < \mu < \mu^*$, in addition to (6.3), we also have the following estimate:*

$$(6.16) \quad \|e^{t\Lambda(0)}\|_{L(C^1_V[0,1])} \leq C e^{-\mu t} \quad \text{for } t \geq 0.$$

Proof. We first show that \mathcal{B} and \mathcal{F} satisfy the following properties: For any $q \in C^1_V[0, 1]$,

$$(6.17) \quad \|\mathcal{B}(u_*q')\|_\infty + \|u_*\mathcal{B}(q')\|_\infty \leq C\|q\|_\infty,$$

$$(6.18) \quad |\mathcal{F}(u_*q')| \leq C\|q\|_\infty.$$

Using the facts that $\lim_{r \rightarrow 0} r p'_*(r) = 0$ and $\lim_{r \rightarrow 0} r^2 p''_*(r) = 0$ (see (3.3)), we can easily prove that $\|u_*\mathcal{B}(q')\|_\infty \leq C\|q\|_\infty$ for $q \in C[0, 1]$. To estimate $\|\mathcal{B}(u_*q')\|_\infty$, we compute

$$\frac{1}{r^3} \int_0^r g_p(\rho) u_*(\rho) q'(\rho) \rho^2 d\rho = \frac{1}{r} u_*(r) g_p(r) q(r) - \frac{1}{r^3} \int_0^r m(\rho) q(\rho) \rho^2 d\rho,$$

where $m(\rho) = g'_p(\rho) u_*(\rho) + g_p(\rho) u'_*(\rho) + \frac{2}{\rho} u_*(\rho) g_p(\rho)$. Taking $r = 1$ we particularly obtain

$$\int_0^1 g_p(\rho) u_*(\rho) q'(\rho) \rho^2 d\rho = u_*(1) g_p(1) q(1) - \int_0^1 m(\rho) q(\rho) \rho^2 d\rho.$$

Since $g_p \in C^1[0, 1]$, $u_* \in C^1[0, 1]$, and $u_*(0) = 0$, we see that $\frac{1}{r} u_* g_p$ and m both belong to $C[0, 1]$. Hence, from the above expressions we see immediately that

$$\begin{aligned} \|\mathcal{B}(u_*q')\|_\infty &= \sup_{0 < r < 1} \left| r p'_*(r) \left[\int_0^1 g_p(\rho) u_*(\rho) q'(\rho) \rho^2 d\rho - \frac{1}{r^3} \int_0^r g_p(\rho) u_*(\rho) q'(\rho) \rho^2 d\rho \right] \right| \\ &\leq C\|q\|_\infty. \end{aligned}$$

Similarly we have $|\mathcal{F}(u_*q')| = \left| \int_0^1 g_p(\rho) u_*(\rho) q'(\rho) \rho^2 d\rho \right| \leq C\|q\|_\infty$. This proves (6.17) and (6.18).

We now proceed to prove (6.16). Let $U_0 \in X_0$ and $U = e^{t\Lambda(0)} U_0$. From the proof of Lemma 4.2 we know that $U \in C([0, \infty), X_0) \cap C^1([0, \infty), X)$. Let $U_0 = (q_0, y_0)$ and $U = (q, y)$. Then (q, y) is the solution of the following problem:

$$\begin{cases} \frac{\partial q}{\partial t} + u_*(r) \frac{\partial q}{\partial r} = a(r)q + \mathcal{B}(q) + b(r)y & \text{for } 0 \leq r \leq 1, \quad t > 0, \\ \frac{dy}{dt} = \mathcal{F}(q) + \kappa y & \text{for } t > 0, \\ q|_{t=0} = q_0(r) & \text{for } 0 \leq r \leq 1, \quad \text{and } y|_{t=0} = y_0. \end{cases}$$

Let $l(r, t) = u_*(r) \frac{\partial q(r, t)}{\partial r}$. By formally differentiating the first equation above in r and multiplying it with $u_*(r)$, we see that (l, y) is a formal solution of the following problem:

$$\begin{cases} \frac{\partial l}{\partial t} + u_*(r) \frac{\partial l}{\partial r} = a(r)l + \mathcal{B}(l) + b(r)y + f_1(r, t) & \text{for } 0 \leq r \leq 1, \quad t > 0, \\ \frac{dy}{dt} = \mathcal{F}(l) + \kappa y + c_1(t) & \text{for } t > 0, \\ l|_{t=0} = l_0(r) & \text{for } 0 \leq r \leq 1, \quad \text{and } y|_{t=0} = y_0, \end{cases}$$

where $l_0(r) = u_*(r)q'_0(r)$, $c_1(t) = \mathcal{F}(q) - \mathcal{F}(u_* \frac{\partial q}{\partial r})$, and

$$f_1(r, t) = u_*(r)a'(r)q(r, t) - \mathcal{B} \left(u_* \frac{\partial q}{\partial r} \right) + u_*(r) \frac{\partial \mathcal{B}(q)}{\partial r} + [u_*(r)b'(r) - b(r)]y(t).$$

We denote $W(t) = (l(\cdot, t), y(t))$, $W_0 = (l_0, y_0)$, and $F_1(t) = (f_1(\cdot, t), c_1(t))$. Then the above problem can be rewritten as follows:

$$\frac{dW}{dt} = \mathbb{A}(0)W + F_1(t) \quad \text{for } t > 0, \quad W(0) = W_0.$$

Using the fact that $U \in C^1([0, \infty), X)$, Corollary 3.2, (6.17), and (6.18), we can easily prove that $F_1 \in C^1([0, \infty), X)$. Thus, by a standard result, we see that the above problem has a unique mild solution, and consequently, the above formal computation makes sense. Moreover, we have

$$W(t) = e^{t\mathbb{A}(0)}W_0 + \int_0^t e^{(t-s)\mathbb{A}(0)}F_1(s)ds \quad \text{for } t \geq 0.$$

From this expression and Lemma 6.2 we see that, for any given $0 < \mu < \mu_*$, we have

$$\|W(t)\|_X \leq Ce^{-\mu t}\|W_0\|_X + C \int_0^t e^{-\mu(t-s)}\|F_1(s)\|_X ds \quad \text{for } t \geq 0.$$

Using (6.17), (6.18), and the fact that $\|U(t)\|_X \leq Ce^{-\mu t}\|U_0\|_X$ ensured by (6.3), we see that

$$\|F_1(t)\|_X \leq C\|U(t)\|_X \leq Ce^{-\mu t}\|U_0\|_X \quad \text{for } t \geq 0.$$

Substituting this estimate into the above inequality and noticing that $W(t) = (u_* \frac{\partial q(\cdot, t)}{\partial r}, y(t))$ and $W_0 = (u_*q'_0, y_0)$, we obtain

$$\|U(t)\|_{X_0} \leq C(1+t)e^{-\mu t}\|U_0\|_{X_0} \quad \text{for } t \geq 0.$$

Now, for any given $0 < \mu < \mu_*$, we arbitrarily take a $\bar{\mu} \in (\mu, \mu_*)$ and first use the above estimate to $\bar{\mu}$ and next use the elementary inequality $(1+t)e^{-\bar{\mu}t} \leq Ce^{-\mu t}$, we see that (6.16) follows. This completes the proof. \square

In what follows we consider the evolution system $\mathbb{U}(t, s, V)$ for a general $V = V(t) \in C([0, \infty), X)$ satisfying the following condition: For some positive constants $\bar{\mu}$, ε and C_0 ,

$$(6.19) \quad \|V(t)\|_X \leq C_0\varepsilon e^{-\bar{\mu}t} \quad \text{for } t \geq 0.$$

LEMMA 6.4. *Assume that $V = V(t) \in C([0, \infty), X)$, and it satisfies (6.19). Let μ^* be as in Lemma 6.2. Then, for any $0 < \mu < \mu^*$, there exists corresponding $\varepsilon_0 > 0$ (depending on μ , $\bar{\mu}$, and C_0) such that if $0 < \varepsilon \leq \varepsilon_0$, then the following estimates hold:*

$$(6.20) \quad \|\mathbb{U}(t, s, V)\|_{L(X)} \leq C_1 e^{-\mu t} \quad \text{for } t \geq 0,$$

$$(6.21) \quad \|\mathbb{U}(t, s, V)\|_{L(X_0)} \leq C_2 e^{-\mu t} \quad \text{for } t \geq 0,$$

where C_1 and C_2 are positive constants depending only on μ and independent of $\bar{\mu}$ and C_0 .

Proof. Given $0 < \mu < \mu^*$ we take a $\mu_1 \in (\mu, \mu^*)$ and fix it. By Lemmas 6.2 and 6.3, we have the following estimates:

$$(6.22) \quad \|e^{t\mathbb{A}(0)}\|_{L(C[0,1])} \leq C_1 e^{-\mu_1 t} \quad \text{for } t \geq 0,$$

$$(6.23) \quad \|e^{t\mathbb{A}(0)}\|_{L(C^1_V[0,1])} \leq C_2 e^{-\mu_1 t} \quad \text{for } t \geq 0.$$

Let $U_0 = (q_0, s_0)$ be an arbitrary point in X , and let $U = \mathbb{U}(t, s, V)U_0$. By definition, U is the solution of the problem (4.12). Let $U = (q, y)$. Then (4.12) can be rewritten as follows:

$$(6.24) \quad \begin{cases} \frac{\partial q}{\partial t} + w_V(r, t) \frac{\partial q}{\partial r} = a(r)q + \mathcal{B}q + b(r)y & \text{for } 0 \leq r \leq 1, \quad t > s, \\ \frac{dy}{dt} = \mathcal{F}(q) + \kappa y & \text{for } t > s, \\ q|_{t=s} = q_0(r) & \text{for } 0 \leq r \leq 1, \quad \text{and } y|_{t=s} = y_0. \end{cases}$$

Let $\bar{r} = T(r, t, s)$ and $r = S(\bar{r}, t, s)$ be as defined in section 5, and let $\tilde{q}(\bar{r}, t, s) = q(S(\bar{r}, t, s), t)$, or, equivalently, $q(r, t) = \tilde{q}(T(r, t, s), t, s)$. Then by using (5.7), we see that (6.24) is transformed into the following problem:

$$(6.25) \quad \begin{cases} \frac{\partial \tilde{q}}{\partial t} + u_*(\bar{r}) \frac{\partial \tilde{q}}{\partial \bar{r}} = \tilde{a}(\bar{r}, t, s)\tilde{q} + \tilde{\mathcal{B}}\tilde{q} + \tilde{b}(\bar{r}, t, s)s & \text{for } 0 \leq \bar{r} \leq 1, \quad t > s, \\ \frac{ds}{dt} = \tilde{\mathcal{F}}(\tilde{q})(t, s) + \kappa s & \text{for } t > s, \\ \tilde{q}|_{t=s} = q_0(\bar{r}) & \text{for } 0 \leq \bar{r} \leq 1, \quad \text{and } s|_{t=s} = s_0, \end{cases}$$

where $\tilde{a}(\bar{r}, t, s) = a(S(\bar{r}, t, s))$, $\tilde{b}(\bar{r}, t, s) = b(S(\bar{r}, t, s))$,

$$\tilde{\mathcal{B}}\tilde{q} = rp'_*(r) \left[\int_0^1 g_p(\rho)\tilde{q}(T(\rho, t, s), t, s)\rho^2 d\rho - \frac{1}{r^3} \int_0^r g_p(\rho)\tilde{q}(T(\rho, t, s), t, s)\rho^2 d\rho \right] \Big|_{r=S(\bar{r}, t, s)},$$

and $\tilde{\mathcal{F}}(\tilde{q})(t, s) = \int_0^1 g_p(\rho)\tilde{q}(T(\rho, t, s), t, s)\rho^2 d\rho$. We define a family of bounded linear operators $\tilde{\mathbb{B}}(t, s, V) : X \rightarrow X$ ($t \geq s \geq 0$) as follows:

$$\tilde{\mathbb{B}}(t, s, V) = \begin{pmatrix} \tilde{a}(\cdot, t, s) + \tilde{\mathcal{B}} & \tilde{b}(\cdot, t, s) \\ \tilde{\mathcal{F}} & \kappa \end{pmatrix}.$$

We also denote $\tilde{U} = (\tilde{q}, y)$. Then (6.25) can be rewritten as follows:

$$(6.26) \quad \frac{d\tilde{U}}{dt} = \mathbb{A}_0(U_*)\tilde{U} + \tilde{\mathbb{B}}(t, s, V)\tilde{U} \quad \text{for } t > s, \quad \tilde{U}|_{t=s} = U_0.$$

Recalling that $\mathbb{A}(0) = \mathbb{A}_0(U_*) + \mathbb{B}$ and denoting

$$\tilde{\mathbb{E}}(t, s, V) = \tilde{\mathbb{B}}(t, s, V) - \mathbb{B} = \begin{pmatrix} \tilde{a}(\cdot, t, s) - a + \tilde{\mathcal{B}} - \mathcal{B} & \tilde{b}(\cdot, t, s) - b \\ \tilde{\mathcal{F}} - \mathcal{F} & 0 \end{pmatrix},$$

we see that

$$\mathbb{A}_0(U_*) + \tilde{\mathbb{E}}(t, s, V) = \mathbb{A}_0(U_*) + \mathbb{B} + \tilde{\mathbb{E}}(t, s, V) = \mathbb{A}(0) + \tilde{\mathbb{E}}(t, s, V).$$

Hence, (6.26) can be further rewritten as follows:

$$(6.27) \quad \frac{d\tilde{U}}{dt} = \mathbb{A}(0)\tilde{U} + \tilde{\mathbb{E}}(t, s, V)\tilde{U} \quad \text{for } t > s, \quad \tilde{U}|_{t=s} = U_0.$$

We know that (6.27) is equivalent to the following integral equation:

$$(6.28) \quad \tilde{U}(t, s) = e^{(t-s)\mathbb{A}(0)}U_0 + \int_s^t e^{(t-\tau)\mathbb{A}(0)}\tilde{\mathbb{E}}(\tau, s, V)\tilde{U}(\tau, s) \quad \text{for } t \geq s.$$

It can be easily shown that, under the assumption (6.19), w_V satisfies the estimates (5.1) and (5.27) (with μ replaced by $\bar{\mu}$). Hence, by Corollary 3.2 and Lemma 5.9, we have

$$\|\tilde{a}(\cdot, t, s) - a\|_{C_V^\downarrow[0,1]} \leq C\varepsilon, \quad \|\tilde{b}(\cdot, t, s) - b\|_{C_V^\downarrow[0,1]} \leq C\varepsilon,$$

and by Corollary 3.2, Lemma 5.9, and Lemma 5.10, we have

$$\|\tilde{\mathcal{B}} - \mathcal{B}\|_{L(C[0,1])} \leq C\varepsilon, \quad \|\tilde{\mathcal{B}} - \mathcal{B}\|_{L(C_V^\downarrow[0,1])} \leq C\varepsilon, \quad \|\tilde{\mathcal{F}} - \mathcal{F}\|_{L(C[0,1],\mathbb{R})} \leq C\varepsilon.$$

It follows that

$$(6.29) \quad \|\tilde{\mathbb{E}}(\tau, s, V)\|_{L(X)} \leq C\varepsilon, \quad \|\tilde{\mathbb{E}}(\tau, s, V)\|_{L(X_0)} \leq C\varepsilon.$$

From (6.22), (6.23), (6.28), and (6.29) we obtain

$$\begin{aligned} \|\tilde{U}(t, s)\|_X &\leq C_1 e^{-\mu_1(t-s)}\|U_0\|_X + C\varepsilon \int_s^t e^{-\mu_1(t-\tau)}\|\tilde{U}(\tau, s)\|_X, \\ \|\tilde{U}(t, s)\|_{X_0} &\leq C_2 e^{-\mu_1(t-s)}\|U_0\|_{X_0} + C\varepsilon \int_s^t e^{-\mu_1(t-\tau)}\|\tilde{U}(\tau, s)\|_{X_0}. \end{aligned}$$

By the Gronwall lemma, these inequalities yield

$$\|\tilde{U}(t, s)\|_X \leq C_1 e^{-(\mu_1 - C\varepsilon)t}\|U_0\|_X, \quad \|\tilde{U}(t, s)\|_{X_0} \leq C_2 e^{-(\mu_1 - C\varepsilon)t}\|U_0\|_{X_0}.$$

Hence, by taking ε sufficiently small such that $\mu_1 - C\varepsilon \geq \mu$, we obtain (6.20) and (6.21). This completes the proof. \square

7. The proof of Theorem 1.1. In order to prove Theorem 1.1, we let μ^* be as in Lemma 6.2 and arbitrarily fix a number $0 < \mu < \mu^*$. Let ε be a positive number to be specified later. For any fixed $U_0 \in X_0$ satisfying $\|U_0\|_{X_0} \leq \varepsilon$, we denote by \mathbf{M} the set of all functions $V = V(t) \in C([0, \infty), X)$ satisfying the following conditions:

$$(7.1) \quad V(0) = U_0, \quad \|V(t)\|_X \leq 2C_1\varepsilon e^{-\mu t} \quad \text{for } t \geq 0,$$

where C_1 is the constant appearing in (6.20). We introduce a metric d on \mathbf{M} by defining $d(V_1, V_2) = \sup_{t \geq 0} e^{\mu t}\|V_1(t) - V_2(t)\|_X$ for $V_1, V_2 \in \mathbf{M}$. It is evident that (\mathbf{M}, d) is a complete metric space. Given $V \in \mathbf{M}$, we consider the following initial value problem:

$$(7.2) \quad \frac{dU(t)}{dt} = \mathbb{A}(V(t))U(t) + \mathbb{G}(U(t)) \quad \text{for } t > 0, \quad U(0) = U_0.$$

LEMMA 7.1. *If ε is sufficiently small, then, for any $V \in \mathbf{M}$, problem (7.2) has a unique solution $U \in C([0, \infty), X_0) \cap C^1([0, \infty), X)$, which satisfies the following estimates:*

$$(7.3) \quad \|U(t)\|_X \leq 2C_1\varepsilon e^{-\mu t}, \quad \|U(t)\|_{X_0} \leq C\varepsilon e^{-\mu t}, \quad \|U'(t)\|_X \leq C\varepsilon e^{-\mu t} \quad \text{for } t \geq 0,$$

where C_1 is as before, and C is another constant independent of V .

Proof. We denote

$$\widetilde{\mathbf{M}} = \{U \in C([0, \infty), X_0) : \|U(t)\|_X \leq 2C_1\varepsilon e^{-\mu t} \text{ and } \|U(t)\|_{X_0} \leq 2C_2\varepsilon e^{-\mu t} \text{ for } t \geq 0\}$$

and introduce a metric d on it by defining $d(U_1, U_2) = \sup_{t \geq 0} e^{\mu t} \|U_1(t) - U_2(t)\|_{X_0}$ for $U_1, U_2 \in \widetilde{\mathbf{M}}$. Here, C_1 and C_2 are positive constants appearing in (6.20) and (6.21), respectively. $(\widetilde{\mathbf{M}}, d)$ is clearly a complete metric space. Given $U \in \widetilde{\mathbf{M}}$, we consider the following initial value problem:

$$(7.4) \quad \frac{d\widetilde{U}(t)}{dt} = \mathbb{A}(V(t))\widetilde{U}(t) + \mathbb{G}(U(t)) \quad \text{for } t > 0, \quad \widetilde{U}(0) = U_0.$$

Since $U(t) \in C([0, \infty), X_0)$, by Corollary 3.3, we have $\mathbb{G}(U(t)) \in C([0, \infty), X_0)$. It follows by Corollary 4.4 that the above problem has a unique solution $\widetilde{U} \in C([0, \infty), X_0) \cap C^1([0, \infty), X)$ and is given by

$$(7.5) \quad \widetilde{U}(t) = \mathbb{U}(t, 0, V)U_0 + \int_0^t \mathbb{U}(t, s, V)\mathbb{G}(U(s))ds.$$

Using this expression and Lemma 6.4 and (2.18) we have

$$\begin{aligned} \|\widetilde{U}(t)\|_X &\leq C_1 e^{-\mu t} \|U_0\|_X + C_1 \int_0^t e^{-\mu(t-s)} \|\mathbb{G}(U(s))\|_X ds \\ &\leq C_1 \varepsilon e^{-\mu t} + C \int_0^t e^{-\mu(t-s)} \|U(s)\|_X^2 ds \leq 2C_1 \varepsilon e^{-\mu t}, \end{aligned}$$

when ε is sufficiently small. Similarly, by using Lemma 6.4 and (3.9), we also have $\|\widetilde{U}(t)\|_{X_0} \leq 2C_2 \varepsilon e^{-\mu t}$, when ε is sufficiently small. Hence, $\widetilde{U} \in \widetilde{\mathbf{M}}$. We now define a mapping $\widetilde{\mathbf{S}} : \widetilde{\mathbf{M}} \rightarrow \widetilde{\mathbf{M}}$ by setting $\widetilde{\mathbf{S}}(U) = \widetilde{U}$ for every $U \in \widetilde{\mathbf{M}}$. We claim that $\widetilde{\mathbf{S}}$ is a contraction mapping. Indeed, for any $U_1, U_2 \in \widetilde{\mathbf{M}}$, let $\widetilde{U}_1 = \mathbf{S}(U_1)$, $\widetilde{U}_2 = \mathbf{S}(U_2)$, and $W = \widetilde{U}_1 - \widetilde{U}_2$. Then W satisfies

$$\frac{dW(t)}{dt} = \mathbb{A}(V(t))W(t) + [\mathbb{G}(U_1(t)) - \mathbb{G}(U_2(t))] \quad \text{for } t > 0, \quad W(0) = 0,$$

so that $W(t) = \int_0^t \mathbb{U}(t, s, V)[\mathbb{G}(U_1(s)) - \mathbb{G}(U_2(s))]ds$. It follows that

$$\begin{aligned} \|W(t)\|_{X_0} &\leq C_2 \int_0^t e^{-\mu(t-s)} \|\mathbb{G}(U_1(s)) - \mathbb{G}(U_2(s))\|_{X_0} ds \\ &\leq C_2 \int_0^t e^{-\mu(t-s)} \|U_1(s) - U_2(s)\|_{X_0} \\ &\quad \left(\int_0^1 \|\mathbb{G}'(\theta U_1(s) + (1-\theta)U_2(s))\|_{L(X_0)} d\theta \right) ds \\ &\leq C \int_0^t e^{-\mu(t-s)} \|U_1(s) - U_2(s)\|_{X_0} \left(\int_0^1 \|\theta U_1(s) + (1-\theta)U_2(s)\|_{X_0} d\theta \right) ds, \end{aligned}$$

which yields $\|W(t)\|_{X_0} \leq C\varepsilon d(U_1, U_2)e^{-\mu t}$. Thus, for ε sufficiently small, we have $d(\widetilde{U}_1, \widetilde{U}_2) = \sup_{t \geq 0} e^{\mu t} \|W(t)\|_{X_0} \leq \frac{1}{2}d(U_1, U_2)$, showing that $\widetilde{\mathbf{S}}$ is a contraction mapping, as we claimed. Thus, by the Banach fixed point theorem, we see that $\widetilde{\mathbf{S}}$ has a unique fixed point in $\widetilde{\mathbf{M}}$, which is clearly a solution of the problem (7.2) in $C([0, \infty), X_0)$. The uniqueness of the solution follows from a standard argument.

From the above argument we see that the solution U of (7.2) satisfies the first two inequalities in (7.3), and $U \in C^1([0, \infty), X)$. It remains to prove that U also satisfies the last inequality in (7.3). The argument is as follows. First, it is straightforward to deduce from condition (7.1) that, for sufficiently small $\varepsilon > 0$, we have $w_V(r, t)/r(1-r) \in C[0, 1]$, and there exist positive constants C_1 and C_2 independent of V such that

$$(7.6) \quad -C_1 r(1-r) \leq w_V(r, t) \leq -C_2 r(1-r) \quad \text{for } 0 \leq r \leq 1 \text{ and } t \geq 0.$$

It follows that, for the solution $U = (q, s)$ of (7.2), we have

$$\sup_{0 \leq r \leq 1} \left| w_V(r, t) \frac{\partial q(r, t)}{\partial r} \right| \leq C \sup_{0 \leq r \leq 1} \left| r(1-r) \frac{\partial q(r, t)}{\partial r} \right| \leq C \|q(\cdot, t)\|_{C_V^1[0,1]}.$$

Using this result and (7.2), we see that

$$\|U'(t)\|_X \leq \|\mathbb{A}(V(t))U(t)\|_X + \|\mathbb{G}(U(t))\|_X \leq C\|U(t)\|_{X_0} + C\|U(t)\|_X^2 \leq C\varepsilon e^{-\mu t}$$

for all $t \geq 0$. This completes the proof of Lemma 7.1. \square

Lemma 7.1, in particular, implies that, for every V in \mathbf{M} , the solution U of (7.2) also belongs to \mathbf{M} . Thus we can define a mapping $\mathbf{S} : \mathbf{M} \rightarrow \mathbf{M}$ as follows: For any $V \in \mathbf{M}$,

$$\mathbf{S}(V) = U = \text{the solution of (7.2)}.$$

LEMMA 7.2. *For ε sufficiently small, \mathbf{S} is a contraction mapping.*

Proof. Let $V_1, V_2 \in \mathbf{M}$ and denote $U_1 = \mathbf{S}(V_1)$, $U_2 = \mathbf{S}(V_2)$, and $W = U_1 - U_2$. Then W satisfies

$$\frac{dW(t)}{dt} = \mathbb{A}(V_1(t))W(t) + [\mathbb{A}(V_1(t)) - \mathbb{A}(V_2(t))]U_2(t) + [\mathbb{G}(U_1(t)) - \mathbb{G}(U_2(t))]$$

for $t > 0$, and $W(0) = 0$. Thus

$$(7.7) \quad \begin{aligned} W(t) &= \int_0^t \mathbb{U}(t, s, V_1) [\mathbb{A}(V_1(s)) - \mathbb{A}(V_2(s))] U_2(s) ds \\ &\quad + \int_0^t \mathbb{U}(t, s, V_1) [\mathbb{G}(U_1(s)) - \mathbb{G}(U_2(s))] ds. \end{aligned}$$

Since the first component of $[\mathbb{A}(V_1(s)) - \mathbb{A}(V_2(s))]U_2(s)$ is $[w_{V_2}(r, s) - w_{V_1}(r, s)]q_2'(r, s)$ and the second component is zero, we have

$$(7.8) \quad \begin{aligned} \|[\mathbb{A}(V_1(s)) - \mathbb{A}(V_2(s))]U_2(s)\|_X &= \max_{0 \leq r \leq 1} |[w_{V_1}(r, s) - w_{V_2}(r, s)]q_2'(r, s)| \\ &\leq \sup_{0 \leq r \leq 1} \left| \frac{w_{V_1}(r, s) - w_{V_2}(r, s)}{r(1-r)} \right| \max_{0 \leq r \leq 1} |r(1-r)q_2'(r, s)| \\ &\leq C\|V_1(s) - V_2(s)\|_X \|U_2(s)\|_{X_0}. \end{aligned}$$

Besides, from (2.18) we have

$$(7.9) \quad \begin{aligned} \|\mathbb{G}(U_1(s)) - \mathbb{G}(U_2(s))\|_X &= \left\| \int_0^1 \mathbb{G}'(\theta U_1(s) + (1-\theta)U_2(s))[U_1(s) - U_2(s)]d\theta \right\|_X \\ &\leq C(\|U_1(s)\|_X + \|U_2(s)\|_X)\|U_1(s) - U_2(s)\|_X. \end{aligned}$$

From (7.7)–(7.9) and Lemma 6.4 we get

$$\begin{aligned} \|U_1(t) - U_2(t)\|_X &\leq C \int_0^t e^{-\mu(t-s)} \|V_1(s) - V_2(s)\|_X \|U_2(s)\|_{X_0} ds \\ &\quad + C \int_0^t e^{-\mu(t-s)} (\|U_1(s)\|_X + \|U_2(s)\|_X) \|U_1(s) - U_2(s)\|_X ds \\ &\leq C\epsilon e^{-\mu t} d(V_1, V_2) + C\epsilon e^{-\mu t} d(U_1, U_2), \end{aligned}$$

which yields $d(U_1, U_2) \leq C\epsilon d(V_1, V_2) + C\epsilon d(U_1, U_2)$. From this inequality the desired assertion easily follows. \square

By Lemma 7.2, if ϵ is sufficiently small, then the mapping \mathbf{S} has a unique fixed point U in \mathbf{M} . Clearly, U is a global solution of (2.16) subject to the initial condition $U(0) = U_0$. Moreover, by Lemma 7.1, we know that the image of \mathbf{S} is contained in $\widetilde{\mathbf{M}}$ so that U satisfies (7.3). From this result and Lemma 2.1 we see that the assertion of Theorem 1.1 follows. This completes the proof of Theorem 1.1.

Acknowledgments. Part of this work was prepared when the author was visiting the Ecole Normale Supérieure (ENS) during June 1–August 28, 2007. He wishes to acknowledge his gratefulness to all staff and faculty of the Department of Mathematics and Applications of ENS, particularly Prof. Perthame and Prof. Rosso, for their hospitality, and the French Ministry of Foreign Affairs, particularly the Section of Science and Technology of the French Consulate at Guangzhou, for financial support for his visit. He is also glad to express thanks to the anonymous referees for helpful suggestions.

REFERENCES

- [1] R. P. ARAUJO AND D. L. MCELWAIN, *A history of the study of solid tumor growth: The contribution of mathematical modeling*, Bull. Math. Biol., 66 (2004), pp. 1039–1091.
- [2] B. BAZALIY AND A. FRIEDMAN, *A free boundary problem for an elliptic-parabolic system: Application to a model of tumor growth*, Comm. Partial Differential Equations, 28 (2003), pp. 517–560.
- [3] B. BAZALIY AND A. FRIEDMAN, *Global existence and asymptotic stability for an elliptic-parabolic free boundary problem: An application to a model of tumor growth*, Indiana Univ. Math. J., 52 (2003), pp. 1265–1304.
- [4] X. CHEN AND A. FRIEDMAN, *A free boundary problem for an elliptic-hyperbolic system: An application to tumor growth*, SIAM J. Math. Anal., 35 (2003), pp. 974–986.
- [5] X. CHEN, S. CUI, AND A. FRIEDMAN, *A hyperbolic free boundary problem modeling tumor growth: Asymptotic behavior*, Trans. Amer. Math. Soc., 357 (2005), pp. 4771–4804.
- [6] S. CUI, *Analysis of a mathematical model for the growth of tumors under the action of external inhibitors*, J. Math. Biol., 44 (2002), pp. 395–426.
- [7] S. CUI, *Analysis of a free boundary problem modelling tumor growth*, Acta Math. Appl. Sin. Engl. Ser., 21 (2005), pp. 1071–1082.
- [8] S. CUI, *Existence of a stationary solution for the modified Ward-King tumor growth model*, Adv. Appl. Math., 36 (2006), pp. 421–445.
- [9] S. CUI, *Well-posedness of a multidimensional free boundary problem modeling the growth of nonnecrotic tumors*, J. Func. Anal., 245 (2007), pp. 1–18.
- [10] S. CUI AND J. ESCHER, *Bifurcation analysis of an elliptic free boundary problem modelling the growth of avascular tumors*, SIAM J. Math. Anal., 39 (2007), pp. 210–235.

- [11] S. CUI AND J. ESCHER, *Asymptotic behavior of solutions of a multidimensional moving boundary problem modeling tumor growth*, Comm. Partial Differential Equations, 33 (2008), pp. 636–655.
- [12] S. CUI AND A. FRIEDMAN, *A hyperbolic free boundary problem modeling tumor growth*, Interfaces Free Bound., 5 (2003), pp. 159–181.
- [13] S. CUI AND A. FRIEDMAN, *A free boundary problem for a singular system of differential equations: An application to a model of tumor growth*, Trans. Amer. Math. Soc., 355 (2003), pp. 3537–3590.
- [14] S. CUI AND X. WEI, *Existence of solutions for a parabolic-hyperbolic free boundary problem*, Acta Math. Appl. Sin. Engl. Ser., 21 (2005), pp. 597–614.
- [15] A. FRIEDMAN, *A hierarchy of cancer models and their mathematical challenges*, Disc. Cont. Dyna. Syst. B, 4 (2004), pp. 147–159.
- [16] A. FRIEDMAN, *Mathematical analysis and challenges arising from models of tumor growth*, Math. Models Methods Appl. Sci., 17 (2007), supplement, pp. 1751–1772.
- [17] A. FRIEDMAN AND B. HU, *Asymptotic stability for a free boundary problem arising in a tumor model*, J. Differential Equations, 227 (2006), pp. 598–639.
- [18] A. FRIEDMAN AND F. REITICH, *Analysis of a mathematical model for the growth of tumors*, J. Math. Biol., 38 (1999), pp. 262–284.
- [19] A. FRIEDMAN AND F. REITICH, *Symmetry-breaking bifurcation of analytic solutions to free boundary problems*, Trans. Amer. Math. Soc., 353 (2000), pp. 1587–1634.
- [20] A. PAZY, *Semigroups of Linear Operators and Applications to Partial Differential Equations*, Springer, New York, 1983.
- [21] G. PETTET, C. PLEASE, AND M. MCELWAIN, *The migration of cells in multicell tumor spheroids*, Bull. Math. Biol., 63 (2001), pp. 231–257.
- [22] T. ROOSE, S. J. CHAPMAN, AND P. K. MAINI, *Mathematical models of avascular tumor growth*, SIAM Rev., 49 (2007), pp. 179–208.
- [23] J. WU AND S. CUI, *Asymptotic behavior of solutions of a free boundary problem modeling the growth of tumors in the presence of inhibitors*, Nonlinearity, 20 (2007), pp. 2389–2408.
- [24] F. ZHOU AND S. CUI, *Well-posedness and stability of a multidimensional moving boundary problem modeling the growth of tumor cord*, Discrete Contin. Dyn. Syst. Ser. A, 21 (2008), pp. 929–943.

WELL-POSEDNESS AND CONVERGENCE TO THE STEADY STATE FOR A MODEL OF MORPHOGEN TRANSPORT*

PIOTR KRZYŻANOWSKI[†], PHILIPPE LAURENÇOT[‡], AND DARIUSZ WRZOSEK[†]

Abstract. Well-posedness and large time convergence to the unique steady state are shown for a model which describes the spreading of morphogens by a nonlinear transport mechanism (transcytosis) and couples a quasilinear parabolic partial differential equation with an ordinary differential equation. A simpler model which assumes linear transport is also investigated for comparison. The analysis of both models requires the construction of specific Liapunov functionals. The study is supplemented by numerical simulations of the sensitivity of the models to the variation of their parameters.

Key words. degenerate parabolic system, Liapunov functional, steady state

AMS subject classifications. 35K55, 35B40, 37L45, 35Q80

DOI. 10.1137/070711608

1. Introduction. Morphogens are signaling molecules that play an important role in the process of cell differentiation, as different morphogen concentrations induce distinct cell fates. In fact, during cell development, spatial gradients of morphogen concentration form while the morphogens move away from a spatially localized source. It is, however, yet unclear whether the spreading of morphogens occurs by passive diffusion in the extracellular medium or by more complex mechanisms such as planar transcytosis (see, e.g., [5, 11] and the references therein). Planar transcytosis explains the transport of morphogens across the tissue as the consequence of repeated cycles of internalization of morphogens inside the cell (endocytosis) and release of morphogens outside the cell (exocytosis) together with the binding and unbinding of morphogens to receptors located at the surface of the cell. A related model accounting only for the latter mechanism was introduced in [6].

A model for the transport of morphogens by passive diffusion was recently derived in [8, Model B] and describes the space and time evolution of the concentrations of free morphogens ℓ and bound receptors s : in dimensionless form, it reads

$$(1.1) \quad \partial_t \ell = D \partial_x^2 \ell + \delta s - \ell(1 - s), \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(1.2) \quad \partial_t s = \ell(1 - s) - (\delta + \varepsilon)s, \quad (t, x) \in (0, \infty) \times (0, 1),$$

supplemented by suitable boundary and initial conditions. In this model, the total number of receptors (free + bound) per cell is assumed to be constant and normalized to one. The reaction terms account for the depletion of free morphogens resulting from the binding with free receptors ($-\ell(1 - s)$), the release of free morphogens after unbinding from the receptors (δs), and the degradation of bound receptors ($-\varepsilon s$), respectively. An attempt to take into account planar transcytosis is made in

*Received by the editors December 21, 2007; accepted for publication (in revised form) July 11, 2008; published electronically December 17, 2008. This work was partially supported by Polish MEiSW grant 1P03A01730 and by “Projet Hubert Curien” POLONIUM 11605VC (2006/2007).

<http://www.siam.org/journals/sima/40-5/71160.html>

[†]Institute of Applied Mathematics, Warsaw University, Banacha 2, 02-097 Warszawa, Poland (p.krzyzanowski@mimuw.edu.pl, d.wrzosek@mimuw.edu.pl).

[‡]Institut de Mathématiques de Toulouse, CNRS (UMR 5219), and Université de Toulouse, F-31062 Toulouse Cedex 9, France (laurenco@math.univ-toulouse.fr).

[3, 7]: it gives a system of two equations similar to (1.1)–(1.2), with the noticeable difference that it includes a nonconstant diffusion coefficient in the equation for the free morphogens. It describes the fact that free receptors have to be available for the transport of morphogens to take place. Depending on the assumptions made during the derivation, the equation for the free morphogens might also contain an additional drift term of the form $\partial_x(\chi(\ell, s) \partial_x s)$. Let us mention here that such a drift term also shows up in the related model proposed in [9]. We will, however, focus in this paper on a model which does not include an additional drift term of the above-mentioned form but accounts for transcytosis by a nonconstant diffusion coefficient which depends on the concentration of bound receptors and vanishes if no free receptor is available ($s = 1$).

More precisely, the model derived in [7] reads, in dimensionless form,

$$(1.3) \quad \partial_t \ell = D \partial_x((1-s) \partial_x \ell) + \delta s - \ell(1-s), \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(1.4) \quad \partial_t s = \ell(1-s) - (\delta + \varepsilon) s, \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(1.5) \quad D(1-s(t, 0)) \partial_x \ell(t, 0) + \nu = \ell(t, 1) = 0, \quad t \in (0, \infty),$$

$$(1.6) \quad (\ell, s)(0, x) = (\ell_0, s_0)(x), \quad x \in (0, 1),$$

with $D > 0$, $\delta > 0$, $\varepsilon \geq 0$, and $\nu > 0$. The strength of the diffusion is clearly related to the concentration of bound receptors s and no diffusion takes place in the absence of free receptors. As for the reaction terms, they are the same as for (1.1)–(1.2) though their biological meaning is different: Indeed, the term $(-\ell(1-s))$ accounts for the endocytosis process, which combines the binding of a free morphogen with a receptor followed by its internalization, and the term (δs) for the exocytosis process, which includes an externalization step and a dissociation step. Finally, the boundary condition at $x = 0$ accounts for a source of morphogens localized on one side of the tissue. We refer to [7] for a complete derivation of the model, and some additional information as well.

The aim of this paper is to supplement the qualitative study of (1.3)–(1.6) performed in [7] with that of mathematical properties, namely, well-posedness and large time behavior. We first point out that, as (1.3)–(1.6) couples a quasilinear parabolic partial differential equation with an ordinary differential equation, it cannot be studied with the usual techniques for systems of parabolic partial differential equations. Nevertheless, a general theory has been developed in [1] to study this kind of system: applying this theory in a suitable functional framework allows us to establish the local well-posedness of (1.3)–(1.6). A Liapunov functional is then constructed from which we deduce the global well-posedness.

More precisely, we assume that

$$(1.7) \quad (\ell_0, s_0) \in V \times V \quad \text{with} \quad \ell_0 \geq 0 \quad \text{and} \quad 0 \leq s_0 < 1,$$

where $V := \{v \in H^1(0, 1) \text{ such that } v(1) = 0\}$ and V' denotes the topological dual of V .

THEOREM 1. *Given initial data (ℓ_0, s_0) satisfying (1.7) the initial-boundary value problem (1.3)–(1.6) possesses a unique weak solution*

$$(\ell, s) \in \mathcal{C}([0, \infty); V \times V) \cap \mathcal{C}^1([0, \infty); V' \times V)$$

such that $\partial_t \ell \in L^2((0, T) \times (0, 1))$ for every $T > 0$,

$$0 \leq \ell(t, x) \quad \text{and} \quad 0 \leq s(t, x) < 1 \quad \text{for} \quad (t, x) \in [0, \infty) \times [0, 1],$$

and

$$\begin{aligned} \langle \partial_t \ell(t), \psi \rangle_{V',V} + D \int_0^1 (1 - s(t)) \partial_x \ell(t) \partial_x \psi \, dx \\ = \nu \psi(0) + \int_0^1 (\delta s(t) - \ell(t)(1 - s(t))) \psi \, dx, \end{aligned}$$

$$\partial_t s(t) = -(\delta + \varepsilon)s(t) + \ell(t)(1 - s(t))$$

for $\psi \in V$ and $t \in [0, \infty)$. Furthermore, for each $t > 0$,

$$(1.8) \quad \mathcal{L}(s(t), \partial_t \sigma(t)) + \int_0^t \mathcal{D}(s(\tau), \partial_t \sigma(\tau)) \, d\tau = \mathcal{L}(s_0, \partial_t \sigma(0)),$$

where $\sigma := -\ln(1 - s)$ and

$$(1.9) \quad \mathcal{L}(v, w) := \mathcal{L}_0(v) + \frac{1}{2} \int_0^1 w(x)^2 \, dx,$$

$$(1.10) \quad \mathcal{L}_0(v) := \int_0^1 \left\{ \frac{D(\delta + \varepsilon)}{2} |\partial_x \Sigma_I(v)(x)|^2 + \varepsilon (\Sigma_I(v) - v)(x) - \nu \Sigma_I(v)(0) \right\} dx,$$

$$(1.11) \quad \begin{aligned} \mathcal{D}(v, w) := \int_0^1 \left(1 - v(x) + \frac{\delta + \varepsilon}{1 - v(x)} \right) w(x)^2 \, dx \\ + D \int_0^1 (1 - v(x)) |\partial_x w(x)|^2 \, dx, \end{aligned}$$

the function Σ_I being defined by $\Sigma_I(r) := -\ln(1 - r)$ for $r \in [0, 1)$.

We now turn to the behavior for large times of weak solutions to (1.3)–(1.6) and first recall from [7] that (1.3)–(1.5) has a single stationary solution.

PROPOSITION 2. *There is a unique stationary solution $(\ell_\infty, s_\infty) \in \mathcal{C}([0, 1]; \mathbb{R}^2)$ to (1.3)–(1.5) satisfying the natural constraints*

$$(1.12) \quad \ell_\infty(x) \geq 0 \quad \text{and} \quad 0 \leq s_\infty(x) < 1 \quad \text{for} \quad x \in [0, 1].$$

The function s_∞ is the unique solution to the boundary value problem

$$(1.13) \quad \begin{aligned} -D \partial_x^2 \Sigma_I(s_\infty) + \frac{\varepsilon}{\delta + \varepsilon} s_\infty &= 0 \quad \text{in} \quad (0, 1), \\ D \partial_x \Sigma_I(s_\infty)(0) + \frac{\nu}{\delta + \varepsilon} &= s_\infty(1) = 0, \end{aligned}$$

while ℓ_∞ is given by $\ell_\infty(x) = (\delta + \varepsilon)s_\infty(x)/(1 - s_\infty(x))$ for $x \in [0, 1]$.

With this notation, our second result reads as follows.

THEOREM 3. *Consider (ℓ_0, s_0) satisfying (1.7) and let (ℓ, s) be the corresponding weak solution to the initial-boundary value problem (1.3)–(1.6). Then $(\ell(t), s(t))$ converges (strongly) towards (ℓ_∞, s_∞) in $L^2(0, 1) \times H^1(0, 1)$ as $t \rightarrow \infty$.*

Roughly speaking, the only piece of information on the convergence that can be retrieved from (1.8) is that $\partial_t \Sigma_I(s)$ vanishes as $t \rightarrow \infty$ and so does $\partial_t s$. Therefore, by (1.4), $\ell(1 - s) - (\delta + \varepsilon)s$ vanishes for large times, which somehow means that s gets closer and closer to $\ell/(\ell + \delta + \varepsilon)$ as time goes by. However, no information is available

on ℓ at this stage. Nevertheless, it seems reasonable to guess that ℓ behaves as the solution λ to

$$(1.14) \quad \partial_t \lambda = \partial_x \left(\frac{\delta + \varepsilon}{\lambda + \delta + \varepsilon} \partial_x \lambda \right) - \frac{\varepsilon}{\delta + \varepsilon} \frac{\lambda}{\lambda + \delta + \varepsilon}, \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(1.15) \quad D \frac{\partial_x \lambda(t, 0)}{\lambda(t, 0) + \delta + \varepsilon} + \frac{\nu}{\delta + \varepsilon} = \lambda(t, 1) = 0, \quad t \in (0, \infty),$$

obtained from (1.3)–(1.5) by replacing s by $\ell/(\ell + \delta + \varepsilon)$. Noting that ℓ_∞ is also a stationary solution to (1.14)–(1.15) allows us to construct another Liapunov functional for (1.3)–(1.6) resembling that available for (1.14)–(1.15) and providing the needed information on ℓ .

The large time convergence towards the steady state is proved in section 3 after establishing the well-posedness in $V \times V$ in section 2. In section 4, we return to the model (1.1)–(1.2) of morphogen transport by passive diffusion and show that it is also well-posed in $V \times V$ once supplemented by the boundary and initial conditions (1.5)–(1.6). To this end we construct a Liapunov function for this problem which turns out to provide much more information than \mathcal{L} . In particular, exponential convergence to the steady state can be proved for (1.1)–(1.2). Finally, we devote section 5 to numerical simulations, where we show that the solutions of both models converge to the steady state at an exponential rate and that this rate cannot be arbitrarily large regardless of the choice of the parameters.

Throughout the paper, C_i , $i \geq 1$, denotes any positive constant depending only on ℓ_0 , s_0 , D , δ , ε , and ν . For a real number $r \in \mathbb{R}$, $r_+ := \max\{r, 0\}$ denotes the positive part of r and we set $\text{sign}_+(r) := 1$ for $r > 0$ and $\text{sign}_+(r) := 0$ for $r \leq 0$.

2. Well-posedness. We first show that the Liapunov functional \mathcal{L} is bounded from below.

LEMMA 4. *For $v \in V$, $0 \leq v < 1$, we have*

$$(2.1) \quad \mathcal{L}_0(v) \geq \mathcal{L}_0(s_\infty) + \frac{D(\delta + \varepsilon)}{2} \int_0^1 |\partial_x (\Sigma_I(v) - \Sigma_I(s_\infty)) (x)|^2 dx.$$

In particular,

$$\mathcal{L}_0(s_\infty) = \min_{v \in V} \{\mathcal{L}_0(v)\}.$$

Proof. Let $v \in V$, $0 \leq v < 1$. It follows from (1.13) by a straightforward computation that

$$\begin{aligned} \mathcal{L}_0(v) - \mathcal{L}_0(s_\infty) &= \frac{D(\delta + \varepsilon)}{2} \int_0^1 |\partial_x (\Sigma_I(v) - \Sigma_I(s_\infty)) (x)|^2 dx \\ &\quad + \varepsilon \int_0^1 \{(1 - s_\infty) (\Sigma_I(v) - \Sigma_I(s_\infty)) - v + s_\infty\} (x) dx. \end{aligned}$$

Since $s_\infty < 1$ by Proposition 2 we observe that

$$(1 - s_\infty) (\Sigma_I(v) - \Sigma_I(s_\infty)) - v + s_\infty = (1 - s_\infty) [\Sigma_I(v) - \Sigma_I(s_\infty) - \Sigma_I'(s_\infty) (v - s_\infty)] \geq 0$$

by the convexity of Σ_I , and (2.1) follows. \square

Proof of Theorem 1. By [1, Theorem 6.4] (with $G = \mathbb{R} \times (-\infty, 1)$) the initial-boundary value problem (1.3)–(1.6) possesses a unique maximal weak solution

$$(2.2) \quad (\ell, s) \in \mathcal{C}([0, T_m]; V \times H^1(0, 1)) \cap \mathcal{C}^1([0, T_m]; V' \times H^1(0, 1))$$

such that

$$(2.3) \quad s(t, x) < 1, \quad (t, x) \in [0, T_m) \times [0, 1],$$

and

$$(2.4) \quad \begin{aligned} \langle \partial_t \ell(t), \psi \rangle_{V', V} + D \int_0^1 (1 - s(t)) \partial_x \ell(t) \partial_x \psi \, dx \\ = \nu \psi(0) + \int_0^1 (\delta s(t) - \ell(t)(1 - s(t))) \psi \, dx, \end{aligned}$$

$$(2.5) \quad \partial_t s(t) = -(\delta + \varepsilon)s(t) + \ell(t)(1 - s(t))$$

for $\psi \in V$ and $t \in [0, T_m)$. In addition, $T_m = \infty$ if, for each $T > 0$,

$$(2.6) \quad (\ell, s) \in \mathcal{BUC}([0, T] \cap [0, T_m); V \times H^1(0, 1))$$

and there is $\vartheta_T \in (0, 1)$ such that

$$(2.7) \quad \max_{x \in [0, 1]} \{s(t, x)\} \leq \vartheta_T \quad \text{for every } t \in [0, T] \cap [0, T_m).$$

We first establish the nonnegativity of (ℓ, s) . For that purpose we notice that since $1 - s \geq 0$, the right-hand side of (1.3)–(1.4) satisfies

$$\begin{aligned} & [\delta s - \ell(1 - s)] \operatorname{sign}_+(-\ell) + [\ell(1 - s) - (\delta + \varepsilon)s] \operatorname{sign}_+(-s) \\ & \geq -\delta(-s)_+ + (1 - s)(-\ell)_+ - (1 - s)(-\ell)_+ + (\delta + \varepsilon)(-s)_+ \\ & \geq 0. \end{aligned}$$

We may then use classical approximations of $r \mapsto r_+$ to deduce from (2.4) and (2.5) that

$$\int_0^1 [(-\ell)_+(t) + (-s)_+(t)] \, dx \leq \int_0^1 [(-\ell_0)_+ + (-s_0)_+] \, dx = 0$$

for $t \in [0, T_m)$. Consequently

$$(2.8) \quad \ell(t, x) \geq 0 \quad \text{and} \quad s(t, x) \geq 0, \quad (t, x) \in [0, T_m) \times [0, 1].$$

It also follows from (1.7), (2.2), and (2.5) that $\partial_t s(t, 1) = -(\delta + \varepsilon)s(t, 1)$ with $s(0, 1) = 0$, whence

$$(2.9) \quad s(t, 1) = 0 \quad \text{and} \quad s(t) \in V \quad \text{for } t \in [0, T_m).$$

We next turn to the identity (1.8): owing to (2.2), (2.3), and (2.9), the function $\sigma := \Sigma_I(s) = -\ln(1 - s)$ is well-defined in $[0, T_m) \times [0, 1]$ and belongs to $C^1([0, T_m); V)$ with

$$\partial_t \sigma = \frac{\partial_t s}{1 - s} \quad \text{and} \quad \partial_x \sigma = \frac{\partial_x s}{1 - s}.$$

Furthermore (2.5) also reads

$$(2.10) \quad \partial_t \sigma = \ell - (\delta + \varepsilon) \frac{s}{1 - s},$$

which, together with (2.2), ensures that $\partial_t \sigma$ belongs to $\mathcal{C}([0, T_m]; V)$ and $\mathcal{C}^1([0, T_m]; V')$. For $t \in [0, T_m)$, $\partial_t \sigma(t)$ is then an admissible test function in (2.4) and it follows from (2.4) and (2.5) that

$$\begin{aligned} & \langle \partial_t \ell(t), \partial_t \sigma(t) \rangle_{V', V} + D \int_0^1 (1 - s(t)) \partial_x \ell(t) \partial_x \partial_t \sigma(t) \, dx \\ &= - \int_0^1 (\partial_t s(t) + \varepsilon s(t)) \partial_t \sigma(t) \, dx + \nu \partial_t \sigma(t, 0). \end{aligned}$$

On the one hand, expressing ℓ in terms of σ and s with the help of (2.10) yields

$$\langle \partial_t \ell(t), \partial_t \sigma(t) \rangle_{V', V} = \frac{1}{2} \frac{d}{dt} \int_0^1 |\partial_t \sigma(t)|^2 \, dx + (\delta + \varepsilon) \int_0^1 \frac{|\partial_t \sigma(t)|^2}{1 - s(t)} \, dx.$$

On the other hand, using once more (2.10), we obtain

$$\begin{aligned} D \int_0^1 (1 - s(t)) \partial_x \ell(t) \partial_x \partial_t \sigma(t) \, dx &= D \int_0^1 (1 - s(t)) |\partial_x \partial_t \sigma(t)|^2 \, dx \\ &\quad + \frac{D(\delta + \varepsilon)}{2} \frac{d}{dt} \int_0^1 |\partial_x \sigma(t)|^2 \, dx. \end{aligned}$$

Combining the above three identities we end up with

$$\frac{d}{dt} \mathcal{L}(s(t), \partial_t \sigma(t)) + \mathcal{D}(s(t), \partial_t \sigma(t)) = 0, \quad t \in [0, T_m),$$

whence

$$(2.11) \quad \mathcal{L}(s(t), \partial_t \sigma(t)) + \int_0^t \mathcal{D}(s(\tau), \partial_t \sigma(\tau)) \, d\tau = \mathcal{L}(s_0, \partial_t \sigma(0)), \quad t \in [0, T_m).$$

We then deduce from (1.9), (2.11), (2.1), and the nonnegativity of $\mathcal{D}(s, \partial_t \sigma)$ that

$$(2.12) \quad \frac{D(\delta + \varepsilon)}{2} \int_0^1 |\partial_x (\sigma(t) - \Sigma_I(s_\infty))|^2 \, dx + \frac{1}{2} \int_0^1 |\partial_t \sigma(t)|^2 \, dx \leq C_0,$$

with $C_0 := -\mathcal{L}_0(s_\infty) + \mathcal{L}(s_0, \partial_t \sigma(0))$. Since $\sigma(t) \in V$ for $t \in [0, T_m)$ it follows from (2.12) and the Poincaré inequality

$$(2.13) \quad \|v\|_2 \leq \|v\|_\infty \leq \|\partial_x v\|_2, \quad v \in V,$$

that there is a constant $C_1 > 0$ such that

$$(2.14) \quad \|\sigma(t)\|_{H^1} \leq C_1, \quad t \in [0, T_m).$$

In particular, $\|\sigma(t)\|_\infty \leq C_1$, from which we readily deduce that

$$(2.15) \quad s(t, x) \leq 1 - e^{-C_1}, \quad (t, x) \in [0, T_m) \times [0, 1].$$

Since $\partial_x s = e^{-\sigma} \partial_x \sigma$ and $\sigma \geq 0$ we further obtain from (2.14) and (2.15) that

$$(2.16) \quad \|s(t)\|_{H^1} \leq 1 + C_1, \quad t \in [0, T_m).$$

We next establish a similar estimate for ℓ and first proceed in a formal way. A rigorous proof based on an approximation argument will be sketched at the end of the proof of Theorem 1. We take $\psi = \partial_t \ell$ in (2.4) and use (2.3), (2.5), (2.12), (2.15), and the nonnegativity of ℓ to obtain

$$\begin{aligned} & \int_0^1 |\partial_t \ell|^2 dx + \frac{D}{2} \frac{d}{dt} \int_0^1 (1-s) |\partial_x \ell|^2 dx - \nu \partial_t \ell(0) \\ &= -\frac{D}{2} \int_0^1 \partial_t s |\partial_x \ell|^2 dx - \int_0^1 (\partial_t s + \varepsilon s) \partial_t \ell dx \\ &\leq \frac{D}{2} (\delta + \varepsilon) \int_0^1 s |\partial_x \ell|^2 dx + \frac{1}{2} \|\partial_t \ell\|_2^2 + \frac{1}{2} \|(1-s) \partial_t \sigma + \varepsilon s\|_2^2 \\ &\leq \frac{D(\delta + \varepsilon)e^{C_1}}{2} \int_0^1 (1-s) |\partial_x \ell|^2 dx + \frac{1}{2} \|\partial_t \ell\|_2^2 + \|\partial_t \sigma\|_2^2 + \varepsilon \\ &\leq \frac{D(\delta + \varepsilon)e^{C_1}}{2} \int_0^1 (1-s) |\partial_x \ell|^2 dx + \frac{1}{2} \|\partial_t \ell\|_2^2 + 2C_0 + \varepsilon. \end{aligned}$$

Consequently, there is a constant $C_2 > 0$ such that

$$\int_0^1 |\partial_t \ell|^2 dx + \frac{d}{dt} \left\{ \int_0^1 D (1-s) |\partial_x \ell|^2 dx - \nu \ell(0) \right\} \leq C_2 \left(1 + \int_0^1 (1-s) |\partial_x \ell|^2 dx \right),$$

whence

$$\begin{aligned} (2.17) \quad & D \int_0^1 (1-s(t)) |\partial_x \ell(t)|^2 dx - \nu \ell(t, 0) + \int_0^t \int_0^1 |\partial_t \ell(\tau)|^2 dx d\tau \\ & \leq (1 + D \|\ell_0\|_{H^1}^2) e^{C_2 t} \end{aligned}$$

for $t \in [0, T_m)$. Recalling (2.15) we readily deduce from (2.17) and the Poincaré inequality (2.13) that

$$(2.18) \quad \|\ell(t)\|_{H^1}^2 + \int_0^t \|\partial_t \ell(\tau)\|_2^2 d\tau \leq C_3 e^{C_3 t}$$

for $t \in [0, T_m)$. According to (2.16) and (2.18), we have shown so far that, given $T > 0$, there is a positive constant M_T depending only on $\ell_0, s_0, D, \delta, \varepsilon, \nu$, and T such that

$$\|\ell(t)\|_{H^1} + \|s(t)\|_{H^1} \leq M_T \quad \text{for } t \in [0, T] \cap [0, T_m).$$

It remains to improve this estimate to a uniform bound on the modulus of continuity with respect to time for ℓ and s in $H^1(0, 1)$ as required by (2.6). To this end, we first notice that (2.5) and (2.18) entail that

$$(2.19) \quad \|\partial_t s(t)\|_{H^1} \leq 2 \|\ell(t)\|_{H^1} \|(1-s)(t)\|_{H^1} + (\delta + \varepsilon) \|s(t)\|_{H^1} \leq C_4 e^{C_4 t}$$

for $t \in [0, T_m)$.

It next follows from (2.4) that $\xi := \partial_t \ell$ formally solves

$$(2.20) \quad \begin{aligned} \partial_t \xi &= D \partial_x ((1-s) \partial_x \xi - \partial_t s \partial_x \ell) \\ &+ (\ell + \delta) \partial_t s - \xi (1-s), \quad (t, x) \in (0, T_m) \times (0, 1), \end{aligned}$$

$$(2.21) \quad (1-s(t, 0)) \partial_x \xi(t, 0) - \partial_t s(t, 0) \partial_x \ell(t, 0) = \xi(t, 1) = 0, \quad t \in (0, T_m).$$

We multiply (2.20) by ξ and integrate over $(0, 1)$ to obtain, thanks to (2.21),

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\xi(t)\|_2^2 &= -D \int_0^1 \partial_x \xi \left((1-s) \partial_x \xi - \partial_t s \partial_x \ell \right) dx \\ &\quad + \int_0^1 (\ell + \delta) \partial_t s \xi dx - \int_0^1 (1-s) \xi^2 dx \\ &\leq -D \int_0^1 (1-s) |\partial_x \xi|^2 dx - \int_0^1 (1-s) \xi^2 dx \\ &\quad + D \|\partial_t s\|_\infty \|\partial_x \ell\|_2 \|\partial_x \xi\|_2 + (\|\ell\|_\infty + \delta) \|\partial_t s\|_2 \|\xi\|_2. \end{aligned}$$

Recalling (2.15), (2.18), (2.19), and the continuous embedding of $H^1(0, 1)$ in $L^\infty(0, 1)$, we deduce from the Young inequality that

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\xi(t)\|_2^2 &\leq -D e^{-C_1} \|\partial_x \xi\|_2^2 - e^{-C_1} \|\xi\|_2^2 + D e^{-C_1} \|\partial_x \xi\|_2^2 \\ &\quad + e^{-C_1} \|\xi\|_2^2 + C_5 \|\partial_t s\|_{H^1}^2 (\|\ell\|_{H^1} + \delta)^2 \\ &\leq C_6 e^{C_6 t}, \end{aligned}$$

whence

$$\|\xi(t)\|_2^2 \leq \|\xi(\tau)\|_2^2 + 2 e^{C_6 t}, \quad 0 \leq \tau \leq t < T_m.$$

Integrating with respect to τ over $(0, t)$ and using (2.18), we end up with

$$\begin{aligned} t \|\xi(t)\|_2^2 &\leq \int_0^t \|\xi(\tau)\|_2^2 d\tau + 2 t e^{C_6 t}, \\ (2.22) \quad \|\xi(t)\|_2^2 &\leq C_7 \left(1 + \frac{1}{t}\right) e^{C_7 t}, \quad t \in (0, T_m). \end{aligned}$$

We next infer from (1.3) and (2.16) that

$$\partial_x^2 \ell = \frac{\partial_t \ell}{D(1-s)} + \frac{\ell}{D} - \frac{\delta}{D} \frac{s}{1-s} + \partial_x \sigma \partial_x \ell.$$

Therefore, by (2.14), (2.15), (2.18), (2.22), and the interpolation inequality

$$\|\partial_x w\|_\infty \leq C_8 \|w\|_{H^2}^{3/4} \|w\|_2^{1/4}, \quad w \in H^2(0, 1),$$

we have

$$\begin{aligned} \|\partial_x^2 \ell(t)\|_2 &\leq \frac{e^{C_1}}{D} \|\partial_t \ell(t)\|_2 + \frac{\|\ell(t)\|_2}{D} + \frac{\delta e^{C_1}}{D} + \|\partial_x \sigma(t)\|_2 \|\partial_x \ell(t)\|_\infty \\ &\leq \frac{\sqrt{C_7} e^{C_1}}{D} \left(1 + \frac{1}{t}\right)^{1/2} e^{C_7 t/2} + \frac{\sqrt{C_3} e^{C_3 t/2}}{D} \\ &\quad + \frac{\delta e^{C_1}}{D} + C_1 C_8 \|\ell(t)\|_{H^2}^{3/4} \|\ell(t)\|_2^{1/4} \\ &\leq C_9 \left(1 + \frac{1}{\sqrt{t}}\right) e^{C_9 t} + \frac{1}{2} \|\ell(t)\|_{H^2} \\ &\leq C_9 \left(1 + \frac{1}{\sqrt{t}}\right) e^{C_9 t} + \frac{1}{2} \sqrt{C_3} e^{C_3 t/2} + \frac{1}{2} \|\partial_x^2 \ell(t)\|_2, \end{aligned}$$

whence

$$(2.23) \quad \|\ell(t)\|_{H^2} \leq C_{10} \left(1 + \frac{1}{\sqrt{t}}\right) e^{C_{10}t}, \quad t \in (0, T_m).$$

Consider next $t_1 \in (0, T_m)$ and $t_2 \in (t_1, T_m)$. It follows from (2.18), (2.23), and the interpolation inequality

$$\|\partial_x w\|_2 \leq C_{11} \|w\|_{H^2}^{1/2} \|w\|_2^{1/2}, \quad w \in H^2(0, 1),$$

that

$$\begin{aligned} \|\ell(t_2) - \ell(t_1)\|_{H^1}^2 &\leq C_{11}^2 \|\ell(t_2) - \ell(t_1)\|_{H^2} \|\ell(t_2) - \ell(t_1)\|_2 \\ &\leq C_{11}^2 (\|\ell(t_2)\|_{H^2} + \|\ell(t_1)\|_{H^2}) \int_{t_1}^{t_2} \|\partial_t \ell(\tau)\|_2 \, d\tau \\ &\leq 2 C_{10} C_{11}^2 \left(1 + \frac{1}{\sqrt{t_1}}\right) e^{C_{10}t_2} (t_2 - t_1)^{1/2} \left(\int_{t_1}^{t_2} \|\partial_t \ell(\tau)\|_2^2 \, d\tau\right)^{1/2} \\ &\leq C_{12} \left(1 + \frac{1}{\sqrt{t_1}}\right) e^{C_{12}t_2} (t_2 - t_1)^{1/2}, \end{aligned}$$

$$(2.24) \quad \|\ell(t_2) - \ell(t_1)\|_{H^1} \leq C_{13} \left(1 + t_1^{-1/4}\right) e^{C_{13}t_2} (t_2 - t_1)^{1/4}.$$

Since ℓ belongs to $\mathcal{BC}([0, T] \cap [0, T_m]; H^1(0, 1))$ for any $T > 0$ by (2.2) and (2.18), it readily follows from (2.24) that

$$(2.25) \quad \ell \in \mathcal{BUC}([0, T] \cap [0, T_m]; H^1(0, 1)) \quad \text{for any } T > 0.$$

Recalling (2.16) and (2.19), we also have

$$s \in \mathcal{BUC}([0, T] \cap [0, T_m]; H^1(0, 1)) \quad \text{for any } T > 0,$$

which, together with (2.15) and (2.25), allows us to conclude that $T_m = \infty$ by (2.6) and (2.7). The proof of Theorem 1 is then complete, provided we can justify the above computations.

To this end we note that s belongs to $\mathcal{C}^1([0, T_m]; H^1(0, 1))$ by (2.2). Setting $\alpha_1 := e^{-C_1}/2$, classical approximation arguments ensure that, for each $\alpha \in (0, \alpha_1)$ and $T \in (0, T_m)$, there is a function $s_{\alpha, T} \in \mathcal{C}^\infty([0, T] \times [0, 1])$ such that

$$(2.26) \quad \|s_{\alpha, T} - s\|_{\mathcal{C}^1([0, T]; H^1(0, 1))} + \|s_{\alpha, T} - s\|_{\mathcal{C}([0, T] \times [0, 1])} \leq \alpha.$$

Owing to (2.15) and (2.26) we have $1 - s_{\alpha, T} \geq \alpha_1 > 0$ and the initial-boundary value problem

$$\begin{aligned} \partial_t \ell_{\alpha, T} - D \partial_x ((1 - s_{\alpha, T}) \partial_x \ell_{\alpha, T}) &= \delta s_{\alpha, T} - (1 - s_{\alpha, T}) \ell_{\alpha, T}, \quad (t, x) \in (0, T) \times (0, 1), \\ D(1 - s_{\alpha, T}(t, 0)) \partial_x \ell_{\alpha, T}(t, 0) + \nu &= \ell_{\alpha, T}(t, 1) = 0, \quad t \in (0, T), \\ \ell_{\alpha, T}(0, x) &= \ell_0(x), \quad x \in (0, 1), \end{aligned}$$

has a unique nonnegative classical solution $\ell_{\alpha, T} \in \mathcal{C}([0, T] \times [0, 1]) \cap \mathcal{C}^\infty((0, T] \times [0, 1])$. On the one hand it is rather straightforward to check that

$$(2.27) \quad \sup_{t \in [0, T]} \{ \|\ell_{\alpha, T} - \ell\|_2^2 \} \leq \alpha \int_0^T (D^2 \|\partial_x \ell(t)\|_2^2 + \delta^2 + \|\ell(t)\|_2^2) \, dt.$$

On the other hand, the smoothness of $\ell_{\alpha,T}$ allows us to proceed as in the derivation of (2.17) and obtain that

$$(2.28) \quad \|\partial_x \ell_{\alpha,T}(t)\|_2^2 + \int_0^t \|\partial_t \ell_{\alpha,T}(\tau)\|_2^2 d\tau \leq C_{14} e^{C_{14}t}, \quad t \in [0, T].$$

Owing to (2.27) we may let $\alpha \rightarrow 0$ in (2.28) and deduce by weak convergence arguments that $\ell \in L^\infty(0, T; V)$ and $\partial_t \ell \in L^2((0, T) \times (0, 1))$ and satisfy an inequality similar to (2.18). Still proceeding as above, we next show that $\ell_{\alpha,T}$ enjoys the properties (2.22) and (2.23), and we use once more weak compactness arguments to provide a rigorous proof of the *BUC*-estimate (2.25). \square

3. Convergence to the steady state. As a first step towards the proof of Theorem 3, we derive the following information on s and $\sigma := \Sigma_I(s)$ from (1.8).

LEMMA 5. *There is a real number $\vartheta \in (0, 1)$ such that*

$$(3.1) \quad 0 \leq s(t, x) \leq 1 - \vartheta \quad \text{for all } (t, x) \in [0, \infty) \times [0, 1]$$

and

$$(3.2) \quad \sigma \in L^\infty(0, \infty; H^1(0, 1)), \quad \partial_t \sigma \in L^2(0, \infty; H^1(0, 1)).$$

Proof. By (1.8) and (2.1) we have

$$\mathcal{L}_0(s_\infty) + \frac{D(\delta + \varepsilon)}{2} \|\partial_x \sigma(t) - \partial_x \Sigma_I(s_\infty)\|_2^2 \leq \mathcal{L}(s(t), \partial_t \sigma(t)) \leq \mathcal{L}(s_0, \partial_t \sigma(0))$$

for all $t \geq 0$, from which we readily conclude that $\partial_x \sigma$ belongs to $L^\infty(0, \infty; L^2(0, 1))$. Since $\sigma(t, 1) = 0$ for every $t \geq 0$, the Poincaré inequality (2.13) entails that σ belongs to $L^\infty(0, \infty; H^1(0, 1))$, which gives the first assertion in (3.2). The space $H^1(0, 1)$ being continuously embedded in $L^\infty(0, 1)$, we further deduce that σ belongs to $L^\infty((0, \infty) \times (0, 1))$, whence (3.1).

Next, since $r + (\delta + \varepsilon)/r \geq 2(\delta + \varepsilon)^{1/2} \geq 2\delta^{1/2}$ for $r \geq 0$ and $1 - s \geq \vartheta > 0$ by (3.1), we have

$$\mathcal{D}(s, \partial_t \sigma) \geq 2\delta^{1/2} \|\partial_t \sigma\|_2^2 + \vartheta D \|\partial_x \partial_t \sigma\|_2^2,$$

and we infer from (1.8) and (2.1) that

$$\mathcal{L}_0(s_\infty) + \int_0^t \left(2\delta^{1/2} \|\partial_t \sigma\|_2^2 + \vartheta D \|\partial_x \partial_t \sigma\|_2^2 \right) d\tau \leq \mathcal{L}(s_0, \partial_t \sigma(0))$$

for all $t \geq 0$, which gives the second assertion in (3.2). \square

We next turn to ℓ and establish the following inequality.

LEMMA 6. *For $t \geq 0$ we have*

$$(3.3) \quad \mathcal{L}_1(\ell(t)) + \int_0^t \mathcal{D}_1(\ell, s) d\tau \leq \mathcal{L}_1(\ell_0) + C_{15} \int_0^t \left\| \partial_x \ln \left(\frac{\ell + \delta + \varepsilon}{\ell_\infty + \delta + \varepsilon} \right) \right\|_2 \|\partial_t \sigma\|_2 d\tau$$

with

$$\mathcal{L}_1(u) := \int_0^1 \left[(u + \delta + \varepsilon) \left(\ln \left(\frac{u + \delta + \varepsilon}{\ell_\infty + \delta + \varepsilon} \right) - 1 \right) + \ell_\infty + \delta + \varepsilon \right] dx \geq 0$$

and

$$\begin{aligned} \mathcal{D}_1(u, v) &:= D \int_0^1 (1 - v) (u + \delta + \varepsilon) \left| \partial_x \ln \left(\frac{u + \delta + \varepsilon}{\ell_\infty + \delta + \varepsilon} \right) \right|^2 dx \\ &\quad + \varepsilon(\delta + \varepsilon) \int_0^1 \frac{u - \ell_\infty}{(u + \delta + \varepsilon)(\ell_\infty + \delta + \varepsilon)} \ln \left(\frac{u + \delta + \varepsilon}{\ell_\infty + \delta + \varepsilon} \right) dx \geq 0. \end{aligned}$$

Proof. We put $\Lambda := \ln(\ell + \delta + \varepsilon)$ and $\Lambda_\infty := \ln(\ell_\infty + \delta + \varepsilon)$. By (1.4) we have

$$(3.4) \quad 1 - s = (\delta + \varepsilon + \partial_t s) e^{-\Lambda} \quad \text{and} \quad s = (\ell - \partial_t s) e^{-\Lambda}.$$

Let $\psi \in V$. Using (1.4) and (3.4) the weak formulation of (1.3) reads

$$\begin{aligned} \langle \partial_t \ell, \psi \rangle_{V',V} + D \int_0^1 (\delta + \varepsilon + \partial_t s) e^{-\Lambda} \partial_x \ell \partial_x \psi \, dx \\ &= \nu \psi(0) - \int_0^1 (\partial_t s + \varepsilon s) \psi \, dx \\ &= \nu \psi(0) - \int_0^1 (\partial_t s + \varepsilon (\ell - \partial_t s) e^{-\Lambda}) \psi \, dx \\ &= \nu \psi(0) - \int_0^1 (1 - \varepsilon e^{-\Lambda}) \partial_t s \psi \, dx - \varepsilon \int_0^1 (1 - (\delta + \varepsilon) e^{-\Lambda}) \psi \, dx. \end{aligned}$$

Similarly, as (ℓ_∞, s_∞) is also a weak solution to (1.3)–(1.5), we have

$$\begin{aligned} D \int_0^1 (\delta + \varepsilon) e^{-\Lambda_\infty} \partial_x \ell_\infty \partial_x \psi \, dx \\ &= \nu \psi(0) - \varepsilon \int_0^1 (1 - (\delta + \varepsilon) e^{-\Lambda_\infty}) \psi \, dx. \end{aligned}$$

Subtracting the above two identities gives

$$\begin{aligned} \langle \partial_t \ell, \psi \rangle_{V',V} + D \int_0^1 (\delta + \varepsilon + \partial_t s) \partial_x (\Lambda - \Lambda_\infty) \partial_x \psi \, dx \\ &\quad + \varepsilon(\delta + \varepsilon) \int_0^1 (e^{-\Lambda_\infty} - e^{-\Lambda}) \psi \, dx \\ (3.5) \quad &= -D \int_0^1 \partial_t s \partial_x \Lambda_\infty \partial_x \psi \, dx - \int_0^1 (1 - \varepsilon e^{-\Lambda}) \partial_t s \psi \, dx. \end{aligned}$$

Clearly, $\Lambda - \Lambda_\infty$ belongs to V and we may thus take $\psi = \Lambda - \Lambda_\infty$ in (3.5) to obtain

$$\begin{aligned} \langle \partial_t \ell, \Lambda - \Lambda_\infty \rangle_{V',V} + D \int_0^1 (\delta + \varepsilon + \partial_t s) |\partial_x (\Lambda - \Lambda_\infty)|^2 \, dx \\ &\quad + \varepsilon(\delta + \varepsilon) \int_0^1 (e^{-\Lambda_\infty} - e^{-\Lambda}) (\Lambda - \Lambda_\infty) \, dx \\ &= -D \int_0^1 \partial_t s \partial_x \Lambda_\infty \partial_x (\Lambda - \Lambda_\infty) \, dx - \int_0^1 (1 - \varepsilon e^{-\Lambda}) \partial_t s (\Lambda - \Lambda_\infty) \, dx. \end{aligned}$$

On the one hand, we have

$$\langle \partial_t \ell, \Lambda - \Lambda_\infty \rangle_{V',V} = \frac{d}{dt} \int_0^1 (\ell + \delta + \varepsilon) (\ln(\ell + \delta + \varepsilon) - 1 - \Lambda_\infty) \, dx = \frac{d}{dt} \mathcal{L}_1(\ell).$$

On the other hand, $\partial_x \Lambda_\infty$ clearly belongs to $L^\infty(0, 1)$ and

$$\begin{aligned} & -D \int_0^1 \partial_t s \partial_x \Lambda_\infty \partial_x (\Lambda - \Lambda_\infty) dx - \int_0^1 (1 - \varepsilon e^{-\Lambda}) \partial_t s (\Lambda - \Lambda_\infty) dx \\ & \leq D \|\partial_x \Lambda_\infty\|_\infty \|\partial_t s\|_2 \|\partial_x (\Lambda - \Lambda_\infty)\|_2 + \|\partial_t s\|_2 \|\Lambda - \Lambda_\infty\|_2 \\ & \leq C_{16} \|\partial_t s\|_2 \|\partial_x (\Lambda - \Lambda_\infty)\|_2, \end{aligned}$$

the last inequality being a consequence of the Poincaré inequality (2.13) which can be applied here since $\Lambda(t) - \Lambda_\infty \in V$ for all $t \geq 0$. Therefore,

$$\frac{d}{dt} \mathcal{L}_1(\ell) + \mathcal{D}_1(\ell, s) \leq C_{16} \|\partial_t s\|_2 \|\partial_x (\Lambda - \Lambda_\infty)\|_2.$$

Integrating the above inequality with respect to time and using the fact that $|\partial_t s| = |(1 - s) \partial_t \sigma| \leq |\partial_t \sigma|$ complete the proof of (3.3). \square

Several bounds on ℓ may be deduced from (3.3) and are listed now.

LEMMA 7. *Introducing the space $W := H^2(0, 1) \cap V$ and denoting its topological dual by W' , we have*

$$(3.6) \quad \ell \ln(1 + \ell) \in L^\infty(0, \infty; L^1(0, 1)),$$

$$(3.7) \quad (\ell + \delta + \varepsilon)^{1/2} \partial_x \ln \left(\frac{\ell + \delta + \varepsilon}{\ell_\infty + \delta + \varepsilon} \right) \in L^2((0, \infty) \times (0, 1)),$$

$$(3.8) \quad \frac{1}{\ell + \delta + \varepsilon} - \frac{1}{\ell_\infty + \delta + \varepsilon} \in L^2((0, \infty) \times (0, 1)),$$

$$(3.9) \quad \partial_t \ell \in L^2(0, \infty; W').$$

Proof. Keeping the notation $\Lambda := \ln(\ell + \delta + \varepsilon)$ and $\Lambda_\infty := \ln(\ell_\infty + \delta + \varepsilon)$, we infer from (3.1), the nonnegativity of ℓ , and the inequalities $|e^{-\Lambda_\infty} - e^{-\Lambda}| \leq |\Lambda - \Lambda_\infty|$ and $\ell + \delta + \varepsilon \geq (\ell + 2\delta + \varepsilon)/2$ that

$$\mathcal{D}_1(\ell, s) \geq \frac{\vartheta D}{2} \int_0^1 (\ell + 2\delta + \varepsilon) |\partial_x (\Lambda - \Lambda_\infty)|_2^2 dx + \varepsilon(\delta + \varepsilon) \|e^{-\Lambda_\infty} - e^{-\Lambda}\|_2^2.$$

It then follows from (3.3) and the Young inequality that

$$\begin{aligned} & \mathcal{L}_1(\ell(t)) + \frac{\vartheta D}{2} \int_0^t \int_0^1 (\ell + 2\delta + \varepsilon) |\partial_x (\Lambda - \Lambda_\infty)|_2^2 dx d\tau \\ & + \varepsilon(\delta + \varepsilon) \int_0^t \|e^{-\Lambda_\infty} - e^{-\Lambda}\|_2^2 d\tau \\ & \leq \mathcal{L}_1(\ell(t)) + \int_0^t \mathcal{D}_1(\ell, s) d\tau \\ & \leq \mathcal{L}_1(\ell_0) + \frac{\vartheta D \delta}{2} \int_0^t \|\partial_x (\Lambda - \Lambda_\infty)\|_2^2 d\tau + C_{17} \int_0^t \|\partial_t \sigma\|_2^2 d\tau, \end{aligned}$$

hence

$$\begin{aligned} & \mathcal{L}_1(\ell(t)) + \frac{\vartheta D}{2} \int_0^t \int_0^1 (\ell + \delta + \varepsilon) |\partial_x (\Lambda - \Lambda_\infty)|_2^2 dx d\tau \\ & + \varepsilon(\delta + \varepsilon) \int_0^t \|e^{-\Lambda_\infty} - e^{-\Lambda}\|_2^2 d\tau \\ & \leq C_{18} \left(1 + \int_0^t \|\partial_t \sigma\|_2^2 d\tau \right). \end{aligned}$$

Recalling that $\partial_t \sigma \in L^2((0, \infty) \times (0, 1))$ by Lemma 5, we readily deduce from the previous inequality, the Poincaré inequality (2.13), and the nonnegativity of \mathcal{L}_1 that $\mathcal{L}_1(\ell)$ belongs to $L^\infty(0, \infty)$ and that (3.7) and (3.8) hold true. In turn, the L^∞ -bound on $\mathcal{L}_1(\ell)$ gives (3.6).

It remains to check (3.9). For that purpose, consider $\psi \in W$. Since $\delta + \varepsilon + \partial_t s = (1 - s)(\ell + \delta + \varepsilon)$ by (1.4), it follows from (3.5) and the Hölder inequality that

$$\begin{aligned} |\langle \partial_t \ell, \psi \rangle_{V', V}| &\leq D \int_0^1 (1 - s)(\ell + \delta + \varepsilon) |\partial_x(\Lambda - \Lambda_\infty)| |\partial_x \psi| dx \\ &\quad + \varepsilon(\delta + \varepsilon) \|e^{-\Lambda_\infty} - e^{-\Lambda}\|_2 \|\psi\|_2 + \|\partial_t s\|_2 (D \|\partial_x \Lambda_\infty\|_\infty + 1) \|\psi\|_{H^1} \\ &\leq D \left(\int_0^1 (\ell + \delta + \varepsilon) |\partial_x(\Lambda - \Lambda_\infty)|^2 dx \right)^{1/2} \left(\int_0^1 (\ell + \delta + \varepsilon) dx \right)^{1/2} \|\partial_x \psi\|_\infty \\ &\quad + \varepsilon(\delta + \varepsilon) \|e^{-\Lambda_\infty} - e^{-\Lambda}\|_2 \|\psi\|_2 + C_{19} \|\partial_t \sigma\|_2 \|\psi\|_{H^1}. \end{aligned}$$

Thanks to (3.2), (3.6), (3.7), (3.8), and the continuous embedding of W in $W^{1, \infty}(0, 1)$, we deduce from the above inequality that

$$\int_0^t |\langle \partial_t \ell, \psi \rangle_{V', V}|^2 d\tau \leq C_{20} \|\psi\|_{H^2},$$

which implies the claim (3.9). \square

Proof of Theorem 3. Let $(t_n)_{n \geq 1}$ be a sequence of positive times such that $t_n \rightarrow \infty$ as $n \rightarrow \infty$. It follows from (3.1), (3.2), (3.6), and the Dunford–Pettis theorem that $(\ell(t_n))_{n \geq 1}$ is weakly compact in $L^1(0, 1)$ and $(s(t_n))_{n \geq 1}$ is weakly compact in $H^1(0, 1)$. There are thus $\ell_* \in L^1(0, 1)$, $s_* \in H^1(0, 1)$, and a subsequence of $(\ell(t_n), s(t_n))_{n \geq 1}$ (not relabeled) such that

$$(3.10) \quad \ell(t_n) \rightharpoonup \ell_* \quad \text{in } L^1(0, 1),$$

$$(3.11) \quad s(t_n) \rightharpoonup s_* \quad \text{in } H^1(0, 1).$$

Owing to the compactness of the embedding of $H^1(0, 1)$ in $\mathcal{C}([0, 1])$, we further obtain

$$(3.12) \quad s(t_n) \longrightarrow s_* \quad \text{in } \mathcal{C}([0, 1]).$$

We now aim at showing that $(\ell_*, s_*) = (\ell_\infty, s_\infty)$. For that purpose, we introduce as usual the functions $(\ell_n, s_n) \in \mathcal{C}([0, 1]; V \times V)$ defined by

$$(\ell_n(t, x), s_n(t, x)) := (\ell(t_n + t, x), s(t_n + t, x)) \quad \text{for } (t, x) \in [0, 1]^2 \quad \text{and } n \geq 1,$$

and claim that

$$(3.13) \quad \ell_n \longrightarrow \ell_* \quad \text{in } \mathcal{C}([0, 1]; L^1_w(0, 1)),$$

$$(3.14) \quad s_n \longrightarrow s_* \quad \text{in } \mathcal{C}([0, 1] \times [0, 1]),$$

$L^1_w(0, 1)$ denoting the usual space $L^1(0, 1)$ endowed with its weak topology. Indeed, on the one hand, the continuous embedding of $H^1(0, 1)$ in $L^\infty(0, 1)$ and the Hölder inequality entail that

$$\begin{aligned} \|s_n(t) - s_n(0)\|_\infty &\leq \int_{t_n}^{t_n+t} \|\partial_t s\|_\infty d\tau \leq \int_{t_n}^{t_n+t} \|\partial_t \sigma\|_\infty d\tau \\ &\leq \left(\int_{t_n}^{t_n+t} \|\partial_t \sigma\|_{H^1}^2 d\tau \right)^{1/2} \end{aligned}$$

for $t \in [0, 1]$, so that

$$\sup_{t \in [0,1]} \|s_n(t) - s_n(0)\|_\infty \leq \left(\int_{t_n}^\infty \|\partial_t \sigma\|_{H^1}^2 d\tau \right)^{1/2}.$$

By (3.2), the right-hand side of the above inequality converges to zero as $n \rightarrow \infty$, which, together with (3.12), gives the claim (3.14). On the other hand, if $\psi \in W = H^2(0, 1) \cap V$ (defined in Lemma 7) and $t \in [0, 1]$,

$$\begin{aligned} \left| \int_0^1 (\ell_n(t) - \ell_n(0)) \psi dx \right| &= \left| \int_{t_n}^{t_n+1} \langle \partial_t \ell(\tau), \psi \rangle_{V',V} d\tau \right| \\ &\leq \int_{t_n}^{t_n+1} \|\partial_t \ell(\tau)\|_{W'} \|\psi\|_W d\tau \\ &\leq \|\psi\|_W \left(\int_{t_n}^{t_n+1} \|\partial_t \ell(\tau)\|_{W'}^2 d\tau \right)^{1/2} \\ &\leq \|\psi\|_W \left(\int_{t_n}^\infty \|\partial_t \ell(\tau)\|_{W'}^2 d\tau \right)^{1/2}. \end{aligned}$$

Owing to (3.9) and (3.10), we pass to the limit as $n \rightarrow \infty$ in the above inequality to conclude that

$$\lim_{n \rightarrow \infty} \sup_{t \in [0,1]} \left| \int_0^1 (\ell_n(t) - \ell_*) \psi dx \right| = 0 \quad \text{for all } \psi \in W.$$

A classical approximation combined with (3.6) allows us to extend the above convergence to every $\psi \in L^\infty(0, 1)$ and completes the proof of (3.13).

We next infer from (3.8) that

$$\begin{aligned} &\lim_{n \rightarrow \infty} \int_0^1 \left\| \frac{1}{\ell_n(t) + \delta + \varepsilon} - \frac{1}{\ell_\infty + \delta + \varepsilon} \right\|_2^2 dt \\ &= \lim_{n \rightarrow \infty} \int_{t_n}^{t_n+1} \left\| \frac{1}{\ell(t) + \delta + \varepsilon} - \frac{1}{\ell_\infty + \delta + \varepsilon} \right\|_2^2 dt \\ &= 0. \end{aligned}$$

Consequently, after extracting a further subsequence if necessary, we may assume that

$$(3.15) \quad \ell_n \longrightarrow \ell_\infty \quad \text{a.e. in } (0, 1) \times (0, 1).$$

Recalling (3.13) we readily conclude from (3.15) that $\ell_* = \ell_\infty$.

Next, (3.1) and (3.14) imply that $0 \leq s_*(x) \leq 1 - \vartheta < 1$ for $x \in [0, 1]$, and we infer from (1.4), (3.2), (3.13), and (3.14) that, if $\psi \in L^\infty(0, 1)$,

$$\begin{aligned} \left| \int_0^1 \left(\ell_* - (\delta + \varepsilon) \frac{s_*}{1 - s_*} \right) \psi dx \right| &= \left| \lim_{n \rightarrow \infty} \int_0^1 \int_0^1 \left(\ell_n(t) - (\delta + \varepsilon) \frac{s_n(t)}{1 - s_n(t)} \right) \psi(x) dx dt \right| \\ &= \left| \lim_{n \rightarrow \infty} \int_0^1 \partial_t \sigma(t_n + t) \psi dx dt \right| \\ &\leq \|\psi\|_\infty \left(\lim_{n \rightarrow \infty} \int_{t_n}^{t_n+1} \|\partial_t \sigma(t)\|_2^2 dt \right)^{1/2} \\ &= 0. \end{aligned}$$

Therefore, $(\delta + \varepsilon)s_*/(1 - s_*) = \ell_* = \ell_\infty$ a.e. in $(0, 1)$ and thus $s_* = s_\infty$.

We have actually shown that (ℓ_∞, s_∞) is the unique cluster point of $(\ell(t), s(t))$ as $t \rightarrow \infty$ for the weak topology of $L^1(0, 1) \times H^1(0, 1)$. The trajectory $\{(\ell(t), s(t)); t \geq 0\}$ being relatively compact for that topology by (3.2), (3.6), and the Dunford–Pettis theorem, we conclude that

$$(\ell(t), s(t)) \rightharpoonup (\ell_\infty, s_\infty) \text{ in } L^1(0, 1) \times H^1(0, 1)$$

as $t \rightarrow \infty$. The embedding of $H^1(0, 1)$ in $L^2(0, 1)$ being compact, we further obtain that

$$(3.16) \quad s(t) \longrightarrow s_\infty \text{ in } L^2(0, 1).$$

Let us finally improve the topology in which the convergence to the steady state takes place. For $\psi \in V$ it follows from the weak formulation of (1.3) and Proposition 2 that

$$\begin{aligned} & \langle \partial_t \ell, \psi \rangle_{V', V} + D \int_0^1 (1 - s) \partial_x (\ell - \ell_\infty) \partial_x \psi \, dx \\ &= D \int_0^1 (s - s_\infty) \partial_x \ell_\infty \partial_x \psi \, dx + \int_0^1 [(\ell_\infty + \delta) (s - s_\infty) - (1 - s) (\ell - \ell_\infty)] \psi \, dx. \end{aligned}$$

Since $\ell - \ell_\infty$ belongs to V , we may take $\psi = \ell - \ell_\infty$ and use the fact that ℓ_∞ does not depend on time to obtain

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\ell - \ell_\infty\|_2^2 + \int_0^1 (1 - s) \left(D |\partial_x (\ell - \ell_\infty)|^2 + |\ell - \ell_\infty|^2 \right) \, dx \\ & \leq \|s - s_\infty\|_2 \left(\|\partial_x \ell_\infty\|_\infty \|\partial_x (\ell - \ell_\infty)\|_2 + (\delta + \|\ell_\infty\|_\infty) \|\ell - \ell_\infty\|_2 \right) \\ & \leq \frac{D\vartheta}{2} \|\partial_x (\ell - \ell_\infty)\|_2^2 + \frac{\vartheta}{2} \|\ell - \ell_\infty\|_2^2 + C_{21} \|s - s_\infty\|_2^2, \end{aligned}$$

the last inequality following from the boundedness of ℓ_∞ and $\partial_x \ell_\infty$ and the parameter ϑ being defined in Lemma 5. Using (3.1) to bound from below $(1 - s)$ in the second term on the left-hand side of the above inequality we get

$$(3.17) \quad \frac{d}{dt} \|\ell - \ell_\infty\|_2^2 + D\vartheta \|\partial_x (\ell - \ell_\infty)\|_2^2 + \vartheta \|\ell - \ell_\infty\|_2^2 \leq 2C_{21} \|s - s_\infty\|_2^2.$$

We next infer from (1.4) and Proposition 2 that

$$\partial_t s = (1 - s) (\ell - \ell_\infty) - (\ell_\infty + \delta + \varepsilon) (s - s_\infty).$$

We differentiate the above identity with respect to x and take the scalar product in $L^2(0, 1)$ of the result with $\partial_x (s - s_\infty)$. Recalling the Poincaré inequality (2.13) and the bound (3.2) we obtain

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\partial_x (s - s_\infty)\|_2^2 + \int_0^1 (\ell_\infty + \delta + \varepsilon) |\partial_x (s - s_\infty)|_2^2 \, dx \\ & \leq \|\partial_x s\|_2 \|\ell - \ell_\infty\|_\infty \|\partial_x (s - s_\infty)\|_2 + \|\partial_x (\ell - \ell_\infty)\|_2 \|\partial_x (s - s_\infty)\|_2 \\ & \quad + \|\partial_x \ell_\infty\|_\infty \|s - s_\infty\|_2 \|\partial_x (s - s_\infty)\|_2 \\ & \leq \|(1 - s) \partial_x \sigma\|_2 \|\partial_x (\ell - \ell_\infty)\|_2 \|\partial_x (s - s_\infty)\|_2 + \|\partial_x (\ell - \ell_\infty)\|_2 \|\partial_x (s - s_\infty)\|_2 \\ & \quad + \|\partial_x \ell_\infty\|_\infty \|s - s_\infty\|_2 \|\partial_x (s - s_\infty)\|_2 \\ & \leq C_{22} \left(\|\partial_x (\ell - \ell_\infty)\|_2 + \|s - s_\infty\|_2 \right) \|\partial_x (s - s_\infty)\|_2 \\ & \leq \frac{\delta}{2} \|\partial_x (s - s_\infty)\|_2^2 + \frac{C_{23}}{2} \left(\|\partial_x (\ell - \ell_\infty)\|_2^2 + \|s - s_\infty\|_2^2 \right). \end{aligned}$$

Since ℓ_∞ and ε are nonnegative we end up with

$$(3.18) \quad \frac{d}{dt} \|\partial_x (s - s_\infty)\|_2^2 + \delta \|\partial_x (s - s_\infty)\|_2^2 \leq C_{23} \left(\|\partial_x (\ell - \ell_\infty)\|_2^2 + \|s - s_\infty\|_2^2 \right).$$

We now multiply (3.18) by $D\vartheta/C_{23}$ and add the result to (3.17) to obtain

$$\begin{aligned} \frac{d}{dt} \left(\|\ell - \ell_\infty\|_2^2 + \frac{D\vartheta}{C_{23}} \|\partial_x (s - s_\infty)\|_2^2 \right) + \vartheta \|\ell - \ell_\infty\|_2^2 \\ + \frac{D\vartheta\delta}{C_{23}} \|\partial_x (s - s_\infty)\|_2^2 \leq C_{24} \|s - s_\infty\|_2^2. \end{aligned}$$

Introducing $\omega := \min \{\vartheta, \delta\} > 0$ and integrating the above differential inequality give

$$\begin{aligned} \left(\|\ell(t) - \ell_\infty\|_2^2 + \frac{D\vartheta}{C_{23}} \|\partial_x (s(t) - s_\infty)\|_2^2 \right) \leq \left(\|\ell_0 - \ell_\infty\|_2^2 + \frac{D\vartheta}{C_{23}} \|\partial_x (s_0 - s_\infty)\|_2^2 \right) e^{-\omega t} \\ + C_{24} \int_0^t \|s(\tau) - s_\infty\|_2^2 e^{\omega(\tau-t)} d\tau \end{aligned}$$

for $t \geq 0$. Since $\|s(\tau) - s_\infty\|_2 \rightarrow 0$ as $\tau \rightarrow \infty$ by (3.16), we deduce from the above inequality that both $\|\ell(t) - \ell_\infty\|_2$ and $\|\partial_x (s(t) - s_\infty)\|_2$ converge to zero as $t \rightarrow \infty$, and the proof of Theorem 3 is complete. \square

4. Well-posedness of (1.1)–(1.2). In this section, we consider the initial-boundary value problem

$$(4.1) \quad \partial_t \ell = D \partial_x^2 \ell + \delta s - \ell (1 - s), \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(4.2) \quad \partial_t s = \ell (1 - s) - (\delta + \varepsilon) s, \quad (t, x) \in (0, \infty) \times (0, 1),$$

$$(4.3) \quad D \partial_x \ell(t, 0) + \nu = \ell(t, 1) = 0, \quad t \in (0, \infty),$$

$$(4.4) \quad (\ell, s)(0, x) = (\ell_0, s_0)(x), \quad x \in (0, 1),$$

with $D > 0$, $\delta > 0$, $\varepsilon \geq 0$, and $\nu > 0$.

THEOREM 8. *Given initial data (ℓ_0, s_0) satisfying (1.7) the initial-boundary value problem (4.1)–(4.4) possesses a unique weak solution*

$$(\ell, s) \in \mathcal{C}([0, \infty); V \times V) \cap \mathcal{C}^1([0, \infty); V' \times V)$$

such that $\partial_t \ell \in L^2((0, T) \times (0, 1))$ for every $T > 0$,

$$(4.5) \quad 0 \leq \ell(t, x) \quad \text{and} \quad 0 \leq s(t, x) < 1 \quad \text{for} \quad (t, x) \in [0, \infty) \times [0, 1],$$

and

$$\begin{aligned} \langle \partial_t \ell(t), \psi \rangle_{V', V} + D \int_0^1 \partial_x \ell(t) \partial_x \psi \, dx \\ = \nu \psi(0) + \int_0^1 (\delta s(t) - \ell(t)(1 - s(t))) \psi \, dx, \end{aligned}$$

$$\partial_t s(t) = -(\delta + \varepsilon)s(t) + \ell(t)(1 - s(t))$$

for $\psi \in V$ and $t \in [0, \infty)$. Furthermore, for each $t > 0$,

$$(4.6) \quad \Lambda(\ell(t), s(t)) + \int_0^t \mathcal{D}_\Lambda(\ell(\tau), s(\tau)) \, d\tau = \Lambda(\ell_0, s_0),$$

where

$$\begin{aligned} \Lambda(u, v) &:= \frac{1}{2} \|u - \ell_\infty\|_2^2 + \Lambda_0(u, v), \\ \Lambda_0(u, v) &:= \int_0^1 (1 - s_\infty) (\ell_\infty + \delta + 2\varepsilon) \left[\Sigma_I(v) - \Sigma_I(s_\infty) - \frac{v - s_\infty}{1 - s_\infty} \right] dx \geq 0, \\ \mathcal{D}_\Lambda(u, v) &:= D \|\partial_x(u - \ell_\infty)\|_2^2 + \int_0^1 \frac{|\partial_t v|^2 + \varepsilon (\ell_\infty + \delta + \varepsilon) (v - s_\infty)^2}{1 - v} dx, \end{aligned}$$

$s_\infty := \ell_\infty / (\ell_\infty + \delta + \varepsilon)$, and ℓ_∞ denotes the unique solution to the boundary-value problem

$$(4.7) \quad -D \partial_x^2 \ell_\infty + \frac{\varepsilon \ell_\infty}{\ell_\infty + \delta + \varepsilon} = 0 \quad \text{in } (0, 1), \quad D \partial_x \ell_\infty(0) + \nu = \ell_\infty(1) = 0.$$

We first point out that the nonnegativity of Λ_0 is a consequence of the convexity of Σ_I . We next observe that Λ differs strongly from \mathcal{L} and provide valuable information on the large time behavior. In fact, exponential convergence to the steady state can be shown; see Proposition 9 below.

Proof. The local well-posedness of (4.1)–(4.4) in $V \times V$ and the bounds (4.5) for (ℓ, s) on the maximal existence time interval $[0, T_m)$ are established as in the proof of Theorem 1 by employing the abstract theory developed in [1]. In addition, since $u \mapsto \varepsilon u / (u + \delta + \varepsilon)$ is a nondecreasing function, the existence and uniqueness of the solution ℓ_∞ to (4.7) are proved by classical arguments [2, 4].

We next turn to the proof of the identity (4.6) for $t \in [0, T_m)$. For that purpose, we take $\psi = \ell - \ell_\infty$ in the weak formulation of (4.1) and (4.7) for ℓ and ℓ_∞ and use (4.2) for s and s_∞ to obtain

$$(4.8) \quad \frac{1}{2} \frac{d}{dt} \|\ell - \ell_\infty\|_2^2 + D \|\partial_x(\ell - \ell_\infty)\|_2^2 = - \int_0^1 (\partial_t s + \varepsilon (s - s_\infty)) (\ell - \ell_\infty) dx.$$

Owing to (4.2) we have

$$\ell - \ell_\infty = \frac{\partial_t s}{1 - s} + (\delta + \varepsilon) \frac{s - s_\infty}{(1 - s)(1 - s_\infty)}.$$

Consequently,

$$\begin{aligned} & - \int_0^1 (\partial_t s + \varepsilon (s - s_\infty)) (\ell - \ell_\infty) dx \\ &= - \int_0^1 \frac{|\partial_t s|^2}{1 - s} dx - \int_0^1 \left[\varepsilon + \frac{\delta + \varepsilon}{1 - s_\infty} \right] (s - s_\infty) \frac{\partial_t s}{1 - s} dx \\ & \quad - \varepsilon (\delta + \varepsilon) \int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx \\ &= - \int_0^1 \frac{|\partial_t s|^2}{1 - s} dx - \int_0^1 (\ell_\infty + \delta + 2\varepsilon) \partial_t ((1 - s_\infty) \Sigma_I(s) - s) dx \\ & \quad - \varepsilon \int_0^1 (\ell_\infty + \delta + \varepsilon) \frac{(s - s_\infty)^2}{1 - s} dx, \end{aligned}$$

where we have used $(\delta + \varepsilon) / (1 - s_\infty) = \ell_\infty + \delta + \varepsilon$ to deduce the last equality. Inserting the previous equality in (4.8) gives (4.6) for $t \in [0, T_m)$ after integration with respect to time.

A useful consequence of (4.6) is that

$$(4.9) \quad \|\ell(t)\|_2 \leq \|\ell_\infty\|_2 + \sqrt{2 \Lambda(\ell_0, s_0)}, \quad t \in [0, T_m).$$

We next proceed as in the proof of (2.18) and take $\psi = \partial_t \ell$ in the weak formulation of (4.1) (recall that a regularization procedure has to be used to justify this step as $\partial_t \ell \notin V$). This gives

$$\|\partial_t \ell\|_2^2 + \frac{D}{2} \frac{d}{dt} \|\partial_x \ell\|_2^2 = \int_0^1 \partial_t \ell (\delta s - \ell (1 - s)) dx \leq \frac{1}{2} \|\partial_t \ell\|_2^2 + (\delta + \|\ell\|_2^2).$$

We then infer from (4.9) that

$$(4.10) \quad \int_0^t \|\partial_t \ell(\tau)\|_2^2 d\tau + D \|\partial_x \ell(t)\|_2^2 \leq C_{25} (1 + t), \quad t \in [0, T_m).$$

We next differentiate (4.2) to obtain the equation solved by $\partial_x s$ and deduce from (4.10) that

$$(4.11) \quad \|\partial_x s(t)\|_2 \leq C_{26} (1 + t), \quad t \in [0, T_m).$$

Finally, we infer from (4.2), (4.5), (4.9), (4.10), (4.11), and the continuous embedding of $H^1(0, 1)$ in $L^\infty(0, 1)$ that

$$s \in W^{1,\infty}([0, T] \cap [0, T_m]; H^1(0, 1)) \quad \text{and} \quad \delta s - \ell (1 - s) \in L^2((0, T) \cap (0, T_m); H^1(0, 1))$$

for all $T > 0$. Combining this last property with (4.1) and [10, Theorem 4.3.1] ensures that

$$\ell \in C^{1/2}((0, T] \times (0, T_m); H^1(0, 1)) \quad \text{for all } T > 0.$$

Collecting the above information allows us to conclude that $T_m = \infty$ and complete the proof. \square

We next turn to the large time behavior and prove the following result.

PROPOSITION 9. *Consider initial data (ℓ_0, s_0) satisfying (1.7) and denote by (ℓ, s) the corresponding weak solution to (4.1)–(4.4) given by Theorem 8. Then, for $t \geq 0$,*

$$(4.12) \quad \|\ell(t) - \ell_\infty\|_2^2 + (\delta + \varepsilon) \|s(t) - s_\infty\|_2^2 \leq 2 \Lambda(\ell_0, s_0) e^{-\chi t}$$

with

$$(4.13) \quad \chi := \min \left\{ D, \frac{D (\delta + \varepsilon)}{2(D + 2)} + \frac{\varepsilon}{2} \right\}.$$

Proposition 9 is a simple consequence of (4.6) and the following two functional inequalities.

LEMMA 10. *Under the assumptions and notation of Proposition 9 we have for all $t \geq 0$*

$$(4.14) \quad \mathcal{D}_\Lambda(\ell(t), s(t)) \geq \chi \Lambda(\ell(t), s(t)),$$

$$(4.15) \quad \Lambda_0(s(t)) \geq \frac{\delta + \varepsilon}{2} \|s(t) - s_\infty\|_2^2.$$

Proof. By the elementary inequality $a^2 \leq ((D + 2)/D) (a - b)^2 + ((D + 2)/2) b^2$ we have

$$\begin{aligned}
 & \int_0^1 \left(\frac{\delta + \varepsilon}{1 - s_\infty} \right)^2 \frac{(s - s_\infty)^2}{1 - s} dx \\
 & \leq \frac{D + 2}{D} \int_0^1 (1 - s) \left[\frac{\delta + \varepsilon}{1 - s_\infty} \frac{s - s_\infty}{1 - s} - (\ell - \ell_\infty) \right]^2 dx \\
 & \quad + \frac{D + 2}{2} \int_0^1 (1 - s) (\ell - \ell_\infty)^2 dx \\
 (4.16) \quad & \leq \frac{D + 2}{D} \left\{ \int_0^1 (1 - s) \left(\frac{\partial_t s}{1 - s} \right)^2 dx + \frac{D}{2} \|\ell - \ell_\infty\|_2^2 \right\},
 \end{aligned}$$

the last inequality following from the identities $\ell + \delta + \varepsilon = (\partial_t s + \delta + \varepsilon)/(1 - s)$ and

$$(4.17) \quad \frac{\delta + \varepsilon}{1 - s_\infty} = \ell_\infty + \delta + \varepsilon \geq \delta + \varepsilon.$$

We next infer from the Poincaré inequality (2.13), (4.16), and (4.17) that

$$\begin{aligned}
 \mathcal{D}_\Lambda(\ell, s) & \geq D \|\ell - \ell_\infty\|_2^2 + \int_0^1 \frac{|\partial_t s|^2}{1 - s} dx + \varepsilon (\delta + \varepsilon) \int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx \\
 & \geq \frac{D}{2} \|\ell - \ell_\infty\|_2^2 + \frac{D}{D + 2} \int_0^1 \left(\frac{\delta + \varepsilon}{1 - s_\infty} \right)^2 \frac{(s - s_\infty)^2}{1 - s} dx \\
 & \quad + \varepsilon (\delta + \varepsilon) \int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx, \\
 (4.18) \quad \mathcal{D}_\Lambda(\ell, s) & \geq \frac{D}{2} \|\ell - \ell_\infty\|_2^2 + (\delta + \varepsilon) \left(\frac{D(\delta + \varepsilon)}{D + 2} + \varepsilon \right) \int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx.
 \end{aligned}$$

We next claim that

$$(4.19) \quad \frac{(z - a)^2}{2} \leq \Sigma_I(z) - \Sigma_I(a) - \frac{z - a}{1 - a} \leq \frac{1}{2} \frac{1 + a}{1 - a} \frac{(z - a)^2}{1 - z}, \quad (a, z) \in [0, 1) \times [0, 1).$$

Indeed, the first inequality in (4.19) follows from the convexity of Σ_I , while the second can be shown by studying the variation of the function

$$z \mapsto (1 - z) \left[\Sigma_I(z) - \Sigma_I(a) - \frac{z - a}{1 - a} \right] - \frac{1}{2} \frac{1 + a}{1 - a} (z - a)^2,$$

the real number a being fixed in $[0, 1)$. Using the second inequality in (4.19) with $z = s(t, x)$ and $a = s_\infty(x)$ gives

$$\int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx \geq \int_0^1 \frac{2(\delta + \varepsilon)}{1 + s_\infty} \left(\Sigma_I(s) - \Sigma_I(s_\infty) - \frac{s - s_\infty}{1 - s_\infty} \right) dx.$$

Using once more (4.17) we realize that

$$\begin{aligned}
 \frac{2(\delta + \varepsilon)}{1 + s_\infty} & = \frac{2(\delta + \varepsilon)}{1 + s_\infty} \frac{(1 - s_\infty) (\ell_\infty + \delta + 2\varepsilon)}{\delta + \varepsilon + \varepsilon(1 - s_\infty)} \\
 & \geq (\delta + \varepsilon) \frac{(1 - s_\infty) (\ell_\infty + \delta + 2\varepsilon)}{\delta + 2\varepsilon} \\
 & \geq \frac{(1 - s_\infty) (\ell_\infty + \delta + 2\varepsilon)}{2},
 \end{aligned}$$

so that

$$\int_0^1 \frac{(s - s_\infty)^2}{(1 - s)(1 - s_\infty)} dx \geq \frac{1}{2} \Lambda_0(s).$$

Inserting the previous inequality in (4.18) gives

$$\mathcal{D}_\Lambda(\ell, s) \geq \frac{D}{2} \|\ell - \ell_\infty\|_2^2 + \frac{1}{2} \left(\frac{D(\delta + \varepsilon)}{D + 2} + \varepsilon \right) \Lambda_0(s) \geq \chi \Lambda(\ell, s),$$

and hence (4.14).

Inequality (4.15) readily follows from (4.17) and the first inequality in (4.19). \square

5. Numerical simulations. Above, we have developed a theory which provides insight into the stationary states and large time behavior of two models of morphogen transport. It is therefore interesting to verify the dependence on the parameters of both the properties of these stationary solutions and the convergence rate to the equilibrium. For this, we have performed numerical simulations on an x86-64 machine under MATLAB 7.2. All source files are available upon request.

We discuss certain qualitative properties of (1.3)–(1.6) and (4.1)–(4.4) for a range of their parameters with default values $D = \nu = \delta = 1$ and $\varepsilon = 0$. For an evaluation of the two models for biologically relevant parameter values, see [7]. In all experiments, we supply zero initial data, $\ell_0(x) = s_0(x) = 0$. For clarity, in what follows we shall refer to (1.3)–(1.6) as the nonlinear diffusion (ND) model and to (4.1)–(4.4) as the linear diffusion (LD) model.

Discretization. To discretize the system (1.3)–(1.6) we employ a finite difference scheme in the spatial variable on a regular mesh with $N = 240$ nodes. Introducing $\ell_i(t) \approx \ell(t, ih)$, $s_i(t) \approx s(t, ih)$, where h is the mesh size and using the first order flux approximation,

$$j_i = \left(1 - \frac{s_{i+1} + s_i}{2} \right) \frac{\ell_{i+1} - \ell_i}{h},$$

we arrive at a system of $2N$ ordinary differential equations (ODEs)

$$(5.1a) \quad \frac{d\ell_i}{dt} = \frac{j_i - j_{i-1}}{h} + \delta s_i - \ell_i (1 - s_i),$$

$$(5.1b) \quad \frac{ds_i}{dt} = -(\delta + \varepsilon)s_i + \ell_i (1 - s_i)$$

for $i = 1, \dots, N$. This discretization has the truncation error of second order with respect to h and is supplemented by a first order approximation of the flux boundary condition. In order to solve (5.1) we use a stiff ODE solver `ode15s` from MATLAB. Since we trace the behavior of the solution up to very small time derivatives, we adjust the absolute and relative error tolerance parameters in `ode15s` to the value of 10^{-13} . Model LD is discretized and solved analogously.

Steady state computation. As observed in [7], without degradation ($\varepsilon = 0$) one has explicit formulas for the stationary states, since from (1.13) it follows that for the ND model there holds

$$(5.2) \quad \ell_\infty(x) = \delta \left(\exp\left(\frac{\nu}{\delta D}(1 - x)\right) - 1 \right),$$

$$(5.3) \quad s_\infty(x) = 1 - \exp\left(-\frac{\nu}{\delta D}(1 - x)\right),$$

while for the LD model, according to Theorem 8,

$$(5.4) \quad \ell_\infty(x) = \frac{\nu}{D}(1 - x),$$

$$(5.5) \quad s_\infty(x) = 1 - \frac{1}{1 + (\nu/(\delta D))(1 - x)}.$$

But when $\varepsilon > 0$ there is no closed formula and we have to approximate the steady states numerically. To this end, we either run the simulation until the L^2 norm of the right-hand side of the discrete ODE system drops (unless otherwise stated) below 10^{-6} or we solve the nonlinear system (5.1), where we replace time derivatives with zeros. Note, however, that although in practice the first procedure turns out to be more robust, none of these procedures guarantees that the obtained solution is sufficiently close to the stationary state, particularly when s_∞ is too close to 1; see [7] for a more detailed discussion.

5.1. Dependence of the convergence rate on the parameters. According to Theorem 3, the solutions of the ND model converge to the stationary state. From Proposition 9 we know that for the LD model the distance between the solution and the equilibrium decays at least at an exponential rate. It may therefore be interesting to identify the actual convergence rate of the two models under consideration for several values of the parameters. In order to do so, we investigate the L^2 norm of the relative difference between the stationary solution (ℓ_∞, s_∞) and the solution at time t ,

$$d_{rel}(t) = \frac{\|(\ell, s)(t) - (\ell_\infty, s_\infty)\|_2}{\|(\ell_\infty, s_\infty)\|_2}.$$

We stop the analysis when d_{rel} drops below 10^{-5} ; otherwise the results would have been polluted with the discretization error.

We consider four cases, when only one of the parameters ε , ν , D , or δ is allowed to vary, while the others are kept at their default values.

The convergence to the stationary state seems to occur indeed at an exponential rate; see the top plots in Figures 1–4. Therefore, we can next compare the dependence of the rates on various parameters of the models. To this end, we introduce an approximate measure of the convergence rate ζ as the mean value of the ratios

$$\zeta_i = -\frac{\ln(d_{rel}(t_{i+1})) - \ln(d_{rel}(t_i))}{t_{i+1} - t_i}$$

for several subsequent t_i at which d_{rel} was measured.

As could be expected by comparing the diffusive parts, the LD model always converges faster than the ND model for the same parameter set. From the modeling point of view, as the constants D and δ have different interpretation in each of the two models [7], it is worth mentioning that for moderate values of D , an approximately tenfold increase of D in the nonlinear diffusion model is needed to recover the convergence speed of the linear diffusion model (see Figure 3). Although in the absence of degradation the stationary states of the two models depend only on the ratio $\nu/(\delta D)$, the convergence rate may be different for the same ratios of $\nu/(\delta D)$, as indicated in Figures 2–4.

Both models show similar dependence of the convergence speed ζ on the degradation rate ε , the unbinding parameter δ , and the diffusion coefficient D . However,

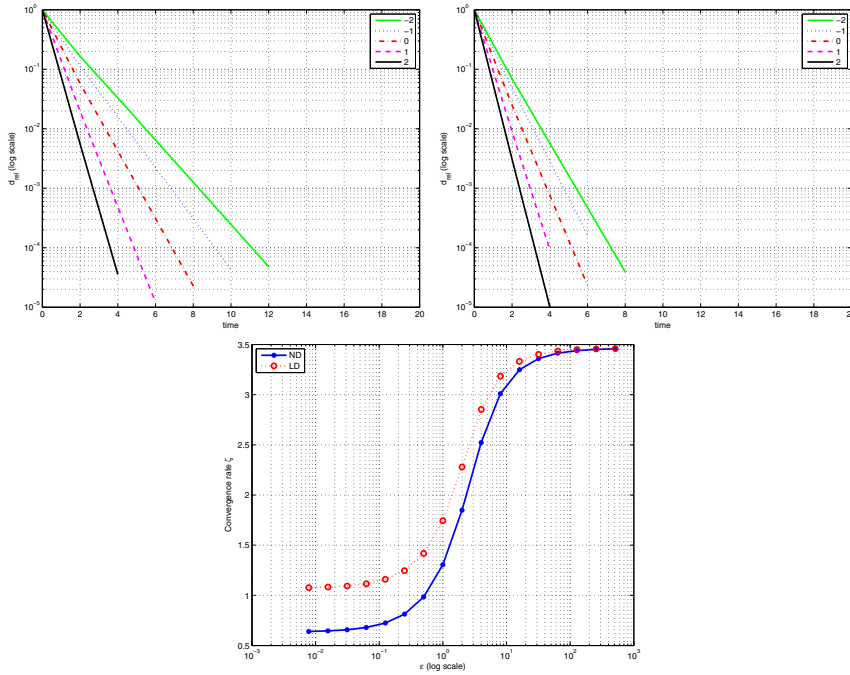


FIG. 1. Convergence to the equilibrium for different values of ε . The ND model (top left) vs. the LD model (top right) for $\varepsilon = 2^k$, $k = -2, \dots, 2$. The comparison of the convergence rates ζ for a broader range of ε (bottom).

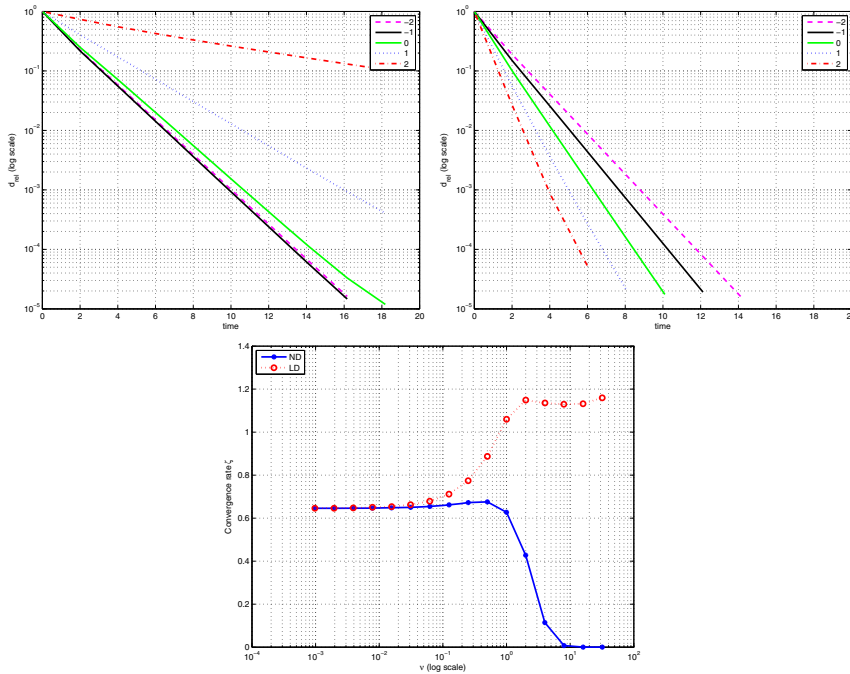


FIG. 2. Convergence to the equilibrium for different values of ν . The ND model (top left) vs. the LD model (top right) for $\nu = 2^k$, $k = -2, \dots, 2$. The comparison of the convergence rates ζ for a broader range of ν (bottom).

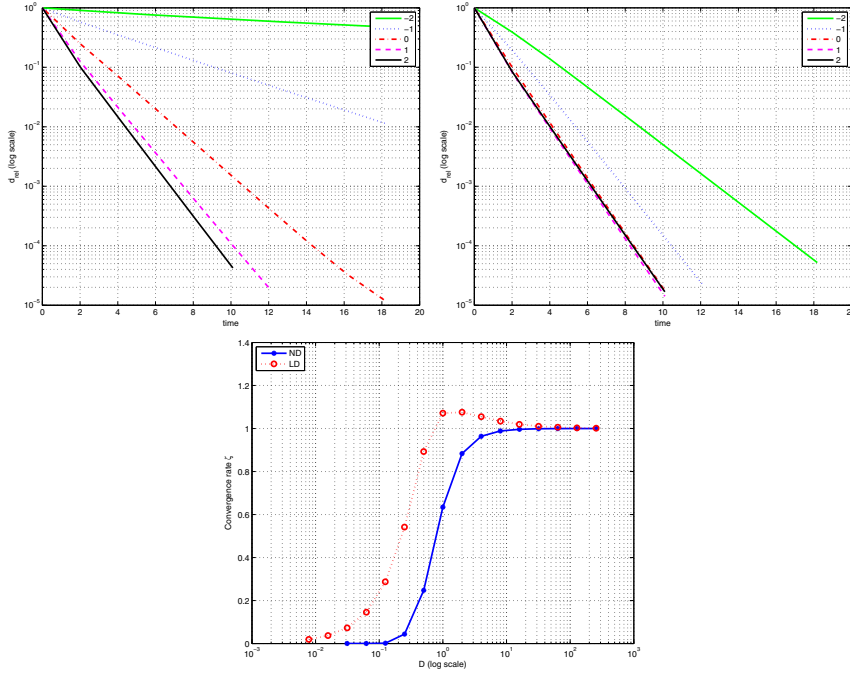


FIG. 3. Convergence to the equilibrium for different values of D . The ND model (top left) vs. the LD model (top right) for $D = 2^k$, $k = -2, \dots, 2$. The comparison of the convergence rates ζ for a broader range of D (bottom).

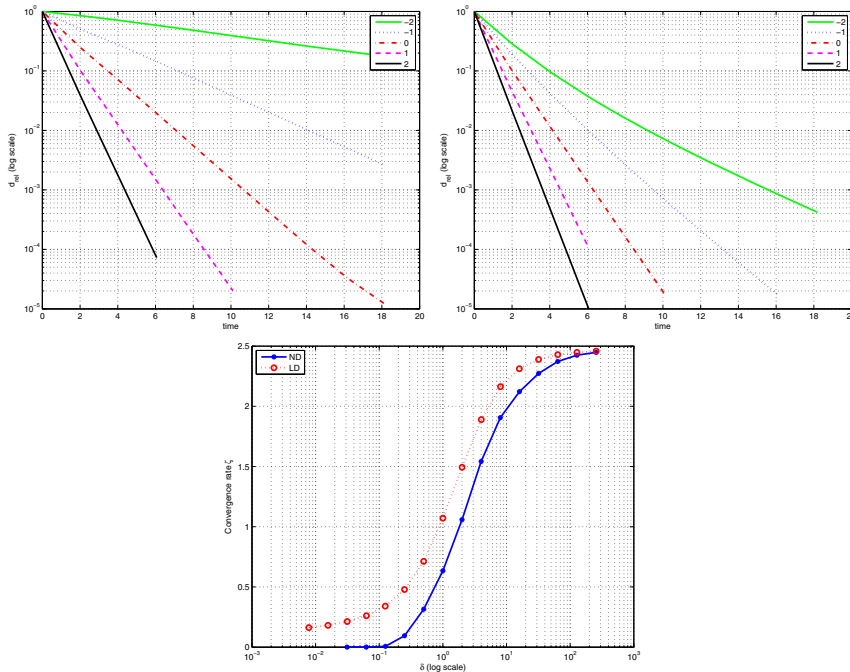


FIG. 4. Convergence to the equilibrium for different values of δ . The ND model (top left) vs. the LD model (top right) for $\delta = 2^k$, $k = -2, \dots, 2$. The comparison of the convergence rates ζ for a broader range of δ (bottom).

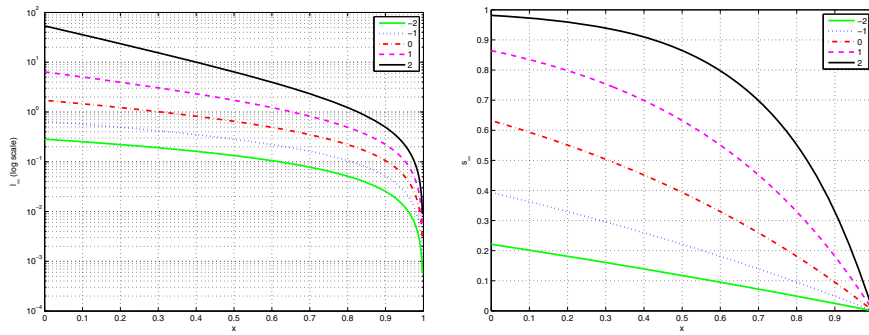


FIG. 5. Stationary states of the ND model, $\nu = 2^k$, $k = -2, -1, 0, 1, 2$, without degradation.

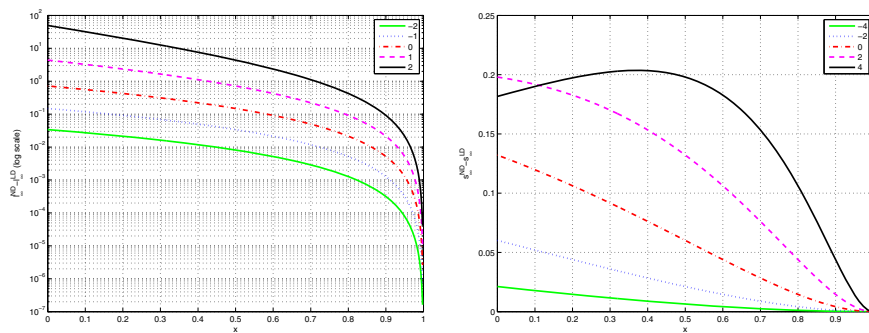


FIG. 6. Difference between the stationary states of the ND and LD models, $\nu = 2^k$, $k = -2, -1, 0, 1, 2$, without degradation.

the two models differ significantly with respect to their sensitivity to the injection parameter ν : the convergence speed increases with ν for the LD model, while the effect of ν on the convergence of the ND model is opposite; cf. Figure 2 (bottom). This is explained by looking at the stationary states of the latter model; see Figures 5 and 6. The stationary solution reaches, in the case of nonlinear diffusion, higher values than in the linear diffusion case. Indeed, from (5.2) and (5.4) it follows that $\|\ell_\infty\|_\infty = \delta (\exp(\nu/(\delta D)) - 1)$ for the ND model which clearly exceeds by far the L^∞ norm ν/D of ℓ_∞ of the LD model. The closer s_∞ gets to 1, the smaller the diffusion, and thus it takes more time for the solution to evolve to the steady state.

It is important to notice that due to the nonlinear effects, both models cannot converge to (ℓ_∞, s_∞) arbitrarily fast, as can be concluded from the bottom parts of Figures 1–4. For each of varying parameters, there is some threshold value of ζ , which cannot be surpassed regardless of the value of the parameter.

Let us observe that for sufficiently large values of D , ε , δ , or $1/\nu$, the convergence speed of both models stabilizes at almost an identical level. On one hand, the similarity in the convergence rates can easily be explained by noticing that the s_∞ -component of the steady state of the nonlinear model is so close to zero that the fluxes for both models are nearly the same. On the other hand, the convergence rate stabilization effect is in full agreement with our previous discussion of the Liapunov function for the LD model in the case of varying D . Indeed, denoting by (ℓ^D, s^D) the solution to the LD model for some prescribed D , it follows from (4.6) and (4.7) that

$\ell^D \rightarrow 0$ in $L^2(0, T; H^1(0, 1))$ as $D \rightarrow \infty$. Hence we deduce that s^D converges towards S , solving $\partial_t S = -(\delta + \varepsilon)S$, which exponentially decays to zero at rate $\delta + \varepsilon$ (cf. Figure 3, where $\delta + \varepsilon = 1$), to be compared with the L^2 norm decay rate $\chi/2 = \delta/4 + \varepsilon/2$ obtained in the limit $D \rightarrow \infty$ from (4.13).

5.2. Dependence of the stationary solutions on the parameters. As confirmed by Figures 5–6 and experiments for other parameters not reported here (see also [7]), for identical set of parameters, the two models produce stationary solutions of quite similar shape. There is a difference in how the ν parameter affects the stationary states of the two models; cf. Figures 5–6. The ℓ_∞ -component of the stationary solution of the nonlinear diffusion model shows more sensitivity to large values of ν than its linear diffusion counterpart. This is because for large ν , the x -derivative of ℓ is much larger in the nonlinear diffusion case, due to higher saturation level of s . The stationary states of the two models are much less sensitive to the changes in the degradation rate ε . There is only a limited range of ε , close to 1, for which large relative change in ε results in a similar large change in the stationary solutions.

Acknowledgments. Part of this research was carried out during P. Krzyżanowski's and D. Wrzosek's visits to the Institut de Mathématiques de Toulouse and during Ph. Laurençot's visit to the Institute of Applied Mathematics at Warsaw University. The authors thank both institutions for their hospitality. They also thank an anonymous referee for pointing out a gap in the original proof of Theorem 1.

REFERENCES

- [1] H. AMANN, *Highly degenerate quasilinear parabolic systems*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 18 (1991), pp. 135–166.
- [2] PH. BÉNILAN, M. G. CRANDALL, AND P. SACKS, *Some L^1 existence and dependence results for semilinear elliptic equations under nonlinear boundary conditions*, Appl. Math. Optim., 17 (1988), pp. 203–224.
- [3] T. BOLLENBACH, K. KRUSE, P. PANTAZIS, M. GONZÁLEZ-GAITÁN, AND F. JÜLICHER, *Robust formation of morphogen gradients*, Phys. Rev. Lett., 94 (2005), article 018103.
- [4] H. BREZIS AND W. A. STRAUSS, *Semi-linear second-order elliptic equations in L^1* , J. Math. Soc. Japan, 25 (1973), pp. 565–590.
- [5] E. V. ENTCHEV AND M. A. GONZÁLEZ-GAITÁN, *Morphogen gradient formation and vesicular trafficking*, Traffic, 3 (2002), pp. 98–109.
- [6] M. KERSZBERG AND L. WOLPERT, *Mechanisms for positional signalling by morphogen transport: A theoretical study*, J. Theoret. Biol., 191 (1998), pp. 103–114.
- [7] P. KRZYŻANOWSKI, PH. LAURENÇOT, AND D. WRZOSEK, *Mathematical models of receptor-mediated transport of morphogens*, submitted.
- [8] A. D. LANDER, Q. NIE, AND F. Y. M. WAN, *Do morphogen gradients arise by diffusion?*, Dev. Cell, 2 (2002), pp. 785–796.
- [9] J. H. MERKIN, D. J. NEEDHAM, AND B. D. SLEEMAN, *A mathematical model for the spread of morphogens with density dependent chemosensitivity*, Nonlinearity, 18 (2005), pp. 2745–2773.
- [10] A. PAZY, *Semigroups of Linear Operators and Applications to Partial Differential Equations*, Appl. Math. Sci. 44, Springer-Verlag, New York, 1983.
- [11] J.-P. VINCENT AND L. DUBOIS, *Morphogen transport along epithelia, an integrated trafficking problem*, Dev. Cell, 3 (2002), pp. 615–623.

ON THE WHITHAM EQUATIONS FOR THE DEFOCUSING COMPLEX MODIFIED KDV EQUATION*

YUJI KODAMA[†], V. U. PIERCE[†], AND FEI-RAN TIAN[†]

Abstract. We study the Whitham equations for the defocusing complex modified KdV (mKdV) equation. These Whitham equations are quasi-linear hyperbolic equations and describe the averaged dynamics of the rapid oscillations which appear in the solution of the mKdV equation when the dispersive parameter is small. The oscillations are referred to as dispersive shocks. The Whitham equations for the mKdV equation are neither strictly hyperbolic nor genuinely nonlinear. We are interested in the solutions of the Whitham equations when the initial values are given by a step-like function. We also compare the results with those of the defocusing nonlinear Schrödinger (NLS) equation. For the NLS equation, the Whitham equations are strictly hyperbolic and genuinely nonlinear. We show that the weak hyperbolicity of the mKdV–Whitham equations is responsible for some new structure in the dispersive shocks which has not been found in the NLS case.

Key words. Whitham equations, non-strictly hyperbolic equations, dispersive shocks

AMS subject classifications. 35L65, 35L67, 35Q05, 35Q15, 35Q53, 35Q58

DOI. 10.1137/070705131

1. Introduction. In [11, 12], Pierce and Tian studied the self-similar solutions of the Whitham equations which describe the zero dispersion limits of the KdV hierarchy. The main feature of the Whitham equations for the higher members of the hierarchy, of which the KdV equation is the first member, is that these Whitham equations are neither strictly hyperbolic nor genuinely nonlinear. This is in sharp contrast to the case of the KdV equation whose Whitham equations are strictly hyperbolic and genuinely nonlinear [8]. In this paper, we extend their studies to the case of the complex modified KdV (mKdV) equation, which is the second member of the defocusing nonlinear Schrödinger (NLS) hierarchy. The Whitham equations for the defocusing NLS equation are strictly hyperbolic and genuinely nonlinear, and they have been studied extensively (see, for example, [4, 6, 7, 10, 14]). However, for the mKdV equation, the Whitham equations are neither strictly hyperbolic nor genuinely nonlinear.

Let us begin with a brief description of the zero dispersion limit of the solution of the NLS equation

$$(1.1) \quad \sqrt{-1} \epsilon \frac{\partial \psi}{\partial t} + 2\epsilon^2 \frac{\partial^2 \psi}{\partial x^2} - 4|\psi|^2 \psi = 0,$$

with the initial data

$$\psi(x, 0) = A_0(x) \exp\left(\sqrt{-1} \frac{S_0(x)}{\epsilon}\right).$$

Here $A_0(x)$ and $S_0(x)$ are real functions that are independent of ϵ . Writing the solution $\psi(x, t; \epsilon) = A(x, t; \epsilon) \exp(\sqrt{-1} \frac{S(x, t; \epsilon)}{\epsilon})$ and using the notation $\rho(x, t; \epsilon) = A^2(x, t; \epsilon)$,

*Received by the editors October 11, 2007; accepted for publication (in revised form) July 8, 2008; published electronically December 17, 2008.

<http://www.siam.org/journals/sima/40-5/70513.html>

[†]Department of Mathematics, Ohio State University, 231 W. 18th Avenue, Columbus, OH 43210 (kodama@math.ohio-state.edu, vpierce@math.ohio-state.edu, tian@math.ohio-state.edu). The first and third authors were supported in part by NSF grant DMS-0404931. The second author was supported in part by NSF grant DMS-0135308.

$v(x, t; \epsilon) = \partial S(x, t; \epsilon) / \partial x$, one obtains the conservation form of the defocusing NLS equation,

$$(1.2) \quad \begin{cases} \frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x}(4\rho v) = 0, \\ \frac{\partial}{\partial t}(\rho v) + \frac{\partial}{\partial x}(4\rho v^2 + 2\rho^2) = \epsilon^2 \frac{\partial}{\partial x} \left(\rho \frac{\partial^2}{\partial x^2} \ln \rho \right). \end{cases}$$

The mass density $\rho = |\psi|^2$ and momentum density $\rho v = \frac{\sqrt{-1}}{2}(\psi\psi_x^* - \psi^*\psi_x)$ have weak limits as $\epsilon \rightarrow 0$ [6]. These limits satisfy a 2×2 system of hyperbolic equations

$$(1.3) \quad \begin{cases} \frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x}(4\rho v) = 0, \\ \frac{\partial}{\partial t}(\rho v) + \frac{\partial}{\partial x}(4\rho v^2 + 2\rho^2) = 0 \end{cases}$$

until its solution develops a shock. System (1.3) can be rewritten in the diagonal form for $\rho \neq 0$,

$$(1.4) \quad \frac{\partial}{\partial t} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + 2 \begin{pmatrix} 3\alpha + \beta & 0 \\ 0 & \alpha + 3\beta \end{pmatrix} \frac{\partial}{\partial x} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = 0,$$

where the Riemann invariants α and β are given by

$$(1.5) \quad \alpha = \frac{v}{2} + \sqrt{\rho}, \quad \beta = \frac{v}{2} - \sqrt{\rho}.$$

As a simple example, we consider the case with $\alpha = a = \text{constant}$. System (1.4) reduces to a single equation

$$(1.6) \quad \frac{\partial \beta}{\partial t} + 2(a + 3\beta) \frac{\partial \beta}{\partial x} = 0.$$

The solution is given by the implicit form

$$\beta(x, t) = h(x - 2(a + 3\beta)t),$$

where $h(x) = \beta(x, 0)$ is the initial function for β . One can easily see that if $\beta(x, 0)$ decreases in some region, then $\beta(x, t)$ develops a shock in a finite time; i.e., $\partial \beta / \partial x$ becomes singular.

After the shock formation in the solution of (1.3) or (1.4), the weak limits are described by the NLS–Whitham equations, which can also be put in the Riemann invariant form [4, 6, 7, 10]

$$(1.7) \quad \frac{\partial u_i}{\partial t} + \lambda_{g,i}(u_1, \dots, u_{2g+2}) \frac{\partial u_i}{\partial x} = 0, \quad i = 1, 2, \dots, 2g + 2,$$

where $\lambda_{g,i}$ are expressed in terms of complete hyperelliptic integrals of genus g [8]. Here the number g is exactly the number of phases in the NLS oscillations with small dispersion. Accordingly, the zero phase $g = 0$ corresponds to no oscillations, and single and higher phases $g \geq 1$ correspond to the NLS oscillations. System (1.4) is viewed as the zero phase Whitham equations. The solution of the Whitham equations (1.7) for $g \geq 1$ then describes the averaged motion of the oscillations appearing in the solution of (1.1) (see, e.g., [7]).

Let us discuss the most important $g = 1$ case in more detail. We note that it is well known that the KdV oscillatory solution, in the single phase regime, can be approximately described by the KdV periodic solution when the dispersive parameter is small [1, 5, 16]. It is very possible to use the method of [1, 16] to show that the solution of the NLS equation (1.1) for small ϵ can be approximately described, in the single phase regime, by the periodic solution of the NLS equation. The NLS periodic solution has the form

$$(1.8) \quad \tilde{\rho}(x, t; \epsilon) = \rho_3 + (\rho_2 - \rho_3) \operatorname{sn}^2 \left(\sqrt{\rho_1 - \rho_3} \frac{\theta(x, t)}{\epsilon}, s \right)$$

with $\theta(x, t) = x - V_1 t - Q$, where $V_1 = V_1(\rho_1, \rho_2, \rho_3)$ and Q is the phase shift. Here ρ_i 's are determined by the equation obtained from (1.2)

$$\frac{\epsilon^2}{4} \left(\frac{d\rho}{d\theta} \right)^2 = (\rho - \rho_1)(\rho - \rho_2)(\rho - \rho_3)$$

with $\rho_1 > \rho_2 > \rho_3$, and $\operatorname{sn}(z, s)$ is the Jacobi elliptic function with the modulus $s = (\rho_2 - \rho_3)/(\rho_1 - \rho_3)$. We can also write ρ_i 's as [3]

$$(1.9) \quad \left\{ \begin{array}{l} \rho_1 = \frac{1}{4}(u_1 + u_2 - u_3 - u_4)^2, \\ \rho_2 = \frac{1}{4}(u_1 - u_2 + u_3 - u_4)^2, \\ \rho_3 = \frac{1}{4}(u_1 - u_2 - u_3 + u_4)^2 \end{array} \right.$$

with $u_1 > u_2 > u_3 > u_4$. The velocity V_1 is then given by

$$V_1 = 2(u_1 + u_2 + u_3 + u_4).$$

It can also be shown that the flow velocity of the periodic solution (1.8) is given by [3]

$$(1.10) \quad \tilde{v}(x, t; \epsilon) = \frac{V_1}{4} + \frac{\sqrt{\rho_1 \rho_2}}{2\rho} (u_1 - u_2 - u_3 + u_4).$$

For constants u_1, u_2, u_3 , and u_4 , formula (1.8) gives the well-known elliptic solution of the NLS equation. To describe the solution $\rho(x, t; \epsilon)$ of the NLS equation (1.2), the quantities u_1, u_2, u_3 , and u_4 , are instead functions of x and t and evolve according to the single phase Whitham equations (1.7) for $g = 1$. The phase shift Q is also a function of u_1, u_2, u_3, u_4 , and depends on the initial values of system (1.2). Suppose the initial function for system (1.4) is given by

$$(1.11) \quad \alpha(x, 0) = a, \quad \beta(x, 0) = \beta_0(x),$$

where a is a constant and $\beta_0(x)$ is a monotone function. The Whitham solution of (1.7) has the property that $u_1 = a$ and that only u_2, u_3, u_4 are nontrivial functions of (x, t) . Then the phase shift $Q(u_1, u_2, u_3, u_4)$ with $u_1 = a$ is the unique solution of the boundary value problem for the Euler–Poisson–Darboux equations

$$(1.12) \quad 2(u_i - u_j) \frac{\partial^2 Q}{\partial u_i \partial u_j} = \frac{\partial Q}{\partial u_i} - \frac{\partial Q}{\partial u_j}, \quad i, j = 2, 3, 4,$$

$$(1.13) \quad Q(a, u, u, u) = f(u),$$

where boundary value $f(u)$ is the inverse of the initial function $\beta_0(x)$ (see Appendix C for details). If, instead, $\beta(x, 0)$ is a constant and $\alpha(x, 0)$ is given by a monotone function, the phase shift Q is determined by a similar boundary value problem for the Euler–Poisson–Darboux equations.

The weak limit of $\rho(x, t; \epsilon)$ of the NLS equation (1.1) as $\epsilon \rightarrow 0$ can be expressed in terms of ρ_1, ρ_2, ρ_3 , and ρ_4 [6] as

$$(1.14) \quad \overline{\rho(x, t)} = \rho_1 - (\rho_1 - \rho_3) \frac{E(s)}{K(s)},$$

where $K(s)$ and $E(s)$ are the complete elliptic integrals of the first and second kind, respectively. This weak limit can also be viewed as the average value of the periodic solution $\tilde{\rho}(x, t; \epsilon)$ of (1.8) over its period $L = 2\epsilon K(s)/\sqrt{\rho_1 - \rho_3}$.

In order to see how a single phase Whitham solution appears, we consider the following step-like initial data for system (1.4):

$$(1.15) \quad \alpha(x, 0) = a, \quad \beta(x, 0) = \begin{cases} b, & x < 0, \\ c, & x > 0, \end{cases}$$

where $a > b, a > c, b \neq c$. The solution of (1.4) develops a shock if and only if $b > c$ (cf. (1.6)). After the formation of a shock, the Whitham equations (1.7) with $g = 1$ kick in. For instance, we consider the Whitham equations with the initial data [7]

$$(1.16) \quad u_1(x, 0) = a, \quad u_2(x, 0) = b, \quad u_3(x, 0) = \begin{cases} b, & x < 0, \\ c, & x > 0, \end{cases} \quad u_4(x, 0) = c.$$

Now notice that the Whitham equations (1.7) for $g = 1$ with the initial data (1.16) can be reduced to a single equation $u_{3t} + \lambda_{1,3}(a, b, u_3, c)u_{3x} = 0$. The equation has a global self-similar solution, which is implicitly given by $x/t = \lambda_3(a, b, u_3, c)$. The x - t plane is then divided into *three* parts:

$$(1) \quad \frac{x}{t} < \gamma, \quad (2) \quad \gamma < \frac{x}{t} < 2a + 4b + 2c, \quad (3) \quad \frac{x}{t} > 2a + 4b + 2c,$$

where $\gamma = 2(a + b + 2c) - 8(a - c)(b - c)/(a + b - 2c)$ (see (2.7) and (2.8) for the derivation). The solution of system (1.4) occupies the first and third parts; i.e.,

(1) for $x/t < \gamma$,

$$\alpha(x, t) = a, \quad \beta(x, t) = b,$$

(3) for $x/t > 2a + 4b + 2c$,

$$\alpha(x, t) = a, \quad \beta(x, t) = c.$$

The Whitham solution of (1.7) with $g = 1$ lives in the second part; i.e.,

(2) for $\gamma < x/t < 2a + 4b + 2c$,

$$u_1(x, t) = a, \quad u_2(x, t) = b, \quad \frac{x}{t} = \lambda_{1,3}(a, b, u_3, c), \quad u_4(x, t) = c,$$

where the solution u_3 can be obtained as a function of the self-similarity variable x/t if

$$\frac{\partial \lambda_{1,3}}{\partial u_3}(a, b, u_3, c) \neq 0.$$

Indeed, it has been shown that the Whitham equations (1.7) are genuinely nonlinear [6, 7]; i.e.,

$$(1.17) \quad \frac{\partial \lambda_{1,i}}{\partial u_i}(u_1, u_2, u_3, u_4) > 0, \quad i = 1, 2, 3, 4,$$

for $u_1 > u_2 > u_3 > u_4$.

In Figure 1.1, we plot the self-similar solution of the Whitham equations (1.7) with $g = 1$ for the NLS equation and the corresponding periodic oscillatory solution (1.8) for the initial data (1.15) with $a = 4$, $b = 1$, and $c = -1$. The oscillations describe a dispersive shock of the NLS equation under a small dispersion. Note here that the oscillations have uniform structure, which is due to an almost linear profile of the Whitham solution u_3 . This will be seen to be in sharp contrast to the case of the mKdV equation, which we will discuss later (cf. Figure 1.2).

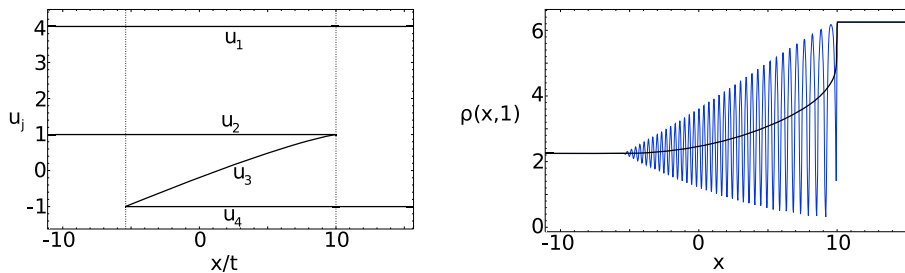


FIG. 1.1. Self-similar solution of the NLS–Whitham equation (1.7) with $g = 1$ and the corresponding oscillatory solution (1.8) of the NLS equation with $\epsilon = 0.3$ at $t = 1$. The dark line in the middle of the oscillations is the weak limit $\bar{\rho}(x, t)$ given by (1.14) at $t = 1$. The initial data are given by (1.15) with $a = 4$, $b = 1$, and $c = -1$.

Most of the figures in this paper have the same form: On the left-hand side is a plot of the solution of the Whitham equations as a function of the self-similarity variable x/t , which is exact, other than a numerical method used to implement the inverse function theorem. On the right-hand side is the oscillatory solution given by (1.8) at $t = 1$ (respectively, (1.25) for mKdV), while the dark plot is the weak limit (1.14) of the oscillatory solution; both plots on the right are also exact. For all the step-like initial data that we study in this paper, the resulting phase shift Q , which is determined by (1.12) and (1.13), is always zero for both NLS and mKdV. This is due to the fact that our step-like initial function has a jump discontinuity at the origin, which implies that $f \equiv 0$.

In the plots on the left-hand sides of Figures 1.1 and 1.2, we demarcate the region where the single phase Whitham equations govern the solution and label the four functions $u_1 > u_2 > u_3 > u_4$. The demarcation and labeling are similar in the other figures, and we will omit them for brevity.

The defocusing NLS equation is just the first member of the defocusing NLS hierarchy; the second is the (defocusing) complex mKdV equation

$$(1.18) \quad \frac{\partial \psi}{\partial t} + \frac{3}{2} |\psi|^2 \frac{\partial \psi}{\partial x} - \frac{\epsilon^2}{4} \frac{\partial^3 \psi}{\partial x^3} = 0.$$

We again use $\psi(x, t; \epsilon) = A(x, t; \epsilon) \exp(\sqrt{-1} \frac{S(x, t; \epsilon)}{\epsilon})$ and notation $\rho(x, t; \epsilon) = A^2(x, t; \epsilon)$, $v(x, t; \epsilon) = \partial S(x, t; \epsilon) / \partial x$ to obtain the conservation form of the mKdV

equation

$$(1.19) \quad \begin{cases} \frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x} \left(\frac{3}{4} \rho^2 + \frac{3}{4} \rho v^2 \right) = \epsilon^2 \frac{\partial}{\partial x} \left(\rho^{3/4} \frac{\partial^2}{\partial x^2} \rho^{1/4} \right), \\ \frac{\partial}{\partial t} (\rho v) + \frac{\partial}{\partial x} \left(\frac{3}{2} \rho^2 v + \frac{3}{4} \rho v^3 \right) = \frac{\epsilon^2}{4} \frac{\partial}{\partial x} \left[\frac{\partial^2}{\partial x^2} (\rho v) - \frac{3}{2} R \right], \end{cases}$$

where

$$R = \frac{3v}{2\rho} \left(\frac{\partial \rho}{\partial x} \right)^2 + \frac{\partial v}{\partial x} \frac{\partial \rho}{\partial x} - v \frac{\partial^2 \rho}{\partial x^2}.$$

The mass density ρ and momentum density ρv for the mKdV equation also have weak limits as $\epsilon \rightarrow 0$ [6]. As in the NLS case, the weak limits satisfy

$$(1.20) \quad \begin{cases} \frac{\partial}{\partial t} \rho + \frac{\partial}{\partial x} \left(\frac{3}{4} \rho^2 + \frac{3}{4} \rho v^2 \right) = 0, \\ \frac{\partial}{\partial t} (\rho v) + \frac{\partial}{\partial x} \left(\frac{3}{2} \rho^2 v + \frac{3}{4} \rho v^3 \right) = 0, \end{cases}$$

until the solution of (1.20) forms a shock. One can rewrite equations (1.20) as

$$(1.21) \quad \frac{\partial}{\partial t} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \frac{3}{8} \begin{pmatrix} 5\alpha^2 + 2\alpha\beta + \beta^2 & 0 \\ 0 & \alpha^2 + 2\alpha\beta + 5\beta^2 \end{pmatrix} \frac{\partial}{\partial x} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = 0,$$

where the Riemann invariants α and β are again given by formula (1.5).

Let us again consider the simplest case $\alpha(x, 0) = a$, where a is constant, to see how the solution of system (1.21) develops a shock. In this case, the system reduces to a single equation

$$\frac{\partial \beta}{\partial t} + \frac{3}{8} (a^2 + 2a\beta + 5\beta^2) \frac{\partial \beta}{\partial x} = 0.$$

As in the NLS case, we consider the initial function given by $\beta(x, 0) = b$ for $x < 0$ and $\beta(x, 0) = c$ for $x > 0$. We recall that, in the NLS case, the zero phase solution of (1.4) develops a shock if and only if $b > c$. However, the solution in the mKdV case develops a shock for $b > c$ if and only if $a + 5b > 0$. In addition, if $b < c$, the solution in the mKdV case develops a shock if and only if $a + 5b < 0$. These differences between the mKdV and NLS cases are due to the weak hyperbolicity of the system (1.21) (note that, for the eigenspeed $\lambda = \frac{3}{8}(\alpha^2 + 2\alpha\beta + 5\beta^2)$ for β , we have $\partial\lambda/\partial\beta = \frac{3}{4}(\alpha + 5\beta)$ which can change sign). As will be shown below, this leads to new structure in the dispersive shocks for the mKdV case.

As in the case of the NLS equation, immediately after the shock formation in the solution of (1.20), the weak limits are described by the mKdV–Whitham equations

$$(1.22) \quad \frac{\partial u_i}{\partial t} + \mu_{g,i}(u_1, \dots, u_{2g+2}) \frac{\partial u_i}{\partial x} = 0, \quad i = 1, 2, \dots, 2g + 2,$$

where $\mu_{g,i}$ can also be expressed in terms of complete hyperelliptic integrals of genus g [6].

In this paper, we study the solution of the Whitham equations (1.22) with $g = 1$ when the initial mass density $\rho(x, 0)$ and momentum density $\rho(x, 0)v(x, 0)$ are

step-like functions. In view of (1.5), this amounts to requiring α and β of system (1.21) to have step-like initial data. We are interested in the following two cases:

(i) $\alpha(x, 0)$ is a constant, and

$$(1.23) \quad \alpha(x, 0) = a, \quad \beta(x, 0) = \begin{cases} b, & x < 0, \\ c, & x > 0, \end{cases} \quad a > b, \ a > c, \ b \neq c,$$

(ii) $\beta(x, 0)$ is a constant, and

$$(1.24) \quad \alpha(x, 0) = \begin{cases} b, & x < 0, \\ c, & x > 0, \end{cases} \quad \beta(x, 0) = a, \quad b > a, \ c > a, \ b \neq c.$$

In the case of the NLS equation, the genuine nonlinearity of the single phase Whitham equations (see (1.17)) warrants that the solution is found by the implicit function theorem. However, the mKdV–Whitham equations (1.22) generally are not genuinely nonlinear; that is, a property like (1.17) is not available (see Appendix B). Our construction of solutions of the Whitham equations (1.22) with $g = 1$ makes use of the non-strict hyperbolicity of the equations. For the NLS case, it is known from [6, 7] that the Whitham equations (1.7) with $g = 1$ are strictly hyperbolic; that is,

$$\lambda_{1,1} > \lambda_{1,2} > \lambda_{1,3} > \lambda_{1,4}$$

for $u_1 > u_2 > u_3 > u_4$. For the mKdV–Whitham equations (1.22) with $g = 1$, the eigen speeds $\mu_{1,i}(u_1, u_2, u_3, u_4)$ may coalesce in the region $u_1 > u_2 > u_3 > u_4$.

Let us now describe one of our main results (see Theorem 3.1) for the single phase mKdV–Whitham equations with step-like initial function (1.23) for $a = 4, b = 1$, and $c = -1$. In this case, the space time is divided into *four* regions (see Figure 1.2) instead of *three* in the case of the NLS equation (cf. Figure 1.1):

$$(1) \ \frac{x}{t} < c_1, \quad (2) \ c_1 < \frac{x}{t} < c_2, \quad (3) \ c_2 < \frac{x}{t} < c_3, \quad (4) \ \frac{x}{t} > c_3,$$

where c_1, c_2 , and c_3 are some constants. In the first and fourth regions, the solution of the 2×2 system (1.21) governs the evolution:

(1) for $x/t < c_1$,

$$\alpha(x, t) = 4, \quad \beta(x, t) = 1;$$

(4) for $x/t > c_3$,

$$\alpha(x, t) = 4, \quad \beta(x, t) = -1.$$

The Whitham solution of the 4×4 system (1.22) with $g = 1$ lives in the second and third regions:

(2) for $c_1 < x/t < c_2$,

$$u_1(x, t) = 4, \quad u_2(x, t) = 1, \quad \frac{x}{t} = \mu_{1,3}(4, 1, u_3, u_4), \quad \frac{x}{t} = \mu_{1,4}(4, 1, u_3, u_4);$$

(3) for $c_2 < x/t < c_3$,

$$u_1(x, t) = 4, \quad u_2(x, t) = 1, \quad \frac{x}{t} = \mu_{1,3}(4, 1, u_3, -1), \quad u_4(x, t) = -1.$$

Note that, in the second region, we have

$$\mu_{1,3}(4, 1, u_3, u_4) = \mu_{1,4}(4, 1, u_3, u_4)$$

on a curve in the region $-1 < u_4 < u_3 < 1$. This implies the non-strict hyperbolicity of the mKdV–Whitham equations (1.22) for $g = 1$.

It is again possible to use the method of [1, 16] to show that the solution of the mKdV equation (1.18) can be approximately described, in the single phase regime, by the periodic solution of the mKdV when ϵ is small. The periodic solution has the same form as (1.8) of the NLS, i.e.,

$$(1.25) \quad \tilde{\rho}(x, t; \epsilon) = \rho_3 + (\rho_2 - \rho_3) \operatorname{sn}^2 \left(\sqrt{\rho_1 - \rho_3} \frac{\theta(x, t)}{\epsilon}, s \right).$$

However, $\theta(x, t)$ is now given by $\theta = x - V_2 t - Q$ with the velocity V_2 (see, e.g., [9]):

$$V_2 = \frac{3}{8}\sigma_1^2 - \frac{1}{2}\sigma_2,$$

where $\sigma_1 = \sum_{j=1}^4 u_j$ and $\sigma_2 = \sum_{i < j} u_i u_j$ are the elementary symmetric functions of degree one and two, respectively. The functions ρ_1 , ρ_2 , and ρ_3 are also given by formula (1.9), and the flow velocity formula (1.10) is still valid for the mKdV. If u_1 , u_2 , u_3 , and u_4 are constants, formula (1.25) gives the periodic solution of the mKdV equation. To describe the solution $\rho(x, t; \epsilon)$ of the mKdV equation (1.18), the quantities u_1 , u_2 , u_3 , and u_4 must satisfy the single phase mKdV–Whitham equations (1.22) for $g = 1$. The phase shift Q is still determined by (1.12) and (1.13) if the initial function of system (1.21) is given by (1.11). The weak limit of $\rho(x, t; \epsilon)$ of the mKdV equation is also given by formula (1.14).

In Figure 1.2, we plot the self-similar solution of the Whitham equations (1.22) for $g = 1$ and the corresponding periodic oscillatory solution (1.25). We note here that the pattern of the oscillations in this case has two different kinds of structure: one corresponds to the region (2), $c_1 < x/t < c_2$, and the other corresponds to the region (3), $c_2 < x/t < c_3$. In Figure 1.2, those regions are separated by a dashed-dotted line. This is in sharp contrast to the NLS case where the oscillations have uniform structure (cf. Figure 1.1). We also note that the weak limit $\overline{\rho(x, t)}$ and the envelope of the oscillations are not C^1 smooth at $x/t = c_2 \approx 3.67$, while they are smooth (even analytic) everywhere within the oscillatory region in the NLS case. Finally, the shapes of the envelopes of the oscillations are quite different: the one on the right of Figure 1.2 looks like a Bordeaux glass and that of Figure 1.1 resembles a martini glass.

As we will show below, for other values of a , b , and c , the solutions of (1.21) and (1.22) with $g = 1$ will be seen to be quite different from the above.

The Whitham equations (1.22) with $g = 1$ for the mKdV equation are analogous to the Whitham equations for the fifth order KdV equation [11]; both Whitham equations are neither strictly hyperbolic nor genuinely nonlinear. For all the step-like initial shock data, the single phase Whitham solutions for the fifth order KdV are also constructed using the non-strict hyperbolicity of the equations. In the case of KdV, the Whitham equations are strictly hyperbolic and genuinely nonlinear, and the oscillations (dispersive shock) have uniform structure. However, in the case of the fifth order KdV, new structure has been found in the dispersive shocks. This structure is similar to the one found here in the dispersive shocks of the mKdV equation.

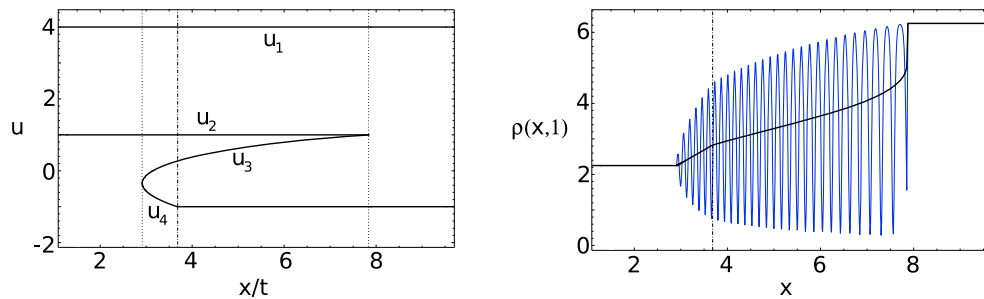


FIG. 1.2. Self-similar solution of the mKdV-Whitham equation (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.1$. The initial data are given by (1.23) with $a = 4$, $b = 1$, and $c = -1$. There are two different kinds of structure in the oscillations, and they are separated by $x/t \approx 3.67$. Both the weak limit and the envelope of the oscillations have a drastic change at $x/t \approx 3.67$.

There are still some significant differences between the mKdV oscillations and the fifth order KdV oscillations. The biggest difference lies in the fact that, in the former case, the zero phase Whitham equations (1.21) form a system of two equations, while in the latter case, the zero phase Whitham equation is a scalar equation. In this paper, we consider initial functions (1.23) and (1.24) in which either $\alpha(x, 0)$ or $\beta(x, 0)$ is a constant; thus equations (1.21) reduce to a single equation. However, the constant initial function $\alpha(x, 0) = a$ of (1.23) or $\beta(x, 0) = a$ of (1.24) still has a significant effect on the behavior of the Whitham solutions (cf. Figures 3.1 and 4.1). Another difference is that mKdV solution ρ can have points (x, t) at which $\rho(x, t; \epsilon) = 0$ for a sequence of vanishing ϵ (see section 5). These points put the mKdV dispersive approximation (1.19) at stake. This phenomenon does not occur in the fifth order KdV case.

The Whitham equations for the higher members of the KdV hierarchy are again neither strictly hyperbolic nor genuinely nonlinear. However, step-like initial shock data will, in general, generate a mixture of single, double, or even higher phases in the higher order KdV oscillations. It has been an open problem to construct such multiphase Whitham solutions. Only rather special step-like initial shock data will produce merely single phase oscillations in the higher order KdV case. Some of these initial data have been studied in [12]. The analogues of inequalities (2.15) and (2.16), which play a crucial role in the fifth order KdV case, are not valid anymore in the higher order KdV case. As a consequence, the approach in [12] is quite different from that in [11]. Indeed, the calculations in [12] are considerably more difficult than in [11].

The Whitham equations for the higher members of the KdV hierarchy are supposed to be analogous to the Whitham equations for the higher members of the (defocusing) NLS hierarchy. It would be very interesting to see how step-like initial shock data generate multiphases in the higher order NLS oscillations.

The organization of the paper is as follows. In section 2, we will study the eigenspeeds $\mu_{g,1}$, $\mu_{g,2}$, $\mu_{g,3}$, and $\mu_{g,4}$ of the Whitham equations (1.22) for $g = 1$. In section 3, we will construct the self-similar solutions of the single phase Whitham equations for the initial function (1.23) with $a > b > c$. In section 4, we will construct the self-similar solution of the Whitham equations for the initial function (1.23) with $a > c > b$. In section 5, we will study those points (x, t) at which the mKdV solution $\rho(x, t; \epsilon) = 0$ for a sequence of vanishing ϵ . In section 6, we will briefly discuss how to handle the other step-like initial data (1.24).

2. The Whitham equations. In this section we define the eigenspeeds $\lambda_{g,i}$ and $\mu_{g,i}$ of the Whitham equations (1.7) and (1.22) with $g = 1$ for the NLS and the mKdV equations. For simplicity, we suppress the subscript $g = 1$ in the notation $\lambda_{g,i}$ and $\mu_{g,i}$ in the rest of the paper.

We first introduce the polynomials of ξ for $n = 0, 1, 2, \dots$ [2, 4, 10]:

$$(2.1) \quad P_n(\xi, u_1, u_2, u_3, u_4) = \xi^{n+2} + a_{n,1}\xi^{n+1} + \dots + a_{n,n+2},$$

where the coefficients $a_{n,1}, a_{n,2}, \dots, a_{n,n+2}$ are uniquely determined by the two conditions

$$\frac{P_n(\xi, u_1, u_2, u_3, u_4)}{\sqrt{(\xi - u_1)(\xi - u_2)(\xi - u_3)(\xi - u_4)}} = \xi^n + \mathcal{O}(\xi^{-2}) \quad \text{for large } |\xi|$$

and

$$\int_{u_2}^{u_1} \frac{P_n(\xi, u_1, u_2, u_3, u_4)}{\sqrt{(u_1 - \xi)(\xi - u_2)(\xi - u_3)(\xi - u_4)}} d\xi = 0.$$

The coefficients of P_n can be expressed in terms of complete elliptic integrals.

The eigenspeeds of the Whitham equations (1.7) with $g = 1$ for the NLS equation are defined in terms of P_0 and P_1 of (2.1) [4, 6, 10],

$$\lambda_i(u_1, u_2, u_3, u_4) = 8 \frac{P_1(u_i, u_1, u_2, u_3, u_4)}{P_0(u_i, u_1, u_2, u_3, u_4)}, \quad i = 1, 2, 3, 4,$$

which give

$$(2.2) \quad \lambda_i(u_1, u_2, u_3, u_4) = 2 \left(\sigma_1(u_1, u_2, u_3, u_4) - \frac{I(u_1, u_2, u_3, u_4)}{\partial_{u_i} I(u_1, u_2, u_3, u_4)} \right).$$

Here $\sigma_1 := \sum_{j=1}^4 u_j$, and $I(u_1, u_2, u_3, u_4)$ is given by a complete elliptic integral [14]

$$(2.3) \quad I(u_1, u_2, u_3, u_4) = \int_{u_2}^{u_1} \frac{d\eta}{\sqrt{(u_1 - \eta)(\eta - u_2)(\eta - u_3)(\eta - u_4)}}.$$

The function I can be rewritten as a contour integral. Hence,

$$(2.4) \quad 2(u_i - u_j) \frac{\partial^2 I}{\partial u_i \partial u_j} = \frac{\partial I}{\partial u_i} - \frac{\partial I}{\partial u_j}, \quad i, j = 1, 2, 3, 4,$$

since the integrand satisfies the same equations for each $\eta \neq u_i$ (cf. (1.12)). This contour integral connection also allows us to give another formulation of I ,

$$(2.5) \quad I(u_1, u_2, u_3, u_4) = \int_{u_4}^{u_3} \frac{d\eta}{\sqrt{(u_1 - \eta)(u_2 - \eta)(u_3 - \eta)(\eta - u_4)}}.$$

It follows from (2.2), (2.3), and (2.5) that

$$(2.6) \quad \lambda_4 - 2\sigma_1 < \lambda_3 - 2\sigma_1 < 0 < \lambda_2 - 2\sigma_1 < \lambda_1 - 2\sigma_1$$

for $u_4 < u_3 < u_2 < u_1$. This implies the strict hyperbolicity of the NLS–Whitham equation (1.7) for $g = 1$.

The eigenspeeds λ_i have the following values [14]: At $u_3 = u_4$, we have

$$(2.7) \quad \begin{cases} \lambda_1 = 6u_1 + 2u_2, \\ \lambda_2 = 2u_1 + 6u_2, \\ \lambda_3 = \lambda_4 = 2(u_1 + u_2 + 2u_4) - \frac{8(u_1 - u_4)(u_2 - u_4)}{u_1 + u_2 - 2u_4}, \end{cases}$$

and at $u_2 = u_3$,

$$(2.8) \quad \begin{cases} \lambda_1 = 6u_1 + 2u_4, \\ \lambda_2 = \lambda_3 = 2u_1 + 4u_3 + 2u_4, \\ \lambda_4 = 2u_1 + 6u_4. \end{cases}$$

Notice that the eigenspeed $\lambda_2 = \lambda_3$ at $u_2 = u_3$ is the same as the velocity of the periodic solution (1.8), i.e., $V_1 = 2\sigma_1 = 2(u_1 + 2u_3 + u_4)$.

The eigenspeeds of the mKdV–Whitham equations (1.22) with $g = 1$ are [6]

$$(2.9) \quad \mu_i(u_1, u_2, u_3, u_4) = 3 \frac{P_2(u_i, u_1, u_2, u_3, u_4)}{P_0(u_i, u_1, u_2, u_3, u_4)}, \quad i = 1, 2, 3, 4.$$

They can be expressed in terms of $\lambda_1, \lambda_2, \lambda_3$, and λ_4 of the NLS–Whitham equations (1.7) with $g = 1$.

LEMMA 2.1. *The eigenspeeds $\mu_i(u_1, u_2, u_3, u_4)$ of (2.9) can be expressed in the form*

$$(2.10) \quad \mu_i = \frac{1}{2}(\lambda_i - 2\sigma_1) \frac{\partial q}{\partial u_i} + q, \quad i = 1, 2, 3, 4,$$

where $\sigma_1 = \sum_{j=1}^4 u_j$ and $q = q(u_1, u_2, u_3, u_4)$ is the solution of the boundary value problem of the Euler–Poisson–Darboux equations (cf. (1.12) and (2.4))

$$(2.11) \quad 2(u_i - u_j) \frac{\partial^2 q}{\partial u_i \partial u_j} = \frac{\partial q}{\partial u_i} - \frac{\partial q}{\partial u_j}, \quad i, j = 1, 2, 3, 4,$$

$$q(u, u, u) = 3u^2.$$

Also the μ_i satisfy the overdetermined systems

$$(2.12) \quad \frac{1}{\mu_i - \mu_j} \frac{\partial \mu_i}{\partial u_j} = \frac{1}{\lambda_i - \lambda_j} \frac{\partial \lambda_i}{\partial u_j}, \quad i \neq j.$$

We omit the proof since it is very similar to the proof of an analogous result for the KdV hierarchy [13].

The boundary value problem (2.11) has a unique solution. The solution is a symmetric quadratic function of u_1, u_2, u_3 , and u_4 :

$$(2.13) \quad q = \frac{3}{8}\sigma_1^2 - \frac{1}{2}\sigma_2,$$

where $\sigma_2 = \sum_{i>j} u_i u_j$ is the elementary symmetric polynomial of degree two. Notice that q gives the velocity of the periodic solution (1.25) for the mKdV equation, i.e., $V_2 = q$.

For NLS, λ_i satisfy [14]

$$(2.14) \quad \frac{\partial \lambda_4}{\partial u_4} < \frac{3 \lambda_3 - \lambda_4}{2 u_3 - u_4} < \frac{\partial \lambda_3}{\partial u_3}$$

for $u_4 < u_3 < u_2 < u_1$. Similar results also hold for the mKdV–Whitham equations (1.22) with $g = 1$.

LEMMA 2.2.

$$(2.15) \quad \frac{\partial \mu_3}{\partial u_3} > \frac{3}{2} \frac{\mu_3 - \mu_4}{u_3 - u_4} \quad \text{if} \quad \frac{\partial q}{\partial u_3} > 0,$$

$$(2.16) \quad \frac{\partial \mu_4}{\partial u_4} < \frac{3}{2} \frac{\mu_3 - \mu_4}{u_3 - u_4} \quad \text{if} \quad \frac{\partial q}{\partial u_4} > 0$$

for $u_4 < u_3 < u_2 < u_1$.

Proof. We use (2.10) and (2.14) to obtain

$$(2.17) \quad \begin{aligned} \frac{\partial \mu_3}{\partial u_3} &= \frac{1}{2} \frac{\partial \lambda_3}{\partial u_3} \frac{\partial q}{\partial u_3} + \frac{1}{2} (\lambda_3 - 2\sigma_1) \frac{\partial^2 q}{\partial u_3^2} \\ &> \frac{3}{4} \frac{\lambda_3 - \lambda_4}{u_3 - u_4} \frac{\partial q}{\partial u_3} + \frac{1}{2} (\lambda_4 - 2\sigma_1) \frac{\partial^2 q}{\partial u_3^2} \end{aligned}$$

and

$$(2.18) \quad \begin{aligned} \mu_3 - \mu_4 &= \frac{1}{2} (\lambda_3 - \lambda_4) \frac{\partial q}{\partial u_3} + \frac{1}{2} (\lambda_4 - 2\sigma_1) \left(\frac{\partial q}{\partial u_3} - \frac{\partial q}{\partial u_4} \right) \\ &= \frac{1}{2} (\lambda_3 - \lambda_4) \frac{\partial q}{\partial u_3} + (\lambda_4 - 2\sigma_1) (u_3 - u_4) \frac{\partial^2 q}{\partial u_3 \partial u_4} \\ &= \frac{2}{3} (u_3 - u_4) \left(\frac{3}{4} \frac{\lambda_3 - \lambda_4}{u_2 - u_3} \frac{\partial q}{\partial u_3} + \frac{3}{2} (\lambda_4 - 2\sigma_1) \frac{\partial^2 q}{\partial u_3 \partial u_4} \right), \end{aligned}$$

where we have used (2.11)

$$\frac{\partial q}{\partial u_3} - \frac{\partial q}{\partial u_4} = 2(u_3 - u_4) \frac{\partial^2 q}{\partial u_3 \partial u_4}.$$

Differentiating this equation with respect to u_3 yields

$$\frac{\partial^2 q}{\partial u_3^2} - 3 \frac{\partial^2 q}{\partial u_3 \partial u_4} = 2(u_3 - u_4) \frac{\partial^3 q}{\partial u_3^2 \partial u_4}.$$

It is here we exploit the fact that q is a quadratic polynomial and conclude that

$$3 \frac{\partial^2 q}{\partial u_3 \partial u_4} = \frac{\partial^2 q}{\partial u_3^2},$$

which, along with (2.17) and (2.18), proves (2.15). Inequality (2.16) can be proved in the same way. \square

The following calculations will be useful in the subsequent sections. Using formula (2.10) for μ_3 and μ_4 and formula (2.2) for λ_3 and λ_4 , we obtain

$$\begin{aligned}
 \mu_3 - \mu_4 &= \frac{I}{(\partial_{u_3} I)(\partial_{u_4} I)} \left[\frac{\partial q}{\partial u_4} \frac{\partial I}{\partial u_3} - \frac{\partial q}{\partial u_3} \frac{\partial I}{\partial u_4} \right] \\
 &= \frac{I}{(\partial_{u_3} I)(\partial_{u_4} I)} \left[\frac{\partial q}{\partial u_4} \left(\frac{\partial I}{\partial u_3} - \frac{\partial I}{\partial u_4} \right) - \left(\frac{\partial q}{\partial u_3} - \frac{\partial q}{\partial u_4} \right) \frac{\partial I}{\partial u_4} \right] \\
 (2.19) \quad &= \frac{2I(u_3 - u_4)}{(\partial_{u_3} I)(\partial_{u_4} I)} M,
 \end{aligned}$$

where

$$M = \frac{\partial q}{\partial u_4} \frac{\partial^2 I}{\partial u_3 \partial u_4} - \frac{\partial^2 q}{\partial u_3 \partial u_4} \frac{\partial I}{\partial u_4}.$$

Here we have used equations (2.4) for I and equations (2.11) for q in equality (2.19). Since q of (2.13) is quadratic, we obtain

$$(2.20) \quad \frac{\partial M}{\partial u_3} = \frac{\partial q}{\partial u_4} \frac{\partial^3 I}{\partial u_3^2 \partial u_4}.$$

We note that another expression for M is

$$M = \frac{\partial q}{\partial u_3} \frac{\partial^2 I}{\partial u_3 \partial u_4} - \frac{\partial^2 q}{\partial u_3 \partial u_4} \frac{\partial I}{\partial u_3}.$$

Hence, we get

$$(2.21) \quad \frac{\partial M}{\partial u_4} = \frac{\partial q}{\partial u_3} \frac{\partial^3 I}{\partial u_3 \partial u_4^2}.$$

We next evaluate $M(u_1, u_2, u_3, u_4)$ when $u_3 = u_4$. Using the integral formula (2.3) for the function I and applying the change of variable $\eta = (u_1 - u_2)\nu + u_2$, we obtain

$$\begin{aligned}
 M \Big|_{u_3=u_4} &= \frac{\frac{\partial q}{\partial u_4}}{4(u_2 - u_4)^3} \int_0^1 \frac{d\nu}{\left(1 + \frac{u_1 - u_2}{u_2 - u_4} \nu\right)^3 \sqrt{\nu(1 - \nu)}} \\
 &\quad - \frac{\frac{\partial^2 q}{\partial u_3 \partial u_4}}{2(u_2 - u_4)^2} \int_0^1 \frac{d\nu}{\left(1 + \frac{u_1 - u_2}{u_2 - u_4} \nu\right)^2 \sqrt{\nu(1 - \nu)}}.
 \end{aligned}$$

The two integrals can be evaluated exactly as

$$\int_0^1 \frac{d\nu}{(1 + \gamma\nu)^3 \sqrt{\nu(1 - \nu)}} = \frac{\pi(8 + 8\gamma + 3\gamma^2)}{8(1 + \gamma)^{\frac{5}{2}}}, \quad \int_0^1 \frac{d\nu}{(1 + \gamma\nu)^2 \sqrt{\nu(1 - \nu)}} = \frac{\pi(2 + \gamma)}{2(1 + \gamma)^{\frac{3}{2}}}$$

for $\gamma > -1$. We finally get

$$(2.22) \quad M \Big|_{u_3=u_4} = \frac{\pi U(u_1, u_2, u_4)}{128[(u_2 - u_4)(u_1 - u_4)]^{\frac{5}{2}}},$$

where

$$\begin{aligned}
 U(u_1, u_2, \xi) &= [8(u_2 - \xi)^2 + 8(u_2 - \xi)(u_1 - u_2) + 3(u_1 - u_2)^2](u_1 + u_2 + 4\xi) \\
 (2.23) \qquad &\quad - 8(u_1 - \xi)(u_2 - \xi)(u_1 + u_2 - 2\xi).
 \end{aligned}$$

Similarly to (2.19) for μ_3 and μ_4 , we have

$$(2.24) \qquad \mu_2 - \mu_3 = \frac{2I(u_2 - u_3)}{(\partial_{u_2} I)(\partial_{u_3} I)} N,$$

where

$$N = \frac{\partial q}{\partial u_2} \frac{\partial^2 I}{\partial u_2 \partial u_3} - \frac{\partial^2 q}{\partial u_2 \partial u_3} \frac{\partial I}{\partial u_2}.$$

Since q of (2.13) is quadratic, we obtain

$$(2.25) \qquad \frac{\partial N}{\partial u_3} = \frac{\partial q}{\partial u_2} \frac{\partial^3 I}{\partial u_2 \partial u_3^2}.$$

Finally, we use (2.7) and (2.10) to calculate

$$\begin{aligned}
 (\mu_2 - \mu_3) \Big|_{u_3=u_4} &= \frac{1}{2} [\lambda_2 - 2(u_1 + u_2 + 2u_4)] \frac{\partial q}{\partial u_2} - \frac{1}{2} [\lambda_3 - 2(u_1 + u_2 + 2u_4)] \frac{\partial q}{\partial u_3} \\
 (2.26) \qquad &= \frac{(u_2 - u_4)}{2(u_1 + u_2 - 2u_4)} V(u_1, u_2, u_4),
 \end{aligned}$$

where

$$(2.27) \qquad V(u_1, u_2, u_4) = 3u_1^2 + 3u_2^2 - 12u_4^2 + 6u_1u_2 + 6u_1u_4 - 6u_2u_4.$$

3. Self-similar solutions. In this section, we construct self-similar solutions of the Whitham equations (1.22) with $g = 1$ for the initial function (1.23) with $a > b > c$. The case with $a > c > b$ will be studied in next section. The solution of the zero phase Whitham equations (1.21) does not develop a shock when $a + 5b \leq 0$. We are therefore interested only in the case $a + 5b > 0$.

We first study the ξ -zero of the cubic polynomial equation

$$(3.1) \qquad U(a, b, \xi) = 0,$$

where U is given by (2.23). It is easy to prove that, for each pair of a and b satisfying $a > b$ and $a + 5b > 0$, $U(a, b, \xi) = 0$ has only one simple real root. Denoting this zero by $\xi(a, b)$, we then deduce that $U(a, b, \xi)$ is positive for $\xi > \xi(a, b)$ and negative for $\xi < \xi(a, b)$. Since $U(a, b, -(a + b)/4) < 0$ in view of (2.23), we must have

$$(3.2) \qquad \xi(a, b) > -\frac{a + b}{4}.$$

For initial function (1.23) with $a > b > c$ and $a + 5b > 0$, we now classify the resulting Whitham solutions into four types:

- I. $\xi(a, b) \leq c$ with any $a > b > c$,
- II. $\xi(a, b) > c$ with $a + 5b > 3(b - c) > 0$,
- III. $\xi(a, b) > c$ with $a + 5b = 3(b - c) > 0$,
- IV. $\xi(a, b) > c$ with $0 < a + 5b < 3(b - c)$.

In Figure 3.1, we illustrate this classification by scaling $a > 0$ to $a = 4$ and plotting regions for types I–IV in the b - c plane.

We will study the second type first.

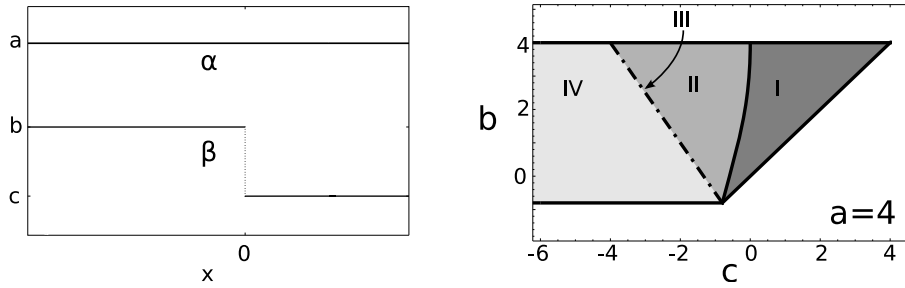


FIG. 3.1. Graph of the step-like initial shock data (1.23) with $a > b > c$ and $a + 5b > 0$, and diagram of regions for types I–IV, where $a > 0$ has been scaled to $a = 4$. The region for type IV is unbounded from the left.

3.1. Type II. Here we consider the step-like initial function (1.23) satisfying $\xi(a, b) > c$ and $a + 5b > 3(b - c) > 0$.

THEOREM 3.1 (see Figure 1.2). *For the step-like initial data (1.23) with $a > b > c$, $a + 5b > 3(b - c)$, and $\xi(a, b) > c$, the solution (α, β) of the zero phase Whitham equations (1.21) and the solution (u_1, u_2, u_3, u_4) of the single phase Whitham equations (1.22) with $g = 1$ are given as follows:*

(1) For $x/t \leq \mu_3(a, b, \xi(a, b), \xi(a, b))$,

$$(3.3) \quad \alpha = a, \quad \beta = b.$$

(2) For $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, u^{**}, c)$,

$$(3.4) \quad u_1 = a, \quad u_2 = b, \quad \frac{x}{t} = \mu_3(a, b, u_3, u_4), \quad \frac{x}{t} = \mu_4(a, b, u_3, u_4),$$

where u^{**} is the unique solution u_3 of $\mu_3(a, b, u_3, c) = \mu_4(a, b, u_3, c)$ in the interval $c < u_3 < b$.

(3) For $\mu_3(a, b, u^{**}, c) \leq x/t < \mu_3(a, b, b, c)$,

$$(3.5) \quad u_1 = a, \quad u_2 = b, \quad \frac{x}{t} = \mu_3(a, b, u_3, c), \quad u_4 = c.$$

(4) For $x/t \geq \mu_3(a, b, b, c)$,

$$(3.6) \quad \alpha = a, \quad \beta = c.$$

The boundaries $x/t = \mu_3(a, b, \xi(a, b), \xi(a, b))$ and $x/t = \mu_3(a, b, b, c)$ are called the trailing and leading edges, respectively. They separate the solutions of the single phase Whitham equations (1.22) with $g = 1$ and the zero phase Whitham equations (1.21). The single phase Whitham solution matches the zero phase Whitham solution in the following fashion (see Figure 1.2):

$$(3.7) \quad (u_1, u_2) = \text{the solution } (\alpha, \beta) \text{ of (1.21) defined outside the region,}$$

$$(3.8) \quad u_3 = u_4,$$

at the trailing edge;

$$(3.9) \quad (u_1, u_4) = \text{the solution } (\alpha, \beta) \text{ of (1.21) defined outside the region,}$$

$$(3.10) \quad u_2 = u_3,$$

at the leading edge.

The proof of Theorem 3.1 is based on a series of lemmas: We first show that the solutions defined by formulae (3.4) and (3.5) indeed satisfy the Whitham equations (1.22) for $g = 1$ [2, 15].

LEMMA 3.2.

- (1) *The functions $u_1, u_2, u_3,$ and u_4 determined by equations (3.4) give a solution of the Whitham equations (1.22) with $g = 1$ as long as u_3 and u_4 can be solved from (3.4) as functions of x and t .*
- (2) *The functions $u_1, u_2, u_3,$ and u_4 determined by equations (3.5) give a solution of the Whitham equations (1.22) with $g = 1$ as long as u_3 can be solved from (3.5) as a function of x and t .*

Proof. (1) u_1 and u_2 obviously satisfy the first two equations of (1.22) for $g = 1$. To verify the third and fourth equations, we observe that

$$(3.11) \quad \frac{\partial \mu_3}{\partial u_4} = \frac{\partial \mu_4}{\partial u_3} = 0$$

on the solution of (3.4). To see this, we use (2.12) to calculate

$$\frac{\partial \mu_3}{\partial u_4} = \frac{\frac{\partial \lambda_3}{\partial u_4}}{\lambda_3 - \lambda_4} (\mu_3 - \mu_4) = 0.$$

The second part of (3.11) can be shown in the same way. We then calculate the partial derivatives of the third equation of (3.4) with respect to x and t ,

$$1 = \frac{\partial \mu_3}{\partial u_3} t u_{3x}, \quad 0 = \frac{\partial \mu_3}{\partial u_3} t u_{3t} + \mu_3,$$

which give the third equation of (1.22) with $g = 1$. The fourth equation of (1.22) with $g = 1$ can be verified in the same way.

- (2) The second part of Lemma 3.2 can easily be proved. □

We now determine the trailing edge. Eliminating x and t from the last two equations of (3.4) yields

$$(3.12) \quad \mu_3(a, b, u_3, u_4) - \mu_4(a, b, u_3, u_4) = 0.$$

Since it degenerates at $u_3 = u_4$, we replace (3.12) by

$$(3.13) \quad F(a, b, u_3, u_4) := \frac{\mu_3(a, b, u_3, u_4) - \mu_4(a, b, u_3, u_4)}{u_3 - u_4} = 0.$$

Therefore, at the trailing edge where $u_3 = u_4$, (3.13), in view of formulae (2.19) and (2.22), reduces to

$$(3.14) \quad U(a, b, u_4) = 0.$$

Noting that $\xi(a, b)$ is the unique solution of (3.1), we then deduce that $u_4 = \xi(a, b)$.

LEMMA 3.3. *Equation (3.13) has a unique solution satisfying $u_3 = u_4$. The solution is $u_3 = u_4 = \xi(a, b)$. The rest of equations (3.4) at the trailing edge are $u_1 = a, u_2 = b,$ and $x/t = \mu_3(a, b, \xi(a, b), \xi(a, b))$.*

Having located the trailing edge, we now solve equations (3.4) in the neighborhood of the trailing edge. We first consider (3.13). We use (2.19) to write F of (3.13) as

$$F(a, b, u_3, u_4) = \frac{2I}{(\partial_{u_3} I)(\partial_{u_3} I)} M(a, b, u_3, u_4).$$

We note that, at the trailing edge $u_3 = u_4 = \xi(a, b)$, we have $M(a, b, \xi(a, b), \xi(a, b)) = 0$ because of (2.22) and (3.14). We then use (2.20) and (2.21) to differentiate F at the trailing edge

$$\begin{aligned} \frac{\partial F(a, b, \xi(a, b), \xi(a, b))}{\partial u_3} &= \frac{\partial F(a, b, \xi(a, b), \xi(a, b))}{\partial u_4} \\ &= \frac{I}{2(\partial_{u_3} I)(\partial_{u_3} I)} [a + b + 4\xi(a, b)] \frac{\partial^3 I}{\partial u_3^2 \partial u_4} > 0, \end{aligned}$$

where we have used the expression (2.13) for q in the last equation and (3.2) in the inequality. These show that (3.13) or, equivalently, (3.12) can be inverted to give u_4 as a decreasing function of u_3 ,

$$(3.15) \quad u_4 = A(u_3),$$

in a neighborhood of $u_3 = u_4 = \xi(a, b)$.

We now extend the solution $A(u_3)$ of (3.12) in the region $c < u_4 < \xi(a, b) < u_3 < b$ as far as possible. We first claim that

$$(3.16) \quad \frac{\partial q(a, b, u_3, u_4)}{\partial u_3} > 0, \quad \frac{\partial q(a, b, u_3, u_4)}{\partial u_4} > 0$$

on the extension. To see this, we first observe that inequalities (3.16) are true at the trailing edge $u_3 = u_4 = \xi(a, b)$. This follows from (2.13) and (3.2). Therefore, inequalities (3.16) hold in a neighborhood of the trailing edge. To prove that (3.16) remains true on the extension, we use formula (2.10) for μ_3 and μ_4 to rewrite (3.12) as

$$\frac{1}{2}[\lambda_3 - 2(a + b + u_3 + u_4)] \frac{\partial q}{\partial u_3} = \frac{1}{2}[\lambda_4 - 2(a + b + u_3 + u_4)] \frac{\partial q}{\partial u_4}.$$

Since the two terms in the two sets of parentheses are both negative in view of (2.6) and since $\frac{\partial q}{\partial u_3} - \frac{\partial q}{\partial u_4} = (u_3 - u_4)/4 > 0$ in view of (2.13), neither $\frac{\partial q}{\partial u_3}$ nor $\frac{\partial q}{\partial u_4}$ can vanish on the extension. This proves inequalities (3.16).

We deduce from Lemma 2.2 that

$$(3.17) \quad \frac{\partial \mu_3}{\partial u_3} > 0, \quad \frac{\partial \mu_4}{\partial u_4} < 0$$

on the solution of (3.12). Because of (3.11) and (3.17), solution (3.15) of (3.12) can be extended as long as $c < u_4 < \xi(a, b) < u_3 < b$.

There are two possibilities: (1) u_3 touches b before or simultaneously as u_4 reaches c and (2) u_4 touches c before u_3 reaches b . It follows from (2.8), (2.10), and (2.13) that

$$(3.18) \quad \mu_3(a, b, b, u_4) - \mu_4(a, b, b, u_4) = \frac{1}{2} (b - u_4)(a + 2b + 3u_4) > 0 \quad \text{for } c \leq u_4 < b,$$

where we have used $a + 2b + 3c > 0$ in the inequality. This shows that (1) is unattainable. Hence, u_4 will touch c before u_3 reaches b . When this happens, (3.12) becomes

$$(3.19) \quad \mu_3(a, b, u_3, c) = \mu_4(a, b, u_3, c).$$

LEMMA 3.4. Equation (3.19) has a simple zero in the interval $c < u_3 < b$, counting multiplicities. If we denote the zero by u^{**} , then $\mu_3(a, b, u_3, c) - \mu_4(a, b, u_3, c)$ is positive for $u_3 > u^{**}$ and negative for $u_3 < u^{**}$.

Proof. We use (2.19) and (2.20) to prove the lemma. In both formulae, $\partial_{u_3} I$, $\partial_{u_4} I$, and $\partial_{u_3 u_4}^2 I$ are all positive functions. By (2.20),

$$(3.20) \quad \frac{\partial M(a, b, u_3, c)}{\partial u_3} = \frac{(a + b + u_3 + 3c)}{4} \frac{\partial^3 I}{\partial u_3^2 \partial u_4} \quad \text{for } c < u_3 < b.$$

We claim that

$$M(a, b, u_3, c) < 0 \quad \text{when } u_3 = c \quad \text{and} \quad M(a, b, u_3, c) > 0 \quad \text{for } u_3 \text{ near } b.$$

The second inequality follows from (2.19) and (3.18). The first inequality can be deduced from formula (2.22):

$$M(a, b, c, c) = \frac{\pi U(a, b, c)}{128[(b - c)(a - c)]^{\frac{5}{2}}} < 0 \quad \text{for } c < \xi(a, b).$$

Therefore, $M(a, b, u_3, c)$ has a zero in the interval $c < u_3 < b$. The uniqueness of the zero follows from (3.20) in that $M(a, b, u_3, c)$ increases or changes from decreasing to increasing as u_3 increases. This zero is exactly u^{**} , and the rest of the theorem can be proved easily. \square

Having solved (3.12) for u_4 as a decreasing function of u_3 for $c < u_4 < \xi(a, b) < u_3 < b$, we turn to equations (3.4). Because of (3.11) and (3.17), the third equation of (3.4) gives u_3 as an increasing function of x/t for $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, u^{**}, c)$. Consequently, u_4 is a decreasing function of x/t in the same interval.

LEMMA 3.5. The last two equations of (3.4) can be inverted to give u_3 and u_4 as increasing and decreasing functions, respectively, of the self-similarity variable x/t in the interval $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, u^{**}, c)$, where u^{**} is given in Lemma 3.4.

We now turn to equations (3.5). We want to solve the third equation when $x/t > \mu_3(a, b, u^{**}, c)$ or, equivalently, when $u_3 > u^{**}$. According to Lemma 3.4, $\mu_3(a, b, u_3, c) - \mu_4(a, b, u_3, c) > 0$ for $u^{**} < u_3 < b$. In view of (3.16), $\partial_{u_3} q(a, b, u_3, c) = (a + b + 3u_3 + c)/4$ is positive at $u_3 = u^{**}$, and, hence, it remains positive for $u_3 > u^{**}$. By (2.15), we have

$$\frac{\partial \mu_3(a, b, u_3, c)}{\partial u_3} > 0.$$

Hence, the third equation of (3.5) can be solved for u_3 as an increasing function of x/t as long as $u^{**} < u_3 < b$. When u_3 reaches b , we have $x/t = \mu_3(a, b, b, c)$. We have therefore proved the following result.

LEMMA 3.6. The third equation of (3.5) can be inverted to give u_3 as an increasing function of x/t in the interval $\mu_3(a, b, u^{**}, c) \leq x/t \leq \mu_3(a, b, b, c)$.

We are ready to conclude the proof of Theorem 3.1. The solutions (3.3) and (3.6) are obvious. According to Lemma 3.5, the last two equations of (3.4) determine u_3 and u_4 as functions of x/t in the region $\mu_3(a, b, \xi(a, b), \xi(a, b)) \leq x/t \leq \mu_3(a, b, u^{**}, c)$. By the first part of Lemma 3.2, the resulting u_1, u_2, u_3 , and u_4 satisfy the Whitham equations (1.22) with $g = 1$. Furthermore, the boundary conditions (3.7) and (3.8) are satisfied at the trailing edge $x = \mu_3(a, b, \xi(a, b), \xi(a, b))$.

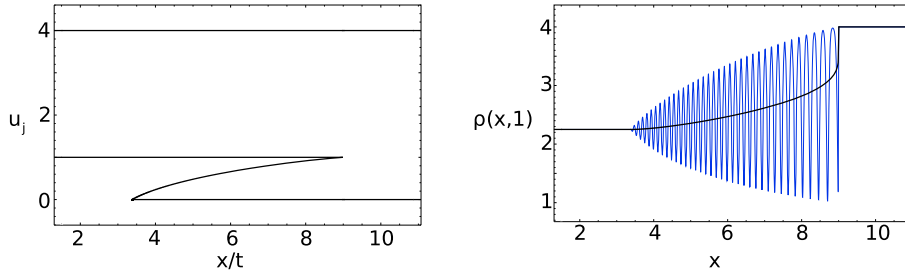


FIG. 3.2. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.07$. The initial data are given by (1.23) with $a = 4$, $b = 1$, and $c = 0$ of type I.

Similarly, by Lemma 3.6, the third equation of (3.5) determines u_3 as a function of x/t in the region $\mu_3(a, b, u^{**}, c) \leq x/t \leq \mu_3(a, b, b, c)$. It then follows from the second part of Lemma 3.2 that u_1, u_2, u_3 , and u_4 of (3.5) satisfy the Whitham equations (1.22) for $g = 1$. They also satisfy the boundary conditions (3.9) and (3.10) at the leading edge $x/t = \mu_3(a, b, b, c)$. We have therefore completed the proof of Theorem 3.1.

A graph of the Whitham solution (u_1, u_2, u_3, u_4) is given in Figure 1.2. It is obtained by plotting the exact solutions of (3.4) and (3.5).

3.2. Type I. Here we consider the initial function (1.23) satisfying $\xi(a, b) \leq c$ with $b > c$ and $a + 5b > 0$.

We will present our proofs only briefly, since they are, more or less, similar to those in section 3.1. The main feature of this case is that the ξ -zero point does not appear in the solution u_3 , and the Whitham equations (1.22) with $g = 1$ are strictly hyperbolic on the solution.

THEOREM 3.7 (see Figure 3.2). *For the step-like initial data (1.23) with $a > b > c$, $a + 5b > 0$, and $\xi(a, b) \leq c$, the solution of the Whitham equations (1.22) with $g = 1$ is given by*

$$u_1 = a, \quad u_2 = b, \quad \frac{x}{t} = \mu_3(a, b, u_3, c), \quad u_4 = c$$

for $\mu_3(a, b, c, c) < x/t < \mu_3(a, b, b, c)$. Outside this interval, the solution of (1.21) is given by

$$\alpha = a, \quad \beta = b \quad \text{for} \quad \frac{x}{t} \leq \mu_3(a, b, c, c)$$

and

$$\alpha = a, \quad \beta = c \quad \text{for} \quad \frac{x}{t} \geq \mu_3(a, b, b, c).$$

Proof. It suffices to show that $\mu_3(a, b, u_3, c)$ is an increasing function of u_3 for $c < u_3 < b$. Substituting (2.13) for q into (2.20) yields

$$\frac{\partial M(a, b, u_3, c)}{\partial u_3} = \frac{1}{4}[a + b + u_3 + 3c] \frac{\partial^3 I}{\partial u_3^2 \partial u_4} \geq \frac{1}{4}[a + b + 4\xi(a, b)] \frac{\partial^3 I}{\partial u_3^2 \partial u_4} > 0$$

for $c < u_3 < b$, where we have used $\xi(a, b) \leq c$ in the first inequality and (3.2) in the second. We now use formula (2.22) to calculate the value of $M(a, b, u_3, c)$ at

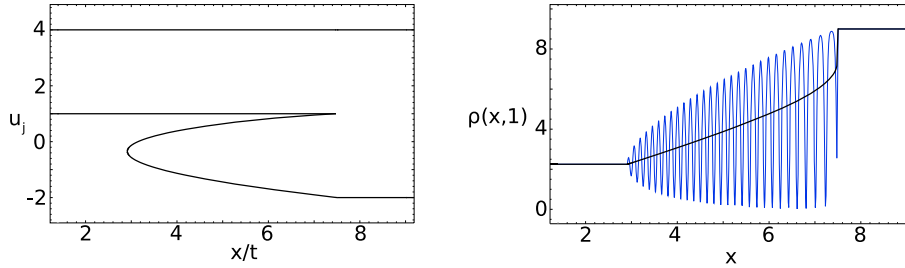


FIG. 3.3. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.1$. The initial data are given by (1.23) with $a = 4, b = 1, c = -2$ of type III.

$u_3 = c$:

$$M(a, b, c, c) = \frac{\pi U(a, b, c)}{128[(b - c)(a - c)]^{\frac{3}{2}}} \geq 0 \quad \text{for } \xi(a, b) \leq c$$

because $U(a, b, \xi) \geq 0$ for $\xi \geq \xi(a, b)$. Therefore, $M(a, b, u_3, c) > 0$ for $c < u_3 < b$. It then follows from (2.19) that $\mu_3(a, b, u_3, c) - \mu_4(a, b, u_3, c) > 0$. Since $\frac{\partial q}{\partial u_3}(a, b, u_3, c) = (a + b + 3u_3 + c)/4 > (a + b + 3\xi(a, b))/4 > 0$ because of (3.2), we conclude from Lemma 2.2 that

$$\frac{d\mu_3(a, b, u_3, c)}{du_3} > 0$$

for $c < u_3 < b$. \square

3.3. Type III. Here we consider the step-like initial function (1.23) satisfying $\xi(a, b) > c$ with $a + 5b = 3(b - c) > 0$.

THEOREM 3.8 (see Figure 3.3). *For the step-like initial data (1.23) with $a > b > c, \xi(a, b) > c$, and $a + 5b = 3(b - c)$, the solution of the $g = 1$ Whitham equations (1.22) with $g = 1$ is given by*

$$u_1 = a, \quad u_2 = b, \quad \frac{x}{t} = \mu_3(a, b, u_3, u_4), \quad \frac{x}{t} = \mu_4(a, b, u_3, u_4)$$

for $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, b, c)$. Outside the region, the solution of (1.21) is given by

$$\alpha = a, \quad \beta = b \quad \text{for } \frac{x}{t} \leq \mu_3(a, b, \xi(a, b), \xi(a, b))$$

and

$$\alpha = a, \quad \beta = c \quad \text{for } \frac{x}{t} \geq \mu_3(a, b, b, c).$$

Proof. It suffices to show that u_3 and u_4 of $\mu_2(a, b, u_3, u_4) - \mu_3(a, b, u_3, u_4) = 0$ reaches b and c , respectively, simultaneously. To see this, we deduce from calculation (3.18) that

$$(3.21) \quad \mu_3(a, b, b, u_4) - \mu_4(a, b, b, u_4) = \frac{1}{2}(b - u_4)(a + 2b + 3u_4)$$

vanishes at $u_4 = (-a - 2b)/3 = c$. \square

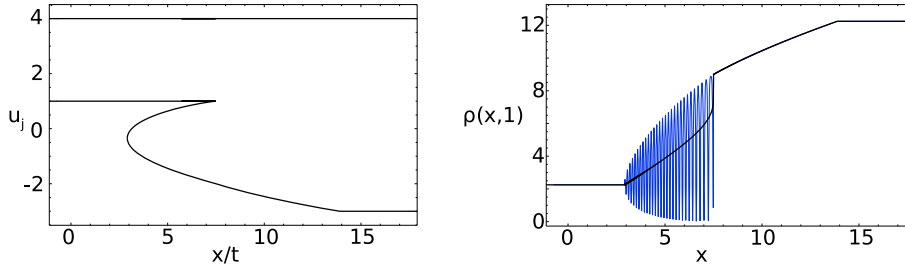


FIG. 3.4. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.1$. The initial data are given by (1.23) with $a = 4$, $b = 1$, and $c = -3$ of type IV. The solution in the region $7.5 < x/t < 111/8$ represents a rarefaction wave. The weak limit $\overline{\rho(x,t)}$ is not C^1 smooth at $x/t = 111/8$.

3.4. Type IV. Here we consider the step-like initial function (1.23) satisfying $\xi(a, b) > c$ with $0 < a + 5b < 3(b - c)$.

THEOREM 3.9 (see Figure 3.4). *For the step-like initial data (1.23) with $a > b > c$, $\xi(a, b) > c$, and $0 < a + 5b < 3(b - c)$, the solution of the Whitham equations (1.18) is given by*

$$u_1 = a, \quad u_2 = b, \quad \frac{x}{t} = \mu_2(a, b, u_3, u_4), \quad x = \mu_3(a, b, u_3, u_4) t$$

for $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, b, -(a + 2b)/3)$. Outside the region, the solution of (1.21) is divided into the three regions:

(1) For $x/t \leq \mu_3(a, b, \xi(a, b), \xi(a, b))$,

$$\alpha = a, \quad \beta = b.$$

(2) For $\mu_3(a, b, b, -(a + 2b)/3) \leq \frac{x}{t} \leq \frac{3}{8}(a^2 + 2ac + 5c^2)$,

$$\alpha = a, \quad \beta = -\frac{1}{5}a - \sqrt{\frac{8}{15} \frac{x}{t} - \frac{4}{25}a^2}.$$

(3) For $x/t \geq \frac{3}{8}(a^2 + 2ac + 5c^2)$,

$$\alpha = a, \quad \beta = c.$$

Proof. By calculation (3.21), when u_3 of $\mu_3(a, b, u_3, u_4) - \mu_4(a, b, u_3, u_4) = 0$ touches b , the corresponding u_4 reaches $-(a + 2b)/3$, which is above c . Hence, the equations

$$\frac{x}{t} = \mu_3(a, b, u_3, u_4), \quad \frac{x}{t} = \mu_4(a, b, u_3, u_4)$$

can be inverted to give u_3 and u_4 as functions of x/t in the region $\mu_3(a, b, \xi(a, b), \xi(a, b)) < x/t < \mu_3(a, b, b, -(a + 2b)/3)$. In region (2), (1.21) has a rarefaction wave solution. \square

4. More self-similar solutions. In this section, we construct self-similar solutions of the $g = 1$ Whitham equations (1.22) for the initial function (1.23) with $a > c > b$. The solution of (1.21) does not develop a shock for $a + 5b \geq 0$. We are therefore interested only in the case $a + 5b < 0$. We classify the resulting Whitham solution into four types:

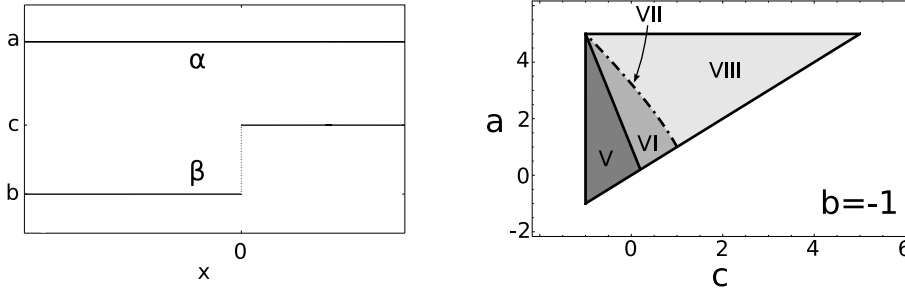


FIG. 4.1. Graph of the step-like initial shock data (1.23) with $a > c > b$ and $a + 5b < 0$, and diagram of regions for types V–VIII, where $b < 0$ has been scaled to $b = -1$.

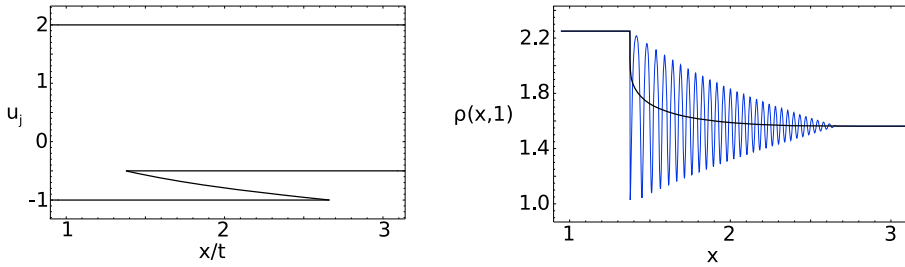


FIG. 4.2. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.014$. The initial data are given by (1.23) with $a = 2$, $b = -1$, and $c = -1/2$ of type V.

V. $a + 5b \leq -4(c - b) < 0$,
 VI. $0 > a + 5b > -4(c - b)$ with $V(a, c, b) < 0$,
 VII. $0 > a + 5b > -4(c - b)$ with $V(a, c, b) = 0$,
 VIII. $0 > a + 5b > -4(c - b)$ with $V(a, c, b) > 0$,
 where V is a quadratic polynomial given by (2.27). In Figure 4.1, we illustrate this classification by scaling $b < 0$ to $b = -1$ and plotting regions for types V–VIII in the a - c plane.

4.1. Type V. Here we consider the step-like initial function (1.23) satisfying $a + 5b \leq -4(c - b) < 0$.

THEOREM 4.1 (see Figure 4.2). *For the step-like initial data (1.23) with $a > c > b$, $a + 5b \leq -4(c - b)$, the solution of the Whitham equations (1.22) with $g = 1$ is given by*

$$u_1 = a, \quad u_2 = c, \quad \frac{x}{t} = \mu_3(a, c, u_3, b), \quad u_4 = b$$

for $\mu_3(a, c, c, b) < x/t < \mu_3(a, c, b, b)$. Outside this interval, the solution of (1.21) is given by

$$\alpha = a, \quad \beta = b \quad \text{for} \quad \frac{x}{t} \leq \mu_3(a, c, c, b)$$

and

$$\alpha = a, \quad \beta = c \quad \text{for} \quad \frac{x}{t} \geq \mu_3(a, c, b, b).$$

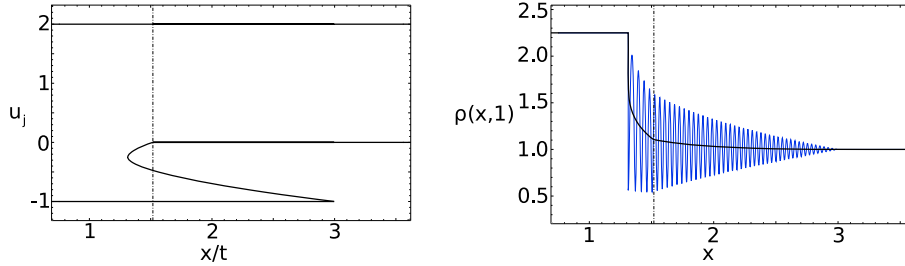


FIG. 4.3. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding periodic oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.018$. The initial data are given by (1.23) with $a = 2$, $b = -1$, and $c = 0$ of type VI. The oscillations have two different kinds of structure, which are separated by $x/t \approx 1.51$. The weak limit $\overline{\rho(x,t)}$ and the envelope of the oscillations have noticeable corners at $x/t \approx 1.51$.

Proof. It suffices to show that $\mu_3(a, c, u_3, b)$ is a decreasing function of u_3 for $b < u_3 < c$. By (2.10), we have

$$\frac{\partial \mu_3(a, c, u_3, b)}{\partial u_3} = \frac{1}{2} \frac{\partial \lambda_3}{\partial u_3} \frac{\partial q}{\partial u_3} + \frac{1}{2} [\lambda_3 - 2(a + c + u_3 + b)] \frac{\partial^2 q}{\partial u_3^2}.$$

The second term is negative because of (2.6) and $\frac{\partial^2 q}{\partial u_3^2} = 3/8 > 0$. The first term is also negative: Its first factor is positive in view of (1.17), while its second factor is

$$\frac{\partial q}{\partial u_3} = \frac{1}{4}(a + c + 3u_3 + b) < 0$$

for $b < u_3 < c$, as we have that $a + b + 4c \leq 0$. \square

4.2. Type VI. Here we consider the step-like initial function (1.23) satisfying $0 > a + 5b > -4(c - b)$ with $V(a, c, b) < 0$.

THEOREM 4.2 (see Figure 4.3). *For the step-like initial data (1.23) with $0 > a + 5b > -4(c - b)$ and $V(a, c, b) < 0$, the solution of the Whitham equations (1.22) with $g = 1$ is given by*

$$(4.1) \quad u_1 = a, \quad \frac{x}{t} = \mu_2(a, u_2, u_3, b), \quad \frac{x}{t} = \mu_3(a, u_2, u_3, b), \quad u_4 = b$$

for $\mu_3(a, -(a + b)/4, -(a + b)/4, b) < x/t \leq \mu_3(a, u^{***}, u^{***}, b)$ and by

$$(4.2) \quad u_1 = a, \quad u_2 = c, \quad \frac{x}{t} = \mu_3(a, c, u_3, b), \quad u_4 = b$$

for $\mu_3(a, u^{***}, u^{***}, b) \leq x/t < \mu_3(a, c, b, b)$, where u^{***} is the unique solution u_3 of $\mu_2(a, c, u_3, b) = \mu_3(a, c, u_3, b)$ in the interval $b < u_3 < c$. Outside the region $\mu_3(a, -(a + b)/4, -(a + b)/4, b) < x/t < \mu_3(a, c, b, b)$, the solution of (1.21) is given by

$$\alpha = a, \quad \beta = b \quad \text{for} \quad \frac{x}{t} \leq \mu_3(a, -(a + b)/4, -(a + b)/4, b)$$

and

$$\alpha = a, \quad \beta = c \quad \text{for} \quad \frac{x}{t} \geq \mu_3(a, c, b, b).$$

Proof. We first locate the “leading” edge, i.e., the solution of (4.1) at $u_2 = u_3$. Eliminating x/t from the first two equations of (4.1) yields

$$(4.3) \quad \mu_2(a, u_2, u_3, b) - \mu_3(a, u_2, u_3, b) = 0.$$

Since it degenerates at $u_2 = u_3$, we replace (4.3) by

$$(4.4) \quad G(a, u_2, u_3, b) := \frac{\mu_2(a, u_2, u_3, b) - \mu_3(a, u_2, u_3, b)}{(u_2 - u_3)\sqrt{(u_1 - u_3)(u_2 - u_4)}I(a, u_2, u_3, b)} = 0.$$

In Appendix A, we show that, at the “leading” edge $u_2 = u_3$, we have

$$G(a, u_3, u_3, b) = 2\left(\frac{\partial q}{\partial u_2} + \frac{\partial q}{\partial u_3}\right) = 0$$

in view of (A.6), which along with (2.13) gives $u_2 = u_3 = -(a + b)/4$. Having located the “leading” edge, we solve (4.4) near $u_2 = u_3 = -(a + b)/4$. We use formula (A.8) to obtain

$$\frac{\partial G(a, -(a + b)/4, -(a + b)/4, b)}{\partial u_2} = \frac{\partial G(a, -(a + b)/4, -(a + b)/4, b)}{\partial u_3} = 2.$$

These show that (4.4) gives u_2 as a decreasing function of u_3 ,

$$(4.5) \quad u_2 = B(u_3),$$

in a neighborhood of $u_2 = u_3 = -(a + b)/4$.

We now extend the solution (4.5) of (4.3) as far as possible in the region $b < u_3 < -(a + b)/4 < u_2 < c$. We use formula (2.10) to obtain

$$\begin{aligned} \frac{\partial \mu_2}{\partial u_2} &= \frac{1}{2} \frac{\partial \lambda_2}{\partial u_2} \frac{\partial q}{\partial u_2} + \frac{1}{2} [\lambda_2 - 2(a + u_2 + u_3 + b)] \frac{\partial^2 q}{\partial u_2^2}, \\ \frac{\partial \mu_3}{\partial u_3} &= \frac{1}{2} \frac{\partial \lambda_3}{\partial u_3} \frac{\partial q}{\partial u_3} + \frac{1}{2} [\lambda_3 - 2(a + u_2 + u_3 + b)] \frac{\partial^2 q}{\partial u_3^2}. \end{aligned}$$

In view of (1.17) and (2.6), we have

$$\begin{aligned} \frac{\partial \mu_2}{\partial u_2} > 0 &\quad \text{if} \quad \frac{\partial q}{\partial u_2} > 0, \\ \frac{\partial \mu_3}{\partial u_3} < 0 &\quad \text{if} \quad \frac{\partial q}{\partial u_3} < 0. \end{aligned}$$

We claim that

$$(4.6) \quad \frac{\partial q}{\partial u_2} > 0, \quad \frac{\partial q}{\partial u_3} < 0$$

on the solution of (4.3) in the region $b < u_3 < -(a + b)/4 < u_2 < c$. To see this, we use formula (2.10) to rewrite (4.3) as

$$\frac{1}{2} [\lambda_2 - 2(a + u_2 + u_3 + b)] \frac{\partial q}{\partial u_2} = \frac{1}{2} [\lambda_3 - 2(a + u_2 + u_3 + b)] \frac{\partial q}{\partial u_3}.$$

This, together with

$$\frac{\partial q}{\partial u_2} - \frac{\partial q}{\partial u_3} = 2(u_2 - u_3) \frac{\partial^2 q}{\partial u_2 \partial u_3} = \frac{1}{2}(u_2 - u_3) > 0$$

for $u_2 > u_3$ and inequalities (2.6), proves (4.6).

Hence, the solution (4.5) can be extended as long as $b < u_3 < -(a+b)/4 < u_2 < c$. There are two possibilities: (1) u_2 touches c before u_3 reaches b , and (2) u_3 touches b before or simultaneously as u_2 reaches c .

Possibility (2) is unattainable. To see this, we use (2.26) to write

$$(4.7) \quad \mu_2(a, u_2, b, b) - \mu_3(a, u_2, b, b) = \frac{(u_2 - b)}{2(a + u_2 - 2b)} V(a, u_2, b),$$

which is negative for $b < u_2 \leq c$ since $V(a, u_2, b)$ of (2.27) is an increasing function of u_2 and since $V(a, c, b) < 0$. Therefore, u_2 will touch c before u_3 reaches b . When this happens, we have

$$(4.8) \quad \mu_2(a, c, u_3, b) - \mu_3(a, c, u_3, b) = 0.$$

LEMMA 4.3. *Equation (4.8) has a simple zero, counting multiplicities, in the interval $b < u_3 < c$. If we denote this zero by u^{***} , then $\mu_2(a, c, u_3, b) - \mu_3(a, c, u_3, b)$ is positive for $u_3 > u^{***}$ and negative for $u_3 < u^{***}$.*

The proof, which involves formulae (2.24) and (2.25), is rather similar to the proof of Lemma 3.19. We will omit it.

We now continue to prove Theorem 4.2. Having solved (4.3) for u_2 as a decreasing function of u_3 for $u^{***} < u_3 < -(a+b)/4$, we can then use the middle two equations of (4.1) to determine u_2 and u_3 as functions of x/t in the interval $\mu_2(a, -(a+b)/4, -(a+b)/4, b) < x/t < \mu_2(a, c, u^{***}, b)$.

We finally turn to equations (4.2). We want to solve the third equation of (4.2), $x/t = \mu_3(a, c, u_3, b)$, for $u_3 < u^{***}$. It is enough to show that $\mu_3(a, c, u_3, b)$ is a decreasing function of u_3 for $u_3 < u^{***}$. According to Lemma 4.3, $\mu_2(a, c, u_3, b) - \mu_3(a, c, u_3, b) < 0$ for $u_3 < u^{***}$. Using formula (2.10) for μ_2 and μ_3 , we have

$$\frac{1}{2}[\lambda_2 - 2(a + c + u_3 + b)] \frac{\partial q}{\partial u_2} < \frac{1}{2}[\lambda_3 - 2(a + c + u_3 + b)] \frac{\partial q}{\partial u_3}.$$

This, together with

$$\frac{\partial q}{\partial u_2} - \frac{\partial q}{\partial u_3} = \frac{1}{2}(c - u_3) > 0$$

for $u_3 < c$ and inequalities (2.6), proves that

$$\frac{\partial q(a, c, u_3, b)}{\partial u_3} < 0$$

for $u_3 < u^{***}$. Hence,

$$\frac{\partial \mu_3}{\partial u_3} = \frac{1}{2} \frac{\partial \lambda_3}{\partial u_3} \frac{\partial q}{\partial u_3} + \frac{1}{2} [\lambda_3 - 2(a + c + u_3 + b)] \frac{\partial^2 q}{\partial u_3^2} < 0,$$

where we have used inequality (1.17). □

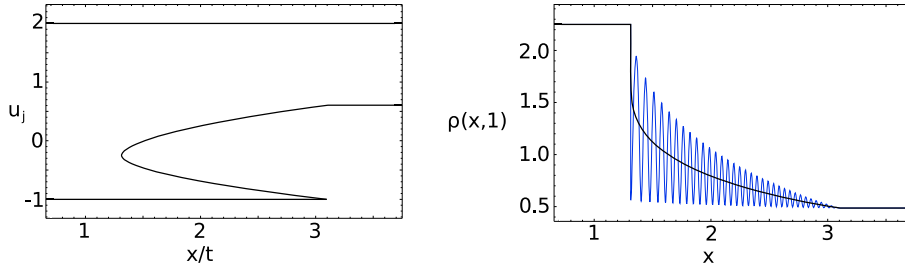


FIG. 4.4. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.03$. The initial data are given by (1.23) with $a = 2$, $b = -1$, and $c = -3 + \sqrt{13}$ of type VII.

4.3. Type VII. Here we consider the step-like initial function (1.23) satisfying $0 > a + 5b > -4(c - b)$ with $V(a, c, b) = 0$.

THEOREM 4.4 (see Figure 4.4). *For the step-like initial data (1.23) with $0 > a + 5b > -4(c - b)$ and $V(a, c, b) = 0$, the solution of the Whitham equations (1.22) with $g = 1$ is given by*

$$u_1 = a, \quad \frac{x}{t} = \mu_2(a, u_2, u_3, b), \quad \frac{x}{t} = \mu_3(a, u_2, u_3, b), \quad u_4 = c$$

for $\mu_3(a, -(a + b)/4, -(a + b)/4, b) < x/t < \mu_3(a, c, b, b)$. Outside the region, the solution of (1.21) is given by

$$\alpha = a, \quad \beta = b \quad \text{for} \quad \frac{x}{t} \leq \mu_3(a, -(a + b)/4, -(a + b)/4, b)$$

and

$$\alpha = a, \quad \beta = c \quad \text{for} \quad \frac{x}{t} \geq \mu_3(a, c, b, b).$$

Proof. It suffices to show that u_2 and u_3 of $\mu_2(a, u_2, u_3, b) - \mu_3(a, u_2, u_3, b) = 0$ reaches c and b , respectively, simultaneously. To see this, we deduce from (4.7) that

$$(4.9) \quad \mu_2(a, c, b, b) - \mu_3(a, c, b, b) = \frac{(c - b)}{2(a + c - 2b)} V(a, c, b)$$

vanishes when $V(a, c, b) = 0$. \square

4.4. Type VIII. Here we consider the step-like initial function (1.23) satisfying $0 > a + 5b > -4(c - b)$ with $V(a, c, b) > 0$.

THEOREM 4.5 (see Figure 4.5). *For the step-like initial data (1.23) with $0 > a + 5b > -4(c - b)$ and $V(a, c, b) > 0$, the solution of the Whitham equations (1.22) with $g = 1$ is given by*

$$u_1 = a, \quad \frac{x}{t} = \mu_2(a, u_2, u_3, b), \quad \frac{x}{t} = \mu_3(a, u_2, u_3, b), \quad u_4 = b$$

for $\mu_3(a, -(a + b)/4, -(a + b)/4, b) < x/t < \mu_3(a, \hat{u}, \hat{u}, b)$, where \hat{u} is the unique u_2 -zero of the quadratic polynomial $V(a, u_2, b)$ in the interval $-(a + b)/4 < u_2 < c$. Outside the region, the solution of (1.21) is divided into the following three regions:

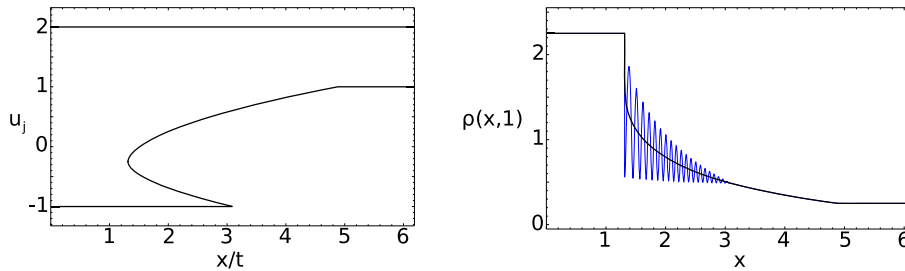


FIG. 4.5. Self-similar solution of the Whitham equations and the corresponding oscillatory solution (1.25) of the mKdV equation with $\epsilon = 0.05$. The initial data are given by (1.23) with $a = 2$, $b = -1$, and $c = 1$ of type VIII. The solution in the region $3.10 < x/t < 39/8$ represents a rarefaction wave. The weak limit $\overline{\rho(x,t)}$ is not C^1 smooth at $x/t = 39/8$.

(1) For $x/t \leq \mu_3(a, -(a+b)/4, -(a+b)/4, b)$,

$$\alpha = a, \quad \beta = b.$$

(2) For $\mu_2(a, \hat{u}, b, b) \leq x/t \leq \frac{3}{8}(a^2 + 2ac + 5c^2)$,

$$\alpha = a, \quad \beta = -\frac{1}{5}a + \sqrt{\frac{8}{15} \frac{x}{t} - \frac{4}{25}a^2}.$$

(3) For $x/t \geq \frac{3}{8}(a^2 + 2ac + 5c^2)$,

$$\alpha = a, \quad \beta = c.$$

Proof. By the calculation (4.9), when u_3 of $\mu_2(a, u_2, u_3, b) - \mu_3(a, u_2, u_3, b) = 0$ touches b , the corresponding u_2 reaches \hat{u} , where $V(a, \hat{u}, b) = 0$. Obviously, $\hat{u} < c$. Hence, equations

$$\frac{x}{t} = \mu_2(a, u_2, u_3, b), \quad \frac{x}{t} = \mu_3(a, u_2, u_3, b)$$

can be inverted to give u_2 and u_3 as functions of x/t in the region $\mu_2(a, -(a+b)/4, -(a+b)/4, b) < x/t < \mu_2(a, \hat{u}, \hat{u}, b)$. To the right of this region, equations (1.21) have a rarefaction wave solution. \square

5. Vacuum points. The nonnegative function $\rho(x, t; \epsilon)$ has to be positive in order for the mKdV dispersive approximation (1.19) (or (1.2) for NLS) to make sense. It is therefore interesting to study those points (x, t) at which $\rho(x, t; \epsilon) = 0$ for a sequence of vanishing ϵ . These points are referred to as the vacuum points in [3].

Since the solution $\rho(x, t; \epsilon)$ of (1.19) can be approximated by the periodic solution (1.25) in the single phase regime, we instead study the vacuum points of the latter solution.

Since $\rho_2 \geq \rho_3$ because of (1.9) for $u_1 \geq u_2 \geq u_3 \geq u_4$, it follows from formula (1.25) that (x, t) is a vacuum point if and only if $\rho_3(x, t) = 0$, which, in view of (1.9), is equivalent to

$$(5.1) \quad u_1(x, t) - u_2(x, t) - u_3(x, t) + u_4(x, t) = 0.$$

For types I–IV, equality (5.1) can occur only if $u_1 - u_2 \leq u_3 - u_4$ at the leading edge of the Whitham solution (see Figures 1.2, 3.2, 3.3, and 3.4). This inequality leads to $a + b \leq 2b$ for types I and II and gives $a \leq 4b$ for types III and IV. Hence,

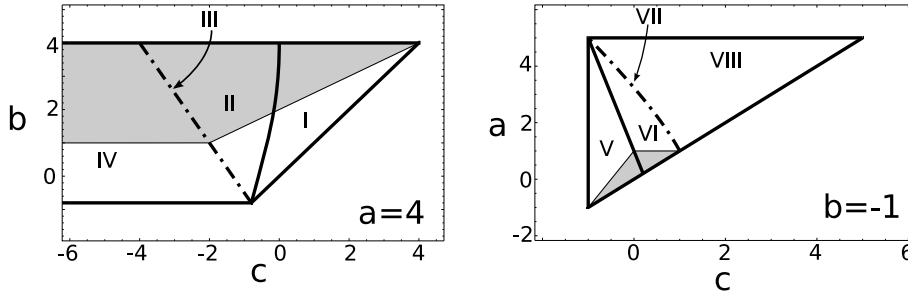


FIG. 5.1. Diagram of regions for the initial data (1.23) which generate oscillatory solutions with vacuum points. The parameters (c, b) or (c, a) form the shaded region.

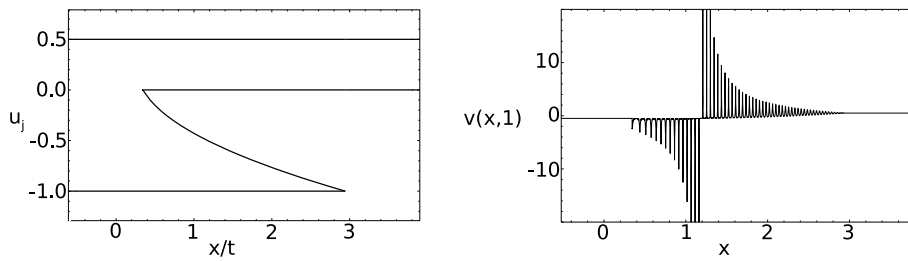


FIG. 5.2. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the flow velocity (1.10) of the oscillatory solution of the mKdV equation with $\epsilon = 0.014$. The initial data are given by (1.24) of type V with parameters $a = 0.5$, $b = -1$, and $c = 0$. A vacuum point is located at $x/t \approx 1.16$.

step-like initial shock data of types I–IV with $4b \geq a > b > c$ and $a + c \leq 2b$ will generate mKdV oscillatory solutions with vacuum points (cf. Figure 5.1).

Similarly, step-like initial shock data of types V and VI with $-b \geq a > c > b$ and $a + b \leq 2c$ will produce mKdV oscillatory solutions with vacuum points. Types VII and VIII oscillatory solutions cannot have any vacuum point (cf. Figure 5.1).

Types I–V oscillatory solutions have at most one vacuum point. To see this, note that the corresponding Whitham solutions have the property that $u_1(x, t)$ and $u_2(x, t)$ are constants and that $u_3(x, t) - u_4(x, t)$ are strictly monotone functions of x/t in the single phase region; so is the left-hand side of (5.1). Hence, (5.1) has at most one solution x/t . Figure 5.2 shows an example of oscillatory solutions with a unique vacuum point. Notice that the velocity $v(x, t; \epsilon)$ has a huge change at this point.

Type VI oscillatory solutions can have more than one vacuum point. They can even have a continuum of vacuum points. Namely, there exists an interval such that (5.1) holds at each point of the interval.

THEOREM 5.1. *For step-like initial data (1.24) of type VI satisfying condition $a + b = 0$, all the points on the closed interval $\mu_3(a, -(a + b)/4, -(a + b)/4, b) \leq x/t \leq \mu_3(a, u^{***}, u^{***}, b)$, where u^{***} is given in Theorem 4.2, are vacuum points.*

Proof. According to Theorem 4.2, the Whitham solution has the property that $u_1 = a$ and $u_4 = b$. Since $a + b = 0$, it follows from (5.1) that it suffices to show that $u_2 = -u_3$ over the closed interval mentioned in Theorem 5.1.

The functions u_2 and u_3 of the Whitham solution satisfies (4.3), which has a unique solution (4.5) with the property that $B(0) = 0$ in the case of $a + b = 0$. It

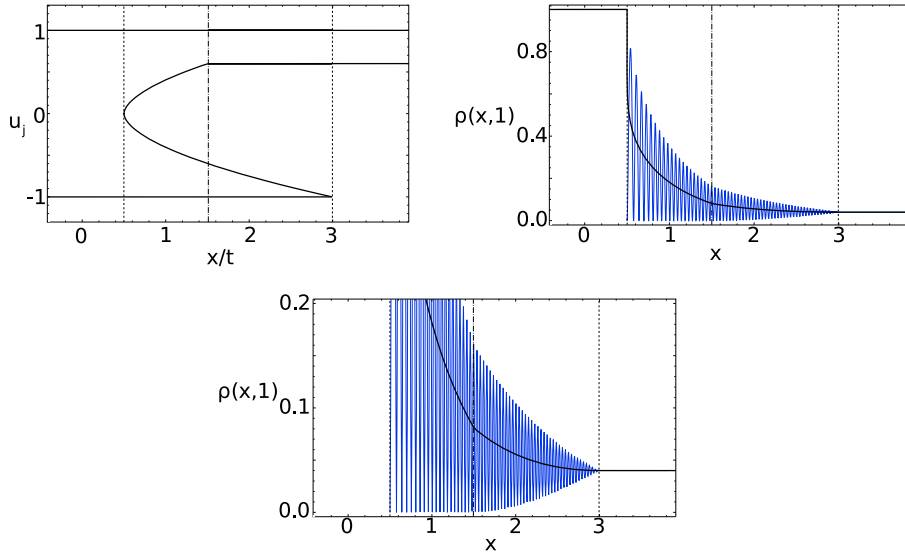


FIG. 5.3. Self-similar solution of the Whitham equations (1.22) with $g = 1$ and the oscillatory solutions (1.25) of the mKdV equation with $\epsilon = 0.04$. The initial data are given by (1.23) with $a = 1$, $b = -1$, and $c = 0.6$ of type VI. The vacuum points occupy the whole interval $0.5 \leq \frac{x}{t} \leq 1.50$.

follows from formula (2.10) for μ_2 and μ_3 , formula (2.13) for q , and formulae (A.2–A.3) for λ_2 and λ_3 that $\mu_2(a, -u_3, u_3, b) = \mu_3(a, -u_3, u_3, b)$. Hence, $u_2 = -u_3$ is the unique solution of (4.3). \square

Figure 5.3 is an illustration of such a solution, where rather than a vacuum point we now have a vacuum interval. It follows from formula (1.10) that $V(x, t; \epsilon) = (u_1 + u_2 + u_3 + u_4)/4 = 0$ at each point of the interval.

We close this section with a remark on the NLS case, which has been used to describe the nonlinear pulse propagation in an optical fiber [7]. The NLS oscillations with step-like initial shock data (1.15), $b > c$, have at most one vacuum point [3]. For more complicated step-like initial shock data, the NLS oscillations can also have a continuum of vacuum points. Indeed, consider the initial data

$$\alpha(x, 0) = \begin{cases} a, & x < 0, \\ b, & x > 0, \end{cases} \quad \beta(x, 0) = \begin{cases} c, & x < 0, \\ d, & x > 0, \end{cases}$$

with $a > c > b > d$ and $a - b = c - d$. It can easily be shown that there is an interval within the single phase regime where all the u_i of the Whitham solution of (1.7) with $g = 1$ are constants. More precisely, $u_1 = a$, $u_2 = c$, $u_3 = b$, and $u_4 = d$ for $\lambda_3(a, c, b, d) \leq x/t \leq \lambda_2(a, c, b, d)$, where λ_2 and λ_3 are the second and third eigenspeeds of the NLS–Whitham equations (1.7) for $g = 1$. Since $a - b = c - d$, (5.1) holds at each point of the interval; hence, this is a vacuum interval. From the point of view of nonlinear optics, the optical wave has trivial frequency chirp in this region, and the region may be considered as a vacuum.

6. Other initial data. We conclude the paper by showing how to handle the initial data (1.24). Inequalities (2.14) are replaced by

$$\frac{\partial \lambda_2}{\partial u_2} < \frac{3}{2} \frac{\lambda_1 - \lambda_2}{u_1 - u_2} < \frac{\partial \lambda_1}{\partial u_1}$$

for $u_4 < u_3 < u_2 < u_1$. Results similar to those of Lemma 2.2 can then be easily proved. The rest of the calculations are similar to those in sections 3 and 4.

Appendix A. Leading edge calculations. The function $I(u_1, u_2, u_3, u_4)$ of (2.3) can be written in terms of the complete elliptic integral of the first kind $K(s)$; i.e.,

$$(A.1) \quad I = \frac{K(s)}{\sqrt{(u_1 - u_3)(u_2 - u_4)}},$$

where

$$s = \frac{(u_1 - u_2)(u_3 - u_4)}{(u_1 - u_3)(u_2 - u_4)}.$$

Using the derivative formula

$$\frac{dK(s)}{ds} = \frac{E(s) - (1 - s)K(s)}{2s(1 - s)},$$

where $E(s)$ is the complete elliptic integral of the second kind, we calculate λ_2 and λ_3 of (2.2):

$$(A.2) \quad \lambda_2 = 2\sigma_1 - 4 \frac{u_1 - u_2}{1 - \frac{u_1 - u_3}{u_2 - u_3} \frac{E}{K}},$$

$$(A.3) \quad \lambda_3 = 2\sigma_1 + 4 \frac{u_3 - u_4}{1 - \frac{u_2 - u_4}{u_2 - u_3} \frac{E}{K}}.$$

We then use (2.10) to write $\mu_2 - \mu_3$ as

$$-2(u_2 - u_3)K(s) \left[\frac{u_1 - u_2}{(u_2 - u_3)K - (u_1 - u_3)E} \frac{\partial q}{\partial u_2} + \frac{u_3 - u_4}{(u_2 - u_3)K - (u_2 - u_4)E} \frac{\partial q}{\partial u_3} \right].$$

Hence, $G(u_1, u_2, u_3, u_4)$ of (4.4) becomes

$$(A.4) \quad G = -2 \left[\frac{u_1 - u_2}{(u_2 - u_3)K - (u_1 - u_3)E} \frac{\partial q}{\partial u_2} + \frac{u_3 - u_4}{(u_2 - u_3)K - (u_2 - u_4)E} \frac{\partial q}{\partial u_3} \right].$$

We now use the asymptotics of $K(s)$ and $E(s)$ as s is close to 1,

$$(A.5) \quad K(s) \approx \frac{1}{2} \log \frac{16}{1 - s}, \quad E(s) \approx 1 + \frac{1}{4}(1 - s) \left(\log \frac{16}{1 - s} - 1 \right),$$

to calculate the $u_2 = u_3$ limit

$$(A.6) \quad G = 2 \left(\frac{\partial q}{\partial u_2} + \frac{\partial q}{\partial u_3} \right).$$

Finally, we can also use the expression (2.13) and the derivative formulae

$$(A.7) \quad \frac{dE(s)}{ds} = \frac{E(s) - K(s)}{2s}, \quad \frac{dK(s)}{ds} = \frac{E(s) - (1 - s)K(s)}{2s(1 - s)}$$

to evaluate the partial derivatives of G in the $u_2 = u_3$ limit

$$(A.8) \quad \begin{aligned} \frac{\partial G}{\partial u_2} \Big|_{u_2=u_3} &= 2 + \frac{(u_1 - u_4)(u_1 + 4u_3 + u_4)}{4(u_1 - u_3)(u_4 - u_3)}, \\ \frac{\partial G}{\partial u_3} \Big|_{u_2=u_3} &= 2 - \frac{(u_1 - u_4)(u_1 + 4u_3 + u_4)}{4(u_1 - u_3)(u_4 - u_3)}. \end{aligned}$$

Appendix B. Loss of genuine nonlinearity. In this appendix, we will show that the Whitham equations (1.22) with $g = 1$ are not genuinely nonlinear. Again, we suppress the subscript $g = 1$ in the notation $\lambda_{g,i}$ and $\mu_{g,i}$.

PROPOSITION B.1. *For fixed $u_1 > u_2 > u_4$ satisfying conditions*

$$u_1 + 4u_2 + u_4 > 0, \quad U(u_1, u_2, u_4) < 0,$$

where the function U is given in (2.23), the derivative

$$\frac{\partial \mu_3(u_1, u_2, u_3, u_4)}{\partial u_3}$$

changes sign as u_3 increases from u_4 to u_2 .

Proof. We first use formula (A.3) for λ_3 , derivative formulae (A.7), and asymptotics (A.5) to calculate

$$(B.1) \quad \frac{\partial \lambda_3}{\partial u_3} \Big|_{u_3=u_2} = +\infty,$$

$$(B.2) \quad \frac{\partial \lambda_3}{\partial u_3} \Big|_{u_3=u_4} = \frac{24(u_2 - u_4)^2 + 24(u_1 - u_2)(u_2 - u_4) + 9(u_1 - u_2)^2}{(u_1 + u_2 - 2u_4)^2}.$$

In the calculation of the last equation, we have also used the asymptotics of $K(s)$ and $E(s)$ near $s = 0$,

$$\begin{aligned} K(s) &= \frac{\pi}{2} \left[1 + \frac{s}{4} + \frac{9}{64}s^2 + \dots + \left(\frac{1 \cdot 3 \cdots (2n-1)}{2 \cdot 4 \cdots 2n} \right)^2 s^n + \dots \right], \\ E(s) &= \frac{\pi}{2} \left[1 - \frac{s}{4} - \frac{3}{64}s^2 - \dots - \frac{1}{2n-1} \left(\frac{1 \cdot 3 \cdots (2n-1)}{2 \cdot 4 \cdots 2n} \right)^2 s^n - \dots \right]. \end{aligned}$$

Differentiating formula (2.10) for μ_3 and using formula (2.13) for q , we obtain

$$\frac{\partial \mu_3(u_1, u_2, u_3, u_4)}{\partial u_3} = \frac{1}{8}(u_1 + u_2 + 3u_3 + u_4) \frac{\partial \lambda_3}{\partial u_3} + \frac{3}{8}[\lambda_3 - 2\sigma_1].$$

We then use the boundary values (2.7) and (2.8) for λ_3 and boundary values (B.1) and (B.2) for $\partial_{u_3} \lambda_3$ to calculate

$$\frac{\partial \mu_3}{\partial u_3} \Big|_{u_3=u_2} = +\infty \quad \text{if } u_1 + 4u_2 + u_4 > 0,$$

and

$$\frac{\partial \mu_3}{\partial u_3} \Big|_{u_3=u_4} = \frac{3U(u_1, u_2, u_4)}{8(u_1 + u_2 - 2u_4)} < 0 \quad \text{if } U(u_1, u_2, u_4) < 0.$$

These two boundary values for $\partial_{u_3}\mu_3$ prove the proposition. \square

The conditions of the proposition can easily be verified. For example, for the initial data (1.23) of types II and III, we have $\xi(a, b) > c$ and $a + 5b \geq 3(b - c) > 0$, which imply $U(a, b, c) < 0$ and $a + 4b + c > 0$. By Proposition B.1, the derivative

$$\frac{\partial\mu_3(a, b, u_3, c)}{\partial u_3}$$

must change sign as u_3 increases from c to b .

Appendix C. Phase shift. In this appendix, we will verify that Q given by (1.12) and (1.13) is the right phase shift for both the NLS and mKdV. We will study only the NLS case; the mKdV case can be handled in the same way.

For the NLS periodic wave train (1.8), the wave number and frequency are

$$k = \frac{\pi\sqrt{\rho_1 - \rho_3}}{\epsilon K(s)}, \quad w = kV_1.$$

The rapid phase is then given by

$$(C.1) \quad \Theta(x, t; \epsilon) = k\theta(x, t) = k[x - V_1 t - Q].$$

To verify that Q of (1.12) and (1.13) is indeed the phase shift, it suffices to show that Θ satisfies the generalized wave number and frequency relations

$$(C.2) \quad \frac{\partial\Theta}{\partial x} = k, \quad \frac{\partial\Theta}{\partial t} = -w.$$

We first observe that the compatibility condition for (C.2) is

$$\frac{\partial k}{\partial t} + \frac{\partial w}{\partial x} = 0,$$

which is the conservation of waves. This equation can be viewed as an additional conservation law satisfied by the solutions of the Whitham equations (1.7). Hence, the eigenspeeds λ_i of the Whitham equations can be calculated using k and w ; i.e.,

$$(C.3) \quad \lambda_i = \frac{\partial_{u_i} w}{\partial_{u_i} k}.$$

We can also rewrite the wave number k in terms of I of (A.1) as

$$k = \frac{\pi}{\epsilon I},$$

which along with formula (2.2) for λ_i gives

$$(C.4) \quad \frac{k}{\partial_{u_i} k} = -\frac{I}{\partial_{u_i} I} = \frac{1}{2}(\lambda_i - 2\sigma_1).$$

Differentiating (C.1) with respect to x and using the property that $u_1 = a$, we obtain

$$(C.5) \quad \begin{aligned} \frac{\partial\Theta}{\partial x} &= k + \sum_{i=2}^4 \left[x \frac{\partial k}{\partial u_i} - t \frac{\partial w}{\partial u_i} - k \frac{\partial Q}{\partial u_i} - Q \frac{\partial k}{\partial u_i} \right] u_{ix} \\ &= k + \sum_{i=2}^4 \frac{\partial k}{\partial u_i} u_{ix} \left[x - \lambda_i t - \frac{1}{2}(\lambda_i - 2\sigma_1) \frac{\partial Q}{\partial u_i} - Q \right], \end{aligned}$$

where we have used formulae (C.3) and (C.4) in the last equality.

It is known that if $Q(a, u_2, u_3, u_4)$ satisfies (1.12) and (1.13), then

$$x = \lambda_i t + \frac{1}{2}(\lambda_i - 2\sigma_1) \frac{\partial Q}{\partial u_i} + Q$$

is the hodograph solution of the Whitham equations (1.7) [14]. Hence, we deduce from (C.5) the first equation of (C.2). The second equation of (C.2) can be shown in the same way.

Acknowledgment. We thank the referees for many valuable suggestions.

REFERENCES

- [1] P. DEIFT, S. VENAKIDES, AND X. ZHOU, *New results in small dispersion KdV by an extension of the steepest descent method for Riemann-Hilbert problems*, Internat. Math. Res. Not., No. 6 (1997), pp. 285–299.
- [2] B. A. DUBROVIN AND S. P. NOVIKOV, *Hydrodynamics of weakly deformed soliton lattices. Differential geometry and Hamiltonian theory*, Russian Math. Surveys, 44 (1989), pp. 35–124.
- [3] G. A. EL, V. V. GEOGJAEV, A. V. GUREVICH, AND A. L. KRYLOV, *Decay of an initial discontinuity in the defocusing NLS hydrodynamics*, Phys. D, 87 (1995), pp. 186–192.
- [4] M. G. FOREST AND J.-E. LEE, *Geometry and modulation theory for the periodic nonlinear Schrödinger equation*, in Oscillation Theory, Computation, and Methods of Compensated Compactness, J. L. Ericksen, D. Kinderlehrer, and M. Slemrod, eds., Springer-Verlag, New York, 1986, pp. 35–69.
- [5] T. GRAVA AND C. KLEIN, *Numerical solution of the small dispersion limit of Korteweg-de Vries and Whitham equations*, Comm. Pure Appl. Math., 60 (2007), pp. 1623–1664.
- [6] S. JIN, C. D. LEVERMORE, AND D. W. MCLAUGHLIN, *The semiclassical limit of the defocusing NLS hierarchy*, Comm. Pure Appl. Math., 52 (1999), pp. 613–654.
- [7] Y. KODAMA, *The Whitham equations for optical communications: Mathematical theory of NRZ*, SIAM J. Appl. Math., 59 (1999), pp. 2162–2192.
- [8] P. D. LAX, C. D. LEVERMORE, AND S. VENAKIDES, *The generation and propagation of oscillations in dispersive initial value problems and their limiting behavior*, in Important Developments in Soliton Theory 1980–1990, T. Fokas and V. E. Zakharov, eds., Springer Ser. Nonlinear Dynam., Springer-Verlag, Berlin, 1993, pp. 205–241.
- [9] V. B. MATVEEV, *Abelian Functions and Solitons*, Preprint 373, University of Wrocław, Wrocław, Poland, 1976.
- [10] M. V. PAVLOV, *Nonlinear Schrödinger equation and the Bogolyubov-Whitham averaging method*, Theoret. and Math. Phys., 71 (1987), pp. 584–588.
- [11] V. U. PIERCE AND F. R. TIAN, *Self-similar solutions of the non-strictly hyperbolic Whitham equations*, Commun. Math. Sci., 4 (2006), pp. 799–822.
- [12] V. U. PIERCE AND F. R. TIAN, *Self-similar solutions of the non-strictly hyperbolic Whitham equations for the KdV hierarchy*, Dyn. Partial Differ. Equ., 4 (2007), pp. 263–282.
- [13] F. R. TIAN, *The Whitham-type equations and linear overdetermined systems of Euler-Poisson-Darboux type*, Duke Math. J., 74 (1994), pp. 203–221.
- [14] F. R. TIAN AND J. YE, *On the Whitham equations for the semiclassical limit of the defocusing nonlinear Schrödinger equation*, Comm. Pure Appl. Math., 52 (1999), pp. 655–692.
- [15] S. P. TSAREV, *Poisson brackets and one-dimensional Hamiltonian systems of hydrodynamic type*, Soviet Math. Dokl., 31 (1985), pp. 488–491.
- [16] S. VENAKIDES, *Higher order Lax-Levermore theory*, Comm. Pure Appl. Math., 43 (1990), pp. 335–362.

PERSISTENCE OF ROLL WAVES FOR THE SAINT VENANT EQUATIONS*

PASCAL NOBLE†

Abstract. The purpose of the article is to study the linear and nonlinear “stability” of roll-waves that are periodic and discontinuous entropic travelling wave solutions of the Saint Venant equations. More precisely, we prove that the Cauchy problem with initial data close to a roll-wave and satisfying suitable compatibility conditions has a solution on a sufficiently small interval.

Key words. roll-waves, shallow water, Cauchy problem, hyperbolic equations

AMS subject classifications. 35L45, 35L67, 35B10, 76H05, 35Q35

DOI. 10.1137/07070810X

1. Introduction. Roll-waves are well-known nonlinear patterns appearing in shallow waters under the effect of gravity and bottom friction. This type of flow is commonly modeled by the “shallow water equations” that can be formally derived from the Navier–Stokes system. This is a hyperbolic system for the fluid height h and the average velocity u that is similar to the isentropic Euler equations. An empiric friction term, so-called Chezy friction term, is added to the momentum equation to model the friction of the bottom; see [7] for more details on the derivation of such models. In [3], Dressler proved the existence of periodic travelling waves that are mathematical solutions of the shallow water equations: those travelling waves are discontinuous, the discontinuity being a Lax shock. Similarly to shocks in hyperbolic systems, several questions arise at this stage: first, the existence of continuous roll-wave solutions of a viscous perturbation of the shallow water equations, the convergence to the inviscid roll-waves in the vanishing viscosity limit and their stability. The existence of continuous travelling waves for viscous shallow water equations follows from a classical Hopf bifurcation argument [9]. More involved is the question of the convergence of those solutions in the vanishing viscosity limit: one can prove, using geometric singular perturbation arguments (Fenichel theorems) that continuous roll-waves are ε close to an inviscid roll-wave with spatial period T and wave speed c_0 for a suitable wave speed $c(\varepsilon, T)$ with ε , the size of the viscosity [10], [11]. Only partial results are known concerning the stability of viscous roll-waves: those patterns are proved to be spectrally stable under large wavelength and small wavelength perturbations and are linearly stable provided that they are strongly spectrally stable; see [12] for more details.

In order to obtain more information on stability of viscous roll-waves, at least in the vanishing viscosity limit, we shall consider the inviscid case. In the case of a shock wave, the connection between the viscous and inviscid shocks has been established by Rousset [15] in the one-dimensional ($1D$) case and Gues and coworkers [6] in the multidimensional case. In the hyperbolic setting, the question of the “stability” of inviscid roll-waves then arises either under $1D$ or multidimensional perturbations.

*Received by the editors November 13, 2007; accepted for publication (in revised form) September 13, 2008; published electronically January 7, 2009.

<http://www.siam.org/journals/sima/40-5/70810.html>

†Université de Lyon, Université Lyon 1, UMR CNRS 5208, Institut Camille Jordan, Batiment du Doyen Jean Braconnier, 43, blvd du 11 novembre 1918, F - 69622 Villeurbanne Cedex, France (noble@math.univ-lyon1.fr).

In that paper and in order to make the connection with the viscous case, we shall focus on small *smooth* perturbations of roll-waves. Due to the presence of an infinite number of shocks, even the formulation of the stability problem is a hard task. In the case of a single shock, one formulates the stability problem among the piecewise smooth functions with a single discontinuity that is close to the shock: working in the reference frame attached to that discontinuity, one can reduce the stability problem to the analysis of an initial boundary value problem. In the case of roll-waves, it is natural to work among the functions that are piecewise regular with discontinuities that are close to the discontinuities of the roll-waves. Here it is not sufficient to work in the reference frame attached to *one* discontinuity; in that case, this would mean that we restrict our attention to *periodic* perturbations of the roll-waves. In that setting, Tamada and Tougou [17] proved that, for suitable large wavelength, roll-waves are spectrally stable. In order to handle more general perturbations, we need to fix *all* the discontinuities; this is done through a Lipschitz change of variable. The author proved in that case that for large wavelength, the spectral problem has no unstable eigenvalue [13]. Due to the hyperbolic nature of the problem, it seems hard to obtain a stability result in that case; we shall consider here the question of the “persistence” of roll-waves.

Similarly to the persistence of shocks, we shall prove here that the Cauchy problem with initial data that are close to a roll-wave and satisfy suitable compatibility conditions is well posed. First, we formulate the “stability” problem among the space of piecewise regular functions with discontinuities close to the roll-waves discontinuities: the shocks are fixed through a Lipschitz change of variable, and we write both the shallow water equations and the Rankine–Hugoniot jump conditions in that setting. Then we linearize that set of equations and analyze the spectral problem. In order to tackle the Cauchy problem, we shall only consider the high-frequency perturbations; we briefly analyze an analogous equation of the Lopatinskii determinant for shocks. Then, following the approach for shocks, we derive a priori estimates on the linear problem and for a linearized problem around an approximate solution. For an initial data satisfying suitable compatibility conditions, we can construct an approximate solution of the shallow water equations on a sufficiently small time interval. We obtain the existence of a solution to the Cauchy problem through a fixed point argument around the approximate solution with the a priori estimates on the linearized problem.

This method introduced by Majda [4] for compact shocks and generalized by Métivier [5] is suitable for *multidimensional perturbations*. Though *valid in the multidimensional case*, the principal results obtained in this paper are written in *the 1D setting*. Indeed, in order to simplify the presentation, we have chosen to emphasize the periodic nature of the problem and the new issues that have to be handled with, in particular, the formulation of the stability problem and the presence of a sonic point in the interior that complicates higher-order estimates. In this particular setting, the energy estimates are derived easily without the microlocal analysis. We briefly discuss at the end of the paper how to extend the analysis proposed here to multidimensional perturbations. Moreover, this approach shall be generalized and developed in a forthcoming paper to handle roll-waves in generalized hyperbolic systems obtained recently by the author in [14].

The paper is organized as follows: in section 2, we recall the Dressler construction of roll-waves and formulate the boundary value problem associated to the persistence problem. In section 3, we linearize the equations around a roll-wave and perform the spectral analysis of that linear problem. We also obtain a priori estimates for

the solutions of that problem. We then consider the linearized problem around an approximate solution: we first obtain a priori estimates on solutions of that linear problem and deduce the well posedness of those equations in suitable functional spaces. Section 4 is devoted to the analysis of the full nonlinear problem: we first obtain compatibility conditions for initial data and construct an approximate solution of the shallow water equations coupled with Rankine–Hugoniot jump conditions. The well posedness of the Cauchy problem is then obtained through a fixed point argument.

2. Formulation of the problem. In this section, we give a short proof of the existence of inviscid roll-waves. Then, we introduce the functional spaces adapted to the stability analysis of that kind of pattern: these are piecewise regular functions with discontinuities close to the discontinuities of a roll-wave. We then perform the Lipschitz change of variable that fixes *all* the shocks and write the shallow water equations and the Rankine–Hugoniot conditions in the reference frame.

2.1. Existence of roll-waves. We briefly recall Dressler analysis for the existence of roll-wave solutions of the Saint Venant system. Let us start with the nondimensional Saint Venant equations

$$(1) \quad \begin{aligned} h_t + (hu)_x &= 0, \\ (hu)_t + \left(\frac{h^2}{2F} + hu^2\right)_x &= h - u^2, \end{aligned}$$

where h is the height of the fluid, u is the average fluid speed, and F is the Froude number. In [3], Dressler looks for periodic travelling wave solutions in the form

$$(h, u)(x, t) = (H, U)(x - ct)$$

of the inviscid Saint Venant equations: these yield a first order differential system. Eliminating U from the system reduces the problem to finding periodic solutions of the scalar equation

$$(2) \quad H' = \frac{H - (c - \frac{q}{H})^2}{\frac{H}{F} - \frac{q^2}{H^2}} = P_1(H),$$

where $q = H(c - U) > 0$, a constant, is the relative discharge rate. Integrating (2) in the form $\xi = \xi_0 + \int \frac{dh}{P_1(h)}$, it is proved that (2) has no continuous periodic solution [3]. Thus, we are looking for periodic solutions with discontinuities that satisfy admissibility conditions. The hyperbolic part of the system being similar to the isentropic Euler equations and the admissibility conditions for a discontinuity located at $x = x(t)$ are the Rankine–Hugoniot conditions

$$(3) \quad [hu] = \dot{x}[h], \quad \left[\frac{h^2}{2F} + hu^2\right] = \dot{x}[hu],$$

where $[u] = \lim_{\epsilon \rightarrow 0} u(x(t) + \epsilon) - u(x(t) - \epsilon)$ and a Lax shock condition

$$u_+ + \sqrt{\frac{h_+}{F}} < \dot{x} < u_- + \sqrt{\frac{h_-}{F}}.$$

Dressler proved the following result.

THEOREM 1. *Given any $F > 4$, $L > 0$, and $c > 0$, there exists a piecewise C^1 periodic travelling wave, with wave speed c and wavelength L , entropic solution of the inviscid Saint Venant equations (1).*

We recall the proof for completeness.

Proof. Let H_- (resp. H_+) be the fluid height before (resp. after) a shock. The entropic shock must satisfy the Rankine–Hugoniot jump conditions

$$(4) \quad \frac{H_+^2}{2F} + \frac{q^2}{H_+} = \frac{H_-^2}{2F} + \frac{q^2}{H_-}.$$

The Saint Venant equations are similar to the isentropic Euler equations for an ideal gas with $\gamma = 2$, where h plays the role of the density ρ . The entropy condition is equivalent to the Lax shock condition (see [18] for more details):

$$U_+ + \sqrt{\frac{H_+}{F}} < c < U_- + \sqrt{\frac{H_-}{F}}.$$

Since the relative discharge rate $H(c - U) = q$ is a constant, we obtain

$$\sqrt{\frac{H_+}{F}} - \frac{q}{H_+} < 0 < \sqrt{\frac{H_-}{F}} - \frac{q}{H_-}.$$

Hence, $H_+ < H_-$ and there exists a “sonic” point H_0 so that $H_+ < H_0 < H_-$ and

$$\frac{H_0}{F} - \frac{q^2}{H_0^2} = 0.$$

In order to satisfy the entropy condition, we must have $\frac{dH}{d\xi}|_{H_0} > 0$ or equivalently, $F > 4$. The denominator of the fraction in (2) vanishes at point H_0 . To pass continuously through this value, the numerator must also vanish at this point H_0 :

$$H_0 - \left(c - \frac{q}{H_0}\right)^2 = 0.$$

We simplify the fraction in (2):

$$(5) \quad \frac{dH}{d\xi} = F \frac{H^2 + (H_0 - c^2)H + \frac{H_0^2}{F}}{H^2 + H_0H + H_0^2} = P_1(H).$$

The construction of a periodic solution with an arbitrary wavelength L is, at this point, straightforward. Let $H(\xi)$ be the special solution of (5) so that $H(0) = H_0$: it is defined implicitly by $H(\xi) = h \Leftrightarrow \xi = f_1(h)$, where f_1 is the primitive of $\frac{1}{P_1}$ such that $f_1(H_0) = 0$. The function f_1 is given by

$$Ff_1(h) = h - H_0 + \frac{H_a^2 + H_0H_a + H_0^2}{H_a - H_b} \ln\left(\frac{h - H_a}{H_0 - H_a}\right) - \frac{H_b^2 + H_0H_b + H_0^2}{H_a - H_b} \ln\left(\frac{h - H_b}{H_0 - H_b}\right),$$

and $H_a > H_b$ are the zeros of P_1 . Define $H_n(\xi) = H(\xi - nL)$. We fit H_n and H_{n+1} together by means of an entropic shock. Eliminating the solution $H_+ = H_-$, the Rankine–Hugoniot condition (4) now reads

$$(6) \quad H_+ + H_- = \frac{2H_0^3}{H_+H_-}.$$

We determine the position of the n th shock ξ_c^n , with $H_n(\xi_c^n) = H_-$ and $H_{n+1}(\xi_c^n) = H_+$. This is equivalent to the relation

$$(7) \quad \int_{H_0}^{H^-} \frac{dh}{P_1(h)} = L + \int_{H_0}^{H^+} \frac{dh}{P_1(h)}.$$

Now, eliminating H_+ between (6) and (7) yields L as a function of H_- . Then we have obtained periodic and entropic solutions of the Saint Venant equations (1). \square

In what follows, we rescale space variables and unknowns by H_0 in the variable $H := H_0 h$ and $\xi := H_0 \xi$: we get the nondimensional equation of profile

$$h' = F \frac{h^2 + \left(1 - \left(1 + \frac{1}{\sqrt{F}}\right)^2\right) h + \frac{1}{F}}{h^2 + h + 1} = P(h).$$

Denote $H_+ = H_0 h_+$ and $H_- = H_0 h_-$. The nondimensional Rankine–Hugoniot conditions reads

$$h_+ + h_- = \frac{2}{h_+ h_-}.$$

We can easily describe the asymptotic profile of the roll-waves as $L \rightarrow \infty$. Indeed, when $L \rightarrow \infty$, the minimum height h_+ converges to h_a given by

$$h_a = \frac{H_a}{H_0} = \frac{1}{2F} \left(1 + 2\sqrt{F} + \sqrt{1 + 4\sqrt{F}}\right).$$

Then the family of roll-waves converges to a solitary wave with a single shock as $L \rightarrow \infty$. Denote $h_m > h_a$ so that $h_a + h_m = \frac{2}{h_a h_m}$. Before the shock, the profile \tilde{h} of the solitary wave is given by the special solution $\tilde{h} = h$, $h \in (h_a, h_m)$, and after the shock, it is a constant $\tilde{h} = h_a$. The description of the roll-wave solutions of the Saint Venant system is complete, and we are now in a position to formulate the stability problem.

2.2. Formulation of the stability problem. In this section, we recall the method, introduced in [13], of studying the stability of Dressler roll-waves that possess an infinite number of shocks. These roll-waves are parametrized by the wavelength L and the wave speed c . When these parameters are fixed, denote H_+ (resp. H_-) the minimum (resp. maximum) height of the roll-wave. This is also the height after (resp. before) a shock. We first write the Saint Venant system (1) into a diagonalized form and introduce the Riemann invariants

$$r = u + 2\sqrt{\frac{h}{F}}, \quad s = u - 2\sqrt{\frac{h}{F}}.$$

The shallow water equations (1) read

$$(8) \quad r_t + \lambda_1(r, s)r_x = Q(r, s), \quad s_t + \lambda_2(r, s)s_x = Q(r, s),$$

where $\lambda_k, k = 1, 2$ and Q are defined by

$$\lambda_1(r, s) = \frac{3r}{4} + \frac{s}{4}, \quad \lambda_2(r, s) = \frac{3s}{4} + \frac{r}{4}, \quad Q(r, s) = 1 - \frac{4}{F} \left(\frac{r+s}{r-s}\right)^2.$$

For a discontinuity located at $x = x(t)$, the Rankine–Hugoniot conditions in these variables reads

$$(9) \quad \begin{aligned} [(r + s)(r - s)^2] &= 2 \dot{x} [(r - s)^2], \\ [(r - s)^4 + 8(r - s)^2(r + s)^2] &= 16 \dot{x} [(r + s)(r - s)^2]. \end{aligned}$$

The derivation of the spectral problem is inspired by the method introduced by Majda for studying the stability of multidimensional shocks in hyperbolic systems [4]. In what follows, we consider a Dressler roll-wave (H, U) (or equivalently, R, S) with wavespeed c and wavelength L : the location of the shocks are given by $x_i(t) = ct + iL, i \in \mathbb{Z}$. We consider small perturbations of that solution: we can suppose it is piecewise C^1 functions with discontinuities at positions $X_i(t) = ct + iL + \varepsilon_i(t), i \in \mathbb{Z}$. The ε_i are supposed to be small. The Rankine–Hugoniot conditions (9) are then given by

$$(10) \quad \begin{aligned} [(r + s)(r - s)^2]_{X_i} &= 2 \dot{X}_i [(r - s)^2]_{X_i}, \\ [(r - s)^4 + 8(r - s)^2(r + s)^2]_{X_i} &= 16 \dot{X}_i [(r + s)(r - s)^2]_{X_i} \end{aligned} \quad \forall i \in \mathbb{Z}.$$

In [4], working in the reference frame attached to the discontinuity made the shock steady. In our case, because of the existence of an *infinite* distribution of shocks, we have to fix *all* the discontinuities in order to make the roll-wave steady, just as Majda did in the case of the stability of planar shocks [4]. In that case, a perturbation of a planar shock has a discontinuity front at points $x_n = \psi(y, t)$, with $y = (x_1, \dots, x_{n-1})$ the transverse variable, and one introduces a new variable $\xi = x_n - \psi(y, t)$ that fixes *all* the front shock and not only a variable $\xi = x_n - \psi(t)$ that fixes only the shock at $y = 0$. In the case of roll-waves, we fix all the shocks by working in the new system of coordinates $(\xi = \xi(x, t), t)$ such that $\xi(X_i(t), t) = iL \forall i \in \mathbb{Z}$. The function ξ has the form

$$(11) \quad \forall x \in]X_i(t), X_{i+1}(t)[, \xi(x, t) = \frac{x - X_i(t)}{X_{i+1}(t) - X_i(t)} L + iL \forall i \in \mathbb{Z}.$$

Remark. Considering $|\varepsilon_i| \ll 1$, we prove that $\xi(x, t) \approx x - ct$ up to first order; it looks like a change of reference frame.

Note that in the case of the stability of roll-waves, we have a sort of multidimensional stability issue where the transverse “variable” is $i \in \mathbb{Z}$. The function ξ is Lipschitz and the system is of order one, so this change of coordinates is licit and does not change the Rankine–Hugoniot jump conditions. We make the change of coordinates $(r, s)(x, t) = (\bar{r}, \bar{s})(\xi(x, t), t)$, where (\bar{r}, \bar{s}) are piecewise C^1 functions with discontinuities located at points $\{iL, i \in \mathbb{Z}\}$. The derivation rules are given by

$$(12) \quad \begin{aligned} \frac{\partial r}{\partial x} &= \frac{\partial \bar{r}}{\partial \xi} \frac{L}{X_{i+1} - X_i}, \\ \frac{\partial r}{\partial t} &= \frac{\partial \bar{r}}{\partial t} - \frac{L}{X_{i+1} - X_i} \left(\dot{X}_i + \frac{\xi - iL}{L} (\dot{X}_{i+1} - \dot{X}_i) \right) \frac{\partial \bar{r}}{\partial \xi}. \end{aligned}$$

Substituting (12) into the shallow water equations (8) and dropping the overlines into the equations, the system (8) reads for all $x \in (iL, (i + 1)L)$,

$$(13) \quad \begin{aligned} r_t + \frac{L}{X_{i+1} - X_i} \left(\lambda_1(r, s) - \left(\dot{X}_i + \gamma_i(x) (\dot{X}_{i+1} - \dot{X}_i) \right) \right) r_x &= Q(r, s), \\ s_t + \frac{L}{X_{i+1} - X_i} \left(\lambda_2(r, s) - \left(\dot{X}_i + \gamma_i(x) (\dot{X}_{i+1} - \dot{X}_i) \right) \right) s_x &= Q(r, s), \end{aligned}$$

with $\gamma_i(x) = \frac{x-iL}{L} \in (0, 1)$. The Rankine–Hugoniot conditions (10) remain unchanged and are given by

$$(14) \quad [hu]_{X_i} = \dot{X}_i[h]_{X_i}, \quad \left[\frac{h^2}{2F} + hu^2 \right]_{X_i} = \dot{X}_i[hu]_{X_i},$$

where h, u have to be understood as functions of r, s . The particular solution $(R, S, X_i(t) = ct + iL)$ computed by Dressler is now a steady solution of (13) and (14).

In this paper, we deal with the problem of the *persistence* of a roll-wave solution: Given any initial condition (r_0, s_0) which has the same structure as a roll-wave solution with discontinuities at points $x_i^0, i \in \mathbb{Z}$, we are going to prove that, for initial conditions close to a roll-wave solution and satisfying suitable compatibility conditions, there exists a solution of the Cauchy problem (13), (14) with the initial conditions $(r, s)(\cdot, 0) = (r_0, s_0)$ on a sufficiently small interval $(0, T^*)$ that conserves a structure similar to the roll-wave structure. This result is in some sense a result of the *stability* of roll waves. Due to the hyperbolic nature of the problem, we expect the apparition of new discontinuities that breaks the roll-wave structure. To deal with this problem, we follow the approach initiated by Majda [4] to prove the existence of compact shocks and developed by Metivier [5] for the stability of multidimensional shocks: In order to deal with the full nonlinear problem, the first task is to prove the well posedness of the linearized equations around an approximate solution, which will be assumed to be close to a roll wave solution. Then, under suitable compatibility conditions, we are able to prove the existence of an approximate solution of (13), (14), and using a fixed point argument in the neighborhood of this approximate solution, we prove the existence of a solution of the full Cauchy problem.

3. Well posedness of the linearized equations. In this section, we prove the well posedness of the linear equations obtained by the linearization of (13), (14) in the neighborhood of an approximate solution close to a roll-wave. This is done by energy estimates on the linear problem and on an adjoint problem. We start with the analysis of the spectral problem associated to the linear equations obtained by linearization of (13), (14) around the *roll-wave solution* and then perform a priori estimates of solutions of that problem and the adjoint of that problem. We prove that these estimates remain true for the solutions of the linear equations obtained by linearization around a function “close” to a roll-wave. We deduce from the energy estimates that the linearized problem around an approximate solution is well posed in suitable functional spaces.

3.1. Energy estimates on the “exact” linearized problem. In what follows, we linearize shallow water equations and Rankine–Hugoniot conditions around an exact roll-wave solution and obtain energy estimates for that linearized problem: we shall deduce estimates for the linearized problem around an approximate solution.

3.1.1. Formulation of the linearized problem. The linearization of the equation in the neighborhood of the roll-wave solution yields the following, after time rescaling $t := H_0 t$:

$$(15) \quad \begin{pmatrix} \partial_t r_i + a(h)\partial_x r_i \\ \partial_t s_i + b(h)\partial_x s_i \end{pmatrix} + L^0(h) \begin{pmatrix} r_i \\ s_i \end{pmatrix} = \mathcal{F}_i \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

with $a(h) = \sqrt{h} - h^{-1}$, $b(h) = -\sqrt{h} - h^{-1}$, $(r_i, s_i) = (r, s)|_{(iL, iL+L)}$,

$$\mathcal{F}_i = \left(\dot{\varepsilon}_i + \frac{x}{L} (\dot{\varepsilon}_{i+1} - \dot{\varepsilon}_i) \right) P(h) \begin{pmatrix} h^{-2} + h^{-\frac{1}{2}} \\ h^{-2} - h^{-\frac{1}{2}} \end{pmatrix} - \frac{c(\varepsilon_{i+1} - \varepsilon_i)P(h)}{L(1 + \sqrt{F})} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

and

$$\begin{aligned} L_{11}^0 &= \frac{3}{4} \left(h^{-2} + h^{-\frac{1}{2}} \right) P(h) - h^{-\frac{3}{2}} \left(1 + \sqrt{F} - h^{-1} \right) \left(1 + \sqrt{F} - h^{-1} - 2h^{\frac{1}{2}} \right), \\ L_{12}^0 &= \frac{1}{4} \left(h^{-2} + h^{-\frac{1}{2}} \right) P(h) + h^{-\frac{3}{2}} \left(1 + \sqrt{F} - h^{-1} \right) \left(1 + \sqrt{F} - h^{-1} + 2h^{\frac{1}{2}} \right), \\ L_{21}^0 &= \frac{1}{4} \left(h^{-2} - h^{-\frac{1}{2}} \right) P(h) - h^{-\frac{3}{2}} \left(1 + \sqrt{F} - h^{-1} \right) \left(1 + \sqrt{F} - h^{-1} - 2h^{\frac{1}{2}} \right), \\ L_{22}^0 &= \frac{3}{4} \left(h^{-2} - h^{-\frac{1}{2}} \right) P(h) + h^{-\frac{3}{2}} \left(1 + \sqrt{F} - h^{-1} \right) \left(1 + \sqrt{F} - h^{-1} + 2h^{\frac{1}{2}} \right). \end{aligned}$$

Here h has to be understood as $h = h(x)$, the equation of profile of the roll-wave. The linearized Rankine–Hugoniot conditions are given by

$$\begin{aligned} 2[h]_{iL} \dot{\varepsilon}_i &= \left[\left(h - h^{-\frac{1}{2}} \right) r + \left(h + h^{-\frac{1}{2}} \right) s \right]_{iL}, \\ 2[h]_{iL} \dot{\varepsilon}_i &= \left[\left(h - h^{-\frac{1}{2}} + \frac{1}{1 + \sqrt{F}} \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2 \right) \right) r \right. \\ (16) \quad &\left. + \left(h + h^{-\frac{1}{2}} - \frac{1}{1 + \sqrt{F}} \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2 \right) \right) s \right]_{iL} \quad \forall i \in \mathbb{Z}, \end{aligned}$$

with $[h]_{iL} = h_+ - h_-$. In order to simplify the notations, we introduce the matrix $M(h)$ so that (16) reads

$$2[h] \dot{\varepsilon} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \left[M(h) \begin{pmatrix} r \\ s \end{pmatrix} \right]_{iL}.$$

Equivalently, system (16) can be written

$$\begin{aligned} 2[h] \dot{\varepsilon}_i &= \left[\left(h - h^{-\frac{1}{2}} \right) r + \left(h + h^{-\frac{1}{2}} \right) s \right]_{iL}, \\ (17) \quad 0 &= \left[\left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2 \right) r - \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2 \right) s \right]_{iL} \quad \forall i \in \mathbb{Z}. \end{aligned}$$

In order to eliminate the derivatives of ε_i in the interior equations (15), one can, as usual, introduce the “good” unknowns:

$$\begin{aligned} \bar{r}_i &= r_i + \left(\varepsilon_i + \frac{x}{L} (\varepsilon_{i+1} - \varepsilon_i) \right) P(h) \left(h^{-2} + h^{-\frac{1}{2}} \right), \\ \bar{s}_i &= s_i + \left(\varepsilon_i + \frac{x}{L} (\varepsilon_{i+1} - \varepsilon_i) \right) P(h) \left(h^{-2} - h^{-\frac{1}{2}} \right). \end{aligned}$$

Dropping the overlines, this yields the system

$$(18) \quad \begin{aligned} \partial_t r_i + a(h) \partial_x r_i &= \bar{\mathcal{F}}_i \quad \forall (x, i) \in (0, L) \times \mathbb{Z}, \\ \partial_t s_i + b(h) \partial_x s_i & \end{aligned}$$

where $\overline{\mathcal{F}}_i$ contains only zeroth order terms. More precisely, $\overline{\mathcal{F}}_i$ is written

$$\begin{aligned} \overline{\mathcal{F}}_i = & -L^0 \begin{pmatrix} r_i \\ s_i \end{pmatrix} - \varepsilon_i P(h) L^0 \begin{pmatrix} h^{-2} + h^{-\frac{1}{2}} \\ h^{-2} - h^{-\frac{1}{2}} \end{pmatrix} \\ & - \frac{x}{L} (\varepsilon_{i+1} - \varepsilon_i) P(h) L^0 \begin{pmatrix} h^{-2} + h^{-\frac{1}{2}} \\ h^{-2} - h^{-\frac{1}{2}} \end{pmatrix} \\ & - \varepsilon_i \text{diag}(a(h), b(h)) \partial_x \begin{pmatrix} P(h)(h^{-2} + h^{-\frac{1}{2}}) \\ P(h)(h^{-2} - h^{-\frac{1}{2}}) \end{pmatrix} \\ & - (\varepsilon_{i+1} - \varepsilon_i) \text{diag}(a(h), b(h)) \partial_x \begin{pmatrix} \frac{x P(h)}{L} (h^{-2} + h^{-\frac{1}{2}}) \\ \frac{x P(h)}{L} (h^{-2} - h^{-\frac{1}{2}}) \end{pmatrix}. \end{aligned}$$

The Rankine–Hugoniot jump conditions are given by

$$\begin{aligned} 2[h]\dot{\varepsilon}_i &= \left[\left(h - h^{-\frac{1}{2}} \right) r + \left(h + h^{-\frac{1}{2}} \right) s \right]_{iL}, \\ (19) \quad - [P(h)(h - h^{-2})] \varepsilon_i &= \left[\left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2 \right) r - \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2 \right) s \right]_{iL}. \end{aligned}$$

As usual, we consider the linear boundary value problem dropping the zeroth order term from (18, 19):

$$\begin{aligned} (20) \quad \partial_t r_i + a(h) \partial_x r_i &= F_i^1, \\ \partial_t s_i + b(h) \partial_x s_i &= F_i^2, \quad \forall i \in \mathbb{Z}, x \in (0, L), \end{aligned}$$

with the linearized Rankine–Hugoniot conditions

$$\begin{aligned} (21) \quad 2[h]\dot{\varepsilon}_i &= \left[\left(h - h^{-\frac{1}{2}} \right) r + \left(h + h^{-\frac{1}{2}} \right) s \right]_{iL} + a_{1,i}, \\ \left[\left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2 \right) r - \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2 \right) s \right]_{iL} &= a_{2,i}, \quad \forall i \in \mathbb{Z}. \end{aligned}$$

These conditions can be interpreted as transmission conditions at the point of discontinuities between (r_i, s_i) and (r_{i+1}, s_{i+1}) . Indeed, a jump of a function f of r, s at point $x = iL$ can be written

$$[f(r, s)]_{iL} = f(r_i(0), s_i(0)) - f(r_{i-1}(L), s_{i-1}(L)).$$

3.1.2. Normal mode analysis. We carry out a similar analysis made in the stability of multidimensional shocks and search for possible unstable eigenmodes. We make a Laplace transform in time and Fourier transform in discrete space (that yields Fourier series) of (20, 21):

$$\lambda \bar{r} + a(h) \bar{r}' = 0, \quad \lambda \bar{s} + b(h) \bar{s}' = 0 \quad \forall x \in (0, L),$$

with

$$\bar{r}(\lambda, \theta, x) = \sum_{k \in \mathbb{Z}} e^{ik\theta} \int_{-\infty}^{\infty} e^{\lambda t} r_k(x, t) dt.$$

The linearized Rankine–Hugoniot conditions reads

$$2\lambda(h_- - h_+) \bar{\varepsilon} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = e^{i\theta} M(h_-) \begin{pmatrix} \bar{r}(L) \\ \bar{s}(L) \end{pmatrix} - M(h_+) \begin{pmatrix} \bar{r}(0) \\ \bar{s}(0) \end{pmatrix}.$$

Here the space \mathbb{E}^s of functions decreasing to 0 at infinity used in the analysis of multidimensional shocks is replaced by the space \mathbb{E}^a of *bounded* functions at the point of singularity x_0 such that $h(x_0) = 1$. A basis of solutions for the differential system (3.1.2) is given by

$$V_1 = \left((x - x_0)^{-\frac{\lambda}{a'(x_0)}} Z(x), 0 \right), \quad V_2 = \left(0, e^{-\lambda \int_{x_0}^x \frac{ds}{b(h(s))}} \right),$$

where Z is an analytic (thus bounded) function on the interval $(0, L)$ that is unique if we choose $Z(0) = 1$. When $\Re(\lambda) > 0$, it is easily seen that $\mathbb{E}^a = \langle V_2 \rangle$ so that the analogous equation of the Lopatinskii determinant is an Evans function given by

$$\Delta(\lambda, \theta) = \frac{2[h]}{1 + \sqrt{F}} \lambda \left(\left(h_+^{\frac{3}{2}} + h_+^{-\frac{3}{2}} + 2 \right) - \left(h_-^{\frac{3}{2}} + h_-^{-\frac{3}{2}} + 2 \right) e^{i\theta - \lambda \int_0^L \frac{ds}{b(h(s))}} \right).$$

The function $h \mapsto h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2$ is decreasing on the interval $(0, 1)$, is increasing on $(1, +\infty)$, and is invariant under the function $\mapsto \frac{1}{h}$. The Rankine–Hugoniot jump conditions for the roll-wave $h_+ + h_- = \frac{2}{h_+^{\frac{3}{2}} h_-}$ imply that $h_- = *10\frac{1}{2}(\sqrt{h_+^2 + \frac{8}{h_+}} - h_+) < \frac{1}{h_+}$ for any $h_+ < 1$. Then we find that $h_+^{\frac{3}{2}} + h_+^{-\frac{3}{2}} + 2 > h_-^{\frac{3}{2}} + h_-^{-\frac{3}{2}} + 2$. As a consequence, on the half plane $\Re(\lambda) > 0$, the Evans function Δ vanishes either when $\lambda = 0$ (for any θ) and for

$$-\Re(\lambda) \int_0^L \frac{ds}{b(h(s))} = \log \left(\frac{h_+^{\frac{3}{2}} + h_+^{-\frac{3}{2}} + 2}{h_-^{\frac{3}{2}} + h_-^{-\frac{3}{2}} + 2} \right), \quad \theta = \Im(\lambda) \int_0^L \frac{ds}{b(h(s))} [2\pi].$$

We obtain a full line of unstable eigenmodes with a positive growth rate $\Re(\lambda)$ in time, since $b(h) < 0$ and $h_+^{\frac{3}{2}} + h_+^{-\frac{3}{2}} + 2 > h_-^{\frac{3}{2}} + h_-^{-\frac{3}{2}} + 2$.

Contrary to the case of multidimensional shocks where there is only one shock, this *does not* imply strong instability since the Evans function is not homogeneous; the growth rate of the unstable eigenmodes is *bounded*, whereas in the case of a single shocks, a full line of eigenvalues with an arbitrary large real part is found. This is due to the fact that in the case of shocks, the equations are invariant under the hyperbolic scaling $(x, t) \mapsto (\mu x, \mu t)$ for any $\mu > 0$. As a consequence, an unstable mode is accompanied with a full line of unstable modes that are more and more unstable as the frequency goes to infinity. This is not the case here due to the periodic nature of the problem. Moreover, it is easily seen that when $L \rightarrow \infty$, $\Re(\lambda) \rightarrow 0$: the zeroth order terms have an importance in the stability issue as pointed out in [13].

3.1.3. Energy estimates. In this section, we obtain energy estimates for solutions of the linear problem (20), (21). We multiply the first (resp. second) equation of (20) by r_i (resp. s_i), integrate on $(0, L)$, and sum for $i \in \mathbb{Z}$:

$$\partial_t \sum_{i \in \mathbb{Z}} \int_0^L (r_i^2 + s_i^2) dx + \sum_{i \in \mathbb{Z}} \int_0^L a(h) \partial_x r_i^2 + b(h) \partial_x s_i^2 = 2 \sum_{i \in \mathbb{Z}} \int_0^L F_i^1 r_i + F_i^2 r_i dx.$$

An integration by parts on the second integral yields

$$\int_0^L a(h) \partial_x r_i^2 + b(h) \partial_x s_i^2 = [a(h) r_i^2 + b(h) s_i^2]_0^L - \int_0^L \partial_x a r_i^2 + \partial_x b s_i^2 dx.$$

We change the indices of summation in the boundary terms and find that

$$(22) \quad \partial_t \sum_{i \in \mathbb{Z}} \int_0^L (r_i^2 + s_i^2) dx - \sum_{i \in \mathbb{Z}} [a(h)r^2 + b(h)s^2]_{iL} \\ = 2 \sum_{i \in \mathbb{Z}} \int_0^L F_i^1 r_i + F_i^2 s_i dx + \sum_{i \in \mathbb{Z}} \int_0^L \partial_x a(h) r_i^2 + \partial_x b(h) s_i^2.$$

Equation (22) can be written in the form

$$(23) \quad \partial_t \int_{\mathbb{R}} r^2 + s^2 dx - \sum_{i \in \mathbb{Z}} [a(r) r^2 + b(h) s^2]_{iL} \\ = 2 \int_{\mathbb{R}} F r + G s dx + \int_{\mathbb{R}} \partial_x a(h) r^2 + \partial_x b(h) s^2 dx.$$

Next, we have to obtain an estimate on the boundary terms. For that purpose, we show that the boundary condition

$$f_- r_i^- - f_+ r_i^+ = g_- s_i^- - g_+ s_i^+,$$

with

$$f_{\pm} = f(h_{\pm}), g_{\pm} = g(h_{\pm}), f(h) = h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2, g(h) = h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2,$$

is maximal dissipative. More precisely, we prove that

$$B(r_i^+, r_i^-, s_i^+, s_i^-) = 0 \Rightarrow \mathcal{Q}(r_i^+, r_i^-, s_i^+, s_i^-) \geq 0,$$

with $Bu = \langle rh, u \rangle$, rh being the vector $rh = {}^t (f_+, -f_-, -g_+, g_-)$ and \mathcal{Q} is the quadratic form

$$\mathcal{Q}(r_i^+, r_i^-, s_i^+, s_i^-) = a_- (r_i^-)^2 + b_- (s_i^-)^2 - a_+ (r_i^+)^2 - b_+ (s_i^+)^2,$$

with $a_{\pm} = a(h_{\pm}), b_{\pm} = b(h_{\pm})$. This property is necessarily maximal, since $\dim(\text{Ker}B) = 3$ and $b_- < 0$, whereas $-b_+, a_-, -a_+ > 0$. The restriction of the quadratic form \mathcal{Q} to $\text{Ker}B$ is represented by the symmetric matrix \mathcal{A} :

$$\mathcal{A} = \begin{pmatrix} a_- g_+^2 - b_+ f_-^2 & b_+ f_- g_- & b_+ f_+ f_- \\ b_+ f_- g_- & b_- g_+^2 - b_+ g_-^2 & -b_+ f_+ g_- \\ b_+ f_+ f_- & -b_+ f_+ g_- & -a_+ g_+^2 - b_+ f_+^2 \end{pmatrix}.$$

The matrix \mathcal{A} is definite positive provided that all the principal determinants $\Delta_i, i = 1, 2, 3$ are strictly nonnegative. Let us check that condition. First, $\Delta_1 = a_- g_+^2 - b_+ f_-^2 \geq 0$ since $a_-, -b_+ \geq 0$, and $\Delta_1 \neq 0$ provided that there is a real shock $h_+ < 1 < h_-$ (otherwise, it is just the stationary solution). A straightforward computation shows that

$$\Delta_2 = g_+^2 (a_- b_- g_+^2 - b_+ (a_- g_-^2 + b_- f_-^2)),$$

and $\Delta_2 > 0$ if and only if the following condition is satisfied:

$$(24) \quad h_+ \left(1 + h_+^{\frac{3}{2}}\right) \left(1 + h_+^{-\frac{3}{2}}\right)^2 < 6h_- + \frac{2}{h_-^2}.$$

It is easily proved that this relation is satisfied for *small amplitude* roll-waves when $h_+ < 1 < h_-$ lie in the neighborhood of $h = 1$. In the large amplitude case, recall that $h_a < h_+ < 1$, where $h_a = \frac{1}{2F}(1 + 2\sqrt{F} + \sqrt{1 + 4\sqrt{F}})$ and $h_+ + h_- = \frac{2}{h_+ h_-}$. We write the Rankine–Hugoniot jump condition for the roll wave as

$$h_- = RH(h_+) = \frac{1}{2} \left(\sqrt{h_+^2 + \frac{8}{h_+}} - h_+ \right)$$

and inserting that relation into (24), the problem is reduced to the analysis of the sign of

$$C(h_+) = 6RH(h_+) + \frac{2}{RH(h_+)^2} - h_+ \left(1 + h_+^{\frac{3}{2}} \right) \left(1 + h_+^{-\frac{3}{2}} \right)^2$$

when $h_+ \in (h_a, 1)$. It is a lengthy but straightforward computation to prove that $C(h_+) > 0$ for all $h_+ \in (h_a, 1)$ provided that $4 < F \leq F_c$, with $F_c \geq 18$. Thus, when $4 < F \leq F_c$, the relation $\Delta_2 > 0$ is satisfied for *all* roll-wave solutions and for larger Froude numbers; there exists $h_c(F) > h_a$ so that $\Delta_2 > 0$ is satisfied for all $h_+ \in (h_c(F), 1)$. Next, we consider the last condition $\Delta_3 > 0$. One can show that

$$\Delta_3 = g_+^4 \left((a_- g_-^2 + b_- f_-^2) a_+ b_+ - (a_+ g_+^2 + b_+ f_+^2) a_- b_- \right)$$

and $\Delta_3 > 0$ if and only if

$$3h_- + \frac{1}{h_-^2} > 3h_+ + \frac{1}{h_+^2}.$$

Then it is an easy calculation to prove that this condition is satisfied for *all* roll-wave solutions if $4 < F \leq F_c$, with $F_c \geq 12$. As a conclusion, the matrix \mathcal{A} is definite positive if $4 < F \leq F_c$ ($F_c \geq 12$). Suppose that this assumption is satisfied. The quadratic form $\mathcal{Q}_{\text{KerB}}$ is definite positive. Denote $w = (r_i^+, r_i^-, s_i^+, s_i^-)$. We deduce that there exist $\kappa > 0$ and $C > 0$ so that

$$\mathcal{Q}(w) \geq \kappa \|w\|^2 - C \|Bw\|^2$$

and substituting that relation into (23), we deduce the estimate

$$\begin{aligned} & \partial_t \|(r, s)\|_2^2 + \kappa \sum_{i \in \mathbb{Z}} |(r_i^\pm, s_i^\pm)|^2 \\ (25) \quad & \leq 2 \int_{\mathbb{R}} F^1 r + F^2 s \, dx + C \|(r, s)\|^2 + \|a_2\|_{l^2(\mathbb{Z})}^2, \end{aligned}$$

where $\|(r, s)\|_2^2 = \int_{\mathbb{R}} r^2 + s^2 \, dx$ is the classical $(L^2(\mathbb{R}))^2$ -norm. From the Rankine–Hugoniot condition, one can easily estimate $\varepsilon_i, \dot{\varepsilon}_i$:

$$|\dot{\varepsilon}_i|^2 \leq C \left(|(r_i^\pm, s_i^\pm)|^2 + a_{1,i}^2 \right).$$

We compute estimates in weighted functional spaces. For $\gamma > 0$ sufficiently large (that depends on $\|(R, S)\|_{W^{1,\infty}(0,L)}$) and using Young’s type inequalities, we obtain the following proposition.

PROPOSITION 1. Assume that the Froude number F satisfies $0 < F < F_c$ (with $F_c > 12$). Then there exists γ_c (depending on $\|(R, S)\|_{W^{1,\infty}(0,L)}$) such that for all $\gamma > \gamma_c$, the following estimate holds true:

$$(26) \quad \gamma \int_{\mathbb{R}} e^{-2\gamma t} \|(r, s)\|_2^2 + \int_{\mathbb{R}} e^{-2\gamma t} \sum_{i \in \mathbb{Z}} |(r_i^\pm, s_i^\pm)|^2 + (\varepsilon_i^2 + \dot{\varepsilon}_i^2) dt \leq C \left(\frac{1}{\gamma} \int_{\mathbb{R}} e^{-2\gamma t} \|(F^1, F^2)\|_2^2 dt + \int_{\mathbb{R}} e^{-2\gamma t} (\|a_1\|_{l^2(\mathbb{Z})}^2 + \|a_2\|_{l^2(\mathbb{Z})}^2) dt \right).$$

We are now in a position to make the energy estimates on the linearized problem around an approximate solution close to a roll-wave. In what follows, we shall assume that the assumption $0 < F < F_c$ is satisfied.

3.2. Energy estimates for the “approximate” linearized problem. In order to tackle the full nonlinear problem, we consider the linearization of the Saint Venant system (13) in the neighborhood of an approximate solution (R_i, S_i, X_i) that is close to the original roll-wave $(R, S, ct + iL)$. In what follows, we shall prove that the estimates (26) obtained in the previous section remain true for (R_i, S_i, X_i) sufficiently “close” to the roll-wave solution $(R, S, ct + iL)$.

3.2.1. Formulation of the problem. We consider an approximate “solution” $(R_i, S_i, X_i) \approx (R, S, ct + iL)$ and linearize the shallow water equations and Rankine–Hugoniot jump conditions around that function:

$$(27) \quad \begin{aligned} \partial_t r_i + \frac{L}{X_{i+1} - X_i} \Lambda_1^i \partial_x r_i &= \bar{F}_i^1, \\ \partial_t s_i + \frac{L}{X_{i+1} - X_i} \Lambda_2^i \partial_x s_i &= \bar{F}_i^2 \end{aligned} \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

where the function $F_k^i, k = 1, 2$ are defined by

$$\begin{aligned} \bar{F}_i^1 &= \frac{\partial Q}{\partial r}(R_i, S_i)r_i + \frac{\partial Q}{\partial s}(R_i, S_i)s_i + \frac{L(\varepsilon_{i+1} - \varepsilon_i)}{(X_{i+1} - X_i)^2} \Lambda_1^i \\ &\quad - \frac{L}{X_{i+1} - X_i} \left(\lambda_1(r_i, s_i) - \left(\dot{\varepsilon}_i + \frac{x}{L} (\dot{\varepsilon}_{i+1} - \dot{\varepsilon}_i) \right) \right) \partial_x R_i \\ \bar{F}_i^2 &= \frac{\partial Q}{\partial r}(R_i, S_i)r_i + \frac{\partial Q}{\partial s}(R_i, S_i)s_i + \frac{L(\varepsilon_{i+1} - \varepsilon_i)}{(X_{i+1} - X_i)^2} \Lambda_2^i \\ &\quad - \frac{L}{X_{i+1} - X_i} \left(\lambda_2(r_i, s_i) - \left(\dot{\varepsilon}_i + \frac{x}{L} (\dot{\varepsilon}_{i+1} - \dot{\varepsilon}_i) \right) \right) \partial_x S_i \end{aligned}$$

and $\Lambda_k^i, k = 1, 2, i \in \mathbb{Z}$ denotes $\Lambda_k^i = \lambda_k(R_i, S_i) - (\dot{X}_i + \frac{x}{L}(\dot{X}_{i+1} - \dot{X}_i))$. The linearized Rankine–Hugoniot jump conditions are given by

$$(28) \quad \begin{aligned} [\mathbb{Q}(R_i, S_i)]\dot{\varepsilon}_i &= \left[\left(\frac{\partial \mathbb{F}}{\partial R} - \dot{X}_i \frac{\partial \mathbb{Q}}{\partial R} \right) r + \left(\frac{\partial \mathbb{F}}{\partial S} - \dot{X}_i \frac{\partial \mathbb{Q}}{\partial S} \right) s \right]_{iL} \\ &\quad + \frac{L\varepsilon_i}{X_{i+1} - X_i} \left[\partial_x (\mathbb{F} - \dot{X}_i \mathbb{Q}) \right]_{iL}, \\ [H(R_i, S_i)]\dot{\varepsilon}_i &= \left[\left(\frac{\partial \mathbb{Q}}{\partial R} - \dot{X}_i \frac{\partial H}{\partial R} \right) r + \left(\frac{\partial \mathbb{Q}}{\partial S} - \dot{X}_i \frac{\partial H}{\partial S} \right) s \right]_{iL} \\ &\quad + \frac{L\varepsilon_i}{X_{i+1} - X_i} \left[\partial_x (\mathbb{Q} - \dot{X}_i H) \right]_{iL}, \end{aligned}$$

where the functions $\mathbb{Q} = HU, \mathbb{F} = \frac{H^2}{2F} + HU^2$ shall be considered as functions of R, S . One can introduce the “good” unknowns to drop the derivatives of ε_i in the interior equation (27). Similarly to the previous section, we analyse the linear problem obtained after dropping zeroth order terms in (27), (28):

$$(29) \quad \begin{aligned} \partial_t r_i + \frac{L}{X_{i+1} - X_i} \Lambda_1^i \partial_x r_i &= F_i^1, \\ \partial_t s_i + \frac{L}{X_{i+1} - X_i} \Lambda_2^i \partial_x s_i &= F_i^2 \end{aligned} \quad \forall x \in (0, L) \quad \forall i \in \mathbb{Z}.$$

The system (29) is completed with the linearized Rankine–Hugoniot conditions

$$(30) \quad \begin{aligned} [\mathbb{Q}(R_i, S_i)]\dot{\varepsilon}_i &= \left[\left(\frac{\partial \mathbb{F}}{\partial R} - \dot{X}_i \frac{\partial \mathbb{Q}}{\partial R} \right) r + \left(\frac{\partial \mathbb{F}}{\partial S} - \dot{X}_i \frac{\partial \mathbb{Q}}{\partial S} \right) s \right]_{iL} + a_{1,i}, \\ [H(R_i, S_i)]\dot{\varepsilon}_i &= \left[\left(\frac{\partial \mathbb{Q}}{\partial R} - \dot{X}_i \frac{\partial H}{\partial R} \right) r + \left(\frac{\partial \mathbb{Q}}{\partial S} - \dot{X}_i \frac{\partial H}{\partial S} \right) s \right]_{iL} + a_{2,i} \end{aligned} \quad \forall i \in \mathbb{Z}.$$

3.2.2. Energy estimates. In what follows, we shall prove the following result.

PROPOSITION 2. *There exists $\eta > 0$ so that for all $m > 0$, there exists a constant C and γ_0 depending on m and for any Lipschitz continuous and bounded function $(R_i, S_i, X_i)_{i \in \mathbb{Z}}$ satisfying*

$$\begin{aligned} \max_{i \in \mathbb{Z}} \left(\|\dot{X}_i - c\|_\infty + \|X_{i+1} - X_i - L\|_\infty + \|R_i - R, S_i - S\|_\infty \right) &\leq \eta, \\ \max_{i \in \mathbb{Z}} \left(|\partial_x(R_i, S_i)|_\infty + \left| \dot{X}_i \right|_\infty \right) &\leq m, \end{aligned}$$

and for all $\gamma \geq \gamma_0, (F^1, F^2) \in \mathcal{D}(\mathbb{R}/\{iL\} \times \mathbb{R}), (a_1, a_2) \in l^2(\mathbb{Z}, \mathcal{D}(\mathbb{R}))$, then a solution $(r_i, s_i, \varepsilon_i)_{i \in \mathbb{Z}}$ of (29), (30) satisfies

$$\begin{aligned} \gamma \left\| e^{-\gamma t} (r, s) \right\|_{L^2(\mathbb{R}/\{iL\} \times \mathbb{R})}^2 + \left\| e^{-\gamma t} (r_i^\pm, s_i^\pm) \right\|_{l^2(\mathbb{Z}, L^2(\mathbb{R}))}^2 + \left\| e^{-\gamma t} \varepsilon_i \right\|_{l^2(\mathbb{Z}, H^1(\mathbb{R}))}^2 \\ \leq C \left(\frac{1}{\gamma} \left\| e^{-\gamma t} (F^1, F^2) \right\|_{L^2(\mathbb{R}/\{iL\} \times \mathbb{R})}^2 + \left\| e^{-\gamma t} (a_1, a_2) \right\|_{l^2(\mathbb{Z}, L^2(\mathbb{R}))}^2 \right). \end{aligned}$$

Proof. Similarly to the previous case, we multiply the first (resp. second) equation of (29) by r_i (resp. s_i), integrate on $(0, L)$, and sum for $i \in \mathbb{Z}$:

$$\begin{aligned} \partial_t \int_{\mathbb{R}} r^2 + s^2 dx + \sum_{i \in \mathbb{Z}} \frac{L}{X_{i+1} - X_i} [\Lambda_1^i r_i^2 + \Lambda_2^i s_i^2]_0^L \\ = 2 \int_{\mathbb{R}} F^1 r + F^2 s dx + \sum_{i \in \mathbb{Z}} \frac{L}{X_{i+1} - X_i} \int_0^L \partial_x \Lambda_1^i r_i^2 + \partial_x \Lambda_2^i s_i^2 dx. \end{aligned}$$

It is a straightforward computation to prove that the boundary terms satisfy

$$\begin{aligned} \sum_{i \in \mathbb{Z}} \frac{L}{X_{i+1} - X_i} [\Lambda_1^i r_i^2 + \Lambda_2^i s_i^2]_0^L \geq - \sum_{i \in \mathbb{Z}} [a(h) r^2 + b(h) s^2]_{iL} \\ - \|\delta\|_{l^\infty(\mathbb{Z})} \sum_i |(r_i^\pm, s_i^\pm)|^2, \end{aligned}$$

where (r_i^\pm, s_i^\pm) denotes the vector $(r_i^+, r_i^-, s_i^+, s_i^-)$ and $(\delta_i)_{i \in \mathbb{Z}}$ satisfies

$$\delta_i = \mathcal{O} \left(\|\dot{X}_i - c\|_\infty + \|X_{i+1} - X_i - L\|_\infty + \|R_i - R, S_i - S\|_\infty \right).$$

On the other hand, for $\max_i (\|R - R_i, S - S_i\|_\infty)$ sufficiently small, we obtain $[H(R_i, S_i)] \neq 0$, and we can write the Rankine–Hugoniot jump conditions in the form

$$\dot{\varepsilon}_i = [\mathcal{L}_1^i(r, s)]_{iL} + L_1^i(a_i, b_i), \quad [\mathcal{L}_2^i(r, s)]_{iL} = L_2^i(a_i, b_i),$$

so that

$$[\mathcal{L}_2^i(r, s)]_{iL} = \left[\left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} - 2 \right) r + \left(h^{\frac{3}{2}} + h^{-\frac{3}{2}} + 2 \right) s \right]_{iL} + \mu_i |(r_i^\pm, s_i^\pm)|,$$

and the real numbers μ_i are estimated by

$$\mu_i = \mathcal{O} \left(\|\dot{X}_i - c\|_\infty + \|X_{i+1} - X_i - L\|_\infty + \|R_i - R, S_i - S\|_\infty \right).$$

Moreover, the linear operators L_k^i are uniformly bounded:

$$|L_k^i(a, b)| \leq C|(a, b)| \quad \forall i \in \mathbb{Z}, \quad k = 1, 2.$$

Choosing μ_i, δ_i sufficiently small and for

$$\max_{i \in \mathbb{Z}} (\|R_i, S_i\|_{W^{1,\infty}(0,L)}) < m, \quad \max_{i \in \mathbb{Z}} (|\dot{X}_i|) < m, \quad m < \infty,$$

we deduce the following estimate:

$$\begin{aligned} & \gamma \int_{\mathbb{R}} e^{-2\gamma t} \|(r, s)\|_2^2 dt + \int_{\mathbb{R}} e^{-2\gamma t} \sum_i |(r_i^\pm, s_i^\pm)|^2 + \varepsilon_i^2 + \varepsilon_i^2 dt \\ & \leq C \left(\frac{1}{\gamma} \int_{\mathbb{R}} e^{-2\gamma t} \|(F^1, F^2)\|_2^2 dt + \int_{\mathbb{R}} e^{-2\gamma t} (\|a_1\|_{l^2(\mathbb{Z})}^2 + \|a_2\|_{l^2(\mathbb{Z})}^2) dt \right) \end{aligned}$$

for any $\gamma \geq \gamma_0 > 0$, where γ_0 depends on $\|\delta\|_{l^\infty(\mathbb{Z})}$ and m . This completes the proof of the proposition. \square

3.3. The adjoint problem. In this section, we formulate the adjoint problem of the linearized equations around the roll-wave and around an approximate solution that is close to a roll-wave. We obtain energy estimates for that adjoint problem that shall be useful to prove the existence of weak solutions to the linearized equations.

3.3.1. The adjoint problem of the “exact” linearized equations. We multiply the first equation (resp. second) of (20) by p_i (resp. q_i), integrate over $(0, L) \times \mathbb{R}$, and sum on $i \in \mathbb{Z}$:

$$\begin{aligned} & \int_{\mathbb{R}} \sum_i \int_0^L \partial_t r_i p_i + \partial_t s_i q_i + a(h) \partial_x r_i p_i + b(h) \partial_x s_i q_i dx dt \\ & = \int_{\mathbb{R}} \sum_i \int_0^L F_i^1 p_i + F_i^2 q_i. \end{aligned}$$

An integration by parts on the interval $(0, L)$ and $t \in \mathbb{R}$ yields

$$\begin{aligned} & - \int_{\mathbb{R}} \sum_i \int_0^L (p_{i,t} + (a(h)p_i)_x) r_i + (q_{i,t} + (b(h)q_i)_x) s_i dx \\ & - \int_{\mathbb{R}} \sum_i [a(h)p r + b(h)q s]_{iL} = \int_{\mathbb{R}} \sum_i \int_0^L F_i^1 p_i + F_i^2 q_i dx dt. \end{aligned}$$

The boundary terms can be written

$$-[a(h)pr + b(h)qs]_{iL} = {}^t(p_i^\pm, q_i^\pm) A(r_i^\pm, s_i^\pm),$$

with $A = \text{diag}(-a_+, a_-, -b_+, b_-)$. We write the Rankine–Hugoniot conditions (21) in the form

$$\begin{pmatrix} \dot{\varepsilon}_i \\ 0 \end{pmatrix} = \underline{B}(r_i^\pm, s_i^\pm) + \begin{pmatrix} a_{1,i} \\ a_{2,i} \end{pmatrix},$$

with $\underline{\varepsilon}_i = 2[h]\varepsilon_i$ and \underline{B} is defined by

$$\underline{B} = \begin{pmatrix} h_+ - h_+^{-\frac{1}{2}} & -h_- + h_-^{-\frac{1}{2}} & h_+ + h_+^{\frac{1}{2}} & -h_- - h_-^{-\frac{1}{2}} \\ 2 - h_+^{\frac{3}{2}} - h_+^{-\frac{3}{2}} & h_-^{\frac{3}{2}} + h_-^{-\frac{3}{2}} - 2 & h_+^{\frac{3}{2}} + h_+^{-\frac{3}{2}} + 2 & -h_-^{\frac{3}{2}} - h_-^{-\frac{3}{2}} - 2 \end{pmatrix}.$$

We decompose \underline{B} in the form $\underline{B} = (A_1 \ A_2) = \begin{pmatrix} B_1 \\ B_2 \end{pmatrix}$, with $A_i \in M_2(\mathbb{R})$ and ${}^t B_i \in \mathbb{R}^4$. It is a straightforward computation to prove that

$$\det(A_1) = (h_- - h_-^{-\frac{1}{2}}) (h_+ - h_+^{-\frac{1}{2}}) \left((h_+^{\frac{1}{2}} - h_+^{-1}) - (h_-^{\frac{1}{2}} - h_-^{-1}) \right) \neq 0,$$

provided that $h_+ < 1 < h_-$. This condition is clearly satisfied for roll-wave solutions (otherwise, there is no roll-wave). Thus, $\text{Ker } \underline{B}$ is two-dimensional (2D) and possesses a basis of vectors: $\langle {}^t B_1^\pm, {}^t B_2^\pm \rangle$. We can apply Lemma 9.4, p. 255 in [1]: there exists $N_1 \in M_{2,4}(\mathbb{R})$ so that $\mathbb{R}^4 = \text{Ker } \underline{B} \oplus \text{Ker } N_1$ and there exist $C, N_2 \in M_{2,4}(\mathbb{R})$ so that $\mathbb{R}^4 = \text{Ker } C \oplus \text{Ker } N_2$ and $A = {}^t N_2 \underline{B} + {}^t C N_1$. As a consequence, the boundary terms can be written in the form

$$\begin{aligned} & \int_{\mathbb{R}} \sum_i {}^t(p_i^\pm, q_i^\pm) A(r_i^\pm, s_i^\pm) \\ &= \int_{\mathbb{R}} \sum_i {}^t(N_2(p_i^\pm, q_i^\pm)) \underline{B}(r_i^\pm, s_i^\pm) + {}^t(C(p_i^\pm, q_i^\pm)) N_1(r_i^\pm, s_i^\pm) \\ &= \int_{\mathbb{R}} \sum_i {}^t(C(p_i^\pm, q_i^\pm)) N_1(r_i^\pm, s_i^\pm) + {}^t(N_2(p_i^\pm, q_i^\pm)) \begin{pmatrix} a_{1,i} - \dot{\varepsilon}_i \\ a_{2,i} \end{pmatrix} \\ &= \int_{\mathbb{R}} \sum_i {}^t(N_2(p_i^\pm, q_i^\pm)) \begin{pmatrix} a_{1,i} \\ a_{2,i} \end{pmatrix} dt + \partial_t (N_2(p_i^\pm, q_i^\pm)_1) \varepsilon_i dt \\ & \quad + \int_{\mathbb{R}} \sum_i {}^t(C(p_i^\pm, q_i^\pm)) N_1(r_i^\pm, s_i^\pm) dt. \end{aligned}$$

The adjoint problem of the linearized equations around a roll-wave is written in the form

$$\begin{aligned} -\partial_t p_i - \partial_x(a(h)p_i) &= 0, \\ -\partial_t q_i - \partial_x(b(h)q_i) &= 0, \end{aligned} \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

with the boundary conditions

$$\partial_t (N_2(p_i^\pm, q_i^\pm)_1) = 0, \quad C(p_i^\pm, q_i^\pm) = 0 \quad \forall i \in \mathbb{Z}.$$

In what follows, we compute energy estimates for solutions of the linear problem:

$$(31) \quad \begin{aligned} -\partial_t p_i - \partial_x(a(h)p_i) &= F_i^1, \\ -\partial_t q_i - \partial_x(b(h)q_i) &= F_i^2, \end{aligned} \quad \forall i \in \mathbb{Z}, \quad x \in (0, L),$$

with the boundary equations

$$(32) \quad \partial_t (N_2 (p_i^\pm, q_i^\pm)_1) = a_{1,i}, \quad C (p_i^\pm, q_i^\pm) = a_{2,i} \quad \forall i \in \mathbb{Z}.$$

We multiply the first equation of (31) by p_i and the second by q_i :

$$(33) \quad -\partial_t \int_{\mathbb{R}} p^2 + q^2 dx + \sum_i [a(h)p^2 + b(h)q^2]_{iL} = 2 \int_{\mathbb{R}} F^1 p + F^2 q dxdt + \int_{\mathbb{R}} \partial_x a(h) p^2 + \partial_x b(h) q^2 dx,$$

where (p, q) is the function $(p, q) = \sum_{i \in \mathbb{Z}} (p_i, q_i) 1_{(iL, (i+1)L)}$. We note $\tilde{a}_{1,i}$, the unique exponentially decreasing solution of $\partial_t \tilde{a}_{1,i} = a_{1,i}$. The boundary conditions (32) now read

$$(34) \quad \begin{aligned} N_2 (p_i^\pm, q_i^\pm)_1 &= \tilde{a}_{1,i}, \\ N_2 (p_i^\pm, q_i^\pm)_2 &= N_2 (p_i^\pm, q_i^\pm)_2 \quad \forall i \in \mathbb{Z}, \\ C (p_i^\pm, q_i^\pm) &= a_{2,i}. \end{aligned}$$

We recall that $A = {}^t N_2 \underline{B} + {}^t C_1$, and (34) is written

$$A (p_i^\pm, q_i^\pm) = (N_2 (p_i^\pm, q_i^\pm)_2)^t B_2 + \tilde{a}_{1,i} {}^t B_1 + {}^t N_1 a_{i,2}.$$

The matrix A is invertible with a bounded inverse so that

$$|(p_i^\pm, q_i^\pm)|^2 \leq C \left(|N_2 (p_i^\pm, q_i^\pm)_2|^2 + |(\tilde{a}_{1,i}, 0, a_{2,i})|^2 \right).$$

As a consequence, there remains to estimate the term $N_2 (p_i^\pm, q_i^\pm)_2$. For that purpose, we write the boundary terms in (33) in the form

$$\begin{aligned} -\mathcal{Q} (p_i^\pm, q_i^\pm) &= -|N_2 (p_i^\pm, q_i^\pm)_2|^2 B_2 A^{-1} {}^t B_2 \\ &\quad - 2N_2 (p_i^\pm, q_i^\pm)_2 \left(A^{-1} {}^t B_2; A^{-1} ({}^t \underline{B}, {}^t N_1) \begin{pmatrix} \tilde{a}_{1,i} \\ 0 \\ a_{2,i} \end{pmatrix} \right) \\ &\quad - \left(A^{-1} ({}^t \underline{B}, {}^t N_1) \begin{pmatrix} \tilde{a}_{1,i} \\ 0 \\ a_{2,i} \end{pmatrix}; ({}^t \underline{B}, {}^t N_1) \begin{pmatrix} \tilde{a}_{1,i} \\ 0 \\ a_{2,i} \end{pmatrix} \right). \end{aligned}$$

Then, we deduce from (33) that

$$\begin{aligned} -\partial_t \int_{\mathbb{R}} p^2 + q^2 - \sum_i |N_2 (p_i^\pm, q_i^\pm)_2|^2 B_2 A^{-1} {}^t B_2 \\ = 2 \int_{\mathbb{R}} F^1 p + F^2 q + \mathcal{O} \left(\sum_i \tilde{a}_{1,i}^2 + |a_{2,i}|^2 \right) \\ + \mathcal{O} \left(\sum_i |N_2 (p_i^\pm, q_i^\pm)_2| |(\tilde{a}_{1,i}, 0, a_{2,i})| \right). \end{aligned}$$

Next, we show that $B_2 A^{-1} {}^t B_2 < 0$. It is an easy computation to prove that

$$B_2 A^{-1} {}^t B_2 = 2 \left(\left(3h_+ + \frac{1}{h_+^2} \right) - \left(3h_- + \frac{1}{h_-^2} \right) \right) < 0.$$

The inequality is satisfied, since we have already assumed that the Froude number F is such that $0 < F < F_c$, and thus

$$3h_- + \frac{1}{h_-} > 3h_+ + \frac{1}{h_+}$$

(see the proof of Proposition 1) so that the boundary conditions are maximal dissipative and obtain energy estimates on the linearized problem. We work in weighted spaces in time: integrating on $t \in \mathbb{R}$ and using Young’s inequalities, we get for $\gamma \geq \gamma_0$ (γ_0 depending on $\|(R, S)\|_{W^{1,\infty}(0,L)}$) that

$$(35) \quad \begin{aligned} & \gamma \int_{\mathbb{R}} e^{2\gamma t} \|(p, q)\|_2^2 dt + \kappa \int_{\mathbb{R}} e^{2\gamma t} \sum_i |(p_i^\pm, q_i^\pm)|^2 \\ & \leq C \left(\frac{1}{\gamma} \int_{\mathbb{R}} e^{2\gamma t} \|(F^1, F^2)\|^2 dt + \int_{\mathbb{R}} e^{2\gamma t} \sum_i |a_{1,i}|^2 + |a_{2,i}|^2 \right), \end{aligned}$$

with $\kappa, C > 0$ constant.

3.3.2. The “approximate” adjoint problem. In this section, we compute the adjoint problem of the linearized problem around an approximate solution and a priori estimates on that linear problem. Recall that the linearized equations around an approximate solution are given by

$$(36) \quad \begin{aligned} \partial_t r_i + \frac{L}{X_{i+1} - X_i} \Lambda_1^i \partial_x r_i &= F_i \\ \partial_t s_i + \frac{L}{X_{i+1} - X_i} \Lambda_2^i \partial_x s_i &= G_i \end{aligned} \quad \forall x \in (0, L), i \in \mathbb{Z},$$

with $\Lambda_k^i = \lambda_k(R^i, S^i) - (\dot{X}_i + \frac{x}{L}(\dot{X}_{i+1} - \dot{X}_i))$. We multiply the first (resp. second) equation of (36) by p_i (resp. q_i), integrate on $(0, L) \times \mathbb{R}$, and sum on $i \in \mathbb{Z}$:

$$(37) \quad \begin{aligned} \int_{\mathbb{R}} \sum_i \int_0^L F_i p_i + G_i q_i &= - \int_{\mathbb{R}} \sum_i \int_0^L \partial_t p_i r_i + \partial_t q_i s_i \\ &\quad - \int_{\mathbb{R}} \sum_i \int_0^L \partial_x (\Lambda_1^i p_i) r_i + \partial_x (\Lambda_2^i q_i) s_i \\ &\quad + \int_{\mathbb{R}} \sum_i {}^t (p_i^\pm, q_i^\pm) A_i (r_i^\pm, s_i^\pm), \end{aligned}$$

with

$$A_i = \text{diag} \left(-\frac{L}{X_{i+1} - X_i} \left(\lambda_1(R_i^+, S_i^+) - \dot{X}_i \right), \frac{L}{X_i - X_{i-1}} \left(\lambda_1(R_i^-, S_i^-) - \dot{X}_i \right), \right. \\ \left. -\frac{L}{X_{i+1} - X_i} \left(\lambda_2(R_i^+, S_i^+) - \dot{X}_i \right), \frac{L}{X_i - X_{i-1}} \left(\lambda_2(R_i^-, S_i^-) - \dot{X}_i \right) \right).$$

It is a straightforward computation to prove that $A_i = A + \mathcal{O}(\eta)$, where η is given by

$$\eta = \max_{i \in \mathbb{Z}} \left(|\dot{X}_i - x|_\infty, |X_{i+1} - X_i - L|_\infty, |(R - R_i, S - S_i)|_\infty \right).$$

We consider the boundary terms in (37): recall that, under the assumption of Proposition 2, the linearized Rankine–Hugoniot conditions can be written in the form

$$\begin{pmatrix} \dot{\xi}_i \\ 0 \end{pmatrix} = \underline{B}_i (r_i^\pm, s_i^\pm) + \begin{pmatrix} a_{1,i} \\ a_{2,i} \end{pmatrix}.$$

We can prove easily that $\underline{B}_i = \underline{B} + \mathcal{O}(\eta)$. Following the computation of the dual problem in the previous section, we decompose \underline{B}_i in the form $\underline{B}_i = (A_{1,i} \ A_{2,i}) = \begin{pmatrix} B_{1,i} \\ B_{2,i} \end{pmatrix}$, with $A_{k,i} \in M_2(\mathbb{R})$ and ${}^t B_{k,i} \in \mathbb{R}^4$. It is a straightforward computation to prove that, for η sufficiently small,

$$\det(A_{1,i}) = \det(A_1) + \mathcal{O}(\eta) \neq 0.$$

Thus, $\text{Ker } \underline{B}_i$ is 2D and possesses a basis of vectors: $\langle {}^t B_{1,i}^\perp, {}^t B_{2,i}^\perp \rangle$. We can apply Lemma 9.4, p. 255 in [1]: there exists $N_{1,i} \in M_{2,4}(\mathbb{R})$ so that $\mathbb{R}^4 = \text{Ker } \underline{B}_i \oplus \text{Ker } N_{1,i}$ and there exist $C_i, N_{2,i} \in M_{2,4}(\mathbb{R})$ so that $\mathbb{R}^4 = \text{Ker } C_i \oplus \text{Ker } N_{2,i}$ and $A_i = {}^t N_{2,i} \underline{B}_i + {}^t C_i N_{1,i}$. At this stage, the formulation of the adjoint problem is straightforward, and one finds

$$\begin{aligned} -\partial_t p_i - \frac{L}{X_{i+1} - X_i} \partial_x (\Lambda_i^1 p_i) &= 0, \\ -\partial_t q_i - \frac{L}{X_{i+1} - X_i} \partial_x (\Lambda_i^2 q_i) &= 0 \end{aligned} \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

with the boundary conditions

$$\partial_t (N_{2,i} (p_i^\pm, q_i^\pm)_1) = 0 \quad C_i (p_i^\pm, q_i^\pm) = 0 \quad \forall i \in \mathbb{Z}.$$

Next we consider the linear problem

$$(38) \quad \begin{aligned} -\partial_t p_i - \frac{L}{X_{i+1} - X_i} \partial_x (\Lambda_i^1 p_i) &= F_i^1, \\ -\partial_t q_i - \frac{L}{X_{i+1} - X_i} \partial_x (\Lambda_i^2 q_i) &= F_i^2 \end{aligned} \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

with the boundary conditions

$$(39) \quad \partial_t (N_{2,i} (p_i^\pm, q_i^\pm)_1) = a_{1,i} \quad C_i (p_i^\pm, q_i^\pm) = a_{2,i} \quad \forall i \in \mathbb{Z}.$$

Following the method of proof of the energy estimates for the adjoint problem associated to the linearization of shallow water equations around roll-waves, one can prove the following result.

PROPOSITION 3. *Under the hypothesis $0 < F < F_c$, there exists $\eta > 0$ so that for all $m > 0$, there exists a constant C and γ_0 depending on m and for any Lipschitz continuous and bounded function $(R_i, S_i, X_i)_{i \in \mathbb{Z}}$ satisfying*

$$\begin{aligned} \max_{i \in \mathbb{Z}} \left(\|\dot{X}_i - c\|_\infty + \|X_{i+1} - X_i - L\|_\infty + \|R_i - R, S_i - S\|_\infty \right) &\leq \eta, \\ \max_{i \in \mathbb{Z}} \left(|\partial_x (R_i, S_i)|_\infty + |\dot{X}_i|_\infty \right) &\leq m, \end{aligned}$$

and for all $\gamma \geq \gamma_0$, $(F^1, F^2) \in \mathcal{D}(\mathbb{R}/\{iL\} \times \mathbb{R})$, $(a_1, a_2) \in l^2(\mathbb{Z}, \mathcal{D}(\mathbb{R}))$, then a solution $(p_i, q_i)_{i \in \mathbb{Z}}$ of (38), (39) satisfies

$$\begin{aligned} &\gamma \|e^{\gamma t} (p, q)\|_{L^2(\mathbb{R}/\{iL\} \times \mathbb{R})}^2 + \|e^{-\gamma t} (p_i^\pm, q_i^\pm)\|_{l^2(\mathbb{Z}, L^2(\mathbb{R}))}^2 \\ &\leq C \left(\frac{1}{\gamma} \|e^{\gamma t} (F^1, F^2)\|_{L^2(\mathbb{R}/\{iL\} \times \mathbb{R})}^2 + \|e^{\gamma t} (a_1, a_2)\|_{l^2(\mathbb{Z}, L^2(\mathbb{R}))} \right). \end{aligned}$$

As a conclusion, we have obtained energy estimates for the linearized problem around an approximate solution, we have formulated an adjoint problem, and we computed estimates on that problem. This is the first step towards the well posedness of the linearized equations.

3.4. Well posedness of the linear differential boundary value problem.

In what follows, we prove the well posedness of the linear boundary value problem without any initial condition: The method of proof is based on a weak formulation of the linearized equations that define an adjoint problem. Using the Riesz representation theorem, we prove the existence of a weak solution that is proved to enjoy more regularity using a mollifying operator in the “transverse” t -direction. At this step, there is a difference with the shock case. In that latter case, the convection matrix in front of the x -derivatives is nonsingular, and we easily get estimates on those derivatives. In the case of roll-waves, the matrix that is in front of the x -derivatives is *singular* in the neighborhood of the sonic point. As a consequence, we also need a regularization step in the interior equations in the neighborhood of the sonic point. In what follows, the presentation of the well posedness results for a linearized problem follows the presentation made in [1] for the persistence of shocks: in the proofs of those results, we shall emphasize the difference between the shock case and the roll-wave case.

In order to simplify the notations, we introduce $\tilde{\Lambda}_i^k = \frac{L}{X_{i+1}-X_i} \Lambda_i^k$ and denote $\Omega = \mathbb{R} \setminus \{iL, i \in \mathbb{Z}\} \times \mathbb{R}$. We shall also note $L_\gamma^2(\Omega)$ the Hilbert space of functions f so that $e^{-\gamma t} f \in L^2(\Omega)$. We first prove the following result (see Theorem 12.4, p. 360 in [1] for the shock case).

PROPOSITION 4. *There exists $\eta > 0$ and $\gamma_0 = \gamma_0(m)$ such that if*

$$\begin{aligned} \max_{i \in \mathbb{Z}} \left(\|\dot{X}_i - c\|_\infty, \|X_{i+1} - X_i - L\|_\infty, \|R_i - \bar{R}\|_\infty, \|S_i - \bar{S}\|_\infty \right) &\leq \eta, \\ \max_{i \in \mathbb{Z}} \left(\|\dot{X}_{i+1} - \dot{X}_i\|_\infty, \|R_i\|_{W^{1,\infty}}, \|S_i\|_{W^{1,\infty}} \right) &\leq m, \end{aligned}$$

then for all $\gamma \geq \gamma_0$, $f \in L_\gamma^2(\Omega)^2$, $g \in l^2(\mathbb{Z}, L_\gamma^2(\mathbb{R}))^2$, there is one and only one $((r, s), \varepsilon) \in L_\gamma^2(\Omega)^2 \times l^2(\mathbb{Z}, H_\gamma^{\frac{1}{2}}(\mathbb{R}))$ such that

$$\mathcal{L}(r, s) = \begin{cases} \partial_t r_i + \tilde{\Lambda}_i^1 \partial_x r_i = f_i^1, \\ \partial_t s_i + \tilde{\Lambda}_i^2 \partial_x s_i = f_i^2 \end{cases} \quad \forall i \in \mathbb{Z}, (x, t) \in]0, L[\times \mathbb{R}$$

and satisfies the boundary conditions

$$\begin{pmatrix} \varepsilon_i \\ 0 \end{pmatrix} = \underline{B}_i (r_i^\pm, s_i^\pm) + \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} \quad \forall i \in \mathbb{Z}, t \in \mathbb{R}.$$

Furthermore, $(r_i^\pm, s_i^\pm)_{i \in \mathbb{Z}} \in l^2(\mathbb{Z}, L_\gamma^2(\mathbb{R})^4)$ and $\varepsilon \in l^2(\mathbb{Z}, H_\gamma^1(\mathbb{R}))$, and we have the estimate

$$\begin{aligned} \gamma \|(r, s)\|_{L_\gamma^2(\Omega)^2}^2 + \|(r^\pm, s^\pm)\|_{l^2(\mathbb{Z}, L_\gamma^2(\mathbb{R})^4)}^2 + \|\varepsilon\|_{l^2(\mathbb{Z}, H_\gamma^1(\mathbb{R}))}^2 \\ \leq C \left(\frac{1}{\gamma} \|f\|_{L_\gamma^2(\Omega)}^2 + \|g\|_{l^2(\mathbb{Z}, L_\gamma^2(\mathbb{R}))}^2 \right). \end{aligned}$$

Proof. Let us first prove the existence of a weak solution with a duality argument. With the notations of Proposition 3, denote \mathcal{E} the space of functions $(p, q) \in \mathcal{D}(\Omega)^2$ so that

$$\partial_t N_{2,i}^t(p_i^\pm, q_i^\pm)_1 - \gamma N_{2,i}^t(p_i^\pm, q_i^\pm)_1 = 0, \quad C_i(p_i^\pm, q_i^\pm) = 0.$$

We define the adjoint operator \mathcal{L}_γ^* of $\mathcal{L}_\gamma = \mathcal{L} + \gamma Id$ as

$$\mathcal{L}_\gamma^*(p, q) = \begin{cases} -\partial_t p_i - \partial_x (\tilde{\Lambda}_i^1 p_i) + \gamma p_i, \\ -\partial_t q_i - \partial_x (\tilde{\Lambda}_i^2 q_i) + \gamma q_i. \end{cases}$$

Let $(f, g) \in L^2(\Omega)^2 \times l^2(\mathbb{Z}, L^2(\mathbb{R}))$ and define a bounded linear form l on $\mathcal{L}_\gamma^* \mathcal{E}$ as

$$l(\mathcal{L}_\gamma^*(p, q)) = \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} \int_0^L p_i f_i^1 + q_i f_i^2 dx + (N_{2,i}(p_i^\pm, q_i^\pm)) \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} dt.$$

We deduce from Proposition 3 on the adjoint problem that $\forall (p, q) \in \mathcal{E}$,

$$\begin{aligned} |l(\mathcal{L}_\gamma^*(p, q))| &\leq C(\|f\|, \|g\|) (\|(p, q)\|_{L^2(\Omega)} + \|(p^\pm, q^\pm)\|_{l^2(\mathbb{Z}, L^2(\mathbb{R}))^4}) \\ &\leq C \left(\frac{1}{\gamma} \|f\| + \frac{1}{\sqrt{\gamma}} \|g\| \right) \|\mathcal{L}_\gamma^*(p, q)\|_{L^2(\Omega)}. \end{aligned}$$

Therefore, by the Hahn–Banach theorem, the linear form l extends to a continuous form on $L^2(\Omega)^2$ and using the Riesz representation theorem, there exists $(r, s) \in L^2(\Omega)^2$ so that

$$l(\mathcal{L}_\gamma^*(p, q)) = \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} \int_0^L (r_i, s_i)^t \mathcal{L}_\gamma^*(p_i, q_i) dx dt \quad \forall (p, q) \in \mathcal{E}.$$

Then, by the definition of l , we find that

$$\sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} \int_0^L (r_i, s_i)^t \mathcal{L}_\gamma^*(p_i, q_i) - (f_i^1, f_i^2)^t (p_i, q_i) dx dt = 0 \quad \forall (p, q) \in \mathcal{D}(\Omega)^2,$$

and $\mathcal{L}_\gamma(r, s) = (f^1, f^2)$ in the sense of distribution. For all $(p, q) \in \mathcal{E}$, one can prove that

$$\begin{aligned} \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} \int_0^L (r_i, s_i)^t \cdot \mathcal{L}_\gamma^*(p_i, q_i) - \mathcal{L}_\gamma(r_i, s_i)^t \cdot \begin{pmatrix} p_i \\ q_i \end{pmatrix} dx dt \\ = \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} N_{2,i}(p_i^\pm, q_i^\pm)^t \cdot \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} dt \\ = \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} N_{2,i}(p_i^\pm, q_i^\pm)^t \cdot \underline{B}_i \begin{pmatrix} r_i^\pm \\ s_i^\pm \end{pmatrix} dt. \end{aligned}$$

As a consequence, we find that

$$(40) \quad \sum_{i \in \mathbb{Z}} \int_{\mathbb{R}} N_{2,i}(p_i^\pm, q_i^\pm)^t \left(\underline{B}_i \begin{pmatrix} r_i^\pm \\ s_i^\pm \end{pmatrix} - \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} \right) dt = 0 \quad \forall (p, q) \in \mathcal{E}.$$

By density, this property holds true for all $(p, q) \in H^1(\Omega)$ such that

$$C_i^t(p_i^\pm, q_i^\pm) = 0, \quad \partial_t N_{2,i}(p_i^\pm, q_i^\pm)_1 - \gamma N_{2,i}(p_i^\pm, q_i^\pm)_1 = 0.$$

Let $(\theta^1, \theta^2) \in l^2(\mathbb{Z}, H^{\frac{1}{2}}(\mathbb{R})^2)$ and $(p_i^\pm, q_i^\pm)_{i \in \mathbb{Z}} \in l^2(\mathbb{Z}, H^{\frac{1}{2}}(\mathbb{R})^2)$ so that

$$N_{2,i}^t(p_i^\pm, q_i^\pm) = \begin{pmatrix} \theta_i^1 \\ \theta_i^2 \end{pmatrix}, \quad C_i^t(p_i^\pm, q_i^\pm) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \forall i \in \mathbb{Z}, t \in \mathbb{R}.$$

The sequence $(p_i^\pm, q_i^\pm)_{i \in \mathbb{Z}}$ exists, is unique, and lies in $l^2(\mathbb{Z}, H^{\frac{1}{2}}(\mathbb{R})^2)$, since the matrix $(N_{2,i}, C_i)$ is invertible under the assumption of Proposition 3. Now with a standard trace-lifting argument, choose $(p, q) \in H^1(\Omega)$ so that $p(iL)^\pm = p_i^\pm, q(iL)^\pm = q_i^\pm$. Then we deduce from (40) that for all $\theta \in l^2(\mathbb{Z}, H^{\frac{1}{2}}(\mathbb{R})^2)$ such that $\frac{d\theta_i^1}{dt} - \gamma\theta_i^1 = 0$,

$$\sum_{i \in \mathbb{Z}} \left\langle \begin{pmatrix} \theta_i^1 \\ \theta_i^2 \end{pmatrix}, \underline{B}_i(r_i^\pm, s_i^\pm) - \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} \right\rangle_{(H^{\frac{1}{2}}(\mathbb{R}), H^{-\frac{1}{2}}(\mathbb{R}))} = 0.$$

Thus, there exists $\varepsilon \in l^2(\mathbb{Z}, H^{\frac{1}{2}}(\mathbb{R}))$ so that

$$\begin{pmatrix} \dot{\varepsilon}_i + \gamma\varepsilon_i \\ 0 \end{pmatrix} = \underline{B}_i(r_i^\pm, s_i^\pm) + \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} \quad \forall i \in \mathbb{Z}.$$

Let us now prove that $(r, s), \varepsilon$ is a strong solution. This is just a matter of smoothing the solution and using the estimates of the first part. We consider a mollifying operator R_δ^t in the t -direction, and we define a regularization $(r^\delta, s^\delta), \varepsilon^\delta$ of (r, s, ε) . Let us first deal with the equation on s . Denote $s_i^\delta = R_\delta^t s_i$, we easily obtain

$$\partial_t s_i^\delta + \tilde{\Lambda}_i^2 \partial_x s_i^\delta = \tilde{\Lambda}_i^2 \left(R_\delta^t \left(\frac{f_i^2}{\tilde{\Lambda}_i^2} \right) + \left[\frac{1}{\tilde{\Lambda}_i^2} \partial_t, R_\delta^t \right] s_i \right).$$

Due to the presence of a *sonic point* $x_i(t)$ so that $\tilde{\Lambda}_1^i(x_i(t), t) = 0$, we cannot directly do the same regularization on r_i . In that case, we do not obtain any information on the derivatives of r_i in the x -direction. Nevertheless, under the assumptions of Proposition 3, there exists $L_1 > 0$ and a constant $C(m)$ so that for all $i \in \mathbb{Z}, iL + 2L_1 \leq x_i(t) \leq (i + 1)L - 2L_1$ and $|\tilde{\Lambda}_1^i(x, t)|^{-1} \leq C(m) \forall x \in [0, L_1] \cup [L - L_1, L]$. Let $\bar{\phi} \in C_c^\infty(0, L)$ so that $\bar{\phi} = 1$ on $[2L_1, L - 2L_1]$ and $\bar{\phi} = 0$ on $[0, L_1] \cup [L - L_1, L]$ and define $\phi = 1 - \bar{\phi}$. We multiply the equation on r_i by ϕ . We find

$$\partial_t(\phi r_i) + \tilde{\Lambda}_i^1 \partial_x(\phi r_i) = \phi f_i^1 + \tilde{\Lambda}_i^1 \partial_x \phi r_i.$$

Then we regularize that equation the same way we did with the equation on s_i :

$$\partial_t(\phi r_i)^\delta + \tilde{\Lambda}_i^1 \partial_x(\phi r_i)^\delta = \tilde{\Lambda}_i^1 \left(R_\delta^t \left(\frac{\phi f_i^1}{\tilde{\Lambda}_i^1} + \partial_x \phi r_i \right) + \left[\frac{1}{\tilde{\Lambda}_i^1} \partial_t, R_\delta^t \right] (\phi r_i) \right).$$

We multiply the equation on r_i by $(1 - \phi)$ and choose a mollifying operator R_δ^x in the x direction with $\delta \leq L_1$:

$$\begin{aligned} \partial_t R_\delta^x((1 - \phi)r_i) + \tilde{\Lambda}_i^1 \partial_x R_\delta^x((1 - \phi)r_i) &= R_\delta^x \left((1 - \phi) f_i^1 - \tilde{\Lambda}_i^1 \partial_x \phi r_i \right) \\ &+ \left[\tilde{\Lambda}_i^1 \partial_x, R_\delta^x \right] ((1 - \phi)r_i). \end{aligned}$$

We define r^δ so that $r_{[[iL, (i+1)L]]}^\delta = R_\delta^t(\phi r_i) + R_\delta^x((1 - \phi)r_i)$. By construction, both r^δ and s^δ lie in $H^1(\Omega)$. On the other hand, we find that r^δ, s^δ , and ε^δ satisfy the boundary conditions

$$\begin{pmatrix} \dot{\varepsilon}_i^\delta \\ 0 \end{pmatrix} = \underline{B}_i(r_i^{\delta, \pm}, s_i^{\delta, \pm}) + \begin{pmatrix} R_\delta^t g_i^1 \\ R_\delta^t g_i^2 \end{pmatrix} + [\underline{B}_i, R_\delta^t](r_i^\pm, s_i^\pm).$$

We can apply Proposition 2 as follows:

$$(41) \quad \begin{aligned} & \gamma \| (r^\delta, s^\delta) \|_{L^2_\gamma(\Omega)^2}^2 + \| (r^{\delta,\pm}, s^{\delta,\pm}) \|_{l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}^4))}^2 + \| \varepsilon^\delta \|_{l^2(\mathbb{Z}, H^1_\gamma(\mathbb{R}))}^2 \\ & \leq C \left(\frac{1}{\gamma} \| f^\delta \|_{L^2_\gamma(\Omega)}^2 + \| g^\delta \|_{l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}))}^2 \right), \end{aligned}$$

with

$$\begin{aligned} f_{i,1}^\delta &= R_\delta^x \left((1 - \phi) f_i^1 - \tilde{\Lambda}_i^1 \partial_x \phi r_i \right) + \left[\tilde{\Lambda}_i^1 \partial_x, R_\delta^x \right] \left((1 - \phi) r_i \right) \\ &\quad + \tilde{\Lambda}_i^1 \left(R_\delta^t \left(\frac{\phi f_i^1}{\tilde{\Lambda}_i^1} + \partial_x \phi r_i \right) + \left[\frac{1}{\tilde{\Lambda}_i^1} \partial_t, R_\delta^t \right] (\phi r_i) \right), \\ f_{i,2}^\delta &= \tilde{\Lambda}_i^2 \left(R_\delta^t \left(\frac{f_i^2}{\tilde{\Lambda}_i^2} \right) + \left[\frac{1}{\tilde{\Lambda}_i^2} \partial_t, R_\delta^t \right] (s) \right) \\ g_i^\delta &= \begin{pmatrix} R_\delta^t g_i^1 \\ R_\delta^t g_i^2 \end{pmatrix} + [\underline{B}_i, R_\delta^t] (r_i^\pm, s_i^\pm). \end{aligned}$$

One can prove the following convergence properties (see [1] for more details):

$$\lim_{\delta \rightarrow 0} \| f^\delta - f \|_{L^2_\gamma(\Omega)} + \| g^\delta - g \|_{l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}))} = 0.$$

By linearity of the equations, these estimates also apply to

$$r^\delta - r^{\delta'}, \quad s^\delta - s^{\delta'}, \quad \varepsilon^\delta - \varepsilon^{\delta'}.$$

Together with the convergence properties, this shows that $r^{\delta,\pm}, s^{\delta,\pm}$, and ε^δ are, respectively, Cauchy sequences in $l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}^4))$ and $l^2(\mathbb{Z}, H^1_\gamma(\mathbb{R}))$. By uniqueness of the limit in the sense of distribution, the sequences $(r^{\delta,\pm}, s^{\delta,\pm})$ and ε^δ converge to r^\pm, s^\pm , and ε in these norms. Then passing to the limit $\delta \rightarrow 0$ in (41), one finds

$$\begin{aligned} & \gamma \| (r, s) \|_{L^2_\gamma(\Omega)^2}^2 + \| (r^\pm, s^\pm) \|_{l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}^4))}^2 + \| \varepsilon \|_{l^2(\mathbb{Z}, H^1_\gamma(\mathbb{R}))}^2 \\ & \leq C \left(\frac{1}{\gamma} \| f \|_{L^2_\gamma(\Omega)}^2 + \| g \|_{l^2(\mathbb{Z}, L^2_\gamma(\mathbb{R}))}^2 \right). \end{aligned}$$

This completes the proof of the proposition. \square

We can prove more regularity on the solution (r, s, ε) under the condition that f, g , and the coefficients R_i, S_i, X_i enjoy more regularity. Similarly to the L^2 estimates, we introduce weighted Sobolev norms:

$$\| u \|_{\mathcal{H}_\gamma^k} = \sum_{|\alpha| \leq k} \gamma^{2(k-|\alpha|)} \| e^{-\gamma t} u \|_{L^2}.$$

The regularity result on (r, s, ε) is formulated as follows (see Theorem 12.5, p. 364 in [1] in the shock case).

PROPOSITION 5. *Under the assumptions of Proposition 4 and that R_i, S_i, X_i enjoy more regularity,*

$$\begin{aligned} & \| R - R_i, S - S_i \|_{l^2(\mathbb{Z}, H^k((0,L) \times \mathbb{R}))} \leq \mu, \\ & \| R_i^\pm - R^\pm, S_i^\pm - S^\pm \|_{l^2(\mathbb{Z}, H^k(\mathbb{R}))} \leq \mu, \\ & \| \dot{X}_i - c, X_{i+1} - X_i - L \|_{l^2(\mathbb{Z}, H^k(\mathbb{R}))} \leq \mu, \end{aligned}$$

with $k > 2$, the solution (r, s, ε) found in Proposition 4 is such that $(r, s) \in H^k(\Omega)$, $(r^\pm, s^\pm) \in l^2(\mathbb{Z}, H^k(\mathbb{R}))$, and $\varepsilon \in l^2(\mathbb{Z}, H^{k+1}(\mathbb{R}))$ and, for any $\gamma \geq \gamma_k(\mu) \geq 1$, satisfy the estimate

$$(42) \quad \begin{aligned} & \gamma \|(r, s)\|_{\mathcal{H}_\gamma^k(\Omega)}^2 + \|(r^\pm, s^\pm)\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^k(\mathbb{R}))}^2 + \|\varepsilon\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^{k+1}(\Omega))}^2 \\ & \leq C(k, \mu) \left(\frac{1}{\gamma} \|f\|_{\mathcal{H}_\gamma^k(\Omega)}^2 + \|g\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^k(\mathbb{R}))}^2 \right). \end{aligned}$$

Proof. We prove that proposition in two steps. The first step is to extend the a priori L^2 estimates obtained in Proposition 2. For all (r, s) and ε compactly supported and smooth functions, the following estimate holds true:

$$(43) \quad \begin{aligned} & \gamma \|(r, s)\|_{\mathcal{H}_\gamma^k(\Omega)}^2 + \|(r^\pm, s^\pm)\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^k(\mathbb{R}))}^2 + \|\varepsilon\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^{k+1}(\Omega))}^2 \\ & \leq C(k, \mu) \left(\frac{1}{\gamma} \|f\|_{\mathcal{H}_\gamma^k(\Omega)}^2 + \|g\|_{l^2(\mathbb{Z}, \mathcal{H}_\gamma^k(\mathbb{R}))}^2 \right), \end{aligned}$$

where $f = \mathcal{L}(r, s)$ and g is defined by

$$g_i = \begin{pmatrix} \dot{\varepsilon}_i \\ 0 \end{pmatrix} - \underline{B}_i(r_i^\pm, s_i^\pm).$$

In order to prove that estimate, we decompose in the form $r_i = \phi r_i + (1 - \phi)r_i$, with ϕ defined in Proposition 4. Deriving the equation satisfied by $\phi r_i, s_i$ in the t -direction and applying Proposition 2, we get estimates for derivatives in the t -direction. We recover the derivatives in the x -direction with

$$\partial_x(\phi r_i) = \frac{1}{\lambda_i^1} (f_i^1 + \partial_x \phi r_i - \partial_t(\phi r_i)), \quad \partial_x s_i = \frac{1}{\Lambda_i^2} (f_i^2 - \partial_t s_i).$$

The obtention of estimates for $(1 - \phi)r_i$ is straightforward, since it is trace-free on the discontinuity points $iL, i \in \mathbb{Z}$.

Now we come back to the solution obtained in Proposition 4. The next step is to prove enough regularity on (r, s, ε) in order to apply the a priori estimate (43). We first consider that the coefficients $R_i - R, S_i - S, X_i - (ct + iL)$ are smooth and compactly supported: the regularity of (r, s, ε) is then obtained by induction (see [1, p. 366] for the details of the proof). We then construct a family $(r^\delta, s^\delta, \varepsilon^\delta)$ of compactly supported functions that converges to (r, s, ε) in Sobolev H_γ^k norms. The regularized functions $(r^\delta, s^\delta, \varepsilon^\delta)$ satisfy the a priori estimate (43) and passing to the limit $\delta \rightarrow 0$ as in Proposition 4, we prove that (r, s, ε) satisfies the energy estimates (42). For the more general case, $R_i - R, S_i - S, X_i - (ct + iL)$ lying in Sobolev spaces, it is just a matter of smoothing the coefficients and passing to the limit. This completes the proof of that proposition. \square

Next, we consider the well posedness of the linearized equation with zero initial data. Let $T \in \mathbb{R}, I_T$ the interval $I_T =]-\infty, T]$, and $\Omega_T = \mathbb{R} \setminus \{iL, i \in \mathbb{Z}\} \times I_T$. One can prove that, under the assumptions of Proposition 4 and for all $f \in L^2(\mathbb{R} \setminus \{iL, i \in \mathbb{Z}\} \times I_T, \mathbb{R}^2)$ and $g \in l^2(\mathbb{Z}, L^2(I_T)^2)$ such that $f|_{t < 0} = 0, g|_{t < 0} = 0$, there exists a unique function $(r, s) \in L^2(\Omega_T)$ and $\varepsilon \in l^2(\mathbb{Z}, H^1(I_T))$ so that

$$\begin{aligned} \partial_t r_i + \tilde{\Lambda}_i^1 \partial_x r_i &= f_i^1, \\ \partial_t s_i + \tilde{\Lambda}_i^2 \partial_x s_i &= f_i^2, \end{aligned} \quad \forall i \in \mathbb{Z}, (x, t) \in]0, L[\times I_T,$$

and

$$\begin{pmatrix} \dot{\varepsilon}_i \\ 0 \end{pmatrix} = \underline{B}_i (r_i^\pm, s_i^\pm) + \begin{pmatrix} g_i^1 \\ g_i^2 \end{pmatrix} \quad \forall i \in \mathbb{Z}, t \in I_T.$$

Furthermore, we have $r|_{t<0} = s|_{t<0} = 0, \quad \varepsilon|_{t<0} = 0, (r_i^\pm, s_i^\pm)_{i \in \mathbb{Z}}$ belongs $l^2(\mathbb{Z}, L_\gamma^2(I_T)^4)$, and the solution (r, s, ε) satisfies the following estimate,

$$\begin{aligned} \gamma \|(r, s)\|_{L_\gamma^2(\Omega_T)^2}^2 + \|(r^\pm, s^\pm)\|_{l^2(\mathbb{Z}, L_\gamma^2(I_T)^4)}^2 + \|\varepsilon\|_{l^2(\mathbb{Z}, H_\gamma^1(I_T))}^2 \\ \leq C \left(\frac{1}{\gamma} \|f\|_{L_\gamma^2(\Omega_T)}^2 + \|g\|_{l^2(\mathbb{Z}, L_\gamma^2(I_T))}^2 \right). \end{aligned}$$

In the case of smoother coefficients, the solution (r, s, ε) of the linearized problem with zero initial data enjoys more regularity. One can prove the following result that shall be useful in the proof of the well posedness of the Cauchy problem for the full nonlinear problem.

PROPOSITION 6. *Assume that the hypotheses of Proposition 4 are satisfied, and the coefficients enjoy the following regularity properties:*

$$\begin{aligned} \|(R_i - R, S_i - S)\|_{l^2(\mathbb{Z}, H^k((0,L) \times I_T))} &\leq \mu, \\ \|(R_i^\pm - R^\pm, S_i^\pm - S^\pm)\|_{l^2(\mathbb{Z}, H^k(I_T))} &\leq \mu, \\ \|\dot{X}_i - c, X_{i+1} - X_i - L\|_{l^2(\mathbb{Z}, H^k(I_T))} &\leq \mu. \end{aligned}$$

Furthermore, suppose that, for some $\tau < T$,

$$R_i - R|_{t<\tau} = S_i - S|_{t<\tau} = 0, \quad (\dot{X}_i - c, X_{i+1} - X_i - L)|_{t<\tau} = 0.$$

Then, for any $f \in H^k(\Omega_T)$ and $g \in l^2(\mathbb{Z}, H^k(I_T))$ such that $f|_{t<0} = 0, \quad g|_{t<0} = 0$, the solution (r, s) belongs to $H^k(\Omega_T)$, (r^\pm, s^\pm) belongs to $l^2(\mathbb{Z}, H^k(I_T))$ and ε to $l^2(\mathbb{Z}, H^{k+1}(I_T))$ with the estimate

$$\begin{aligned} \frac{1}{T} \|(r, s)\|_{H^k(\Omega_T)}^2 + \|r^\pm, s^\pm\|_{l^2(\mathbb{Z}, H^k(I_T))}^2 + \|\varepsilon\|_{l^2(\mathbb{Z}, H^{k+1}(I_T))}^2 \\ \leq C(k, \mu) \left(T \|f\|_{H^k(\Omega_T)}^2 + \|g\|_{l^2(\mathbb{Z}, H^k(I_T))}^2 \right). \end{aligned}$$

4. The nonlinear problem. In this section, we solve the full nonlinear Cauchy problem with prescribed initial data (r_0, s_0, X_0) . To actually solve that problem, we need compatibility conditions on (r_0, s_0, X_0) . Under suitable compatibility conditions, we shall prove the existence of an approximate solution $(r^{(a)}, s^{(a)}, X^{(a)})$ on a time interval $(0, T^0)$ for a sufficiently small interval. We can prove the well posedness of the full Cauchy problem through an iterative scheme based on the approximate linear problem introduced in the previous section. That scheme is proved to converge to a solution of the Cauchy problem with the energy estimates obtained previously.

4.1. Construction of an approximate solution. In what follows, we shall construct an approximate solution of the shallow water equations coupled with the Rankine–Hugoniot jump condition for an initial data (r_0, s_0, X_0) that satisfies some compatibility conditions. In what follows, we write compatibility conditions for the shallow water equations.

4.1.1. Compatibility conditions. After the change of variable introduced in section 2, the shallow water system is written as

$$(44) \quad \begin{aligned} \partial_t r_i + \Lambda_i^1 \partial_x r_i &= Q(r_i, s_i), \\ \partial_t s_i + \Lambda_i^2 \partial_x s_i &= Q(r_i, s_i) \end{aligned} \quad \forall (x, i) \in (0, L) \times \mathbb{Z},$$

where the Λ_i^k are defined by

$$\Lambda_i^k = \frac{L}{X_{i+1} - X_i} \left(\lambda_k(r_i, s_i) - \left(\dot{X}_i + \frac{x}{L} (\dot{X}_{i+1} - \dot{X}_i) \right) \right),$$

whereas the Rankine–Hugoniot jump conditions are given by

$$(45) \quad \dot{X}_i = F(r_i^+, r_i^-, s_i^+, s_i^-), \quad G(r_i^+, r_i^-, s_i^+, s_i^-) = 0 \quad \forall i \in \mathbb{Z},$$

with

$$\begin{aligned} F(r_i^\pm, s_i^\pm) &= [(r - s)^4 + 8(r - s)^2(r + s)^2]_{iL} [(r - s)^2]_{iL} - 8[(r + s)(r - s)^2]_{iL}^2, \\ G(r_i^\pm, s_i^\pm) &= \frac{1}{2} [(r + s)(r - s)^2]_{iL} [(r - s)^2]_{iL}^{-1}. \end{aligned}$$

Suppose that both (r, s) and X are smooth enough: applying Faa di Bruno’s formula to the jump conditions (45), we obtain

$$(46) \quad \frac{d^{p+1} X_i}{dt^{p+1}} = \sum_{i=1}^p \sum_{i_1 + \dots + i_m = p} c_{i_1 \dots i_m} d^m F(r_i^\pm, s_i^\pm) \cdot \left(\frac{d^{i_1}}{dt^{i_1}} \begin{pmatrix} r_i^\pm \\ s_i^\pm \end{pmatrix}, \dots, \frac{d^{i_m}}{dt^{i_m}} \begin{pmatrix} r_i^\pm \\ s_i^\pm \end{pmatrix} \right).$$

Furthermore, we derive the interior equations (44) with the Faa di Bruno formula:

$$(47) \quad \begin{aligned} \partial_t^{p+1} r_i &= \partial_t^p Q(r_i, s_i) - \sum_{l=0}^p C_p^l \partial_t^l \Lambda_i^1 \partial_x \partial_t^{p-l} r_i, \\ \partial_t^{p+1} s_i &= \partial_t^p Q(r_i, s_i) - \sum_{l=0}^p C_p^l \partial_t^l \Lambda_i^2 \partial_x \partial_t^{p-l} s_i. \end{aligned}$$

It is a straightforward computation to prove that

$$\partial_t^p Q(r_i, s_i) = \sum_{m=1}^p \sum_{i_1 + \dots + i_m = p} c_{i_1 \dots i_m} d^m Q(r_i, s_i) \cdot \left(\partial_t^{i_1} \begin{pmatrix} r_i \\ s_i \end{pmatrix}, \dots, \partial_t^{i_m} \begin{pmatrix} r_i \\ s_i \end{pmatrix} \right).$$

Moreover, for $l = 1, \dots, p - 1$, we easily obtain

$$\begin{aligned} \partial_t^l \Lambda_i^k &= \sum_{m=1}^l \sum_{i_1 + \dots + i_m = l} c_{i_1 \dots i_m} d_{1,4}^m \Lambda_i^k \cdot \left(\partial_t^{i_1} \begin{pmatrix} r_i \\ s_i \\ X_i \\ X_{i+1} \end{pmatrix}, \dots, \partial_t^{i_m} \begin{pmatrix} r_i \\ s_i \\ X_i \\ X_{i+1} \end{pmatrix} \right) \\ &+ \sum_{m=1}^l \sum_{i_1 + \dots + i_m = l} c_{i_1 \dots i_m} d_{5,6}^m \Lambda_i^k \cdot \left(\partial_t^{i_1} \begin{pmatrix} \dot{X}_i \\ \dot{X}_{i+1} \end{pmatrix}, \dots, \partial_t^{i_m} \begin{pmatrix} \dot{X}_i \\ \dot{X}_{i+1} \end{pmatrix} \right), \end{aligned}$$

where $d_{1,4} \Lambda_i^k$ denotes the differential of Λ_i^k with respect to (r, s, X_i, X_{i+1}) and $d_{5,6} \Lambda_i^k$ denotes the differential with respect to $(\dot{X}_i, \dot{X}_{i+1})$. Here, the derivatives in X_i are no

larger than p . For $l = p$, one finds

$$\begin{aligned} \partial_t^p \Lambda_k^i &= \sum_{m=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d_{1,4}^m \Lambda_k \cdot \left(\partial_t^{i_1} \begin{pmatrix} r_i \\ s_i \\ X_i \\ X_{i+1} \end{pmatrix}, \dots, \partial_t^{i_m} \begin{pmatrix} r_i \\ s_i \\ X_i \\ X_{i+1} \end{pmatrix} \right) \\ &+ \sum_{m=2}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d_{5,6}^m \Lambda_k \cdot \left(\partial_t^{i_1} \begin{pmatrix} \dot{X}_i \\ \dot{X}_{i+1} \end{pmatrix}, \dots, \partial_t^{i_m} \begin{pmatrix} \dot{X}_i \\ \dot{X}_{i+1} \end{pmatrix} \right) \\ &+ \nabla_{5,6} \Lambda_k \cdot \partial_t^{p+1} \begin{pmatrix} X_i \\ X_{i+1} \end{pmatrix}. \end{aligned}$$

The derivatives of X_i up to order $p + 1$ are replaced by the expressions found in (46) involving derivatives of order $j \leq p$. Following the definitions of compatibility conditions introduced in [1, p. 370] (see the definition of “compatibility conditions to order s ”), we define the compatibility conditions for a pair of initial data (r_0, s_0, X_0) .

We fix (r_0, s_0, X_0) an initial data and define the functions (r_p, s_p) and X_p so that

$$(48) \quad \begin{aligned} X_{1,i} &= F(r_{0,i}^\pm, s_{0,i}^\pm), \\ r_{1,i} &= Q(r_{0,i}, s_{0,i}) - \Lambda_{0,i}^1 \partial_x r_{0,i}, \\ s_{1,i} &= Q(r_{0,i}, s_{0,i}) - \Lambda_{0,i}^2 \partial_x s_{0,i}, \end{aligned}$$

and

$$(49) \quad \begin{aligned} X_{p+1,i} &= \sum_{i=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d^m F(r_{0,i}^\pm, s_{0,i}^\pm) \cdot \left(\begin{pmatrix} r_{i_1,i}^\pm \\ s_{i_1,i}^\pm \end{pmatrix}, \dots, \begin{pmatrix} r_{i_m,i}^\pm \\ s_{i_m,i}^\pm \end{pmatrix} \right), \\ r_{p+1,i} &= \sum_{m=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d^m Q(r_{0,i}, s_{0,i}) \cdot \left(\begin{pmatrix} r_{i_1,i} \\ s_{i_1,i} \end{pmatrix}, \dots, \begin{pmatrix} r_{i_m,i} \\ s_{i_m,i} \end{pmatrix} \right) \\ &\quad - \sum_{l=0}^p C_p^l \Lambda_{l,i}^1 \partial_x r_{p-l,i}, \\ s_{p+1,i} &= \sum_{m=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d^m Q(r_{0,i}, s_{0,i}) \cdot \left(\begin{pmatrix} r_{i_1,i} \\ s_{i_1,i} \end{pmatrix}, \dots, \begin{pmatrix} r_{i_m,i} \\ s_{i_m,i} \end{pmatrix} \right) \\ &\quad - \sum_{l=0}^p C_p^l \Lambda_{l,i}^2 \partial_x s_{p-l,i}, \end{aligned}$$

where $\Lambda_{l,i}^k$ are defined, for any $l = 0, \dots, p - 1$, by

$$\begin{aligned} \Lambda_{0,i}^k &= \frac{L}{X_{0,i+1} - X_{0,i}} \left(\lambda_k(r_{0,i}, s_{0,i}) - \left(X_{1,i} + \frac{x}{L} (X_{1,i+1} - X_{1,i}) \right) \right), \\ \Lambda_{l,i}^k &= \sum_{m=1}^l \sum_{i_1+\dots+i_m=l} c_{i_1\dots i_m} d_{1,4}^m \Lambda_i^k \cdot \left(\begin{pmatrix} r_{i_1,i} \\ s_{i_1,i} \\ X_{i_1,i} \\ X_{i_1,i+1} \end{pmatrix}, \dots, \begin{pmatrix} r_{i_m,i} \\ s_{i_m,i} \\ X_{i_m,i} \\ X_{i_m,i+1} \end{pmatrix} \right) \\ &\quad + \sum_{m=1}^l \sum_{i_1+\dots+i_m=l} c_{i_1\dots i_m} d_{5,6}^m \Lambda_i^k \cdot \left(\begin{pmatrix} X_{i_1+1,i} \\ X_{i_1+1,i+1} \end{pmatrix}, \dots, \begin{pmatrix} X_{i_m+1,i} \\ X_{i_m+1,i+1} \end{pmatrix} \right), \end{aligned}$$

and, for $l = p$,

$$\begin{aligned} \Lambda_{p,i}^k &= \sum_{m=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d_{1,4}^m \Lambda_i^k \cdot \left(\left(\begin{matrix} r_{i_1,i} \\ s_{i_1,i} \\ X_{i_1,i} \\ X_{i_1,i+1} \end{matrix} \right), \dots, \left(\begin{matrix} r_{i_m,i} \\ s_{i_m,i} \\ X_{i_m,i} \\ X_{i_m,i+1} \end{matrix} \right) \right) \\ &+ \sum_{m=2}^p \sum_{i_1+\dots+i_m=p} c_{i_1\dots i_m} d_{5,6}^m \Lambda_i^k \cdot \left(\left(\begin{matrix} X_{i_1+1,i} \\ X_{i_1+1,i+1} \end{matrix} \right), \dots, \left(\begin{matrix} X_{i_m+1,i} \\ X_{i_m+1,i+1} \end{matrix} \right) \right) \\ &+ \nabla_{5,6} \Lambda_i^k \cdot \left(\begin{matrix} X_{p+1,i} \\ X_{p+1,i+1} \end{matrix} \right). \end{aligned}$$

For the definition of the sequence (r_p, s_p, X_p) , we have not used the jump condition $G(r_i^\pm, s_i^\pm) = 0$. In order to obtain a solution to the Cauchy problem, the sequence (r_p, s_p, X_p) shall satisfy some compatible conditions. We define here those compatibility conditions.

DEFINITION 1. A pair of initial data (r_0, s_0) and X_0 is said to be compatible up to order s if, for all $p \in \{0, \dots, s\}$, the functions (r_p, s_p, X_p) defined by (49) satisfy

$$0 = \sum_{m=1}^p \sum_{i_1+\dots+i_m=p} c_{i_1,\dots,i_m} d^m G(r_{0,i}^\pm, s_{0,i}^\pm) \cdot \left(\left(\begin{matrix} r_{i_1,i}^\pm \\ s_{i_1,i}^\pm \end{matrix} \right), \dots, \left(\begin{matrix} r_{i_m,i}^\pm \\ s_{i_m,i}^\pm \end{matrix} \right) \right).$$

4.1.2. Approximate solution. The purpose of this section is to prove that, under the compatibility conditions defined above, one can construct an approximate solution of the full Cauchy problem. Denote $\Omega = \mathbb{R} \setminus \{iL, i \in \mathbb{Z}\}$. We show the following result.

PROPOSITION 7. Let (r_0, s_0) and X_0 be an initial data so that $(r_0, s_0) \in (R, S) + H^{m+\frac{1}{2}}(\Omega)$ and $(X_{0,i} - iL)_{i \in \mathbb{Z}} \in l^2(\mathbb{Z})$, with $m > 2$ an integer that is compatible up to order $m - 1$ and such that

$$\|(r_0, s_0) - (R, S)\|_{L^\infty(\Omega)} \leq \rho, \quad \max_{i \in \mathbb{Z}} (|X_{0,i+1} - X_{0,i} - L|) \leq \rho.$$

Then for ρ sufficiently small, there exists $T_0 > 0$ and

$$(r_a, s_a) \in (R, S) + H^{m+1}(\Omega), \quad (X_{a,i} - (ct + iL)) \in l^2(\mathbb{Z}, H^{m+1}(\mathbb{R})),$$

with both $(r_a, s_a) - (R, S)$ and $X_{a,i} - (ct + iL)$ vanishing for $|t| \geq 2T_0$, so that $(r_a, s_a)|_{t=0} = (r_0, s_0)$, $X_{a,i}(0) = X_{0,i}$ and for all $(x, t) \in \Omega \times [-T_0, T_0]$,

$$\|(r_a - r_0, s_a - s_0)\| \leq \frac{\rho}{2}, \quad |X_{a,i} - X_{0,i}| \leq \frac{\rho}{2}.$$

Furthermore, the functions f_a and g_a defined as

$$\begin{aligned} f_{a,i}^1 &= \partial_t r_{a,i} + \Lambda_{a,i}^1 \partial_x r_{a,i} - Q(r_{a,i}, s_{a,i}), \\ f_{a,i}^2 &= \partial_t s_{a,i} + \Lambda_{a,i}^2 \partial_x s_{a,i} - Q(r_{a,i}, s_{a,i}), \end{aligned}$$

and $g_{a,i}^1 = \dot{X}_{a,i} - F(r_{a,i}^\pm, s_{a,i}^\pm)$, $g_{a,i}^2 = G(r_{a,i}^\pm, s_{a,i}^\pm)$ are such that

$$\partial_t^p f_a|_{t=0} = 0, \quad \partial_t^p g_a|_{t=0} = 0 \quad \forall p \in \{0, \dots, m - 1\}.$$

The function f_a belongs to $H^m(\Omega \times \mathbb{R})$, g_a lies in $l^2(\mathbb{Z}, H^m(\mathbb{R}))$, and both vanish for $|t| \geq 2T_0$.

Proof. For $k = 1 \dots m - 1$, we construct a sequence $(r_k, s_k) \in H^{m+\frac{1}{2}-k}(\Omega)$, X_k so that $X_{1,i} - c \in l^2(\mathbb{Z})$ and $X_{k,i} \in l^2(\mathbb{Z})$ ($k \geq 2$) that satisfies the compatibility conditions (49). Then, by trace lifting, we find $(r_a, s_a) \in (R, S) + H^{m+1}(\Omega \times \mathbb{R})$ and $X_a \in (ct + iL) + l^2(\mathbb{Z}, H^{m+1}(\mathbb{R}))$ so that

$$\begin{aligned} \|r_a - R, s_a - S\|_{H^{m+1}(\Omega \times \mathbb{R})} &\leq C\|r_0 - R, s_0 - S\|_{H^{m+1}(\Omega)}, \\ \|X_a - (ct + iL)\|_{l^2(\mathbb{Z}, H^{m+1}(\mathbb{R}))} &\leq C\|X_0 - iL\|_{l^2(\mathbb{Z})}, \end{aligned}$$

and $\partial_t^k(r_{a,i}, s_{a,i})|_{t=0} = (r_{k,i}, s_{k,i})$, $\partial_t^k(X_{a,i})|_{t=0} = X_{k,i} \forall k \in \{0, \dots, m-1\}$. By Sobolev embeddings, we find $T_0 > 0$ so that

$$\|(r_a - r_0, s_a - s_0)\| \leq \frac{\rho}{2}, \quad \|X_a - X_0\| \leq \frac{\rho}{2}$$

for $|t| \leq T_0$. In order to obtain functions f_a, g_a that vanish for $|t| \geq 2T_0$, we choose $\phi \in \mathcal{D}(\mathbb{R})$ a cut-off function so that $\phi(t) = 1$ if $|t| \leq T_0$ and $\phi(t) = 0$ if $|t| \geq 2T_0$. Then, substitute (r_a, s_a) and X_a , respectively, by $\phi(r_a, s_a) + (1 - \phi)(R, S)$ and $\phi X_{a,i} + (1 - \phi)(ct + iL)$. Then, one can prove that f_a belongs to $H^m(\Omega \times \mathbb{R})$ and $g^{(a)}$ lies in $l^2(\mathbb{Z}, H^m(\mathbb{R}))$, and both vanish for $|t| \geq 2T_0$. Finally, according to the compatibility conditions that are satisfied up to order $m - 1$, the derivatives of these functions vanish at $t = 0$. The proof of the proposition is then completed. \square

4.2. The fixed point argument. In this section, we prove the main theorem of this paper on the persistence of roll-waves.

THEOREM 2. *Assume that the Froude number F satisfies $0 < F < F_c$. There exists $\rho > 0$ so that for any $(r_0, s_0) \in (R, S) + H^{m+\frac{1}{2}}(\Omega)$, with $m > 2$ and $(X_{0,i} - iL)_{i \in \mathbb{Z}} \in l^2(\mathbb{Z})$, compatible up to order $m - 1$ and such that*

$$\|(r_0, s_0) - (R, S)\|_{L^\infty(\Omega)} \leq \rho, \quad \|X_{0,i+1} - X_{0,i} - L\|_{l^\infty(\mathbb{Z})} \leq \rho,$$

there is $T > 0$ and a solution (r, s) and X of the shallow water system (8) coupled with the Rankine-Hugoniot jump conditions (9) so that

$$(r, s)|_{t=0} = (r_0, s_0), \quad X_i(t = 0) = X_{0,i},$$

and (r, s) belongs to $(R, S) + H^m(\Omega \times [0, T])$, whereas $(X_i - (ct + iL))_{i \in \mathbb{Z}}$ lies in $l^2(\mathbb{Z}, H^{m+1}(0, T))$.

Proof. The method of proof is based on an iterative scheme: We search the solution in the form $(r, s) = (r_a, s_a) + (\bar{r}, \bar{s})$, $X = X_a + \varepsilon$, where (r_a, s_a, X_a) is the approximate solution associated to the initial data (r_0, s_0, X_0) and defined in Proposition (7). For that purpose, we introduce the iterative scheme $r^{(0)} = s^{(0)} = 0$, $\varepsilon^{(0)} = 0$, and $(r^{(k)}, s^{(k)})$, $\varepsilon^{(k)}$ is defined inductively as the unique solution, with zero initial data, of the system

$$(50) \quad \partial_t r_i^{(k)} + \Lambda_i^{1,(k-1)} \partial_x r_i^{(k)} = f_i^{1,(k-1)}, \quad \partial_t s_i^{(k)} + \Lambda_i^{2,(k-1)} \partial_x s_i^{(k)} = f_i^{2,(k-1)},$$

where the functions $f_i^{n,(k-1)}$, $n = 1, 2$ are defined by

$$\begin{aligned} f_i^{1,(k-1)} &= Q_i^{(k-1)} - \partial_t r_i^{(k-1)} - \Lambda_i^{1,(k-1)} \partial_x r_i^{(k-1)}, \\ f_i^{2,(k-1)} &= Q_i^{(k-1)} - \partial_t s_i^{(k-1)} - \Lambda_i^{2,(k-1)} \partial_x s_i^{(k-1)}, \end{aligned}$$

with $Q_i^{(k)} = Q(r_{a,i} + r^{(k)}, s_{a,i} + s^{(k)})$ and

$$\Lambda_i^{n,(k)} = \frac{L}{X_{a,i+1} - X_{a,i} + \varepsilon_{i+1}^{(k)} - \varepsilon_i^{(k)}} \left(\lambda^n \left(r_{a,i} + r^{(k)} s_{a,i} + s^{(k)} \right) - \left(\dot{X}_{a,i} + \dot{\varepsilon}_i^{(k)} + \frac{x}{L} \left(\dot{X}_{a,i+1} - \dot{X}_{a,i} + \dot{\varepsilon}_{i+1}^{(k)} - \dot{\varepsilon}_i^{(k)} \right) \right) \right).$$

Equations (50) are supplemented by boundary conditions that read

$$(51) \quad \begin{pmatrix} \dot{\varepsilon}_i^{(k)} \\ 0 \end{pmatrix} = \underline{B}_i^{(k-1)} \left(r_i^{(k),\pm}, s_i^{(k),\pm} \right) + \begin{pmatrix} g_i^{1,(k-1)} \\ g_i^{2,(k-1)} \end{pmatrix},$$

where the functions $g_i^{n,(k)}$ are defined by

$$\begin{pmatrix} g_i^{1,(k)} \\ g_i^{2,(k)} \end{pmatrix} = \begin{pmatrix} -\dot{X}_{a,i} + F \left(r_{a,i}^\pm + r_i^{(k),\pm}, s_{a,i}^\pm + s_i^{(k),\pm} \right) \\ G \left(r_{a,i}^\pm + r_i^{(k),\pm}, s_{a,i}^\pm + s_i^{(k),\pm} \right) \end{pmatrix} - \underline{B}_i^{(k)} \left(r_i^{(k),\pm}, s_i^{(k),\pm} \right).$$

We estimate the nonlinear terms $f_i^{n,(k-1)}$ and $g_i^{n,(k-1)}$. One can prove that there exists $M > 0$ so that for any $T \in (0, T_0)$, and if $r^{(k-1)}, s^{(k-1)}$ belongs to $H^m(\Omega) \times I_T$, $r^{(k-1),\pm}, s^{(k-1),\pm}$ lies in $l^2(\mathbb{Z}, H^m(I_T))$ and $\varepsilon^{(k-1)}$ to $l^2(\mathbb{Z}, H^{m+1}(I_T))$, vanish for $t < 0$ and are such that

$$(52) \quad \begin{aligned} \|r^{(k-1)}, s^{(k-1)}\|_{H^m(\Omega \times I_T)} &\leq M, \\ \|r^{(k-1),\pm}, s^{(k-1),\pm}\|_{l^2(\mathbb{Z}, H^m(I_T))} &\leq M, \\ \|\varepsilon^{(k-1)}\|_{l^2(\mathbb{Z}, H^{m+1}(I_T))} &\leq M, \end{aligned}$$

then $f^{n,(k-1)}$ lies in $H^m(\Omega \times I_T)$, $g^{n,(k-1)}$ in $l^2(\mathbb{Z}, H^m(I_T))$ and satisfy the estimates

$$\|f^{n,(k-1)}\|_{H^m(\Omega \times I_T)} \leq C(M), \quad \|g^{n,(k-1)}\|_{l^2(\mathbb{Z}, H^m(I_T))} \leq C(M)T + \eta(T),$$

with $\eta(T) \rightarrow 0$ as $T \rightarrow 0$ and $M \mapsto C(M)$ is a continuous function. The function $(r^{(0)}, s^{(0)}, \varepsilon^{(0)})$ satisfies all the assumptions. We can apply the Proposition 6: There exists a unique solution $(r^{(1)}, s^{(1)}) \in H^m(\Omega \times I_T)$ and $\varepsilon^{(1)} \in l^2(\mathbb{Z}, H^{m+1}(I_T))$ solution of (50), (51) with $k = 1$. Moreover, diminishing T if necessary, one can prove, using the energy estimate in Proposition 6, that

$$\begin{aligned} \|r^{(1)}, s^{(1)}\|_{H^m(\Omega \times I_T)} &\leq M, \\ \|r^{(1),\pm}, s^{(1),\pm}\|_{l^2(\mathbb{Z}, H^m(I_T))} &\leq M, \\ \|\varepsilon^{(1)}\|_{l^2(\mathbb{Z}, H^{m+1}(I_T))} &\leq M. \end{aligned}$$

By induction, one can construct $(r^{(k)}, s^{(k)}) \in H^m(\Omega \times I_T)$ and $\varepsilon^{(k)} \in l^2(\mathbb{Z}, H^{m+1}(I_T))$ that satisfy (52) for $k \geq 1$. Now using the fact that $W^{1,\infty}(\mathbb{R}) \hookrightarrow H^m(\mathbb{R})$, we obtain uniform estimates on $\|(r^{(k)}, s^{(k)})\|_{W^{1,\infty}(\mathbb{R})}$. We write the equations satisfied by $(r^{(k)} - r^{(k-1)}, s^{(k)} - s^{(k-1)})$ that is denoted $u^{(k)}$ in what follows, and $\varepsilon^{(k)} - \varepsilon^{(k-1)}$. Using the L^2 energy estimate of Proposition 4, we prove that, for T sufficiently small,

$$\begin{aligned} \|u^{(k+1)}\|_{L^2(\Omega \times I_T)} &\leq \frac{1}{2} \|u^{(k)}\|_{L^2(\Omega \times I_T)}, \\ \|u^{(k+1),\pm}\|_{l^2(\mathbb{Z}, L^2(I_T))} &\leq \frac{1}{2} \|u^{(k),\pm}\|_{l^2(\mathbb{Z}, L^2(I_T))}, \\ \|\varepsilon^{(k+1)} - \varepsilon^{(k)}\|_{l^2(\mathbb{Z}, H^1(I_T))} &\leq \frac{1}{2} \|\varepsilon^{(k)} - \varepsilon^{(k-1)}\|_{l^2(\mathbb{Z}, H^1(I_T))}. \end{aligned}$$

Then, $(r^{(k)}, s^{(k)})$, $(r^{(k),\pm}, s^{(k),\pm})$, and $\varepsilon^{(k)}$ are, respectively, Cauchy sequences in $L^2(\Omega \times I_T)$, $l^2(\mathbb{Z}, L^2(I_T))$, and $l^2(\mathbb{Z}, H^1(I_T))$ and converge in that space. We denote \bar{r} , \bar{s} , and $\bar{\varepsilon}$ the limits of $(r^{(k)}, s^{(k)})$ and $\varepsilon^{(k)}$ as $k \rightarrow \infty$. Necessarily, (\bar{r}, \bar{s}) lies in $H^m(\Omega \times I_T)$, $(\bar{r}^\pm, \bar{s}^\pm)$ in $l^2(\mathbb{Z}, H^m(I_T))$, and $\bar{\varepsilon}$ belongs to $l^2(\mathbb{Z}, H^{m+1}(I_T))$. By standard interpolation arguments, we obtain strong convergence in appropriate Sobolev spaces so that we can make $k \rightarrow \infty$ in (50),(51). As a result, the function $(r_a + \bar{r}, s_a + \bar{s})$ and $X_a + \bar{\varepsilon}$ solves the full Cauchy problem (8), (9) on the interval $(0, T)$. This concludes the proof of the theorem and the persistence of roll-waves for suitable initial data. \square

5. Conclusion. In this paper, we proved the well posedness of the shallow water equations in the neighborhood of roll-wave solutions for initial conditions that satisfy suitable compatibility conditions. The main issue here is to formulate the stability problem in the presence of an infinite number of shocks. This is done through a Lipschitz change of variable that fixes *all* the discontinuities of functions that are close to roll-waves. Writing the shallow water equations with Riemann invariants, we diagonalize the shallow water system; in that case, it is easier to obtain estimates on the linearized problems. We prove the existence of weak solutions to the linearized problem through the formulation of an adjoint problem just as in the shock case. We then prove regularity properties on that solution. The main difference with the shock case is the presence of a singularity in the interior equations due to the presence of a sonic point: a regularization in the transverse direction is not sufficient here, and we regularize the solution also in the direction of propagation. The existence of a solution to the full nonlinear problem is then obtain through a fixed point argument. A key assumption is that the Froude number F (that measures whether a the flow of the channel is fluvial or torrential) has to be smaller than a critical value F_c : this suggests that in the super critical regime $F > F_c$, the roll-wave structure may be destroyed. This would be an interesting question to see what happens in that regime even from a numerical point of view.

Though written in a 1D setting, we claim that the results obtained in this paper hold true in the *multidimensional setting*. We briefly discuss the method of proof. Indeed, in the 2D case, one shall consider a transverse coordinate y and a transverse speed v ; the hyperbolic part of the equations is nothing but the isentropic Euler equations. Here we shall assume that the shocks are located at the positions $\Psi_i(t, y) \approx iT$, $i \in \mathbb{Z}$. Next we perform a change of variable similar to the one used in this paper, and one can easily formulate the linearized equations. The main difference with the 1D case is the obtention of energy estimates for the linearized (exact and approximated) equations. However, once the linearization procedure is performed, we are left with an infinite but discrete number of Lax shocks that can be treated separately. For each Lax shock, one can use Kreiss symmetrizers to obtain energy estimates and then sum all the contributions over $i \in \mathbb{Z}$. The end of the proof is then completely similar to the 1D case presented here. This method shall be developed in a forthcoming paper on the persistence of roll-waves in general hyperbolic systems with discontinuities satisfying Lax shock conditions obtained recently by the author[14]. We shall also consider the persistence of roll-waves with characteristic discontinuities. This is of interest from a physical point of view, since they may appear when the pressure term in shallow water equations is not convex like in stratified flows (see [2] for more details).

At this stage, several questions arise. First, we have restricted our analysis to *smooth* perturbations: in the 1D case, because of the progress made in the analysis of the 1D problem, it seems relevant to consider more general perturbations. In the case of shock waves, the persistence of shocks satisfying the Majda–Liu stability condition

under small perturbations with a bounded variation (BV) using the Glimm scheme or the front tracking method has been established (see [16], [8] for more details). It is thus a natural question whether roll-waves persist under BV perturbations, though it may be a very difficult problem in the case of roll-waves. Another interesting question is the stability of viscous roll-waves that are close to Dressler roll-waves in the vanishing viscosity limit [10, 11]; that question has been treated recently in the shock case for multidimensional perturbations [6]. Moreover, the linear stability of viscous roll-waves has been proved under strong spectral stability [12]. Energy estimates have been obtained through Green functions that blow up in the vanishing viscosity case, just as in the shock case. It would be interesting here to adapt the method introduced by Metivier and coworkers [6] to the roll-wave case: the main issue would be the presence of an infinite number of viscous shocks. This question is an open problem that is postponed to a forthcoming work.

Acknowledgment. The author wishes to thank Prof. Sylvie Benzoni for many useful discussions on the multidimensional stability of shock waves.

REFERENCES

- [1] S. BENZONI-GAVAGE AND D. SERRE, *Multidimensional Hyperbolic Partial Differential Equations: First-order Systems and Applications*, Oxford Math. Monogr., Clarendon Press, Oxford University Press, Oxford, 2007.
- [2] A. BOUDDLAL AND V. LIAPIDEVSKII, *Multi-shock structure of roll-waves*, C.R.A.S. Mécanique, 332 (2004), pp. 659–664.
- [3] R. DRESSLER, *Mathematical solution of the problem of roll waves in inclined open channels*, Comm. Pure Appl. Math., 2 (1949), pp. 149–190.
- [4] A. MAJDA, *The stability of multidimensional shock fronts*, Mem. Amer. Math. Soc., 41 (1983), No. 275.
- [5] G. MÉTIVIER, *Stability of multidimensional shocks*, in Advances in the Theory of Shock Waves, pp. 25–103, Progr. Nonlinear Differential Equations Appl. 47, Birkhäuser Boston, Boston 2001, pp. 25–103.
- [6] O. GUÈS, G. MÉTIVIER, M. WILLIAMS, AND K. ZUMBRUN, *Existence and stability of multidimensional shock fronts in the vanishing viscosity limit*, Arch. Ration. Mech. Anal., 175 (2005), pp. 151–244.
- [7] J.-F. GERBEAU AND B. PERTHAME, *Derivation of viscous Saint-Venant system for Lamina shallow-water: Numerical validation*, Discrete Contin. Dyn. Syst. Ser. B, 1 (2001), pp. 89–102.
- [8] M. LEWICKA AND K. TRIVISA, *On the L^1 well posedness of systems of conservation laws near solutions containing two large shocks*, J. Differential Equations, 179 (2002), pp. 133–177.
- [9] D.J. NEEDHAM AND J.H. MERKIN, *On roll waves down an open inclined channel*, Proc. Roy. Soc. Lond. Ser. A Math. Phys. Eng. Sci., 394 (1984), pp. 259–278.
- [10] J. HÄRTERICH, *Existence of rollwaves in a viscous shallow water equation*, in Proceedings of EQUADIFF 2003, World Scientific, Hackensack, NJ, 2005, pp. 511–516.
- [11] P. NOBLE, *Méthodes de Variétés Invariantes Pour les Équations de Saint Venant et les Systèmes Hamiltoniens Discrets*, Thèse Université Toulouse, Toulouse, France, 2003.
- [12] P. NOBLE, *Linear stability of viscous roll-waves*, Comm. Partial Differential Equations, 32 (2007), pp. 1681–1713.
- [13] P. NOBLE, *Linear stability of roll waves* Indiana Univ. Math. J., 55 (2006), pp. 795–848.
- [14] P. NOBLE, *Roll-waves in general hyperbolic systems with source terms*, SIAM J. Appl. Math., 67 (2007), pp. 1202–1212.
- [15] F. ROUSSET, *Viscous approximation of strong shocks of systems of conservation laws*, SIAM J. Math. Anal., 35 (2003), pp. 492–519.
- [16] S. SCHOCHET, *Sufficient conditions for local existence via Glimm’s scheme for large BV data*, J. Differential Equations, 89 (1991), pp. 317–354.
- [17] K. TAMADA AND H. TOUGOU, *Stability of roll-waves on thin Lamina flow down an inclined plane wall*, J. Phys. Soc. Japan, 47 (1979), pp. 1992–1998.
- [18] J.P. VILA, *Sur la Théorie et l’Approximation Numérique des Problèmes Hyperboliques Non Linéaires. Application aux Équations de Saint Venant et à la Modélisation des Avalanches de neige Dense*, Thèse Univ. Paris VI, Paris, France, (1986).

A NONLOCAL p -LAPLACIAN EVOLUTION EQUATION WITH NONHOMOGENEOUS DIRICHLET BOUNDARY CONDITIONS*

F. ANDREU[†], J. M. MAZÓN[‡], J. D. ROSSI[§], AND J. TOLEDO[‡]

Abstract. In this paper we study the nonlocal p -Laplacian-type diffusion equation $u_t(t, x) = \int_{\mathbb{R}^N} J(x-y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy$, $(t, x) \in]0, T[\times \Omega$, with $u(t, x) = \psi(x)$ for $(t, x) \in]0, T[\times (\mathbb{R}^N \setminus \Omega)$. If $p > 1$, this is the nonlocal analogous problem to the well-known local p -Laplacian evolution equation $u_t = \operatorname{div}(|\nabla u|^{p-2} \nabla u)$ with Dirichlet boundary condition $u(t, x) = \psi(x)$ on $(t, x) \in]0, T[\times \partial\Omega$. If $p = 1$, this is the nonlocal analogous to the total variation flow. When $p = +\infty$ (this has to be interpreted as the limit as $p \rightarrow +\infty$ in the previous model) we find an evolution problem that can be seen as a nonlocal model for the formation of sandpiles (here $u(t, x)$ stands for the height of the sandpile) with prescribed height of sand outside of Ω . We prove, as main results, existence, uniqueness, a contraction property that gives well posedness of the problem, and the convergence of the solutions to solutions of the local analogous problem when a rescaling parameter goes to zero.

Key words. nonlocal diffusion, p -Laplacian, nonhomogeneous Dirichlet boundary conditions, total variation flow, sandpiles

AMS subject classifications. 35B40, 45A07, 45G10

DOI. 10.1137/080720991

1. Introduction. In this paper we study the nonlocal diffusion equation

$$u_t(t, x) = \int_{\mathbb{R}^N} J(x-y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy \quad (t, x) \in]0, T[\times \Omega,$$

where Ω is a bounded domain and u is prescribed in $\mathbb{R}^N \setminus \Omega$ as $u(t, x) = \psi(x)$ for $(t, x) \in]0, T[\times (\mathbb{R}^N \setminus \Omega)$. We consider $1 < p < +\infty$ as well as the extreme cases $p = 1$ and the limit $p \nearrow +\infty$. Throughout the paper, we assume that $J : \mathbb{R}^N \rightarrow \mathbb{R}$ is a nonnegative, radial, continuous function, strictly positive in $B(0, 1)$, vanishing in $\mathbb{R}^N \setminus B(0, 1)$ and such that $\int_{\mathbb{R}^N} J(z) dz = 1$.

First, let us briefly introduce the prototype of nonlocal problem that will be considered along this work. Nonlocal evolution equations of the form

$$(1.1) \quad u_t(t, x) = (J * u - u)(t, x) = \int_{\mathbb{R}^N} J(x-y)u(t, y) dy - u(t, x),$$

and variations of it, have been recently widely used to model diffusion processes. More precisely, as stated in [31], if $u(t, x)$ is thought of as a density at the point x at time t and $J(x-y)$ is thought of as the probability distribution of jumping from location y to

*Received by the editors April 11, 2008; accepted for publication (in revised form) September 13, 2008; published electronically January 7, 2009.

<http://www.siam.org/journals/sima/40-5/72099.html>

[†]Departament de Matemàtica Aplicada, Universitat de València, Valencia, Spain (fuensanta.andreu@uv.es). This author was partially supported by the PNPGC project, reference MTM2008-03176.

[‡]Departament d'Anàlisi Matemàtica, Universitat de València, Valencia, Spain (mazon@uv.es, toledojj@uv.es). These authors were partially supported by the PNPGC project, reference MTM2008-03176.

[§]IMDEA Matemáticas, C-IX, Campus Cantoblanco UAM, Madrid, Spain. On leave from Dpto. de Matemáticas, FCEyN Universidad de Buenos Aires, 1428 Buenos Aires, Argentina (jrossi@dm.uba.ar). This author was partially supported by ANPCyT PICT 5009, UBA X066, CONICET (Argentina).

location x , then $\int_{\mathbb{R}^N} J(y-x)u(t,y) dy = (J*u)(t,x)$ is the rate at which individuals are arriving at position x from all other places and $-u(t,x) = -\int_{\mathbb{R}^N} J(y-x)u(t,x) dy$ is the rate at which they are leaving location x to travel to all other sites. This consideration, in the absence of external or internal sources, leads immediately to the fact that the density u satisfies (1.1). For recent references on nonlocal diffusion, see [4], [5], [6], [9], [11], [12], [20], [21], [22], [23], [24], [25], [26], [27], [31], [33], [36] and references therein.

The first goal of this paper is to study the following nonlocal nonlinear diffusion problem:

$$P_p^J(u_0, \psi) \begin{cases} u_t(t, x) = \int_{\Omega} J(x-y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy \\ \quad + \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)|\psi(y) - u(t, x)|^{p-2}(\psi(y) - u(t, x)) dy, \\ \quad \quad \quad (t, x) \in]0, T[\times \Omega, \\ u(0, x) = u_0(x), \quad x \in \Omega. \end{cases}$$

Here $\Omega_J = \Omega + \text{supp}(J)$ and ψ is a given function $\psi : \Omega_J \setminus \bar{\Omega} \rightarrow \mathbb{R}$.

Observe that we can rewrite $P_p^J(u_0, \psi)$, setting $u(t, x) = \psi(x)$ in $\Omega_J \setminus \bar{\Omega}$, as

$$\begin{cases} u_t(t, x) = \int_{\Omega_J} J(x-y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy, \quad (t, x) \in]0, T[\times \Omega, \\ u(t, x) = \psi(x), \quad (t, x) \in]0, T[\times (\Omega_J \setminus \bar{\Omega}), \\ u(0, x) = u_0(x), \quad x \in \Omega, \end{cases}$$

and we call it the *nonlocal p -Laplacian problem with Dirichlet boundary condition*. Note that we are prescribing the values of u outside the domain Ω and not only on its boundary. This is due to the nonlocal character of the problem.

Let us state the precise definition of solution. Solutions to $P_p^J(u_0, \psi)$ will be understood in the following sense.

DEFINITION 1.1. *Let $1 < p < +\infty$. A solution of $P_p^J(u_0, \psi)$ in $[0, T]$ is a function*

$$u \in C([0, T]; L^1(\Omega)) \cap W^{1,1}(]0, T[; L^1(\Omega)),$$

which satisfies $u(0, x) = u_0(x)$ a.e. $x \in \Omega$ and

$$u_t(t, x) = \int_{\Omega} J(x-y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy \\ + \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)|\psi(y) - u(t, x)|^{p-2}(\psi(y) - u(t, x)) dy,$$

for a.e. $t \in]0, T[$ and a.e. $x \in \Omega$.

Our first result shows existence and uniqueness of a global solution for this problem. Moreover, a contraction principle holds.

THEOREM 1.2. *Assume $p > 1$ and let $u_0 \in L^p(\Omega)$, $\psi \in L^p(\Omega_J \setminus \bar{\Omega})$. Then, there exists a unique solution to $P_p^J(u_0, \psi)$ in the sense of Definition 1.1. Moreover, if $u_{i0} \in L^1(\Omega)$ and u_i is a solution in $[0, T]$ of $P_p^J(u_{i0}, \psi)$, $i = 1, 2$, respectively. Then*

$$\int_{\Omega} (u_1(t) - u_2(t))^+ \leq \int_{\Omega} (u_{10} - u_{20})^+ \quad \text{for every } t \in [0, T].$$

If $u_{i0} \in L^p(\Omega)$, $i = 1, 2$, then

$$\|u_1(t) - u_2(t)\|_{L^p(\Omega)} \leq \|u_{10} - u_{20}\|_{L^p(\Omega)} \quad \text{for every } t \in [0, T].$$

Our next step is to rescale the kernel J appropriately and take the limit as the scaling parameter goes to zero. To be more precise, for $p > 1$, we consider the local p -Laplace evolution equation with Dirichlet boundary condition

$$D_p(u_0, \tilde{\psi}) \quad \begin{cases} u_t = \Delta_p u & \text{in }]0, T[\times \Omega, \\ u = \tilde{\psi} & \text{on }]0, T[\times \partial\Omega, \\ u(0, x) = u_0(x) & \text{in } \Omega, \end{cases}$$

where the boundary datum $\tilde{\psi}$ is assumed to be the trace of a function defined in a larger domain and the operator in the equation, $\Delta_p u = \operatorname{div}(|\nabla u|^{p-2} \nabla u)$, is the usual local p -Laplacian.

We prove that the solutions of this local problem can be approximated by solutions of a sequence of nonlocal p -Laplacian problems of the form P_p^J . Indeed, for given $p \geq 1$ and J we consider the rescaled kernels

$$(1.2) \quad J_{p,\varepsilon}(x) := \frac{C_{J,p}}{\varepsilon^{p+N}} J\left(\frac{x}{\varepsilon}\right), \quad \text{where} \quad C_{J,p}^{-1} := \frac{1}{2} \int_{\mathbb{R}^N} J(z) |z_N|^p dz$$

is a normalizing constant in order to obtain the p -Laplacian in the limit instead of a multiple of it, and we obtain the following result.

THEOREM 1.3. *Let Ω be a smooth bounded domain in \mathbb{R}^N and $\tilde{\psi} \in W^{1/p', p}(\partial\Omega) \cap L^\infty(\partial\Omega)$. Let $\psi \in W^{1,p}(\Omega_J) \cap L^\infty(\Omega_J)$ such that the trace $\psi|_{\partial\Omega} = \tilde{\psi}$. Assume $J(x) \geq J(y)$ if $|x| \leq |y|$. Let $T > 0$ and $u_0 \in L^p(\Omega)$. Let u_ε be the unique solution of $P_p^{J_{p,\varepsilon}}(u_0, \psi)$ and u the unique solution of $D_p(u_0, \tilde{\psi})$ (see section 2.2). Then*

$$(1.3) \quad \lim_{\varepsilon \rightarrow 0} \sup_{t \in [0, T]} \|u_\varepsilon(t, \cdot) - u(t, \cdot)\|_{L^p(\Omega)} = 0.$$

Note that the above result says that P_p^J is a nonlocal problem analogous to the p -Laplacian with Dirichlet boundary condition.

The second goal of this paper is to study the Dirichlet problem for $p = 1$, called the nonlocal total variation flow, which can be written formally as

$$P_1^J(u_0, \psi) \quad \begin{cases} u_t(t, x) = \int_{\Omega} J(x-y) \frac{u(t, y) - u(t, x)}{|u(t, y) - u(t, x)|} dy \\ \quad + \int_{\Omega_J \setminus \bar{\Omega}} J(x-y) \frac{\psi(t, y) - u(t, x)}{|\psi(t, y) - u(t, x)|} dy, & (t, x) \in]0, T[\times \Omega, \\ u(0, x) = u_0(x), & x \in \Omega. \end{cases}$$

We give the following definition of what we understand by a solution of $P_1^J(u_0, \psi)$.

DEFINITION 1.4. *A solution of $P_1^J(u_0, \psi)$ in $[0, T]$ is a function*

$$u \in C([0, T]; L^1(\Omega)) \cap W^{1,1}([0, T]; L^1(\Omega)),$$

which satisfies $u(0, x) = u_0(x)$ a.e. $x \in \Omega$ and

$$u_t(t, x) = \int_{\Omega_J} J(x-y) g(t, x, y) dy \quad \text{a.e. in }]0, T[\times \Omega,$$

for some $g \in L^\infty(0, T; L^\infty(\Omega_J \times \Omega))$ with $\|g\|_\infty \leq 1$ such that for almost every $t \in]0, T[$, $g(t, x, y) = -g(t, y, x)$ and

$$J(x - y)g(t, x, y) \in J(x - y)\text{sign}(u(t, y) - u(t, x)), \quad (x, y) \in \Omega \times \Omega,$$

$$J(x - y)g(t, x, y) \in J(x - y)\text{sign}(\psi(y) - u(t, x)), \quad (x, y) \in \Omega \times (\Omega_J \setminus \overline{\Omega}).$$

Here, sign is the multivalued function defined by

$$\text{sign}(r) := \begin{cases} 1 & \text{if } r > 0 \\ [-1, 1] & \text{if } r = 0 \\ -1 & \text{if } r < 0. \end{cases}$$

We use sign_0 to denote the univalued function

$$\text{sign}_0(r) := \begin{cases} 1 & \text{if } r > 0 \\ 0 & \text{if } r = 0 \\ -1 & \text{if } r < 0. \end{cases}$$

To get the existence and uniqueness of these kinds of solutions, the idea is to take the limit as $p \searrow 1$ of solutions to P_p^J with $p > 1$.

THEOREM 1.5. *Let $u_0 \in L^1(\Omega)$ and $\psi \in L^1(\Omega_J \setminus \overline{\Omega})$. Then, there exists a unique solution to $P_1^J(u_0)$ in the sense of Definition 1.4. Moreover, if $u_{i_0} \in L^1(\Omega)$ and u_i are solutions in $[0, T]$ of $P_1^J(u_{i_0})$, $i = 1, 2$. Then*

$$\int_{\Omega} (u_1(t) - u_2(t))^+ \leq \int_{\Omega} (u_{10} - u_{20})^+ \quad \text{for every } t \in [0, T].$$

In this case we can rescale the kernel as in (1.2) in order to obtain convergence of the solutions of the corresponding rescaled problem to the solution of the Dirichlet problem for the total variational flow, that is,

$$D_1(u_0, \tilde{\psi}) \quad \begin{cases} u_t = \text{div} \left(\frac{Du}{|Du|} \right) & \text{in }]0, T[\times \Omega, \\ u = \tilde{\psi} & \text{on }]0, T[\times \partial\Omega, \\ u(0, x) = u_0(x) & \text{in } \Omega. \end{cases}$$

THEOREM 1.6. *Let Ω be a smooth bounded domain in \mathbb{R}^N . Assume $J(x) \geq J(y)$ if $|x| \leq |y|$. Let $T > 0$, $u_0 \in L^1(\Omega)$, $\tilde{\psi} \in L^\infty(\partial\Omega)$, and $\psi \in W^{1,1}(\Omega_J \setminus \overline{\Omega}) \cap L^\infty(\Omega_J \setminus \overline{\Omega})$ such that the trace $\psi|_{\partial\Omega} = \tilde{\psi}$. Let u_ε be the unique solution of $P_1^{J, \varepsilon}(u_0, \psi)$. Then, if u is the unique solution of $D_1(u_0, \tilde{\psi})$ (see section 3.2),*

$$\lim_{\varepsilon \rightarrow 0} \sup_{t \in [0, T]} \|u_\varepsilon(t, \cdot) - u(t, \cdot)\|_{L^1(\Omega)} = 0.$$

Finally, the third goal of this paper is to study the limit case $p = +\infty$, which has to be understood as the limit of our nonlocal evolution problems as $p \rightarrow +\infty$ (see section 4). In this case we recover a nonlocal model for the evolution of sandpiles which is the nonlocal version of the Prigozhin model [35]. Then, the nonlocal limit problem with source for $p = +\infty$ can be written as

$$P_\infty^J(u_0, \psi, f) \quad \begin{cases} f(t, \cdot) - u_t(t, \cdot) \in \partial G_{\infty, \psi}^J(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x), \end{cases}$$

where $G_{\infty,\psi}^J$ is the functional

$$G_{\infty,\psi}^J(u) = \begin{cases} 0 & \text{if } |u(x) - u(y)| \leq 1, \text{ for } x, y \in \Omega \\ & \text{and } |\psi(y) - u(x)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ & \text{with } x - y \in \text{supp}(J) \\ +\infty & \text{in the other case,} \end{cases}$$

that is, $G_{\infty,\psi}^J = I_{K_{\infty,\psi}^J}$, the indicator function of the set

$$K_{\infty,\psi}^J := \left\{ u \in L^2(\Omega) : \begin{array}{l} |u(x) - u(y)| \leq 1, x, y \in \Omega \\ \text{and } |\psi(y) - u(x)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ \text{with } x - y \in \text{supp}(J) \end{array} \right\}.$$

More precisely, we obtain the following result.

THEOREM 1.7. *Let $\psi \in L^\infty(\Omega_J \setminus \bar{\Omega})$ such that $K_{\infty,\psi}^J \neq \emptyset$. Let $T > 0$, $f \in L^2(0, T; \cap_{q \geq 2} L^q(\Omega))$, $u_0 \in \cap_{q \geq 2} L^q(\Omega)$ such that $u_0 \in K_{\infty,\psi}^J$, and u_p , $p \geq 2$, the unique solution of the nonlocal p -Laplacian with a source term f , $P_p^J(u_0, \psi, f)$ (see section 4). Then, if u_∞ is the unique solution to $P_\infty^J(u_0, \psi, f)$,*

$$\lim_{p \rightarrow \infty} \sup_{t \in [0, T]} \|u_p(t, \cdot) - u_\infty(t, \cdot)\|_{L^2(\Omega)} = 0.$$

Our next step is to rescale the kernel J appropriately and take the limit as the scaling parameter goes to zero. We will suppose that Ω is convex and ψ verifies $\|\nabla\psi\|_\infty \leq 1$. For $\varepsilon > 0$, we rescale the functional $G_{\infty,\psi}^J$ as follows:

$$G_{\infty,\psi}^\varepsilon(u) = \begin{cases} 0 & \text{if } |u(x) - u(y)| \leq \varepsilon, \text{ for } x, y \in \Omega \\ & \text{and } |\psi(y) - u(x)| \leq \varepsilon, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ & \text{with } |x - y| \leq \varepsilon \\ +\infty & \text{in the other case,} \end{cases}$$

that is, $G_{\infty,\psi}^\varepsilon = I_{K_{\infty,\psi}^\varepsilon}$, where

$$K_{\infty,\psi}^\varepsilon := \left\{ u \in L^2(\Omega) : \begin{array}{l} |u(x) - u(y)| \leq \varepsilon, x, y \in \Omega \\ \text{and } |\psi(y) - u(x)| \leq \varepsilon, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ \text{with } |x - y| \leq \varepsilon \end{array} \right\}.$$

Consider the gradient flow associated to the functional $G_{\infty,\psi}^\varepsilon$

$$P_\infty^\varepsilon(u_0, \psi, f) \begin{cases} f(t, \cdot) - u_t(t, \cdot) \in \partial I_{K_{\infty,\psi}^\varepsilon}(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x), & \text{in } \Omega, \end{cases}$$

and the limit problem

$$P_\infty(u_0, \psi, f) \begin{cases} f(t, \cdot) - u_{\infty,t} \in \partial I_{K_\psi}(u_\infty), & \text{a.e. } t \in]0, T[, \\ u_\infty(0, x) = u_0(x), & \text{in } \Omega, \end{cases}$$

where

$$K_\psi := \{u \in W^{1,\infty}(\Omega) : \|\nabla u\|_\infty \leq 1, u|_{\partial\Omega} = \psi|_{\partial\Omega}\}.$$

Now we state our result concerning the limit as $\varepsilon \rightarrow 0$ for the sandpile model ($p = +\infty$).

THEOREM 1.8. *Assume Ω is a convex bounded domain in \mathbb{R}^N . Let $T > 0$, $f \in L^2(0, T; L^2(\Omega))$, $\psi \in W^{1,\infty}(\Omega_J \setminus \bar{\Omega})$ such that $\|\nabla\psi\|_\infty \leq 1$, $u_0 \in W^{1,\infty}(\Omega)$ such that $\|\nabla u_0\|_\infty \leq 1$ and $u_0|_{\partial\Omega} = \psi|_{\partial\Omega}$ (this means $u_0 \in K_\psi$), and consider $u_{\infty,\varepsilon}$ the unique solution of $P_\infty^\varepsilon(u_0, \psi, f)$. Then, if v_∞ is the unique solution of $P_\infty(u_0, \psi, f)$, we have*

$$\lim_{\varepsilon \rightarrow 0} \sup_{t \in [0, T]} \|u_{\infty,\varepsilon}(t, \cdot) - v_\infty(t, \cdot)\|_{L^2(\Omega)} = 0.$$

Closely related to the present work are [5] and [6] where the homogeneous Neumann problem and its limit as p goes to infinity or to one are considered. The difference here is that we are now considering Dirichlet boundary conditions, not only the homogeneous case, but also the nonhomogeneous case, and this introduces new difficulties specially when one tries to recover the local models when $\varepsilon \rightarrow 0$. Remark that in our nonlocal formulation we are not imposing any continuity between the values of u inside Ω and outside it, ψ . However, when dealing with local problems usually the boundary datum is taken in the sense of traces, that is, $u|_{\partial\Omega} = \psi$. Recovering this condition as $\varepsilon \rightarrow 0$ is one of the main contributions of the present work.

Note that, as it happens for the local p -Laplacian, the Dirichlet problem can be written as a Neumann problem with a particular flux that depends on the solution itself. Indeed, the problem $P_p^J(u_0, \psi)$ can be written as

$$\begin{cases} u_t(t, x) = \int_\Omega J(x - y)|u(t, y) - u(t, x)|^{p-2}(u(t, y) - u(t, x)) dy + \varphi(x, u(x)) \\ u(0, x) = u_0(x), \quad x \in \Omega, \end{cases} \quad (t, x) \in]0, T[\times \Omega,$$

where

$$\varphi(x, u(x)) = \int_{\Omega_J \setminus \bar{\Omega}} J(x - y)|\psi(t, y) - u(t, x)|^{p-2}(\psi(t, y) - u(t, x)) dy.$$

In the homogeneous case, $\psi \equiv 0$,

$$\varphi(x, u(x)) = - \left(\int_{\Omega_J \setminus \bar{\Omega}} J(x - y) dy \right) |u(t, x)|^{p-2}u(t, x).$$

This problem is a nonhomogeneous Neumann problem (see [5]) with a prescribed flux given by φ .

Let us finish the introduction by collecting some notations and results that will be used in the sequel. Following [7] (see also [2]), let

$$(1.4) \quad X(\Omega) = \{z \in L^\infty(\Omega, \mathbb{R}^n) : \operatorname{div}(z) \in L^1(\Omega)\}.$$

If $z \in X(\Omega)$ and $w \in BV(\Omega) \cap L^\infty(\Omega)$, we define the functional $(z, Dw) : C_0^\infty(\Omega) \rightarrow \mathbb{R}$ by the formula

$$(1.5) \quad \langle (z, Dw), \varphi \rangle = - \int_\Omega w \varphi \operatorname{div}(z) dx - \int_\Omega w z \cdot \nabla \varphi dx.$$

Then (z, Dw) is a Radon measure in Ω ,

$$(1.6) \quad \int_{\Omega} (z, Dw) = \int_{\Omega} z \cdot \nabla w \, dx$$

for all $w \in W^{1,1}(\Omega) \cap L^{\infty}(\Omega)$ and

$$(1.7) \quad \left| \int_B (z, Dw) \right| \leq \int_B |(z, Dw)| \leq \|z\|_{\infty} \int_B \|Dw\|$$

for any Borel set $B \subseteq \Omega$.

In [7], a weak trace on $\partial\Omega$ of the normal component of $z \in X(\Omega)$ is defined. Concretely, it is proved that there exists a linear operator $\gamma : X(\Omega) \rightarrow L^{\infty}(\partial\Omega)$ such that

$$\|\gamma(z)\|_{\infty} \leq \|z\|_{\infty},$$

and

$$\gamma(z)(x) = z(x) \cdot \nu(x) \quad \text{for all } x \in \partial\Omega \text{ if } z \in C^1(\overline{\Omega}, \mathbb{R}^N).$$

We shall denote $\gamma(z)(x)$ by $[z, \nu](x)$. Moreover, the following *Green's formula*, relating the function $[z, \nu]$ and the measure (z, Dw) , for $z \in X(\Omega)$ and $w \in BV(\Omega) \cap L^{\infty}(\Omega)$, is established:

$$(1.8) \quad \int_{\Omega} w \operatorname{div}(z) \, dx + \int_{\Omega} (z, Dw) = \int_{\partial\Omega} [z, \nu] w \, d\mathcal{H}^{N-1}.$$

Organization of the paper. The rest of the paper is organized as follows. In the second section we prove the existence and uniqueness of strong solutions for the nonlocal p -Laplacian problem with Dirichlet boundary conditions for $p > 1$ and we show that our model approaches local p -Laplacian evolution equation with Dirichlet boundary condition. In section 3 we study the Dirichlet problem for the nonlocal total variation flow, proving convergence to the local model when the problem is rescaled appropriately as well. Finally, in section 4 we study the case $p = \infty$, obtaining a model for sandpiles with Dirichlet boundary conditions.

2. The case $p > 1$.

2.1. Existence of solutions for the nonlocal problems. We first study $P_p^J(u_0, \psi)$ from the point of view of nonlinear semigroup theory ([15], [28]). For that we introduce in $L^1(\Omega)$ the following operator associated with our problem.

DEFINITION 2.1. For $1 < p < +\infty$ and $\psi : \Omega_J \setminus \overline{\Omega} \rightarrow \mathbb{R}$, such that $|\psi|^{p-1} \in L^1(\Omega_J \setminus \overline{\Omega})$, we define in $L^1(\Omega)$ the operator $B_{p,\psi}^J$ by

$$\begin{aligned} B_{p,\psi}^J(u)(x) &= - \int_{\Omega} J(x-y) |u(y) - u(x)|^{p-2} (u(y) - u(x)) \, dy \\ &\quad - \int_{\Omega_J \setminus \overline{\Omega}} J(x-y) |\psi(y) - u(x)|^{p-2} (\psi(y) - u(x)) \, dy, \quad x \in \Omega. \end{aligned}$$

Remark 2.2. (i). We will set overall the section,

$$u_{\psi}(x) := \begin{cases} u(x) & \text{if } x \in \Omega, \\ \psi(x) & \text{if } x \in \Omega_J \setminus \overline{\Omega}, \\ 0 & \text{if } x \notin \Omega_J. \end{cases}$$

Therefore, we can rewrite

$$B_{p,\psi}^J(u)(x) = - \int_{\Omega_J} J(x-y)|u_\psi(y) - u(x)|^{p-2}(u_\psi(y) - u(x)) dy, \quad x \in \Omega.$$

(ii) If $\psi = 0$, then

$$B_{p,0}^J(u)(x) = - \int_{\Omega} J(x-y)|u(y) - u(x)|^{p-2}(u(y) - u(x)) dy + \left(\int_{\Omega_J \setminus \bar{\Omega}} J(x-y)dy \right) |u(x)|^{p-2}u(x), \quad x \in \Omega.$$

Remark 2.3. It is easy to see that

(i) If $\psi = 0$, $B_{p,0}^J$ is positively homogeneous of degree $p - 1$,

(ii) $L^{p-1}(\Omega) \subset \text{Dom}(B_{p,\psi}^J)$, if $p > 2$.

(iii) For $1 < p \leq 2$, $\text{Dom}(B_{p,\psi}^J) = L^1(\Omega)$ and $B_{p,\psi}^J$ is closed in $L^1(\Omega) \times L^1(\Omega)$.

We have the following monotonicity lemma, whose proof is straightforward.

LEMMA 2.4. *Let $1 < p < +\infty$, $\psi : \Omega_J \setminus \bar{\Omega} \rightarrow \mathbb{R}$, $|\psi|^{p-1} \in L^1(\Omega_J \setminus \bar{\Omega})$, and $T : \mathbb{R} \rightarrow \mathbb{R}$ a nondecreasing function. Then,*

(i) *for every $u, v \in L^p(\Omega)$ such that $T(u - v) \in L^p(\Omega)$, it holds*

$$(2.1) \quad \begin{aligned} & \int_{\Omega} (B_{p,\psi}^J u(x) - B_{p,\psi}^J v(x)) T(u(x) - v(x)) dx \\ &= \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) (T(u_\psi(y) - v_\psi(y)) - T(u_\psi(x) - v_\psi(x))) \\ & \quad \times (|u_\psi(y) - u_\psi(x)|^{p-2}(u_\psi(y) - u_\psi(x)) \\ & \quad - |v_\psi(y) - v_\psi(x)|^{p-2}(v_\psi(y) - v_\psi(x))) dy dx. \end{aligned}$$

(ii) *Moreover, if T is bounded, (2.1) holds for $u, v \in \text{Dom}(B_{p,\psi}^J)$.*

We have the following Poincaré’s type inequality.

PROPOSITION 2.5. *Given Ω a bounded domain in \mathbb{R}^N , $J : \mathbb{R}^N \rightarrow \mathbb{R}$ a nonnegative, radial, continuous function, such that $\int_{\mathbb{R}^N} J(z) dz > 0$, $p \geq 1$ and $\psi \in L^p(\Omega_J \setminus \bar{\Omega})$, there exists $\lambda = \lambda(J, \Omega, p) > 0$ such that*

$$(2.2) \quad \lambda \int_{\Omega} |u(x)|^p dx \leq \int_{\Omega} \int_{\Omega_J} J(x-y)|u_\psi(y) - u(x)|^p dy dx + \int_{\Omega_J \setminus \bar{\Omega}} |\psi(y)|^p dy$$

for all $u \in L^p(\Omega)$.

Proof. First, let us assume that there exist $r, \alpha > 0$ such that $J(x) \geq \alpha$ in $B(0, r)$.

Let

$$B_0 = \{x \in \Omega_J \setminus \bar{\Omega} : d(x, \Omega) \leq r/2\},$$

$$B_1 = \{x \in \Omega : d(x, B_0) \leq r/2\},$$

$$B_j = \left\{ x \in \Omega \setminus \cup_{k=1}^{j-1} B_k : d(x, B_{j-1}) \leq r/2 \right\}, \quad j = 2, 3, \dots$$

Observe that we can cover Ω by a finite number of nonnull sets $\{B_j\}_{j=1}^{l_r}$. Now

$$\int_{\Omega} \int_{\Omega_J} J(x-y)|u_\psi(y) - u(x)|^p dy dx \geq \int_{B_j} \int_{B_{j-1}} J(x-y)|u_\psi(y) - u(x)|^p dy dx,$$

$j = 1, \dots, l_r$, and

$$\begin{aligned} & \int_{B_j} \int_{B_{j-1}} J(x-y)|u_\psi(y) - u(x)|^p dy dx \\ & \geq \frac{1}{2^p} \int_{B_j} \int_{B_{j-1}} J(x-y)|u(x)|^p dy dx - \int_{B_j} \int_{B_{j-1}} J(x-y)|u_\psi(y)|^p dy dx \\ & = \frac{1}{2^p} \int_{B_j} \left(\int_{B_{j-1}} J(x-y) dy \right) |u(x)|^p dx - \int_{B_{j-1}} \left(\int_{B_j} J(x-y) dx \right) |u_\psi(y)|^p dy \\ & \geq \frac{1}{2^p} \min_{x \in \overline{B_j}} \int_{B_{j-1}} J(x-y) dy \int_{B_j} |u(x)|^p dx - \beta \int_{B_{j-1}} |u_\psi(y)|^p dy, \end{aligned}$$

where $\beta = \int_{\mathbb{R}^N} J(x) dx$. Hence

$$\int_{\Omega} \int_{\Omega_j} J(x-y)|u_\psi(y) - u(x)|^p dy dx \geq \alpha_j \int_{B_j} |u(x)|^p dx - \beta \int_{B_{j-1}} |u_\psi(y)|^p dy,$$

where

$$\alpha_j = \frac{1}{2^p} \min_{x \in \overline{B_j}} \int_{B_{j-1}} J(x-y) dy > 0.$$

Therefore, since $u_\psi(y) = \psi(y)$ if $y \in B_0$, $u_\psi(y) = u(y)$ if $y \in B_j$, $j = 1, \dots, l_r$, $B_j \cap B_i = \emptyset$, for all $i \neq j$ and $|\Omega \setminus \cup_{j=1}^{l_r} B_j| = 0$, it is easy to see that there exists $\hat{\lambda} = \hat{\lambda}(J, \Omega, p) > 0$ such that

$$\int_{\Omega} |u|^p \leq \hat{\lambda} \int_{\Omega} \int_{\Omega_j} J(x-y)|u_\psi(y) - u(x)|^p dy dx + \hat{\lambda} \int_{B_0} |\psi|^p.$$

The proof is finished by taking $\lambda = \hat{\lambda}^{-1}$.

In the general case we have that there exist $\mathbf{a} \geq 0$ and $r, \alpha > 0$ such that

$$(2.3) \quad J(x) \geq \alpha \text{ in the annulus } A(0, \mathbf{a}, r).$$

In this case we proceed as before with the same choice of the sets B_j for $j \geq 0$ and

$$B_{-j} = \left\{ x \in \Omega_j \setminus \left(\Omega \cup \cup_{k=0}^{j-1} B_{-k} \right) : d(x, B_{-j+1}) \leq r/2 \right\}, \quad j = 1, 2, 3, \dots$$

Observe that for each B_j , $j \geq 1$, there exists B_{j^e} with $j^e < j$ and such that

$$(2.4) \quad |(x + A(0, \mathbf{a}, r)) \cap B_{j^e}| > 0 \quad \forall x \in \overline{B_j}.$$

With this choice of B_j and taking into account (2.3) and (2.4), as before, we obtain

$$\begin{aligned} \int_{\Omega} \int_{\Omega_j} J(x-y)|u_\psi(y) - u(x)|^p dy dx & \geq \int_{B_j} \int_{B_{j^e}} J(x-y)|u_\psi(y) - u(x)|^p dy dx \\ & \geq \alpha_j \int_{B_j} |u(x)|^p dx - \beta \int_{B_{j^e}} |u_\psi(y)|^p dy, \end{aligned}$$

$j = 1, \dots, l_r$, where

$$\alpha_j = \frac{1}{2^p} \min_{x \in \overline{B_j}} \int_{B_{j^e}} J(x-y) dy > 0$$

and $\beta = \int_{\mathbb{R}^N} J(x) dx$. And we conclude as before. \square

Remark 2.6. Note that in [5] it is proved a Poincaré’s type inequality for Neumann boundary conditions, but assuming that $J(0) > 0$ (otherwise there is a counterexample). Surprisingly, for the Dirichlet problem we do not need positivity at the origin for J . This is due to the fact that for the Dirichlet problem the outside values influence the inside values.

In the next result we prove that $B_{p,\psi}^J$ is a completely accretive operator (see [14]) and verifies a range condition. In short, this means that for any $\phi \in L^p(\Omega)$ there is a unique solution of the problem $u + B_{p,\psi}^J u = \phi$ and the resolvent $(I + B_{p,\psi}^J)^{-1}$ is a contraction in $L^q(\Omega)$ for all $1 \leq q \leq +\infty$.

THEOREM 2.7. *Let $1 < p < +\infty$. For $\psi \in L^p(\Omega_J \setminus \overline{\Omega})$, the operator $B_{p,\psi}^J$ is completely accretive and verifies the range condition*

$$(2.5) \quad L^p(\Omega) \subset \text{Ran}(I + B_{p,\psi}^J).$$

Proof. Given $u_i \in \text{Dom}(B_{p,\psi}^J)$, $i = 1, 2$, by the monotonicity Lemma 2.4, for any $q \in C^\infty(\mathbb{R})$, $0 \leq q' \leq 1$, $\text{supp}(q')$ compact, $0 \notin \text{supp}(q)$, we have that

$$\int_{\Omega} (B_{p,\psi}^J u_1(x) - B_{p,\psi}^J u_2(x)) q(u_1(x) - u_2(x)) dx \geq 0,$$

from where it follows that $B_{p,\psi}^J$ is a completely accretive operator (see [14]).

To show that $B_{p,\psi}^J$ satisfies the range condition we have to prove that for any $\phi \in L^p(\Omega)$ there exists $u \in \text{Dom}(B_{p,\psi}^J)$ such that $\phi = u + B_{p,\psi}^J u$.

Assume first $p \geq 2$. Let $\phi \in L^p(\Omega)$ and set

$$K = \{w \in L^p(\Omega_J) : w = \psi \text{ in } \Omega_J \setminus \overline{\Omega}\}.$$

We consider the continuous monotone operator $A : K \rightarrow L^{p'}(\Omega_J)$ defined by

$$A(w)(x) := w(x) - \int_{\Omega_J} J(x - y) |w(y) - w(x)|^{p-2} (w(y) - w(x)) dy.$$

A is coercive in $L^p(\Omega_J)$. In fact, by Proposition 2.5, for any $w \in K$,

$$\begin{aligned} \int_{\Omega_J} A(w)w &= \int_{\Omega_J} w^2 - \int_{\Omega_J} \int_{\Omega_J} J(x - y) |w(y) - w(x)|^{p-2} (w(y) - w(x)) dy w(x) dx \\ &\geq \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x - y) |w(y) - w(x)|^p dy dx \\ &\geq \frac{1}{2} \int_{\Omega} \int_{\Omega_J} J(x - y) |w_\psi(y) - w(x)|^p dy dx \geq \frac{\lambda}{2} \|w\|_{L^p(\Omega)}^p - \frac{1}{2} \int_{\Omega_J \setminus \overline{\Omega}} |\psi|^p. \end{aligned}$$

Therefore,

$$\lim_{\substack{\|w\|_{L^p(\Omega_J)} \rightarrow +\infty \\ w \in K}} \frac{\int_{\Omega_J} A(w)w}{\|w\|_{L^p(\Omega_J)}} = +\infty.$$

Now, since $p \geq 2$, we have the function $\phi_\psi \in L^{p'}(\Omega_J)$. Then, applying [32, Corollary III.1.8] to the operator $B(w) := A(w) - \phi_\psi$, we get there exists $w \in K$, such that

$$w(x) - \int_{\Omega_J} J(x - y) |w(y) - w(x)|^{p-2} (w(y) - w(x)) dy = \phi_\psi(x) \quad \text{for all } x \in \Omega_J.$$

Hence, $u := w|_{\Omega}$ satisfies

$$u(x) - \int_{\Omega_J} J(x-y)|u_{\psi}(y) - u(x)|^{p-2}(u_{\psi}(y) - u(x)) dy = \phi(x) \quad \text{for all } x \in \Omega,$$

and, consequently, $\phi = u + B_{p,\psi}^J u$.

Suppose now $1 < p < 2$. By the results in [5], we know that the operator

$$B_p^J u(x) = - \int_{\Omega} J(x-y)|u(y) - u(x)|^{p-2}(u(y) - u(x)) dy$$

is m -accretive in $L^1(\Omega)$ and satisfies what is called property (M_0) ; that is, for any $q \in C^\infty(\mathbb{R})$, $0 \leq q' \leq 1$, $\text{supp}(q')$ compact, $0 \notin \text{supp}(q)$, and $(u, v) \in B_p^J$,

$$\int_{\Omega} q(u)v \geq 0.$$

On the other hand,

$$\varphi(x, r) = - \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)|\psi(y) - r|^{p-2}(\psi(y) - r) dy$$

is continuous and nondecreasing in r for almost every $x \in \Omega$, and an $L^1(\Omega)$ function for all r . Therefore, by [3, Theorem 3.1], $B_{p,\psi}^J u(x) = B_p^J u(x) + \varphi(x, u(x))$ is m -accretive in $L^1(\Omega)$. \square

Remark 2.8. If $\mathcal{B}_{p,\psi}^J$ denotes the closure of $B_{p,\psi}^J$ in $L^1(\Omega)$, by Theorem 2.7, we have $\mathcal{B}_{p,\psi}^J$ is m -completely accretive in $L^1(\Omega)$ (see [14]). Therefore, by the nonlinear semigroup theory (see [15] and [14]), there exists an unique mild-solution of the abstract Cauchy problem

$$(2.6) \quad \begin{cases} u'(t) + B_{p,\psi}^J u(t) = 0, & t \in (0, T), \\ u(0) = u_0, \end{cases}$$

given by the Crandall–Liggett exponential formula

$$e^{-t\mathcal{B}_{p,\psi}^J} u_0 = \lim_n \left(I + \frac{t}{n} \mathcal{B}_{p,\psi}^J \right)^{-n} u_0.$$

Now, due to regularity results for mild solutions, under certain hypothesis, this mild solution is a strong solution of the abstract Cauchy problem (2.6) (see [14]) which means, for our problem $P_p^J(u_0, \psi)$, a solution in the sense of Definition 1.1.

The following result states the existence and uniqueness results for $P_p^J(u_0, \psi)$. From it, Theorem 1.2 can be derived.

THEOREM 2.9. *Assume $p > 1$. Let $T > 0$, $\psi \in L^p(\Omega_J \setminus \bar{\Omega})$, and $u_0 \in L^1(\Omega)$. Then, there exists a unique mild-solution u of (2.6). Moreover,*

(1) *if $u_0 \in L^p(\Omega)$, the unique mild solution u of (2.6) is a solution of $P_p^J(u_0, \psi)$ in the sense of Definition 1.1. If $1 < p \leq 2$, this is true for any $u_0 \in L^1(\Omega)$ and any ψ such that $|\psi|^{p-1} \in L^1(\Omega_J \setminus \bar{\Omega})$.*

(2) *Let $u_{i0} \in L^1(\Omega)$ and u_i a solution in $[0, T]$ of $P_p^J(u_{i0})$, $i = 1, 2$. Then*

$$\int_{\Omega} (u_1(t) - u_2(t))^+ \leq \int_{\Omega} (u_{10} - u_{20})^+ \quad \text{for every } t \in]0, T[.$$

Moreover, for $q \in [1, +\infty]$, if $u_{i0} \in L^q(\Omega)$, $i = 1, 2$, then

$$\|u_1(t) - u_2(t)\|_{L^q(\Omega)} \leq \|u_{10} - u_{20}\|_{L^q(\Omega)} \quad \text{for every } t \in]0, T[.$$

Proof. As a consequence of Theorem 2.7 we get the existence of mild solution of (2.6) (see Remark 2.8). Now, due to the complete accretivity of $B_{p,\psi}^J$ and the range condition (2.5), by regularity results for mild solutions (see [14]), $u(t)$ is a strong solution, that is, a solution of $P_p^J(u_0, \psi)$ in the sense of Definition 1.1. Moreover, in the case $1 < p \leq 2$, since $\text{Dom}(B_{p,\psi}^J) = L^1(\Omega)$ and $B_{p,\psi}^J$ is closed in $L^1(\Omega) \times L^1(\Omega)$, the result holds for L^1 -data. Finally, the contraction principle is a consequence of the general nonlinear semigroup theory ([15], [28]). \square

2.2. Convergence to the p -Laplacian. Our main goal in this section is to show that the solution to the Dirichlet problem for the p -Laplacian equation $D_p(u_0, \tilde{\psi})$ can be approximated by solutions to suitable nonlocal Dirichlet problems $P_p^J(u_0, \psi)$.

Let us first recall the following result from [5]. For a function g defined in a set D , we define

$$\bar{g}(x) = \begin{cases} g(x) & \text{if } x \in D, \\ 0 & \text{otherwise,} \end{cases}$$

and we denote by χ_D the characteristic function of D .

PROPOSITION 2.10 ([5]). *Let $1 \leq q < +\infty$, D a bounded domain in \mathbb{R}^N , $\rho : \mathbb{R}^N \rightarrow \mathbb{R}$ a nonnegative continuous radial function with compact support, nonidentically zero, and $\rho_n(x) := n^N \rho(nx)$. Let $\{f_n\}$ be a sequence of functions in $L^q(D)$ such that*

$$(2.7) \quad \int_D \int_D |f_n(y) - f_n(x)|^q \rho_n(y - x) \, dx \, dy \leq M \frac{1}{n^q}.$$

1. *If $\{f_n\}$ is weakly convergent in $L^q(D)$ to f , then*

(i) *if $q > 1$, $f \in W^{1,q}(D)$ and moreover*

$$(\rho(z))^{1/q} \chi_D \left(x + \frac{1}{n} z \right) \frac{\bar{f}_n \left(x + \frac{1}{n} z \right) - f_n(x)}{1/n} \rightharpoonup (\rho(z))^{1/q} z \cdot \nabla f(x)$$

weakly in $L^q(D) \times L^q(\mathbb{R}^N)$;

(ii) *if $q = 1$, $f \in BV(D)$ and moreover*

$$\rho(z) \chi_D \left(\cdot + \frac{1}{n} z \right) \frac{\bar{f}_n \left(\cdot + \frac{1}{n} z \right) - f_n(\cdot)}{1/n} \rightharpoonup \rho(z) z \cdot Df$$

weakly as measures.

2. *Assume D is a smooth bounded domain in \mathbb{R}^N and $\rho(x) \geq \rho(y)$ if $|x| \leq |y|$. Then $\{f_n\}$ is relatively compact in $L^q(D)$ and, consequently, there exists a subsequence $\{f_{n_k}\}$ such that*

(i) *if $q > 1$, $f_{n_k} \rightarrow f$ in $L^q(D)$ with $f \in W^{1,q}(D)$;*

(ii) *If $q = 1$, $f_{n_k} \rightarrow f$ in $L^1(D)$ with $f \in BV(D)$.*

Let us now recall some results about the p -Laplacian equation

$$D_p(u_0, \tilde{\psi}) \quad \begin{cases} u_t = \Delta_p u & \text{in }]0, T[\times \Omega, \\ u = \tilde{\psi} & \text{on }]0, T[\times \partial\Omega, \\ u(0, x) = u_0(x) & \text{in } \Omega. \end{cases}$$

In the case $\tilde{\psi} \in W^{1/p',p}(\partial\Omega)$, associated to the p -Laplacian with nonhomogeneous Dirichlet boundary condition, in [2] it is defined the operator $A_{p,\tilde{\psi}} \subset L^1(\Omega) \times L^1(\Omega)$ as $(u, \hat{u}) \in A_{p,\tilde{\psi}}$ if and only if $\hat{u} \in L^1(\Omega)$, $u \in W^{1,p}_{\tilde{\psi}}(\Omega) := \{u \in W^{1,p}(\Omega) : u|_{\partial\Omega} = \tilde{\psi} \mathcal{H}^{N-1} - a.e. \text{ on } \partial\Omega\}$ and

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \cdot \nabla(u - v) \leq \int_{\Omega} \hat{u}(u - v) \quad \text{for every } v \in W^{1,p}_{\tilde{\psi}}(\Omega) \cap L^{\infty}(\Omega).$$

This inequality is equivalent to

$$\int_{\Omega} |\nabla u|^{p-2} \nabla u \cdot \nabla w = \int_{\Omega} \hat{u}w \quad \text{for every } w \in W^{1,p}_0(\Omega) \cap L^{\infty}(\Omega).$$

Moreover, for $\tilde{\psi} \in W^{1/p',p}(\partial\Omega) \cap L^{\infty}(\partial\Omega)$, $A_{p,\tilde{\psi}}$ is proved to be a completely accretive operator in $L^1(\Omega)$, satisfying the range condition $L^{\infty}(\Omega) \subset \text{Ran}(I + A_{p,\tilde{\psi}})$, and it is easy to see that $\overline{D(A_{p,\tilde{\psi}})}^{L^1(\Omega)} = L^1(\Omega)$. Therefore, its closure $\mathcal{A}_{p,\tilde{\psi}}$ in $L^1(\Omega) \times L^1(\Omega)$ is an m -completely accretive operator in $L^1(\Omega)$. Consequently, for any $u_0 \in L^1(\Omega)$ there exists a unique mild solution $u(t) = e^{-t\mathcal{A}_{p,\tilde{\psi}}}u_0$ of the abstract Cauchy problem associated to $D_p(u_0, \tilde{\psi})$, given by Crandall–Liggett’s exponential formula. Due to the complete accretivity of the operator $\mathcal{A}_{p,\tilde{\psi}}$, in the case $u_0 \in D(\mathcal{A}_{p,\tilde{\psi}})$ this mild solution is the unique strong solution of problem $D_p(u_0, \tilde{\psi})$.

In the homogeneous case $\tilde{\psi} = 0$, due to the results in [13], we can say that for any $u_0 \in L^1(\Omega)$, the mild solution $u(t) = e^{-t\mathcal{A}_{p,0}}u_0$ is the unique entropy solution of problem $D_p(u_0, 0)$.

For given $p > 1$ and J , we consider the rescaled kernels

$$J_{p,\varepsilon}(x) := \frac{C_{J,p}}{\varepsilon^{p+N}} J\left(\frac{x}{\varepsilon}\right), \quad \text{where} \quad C_{J,p}^{-1} := \frac{1}{2} \int_{\mathbb{R}^N} J(z)|z_N|^p dz$$

is a normalizing constant in order to obtain the p -Laplacian in the limit instead of a multiple of it.

PROPOSITION 2.11. *Let Ω be a smooth bounded domain in \mathbb{R}^N and let $\tilde{\psi} \in W^{1/p',p}(\partial\Omega) \cap L^{\infty}(\partial\Omega)$. Let $\psi \in W^{1,p}(\Omega_J) \cap L^{\infty}(\Omega_J)$ such that $\psi|_{\partial\Omega} = \tilde{\psi}$. Assume $J(x) \geq J(y)$ if $|x| \leq |y|$. Then, for any $\phi \in L^{\infty}(\Omega)$,*

$$(2.8) \quad \left(I + B_{p,\psi}^{J_{p,\varepsilon}}\right)^{-1} \phi \rightarrow \left(I + A_{p,\tilde{\psi}}\right)^{-1} \phi \quad \text{in } L^p(\Omega) \text{ as } \varepsilon \rightarrow 0.$$

Proof. We denote

$$\Omega_{\varepsilon} := \Omega_{J_{p,\varepsilon}} = \Omega + \text{supp}(J_{p,\varepsilon}).$$

For $\varepsilon > 0$ small, let $u_{\varepsilon} = \left(I + B_{p,\psi}^{J_{p,\varepsilon}}\right)^{-1} \phi$. Then,

$$(2.9) \quad \int_{\Omega} u_{\varepsilon} v - \frac{C_{J,p}}{\varepsilon^{p+N}} \int_{\Omega} \int_{\Omega_{\varepsilon}} J\left(\frac{x-y}{\varepsilon}\right) |(u_{\varepsilon})_{\psi}(y) - u_{\varepsilon}(x)|^{p-2} \times ((u_{\varepsilon})_{\psi}(y) - u_{\varepsilon}(x)) dy v(x) dx = \int_{\Omega} \phi v$$

for every $v \in L^{\infty}(\Omega)$.

Let $M := \max\{\|\phi\|_{L^\infty(\Omega)}, \|\psi\|_{L^\infty(\Omega_J)}\}$. Taking $v = (u_\epsilon - M)^+$ in (2.9), we get

$$\begin{aligned} & \int_{\Omega} u_\epsilon(x)(u_\epsilon(x) - M)^+ dx - \frac{C_{J,p}}{\epsilon^{p+N}} \int_{\Omega} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) |(u_\epsilon)_\psi(y) - u_\epsilon(x)|^{p-2} \\ & \quad \times ((u_\epsilon)_\psi(y) - u_\epsilon(x)) dy (u_\epsilon(x) - M)^+ dx \\ & = \int_{\Omega} \phi(x)(u_\epsilon(x) - M)^+ dx. \end{aligned}$$

Now,

$$\begin{aligned} & -\frac{C_{J,p}}{\epsilon^{p+N}} \int_{\Omega} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) |(u_\epsilon)_\psi(y) - u_\epsilon(x)|^{p-2} ((u_\epsilon)_\psi(y) - u_\epsilon(x)) dy \\ & \quad \times (u_\epsilon(x) - M)^+ dx \\ & = -\frac{C_{J,p}}{\epsilon^{p+N}} \int_{\Omega_\epsilon} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) |(u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)|^{p-2} ((u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)) dy \\ & \quad \times ((u_\epsilon)_\psi(x) - M)^+ dx \\ & = \frac{C_{J,p}}{2\epsilon^{p+N}} \int_{\Omega_\epsilon} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) |(u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)|^{p-2} ((u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)) \\ & \quad \times (((u_\epsilon)_\psi(y) - M)^+ - ((u_\epsilon)_\psi(x) - M)^+) dy dx \\ & \geq 0. \end{aligned}$$

Therefore,

$$\int_{\Omega} u_\epsilon(x)(u_\epsilon(x) - M)^+ dx \leq \int_{\Omega} \phi(x)(u_\epsilon(x) - M)^+ dx.$$

Consequently, we have

$$\int_{\Omega} (u_\epsilon(x) - M)(u_\epsilon(x) - M)^+ dx \leq \int_{\Omega} (\phi(x) - M)(u_\epsilon(x) - M)^+ dx \leq 0,$$

and $u_\epsilon(x) \leq M$ for almost all $x \in \Omega$. Analogously, we can obtain $-M \leq u_\epsilon(x)$ for almost all $x \in \Omega$. Thus

$$(2.10) \quad \|u_\epsilon\|_{L^\infty(\Omega)} \leq M \quad \text{for all } \epsilon > 0,$$

and, therefore, there exists a sequence $\epsilon_n \rightarrow 0$ such that

$$u_{\epsilon_n} \rightharpoonup u \quad \text{weakly in } L^1(\Omega).$$

Taking $v = u_\epsilon - \psi$ in (2.9) we get

$$(2.11) \quad \begin{aligned} & \int_{\Omega} u_\epsilon(u_\epsilon - \psi) - \frac{C_{J,p}}{\epsilon^{p+N}} \int_{\Omega_\epsilon} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) |(u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)|^{p-2} \\ & \quad \times ((u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)) dy ((u_\epsilon)_\psi(x) - \psi(x)) dx = \int_{\Omega} \phi(u_\epsilon - \psi). \end{aligned}$$

Now, by (2.11) and (2.10),

$$\begin{aligned} & \frac{C_{J,p}}{2\epsilon^N} \int_{\Omega_\epsilon} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) \frac{|(u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)|^p}{\epsilon^p} dy dx \\ & \leq \frac{C_{J,p}}{2\epsilon^N} \int_{\Omega_\epsilon} \int_{\Omega_\epsilon} J\left(\frac{x-y}{\epsilon}\right) \frac{|(u_\epsilon)_\psi(y) - (u_\epsilon)_\psi(x)|^{p-1}}{\epsilon^{p-1}} \frac{|\psi(y) - \psi(x)|}{\epsilon} dy dx + M_1. \end{aligned}$$

Since $\psi \in W^{1,p}(\Omega_J)$, using Young's inequality, we obtain

$$\frac{1}{\varepsilon^N} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) \frac{|(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)|^p}{\varepsilon^p} dy dx \leq M_2.$$

Moreover,

$$\begin{aligned} & \int_{\Omega_J} \int_{\Omega_J} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)}{\varepsilon} \right|^p dx dy \\ &= \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)}{\varepsilon} \right|^p dx dy \\ & \quad + 2 \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \int_{\Omega_\varepsilon} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{\psi(y) - (u_\varepsilon)_\psi(x)}{\varepsilon} \right|^p dx dy \\ (2.12) \quad & \quad + \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{\psi(y) - \psi(x)}{\varepsilon} \right|^p dx dy \\ &= \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)}{\varepsilon} \right|^p dx dy \\ & \quad + 2 \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \int_{\Omega_\varepsilon \setminus \overline{\Omega}} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{\psi(y) - \psi(x)}{\varepsilon} \right|^p dx dy \\ & \quad + \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \int_{\Omega_J \setminus \overline{\Omega_\varepsilon}} \frac{1}{\varepsilon^N} J\left(\frac{x-y}{\varepsilon}\right) \left| \frac{\psi(y) - \psi(x)}{\varepsilon} \right|^p dx dy \leq M_3. \end{aligned}$$

Therefore, by Proposition 2.10, there exists a subsequence, denoted as above, and $w \in W^{1,p}(\Omega_J)$ such that

$$(u_{\varepsilon_n})_\psi \rightarrow w \quad \text{strongly in } L^p(\Omega_J).$$

Hence, $w = u$ in Ω and, by [18, Proposition IX.18] and the properties of the trace, $u \in W^{1,p}_\psi(\Omega)$. Moreover, by Proposition 2.10,

$$(2.13) \quad \left(\frac{C_{J,p}}{2} J(z)\right)^{1/p} \chi_\Omega(x + \varepsilon_n z) \frac{(u_{\varepsilon_n})_\psi(x + \varepsilon_n z) - (u_{\varepsilon_n})_\psi(x)}{\varepsilon_n} \rightharpoonup \left(\frac{C_{J,p}}{2} J(z)\right)^{1/p} z \cdot \nabla u(x)$$

weakly in $L^p(\Omega) \times L^p(\mathbb{R}^N)$ (observe that $\chi_\Omega(x + \varepsilon_n z)(u_{\varepsilon_n})_\psi(x + \varepsilon_n z) = \chi_\Omega(x + \varepsilon_n z)\bar{u}_{\varepsilon_n}(x + \varepsilon_n z)$). We can also assume that

$$\begin{aligned} & (J(z))^{1/p'} \left| \frac{(u_{\varepsilon_n})_\psi(x + \varepsilon_n z) - (u_{\varepsilon_n})_\psi(x)}{\varepsilon_n} \right|^{p-2} \chi_{\Omega_{\varepsilon_n}}(x + \varepsilon_n z) \\ & \quad \times \frac{(u_{\varepsilon_n})_\psi(x + \varepsilon_n z) - (u_{\varepsilon_n})_\psi(x)}{\varepsilon_n} \rightharpoonup (J(z))^{1/p'} \chi(x, z) \end{aligned}$$

weakly in $L^{p'}(\Omega_J) \times L^{p'}(\mathbb{R}^N)$, for some function $\chi \in L^{p'}(\Omega_J) \times L^{p'}(\mathbb{R}^N)$.

Passing to the limit in (2.9) for $\varepsilon = \varepsilon_n$, we get

$$(2.14) \quad \int_{\Omega} uv + \int_{\mathbb{R}^N} \int_{\Omega} \frac{C_{J,p}}{2} J(z) \chi(x, z) z \cdot \nabla v(x) dx dz = \int_{\Omega} \phi v$$

for every v smooth with support in Ω and by approximation for every $v \in W^{1,p}_0(\Omega)$.

Finally, working as in Proposition 3.3. of [5], we can prove

$$(2.15) \quad \int_{\mathbb{R}^N} \int_{\Omega} \frac{C_{J,p}}{2} J(z) \chi(x, z) z \cdot \nabla v(x) \, dx \, dz = \int_{\Omega} |\nabla u|^{p-2} \nabla u \cdot \nabla v$$

and the proof is finished. \square

From the above Proposition, by the standard results of the nonlinear semigroup theory (see [19] or [15]), we obtain Theorem 1.3.

3. The nonlocal total variation flow. The case $p = 1$.

3.1. Existence of solutions for the nonlocal problem. This section deals with the existence and uniqueness of solutions for the nonlocal 1-Laplacian problem with Dirichlet boundary condition,

$$P_1^J(u_0, \psi) \begin{cases} u_t(t, x) = \int_{\Omega} J(x-y) \frac{u(t, y) - u(t, x)}{|u(t, y) - u(t, x)|} \, dy \\ \quad + \int_{\Omega_J \setminus \bar{\Omega}} J(x-y) \frac{\psi(y) - u(t, x)}{|\psi(y) - u(t, x)|} \, dy, & x \in \Omega. \\ u(0, x) = u_0(x). \end{cases}$$

As in the case $p > 1$, to prove existence and uniqueness of solutions of $P_1^J(u_0, \psi)$ we use the Nonlinear Semigroup Theory, so we start by introducing the following operator in $L^1(\Omega)$.

DEFINITION 3.1. *Given $\psi \in L^1(\Omega_J \setminus \bar{\Omega})$, we define the operator $B_{1,\psi}^J$ in $L^1(\Omega) \times L^1(\Omega)$ by $\hat{u} \in B_{1,\psi}^J u$ if and only if $u, \hat{u} \in L^1(\Omega)$, there exists $g \in L^\infty(\Omega_J \times \Omega_J)$, $g(x, y) = -g(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, $\|g\|_\infty \leq 1$,*

$$(3.1) \quad \hat{u}(x) = - \int_{\Omega_J} J(x-y) g(x, y) \, dy \quad \text{a.e. } x \in \Omega,$$

and

$$(3.2) \quad J(x-y) g(x, y) \in J(x-y) \text{sign}(u(y) - u(x)) \quad \text{a.e. } (x, y) \in \Omega \times \Omega,$$

$$(3.3) \quad J(x-y) g(x, y) \in J(x-y) \text{sign}(\psi(y) - u(x)) \quad \text{a.e. } (x, y) \in \Omega \times (\Omega_J \setminus \bar{\Omega}).$$

Remark 3.2. Observe that

(i) we can rewrite (3.2) + (3.3) as

$$(3.4) \quad J(x-y) g(x, y) \in J(x-y) \text{sign}(u_\psi(y) - u(x)) \quad \text{a.e. } (x, y) \in \Omega \times \Omega_J,$$

where we set as above, and overall the section,

$$u_\psi(x) := \begin{cases} u(x) & \text{if } x \in \Omega, \\ \psi(x) & \text{if } x \in \Omega_J \setminus \bar{\Omega}, \\ 0 & \text{if } x \notin \Omega_J. \end{cases}$$

(ii) It holds $L^1(\Omega) = \text{Dom}(B_{1,\psi}^J)$ and $B_{1,\psi}^J$ is closed in $L^1(\Omega) \times L^1(\Omega)$.

(iii) It is not difficult to see that, if $g \in L^\infty(\Omega_J \times \Omega_J)$, $g(x, y) = -g(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, $\|g\|_\infty \leq 1$,

$$J(x-y) g(x, y) \in J(x-y) \text{sign}(z(y) - z(x)) \quad \text{a.e. } (x, y) \in \Omega_J \times \Omega_J$$

is equivalent to

$$-\int_{\Omega_J} \int_{\Omega_J} J(x-y)g(x,y) dy z(x) dx = \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y)|z(y) - z(x)| dy dx.$$

THEOREM 3.3. *Let $\psi \in L^1(\Omega_J \setminus \bar{\Omega})$. The operator $B_{1,\psi}^J$ is completely accretive and satisfies the range condition*

$$L^\infty(\Omega) \subset \text{Ran}(I + B_{1,\psi}^J).$$

Proof. Let $\hat{u}_i \in B_{1,\psi}^J u_i$, $i = 1, 2$, and set $u_i(y) = \psi(y)$ in $\Omega_J \setminus \bar{\Omega}$. Then, there exist $g_i \in L^\infty(\Omega_J \times \Omega_J)$, $\|g_i\|_\infty \leq 1$, $g_i(x, y) = -g_i(y, x)$, $J(x-y)g_i(x, y) \in J(x-y)\text{sign}(u_i(y) - u_i(x))$ for almost all $(x, y) \in \Omega \times \Omega_J$, such that

$$\hat{u}_i(x) = -\int_{\Omega_J} J(x-y)g_i(x, y) dy \quad \text{a.e. } x \in \Omega$$

for $i = 1, 2$. Given $q \in C^\infty(\mathbb{R})$, $0 \leq q' \leq 1$, $\text{supp}(q')$ compact, $0 \notin \text{supp}(q)$, we have

$$\begin{aligned} & \int_{\Omega} (\hat{u}_1(x) - \hat{u}_2(x))q(u_1(x) - u_2(x)) dx \\ &= \frac{1}{2} \int_{\Omega} \int_{\Omega} J(x-y)(g_1(x, y) - g_2(x, y)) (q(u_1(y) - u_2(y)) - q(u_1(x) - u_2(x))) dx dy \\ & \quad - \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)(g_1(x, y) - g_2(x, y)) (q(u_1(x) - u_2(x))) dx dy \\ & \geq \frac{1}{2} \int_{\Omega} \int_{\Omega} J(x-y)(g_1(x, y) - g_2(x, y)) (q(u_1(y) - u_2(y)) - q(u_1(x) - u_2(x))) dx dy. \end{aligned}$$

Now, by the mean value theorem

$$\begin{aligned} & J(x-y)(g_1(x, y) - g_2(x, y)) [q(u_1(y) - u_2(y)) - q(u_1(x) - u_2(x))] \\ &= J(x-y)(g_1(x, y) - g_2(x, y))q'(\xi) [(u_1(y) - u_2(y)) - (u_1(x) - u_2(x))] \\ &= J(x-y)q'(\xi) [g_1(x, y)(u_1(y) - u_1(x)) - g_1(x, y)(u_2(y) - u_2(x))] \\ & \quad - J(x-y)q'(\xi) [g_2(x, y)(u_1(y) - u_1(x)) - g_1(x, y)(u_2(y) - u_2(x))] \geq 0, \end{aligned}$$

since

$$J(x-y)g_i(x, y)(u_i(y) - u_i(x)) = J(x-y)|u_i(y) - u_i(x)|, \quad i = 1, 2,$$

and

$$-J(x-y)g_i(x, y)(u_j(y) - u_j(x)) \geq -J(x-y)|u_j(y) - u_j(x)|, \quad i \neq j.$$

Hence

$$\int_{\Omega} (\hat{u}_1(x) - \hat{u}_2(x))q(u_1(x) - u_2(x)) dx \geq 0,$$

from which it follows that $B_{1,\psi}^J$ is a completely accretive operator.

To show that $B_{1,\psi}^J$ satisfies the range condition, let us see that for any $\phi \in L^\infty(\Omega)$,

$$\lim_{p \rightarrow 1^+} (I + B_{p,\psi}^J)^{-1} \phi = (I + B_{1,\psi}^J)^{-1} \phi \quad \text{weakly in } L^1(\Omega).$$

We prove this in several steps.

Step 1. Let us first suppose that $\psi \in L^\infty(\Omega_J \setminus \bar{\Omega})$. For $1 < p < +\infty$, by Theorem 2.7, there is u_p such that $u_p = (I + B_{p,\psi}^J)^{-1}\phi$, that is,

$$(3.5) \quad u_p(x) - \int_{\Omega_J} J(x-y) |(u_p)_\psi(y) - u_p(x)|^{p-2} ((u_p)_\psi(y) - u_p(x)) dy = \phi(x),$$

a.e. $x \in \Omega$. It is easy to see that $\|u_p\|_\infty \leq \sup\{\|\phi\|_\infty, \|\psi\|_\infty\}$. Therefore, there exists a sequence $p_n \rightarrow 1$ such that

$$u_{p_n} \rightharpoonup u \quad \text{weakly in } L^2(\Omega).$$

On the other hand, we also have

$$\frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) |(u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)|^{p_n} dy dx \leq M_2, \quad \forall n \in \mathbb{N}.$$

Consequently, for any measurable subset $E \subset \Omega_J \times \Omega_J$, we have

$$\begin{aligned} & \left| \int \int_E J(x-y) |(u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)|^{p_n-2} ((u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)) \right| \\ & \leq \int \int_E J(x-y) |(u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)|^{p_n-1} \leq M_2 |E|^{\frac{1}{p_n}}. \end{aligned}$$

Hence, by the Dunford–Pettis theorem we may assume that there exists $g(x, y)$ such that

$$J(x-y) |(u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)|^{p_n-2} ((u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)) \rightharpoonup J(x-y)g(x, y),$$

weakly in $L^1(\Omega_J \times \Omega_J)$, $g(x, y) = -g(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, and $\|g\|_\infty \leq 1$.

Therefore, by (3.5),

$$(3.6) \quad u(x) - \int_{\Omega_J} J(x-y)g(x, y) dy = \phi(x) \quad \text{a.e. } x \in \Omega.$$

Then, to finish the proof it is enough to show that

$$(3.7) \quad \begin{aligned} & - \int_{\Omega_J} \int_{\Omega_J} J(x-y)g(x, y) dy u_\psi(x) dx \\ & = \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) |u_\psi(y) - u_\psi(x)| dy dx. \end{aligned}$$

In fact, by (3.5) and (3.6),

$$\begin{aligned} & \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) |(u_{p_n})_\psi(y) - (u_{p_n})_\psi(x)|^{p_n} dy dx = \int_{\Omega} \phi u_{p_n} - \int_{\Omega} u_{p_n} u_{p_n} \\ & - \int_{\Omega_J \setminus \bar{\Omega}} \int_{\Omega_J} J(x-y) |\psi(y) - (u_{p_n})_\psi(x)|^{p_n-2} (\psi(y) - (u_{p_n})_\psi(x)) dy \psi(x) dx \\ & = \int_{\Omega} \phi u - \int_{\Omega} uu - \int_{\Omega} \phi(u - u_{p_n}) + \int_{\Omega} 2u(u - u_{p_n}) - \int_{\Omega} (u - u_{p_n})(u - u_{p_n}) \\ & \quad - \int_{\Omega_J \setminus \bar{\Omega}} \int_{\Omega_J} J(x-y) |\psi(y) - (u_{p_n})_\psi(x)|^{p_n-2} (\psi(y) - (u_{p_n})_\psi(x)) dy \psi(x) dx \\ & \leq - \int_{\Omega_J} \int_{\Omega_J} J(x-y) g(x, y) dy u(x) dx - \int_{\Omega} \phi(u - u_{p_n}) + \int_{\Omega} 2u(u - u_{p_n}) \\ & \quad + \int_{\Omega_J \setminus \bar{\Omega}} \int_{\Omega_J} J(x-y) g(x, y) dy \psi(x) dx \\ & \quad - \int_{\Omega_J \setminus \bar{\Omega}} \int_{\Omega_J} J(x-y) |\psi(y) - (u_{p_n})_\psi(x)|^{p_n-2} (\psi(y) - (u_{p_n})_\psi(x)) dy \psi(x) dx, \end{aligned}$$

and so,

$$\begin{aligned} & \limsup_{n \rightarrow +\infty} \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) |u_{p_n}(y) - u_{p_n}(x)|^{p_n} dy dx \\ & \leq - \int_{\Omega} \int_{\Omega} J(x-y) g(x, y) dy u(x) dx. \end{aligned}$$

Now, by the monotonicity Lemma 2.4, for all $\rho \in L^\infty(\Omega)$,

$$\begin{aligned} & - \int_{\Omega_J} \int_{\Omega_J} J(x-y) |\rho(y) - \rho(x)|^{p_n-2} (\rho(y) - \rho(x)) dy (u_{p_n}(x) - \rho(x)) dx \\ & \leq - \int_{\Omega_J} \int_{\Omega_J} J(x-y) |u_{p_n}(y) - u_{p_n}(x)|^{p_n-2} (u_{p_n}(y) - u_{p_n}(x)) dy (u_{p_n}(x) - \rho(x)) dx. \end{aligned}$$

Taking limits,

$$\begin{aligned} & - \int_{\Omega_J} \int_{\Omega_J} J(x-y) \text{sign}_0(\rho(y) - \rho(x)) dy (u(x) - \rho(x)) dx \\ & \leq - \int_{\Omega_J} \int_{\Omega_J} J(x-y) g(x, y) dy (u(x) - \rho(x)) dx. \end{aligned}$$

Taking now, $\rho = u \pm \lambda u$, $\lambda > 0$, and letting $\lambda \rightarrow 0$, we get (3.7), and the proof is finished for this class of data.

Step 2. Let us now suppose that ψ^- is bounded. Let $\psi_n = T_n(\psi)$, n large enough such that $\psi_n^- = \psi^-$. Then, $\{\psi_n\}$ is a nondecreasing sequence that converges in L^1 to ψ . By *Step 1*, there exists $u_n = (I + B_{1, \psi_n}^J)^{-1} \phi$, that is, there exists $g_n \in L^\infty(\Omega_J \times \Omega_J)$, $g_n(x, y) = -g_n(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, $\|g_n\|_\infty \leq 1$,

$$(3.8) \quad u_n(x) - \int_{\Omega_J} J(x-y) g_n(x, y) dy = \phi(x) \quad \text{a.e. } x \in \Omega$$

and

$$\begin{aligned} (3.9) \quad & - \int_{\Omega_J} \int_{\Omega_J} J(x-y) g_n(x, y) dy (u_n)_{\psi_n}(x) dx \\ & = \frac{1}{2} \int_{\Omega_J} \int_{\Omega_J} J(x-y) |(u_n)_{\psi_n}(y) - (u_n)_{\psi_n}(x)| dy dx. \end{aligned}$$

Therefore, by monotonicity,

$$\int_{\Omega_J} \int_{\Omega_J} ((u_n)_{\psi_n} - (u_{n+1})_{\psi_{n+1}}) ((u_n)_{\psi_n} - (u_{n+1})_{\psi_{n+1}})^+ \leq 0,$$

which implies $u_n \leq u_{n+1}$. Since $\{u_n\}$ is bounded in L^∞ we have $\{u_n\}$ converges to a function u in L^2 . On the other hand, we can suppose that $J(x - y)g_n(x, y)$ converges weakly in L^2 to $J(x - y)g(x, y)$, $g(x, y) = -g(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, and $\|g\|_\infty \leq 1$. Hence, passing to the limit in (3.8) and (3.9) we obtain $u = (I + B_{1,\psi}^J)^{-1}\phi$.

Step 3. For a general $\psi \in L^1(\Omega_J \setminus \bar{\Omega})$, apply *Step 2* to $\psi_n = \sup\{\psi, -n\}$ and use monotonicity in a similar way to finish the proof. \square

Proof of Theorem 1.5. As a consequence of the above results, we have that the abstract Cauchy problem

$$(3.10) \quad \begin{cases} u'(t) + B_{1,\psi}^J u(t) \ni 0, & t \in (0, T), \\ u(0) = u_0 \end{cases}$$

has a unique mild solution u for every initial datum $u_0 \in L^1(\Omega)$ and $T > 0$ (see [15]). Moreover, due to the complete accretivity of the operator $B_{1,\psi}^J$, the mild solution of (3.10) is a strong solution ([14]). Consequently, the proof is concluded. \square

3.2. Convergence to the total variation flow. Let us start recalling some results from [1] (see also [2]) about the Dirichlet problem for the total variational flow, that is,

$$D_1(u_0, \tilde{\psi}) \quad \begin{cases} u_t = \operatorname{div} \left(\frac{Du}{|Du|} \right) & \text{in }]0, T[\times \Omega, \\ u = \tilde{\psi} & \text{on }]0, T[\times \partial\Omega, \\ u(0, x) = u_0(x) & \text{in } \Omega, \end{cases}$$

with $\tilde{\psi} \in L^1(\partial\Omega)$.

THEOREM 3.4 ([1]). *Let $T > 0$ and $\tilde{\psi} \in L^1(\partial\Omega)$. For any $u_0 \in L^1(\Omega)$ ($L^2(\Omega)$) there exists a unique entropy (strong) solution $u(t)$ of $D_1(u_0, \tilde{\psi})$.*

Associated to $-\operatorname{div}(\frac{Du}{|Du|})$ with Dirichlet boundary conditions, in [1] it is defined the operator $\mathcal{A}_{\tilde{\psi}} \subset L^1(\Omega) \times L^1(\Omega)$ as follows: $(u, v) \in \mathcal{A}_{\tilde{\psi}}$ if and only if $u, v \in L^1(\Omega)$, $q(u) \in BV(\Omega)$ for all $q \in \mathcal{P} := \{q \in W^{1,\infty}(\mathbb{R}) : q' \geq 0, \operatorname{supp}(q') \text{ is compact}\}$, and there exists $\zeta \in X(\Omega)$ (where $X(\Omega)$ is defined by (1.4)), with $\|\zeta\|_\infty \leq 1$, $v = -\operatorname{div}(\zeta)$ in $\mathcal{D}'(\Omega)$ such that

$$(3.11) \quad \int_{\Omega} (w - q(u))v \leq \int_{\Omega} (\zeta, Dw) - |Dq(u)| + \int_{\partial\Omega} |w - q(\tilde{\psi})| - \int_{\partial\Omega} |q(u) - q(\tilde{\psi})|$$

for every $w \in BV(\Omega) \cap L^\infty(\Omega)$ and every $q \in \mathcal{P}$. Also in [1] it is proved that the following assertions are equivalent:

- (a) $(u, v) \in \mathcal{A}_{\tilde{\psi}}$,
- (b) $u, v \in L^1(\Omega)$, $q(u) \in BV(\Omega)$ for all $q \in \mathcal{P}$, and there exists $\zeta \in X(\Omega)$, with $\|\zeta\|_\infty \leq 1$, $v = -\operatorname{div}(\zeta)$ in $\mathcal{D}'(\Omega)$ such that

$$(3.12) \quad \int_{\Omega} (\zeta, Dq(u)) = |Dq(u)| \quad \forall q \in \mathcal{P},$$

$$(3.13) \quad [\zeta, \nu] \in \operatorname{sign} \left(q(\tilde{\psi}) - q(u) \right) \quad \mathcal{H}^{N-1} - a.e. \text{ on } \partial\Omega, \quad \forall q \in \mathcal{P}.$$

Moreover, it is shown that $\mathcal{A}_{\tilde{\psi}}$ is an m -completely accretive operator in $L^1(\Omega)$ with dense domain and that for any $u_0 \in L^1(\Omega)$, the unique entropy solution $u(t)$ of problem $D_1(u_0, \tilde{\psi})$ coincides with the unique mild solution $e^{-t\mathcal{A}_{\tilde{\psi}}}u_0$ given by Crandall–Liggett’s exponential formula.

Now, given J , we consider the rescaled kernels

$$J_{1,\varepsilon}(x) := \frac{C_{J,1}}{\varepsilon^{1+N}} J\left(\frac{x}{\varepsilon}\right), \quad \text{with} \quad C_{J,1}^{-1} := \frac{1}{2} \int_{\mathbb{R}^N} J(z)|z_N| dz,$$

that is, a normalizing constant in order to obtain the 1-Laplacian in the limit instead of a multiple of it.

PROPOSITION 3.5. *Let Ω be a smooth bounded domain in \mathbb{R}^N and $\tilde{\psi} \in L^\infty(\partial\Omega)$. Let $\psi \in W^{1,1}(\Omega_J \setminus \bar{\Omega}) \cap L^\infty(\Omega_J \setminus \bar{\Omega})$ such that $\psi|_{\partial\Omega} = \tilde{\psi}$. Assume $J(x) \geq J(y)$ if $|x| \leq |y|$. Then, for any $\phi \in L^\infty(\Omega)$,*

$$(3.14) \quad \left(I + B_{1,\psi}^{J_{1,\varepsilon}}\right)^{-1} \phi \rightarrow \left(I + \mathcal{A}_{\tilde{\psi}}\right)^{-1} \phi \quad \text{strongly in } L^1(\Omega) \text{ as } \varepsilon \rightarrow 0.$$

Proof. Given $\varepsilon > 0$ small, we set $u_\varepsilon = (I + B_{1,\psi}^{J_{1,\varepsilon}})^{-1}\phi$ and denote

$$\Omega_\varepsilon := \Omega_{J_{1,\varepsilon}} = \Omega + \text{supp}(J_{1,\varepsilon}).$$

Then, there exists $g_\varepsilon \in L^\infty(\Omega_\varepsilon \times \Omega_\varepsilon)$, $g_\varepsilon(x, y) = -g_\varepsilon(y, x)$ for almost all $(x, y) \in \Omega_\varepsilon \times \Omega_\varepsilon$, $\|g_\varepsilon\|_\infty \leq 1$, such that

$$\begin{aligned} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y) &\in J\left(\frac{x-y}{\varepsilon}\right) \text{sign}(u_\varepsilon(y) - u_\varepsilon(x)) && \text{a.e. } (x, y) \in \Omega \times \Omega, \\ J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y) &\in J\left(\frac{x-y}{\varepsilon}\right) \text{sign}(\tilde{\psi}(y) - u_\varepsilon(x)) && \text{a.e. } (x, y) \in \Omega \times (\Omega_\varepsilon \setminus \bar{\Omega}) \end{aligned}$$

and

$$(3.15) \quad u_\varepsilon(x) - \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y) dy = \phi(x) \quad \text{a.e. } x \in \Omega.$$

Therefore, for $v \in L^\infty(\Omega_J)$, we can write

$$(3.16) \quad \begin{aligned} \int_{\Omega} u_\varepsilon(x)v(x) dx - \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y)v(x) dy dx \\ = \int_{\Omega} \phi(x)v(x) dx. \end{aligned}$$

Observe that we can extend g_ε to a function in $L^\infty(\Omega_J \times \Omega_J)$, $g_\varepsilon(x, y) = -g_\varepsilon(y, x)$ for almost all $(x, y) \in \Omega_J \times \Omega_J$, $\|g_\varepsilon\|_{L^\infty(\Omega_J)} \leq 1$, such that

$$J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y) \in J\left(\frac{x-y}{\varepsilon}\right) \text{sign}((u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)) \quad \text{a.e. } (x, y) \in \Omega_J \times \Omega_J.$$

Let $M := \max\{\|\phi\|_{L^\infty(\Omega)}, \|\psi\|_{L^\infty(\Omega_J \setminus \bar{\Omega})}\}$. Taking $v = (u_\varepsilon - M)^+$ in (3.16), we get

$$\begin{aligned} \int_{\Omega} u_\varepsilon(x)(u_\varepsilon(x) - M)^+ dx - \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x, y)(u_\varepsilon(x) - M)^+ dy dx \\ = \int_{\Omega} \phi(x)(u_\varepsilon(x) - M)^+ dx. \end{aligned}$$

Now

$$\begin{aligned} & -\frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) ((u_\varepsilon)_\psi(x) - M)^+ dy dx \\ &= \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) (((u_\varepsilon)_\psi(y) - M)^+ - ((u_\varepsilon)_\psi(x) - M)^+) dy dx \\ &\geq 0. \end{aligned}$$

Hence, we get

$$\int_{\Omega} u_\varepsilon(x)(u_\varepsilon(x) - M)^+ dx \leq \int_{\Omega} \phi(x)(u_\varepsilon(x) - M)^+ dx.$$

Consequently,

$$0 \leq \int_{\Omega} (u_\varepsilon(x) - M)(u_\varepsilon(x) - M)^+ dx \leq \int_{\Omega} (\phi(x) - M)(u_\varepsilon(x) - M)^+ dx \leq 0,$$

and we deduce $u_\varepsilon(x) \leq M$ for almost all $x \in \Omega$. Analogously, we can obtain $-M \leq u_\varepsilon(x)$ for almost all $x \in \Omega$. Thus

$$(3.17) \quad \|u_\varepsilon\|_{L^\infty(\Omega)} \leq M \quad \text{for all } \varepsilon > 0;$$

from here, we can assume there exists a sequence $\varepsilon_n \rightarrow 0$ such that

$$u_{\varepsilon_n} \rightharpoonup u \quad \text{weakly in } L^1(\Omega).$$

Taking $v = u_\varepsilon$ in (3.16), we have

$$(3.18) \quad \int_{\Omega} u_\varepsilon(x)u_\varepsilon(x) dx - \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy u_\varepsilon(x) dx = \int_{\Omega} \phi(x)u_\varepsilon(x) dx.$$

Observe that

$$\begin{aligned} & -\frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy u_\varepsilon(x) dx \\ &= -\frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy (u_\varepsilon)_\psi(x) dx \\ &\quad + \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon \setminus \bar{\Omega}} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy \psi(x) dx. \end{aligned}$$

Then

$$\begin{aligned} & \left| \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon \setminus \bar{\Omega}} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy \psi(x) dx \right| \\ &\leq \frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon \setminus \bar{\Omega}} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) dy |\psi(x)| dx \\ &\leq \frac{C_{J,1}}{\varepsilon} M \int_{\Omega_\varepsilon \setminus \bar{\Omega}} \left(\frac{1}{\varepsilon^N} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) dy \right) dx \\ &\leq \frac{C_{J,1}}{\varepsilon} M |\Omega_\varepsilon \setminus \Omega| \leq M_1. \end{aligned}$$

On the other hand,

$$\begin{aligned} & -\frac{C_{J,1}}{\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) g_\varepsilon(x,y) dy (u_\varepsilon)_\psi(x) dx \\ & = \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx. \end{aligned}$$

Consequently, from (3.17) and (3.18), it follows that

$$(3.19) \quad \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \leq M_2.$$

Let us compute,

$$\begin{aligned} & \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J} \int_{\Omega_J} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & = \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & \quad + \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & \quad + \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & \quad + \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx. \end{aligned}$$

Now, since $\psi \in W^{1,1}(\Omega_J \setminus \bar{\Omega})$, we get

$$\begin{aligned} & \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & = \frac{C_{J,1}}{2\varepsilon^N} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) \frac{|\psi(y) - \psi(x)|}{\varepsilon} dy dx \leq M_3. \end{aligned}$$

On the other hand, we have

$$\begin{aligned} & \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \\ & = \frac{C_{J,1}}{2\varepsilon^N} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \int_{\Omega_\varepsilon \setminus \bar{\Omega}} J\left(\frac{x-y}{\varepsilon}\right) \frac{|\psi(y) - \psi(x)|}{\varepsilon} dy dx \\ & \leq M_4 \frac{C_{J,1}}{2} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} \left(\frac{1}{\varepsilon^N} \int_{\Omega_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) dy \right) dx \leq M_5. \end{aligned}$$

With similar arguments we obtain

$$\frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_\varepsilon} \int_{\Omega_J \setminus \bar{\Omega}_\varepsilon} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \leq M_6.$$

Therefore,

$$(3.20) \quad \frac{C_{J,1}}{2\varepsilon^{1+N}} \int_{\Omega_J} \int_{\Omega_J} J\left(\frac{x-y}{\varepsilon}\right) |(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)| dy dx \leq M_7.$$

In particular, we get

$$\int_{\Omega_J} \int_{\Omega_J} \frac{1}{2} \frac{C_{J,1}}{\varepsilon^N} J \left(\frac{x-y}{\varepsilon} \right) \left| \frac{(u_\varepsilon)_\psi(y) - (u_\varepsilon)_\psi(x)}{\varepsilon} \right| dx dy \leq M_7 \quad \forall n \in \mathbb{N}.$$

By Proposition 2.10, there exists a subsequence, denote equal, and $w \in BV(\Omega_J)$ such that

$$(u_{\varepsilon_n})_\psi \rightarrow w \quad \text{strongly in } L^1(\Omega_J)$$

and

$$(3.21) \quad \frac{C_{J,1}}{2} J(z) \chi_{\Omega}(\cdot + \varepsilon_n z) \frac{(u_{\varepsilon_n})_\psi(\cdot + \varepsilon_n z) - (u_{\varepsilon_n})_\psi(\cdot)}{\varepsilon_n} \rightharpoonup \frac{C_{J,1}}{2} J(z) z \cdot Dw$$

weakly as measures. Hence, it is easy to obtain that

$$w(x) = u_\psi(x) = \begin{cases} u(x) & \text{in } x \in \Omega, \\ \psi(x) & \text{in } x \in \Omega_J \setminus \overline{\Omega}, \end{cases}$$

and $u \in BV(\Omega)$.

Moreover, we can also assume that

$$(3.22) \quad J(z) \chi_{\Omega_J}(x + \varepsilon_n z) \bar{g}_{\varepsilon_n}(x, x + \varepsilon_n z) \rightharpoonup \Lambda(x, z)$$

weakly* in $L^\infty(\Omega_J) \times L^\infty(\mathbb{R}^N)$ for some function $\Lambda \in L^\infty(\Omega_J) \times L^\infty(\mathbb{R}^N)$, $\Lambda(x, z) \leq J(z)$ almost everywhere in $\Omega_J \times \mathbb{R}^N$. Taking in (3.16) $v \in \mathcal{D}(\Omega)$, we get for $\varepsilon = \varepsilon_n$ small enough

$$(3.23) \quad \begin{aligned} \int_{\Omega} u_{\varepsilon_n}(x) v(x) dx - \frac{C_{J,1}}{\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J \left(\frac{x-y}{\varepsilon_n} \right) g_{\varepsilon_n}(x, y) v(x) dy dx \\ = \int_{\Omega} \phi(x) v(x) dx. \end{aligned}$$

Changing variables and taking into account (3.23), we can write

$$(3.24) \quad \begin{aligned} & \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \int_{\Omega} J(z) \chi_{\Omega}(x + \varepsilon_n z) \bar{g}_{\varepsilon_n}(x, x + \varepsilon_n z) dz \frac{\bar{v}(x + \varepsilon_n z) - v(x)}{\varepsilon_n} dx \\ & = - \frac{C_{J,1}}{\varepsilon_n} \int_{\mathbb{R}^N} \int_{\Omega} J(z) \chi_{\Omega}(x + \varepsilon_n z) \bar{g}_{\varepsilon_n}(x, x + \varepsilon_n z) dz v(x) dx \\ & = \int_{\Omega} (\phi(x) - u_{\varepsilon_n}(x)) v(x) dx. \end{aligned}$$

By (3.22), passing to the limit in (3.24), we get

$$(3.25) \quad \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \int_{\Omega} \Lambda(x, z) z \cdot \nabla v(x) dx dz = \int_{\Omega} (\phi(x) - u(x)) v(x) dx$$

for all $v \in \mathcal{D}(\Omega)$. We set $\zeta = (\zeta_1, \dots, \zeta_N)$, the vector field defined by

$$\zeta_i(x) := \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \Lambda(x, z) z_i dz, \quad i = 1, \dots, N.$$

Then, $\zeta \in L^\infty(\Omega_J, \mathbb{R}^N)$, and from (3.25),

$$-\operatorname{div}(\zeta) = \phi - u \quad \text{in } \mathcal{D}'(\Omega).$$

Let us see that

$$\|\zeta\|_{L^\infty(\Omega_J)} \leq 1.$$

Given $\xi \in \mathbb{R}^N \setminus \{0\}$, let R_ξ be the rotation such that $R_\xi^t(\xi) = \mathbf{e}_1|\xi|$. If we make the change of variables $z = R_\xi(y)$, we obtain

$$\begin{aligned} \zeta(x) \cdot \xi &= \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \Lambda(x, z) z \cdot \xi \, dz = \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \Lambda(x, R_\xi(y)) R_\xi(y) \cdot \xi \, dy \\ &= \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} \Lambda(x, R_\xi(y)) y_1 |\xi| \, dy. \end{aligned}$$

On the other hand, since J is a radial function and $\Lambda(x, z) \leq J(z)$ almost everywhere,

$$C_{J,1}^{-1} = \frac{1}{2} \int_{\mathbb{R}^N} J(z) |z_1| \, dz$$

and

$$|\zeta(x) \cdot \xi| \leq \frac{C_{J,1}}{2} \int_{\mathbb{R}^N} J(y) |y_1| \, dy |\xi| = |\xi| \quad \text{a.e. } x \in \Omega_J.$$

Therefore, $\|\zeta\|_{L^\infty(\Omega_J)} \leq 1$.

To finish the proof, that is, to show that $u = (I + \mathcal{A}_{\tilde{\psi}})^{-1}\phi$, since $u \in L^\infty(\Omega)$ and $\tilde{\psi} \in L^\infty(\partial\Omega)$, we need only to prove that

$$(3.26) \quad (\zeta, Du) = |Du| \quad \text{as measures in } \Omega$$

and

$$(3.27) \quad [\zeta, \nu] \in \operatorname{sign}(\tilde{\psi} - u) \quad \mathcal{H}^{N-1} - \text{a.e. on } \partial\Omega.$$

Given $0 \leq \varphi \in \mathcal{D}(\Omega)$, taking $\varepsilon = \varepsilon_n$ and $v = \varphi u_{\varepsilon_n}$ in (3.16), we get

$$\begin{aligned} & -\frac{C_{J,1}}{\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x, y) u_{\varepsilon_n}(x) \varphi(x) \, dy \, dx \\ (3.28) \quad & = \frac{C_{J,1}}{2\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x, y) (u_{\varepsilon_n}(y)\varphi(y) - u_{\varepsilon_n}(x)\varphi(x)) \, dy \, dx \\ & = \int_{\Omega} (\phi(x) - u_{\varepsilon_n}(x)) u_{\varepsilon_n}(x) \varphi(x) \, dx. \end{aligned}$$

Now, we decompose the double integral as follows,

$$I_n := \frac{C_{J,1}}{2\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x, y) (u_{\varepsilon_n}(y)\varphi(y) - u_{\varepsilon_n}(x)\varphi(x)) \, dy \, dx = I_n^1 + I_n^2,$$

where

$$\begin{aligned} I_n^1 &:= \frac{C_{J,1}}{2\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J\left(\frac{x-y}{\varepsilon_n}\right) |u_{\varepsilon_n}(y) - u_{\varepsilon_n}(x)| \varphi(y) \, dy \, dx \\ &= \frac{C_{J,1}}{2} \int_{\Omega} \int_{\Omega} J(z) \chi_{\Omega}(x + \varepsilon_n z) \frac{|\bar{u}_{\varepsilon_n}(x + \varepsilon_n z) - u_{\varepsilon_n}(x)|}{\varepsilon_n} \varphi(x + \varepsilon_n z) \, dz \, dx \end{aligned}$$

and

$$\begin{aligned} I_n^2 &:= \frac{C_{J,1}}{2\varepsilon_n^{1+N}} \int_{\Omega} \int_{\Omega} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x,y) u_{\varepsilon_n}(x) (\varphi(y) - \varphi(x)) dy dx \\ &= \frac{C_{J,1}}{2} \int_{\Omega} \int_{\Omega} J(z) \chi_{\Omega}(x + \varepsilon_n z) \bar{g}_{\varepsilon_n}(x, x + \varepsilon_n z) u_{\varepsilon_n}(x) \frac{\bar{\varphi}(x + \varepsilon_n z) - \varphi(x)}{\varepsilon_n} dz dx. \end{aligned}$$

Having in mind (3.21), it follows that

$$\lim_{n \rightarrow \infty} I_n^1 \geq \frac{C_{J,1}}{2} \int_{\Omega} \int_{\Omega} J(z) \varphi(x) |z \cdot Du| = \int_{\Omega} \varphi |Du|.$$

On the other hand, since

$$u_{\varepsilon_n} \rightarrow u \quad \text{strongly in } L^1(\Omega),$$

by (3.22), we get

$$\lim_{n \rightarrow \infty} I_n^2 = \frac{C_{J,1}}{2} \int_{\Omega} \int_{\mathbb{R}^N} u(x) \Lambda(x, z) z \cdot \nabla \varphi(x) dz dx = \int_{\Omega} u(x) \zeta(x) \cdot \nabla \varphi(x) dx.$$

Therefore, taking $n \rightarrow +\infty$ in (3.28), we obtain

$$(3.29) \quad \int_{\Omega} \varphi |Du| + \int_{\Omega} u(x) \zeta(x) \cdot \nabla \varphi(x) dx \leq \int_{\Omega} (\phi(x) - u(x)) u(x) \varphi(x) dx.$$

By Green's formula,

$$\begin{aligned} \int_{\Omega} (\phi(x) - u(x)) u(x) \varphi(x) dx &= - \int_{\Omega} \operatorname{div}(\zeta) u \varphi dx = \int_{\Omega} (\zeta, D(\varphi u)) \\ &= \int_{\Omega} \varphi(\zeta, Du) + \int_{\Omega} u(x) \zeta(x) \cdot \nabla \varphi(x) dx. \end{aligned}$$

Since $|(\zeta, Du)| \leq |Du|$, the last identity and (3.29) give (3.26).

Finally, we show that (3.27) holds. We take $w_m \in W^{1,1}(\Omega) \cap C(\Omega)$ such that $w_m = \tilde{\psi} \mathcal{H}^{N-1}$ -a.e. on $\partial\Omega$, and $w_m \rightarrow u$ in $L^1(\Omega)$. Taking $v = v_{m,n} := (u_{\varepsilon_n})_{\tilde{\psi}} - (w_m)_{\tilde{\psi}}$ in (3.16), we get

$$\begin{aligned} (3.30) \quad & \int_{\Omega} (\phi(x) - u_{\varepsilon_n}(x))(u_{\varepsilon_n}(x) - w_m(x)) dx \\ &= - \frac{C_{J,1}}{\varepsilon_n^{1+N}} \int_{\Omega_J} \int_{\Omega_J} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x,y) v_{m,n}(x) dy dx \\ &= \frac{C_{J,1}}{2\varepsilon_n^{1+N}} \int_{\Omega_J} \int_{\Omega_J} J\left(\frac{x-y}{\varepsilon_n}\right) g_{\varepsilon_n}(x,y) (v_{m,n}(x) - v_{m,n}(y)) dy dx \\ &= H_n^1 + H_{m,n}^1, \end{aligned}$$

where

$$H_n^1 = \frac{C_{J,1}}{2} \int_{\Omega_J} \int_{\mathbb{R}^N} J(z) \chi_{\Omega_J}(x + \varepsilon_n z) \left| \frac{(u_{\varepsilon_n})_{\tilde{\psi}}(x + \varepsilon_n z) - (u_{\varepsilon_n})_{\tilde{\psi}}(x)}{\varepsilon_n} \right| dz dx$$

and

$$\begin{aligned} H_{m,n}^2 &= - \frac{C_{J,1}}{2} \int_{\Omega_J} \int_{\mathbb{R}^N} J(z) \chi_{\Omega_J}(x + \varepsilon_n z) \bar{g}_{\varepsilon_n}(x, x + \varepsilon_n z) \\ &\quad \times \frac{(w_m)_{\tilde{\psi}}(x + \varepsilon_n z) - (w_m)_{\tilde{\psi}}(x)}{\varepsilon_n} dz dx. \end{aligned}$$

Arguing as before,

$$\lim_{n \rightarrow \infty} H_n^1 \geq \int_{\Omega_J} |Du_\psi| = \int_{\Omega} |Du| + \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} + \int_{\Omega_J \setminus \bar{\Omega}} |\nabla\psi|.$$

On the other hand, since $(w_m)_\psi \in W^{1,1}(\Omega_J)$, by (3.22),

$$\lim_{n \rightarrow \infty} H_{m,n}^2 = -\frac{C_{J,1}}{2} \int_{\Omega_J} \int_{\mathbb{R}^N} \Lambda(x, z) z \cdot \nabla(w_m)_\psi(x) dz dx = -\int_{\Omega_J} \zeta(x) \cdot \nabla(w_m)_\psi(x) dx.$$

Consequently, taking $n \rightarrow \infty$ in (3.30), we get

$$(3.31) \quad \begin{aligned} & \int_{\Omega} (\phi(x) - u(x))(u(x) - w_m(x)) dx \\ & \geq \int_{\Omega} |Du| + \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} + \int_{\Omega_J \setminus \bar{\Omega}} |\nabla\psi| - \int_{\Omega_J} \zeta(x) \cdot \nabla(w_m)_\psi(x) dx. \end{aligned}$$

Now,

$$\begin{aligned} -\int_{\Omega_J} \zeta(x) \cdot \nabla(w_m)_\psi(x) dx &= -\int_{\Omega} \zeta(x) \cdot \nabla w_m(x) dx - \int_{\Omega_J \setminus \bar{\Omega}} \zeta(x) \cdot \nabla\psi(x) dx \\ &= \int_{\Omega} \operatorname{div}\zeta(x)w_m(x) dx - \int_{\partial\Omega} [\zeta, \nu]\tilde{\psi} d\mathcal{H}^{N-1} - \int_{\Omega_J \setminus \bar{\Omega}} \zeta(x) \cdot \nabla\psi(x) dx. \end{aligned}$$

Since

$$\int_{\Omega_J \setminus \bar{\Omega}} |\nabla\psi| - \int_{\Omega_J \setminus \bar{\Omega}} \zeta(x) \cdot \nabla\psi(x) dx \geq 0,$$

from (3.31), we have

$$\begin{aligned} & \int_{\Omega} (\phi(x) - u(x))(u(x) - w_m(x)) dx \\ & \geq \int_{\Omega} |Du| + \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} + \int_{\Omega} \operatorname{div}\zeta(x)w_m(x) dx - \int_{\partial\Omega} [\zeta, \nu]\tilde{\psi} d\mathcal{H}^{N-1}. \end{aligned}$$

Letting $m \rightarrow \infty$, and using Green's formula, we deduce

$$\begin{aligned} 0 & \geq \int_{\Omega} |Du| + \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} + \int_{\Omega} \operatorname{div}\zeta(x)u(x) dx - \int_{\partial\Omega} [\zeta, \nu]\tilde{\psi} d\mathcal{H}^{N-1} \\ & = \int_{\Omega} |Du| + \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} - \int_{\Omega} (\zeta, Du) + \int_{\partial\Omega} [\zeta, \nu]u d\mathcal{H}^{N-1} \\ & \quad - \int_{\partial\Omega} [\zeta, \nu]\tilde{\psi} d\mathcal{H}^{N-1}. \end{aligned}$$

By (3.26), we obtain

$$\int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1} \leq \int_{\partial\Omega} [\zeta, \nu](\tilde{\psi} - u) d\mathcal{H}^{N-1} \leq \int_{\partial\Omega} |u - \tilde{\psi}| d\mathcal{H}^{N-1}.$$

Therefore,

$$[\zeta, \nu] \in \operatorname{sign}(\tilde{\psi} - u) \quad \mathcal{H}^{N-1} - \text{a.e. on } \partial\Omega,$$

and the proof is finished. \square

From the above Proposition, by standard results of the Nonlinear Semigroup Theory (see, [19] or [15]), we obtain Theorem 1.6.

4. Limit as $p \rightarrow +\infty$. A model for sandpiles.

4.1. A model for sandpiles. Let sand be poured out onto a rigid surface, $y = u_0(x)$, given in a bounded open subset Ω of \mathbb{R}^2 with Lipschitz boundary $\partial\Omega$. If the support boundary is open and we assume that the angle of stability is equal to $\frac{\pi}{4}$, a model for pile surface evolution was proposed by Prigozhin [35] as

$$(4.1) \quad \partial_t u + \operatorname{div} \mathbf{q} = f, \quad u|_{t=0} = u_0, \quad u|_{\partial\Omega} = u_0|_{\partial\Omega},$$

where $u(t, x)$ is the unknown pile surface, $f(t, x) \geq 0$ is the given source density, and $\mathbf{q}(t, x)$ is the unknown horizontal projection of the flux of sand pouring down the pile surface. If the support has no slopes steeper than the sand angle of repose, $\|\nabla u_0\|_\infty \leq 1$, Prigozhin ([35], see also [10], [29], and the references therein) proposed to take $\mathbf{q} = -m\nabla u$, where $m \geq 0$ is the Lagrange multiplier related to the constraint $\|\nabla u\|_\infty \leq 1$ and satisfies $m(\|\nabla u\|^2 - 1) = 0$ and reformulated this model as the following variational inequality:

$$(4.2) \quad \begin{cases} f(t, \cdot) - u_t(t) \in \partial I_{K(u_0)}(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x), \end{cases}$$

where

$$K(u_0) := \{v \in W^{1,\infty}(\Omega) : \|\nabla v\|_\infty \leq 1, v|_{\partial\Omega} = u_0|_{\partial\Omega}\}.$$

Our aim is to approximate the Prigozhin model for the sandpile by a nonlocal model (Theorem 1.8) obtained as the limit as $p \rightarrow +\infty$ of the nonlocal p -Laplacian problem with Dirichlet boundary condition (Theorem 1.7).

To identify the limit as $p \rightarrow +\infty$ of the solutions u_p of problem $P_p^J(u_0, \psi)$ we will use the methods of convex analysis, and so we first recall some terminology (see [30], [17], and [8]). If H is a real Hilbert space with inner product (\cdot, \cdot) and $\Psi : H \rightarrow (-\infty, +\infty]$ is convex, then the subdifferential of Ψ is defined as the multivalued operator $\partial\Psi$ given by

$$v \in \partial\Psi(u) \iff \Psi(w) - \Psi(u) \geq (v, w - u) \quad \forall w \in H.$$

The epigraph of Ψ is defined by $\operatorname{Epi}(\Psi) = \{(u, \lambda) \in H \times \mathbb{R} : \lambda \geq \Psi(u)\}$.

Given K a closed convex subset of H , the indicator function of K is defined by

$$I_K(u) = \begin{cases} 0 & \text{if } u \in K, \\ +\infty & \text{if } u \notin K. \end{cases}$$

Then it is easy to see that the subdifferential is characterized as follows:

$$v \in \partial I_K(u) \iff u \in K \text{ and } (v, w - u) \leq 0 \quad \forall w \in K.$$

In case the convex functional $\Psi : H \rightarrow (-\infty, +\infty]$ is proper, lower-semicontinuous, and $\min \Psi = 0$, it is well known (see [17]) that the abstract Cauchy problem

$$\begin{cases} u'(t) + \partial\Psi(u(t)) \ni f(t), & \text{a.e. } t \in]0, T[, \\ u(0) = u_0, \end{cases}$$

has a unique strong solution for any $f \in L^2(0, T; H)$ and $u_0 \in \overline{D(\partial\Psi)}$.

The following convergence was studied by Mosco in [34] (see [8]). Suppose X is a metric space and $A_n \subset X$. We define

$$\liminf_{n \rightarrow \infty} A_n = \{x \in X : \exists x_n \in A_n, x_n \rightarrow x\}$$

and

$$\limsup_{n \rightarrow \infty} A_n = \{x \in X : \exists x_{n_k} \in A_{n_k}, x_{n_k} \rightarrow x\}.$$

In the case X is a normed space, we note by $s - \lim$ and $w - \lim$ the above limits associated, respectively, to the strong and to the weak topology of X .

Given a sequence $\Psi_n, \Psi : H \rightarrow (-\infty, +\infty]$ of convex lower-semicontinuous functionals, we say that Ψ_n converges to Ψ in the sense of Mosco if

$$(4.3) \quad w - \limsup_{n \rightarrow \infty} \text{Epi}(\Psi_n) \subset \text{Epi}(\Psi) \subset s - \liminf_{n \rightarrow \infty} \text{Epi}(\Psi_n).$$

It is easy to see that (4.3) is equivalent to the two following conditions:

$$(4.4) \quad \forall u \in D(\Psi) \exists u_n \in D(\Psi_n) : u_n \rightarrow u \text{ and } \Psi(u) \geq \limsup_{n \rightarrow \infty} \Psi_n(u_n);$$

$$(4.5) \quad \text{for every subsequence } n_k, \text{ when } u_k \rightharpoonup u, \text{ it holds } \Psi(u) \leq \liminf_k \Psi_{n_k}(u_k).$$

As a consequence of the results in [19] and [8] we can write the following result.

THEOREM 4.1. *Let $\Psi_n, \Psi : H \rightarrow (-\infty, +\infty]$ convex lower-semicontinuous functionals. Then the following statements are equivalent.*

- (i) Ψ_n converges to Ψ in the sense of Mosco.
- (ii) $(I + \lambda \partial \Psi_n)^{-1} u \rightarrow (I + \lambda \partial \Psi)^{-1} u, \quad \forall \lambda > 0, u \in H.$

Moreover, any of these two conditions (i) or (ii) imply that

- (iii) for every $u_0 \in D(\partial \Psi)$ and $u_{0,n} \in \overline{D(\partial \Psi_n)}$ such that $u_{0,n} \rightarrow u_0$, and every $f_n, f \in L^2(0, T; H)$ with $f_n \rightarrow f$, if $u_n(t), u(t)$ are the strong solutions of the abstract Cauchy problems

$$\begin{cases} u'_n(t) + \partial \Psi_n(u_n(t)) \ni f_n, & a.e. \ t \in]0, T[, \\ u_n(0) = u_{0,n}, \end{cases}$$

and

$$\begin{cases} u'(t) + \partial \Psi(u(t)) \ni f, & a.e. \ t \in]0, T[, \\ u(0) = u_0, \end{cases}$$

respectively, then

$$u_n \rightarrow u \quad \text{in } C([0, T] : H).$$

4.2. Limit as $p \rightarrow +\infty$. Let us consider the nonlocal p -Laplacian evolution problem with source

$$P_p^J(u_0, \psi, f) \begin{cases} u_t(t, x) = \int_{\Omega} J(x - y) |u(t, y) - u(t, x)|^{p-2} (u(t, y) - u(t, x)) dy + f(t, x), \\ \hspace{15em} (t, x) \in]0, T[\times \Omega, \\ u(t, x) = \psi(x), \quad (t, x) \in]0, T[\times (\Omega_J \setminus \overline{\Omega}), \\ u(0, x) = u_0(x), \quad x \in \Omega. \end{cases}$$

This problem is associated to the energy functional

$$G_{p,\psi}^J(u) = \frac{1}{2p} \int_{\Omega} \int_{\Omega} J(x-y)|u(y) - u(x)|^p dy dx + \frac{1}{p} \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)|\psi(y) - u(x)|^p dy dx.$$

With a formal calculation, taking limit as $p \rightarrow +\infty$, we arrive to the functional

$$G_{\infty,\psi}^J(u) = \begin{cases} 0 & \text{if } |u(x) - u(y)| \leq 1, \text{ for } x, y \in \bar{\Omega} \\ & \text{and } |\psi(y) - u(x)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ & \text{with } x - y \in \text{supp}(J) \\ +\infty & \text{in the other case.} \end{cases}$$

Hence, if we define

$$K_{\infty,\psi}^J := \left\{ u \in L^2(\Omega) : \begin{array}{l} |u(x) - u(y)| \leq 1, x, y \in \Omega \\ \text{and } |\psi(y) - u(x)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ \text{with } x - y \in \text{supp}(J) \end{array} \right\},$$

we have that the functional $G_{\infty,\psi}^J$ is given by the indicator function of $K_{\infty,\psi}^J$; that is, $G_{\infty,\psi}^J = I_{K_{\infty,\psi}^J}$. Then, the *nonlocal limit problem* can be written as

$$P_{\infty}^J(u_0, \psi, f) \quad \begin{cases} f(t, \cdot) - u_t(t) \in \partial I_{K_{\infty,\psi}^J}(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x). \end{cases}$$

Proof of Theorem 1.7. Let $T > 0$. By Theorem 4.1, to prove the result it is enough to show that the functionals

$$G_{p,\psi}^J(u) = \frac{1}{2p} \int_{\Omega} \int_{\Omega} J(x-y)|u(y) - u(x)|^p dy dx + \frac{1}{p} \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y)|\psi(y) - u(x)|^p dy dx$$

converge to

$$G_{\infty,\psi}^J(u) = \begin{cases} 0 & \text{if } |u(x) - u(y)| \leq 1, \text{ for } x, y \in \Omega \\ & \text{and } |\psi(y) - u(x)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \bar{\Omega}, \\ & \text{with } x - y \in \text{supp}(J) \\ +\infty & \text{in the other case} \end{cases}$$

as $p \rightarrow +\infty$, in the sense of Mosco. First, let us check that

$$(4.6) \quad \text{Epi}(G_{\infty,\psi}^J) \subset s - \liminf_{p \rightarrow +\infty} \text{Epi}(G_{p,\psi}^J).$$

To this end let $(u, \lambda) \in \text{Epi}(G_{\infty,\psi}^J)$. We can assume that $u \in K_{\infty,\psi}^J$ and $\lambda \geq 0$ (as $G_{\infty,\psi}^J(u) = 0$). Now take

$$(4.7) \quad v_p = u \quad \text{and} \quad \lambda_p = G_{p,\psi}^J(u) + \lambda.$$

Then, as $\lambda \geq 0$ we have $(v_p, \lambda_p) \in \text{Epi}(G_{p,\psi}^J)$. Obviously, $v_p = u \rightarrow u$ in $L^2(\Omega)$, and as $u \in K_{\infty,\psi}^J$,

$$\begin{aligned} G_{p,\psi}^J(u) &= \frac{1}{2p} \int_{\Omega} \int_{\Omega} J(x-y) |u(y) - u(x)|^p \, dy \, dx \\ &\quad + \frac{1}{p} \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y) |\psi(y) - u(x)|^p \, dy \, dx \\ &\leq \frac{1}{2p} \int_{\Omega} \int_{\Omega} J(x-y) \, dy \, dx + \frac{1}{p} \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y) \, dy \, dx \rightarrow 0 \end{aligned}$$

as $p \rightarrow +\infty$. Therefore, we get (4.6). Finally, let us prove that

$$(4.8) \quad w - \limsup_{p \rightarrow +\infty} \text{Epi}(G_{p,\psi}^J) \subset \text{Epi}(G_{\infty,\psi}^J).$$

To this end, let us consider a sequence $(u_{p_j}, \lambda_{p_j}) \in \text{Epi}(G_{p_j,\psi}^J)$; that is, $G_{p_j,\psi}^J(u_{p_j}) \leq \lambda_{p_j}$, with

$$u_{p_j} \rightharpoonup u, \quad \text{and} \quad \lambda_{p_j} \rightarrow \lambda.$$

Since, $0 \leq G_{p_j,\psi}^J(u_{p_j}) \leq \lambda_{p_j} \rightarrow \lambda$, $0 \leq \lambda$. On the other hand, we have that there exists a constant $C > 0$ such that

$$\begin{aligned} (p_j C)^{1/p_j} &\geq (p_j G_{p_j,\psi}^J(u_{p_j}))^{1/p_j} = \left(\frac{1}{2} \int_{\Omega} \int_{\Omega} J(x-y) |u_{p_j}(y) - u_{p_j}(x)|^{p_j} \, dy \, dx \right. \\ &\quad \left. + \int_{\Omega} \int_{\Omega_J \setminus \bar{\Omega}} J(x-y) |\psi(y) - u_{p_j}(x)|^{p_j} \, dy \, dx \right)^{1/p_j}. \end{aligned}$$

Then, by the above inequality,

$$\begin{aligned} &\left(\int_{\Omega} \int_{\Omega} J(x-y) |u_{p_j}(y) - u_{p_j}(x)|^q \, dy \, dx \right)^{1/q} \\ &\leq \left(\int_{\Omega} \int_{\Omega} J(x-y) \, dy \, dx \right)^{(p_j-q)/p_j q} \\ &\quad \times \left(\int_{\Omega} \int_{\Omega} J(x-y) |u_{p_j}(y) - u_{p_j}(x)|^{p_j} \, dy \, dx \right)^{1/p_j} \\ &\leq \left(\int_{\Omega} \int_{\Omega} J(x-y) \, dy \, dx \right)^{(p_j-q)/p_j q} (C p_j)^{1/p_j}. \end{aligned}$$

Hence, we can extract a subsequence (if necessary) and let $p_j \rightarrow +\infty$ to obtain

$$\left(\int_{\Omega} \int_{\Omega} J(x-y) |u(y) - u(x)|^q \, dy \, dx \right)^{1/q} \leq \left(\int_{\Omega} \int_{\Omega} J(x-y) \, dy \, dx \right)^{1/q}.$$

Now, just taking $q \rightarrow +\infty$, we get

$$|u(x) - u(y)| \leq 1 \quad \text{a.e. } (x, y) \in \Omega \times \Omega, \quad x - y \in \text{supp}(J).$$

With a similar argument we obtain

$$|u(x) - \psi(y)| \leq 1 \quad \text{a.e. } x \in \Omega, \quad y \in \Omega_J \setminus \bar{\Omega}, \quad \text{with } x - y \in \text{supp}(J).$$

Hence, we conclude that $u \in K_{\infty,\psi}^J$. This ends the proof. \square

4.3. Rescaling. We will assume now that Ω is convex and ψ verifies $\|\nabla\psi\|_\infty \leq 1$. For $\varepsilon > 0$, we rescale the functional $G_{\infty,\psi}^J$, as follows:

$$G_{\infty,\psi}^\varepsilon(u) = \begin{cases} 0 & \text{if } |u(x) - u(y)| \leq \varepsilon, \text{ for } x, y \in \Omega \\ & \text{and } |\psi(y) - u(x)| \leq \varepsilon, \text{ for } x \in \Omega, y \in \Omega_J \setminus \overline{\Omega}, \\ & \text{with } |x - y| \leq \varepsilon \\ +\infty & \text{in the other case.} \end{cases}$$

In other words, $G_{\infty,\psi}^\varepsilon = I_{K_{\infty,\psi}^\varepsilon}$, where

$$K_{\infty,\psi}^\varepsilon := \left\{ u \in L^2(\Omega) : \begin{array}{l} |u(x) - u(y)| \leq \varepsilon, x, y \in \Omega \\ \text{and } |\psi(y) - u(x)| \leq \varepsilon, \text{ for } x \in \Omega, y \in \Omega_J \setminus \overline{\Omega}, \\ \text{with } |x - y| \leq \varepsilon \end{array} \right\}.$$

Consider the gradient flow associated to the functional $G_{\infty,\psi}^\varepsilon$

$$P_\infty^\varepsilon(u_0, \psi, f) \begin{cases} f(t, \cdot) - u_t(t, \cdot) \in \partial I_{K_{\infty,\psi}^\varepsilon}(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x), & \text{in } \Omega, \end{cases}$$

and the problem

$$P_\infty(u_0, \psi, f) \begin{cases} f(t, \cdot) - u_{\infty,t} \in \partial I_{K_\psi}(u_\infty), & \text{a.e. } t \in]0, T[, \\ u_\infty(0, x) = u_0(x), & \text{in } \Omega, \end{cases}$$

where

$$K_\psi := \{u \in W^{1,\infty}(\Omega) : \|\nabla u\|_\infty \leq 1, u|_{\partial\Omega} = \psi|_{\partial\Omega}\}.$$

Observe that if $u \in K_\psi$, $\|\nabla u\|_\infty \leq 1$. Then, since $\|\nabla\psi\|_\infty \leq 1$ and Ω is convex, we have $|u(x) - u(y)| \leq |x - y|$ and $|u(x) - \psi(y)| \leq |x - y|$, from where it follows that $u \in K_{\infty,\psi}^\varepsilon$, that is, $K_\psi \subset K_{\infty,\psi}^\varepsilon$.

With all these definitions and notations, we can proceed with the limit as $\varepsilon \rightarrow 0$ for the sandpile model ($p = +\infty$).

Proof of Theorem 1.8. Since $u_0 \in K_\psi$, $u_0 \in K_{\infty,\psi}^\varepsilon$ for all $\varepsilon > 0$. Again we are using that $\|\nabla\psi\|_\infty \leq 1$. Consequently, there exists $u_{\infty,\varepsilon}$ the unique solution of $P_\infty^\varepsilon(u_0, \psi, f)$.

By Theorem 4.1, to prove the result it is enough to show that $I_{K_{\infty,\psi}^\varepsilon}$ converges to I_{K_ψ} in the sense of Mosco. Using that $\|\nabla\psi\|_\infty \leq 1$ it is easy to obtain that

$$(4.9) \quad K_{\infty,\psi}^{\varepsilon_1} \subset K_{\infty,\psi}^{\varepsilon_2}, \quad \text{if } \varepsilon_1 \leq \varepsilon_2.$$

Since $K_\psi \subset K_{\infty,\psi}^\varepsilon$ for all $\varepsilon > 0$, we have

$$K_\psi \subset \bigcap_{\varepsilon>0} K_{\infty,\psi}^\varepsilon.$$

On the other hand, if

$$u \in \bigcap_{\varepsilon>0} K_{\infty,\psi}^\varepsilon,$$

we have

$$|u(y) - u(x)| \leq |y - x|, \quad \text{a.e. } x, y \in \Omega,$$

and moreover

$$|u(y) - \psi(x)| \leq |y - x|, \quad \text{a.e. } x \in \Omega_J \setminus \overline{\Omega}, y \in \Omega,$$

from where it follows that $u \in K_\psi$. Therefore, we have

$$(4.10) \quad K_\psi = \bigcap_{\varepsilon > 0} K_{\infty, \psi}^\varepsilon.$$

Note that

$$(4.11) \quad \text{Epi}(I_{K_\psi}) = K_\psi \times [0, \infty[, \quad \text{Epi}(I_{K_{\infty, \psi}^\varepsilon}) = K_{\infty, \psi}^\varepsilon \times [0, \infty[\quad \forall \varepsilon > 0.$$

By (4.10) and (4.11),

$$(4.12) \quad \text{Epi}(I_{K_\psi}) \subset s - \liminf_{\varepsilon \rightarrow 0} \text{Epi}(I_{K_{\infty, \psi}^\varepsilon}).$$

On the other hand, given $(u, \lambda) \in w - \limsup_{\varepsilon \rightarrow 0} \text{Epi}(I_{K_{\infty, \psi}^\varepsilon})$ there exists $(u_{\varepsilon_k}, \lambda_k) \in K_{\varepsilon_k, \psi} \times [0, \infty[$, such that $\varepsilon_k \rightarrow 0$ and

$$u_{\varepsilon_k} \rightharpoonup u \quad \text{in } L^2(\Omega), \quad \lambda_k \rightarrow \lambda \quad \text{in } \mathbb{R}.$$

By (4.9), given $\varepsilon > 0$, there exists k_0 , such that $u_{\varepsilon_k} \in K_{\infty, \psi}^\varepsilon$ for all $k \geq k_0$. Then, since $K_{\infty, \psi}^\varepsilon$ is a closed convex set, we get $u \in K_{\infty, \psi}^\varepsilon$, and, by (4.10), we obtain that $u \in K_0$. Consequently,

$$(4.13) \quad w - \limsup_{n \rightarrow \infty} \text{Epi}(I_{K_{\infty, \psi}^\varepsilon}) \subset \text{Epi}(I_{K_\psi}).$$

Finally, by (4.12), (4.13), and having in mind (4.3), we obtain that $I_{K_{\infty, \psi}^\varepsilon}$ converges to I_{K_ψ} in the sense of Mosco. \square

4.4. Explicit solutions. Our goal now is to show some explicit examples that illustrate the behavior of the solutions when $p = +\infty$.

Remark 4.2. There is a natural upper bound (and of course also a natural lower bound) for the solutions with boundary datum ψ outside Ω (regardless the source term f). Indeed, given a bounded domain $\Omega \subset \mathbb{R}^N$ let us define inductively

$$\Omega_1 = \{x \in \Omega : |x - y| < 1 \text{ for some } y \in \Omega_J \setminus \overline{\Omega}\}$$

and, for $j \geq 2$,

$$\Omega_j = \left\{ x \in \Omega \setminus \bigcup_{i=1}^{j-1} \Omega_i : |x - y| < 1 \text{ for some } y \in \Omega_{j-1} \right\}.$$

Then, since $u(t) \in K_{\infty, \psi}^J$ we must have

$$u(t, x) \leq \psi(y) + 1 \quad \text{if } |x - y| \leq 1, x \in \Omega_1, y \in \Omega_J \setminus \overline{\Omega},$$

and for any $j \geq 2$

$$u(t, x) \leq u(t, y) + 1 \quad \text{if } |x - y| \leq 1, x \in \Omega_j, y \in \Omega_{j-1} \setminus \Omega_j.$$

Therefore we have an upper bound for $u(t, x)$ in the whole Ω ,

$$u(t, x) \leq \Psi_1(x),$$

where Ψ_1 is defined by the inductive formula,

$$\Psi_1(x) = \max \{ \psi(y) + 1 : y \in \Omega_J \setminus \overline{\Omega}, |x - y| \leq 1 \}, \text{ for } x \in \Omega_1,$$

and

$$\Psi_1(x) = \max \{ \Psi_1(y) + 1 : y \in \Omega_{j-1}, |x - y| \leq 1 \}, \text{ for } x \in \Omega_j, \text{ if } j \geq 2.$$

Analogously, we can obtain a lower bound for $u(t, x)$,

$$u(t, x) \geq \Phi_1(x),$$

where Φ_1 is defined by the inductive formula,

$$\Phi_1(x) = \min \{ \psi(y) - 1 : y \in \Omega_J \setminus \overline{\Omega}, |x - y| \leq 1 \}, \text{ for } x \in \Omega_1,$$

and

$$\Phi_1(x) = \min \{ \Phi_1(y) - 1 : y \in \Omega_{j-1}, |x - y| \leq 1 \}, \text{ for } x \in \Omega_j, \text{ if } j \geq 2.$$

With this remark in mind we show some explicit examples of solutions to

$$P_\infty^J(u_0, \psi, f) \begin{cases} f(t, x) - u_t(t, x) \in \partial G_{\infty, \psi}^J(u(t)), & \text{a.e. } t \in]0, T[, \\ u(0, x) = u_0(x), & \text{in } \Omega, \end{cases}$$

where

$$G_{\infty, \psi}^J(u) = \begin{cases} 0 & \text{if } u \in L^2(\Omega), |u(x) - u(y)| \leq 1, \text{ for } x, y \in \Omega, |x - y| \leq 1, \\ & \text{and } |u(x) - \psi(y)| \leq 1, \text{ for } x \in \Omega, y \in \Omega_J \setminus \overline{\Omega}, |x - y| \leq 1, \\ +\infty & \text{in the other case.} \end{cases}$$

In order to verify that a function $u(t, x)$ is a solution to $P_\infty^J(u_0, \psi, f)$, we need to check that

$$(4.14) \quad G_{\infty, \psi}^J(v) \geq G_{\infty, \psi}^J(u) + \langle f - u_t, v - u \rangle, \quad \text{for all } v \in L^2(\Omega).$$

To this end we can assume that $v \in K_{\infty, \psi}^J$ (otherwise $G_{\infty, \psi}^J(v) = +\infty$ and then (4.14) becomes trivial). Therefore, we need to check that

$$(4.15) \quad u(t, \cdot) \in K_{\infty, \psi}^J$$

and, by (4.14), that

$$(4.16) \quad \int_{\Omega} (f(t, x) - u_t(t, x))(v(x) - u(t, x)) dx \leq 0$$

for every $v \in K_{\infty, \psi}^J$.

Example 1. Let us consider a nonnegative source f and as initial condition the upper bound defined in the previous remark, $u_0(x) = \Psi_1(x)$. Then the solution to $P_\infty^J(u_0, \psi, f)$ is given by

$$u(t, x) \equiv \Psi_1(x)$$

for every $t > 0$. Indeed, $\Psi_1(x) \in K_{\infty, \psi}^J$ and for every $v \in K_{\infty, \psi}^J$ we have that $v(x) \leq \Psi_1(x)$, and therefore

$$\int_{\Omega} (f(t, x) - u_t(t, x))(v(x) - u(t, x)) \, dx = \int_{\Omega} f(t, x)(v(x) - \Psi_1(x)) \, dx \leq 0,$$

as we have to show.

In general, given a nonnegative source f supported in $D \subset \Omega$, any initial condition $u_0 \in K_{\infty, \psi}^J$ that verifies $u_0(x) = \Psi_1(x)$ in D produces a stationary solution $u(t, x) \equiv u_0(x)$.

Analogously, it can be shown that $u(t, x) \equiv \Phi_1(x)$ when $u_0(x) = \Phi_1(x)$ and $f(t, x) \leq 0$.

Example 2. Now, let us assume that we are in an interval $\Omega = (-L, L)$, $\psi = 0$, $\varepsilon = L/n$, $n \in \mathbf{N}$, $u_0 = 0$ which belongs to $K_{\varepsilon, 0}$, and the source f is an approximation of a delta function,

$$f(t, x) = f_{\eta}(t, x) = \frac{1}{\eta} \chi_{[-\frac{\eta}{2}, \frac{\eta}{2}]}(x), \quad 0 < \eta \leq 2\varepsilon.$$

Using the same ideas of [6], it is easy to verify the following general formula that describes the solution of $P_{\infty}^{\varepsilon}(u_0, \psi, f)$ for every $t \geq 0$. For any given integer $l \geq 0$ we have

$$u(t, x) = \begin{cases} l\varepsilon + k_l(t - t_l), & x \in [-\frac{\eta}{2}, \frac{\eta}{2}], \\ (l - 1)\varepsilon + k_l(t - t_l), & x \in [-\frac{\eta}{2} - \varepsilon, \frac{\eta}{2} + \varepsilon] \setminus [-\frac{\eta}{2}, \frac{\eta}{2}], \\ \dots \\ k_l(t - t_l), & x \in [-\frac{\eta}{2} - l\varepsilon, \frac{\eta}{2} + l\varepsilon] \setminus [-\frac{\eta}{2} - (l - 1)\varepsilon, \frac{\eta}{2} + (l - 1)\varepsilon], \\ 0, & x \notin [-\frac{\eta}{2} - l\varepsilon, \frac{\eta}{2} + l\varepsilon], \end{cases}$$

for $t \in [t_l, t_{l+1})$, where

$$k_l = \frac{1}{2l\varepsilon + \eta} \quad \text{and} \quad t_{l+1} = t_l + \frac{\varepsilon}{k_l}, \quad t_0 = 0.$$

This general formula is valid until the time at which the solution verifies $u(t, x) = \Psi_{\varepsilon}(x)$ for $x \in [-\frac{\eta}{2}, \frac{\eta}{2}]$ (the support of f), that is, until $T = t_{l^*+1}$, where

$$l^* \text{ is the first } l \text{ such that } l\varepsilon + k_l(t_{l+1} - t_l) = \Psi_{\varepsilon}(0)$$

and

Ψ_{ε} is the natural upper bound defined in Remark 4.2

for the corresponding rescaled kernel. Observe that for this l^* , $\frac{\eta}{2} + l^*\varepsilon \leq L$. From that time on the solution is stationary, that is, $u(t, x) = u(T, x)$ for all $t > T$.

From the above formula, taking limits as $\eta \rightarrow 0$, we get that the expected solution to $P_{\infty}^{\varepsilon}(u_0, \psi, \delta_0)$ is given, for any given integer $l \geq 1$, by

$$(4.17) \quad u(t, x) = \begin{cases} (l - 1)\varepsilon + k_l(t - t_l), & x \in [-\varepsilon, \varepsilon], \\ (l - 2)\varepsilon + k_l(t - t_l), & x \in [-2\varepsilon, 2\varepsilon] \setminus [-\varepsilon, \varepsilon], \\ \dots \\ k_l(t - t_l), & x \in [-l\varepsilon, l\varepsilon] \setminus [-(l - 1)\varepsilon, (l - 1)\varepsilon], \\ 0, & x \notin [-l\varepsilon, l\varepsilon], \end{cases}$$

for $t \in [t_l, t_{l+1})$, where $k_l = \frac{1}{2l\varepsilon}$, $t_{l+1} = t_l + \frac{\varepsilon}{k_l}$, $t_1 = 0$, until $T = t_{l^*+1}$, where

$$l^* \text{ is the first } l \text{ such that } l\varepsilon + k_l(t_{l+1} - t_l) = \Psi_\varepsilon(0).$$

And from that time on the solution is stationary, that is, $u(t, x) = u(T, x)$ for all $t > T$.

Remark that, since the space of functions $K_{\infty, \psi}^\varepsilon$ is not contained into $C(\mathbb{R})$, the formulation (4.16) with $f = \delta_0$ does not make sense. Hence the function $u(t, x)$ described by (4.17) is to be understood as a *generalized solution* to $P_\infty^\varepsilon(u_0, \psi, \delta_0)$ (it is obtained as a limit of solutions to approximating problems).

Note that the function $u(T, x)$ is a “regular and symmetric pyramid” composed by squares of side ε which is one step below the upper profile Ψ_ε .

Recovering the sandpile model as $\varepsilon \rightarrow 0$. Now, to recover the sandpile model, take the limit as $\varepsilon \rightarrow 0$ in the previous example to get that $u(t, x) \rightarrow v(t, x)$, where

$$v(t, x) = (l - |x|)^+ \quad \text{for } t = l^2,$$

until the time at which $t = L^2$, and from that time the solution is stationary.

A similar argument shows that, for any $a \in (0, L)$, the *generalized solution* to $P_\infty^\varepsilon(0, 0, \delta_a)$ converges as $\varepsilon \rightarrow 0$ to $v(t, x)$, where

$$v(t, x) = (l - |x - a|)^+ \quad \text{for } t = l^2,$$

until the time at which $t = (L - a)^2$, and from that time the solution is stationary.

These concrete examples illustrate the general convergence result in Theorem 1.8.

Acknowledgments. We want to thank the referees for their help to improve the manuscript. Part of this work was performed during a visit of JDR to Univ. de Valencia. He is thankful for the warm hospitality and the stimulating working atmosphere found there.

REFERENCES

- [1] F. ANDREU, C. BALLESTER, V. CASELLES, AND J. M. MAZÓN, *The Dirichlet problem for the total variation flow*, J. Funct. Anal., 180 (2001), pp. 347–403.
- [2] F. ANDREU, V. CASELLES, AND J. M. MAZÓN, *Parabolic quasilinear equations minimizing linear growth functionals*, Progress in Mathematics, Vol. 223, Birkhauser, 2004.
- [3] F. ANDREU, J. M. MAZÓN, AND J. TOLEDO, *Stabilization of solutions of the filtration equation with absorption and non-linear flux*, NoDEA, 2 (1995), pp. 267–289.
- [4] F. ANDREU, J. M. MAZÓN, J. D. ROSSI, AND J. TOLEDO, *The Neumann problem for nonlocal nonlinear diffusion equations*, J. Evol. Eqn., 8 (2008), pp. 189–215.
- [5] F. ANDREU, J. M. MAZÓN, J. D. ROSSI, AND J. TOLEDO, *A nonlocal p -Laplacian evolution equation with Neumann boundary conditions*, J. Math. Pures Appl., 90 (2008), pp. 201–227.
- [6] F. ANDREU, J. M. MAZÓN, J. D. ROSSI, AND J. TOLEDO, *The limit as $p \rightarrow \infty$ in a nonlocal p -Laplacian evolution equation. A nonlocal approximation of a model for sandpiles*, Calc. Var. Partial Differential Equations, to appear.
- [7] G. ANZELLOTTI, *Pairings between measures and bounded functions and compensated compactness*, Ann. Mat. Pura Appl. IV, 135 (1983), pp. 293–318.
- [8] H. ATTOUCH, *Familles d’opérateurs maximaux monotones et mesurabilité*, Ann. Mat. Pura Appl., 120 (1979), pp. 35–111.
- [9] G. BARLES AND C. IMBERT, *Second-order elliptic integro-differential equations: Viscosity solutions theory revisited*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 25 (2008), pp. 567–585.
- [10] J. W. BARRETT AND L. PRIGOZHIN, *Dual formulations in critical state problems*, Interfaces Free Bound., 8 (2006), pp. 349–370.

- [11] P. BATES AND A. CHMAJ, *A discrete convolution model for phase transitions*, Arch. Ration. Mech. Anal., 150 (1999), pp. 281–305.
- [12] P. BATES, P. FIFE, X. REN, AND X. WANG, *Travelling waves in a convolution model for phase transitions*, Arch. Ration. Mech. Anal., 138 (1997), pp. 105–136.
- [13] PH. BÉNILAN, L. BOCCARDO, T. GALLOUET, R. GARIÉPY, M. PIERRE, AND J. L. VAZQUEZ, *An L^1 -theory of existence and uniqueness of solutions of nonlinear elliptic equations*, Ann. Sc. Norm. Super. Pisa, IV, XXII (1995), pp. 241–273.
- [14] PH. BÉNILAN AND M. G. CRANDALL, *Completely accretive operators*, In Semigroup Theory and Evolution Equations (Delft, 1989), Lecture Notes in Pure and Appl. Math. 135, Dekker, New York, 1991, pp. 41–75.
- [15] PH. BÉNILAN, M. G. CRANDALL, AND A. PAZY, *Evolution Equations Governed by Accretive Operators*, book to appear.
- [16] H. BREZIS, *Équations et inéquations non linéaires dans les espaces vectoriels en dualité*, Ann. Inst. Fourier, 18 (1968), pp. 115–175.
- [17] H. BREZIS, *Opérateurs Maximaux Monotones et Semi-groupes de Contractions dans les Espaces de Hilbert*, North-Holland, Amsterdam, 1973.
- [18] H. BREZIS, *Analyse fonctionnelle. Théorie et applications*, Masson, 1978.
- [19] H. BREZIS AND A. PAZY, *Convergence and approximation of semigroups of nonlinear operators in Banach spaces*, J. Funct. Anal., 9 (1972), pp. 63–74.
- [20] L. CAFFARELLI AND L. SILVESTRE, *Regularity Theory for Fully Nonlinear Integro-Differential Equations*, preprint.
- [21] L. CAFFARELLI, S. SALSA, AND L. SILVESTRE, *Regularity estimates for the solution and the free boundary of the obstacle problem for the fractional Laplacian*, Inventiones Mathematicae, 171 (2008), pp. 425–461.
- [22] C. CARRILLO AND P. FIFE, *Spatial effects in discrete generation population models*, J. Math. Biol., 50 (2005), pp. 161–188.
- [23] E. CHASSEIGNE, M. CHAVES, AND J. D. ROSSI, *Asymptotic behaviour for nonlocal diffusion equations*, J. Math. Pures Appl., 86 (2006), pp. 271–291.
- [24] X. CHEN, *Existence, uniqueness and asymptotic stability of travelling waves in nonlocal evolution equations*, Adv. Differential Equations, 2 (1997), pp. 125–160.
- [25] C. CORTAZAR, M. ELGUETA, AND J. D. ROSSI, *A non-local diffusion equation whose solutions develop a free boundary*, Ann. Henri Poincaré, 6 (2005), pp. 269–281.
- [26] C. CORTAZAR, M. ELGUETA, J. D. ROSSI, AND N. WOLANSKI, *Boundary fluxes for non-local diffusion*, J. Differential Equations, 234 (2007), pp. 360–390.
- [27] C. CORTAZAR, M. ELGUETA, J. D. ROSSI, AND N. WOLANSKI, *How to approximate the heat equation with Neumann boundary conditions by nonlocal diffusion problems*, Arch. Ration. Mech. Anal., 187 (2008), pp. 137–156.
- [28] M. G. CRANDALL, *Nonlinear semigroups and evolution governed by accretive operators*, in Proceedings of the Symposium in Pure Mathematics, Part I, Vol. 45, F. Browder ed., AMS, Providence, RI, 1986, pp. 305–338.
- [29] S. DUMONT AND N. IGBIDA, *On a Dual Formulation for the Growing Sandpile Model*, European J. Appl. Math., to appear.
- [30] I. EKELAND AND R. TEMAM, *Convex Analysis and Variational Problems*, North-Holland, Amsterdam, 1972.
- [31] P. FIFE, *Some nonclassical trends in parabolic and parabolic-like evolutions*, Trends in nonlinear analysis, Springer, Berlin, 2003, pp. 153–191.
- [32] D. KINDERLEHRER AND G. STAMPACCHIA, *An Introduction to Variational Inequalities and their Applications*, Academic Press, New York, 1980.
- [33] S. KINDERMANN, S. OSHER, AND P. W. JONES, *Deblurring and denoising of images by nonlocal functionals*, Multiscale Model. Simul., 4 (2005), pp. 1091–1115.
- [34] U. MOSCO, *Convergence of convex sets and solutions of variational inequalities*, Adv. Math., 3 (1969), pp. 510–585.
- [35] L. PRIGOZHIN, *Variational models of sandpile growth*, Eur. J. Appl. Math., 7 (1996), pp. 225–236.
- [36] L. SILVESTRE, *Hölder estimates for solutions of integro differential equations like the fractional Laplace*, Indiana Univ. Math. J., 55 (2006), pp. 1155–1174.

GROW-UP RATE AND REFINED ASYMPTOTICS FOR A TWO-DIMENSIONAL PATLAK–KELLER–SEGEL MODEL IN A DISK*

NIKOS I. KAVALLARIS[†] AND PHILIPPE SOUPLET[‡]

Abstract. We consider a special case of the Patlak–Keller–Segel system in a disc, which arises in the modeling of chemotaxis phenomena. For a critical value of the total mass, the solutions are known to be global in time but with density becoming unbounded, leading to a phenomenon of mass-concentration in infinite time. We establish the precise grow-up rate and obtain refined asymptotic estimates of the solutions. Unlike in most of the similar, recently studied, grow-up problems, the rate is neither polynomial nor exponential. In fact, the maximum of the density behaves like $e^{\sqrt{2t}}$ for large time. In particular, our study provides a rigorous proof of a behavior suggested by Sire and Chavanis [*Phys. Rev. E* (3), 66 (2002), 046133] on the basis of formal arguments.

Key words. chemotaxis system, critical mass, grow-up, sub-/supersolutions

AMS subject classifications. Primary, 35Q, 35K60, 35B40, 92C17; Secondary, 35Q72

DOI. 10.1137/080722229

1. Introduction.

1.1. The complete Patlak–Keller–Segel model. Out of the many mathematical models that have been proposed to deal with particular aspects of chemotaxis, that proposed by Patlak in 1953 (cf. [40]) and Keller and Segel in 1970 (cf. [32]) has received particular attention. The so-called (two-dimensional) Patlak–Keller–Segel model consists of two equations, describing the evolution of the population density $\rho(x, t)$ of bacteria and the concentration $c(x, t)$ of a chemical attracting substance in a bounded domain $\Omega \subset \mathbb{R}^2$ and in a time interval $[0, T]$:

$$(1.1) \quad \frac{\partial \rho}{\partial t} = \nabla \cdot (D_1 \nabla \rho - \chi \rho \nabla c),$$

$$(1.2) \quad \theta \frac{\partial c}{\partial t} = \Delta c - ac + \rho.$$

More precisely, the first equation describes the random (Brownian) diffusion of the population of cells, which is biased in the direction of a drift velocity, proportional to the gradient of the concentration of the chemoattractant. The diffusion coefficient is denoted by $D_1 > 0$, and the proportionality coefficient of the drift (mobility parameter) is denoted by $\chi > 0$. According to the second equation, the chemoattractant, which is directly emitted by the cells, diffuses with a diffusion coefficient $D_2 = 1/\theta > 0$ on the substrate, while it is generated proportionally to the density of cells and at the same time is degraded with a rate equal to $a/\theta \geq 0$. In order for system (1.1)–(1.2) to be well posed, it should be supplemented with some initial conditions

$$(1.3) \quad \rho(x, 0) = \rho_0(x) \geq 0, \quad c(x, 0) = c_0(x) \geq 0,$$

*Received by the editors April 23, 2008; accepted for publication (in revised form) September 15, 2008; published electronically January 7, 2009.

<http://www.siam.org/journals/sima/40-5/72222.html>

[†]Department of Statistics and Actuarial-Financial Mathematics, University of the Aegean, Vourlioti Building, Gr-83200 Karlovassi, Samos, Greece (nkaval@aegean.gr).

[‡]Laboratoire Analyse Géométrie et Applications, UMR CNRS 7539, Institut Galilée, Université Paris-Nord, 99 av. J.-B. Clément, 93430 Villetaneuse, France (souple@math.univ-paris13.fr).

along with conditions on the boundary $\partial\Omega$. A natural boundary condition, since it guarantees the conservation of total mass, is the no-flux-type condition for ρ , namely,

$$(1.4) \quad \frac{\partial\rho}{\partial\nu} - \rho\frac{\partial c}{\partial\nu} = 0 \quad \text{on } \partial\Omega,$$

where ν stands for the outer unit normal vector at $\partial\Omega$. As for c , a Dirichlet-type boundary condition is assumed, i.e.,

$$(1.5) \quad c = 0 \quad \text{on } \partial\Omega;$$

cf. [7, 46]. Note that the parabolic system (1.1)–(1.2) preserves the nonnegativity of the initial conditions, i.e., $\rho, c \geq 0$ for $t > 0$, which is also expected to be true for the physical problem. For simplicity, D_1, χ are considered to be constant and under suitable scaling can be taken as $D_1 = \chi = 1$.

In view of experimental facts, the coefficients θ and a are assumed to be small, and a simplified form of the Patlak–Keller–Segel system is obtained (in fact, this corresponds to the case when the diffusion and production of c are much faster than the dynamics of ρ and the degradation of c). Namely, by considering the limiting case $\theta, a \rightarrow 0+$, the parabolic-parabolic system (1.1)–(1.5) is reduced to the elliptic-parabolic system

$$(1.6) \quad \frac{\partial\rho}{\partial t} = \nabla \cdot (\nabla\rho - \rho\nabla c), \quad x \in \Omega, \quad t \in (0, T),$$

$$(1.7) \quad -\Delta c = \rho, \quad x \in \Omega, \quad t \in (0, T),$$

$$(1.8) \quad \frac{\partial\rho}{\partial\nu} - \rho\frac{\partial c}{\partial\nu} = 0, \quad x \in \partial\Omega, \quad t \in (0, T),$$

$$(1.9) \quad c = 0, \quad x \in \partial\Omega, \quad t \in (0, T),$$

$$(1.10) \quad \rho(x, 0) = \rho_0(x) \geq 0, \quad x \in \Omega.$$

Note that in order for (1.6)–(1.10) to be well posed, only the initial data $\rho(x, 0) = \rho_0(x)$ must be prescribed. Moreover, owing to the boundary condition (1.8), the total (mass) population of cells is conserved; that is,

$$\|\rho(\cdot, t)\|_1 = \|\rho_0\|_1 =: \Lambda \quad \text{for } t > 0.$$

For a more detailed analysis regarding the modeling as well as the behavior of solutions of chemotaxis systems, see the review papers [24, 28, 29] as well as the monograph [48]. Here, it should be noticed that system (1.6)–(1.10) is also known as the Smoluchowski–Poisson system and can describe the motion of the mean field of many self-gravitating particles [13, 1, 2, 7, 46, 49, 50] or that of polymer molecules [16]. The behavior of the solution to (1.6)–(1.10) strongly depends on the parameter Λ . In fact, if $\Lambda > 8\pi$ and $\Omega = B(0, R)$, $R > 0$, then solutions of (1.6)–(1.10) blow up in a finite time $T^*(\rho_0)$; that is,

$$\lim_{t \rightarrow T^*} \|\rho(\cdot, t)\|_{H^1} = \lim_{t \rightarrow T^*} \|\rho(\cdot, t)\|_{L^p} = \lim_{t \rightarrow T^*} \int_{\Omega} (\rho \log \rho)(x, t) \, dx = \infty$$

for every $p > 1$; see [9, Theorem 2(i)], [2, Theorem 2]. On the other hand, for $\Lambda < 8\pi$, all solutions of system (1.6)–(1.10) are global in time; cf. [7, Theorem 2 (iv)]. In the critical case $\Lambda = 8\pi$ an infinite-time blow-up (grow-up) occurs, i.e., $\|\rho(\cdot, t)\|_{\infty} \rightarrow \infty$ as $t \rightarrow \infty$; cf. [39, Theorem 3], [5, Proposition 3.2]. Finite- or infinite-time blow-up can

be accompanied by the occurrence of a δ -function formation in the blow-up set (this represents the trend of populations to concentrate to form sporae) and is known in the literature as *chemotactic collapse*. This phenomenon was conjectured by Childress and Percus [14] and Nanjundiah [38] and was first verified, via matched asymptotics arguments, for a radially symmetric simplified Patlak–Keller–Segel system in [25]. A result regarding the infinite-time Dirac mass formation for ρ can be found in [39], where some characterization of grow-up (mass-concentration) points together with more grow-up results for different types of boundary conditions are also obtained [39, Theorems 2 and 3]. For blow-up results concerning a variation of system (1.6)–(1.10) but with Neumann boundary conditions for both ρ and c , see [30, 35, 3, 36, 44].

1.2. The simplified Patlak–Keller–Segel model in the radial case. In the case when Ω is the ball $B(0, R)$, $R > 0$, and the initial data $\rho_0(x) = \rho_0(r)$ is radially symmetric, the solution of system (1.6)–(1.10) is radially symmetric, i.e., $\rho(x, t) = \rho(r, t)$, with $r = |x|$. In this case the elliptic-parabolic system (1.6)–(1.10) can be greatly simplified. Namely, by introducing the cumulative mass distribution

$$Q(r, t) := \int_{B(0, r)} \rho(x, t) dx = 2\pi \int_0^r s\rho(s, t) ds,$$

which is equal to the mass contained in the sphere $B(0, r)$, the system reduces to a single equation

$$(1.11) \quad Q_t = Q_{rr} - \frac{1}{r}Q_r + \frac{1}{2\pi r}QQ_r, \quad 0 < r < R, \quad t > 0,$$

$$(1.12) \quad Q(0, t) = 0, \quad Q(R, t) = \Lambda.$$

By the definition of Q , the function

$$(1.13) \quad Q(r, 0) = Q_0(r)$$

is positive nondecreasing and satisfies the compatibility conditions $Q_0(0) = 0$ and $Q_0(R) = \Lambda$.

As mentioned in [4, 7], the formulation (1.11)–(1.13) allows the consideration of some initial data for the density ρ which could be either unbounded or singular (such as measures)—a case that seems rather realistic. This means that the initial data Q_0 for problem (1.11)–(1.13) could have unbounded derivatives $Q_{0,r}$ or even be discontinuous. Moreover, using formulation (1.11)–(1.13), we have the comparison principle at hand, which is not available for system (1.6)–(1.10). Due to the scaling properties of (1.11), we can assume without loss of generality that problem (1.11)–(1.13) is posed in the unit ball $B(0, 1)$. (Indeed, it is easily seen that if $Q(r, t)$ is a solution of (1.11)–(1.13), then $Q(Rr, R^2t)$ is also a solution.)

Using the new variable $x = r^2$ (no confusion with the original variable x in (1.6)–(1.10) should occur) and defining $N(x, t) = Q(r, t)$, we are led to the problem

$$(1.14) \quad N_t = 4xN_{xx} + \frac{1}{\pi}NN_x, \quad 0 < x < 1, \quad t > 0,$$

$$(1.15) \quad N(0, t) = 0, \quad N(1, t) = \Lambda,$$

$$(1.16) \quad N(x, 0) = N_0(x), \quad 0 < x < 1.$$

Note that (1.14) differs from the Burgers equation only by the variable coefficient x in the diffusion term. The above problem, although it degenerates at $x = 0$, may be

handled more easily than (1.11)–(1.13) since it contains fewer terms and at the same time does not have any singular coefficients in the first order terms.

As expected, the behavior of the solution of problem (1.14)–(1.16) (which is well defined for suitable initial data) depends on Λ . For $\Lambda > 8\pi$ the solution N ceases to exist in a finite time T^* ; more precisely, the boundary condition $N(0, t) = 0$ is no longer fulfilled at $t = T^*$. Moreover, a “gradient blow-up” occurs at $t = T^*$ in the sense that $N_x(0, t) \rightarrow \infty$ as $t \rightarrow T^*$ and the density ρ also becomes unbounded at time T^* ; cf. [8, Theorem 2(i)]. On the other hand, for $0 < \Lambda < 8\pi$ and any (admissible) initial data there is a unique global-in-time solution $N \in C([0, \infty); L^2(0, 1)) \cap C^{2,1}((0, 1) \times (0, \infty))$. Furthermore, N converges to the unique steady state solution:

$$(1.17) \quad N(\cdot, t) \rightarrow N_d = 8\pi \frac{x}{x+d} \quad \text{as } t \rightarrow \infty$$

in $L^p(\Omega)$, $p \geq 1$, and even in $L^\infty(\Omega)$ provided that $\sup_{t \geq 0} |N_x|_\infty < \infty$, where $d = \frac{8\pi}{\Lambda} - 1 > 0$; cf. [5]. In this case the rate of the L^1 -convergence of $N(\cdot, t)$ to N_d is shown to be exponential.

In the borderline case $\Lambda = 8\pi$, the situation is still different and the problem exhibits a typical critical behavior. Namely, it is proved in [5] that there exists a global-in-time solution N , which converges to the “singular” steady state $\tilde{N}(x) \equiv 8\pi$ (note that \tilde{N} does not satisfy the boundary condition at $x = 0$). Actually, as was proven in [39, Theorem 3], an infinite-time Dirac mass formation at the origin $r = 0$ of the ball occurs in this case. However, neither an estimation of the grow-up rate nor the asymptotic profile of the grow-up are provided in [39]. On the other hand, the authors in [5, Proposition 3.2] obtain the decay estimate

$$(1.18) \quad \|N(\cdot, t) - 8\pi\|_{L^1} \leq \frac{8\pi}{t} \quad \text{for } t \geq 1.$$

Estimate (1.18) seems to be far from optimal since formal asymptotics performed in [46] suggest a temporal decay estimate of order

$$(1.19) \quad \|N(\cdot, t) - 8\pi\|_{L^1} \approx O(e^{-\sqrt{2t}}) \quad \text{as } t \rightarrow \infty.$$

Remarks 1.1. (a) The (nonradial) Patlak–Keller–Segel system has also been studied in the whole plane \mathbb{R}^2 . Again, the behavior of solutions depends on the initial mass of the system, and a dichotomy is found [11, 17]. More precisely, assuming $0 \leq (1 + |x|^2)\rho_0$ and $\rho_0 \log \rho_0 \in L^1$, there exists a critical value of the mass $N_c := 8\pi$ such that if $0 < \|\rho_0\|_1 < N_c$ (subcritical case), only global-in-time solutions exist, while if $\|\rho_0\|_1 > N_c$ (supercritical case), the solutions blow up in finite time [11, 41]. Moreover, in the subcritical case, solutions converge to a self-similar profile as $t \rightarrow \infty$ [6, 11]. Finally, for the critical case $N = N_c$, which was studied in [10], the solution is global-in-time and grows up as a Dirac mass at the center of mass as $t \rightarrow \infty$.

(b) The only previous mathematical study of grow-up rates for a system of Patlak–Keller–Segel type concerns high dimensions, namely, $n \geq 11$, and was performed recently in [45]. There, some radial global unbounded solutions were constructed in a ball and an infinite sequence of polynomial grow-up rates was obtained (for a suitable sequence of initial data). On the contrary, our results in the present paper for $n = 2$ exhibit a grow-up rate independent of the initial data.

(c) Concerning the parabolic-parabolic Patlak–Keller–Segel system in a bounded domain with Neumann conditions, interesting results can be found in [26, 27] and [19,

37], respectively, on the asymptotics of finite-time blow-up and on global existence and the convergence of global bounded solutions.

(d) Problem (1.14)–(1.16) is in fact independent of the boundary condition (1.9) for c . Of course the boundary condition for c has to be taken into account if one wants to determine c , and not only ρ , from N .

(e) In the case where the diffusion of cells is very slow as compared to the diffusion of the chemoattractant, i.e., when $D_1 \ll D_2$ in (1.1)–(1.2), the complete system (1.1)–(1.5) is reduced to a single but nonlocal equation; cf. [51]. The global existence (subcritical case) and the finite-time blow-up (supercritical case) of solutions of the derived nonlocal equation in the two-dimensional case are studied in [31].

1.3. Heuristic description of the results and methods. Our aim is to prove rigorously the decay rate (1.19) as well as to provide a refined asymptotic profile for $N(x, t)$ as $t \rightarrow \infty$. To normalize the constants arising in the calculations we shall work with the equivalent problem

$$(1.20) \quad u_t = xu_{xx} + 2uu_x, \quad 0 < x < 1, \quad t > 0,$$

$$(1.21) \quad u(0, t) = 0, \quad u(1, t) = \xi := \frac{\Lambda}{8\pi}, \quad t > 0,$$

$$(1.22) \quad u(x, 0) = u_0(x), \quad 0 < x < 1,$$

which is obtained from (1.14)–(1.16) by setting $N(x, t) = 8\pi u(x, 4t)$.

It is clear that the stationary part $xu_{xx} + 2uu_x = 0$ of the parabolic equation (1.20) is invariant under the rescaling $u(x) \mapsto u(kx)$ ($k > 0$). Moreover, the steady state solution for $\xi < 1$ is given by

$$U_a(x) = 1 - \frac{1}{ax + 1} = \frac{ax}{ax + 1}$$

for $a = U'_a(0) = \frac{\xi}{1-\xi} > 0$. Rewriting (1.17) in terms of u , we have

$$u(x, t) \rightarrow U_a(x) \quad \text{as } t \rightarrow \infty.$$

Also, observe that $U_a(x)$ converges in a monotone increasing way to the “singular” steady state $U \equiv 1$ as $a \rightarrow \infty$ ($\xi \rightarrow 1$).

Motivated by the above considerations, we shall look for sub- and supersolutions of problem (1.20)–(1.22), which are perturbations of a moving family of steady states

$$(1.23) \quad U_{a(t)}(x),$$

the perturbation being defined in terms of the self-similar variable $y = a(t)x$. Here, a is a function of time, diverging to infinity, which is a priori unknown and will be eventually identified by a suitable “matching” procedure (see Remark 1.3). Note that such a form was used in [46] to construct “approximate solutions” leading to the formal asymptotics mentioned above. However, the expansions in [46] contained only a correction term at the first order, while those that we here construct involve first and second order terms (see Lemmas 4.1 and 4.2). This seems necessary to obtain (rigorous) sub- and supersolutions living “close” to the actual solution and eventually provides us with the desired decay rate as well as a good description of the asymptotic profile.

As a consequence of this construction, we shall derive an infinite-time boundary gradient grow-up result, with the grow-up rate

$$(1.24) \quad u_x(0, t) = A(t) \left(1 + O(t^{-1/2} \log t) \right) \quad \text{as } t \rightarrow \infty, \quad \text{where } A(t) = \exp \left[\frac{5}{2} + \sqrt{2t} \right].$$

In terms of the solution of the original system (1.6)–(1.10), (1.24) gives an estimate of the central density of bacteria, since $\rho(0, t) = 8u_x(0, 4t)$. Note that (1.24) was also predicted by formal arguments in [46]. At the same time, we show the C^1 regularity of u up to the boundary for all finite time intervals, hence ruling out the possibility of $\sup_{x \in [0, 1]} u_x(x, t)$ blowing up in finite time. This problem was left open in [5]. Moreover, we obtain a precise asymptotic expansion of the solution (see formula (2.8) in Theorem 2.1). It expresses the solution as the sum of a quasi-stationary profile (cf. (1.23)) and a correction term which becomes significant only for x bounded away from 0. As a consequence, we obtain the decay

$$\|u(\cdot, t) - 1\|_{L^1(0, 1)} = \left(1 + O(t^{-1/2} \log t) \right) \sqrt{2t} \exp \left[-\frac{5}{2} - \sqrt{2t} \right] \quad \text{as } t \rightarrow \infty,$$

again in accordance with the predictions in [46].

Remark 1.2. The study of unbounded global classical solutions of superlinear parabolic problems and their asymptotic behavior has recently attracted substantial mathematical interest. Particular effort has been devoted to the reaction-diffusion equation

$$(1.25) \quad u_t - \Delta u = u^p,$$

where such solutions are known to exist for suitable $p > 1$. Let us mention the works [33, 18, 22] for the Dirichlet problem in a ball and [42, 20, 21, 34] for the Cauchy problem. See also [43, sections 22 and 29] for related questions. The case of the Frank–Kamenetskii equation (with nonlinearity e^u instead of u^p) is also studied in [18]. As for the diffusive Hamilton–Jacobi equation

$$(1.26) \quad u_t - u_{xx} = |u_x|^p,$$

with $p > 2$, results of this kind can be found in [47]. A common feature in all these examples is the stabilization of the solution to a singular steady state, the growing-up quantity being $\|u(t)\|_\infty$ (resp., $\|u_x(t)\|_\infty$) for (1.25) (resp., (1.26)). It has to be noted that the grow-up rates, as $t \rightarrow \infty$, usually behave like either $e^{\mu t}$ or t^k . The only known exception (see [22]) seems to be the case of (1.25) with zero boundary conditions, in a ball of \mathbb{R}^4 with critical Sobolev exponent ($p = 3$). In this situation, unlike in spatial dimensions $n \neq 4$, there holds $\log \|u(t)\|_\infty \sim 2\sqrt{t}$, thus leading to a rate similar to that in our problem.

Remark 1.3. Let us point out that our determination of the grow-up rate for problem (1.20)–(1.22) is achieved through the matching of (sub- or super-) solutions with the imposed boundary condition at the right endpoint $x = 1$, an idea also present in [22]. This is different from what is done in [18, 20, 21, 47], where the grow-up rate is determined by the matching between inner and outer (sub-/super-) solutions. In those works, the inner solution corresponds to a self-similar, quasi-stationary evolution along a continuum of regular steady states (similar to (1.23)), but the outer solution is obtained by a linearization around the singular steady state. Here, on the contrary,

the behavior of u in the inner and outer regions is unified through a single self-similar variable $y = a(t)x$ (cf. formulae (4.1) and (4.27) below). Moreover, we point out that, in our case, a linearization around the “singular” steady state $U \equiv 1$ would not give the desired grow-up rate (1.24) but only a nonoptimal exponential upper bound based on an associated eigenvalue problem.

The paper is organized as follows. In section 2 we state the main results. Section 3 contains a number of preliminary results: In subsection 3.1, we recall the basic facts concerning local existence and comparison. In subsections 3.2 and 3.3, we show that a control on the slope at $x = 0$ is enough to prevent gradient blow-up (Lemma 3.2), and we obtain preliminary estimates of solutions for small time, which will be useful in section 4 to initialize the comparison with the main sub-/supersolutions (Lemmas 3.3 and 3.4). Subsection 3.4 is devoted to the study of a second order ordinary differential operator which plays a key role in the subsequent construction of sub-/supersolutions. The proofs of the main results are given in section 4: Subsections 4.1 and 4.2 are devoted to the construction of the main sub- and supersolutions, respectively; the proofs of Theorem 2.1 and Corollary 2.2 are finally completed in subsection 4.3.

2. Main results. Consider the problem

$$(2.1) \quad u_t - xu_{xx} = 2uu_x, \quad 0 < x < 1, \quad t > 0,$$

$$(2.2) \quad u(0, t) = 0, \quad t > 0,$$

$$(2.3) \quad u(1, t) = 1, \quad t > 0,$$

$$(2.4) \quad u(x, 0) = u_0(x), \quad 0 \leq x \leq 1.$$

Concerning the initial data, we assume that

$$(2.5) \quad u_0 \in C([0, 1]), \quad u_0(0) = 0, \quad u_0(1) = 1, \quad u_0 \text{ is nondecreasing,}$$

and that

$$(2.6) \quad u_0(x) \leq Kx, \quad 0 < x < 1, \quad \text{for some } K \geq 1.$$

Problem (2.1)–(2.4) admits a unique global solution, with $u \in C([0, 1] \times [0, \infty))$, $u \in C^{2,1}((0, 1] \times (0, \infty))$ (see [5] and section 3.1 below). Moreover, it was shown in [5] that $0 \leq u \leq 1$, u is nondecreasing in x , and

$$(2.7) \quad \lim_{t \rightarrow \infty} u(x, t) = 1, \quad \text{uniformly for } x \text{ in compact subsets of } (0, 1].$$

Our main results are the following.

THEOREM 2.1. *Let u_0 satisfy (2.5) and (2.6), and denote by u the global solution of problem (2.1)–(2.4). Then there holds*

$$(2.8) \quad 1 - u(x, t) = \frac{1 - x + O(t^{-1/2} \log t)}{1 + A(t)x} \quad \text{uniformly in } [0, 1], \text{ as } t \rightarrow \infty,$$

with

$$(2.9) \quad A(t) = \exp \left[\frac{5}{2} + \sqrt{2t} \right].$$

Moreover, we have the regularity property

$$(2.10) \quad u_x \in C([0, 1] \times (0, \infty))$$

and the estimate

$$(2.11) \quad u_x(0, t) = A(t) \left(1 + O(t^{-1/2} \log t) \right) \quad \text{as } t \rightarrow \infty.$$

COROLLARY 2.2. *Let u_0 satisfy (2.5) and (2.6), denote by u the global solution of problem (2.1)–(2.4), and let $A(t)$ be given by (2.9).*

(i) *The solution u satisfies the inner layer expansion (quasi-stationary behavior)*

$$1 - u(x, t) = \frac{1 + o(1)}{1 + A(t)x} \quad \text{uniformly in any region } x \leq o(1), \text{ as } t \rightarrow \infty.$$

(ii) *We have the L^1 -decay rate*

$$\|u(\cdot, t) - 1\|_{L^1(0,1)} = \left(1 + O(t^{-1/2} \log t) \right) \sqrt{2t} \exp \left[-\frac{5}{2} - \sqrt{2t} \right] \quad \text{as } t \rightarrow \infty.$$

3. Preliminaries.

3.1. Local existence and comparison principle. By [5, Theorem 2.1], we know that for any u_0 satisfying (2.5), problem (2.1)–(2.4) admits a global solution u in the following sense:

$$(3.1) \quad u \in C([0, \infty); L^1(0, 1)),$$

$$(3.2) \quad u \in C^{2,1}((0, 1] \times (0, \infty)),$$

$$(3.3) \quad u_x \geq 0, \quad 0 < x \leq 1, \quad t > 0,$$

and

$$(3.4) \quad u_t - xu_{xx} = 2uu_x, \quad 0 < x < 1, \quad t > 0,$$

$$(3.5) \quad \lim_{x \rightarrow 0^+} u(x, t) = 0 \quad \text{for a.e. } t > 0,$$

$$(3.6) \quad u(1, t) = 1, \quad t > 0,$$

$$(3.7) \quad u(\cdot, 0) = u_0 \quad \text{in } L^1(0, 1).$$

Note that the solution is obtained in [5] as a limit of solutions of regularized problems, where xu_{xx} is replaced by $(x + \varepsilon)u_{xx}$, $\varepsilon > 0$. Also, for each $T > 0$, u is the unique local solution of (3.1)–(3.7) on $(0, T)$. If, moreover, u_0 satisfies (2.6), then

$$u \in C([0, 1] \times [0, \infty));$$

see also [8, Theorem 1 (i)]. The continuity for $t > 0$ and $x = 0$ follows from [5, Propositions 2.4 and 2.5]. For $t = 0$ and $x \in [0, 1]$, the continuity can be established by comparison with simple barrier functions.

The following proposition provides a comparison principle suitable to our needs.

PROPOSITION 3.1. *Let $\tau > 0$ and the functions u, v satisfy the following regularity conditions:*

$$(3.8) \quad u, v \in C([0, \tau]; L^1(0, 1)),$$

$$(3.9) \quad u, v \in C^{2,1}((0, 1] \times (0, \tau)),$$

$$(3.10) \quad u_x, v_x \in L^1_{loc}([0, 1] \times (0, \tau)),$$

$$(3.11) \quad u, v \in L^\infty_{loc}([0, 1] \times (0, \tau)).$$

Assume that

$$(3.12) \quad u_t - xu_{xx} - 2uu_x \leq v_t - xv_{xx} - 2vv_x, \quad 0 < x < 1, \quad 0 < t < \tau,$$

$$(3.13) \quad \lim_{x \rightarrow 0^+} u(x, t) \leq \lim_{x \rightarrow 0^+} v(x, t) \quad \text{for a.e. } t \in (0, \tau),$$

$$(3.14) \quad u(1, t) \leq v(1, t), \quad 0 < t < \tau,$$

$$(3.15) \quad u(\cdot, 0) \leq v(\cdot, 0) \quad \text{a.e. in } (0, 1).$$

Then $u \leq v$ in $(0, 1) \times (0, \tau)$.

Observe that (3.10) implies that the limits in (3.13) exist for a.e. $t \in (0, \tau)$. Note also that for any u_0 satisfying (2.5), the solution u of problem (3.1)–(3.7) satisfies conditions (3.10) and (3.11) as a consequence of (3.3), (3.5), and (3.6).

Proof of Proposition 3.1. The proof is a modification of the stability proof in [5, Theorem 3.1]. Let $z = u - v$. By (3.12), we have

$$(3.16) \quad z_t \leq xz_{xx} + 2uu_x - 2vv_x = \partial_x(xz_x + z(u + v - 1)), \quad 0 < x < 1, \quad 0 < t < \tau.$$

For $\delta \in (0, 1)$ we define the following C^1 (and piecewise C^2) convex approximations of the function $s \rightarrow s_+ = \max(s, 0)$:

$$\phi_\delta(s) = \begin{cases} 0 & \text{if } -\infty < s \leq \delta, \\ (2\delta)^{-1}(s - \delta)^2 & \text{if } \delta \leq s \leq 2\delta, \\ s - 3\delta/2 & \text{if } 2\delta < s < \infty. \end{cases}$$

Fix $0 < t_1 < t_2 < \tau$ and $\delta \in (0, 1)$. Then, for any $\varepsilon \in (0, 1)$, multiplying (3.16) by $\phi'_\delta(z)$, integrating by parts, and using (3.11), (3.14), $0 \leq \phi'_\delta \leq 1$, and $\phi''_\delta \geq 0$, we obtain

$$\begin{aligned} & \int_\varepsilon^1 \phi_\delta(z(x, t_2)) \, dx - \int_\varepsilon^1 \phi_\delta(z(x, t_1)) \, dx \\ &= \int_{t_1}^{t_2} \int_\varepsilon^1 \phi'_\delta(z) z_t \, dx \, dt \leq \int_{t_1}^{t_2} \int_\varepsilon^1 \phi'_\delta(z) \partial_x(xz_x + z(u + v - 1)) \, dx \, dt \\ &= \int_{t_1}^{t_2} \left[\phi'_\delta(z)(xz_x + z(u + v - 1)) \right]_\varepsilon^1 \, dt - \int_{t_1}^{t_2} \int_\varepsilon^1 \phi''_\delta(z)(xz_x + z(u + v - 1))z_x \, dx \, dt \\ &\leq \varepsilon \int_{t_1}^{t_2} |z_x(\varepsilon, t)| \, dt + C \int_{t_1}^{t_2} \phi'_\delta(z(\varepsilon, t)) \, dt - \int_{t_1}^{t_2} \int_\varepsilon^1 \phi''_\delta(z)(u + v - 1)zz_x \, dx \, dt \\ &\equiv I_\varepsilon + J_\varepsilon + K_\varepsilon, \end{aligned}$$

where C depends on t_1, t_2 but is independent of ε (and δ). Owing to (3.10), there exists a sequence $\varepsilon_n \rightarrow 0+$ such that $\lim_{n \rightarrow \infty} I_{\varepsilon_n} = 0$. Next, since $\lim_{n \rightarrow \infty} \phi'_\delta(z(\varepsilon_n, t)) = 0$ for a.e. $t \in (t_1, t_2)$ due to (3.13) and the definition of ϕ_δ , we deduce that $\lim_{n \rightarrow \infty} J_{\varepsilon_n} = 0$ by dominated convergence. Consequently,

$$\begin{aligned} \int_0^1 \phi_\delta(z(x, t_2)) \, dx - \int_0^1 \phi_\delta(z(x, t_1)) \, dx &\leq - \int_{t_1}^{t_2} \int_0^1 \phi''_\delta(z)(u + v - 1)zz_x \, dx \, dt \\ &\leq C \int_{t_1}^{t_2} \int_0^1 |z\phi''_\delta(z)||z_x| \, dx \, dt. \end{aligned}$$

Now, observe that $\lim_{\delta \rightarrow 0} \phi_\delta(s) = s_+$ and $\lim_{\delta \rightarrow 0} s\phi''_\delta(s) = 0$ for each $s \in \mathbb{R}$. Using $0 \leq \phi_\delta(s) \leq s_+$, (3.10), and (3.11), we may pass to the limit $\delta \rightarrow 0$ by dominated

convergence in the preceding inequality, and we obtain

$$\int_0^1 z_+(x, t_2) dx - \int_0^1 z_+(x, t_1) dx \leq 0.$$

Letting $t_1 \rightarrow 0+$ and using (3.8) and (3.15), we conclude that $\int_0^1 z_+(x, t_2) dx = 0$ for all $t_2 \in (0, \tau)$; hence $u \leq v$ in $(0, 1) \times (0, \tau)$. \square

3.2. Sufficient condition for C^1 regularity. As noted in, e.g., [12, section 2.2] and [23, section 2], by means of the transformation

$$w(r, t) := \frac{8}{r^2} u(r^2, 4t) = \frac{1}{\pi r^2} \int_{B_r} \rho(y, 4t) dy$$

(and $w_0(r) := 8r^{-2}u_0(r^2)$), problem (2.1)–(2.4) becomes equivalent to

$$(3.17) \quad w_t - \tilde{\Delta}w = w^2 + \frac{r}{2} w w_r, \quad 0 < r < 1, \quad t > 0,$$

$$(3.18) \quad w_r(0, t) = 0, \quad t > 0,$$

$$(3.19) \quad w(1, t) = 8, \quad t > 0,$$

$$(3.20) \quad w(r, 0) = w_0(r), \quad 0 < r < 1,$$

where $\tilde{\Delta}w := w_{rr} + \frac{3}{r}w_r$, which in turn corresponds to the radial Laplacian in four space dimensions. It should be noticed that w has the same scale invariance as ρ but is smoother than ρ . Problem (3.17)–(3.20) turns out to be convenient regarding the study of C^1 regularity of u up to the boundary, which was left open in [5]. Namely, using this transformation, one can show the following additional properties for u .

LEMMA 3.2. *Let u_0 satisfy (2.5) and (2.6), and denote by u the global solution of problem (2.1)–(2.4). For any given $T > 0$, if*

$$(3.21) \quad \sup_{0 < x < 1, 0 < t < T} \frac{u(x, t)}{x} < \infty,$$

then

$$(3.22) \quad u_x \in C([0, 1] \times (0, T])$$

and

$$(3.23) \quad u_x(0, t) > 0, \quad 0 < t \leq T.$$

Note that the assumption (3.21) (for all $T > 0$) will be shown in section 4 (see Lemma 4.3(ii)) as a consequence of our main supersolution construction, hence leading to the global C^1 regularity property (2.10) in Theorem 2.1.

Proof of Lemma 3.2. Since the semilinear parabolic equation in (3.17) has only linear growth with respect to the gradient, standard arguments based on the variation-of-constants formula (see, e.g., [43, Example 51.30] or [15, p. 889]) show that problem (3.17)–(3.20) is locally well posed in the space of (radial, nonnegative) L^∞ -functions. More precisely, for any $0 \leq w_0 \in L^\infty(0, 1)$, there exists a unique, maximal (radial, nonnegative) classical solution w of (3.17)–(3.20), with $w \in C^{2,1}([0, 1] \times (0, T_m))$ and $w \in C([0, T_m]; L^q(0, 1))$ for all finite $q \leq 1$. Moreover,

$$(3.24) \quad T_m < \infty \implies \lim_{t \rightarrow T_m} \|w(t)\|_\infty = \infty.$$

Now, if u_0 satisfies (2.5)–(2.6), then $w_0(r) := 8r^{-2}u_0(r^2)$ verifies $0 \leq w_0 \in L^\infty(0, 1)$. Denote by w the corresponding maximal solution of (3.17)–(3.20), and let

$$\tilde{u}(x, t) = \frac{x}{8}w(\sqrt{x}, t/4), \quad 0 \leq x \leq 1, \quad 0 \leq t < 4T_m.$$

We see that \tilde{u} satisfies the regularity conditions in (3.8), (3.9), and (3.11) with $\tau = 4T_m$. Also, since

$$\tilde{u}_x(x, t) = \frac{1}{8}w(\sqrt{x}, t/4) + \frac{\sqrt{x}}{16}w_r(\sqrt{x}, t/4), \quad 0 < x \leq 1, \quad 0 < t < 4T_m,$$

we have

$$(3.25) \quad \tilde{u} \in C^{1,0}([0, 1] \times (0, 4T_m))$$

and hence in particular (3.10) with $\tau = 4T_m$. Then one easily checks that \tilde{u} solves (2.1)–(2.3) on $[0, 4T_m]$, along with (3.7). We may thus apply Proposition 3.1 to deduce that $\tilde{u} = u$ on $(0, T_0)$ with $T_0 = \min(T, 4T_m)$.

We claim that $T_m > T/4$. Indeed, if $T_m \leq T/4$, then (3.24) implies

$$\lim_{t \rightarrow 4T_m} \sup_{0 < x < 1} \frac{u(x, t)}{x} = \infty,$$

contradicting (3.21). Consequently, property (3.22) follows from (3.25). Moreover, $w > 0$ in $[0, 1] \times (0, T_m)$ by the strong maximum principle, which readily implies (3.23). \square

3.3. Small time estimates. Throughout the paper, we denote by \mathcal{P} the parabolic operator defined by

$$(3.26) \quad \mathcal{P}v := v_t - xv_{xx} - 2vv_x.$$

The following two lemmas will be useful to initialize the comparison between u and the main sub-/supersolutions constructed in section 4.

LEMMA 3.3. *Let u_0 satisfy (2.5) and (2.6), and denote by u the global solution of problem (2.1)–(2.4). Then there exist $\tau, \eta > 0$ such that*

$$(3.27) \quad u(x, t) \leq 2Kx, \quad 0 \leq x \leq 1, \quad 0 \leq t \leq \tau,$$

and

$$u(x, \tau) \leq 1 - \eta(1 - x), \quad 0 \leq x \leq 1.$$

Proof. Define

$$\bar{v}(x, t) = \frac{Kx}{1 - 2Kt}, \quad 0 \leq x \leq 1, \quad 0 \leq t < 1/2K.$$

Since

$$\mathcal{P}\bar{v} = \frac{2K^2x}{(1 - 2Kt)^2} - \frac{2K^2x}{(1 - 2Kt)^2} = 0$$

and $\bar{v}(1, t) \geq K \geq 1$, the comparison principle guarantees that

$$u(x, t) \leq \bar{v}(x, t) \leq 2Kx, \quad 0 \leq x \leq 1, \quad 0 < t \leq 1/4K.$$

Fix $\tau = 1/4K$. By Hopf’s lemma, we have $u_x(1, \tau) > 0$. Since $u(\cdot, \tau)$ is nondecreasing in x , this implies $u(x, \tau) \leq 1 - \eta(1 - x)$, $0 \leq x \leq 1$, for $\eta > 0$ sufficiently small, and the conclusion follows. \square

LEMMA 3.4. *Let u_0 satisfy (2.5) and (2.6), and denote by u the global solution of problem (2.1)–(2.4). For any given $\delta \in (0, 1)$, there exists $T_\delta > 0$ such that*

$$u(x, T_\delta) \geq \min(1 - \delta, x/\delta), \quad 0 \leq x \leq 1.$$

Proof. By Lemma 3.2, we may fix a small $\tau > 0$ such that $u_x(\cdot, \tau) \in C([0, 1])$ and $u_x(0, \tau) > 0$. Therefore, $u(x, \tau) \geq \eta x$ for all $x \in [0, 1]$ and some $\eta \in (0, 1)$. Since $\mathcal{P}[\eta x] \leq 0$ and $\eta < 1 = u(1, t)$, the comparison principle implies that

$$(3.28) \quad u(x, t) \geq \eta x, \quad 0 \leq x \leq 1, t \geq \tau.$$

On the other hand, by (2.7), there exists $T > \tau$ such that

$$(3.29) \quad u(x, t) \geq 1 - \delta, \quad (1 - \delta)\delta \leq x \leq 1, t \geq T.$$

Define

$$(3.30) \quad \underline{v}(x, t) = (\eta + 2\eta^2(t - T))x, \quad 0 \leq x \leq 1, t \geq T.$$

We have

$$(3.31) \quad \mathcal{P}\underline{v} = 2\eta^2x - 2(\eta + 2\eta^2(t - T))^2x \leq 0$$

and, due to (3.28),

$$(3.32) \quad u(x, T) \geq \underline{v}(x, T), \quad 0 \leq x \leq 1.$$

Since $\eta < 1 < 1/\delta$, we may find $T_\delta > T$ such that

$$(3.33) \quad \eta + 2\eta^2(T_\delta - T) = 1/\delta.$$

In view of (3.29), (3.30), and (3.33), we have

$$(3.34) \quad u((1 - \delta)\delta, t) \geq 1 - \delta \geq \underline{v}((1 - \delta)\delta, t), \quad T \leq t \leq T_\delta.$$

It then follows from (3.31), (3.32), (3.34), and the comparison principle that

$$u(x, T_\delta) \geq x/\delta, \quad 0 \leq x \leq (1 - \delta)\delta.$$

Using $\mathcal{P}[x/\delta] \leq 0$ and (3.29), we deduce that

$$u(x, t) \geq x/\delta, \quad 0 \leq x \leq (1 - \delta)\delta, t \geq T_\delta.$$

This combined with (3.29) yields the desired conclusion. \square

3.4. ODE lemmas. We consider the differential operator

$$\mathcal{L}w := yw'' + \frac{2yw'}{1 + y} + \frac{2w}{(1 + y)^2}.$$

First, the expression for the operator \mathcal{L} reads

$$\mathcal{L}w = \left[\left(\frac{y}{y + 1} \right)^2 \left(\frac{(y + 1)^2}{y} w \right)' \right]'$$

If $\psi \in C([0, \infty))$ satisfies $\psi(y) = O(y)$ as $y \rightarrow 0$, then the problem

$$\begin{aligned} \mathcal{L}w &= \psi, \quad y > 0, \\ w(0) &= 0, \quad w'(0) = 0 \end{aligned}$$

admits a unique solution w , and w can be represented as

$$w(y) = \mathcal{L}_0^{-1}\psi := \frac{y}{(y+1)^2} \int_0^y \left(\frac{t+1}{t}\right)^2 \int_0^t \psi(s) ds dt$$

(note that the integral is convergent due to the assumption on ψ). In particular,

$$(3.35) \quad \psi \geq 0 \text{ on } [0, \infty) \implies \mathcal{L}_0^{-1}\psi \geq 0 \text{ on } [0, \infty).$$

On the other hand, $w_0(y) = \frac{y}{(y+1)^2}$ solves

$$\begin{aligned} \mathcal{L}w &= 0, \quad y > 0, \\ w(0) &= 0, \quad w'(0) = 1. \end{aligned}$$

In the next section, to construct our main sub- and supersolutions, we will need to know the asymptotic behavior of the action of the operator \mathcal{L}_0 on some particular functions as $y \rightarrow \infty$. More precisely, let us define

$$\begin{aligned} f &= (I + \mathcal{L}_0^{-1}) \left(\frac{y}{(y+1)^2} \right), \\ \tilde{f} &= 2ff' - yf' + f, \\ g &= \mathcal{L}_0^{-1}\tilde{f}, \end{aligned}$$

and

$$h = \mathcal{L}_0^{-1}(\tilde{f} + M\varphi) = g + M\mathcal{L}_0^{-1}\varphi,$$

where $M > 0$ and

$$\varphi \in C^1([0, \infty)), \quad \varphi(0) = 0, \quad \varphi(y) = \frac{1}{\log y}, \quad y \geq 2.$$

We have the following lemmas.

LEMMA 3.5. *As $y \rightarrow \infty$, the function f satisfies*

$$\begin{aligned} \text{(i)} \quad f(y) &= \log y - 2 + O\left(\frac{\log^2 y}{y}\right); \\ \text{(ii)} \quad f'(y) &= \frac{1}{y} + O\left(\frac{\log^2 y}{y^2}\right). \end{aligned}$$

LEMMA 3.6. *As $y \rightarrow \infty$, the functions g and h satisfy*

$$\begin{aligned} \text{(i)} \quad g(y) &= \frac{y \log y}{2} - \frac{9y}{4} + O(\log^3 y); \\ \text{(ii)} \quad g'(y) &= \frac{\log y}{2} - \frac{7}{4} + O\left(\frac{\log^3 y}{y}\right); \\ \text{(iii)} \quad h(y) &= \frac{y \log y}{2} - \frac{9y}{4} + O\left(\frac{y}{\log y}\right); \\ \text{(iv)} \quad h'(y) &= \frac{\log y}{2} - \frac{7}{4} + O\left(\frac{1}{\log y}\right). \end{aligned}$$

Proof of Lemma 3.5.

(i) Using

$$(3.36) \quad \log(y + 1) = \log y + O(1/y) \quad \text{as } y \rightarrow \infty,$$

we obtain

$$\begin{aligned} f(y) &= \frac{y}{(y + 1)^2} \left[1 + \int_0^y \left(\frac{t + 1}{t} \right)^2 \int_0^t \left(\frac{1}{s + 1} - \frac{1}{(s + 1)^2} \right) ds dt \right] \\ &= \frac{y}{(1 + y)^2} \left[1 + \int_0^y \left\{ \left(1 + \frac{2}{t} + \frac{1}{t^2} \right) \left(\log(t + 1) - \frac{t}{t + 1} \right) \right\} dt \right] \\ &= \frac{y}{(1 + y)^2} [y \log y - 2y + O(\log^2 y)] \\ &= \log y - 2 + O\left(\frac{\log^2 y}{y}\right) \quad \text{as } y \rightarrow \infty. \end{aligned}$$

(ii) Now, using (i), we deduce

$$\begin{aligned} f'(y) &= \left[\frac{1}{y} - \frac{2}{y + 1} \right] f(y) + \frac{1}{y} \left(\log(y + 1) - \frac{y}{y + 1} \right) \\ &= \left(-\frac{1}{y} + O\left(\frac{1}{y^2}\right) \right) \left(\log y - 2 + O\left(\frac{\log^2 y}{y}\right) \right) + \frac{\log(y + 1)}{y} - \frac{1}{y + 1} \\ &= \frac{1}{y} + O\left(\frac{\log^2 y}{y^2}\right) \quad \text{as } y \rightarrow \infty. \quad \square \end{aligned}$$

To show Lemma 3.6, we first note that, due to Lemma 3.5 and (3.36), we have

$$\tilde{f}(y) = f_1(y) - 3f_2(y) + O(f_3(y)) \quad \text{as } y \rightarrow \infty,$$

where

$$f_1(y) = \log(y + 1), \quad f_3(y) = \frac{\log^2(y + 1)}{y + 1},$$

and

$$f_2 \in C^1([0, \infty)), \quad f_2(0) = 0, \quad f_2(y) = 1, \quad y \geq 1.$$

Denote $g_i = \mathcal{L}_0^{-1} f_i$ for $i = 1, 2, 3$ and $g_4 = \mathcal{L}_0^{-1} \varphi$. Lemma 3.6 will then be an immediate consequence of the following lemma.

LEMMA 3.7. *As $y \rightarrow \infty$, the functions g_i satisfy*

- (i) $g_1(y) = \frac{y \log y}{2} - \frac{3y}{4} + O(\log y)$ and $g'_1(y) = \frac{\log y}{2} - \frac{1}{4} + O\left(\frac{\log y}{y}\right)$;
- (ii) $g_2(y) = \frac{y}{2} + O(1)$ and $g'_2(y) = \frac{1}{2} + O\left(\frac{1}{y}\right)$;
- (iii) $g_3(y) = O(\log^3 y)$ and $g'_3(y) = O\left(\frac{\log^3 y}{y}\right)$;
- (iv) $g_4(y) = O\left(\frac{y}{\log y}\right)$ and $g'_4(y) = O\left(\frac{1}{\log y}\right)$.

Proof.

(i) As $y \rightarrow \infty$, we have

$$\begin{aligned} g_1(y) &= \frac{y}{(y+1)^2} \int_0^y \left(\frac{t+1}{t}\right)^2 ((t+1)\log(t+1) - t) dt \\ &= \frac{y}{(1+y)^2} \int_0^y [(t+1)\log(t+1) - t + O(\log(t+1))] dt \\ &= \left[\frac{1}{y} + O\left(\frac{1}{y^2}\right)\right] \left[\frac{y^2 \log y}{2} - \frac{3y^2}{4} + O(y \log y)\right] \\ &= \frac{y \log y}{2} - \frac{3y}{4} + O(\log y); \end{aligned}$$

hence

$$\begin{aligned} g'_1(y) &= \left[\frac{1}{y} - \frac{2}{y+1}\right] g_1(y) + \frac{1}{y} \int_0^y \log(s+1) ds \\ &= \left(-\frac{1}{y} + O\left(\frac{1}{y^2}\right)\right) \left(\frac{y \log y}{2} - \frac{3y}{4} + O(\log y)\right) + \frac{(y+1)\log(y+1) - y}{y} \\ &= \frac{\log y}{2} - \frac{1}{4} + O\left(\frac{\log y}{y}\right). \end{aligned}$$

(ii) Due to the definition of g_2 , there exist constants $C_1, C_2 \in \mathbb{R}$ such that, for $y > 1$,

$$\begin{aligned} g_2(y) &= \frac{y}{(y+1)^2} \left[\int_0^1 \left(\frac{t+1}{t}\right)^2 \int_0^t f_2(s) ds dt + \int_1^y \left(\frac{t+1}{t}\right)^2 (C_1 + t) dt \right] \\ &= \frac{y}{(y+1)^2} \left[C_2 + \frac{y^2}{2} + O(y) \right] = \frac{y}{2} + O(1) \quad \text{as } y \rightarrow \infty. \end{aligned}$$

Therefore, for $y > 1$,

$$\begin{aligned} g'_2(y) &= \left[\frac{1}{y} - \frac{2}{y+1}\right] g_2(y) + \frac{y}{(y+1)^2} \left(\frac{y+1}{y}\right)^2 (C_1 + y) \\ &= \left[\frac{1}{y} - \frac{2}{y+1}\right] g_2(y) + 1 + O\left(\frac{1}{y}\right) = \frac{1}{2} + O\left(\frac{1}{y}\right) \quad \text{as } y \rightarrow \infty. \end{aligned}$$

(iii) As $y \rightarrow \infty$, we have

$$\begin{aligned} g_3(y) &= \frac{y}{3(y+1)^2} \int_0^y \left(\frac{t+1}{t}\right)^2 \log^3(t+1) dt \\ &= \frac{y}{3(1+y)^2} \int_0^y O(\log^3(t+1)) dt \\ &= \frac{y}{3(1+y)^2} O((y+1)\log^3(y+1)) = O(\log^3 y); \end{aligned}$$

hence

$$\begin{aligned} g'_3(y) &= \left[\frac{1}{y} - \frac{2}{y+1}\right] g_3(y) + \frac{1}{y} \int_0^y \frac{\log^2(t+1)}{(t+1)} dt \\ &= \left(-\frac{1}{y} + O\left(\frac{1}{y^2}\right)\right) O(\log^3 y) + \frac{\log^3(y+1)}{3y} = O\left(\frac{\log^3 y}{y}\right). \end{aligned}$$

(iv) Similarly to (ii), in this case we have, for $y > 2$,

$$\begin{aligned} g_4(y) &= \frac{y}{(y+1)^2} \left[\int_0^2 \left(\frac{t+1}{t}\right)^2 \int_0^t \varphi(s) ds dt + \int_2^y \left(\frac{t+1}{t}\right)^2 \left(C_1 + \int_2^t \frac{ds}{\log s}\right) dt \right] \\ &= \frac{y}{(y+1)^2} \left[C_2 + \int_2^y \left(\frac{t+1}{t}\right)^2 O\left(\frac{t}{\log t}\right) dt \right] \\ &= \frac{y}{(y+1)^2} O\left(\frac{y^2}{\log y}\right) = O\left(\frac{y}{\log y}\right) \quad \text{as } y \rightarrow \infty. \end{aligned}$$

Therefore, for $y > 2$,

$$\begin{aligned} g'_4(y) &= \left[\frac{1}{y} - \frac{2}{y+1} \right] g_4(y) + \frac{y}{(y+1)^2} \left(\frac{y+1}{y}\right)^2 \left(C_1 + \int_2^y \frac{ds}{\log s}\right) \\ &= O\left(\frac{1}{\log y}\right) \quad \text{as } y \rightarrow \infty. \quad \square \end{aligned}$$

4. Proof of the main results.

4.1. Construction of a subsolution. Motivated by the idea of an asymptotic expansion around (moving) steady states and based on a self-similar variable (see subsection 1.3), we make the following ansatz:

$$\underline{u}(x, t) = 1 - \frac{1}{y+1} + b(t)f(y) - b^2(t)g(y), \quad y = a(t)x,$$

where the functions a, b, f, g have to be determined. Note that the variable y now ranges into the time-dependent interval $[0, a(t)]$. Here, a and b are expected to satisfy

$$a(t) \sim u_x(0, t), \quad \lim_{t \rightarrow \infty} a(t) = \infty, \quad \text{and} \quad \lim_{t \rightarrow \infty} b(t) = 0.$$

LEMMA 4.1. *The problem*

$$\begin{aligned} \underline{u}_t - x\underline{u}_{xx} &\leq 2\underline{u}\underline{u}_x, \quad 0 < x < 1, \quad t > 0, \\ \underline{u}(0, t) &= 0, \quad t \geq 0, \\ \underline{u}(1, t) &< 1, \quad t \geq 0, \end{aligned}$$

admits a solution of the form

$$(4.1) \quad \underline{u}(x, t) = 1 - \frac{1}{y+1} + b(t)f(y) - b^2(t)g(y), \quad y = a(t)x,$$

where the smooth functions $a > 0, b > 0, f \geq 0$, and g have the following properties:

$$f(y) \sim \log y, \quad g(y) \sim \frac{y \log y}{2} \quad \text{as } y \rightarrow \infty,$$

$$(4.2) \quad f(0) = g(0) = f'(0) = g'(0) = 0,$$

$$(4.3) \quad a(t) = \left(1 + O(t^{-1/2} \log t)\right) \exp\left[\frac{5}{2} + \sqrt{2t}\right] \quad \text{as } t \rightarrow \infty,$$

$$(4.4) \quad b(t) = \frac{1 + O(t^{-1/2})}{a(t) \log a(t)} \quad \text{as } t \rightarrow \infty,$$

and, moreover,

$$(4.5) \quad \underline{u}_x > 0, \quad 0 \leq x \leq 1, \quad t \geq 0.$$

Proof. Step 1. Construction of the subsolution. In what follows we shall omit the variables t and/or y when no confusion is likely. We take \underline{u} as in (4.1) where we assume

$$(4.6) \quad a > 0 \quad \text{and} \quad \lim_{t \rightarrow \infty} a(t) = \infty.$$

We compute

$$(4.7) \quad \begin{aligned} \underline{u}_x &= \frac{a}{(1+y)^2} + baf'(y) - b^2ag'(y), \\ \underline{u}_{xx} &= \frac{-2a^2}{(1+y)^3} + ba^2f''(y) - b^2a^2g''(y), \end{aligned}$$

and

$$\underline{u}_t = \frac{a'x}{(1+y)^2} + \left[b'f(y) + ba'xf'(y) - 2bb'g(y) - b^2a'xg'(y) \right];$$

hence

$$\underline{u}_t = \frac{a'}{a} \frac{y}{(1+y)^2} + \left[b'f(y) + \frac{ba'}{a}yf'(y) - 2bb'g(y) - \frac{b^2a'}{a}yg'(y) \right].$$

Recall that the operator \mathcal{P} is defined in (3.26). It follows that

$$\begin{aligned} \mathcal{P}\underline{u} &= \frac{a'}{a} \frac{y}{(1+y)^2} + \left[b'f + \frac{ba'}{a}yf' - 2bb'g - \frac{b^2a'}{a}yg' \right] \\ &\quad + 2a \frac{y}{(1+y)^3} - bayf'' + b^2ayg'' \\ &\quad - 2a \left[\frac{y}{1+y} + bf - b^2g \right] \left[\frac{1}{(1+y)^2} + bf' - b^2g' \right]. \end{aligned}$$

Collecting terms of the same order in b yields

$$(4.8) \quad \begin{aligned} \mathcal{P}\underline{u} &= \frac{a'}{a} \frac{y}{(1+y)^2} + \left[b'f + \frac{ba'}{a}yf' - 2bb'g - \frac{b^2a'}{a}yg' \right] - ab \left[yf'' + \frac{2f}{(1+y)^2} + \frac{2yf'}{1+y} \right] \\ &\quad + ab^2 \left[yg'' + \frac{2g}{(1+y)^2} + \frac{2yg'}{1+y} - 2ff' \right] + 2ab^3 \left[f'g + fg' \right] - 2ab^4gg'. \end{aligned}$$

The natural scaling of the equation leads to the choice

$$(4.9) \quad b := \frac{a'}{a^2},$$

so that in the right-hand side (RHS) of (4.8), the first term will be of the same order as the terms in the second bracket. Assuming

$$(4.10) \quad a' > 0,$$

we also denote

$$(4.11) \quad \gamma := \left(\frac{a}{a'}\right)'$$

and observe that

$$(4.12) \quad \frac{a'}{a} = ab, \quad b' = -(1 + \gamma)ab^2,$$

where the last equality comes from

$$\gamma = \left(\frac{a}{a'}\right)' = \left(\frac{1}{ab}\right)' = -\frac{b'}{ab^2} - \frac{a'}{a^2b} = -\frac{b'}{ab^2} - 1.$$

To make things clear, let us already stress that the final choice of a will guarantee

$$\gamma \geq 0 \quad \text{and} \quad \lim_{t \rightarrow \infty} \gamma(t) = 0.$$

Using (4.12), we can recast identity (4.8) in the form

$$\begin{aligned} \mathcal{P}\underline{u} &= ab \left[\frac{y}{(1+y)^2} - yf'' - \frac{2f}{(1+y)^2} - \frac{2yf'}{1+y} \right] \\ &+ ab^2 \left[yg'' + \frac{2g}{(1+y)^2} + \frac{2yg'}{1+y} - 2ff' + yf' - (1+\gamma)f \right] \\ &+ ab^3 \left[2f'g + 2fg' - yg' + 2(1+\gamma)g \right] - 2ab^4gg'; \end{aligned}$$

hence

$$(4.13) \quad \begin{aligned} \mathcal{P}\underline{u} &= ab \left[\frac{y}{(1+y)^2} - \mathcal{L}f \right] + ab^2 \left[\mathcal{L}g - 2ff' + yf' - (1+\gamma)f \right] \\ &+ ab^3 \left[2f'g + 2fg' - yg' + 2(1+\gamma)g \right] - 2ab^4gg'. \end{aligned}$$

Let us now choose

$$(4.14) \quad f := (I + \mathcal{L}_0^{-1}) \left(\frac{y}{(1+y)^2} \right) \geq \frac{y}{(1+y)^2} \geq 0,$$

which solves $\mathcal{L}f = y/(1+y)^2$ for $y > 0$, with $f(0) = 0$ and $f'(0) = 1$ (cf. subsection 3.4). Dividing by ab^2 , we see that $\mathcal{P}\underline{u}$ has the same sign as the quantity

$$A := \left[\mathcal{L}g - 2ff' + yf' - (1+\gamma)f \right] + b \left[2f'g + 2fg' - yg' + 2(1+\gamma)g \right] - 2b^2gg'.$$

Next we choose

$$g := \mathcal{L}_0^{-1}(2ff' - yf' + f).$$

Therefore,

$$A = -\gamma f + b \left[2f'g + 2fg' - yg' + 2(1+\gamma)g \right] - 2b^2gg'.$$

We now proceed to show that the quantity A is nonpositive for large t by considering separately the regions $y_0 \leq y \leq a(t)$ and $0 \leq y \leq y_0$ for some large y_0 independent of t . At this point, we make the additional assumptions that

$$(4.15) \quad \gamma \sim \frac{1}{\log a} \quad \text{as } t \rightarrow \infty$$

and

$$(4.16) \quad a' \sim \frac{a}{\log a}; \quad \text{hence } b \sim \frac{1}{a \log a} \quad \text{as } t \rightarrow \infty$$

(which will be verified on the final choice of the function a ; actually more precise expansions of γ and b will be needed in the final matching process). By Lemmas 3.5 and 3.6, we have

$$f \sim \log y, \quad f'g = o(g), \quad fg' = o(g), \quad \text{and } yg' \sim g \sim \frac{y \log y}{2} \quad \text{as } y \rightarrow \infty.$$

Consequently, fixing $\delta > 0$ and taking y_0 and t_0 large enough, we have, for $y_0 \leq y \leq a(t)$ and $t \geq t_0$,

$$-\gamma f \leq (-1 + \delta) \frac{\log y}{\log a}, \quad g, g' \geq 0,$$

and

$$\begin{aligned} 2f'g + 2fg' - yg' + 2(1 + \gamma)g &\leq \delta y \log y - \frac{1}{2}y \log y + 2(1 + \gamma) \left(\frac{1}{2} + \delta\right) y \log y \\ &\leq \left(\frac{1}{2} + 4\delta\right) y \log y. \end{aligned}$$

Taking $\delta = 1/12$ and also using (4.16), we thus obtain

$$A \leq (-1 + \delta) \frac{\log y}{\log a} + \frac{1}{a \log a} \left(\frac{1}{2} + 5\delta\right) y \log y = \frac{\log y}{\log a} \left(-1 + \delta + \left(\frac{1}{2} + 5\delta\right) \frac{y}{a}\right) \leq 0$$

for $y_0 \leq y \leq a(t)$ and $t \geq t_0$ (possibly larger). Next, for $0 \leq y \leq y_0$, (4.14) implies

$$f(y) \geq c_1 y, \quad 0 \leq y \leq y_0,$$

whereas $f(0) = g(0) = 0$ yields

$$|2f'g + 2fg' - yg'| + 4|g| + |gg'| \leq c_2 y, \quad 0 \leq y \leq y_0,$$

with $c_1, c_2 > 0$. Consequently,

$$A \leq -c_1 \gamma(t)y + c_2 b(t)y = \left[-c_1 + c_2 \frac{b(t)}{\gamma(t)}\right] \gamma(t)y \leq 0$$

on $[0, y_0]$ for t large enough, due to (4.15) and (4.16).

We have thus proved that, under conditions (4.6), (4.10), (4.15), and (4.16), there holds $\mathcal{P}\underline{u} \leq 0$ in $(0, 1) \times (T, \infty)$ for T large enough. By a time shift we may obviously take $T = 0$.

Step 2. Determination of $a(t)$ by matching at the outer boundary. Next, the determination of $a(t)$ will be done by “matching” with the boundary condition at $x = 1$, i.e., by writing

$$\underline{u}(1, t) < 1,$$

which is equivalent to

$$(4.17) \quad bf(a) - b^2g(a) < \frac{1}{a + 1}.$$

It is of course sufficient to check (4.17) for large t (thanks to the possibility of shifting time). Let us first sketch the resolution of (4.17) in a rough way. Using (4.9) and applying Lemmas 3.5(i) and 3.6(i) at leading order, we are left with

$$(4.18) \quad \frac{a'}{a^2} \left(\log a - \frac{a'}{a^2} \frac{a \log a}{2} \right) \lesssim \frac{1}{a + 1} \quad \text{as } t \rightarrow \infty.$$

We expect the second term in the bracket of the left-hand side (LHS) of the preceding relation to be much smaller than the first one as $a \rightarrow \infty$. If we ignore it, we obtain the differential inequality

$$(4.19) \quad a' \lesssim \frac{a}{\log a} \quad \text{as } t \rightarrow \infty,$$

which implies that

$$a(t) \lesssim e^{\sqrt{2t}} \quad \text{as } t \rightarrow \infty.$$

However, the latter estimation is not accurate enough to show the desired estimate, and we thus add a correction term in (4.19). More precisely, we look for $a(t)$ as the solution of

$$(4.20) \quad a' = \frac{a}{\log a}(1 + \eta), \quad t > 0, \quad a(0) = 2,$$

where the correction term $\eta = \eta(a)$ has the form

$$(4.21) \quad \eta = \frac{5}{2 \log a} + \frac{K}{\log^2 a}, \quad K > 0.$$

Now plugging (4.20) into (4.17), recalling also (4.9), and using the exact asymptotic behavior of $f(a), g(a)$ as $a \rightarrow \infty$, provided by Lemmas 3.5(i) and 3.6(i), we are reduced to the condition

$$\frac{1 + \eta}{a \log a} \left[\log a - 2 + O\left(\frac{\log^2 a}{a}\right) - \frac{1 + \eta}{a \log a} \left(\frac{a \log a}{2} - \frac{9a}{4} + O(\log^3 a) \right) \right] < \frac{1}{a + 1}$$

as $t \rightarrow \infty$ or

$$(4.22) \quad (1 + \eta) \left[1 - \frac{5 + \eta}{2 \log a} + \frac{9}{4 \log^2 a} + \frac{9\eta}{4 \log^2 a} + O\left(\frac{\log a}{a}\right) \right] < \frac{a}{a + 1} \quad \text{as } t \rightarrow \infty.$$

Plugging (4.21) into (4.22), we first obtain

$$\left(1 + \frac{5}{2 \log a} + \frac{K}{\log^2 a} \right) \left[1 - \frac{5}{2 \log a} + \frac{1}{\log^2 a} + O\left(\frac{1}{\log^3 a}\right) \right] < 1 - \frac{1}{a + 1} \quad \text{as } t \rightarrow \infty$$

and finally

$$(4.23) \quad 1 + \frac{4K - 21}{4 \log^2 a} + O\left(\frac{1}{\log^3 a}\right) < 1 - \frac{1}{a + 1} \quad \text{as } t \rightarrow \infty.$$

In order for (4.23) to be satisfied, we choose $K < \frac{21}{4}$, and we then obtain the following ODE for a :

$$(4.24) \quad a' = \frac{a}{\log a} \left(1 + \frac{5}{2 \log a} + \frac{K}{\log^2 a}\right), \quad t > 0, \quad a(0) = 2.$$

Here we should mention that the same form for η , given by (4.21), will also be considered for the supersolution, with a different constant K ; see the next subsection. Equation (4.24) implies that

$$a' = \frac{a}{\log a - 5/2} (1 + O(\log^{-2} a)) \quad \text{as } t \rightarrow \infty,$$

and, integrating with respect to t , we get

$$\log^2 a - 5 \log a = 2t + O\left(\int_0^t \frac{ds}{\log^2(a(s))}\right) \quad \text{as } t \rightarrow \infty.$$

Solving the quadratic polynomial in $\log a$ and noting that $\log a(s) \geq \sqrt{2s}$ by (4.24), we end up with (4.3). As for (4.4), it follows from

$$(4.25) \quad \begin{aligned} b &= \frac{a'}{a^2} = \frac{1}{a \log a} \left(1 + \frac{5}{2 \log a} + \frac{K}{\log^2 a}\right) = \frac{1 + O(\log^{-1} a)}{a \log a} \\ &= \frac{1 + O(t^{-1/2})}{a \log a} \quad \text{as } t \rightarrow \infty, \end{aligned}$$

where we used (4.24) and $\log a \geq \sqrt{2t}$ as $t \rightarrow \infty$. On the other hand, denoting $G(s) = s + (5/2)s^2 + Ks^3$ and using (4.24), we see that $\gamma = (a/a')'$ satisfies

$$(4.26) \quad \gamma = \left[\frac{1}{G(1/\log a)} \right]' = \frac{a'}{a \log^2 a} \frac{G'(1/\log a)}{G^2(1/\log a)} = \frac{G'(1/\log a)}{\log^2 a G(1/\log a)} = H(1/\log a),$$

where $H(s) = s(1 + 5s + 3Ks^2)(1 + (5/2)s + Ks^2)^{-1}$. Finally, the assumed properties (4.6), (4.10), (4.15), and (4.16) of a, b, γ are immediate consequences of (4.24), (4.25), and (4.26).

Step 3. Proof of (4.5). By (4.7) we have

$$a^{-1}u_x = \frac{1}{(1 + y)^2} + bf'(y) - b^2g'(y).$$

By Lemma 3.6 and (4.16), taking y_1 and t_1 large enough, we have, for $t \geq t_1$ and $y_1 \leq y \leq a(t)$,

$$a^{-1}u_x \geq \frac{1}{2a^2} - \frac{\log y}{a^2 \log^2 a} \geq \frac{1}{a^2} \left(\frac{1}{2} - \frac{1}{\log a}\right) > 0.$$

Now, for $0 \leq y \leq y_1$ and $t \geq t_1$ possibly larger, we get

$$a^{-1}u_x \geq \frac{1}{(1 + y_1)^2} - Cb(t) - Cb^2(t) > 0.$$

By a time shift we may obviously take $t_1 = 0$, and (4.5) is proved. \square

Remark 4.1. It is still possible to obtain a qualitatively correct subsolution with just a two-term expansion, for instance, by making the simple choice $f(y) = \log(1+y)$, $g = 0$ in (4.1). However, this yields only the lower grow-up rate up to a multiplicative constant, i.e., $u_x(0, t) \geq Ce^{\sqrt{2t}}$, and does not enable one to deduce an expansion of the form (2.8).

4.2. Construction of a supersolution. The form (4.1) does not seem to be sufficient to construct an accurate supersolution (i.e., leading to a function $a(t)$ fulfilling (4.3)). We need a slight perturbation, corresponding to the modified ansatz

$$(4.27) \quad \bar{u}(x, t) = 1 - \frac{1}{y + 1} + b(t)f(y) - b^2(t)\tilde{g}(y, t), \quad y = a(t)x,$$

where

$$(4.28) \quad \tilde{g}(y, t) = (1 + \varepsilon(t))h(y),$$

and $\varepsilon(t)$ goes to 0 as $t \rightarrow \infty$.

LEMMA 4.2. *The problem*

$$\begin{aligned} \bar{u}_t - x\bar{u}_{xx} &\geq 2\bar{u}\bar{u}_x, & 0 < x < 1, t > 0, \\ \bar{u}(0, t) &= 0, & t \geq 0, \\ \bar{u}(1, t) &\geq 1, & t \geq 0, \end{aligned}$$

admits a solution of the form

$$\bar{u}(x, t) = 1 - \frac{1}{y + 1} + b(t)f(y) - (1 + \varepsilon(t))b^2(t)h(y), \quad y = a(t)x,$$

where the smooth functions $a(t), b(t), f(y), h(y)$ have the following properties:

$$f(y) \sim \log y, \quad h(y) \sim \frac{y \log y}{2} \quad \text{as } y \rightarrow \infty,$$

$$f(0) = h(0) = f'(0) = h'(0) = 0,$$

$$(4.29) \quad a(t) = \left(1 + O(t^{-1/2} \log t)\right) \exp \left[\frac{5}{2} + \sqrt{2t}\right] \quad \text{as } t \rightarrow \infty,$$

$$(4.30) \quad b(t) = \frac{1 + O(t^{-1/2})}{a(t) \log a(t)} \quad \text{as } t \rightarrow \infty,$$

$$(4.31) \quad \varepsilon(t) \sim (2t)^{-1/2} \quad \text{as } t \rightarrow \infty.$$

Proof. Step 1. Construction of the supersolution. Taking \bar{u} as defined in (4.27), (4.28), the expression for $\mathcal{P}\bar{u}$ is similar to (4.13), except that g, g', g'' are now replaced with $\tilde{g}, \tilde{g}_y, \tilde{g}_{yy}$ and an additional term $-b^2\varepsilon'h$ is added, which is inherited from \bar{u}_t . As in the proof of Lemma 4.1, b and γ are defined through (4.9) and (4.11), and we assume (4.6), (4.10), (4.15), and (4.16). This leads to

$$\begin{aligned} \mathcal{P}\bar{u} &= ab \left[\frac{y}{(1+y)^2} - \mathcal{L}f \right] + ab^2 \left[\mathcal{L}\tilde{g} - 2ff' + yf' - (1+\gamma)f - \frac{\varepsilon'}{a}h \right] \\ &\quad + ab^3 \left[2f'\tilde{g} + 2f\tilde{g}_y - y\tilde{g}_y + 2(1+\gamma)\tilde{g} \right] - 2ab^4\tilde{g}\tilde{g}_y. \end{aligned}$$

Replacing $\tilde{g}(y, t)$ with $(1 + \varepsilon(t))h(y)$, we obtain

$$(4.32) \quad \mathcal{P}\bar{u} = ab \left[\frac{y}{(1+y)^2} - \mathcal{L}f \right] + ab^2 \left[(1 + \varepsilon)\mathcal{L}h - 2ff' + yf' - (1 + \gamma)f - \frac{\varepsilon'}{a}h \right] \\ + (1 + \varepsilon)ab^3 \left[2f'h + 2fh' - yh' + 2(1 + \gamma)h \right] - 2(1 + \varepsilon)^2ab^4hh'.$$

Denote

$$B_0 := (1 + \varepsilon)\mathcal{L}h - 2ff' + yf' - (1 + \gamma)f - \frac{\varepsilon'}{a}h.$$

As in Lemma 4.1, we first choose

$$f := (I + \mathcal{L}_0^{-1}) \left(\frac{y}{(1+y)^2} \right).$$

Dividing identity (4.32) by ab^2 , we see that $\mathcal{P}\bar{u}$ has the same sign as the quantity

$$B := B_0 + (1 + \varepsilon)b \left[2f'h + 2fh' - yh' + 2(1 + \gamma)h \right] - 2(1 + \varepsilon)^2b^2hh'.$$

Next we choose

$$h = \mathcal{L}_0^{-1}(2ff' - yf' + f + M\varphi),$$

where φ satisfies

$$(4.33) \quad \varphi(0) = 0, \quad \varphi'(0) > 0, \quad \varphi(y) > 0 \text{ for } 0 < y < 2, \quad \varphi(y) = 1/\log y \text{ for } y \geq 2.$$

Note that, taking $M > 2$ suitably large, we have

$$2ff' - yf' + f + M\varphi \geq 0, \quad y \geq 0,$$

by Lemma 3.5; hence

$$(4.34) \quad h(y) \geq 0, \quad y \geq 0,$$

due to (3.35). We compute

$$B_0 = (1 + \varepsilon)(2ff' - yf' + f + M\varphi) - 2ff' + yf' - (1 + \gamma)f - \frac{\varepsilon'}{a}h \\ = \varepsilon(2ff' - yf') + M(1 + \varepsilon)\varphi + (\varepsilon - \gamma)f - \frac{\varepsilon'}{a}h.$$

At this point, we choose

$$(4.35) \quad \varepsilon = \gamma,$$

where γ is defined by (4.11), and we assume again (4.15) and (4.16) along with

$$(4.36) \quad \gamma' \leq 0$$

(these assumptions will be verified on the final choice of the function a). By (4.34) and (4.36) we have

$$B_0 \geq \gamma(2ff' - yf') + M(1 + \gamma)\varphi,$$

and it follows that

$$(4.37) \quad (1+\gamma)^{-1}B \geq \left[\frac{\gamma}{1+\gamma}(2ff' - yf') + M\varphi \right] + b \left[2f'h + 2fh' - yh' + 2(1+\gamma)h \right] - 2(1+\gamma)b^2hh'.$$

To show that the RHS of (4.37) is nonnegative for large t , we again consider separately the regions $y_0 \leq y \leq a(t)$ and $0 \leq y \leq y_0$ for some large y_0 independent of t . By the estimates in Lemma 3.5, we have

$$(4.38) \quad f \sim \log y, \quad f' \sim \frac{1}{y}, \quad \text{and} \quad yh' \sim h \sim \frac{y \log y}{2} \quad \text{as } y \rightarrow \infty.$$

Consequently, fixing $\delta > 0$, using (4.15), and taking y_0 and t_0 large enough, we have, for $y_0 \leq y \leq a(t)$ and $t \geq t_0$,

$$\begin{aligned} \frac{\gamma}{1+\gamma}(2ff' - yf') + M\varphi &\geq -\frac{3}{2}\gamma + \frac{M}{\log y} \geq -\frac{2}{\log a} + \frac{M}{\log y} \geq \frac{M-2}{\log y}, \\ 2f'h + 2fh' - yh' + 2(1+\gamma)h &\geq -\left(\frac{1}{2} + \delta\right)y \log y + 2(1+\gamma)\left(\frac{1}{2} - \delta\right)y \log y \geq \left(\frac{1}{2} - 4\delta\right)y \log y, \end{aligned}$$

and

$$2(1+\gamma)hh' \leq y \log^2 y.$$

Assuming $M \geq 3$, taking $\delta = 1/8$, and also using (4.16), we infer that

$$(1+\gamma)^{-1}B \geq \frac{1}{\log a} - \frac{y \log^2 y}{a^2 \log^2 a} = \frac{1}{\log a} \left(1 - \frac{y \log^2 y}{a^2 \log a} \right) \geq \frac{1}{\log a} \left(1 - \frac{\log a}{a} \right) \geq 0$$

for $y_0 \leq y \leq a(t)$ and $t \geq t_0$. Next, for $0 \leq y \leq y_0$, (4.33) implies

$$M\varphi(y) \geq c_1 y,$$

whereas $f(0) = h(0) = 0$ yields

$$2ff' - yf' \geq -c_2 y, \quad 2f'h + 2fh' - yh' + 2(1+\gamma)h \geq -c_2 y, \quad 2(1+\gamma)hh' \leq c_2 y,$$

with $c_1, c_2 > 0$. Therefore,

$$(1+\gamma)^{-1}B \geq -\gamma(t)c_2 y + c_1 y - b(t)c_2 y - b^2(t)c_2 y = \left[c_1 - c_2(\gamma(t) + b(t) + b^2(t)) \right] y \geq 0$$

on $[0, y_0]$ for t large enough.

We have thus proved that $\mathcal{P}\bar{u} \geq 0$ in $(0, 1) \times (T, \infty)$ for T large enough. By a time shift we may obviously take $T = 0$.

Step 2. Determination of $a(t)$ by matching at the outer boundary. The determination of $a(t)$ will be done again by “matching” with the boundary condition at $x = 1$. Imposing that

$$\underline{u}(1, t) \geq 1,$$

we obtain

$$bf(a) - b^2(1 + \varepsilon)h(a) \geq \frac{1}{a + 1},$$

and taking (4.9) into account, we end up with

$$(4.39) \quad \frac{a'}{a^2} \left(f(a) - \frac{a'}{a^2}(1 + \varepsilon)h(a) \right) \geq \frac{1}{a + 1}.$$

Again it suffices to check (4.39) for large t . Following the same reasoning as in the case of a subsolution, we again look for $a(t)$ as a solution of

$$(4.40) \quad a' = \frac{a}{\log a}(1 + \eta) \quad \text{as } t \rightarrow \infty,$$

where the correction term $\eta = \eta(a)$ is given by (4.21) with K a constant to be determined.

Plugging (4.40) into (4.39) and using the exact asymptotic behavior of $f(a), h(a)$ as $a \rightarrow \infty$, given by Lemmas 3.5 and 3.6, we arrive at

$$\frac{1 + \eta}{a \log a} \left[\log a - 2 + O\left(\frac{\log^2 a}{a}\right) - \frac{(1 + \eta)(1 + \varepsilon)}{a \log a} \left(\frac{a \log a}{2} - \frac{9a}{4} + O\left(\frac{a}{\log a}\right) \right) \right] \geq \frac{1}{a + 1}$$

or

$$(4.41) \quad (1 + \eta) \left[1 - \frac{2}{\log a} + O\left(\frac{\log a}{a}\right) - \frac{(1 + \eta)(1 + \varepsilon)}{\log a} \left(\frac{1}{2} - \frac{9}{4 \log a} + O\left(\frac{1}{\log^2 a}\right) \right) \right] \geq \frac{a}{a + 1}$$

as $t \rightarrow \infty$. Denote by Γ the quantity in the bracket in the LHS of (4.41). Using (4.21) and (4.41), we obtain

$$\begin{aligned} \Gamma &= 1 - \frac{2}{\log a} + O\left(\frac{\log a}{a}\right) \\ &\quad - \left(1 + \frac{5}{2 \log a} + \frac{K}{\log^2 a} \right) \left(1 + \frac{1}{\log a} + \frac{5}{2 \log^2 a} \right) \left(\frac{1}{2 \log a} - \frac{9}{4 \log^2 a} + O\left(\frac{1}{\log^3 a}\right) \right) \\ &= 1 - \frac{5}{2 \log a} + \frac{1}{2 \log^2 a} + O\left(\frac{1}{\log^3 a}\right) \quad \text{as } t \rightarrow \infty. \end{aligned}$$

Then (4.41) becomes equivalent to

$$\left(1 + \frac{5}{2 \log a} + \frac{K}{\log^2 a} \right) \left(1 - \frac{5}{2 \log a} + \frac{1}{2 \log^2 a} + O\left(\frac{1}{\log^3 a}\right) \right) \geq 1 - \frac{1}{a + 1} \quad \text{as } t \rightarrow \infty;$$

that is,

$$(4.42) \quad 1 + \frac{4K - 23}{4 \log^2 a} + O\left(\frac{1}{\log^3 a}\right) \geq 1 - \frac{1}{a + 1} \quad \text{as } t \rightarrow \infty.$$

For (4.42) to be satisfied we choose $K > \frac{23}{4}$, and we again take a to be the solution of the ODE (4.24). Then, by the end of step 2 of the proof of Lemma 4.1, we obtain (4.29) and (4.30) as well as the assumed properties (4.6), (4.10), (4.15), and (4.16) of a, b, γ . Finally, we note that (4.15), (4.29), and (4.35) guarantee (4.31), and that (4.26) implies

$$\gamma' = \frac{-a'}{a \log^2 a} H'(1/\log a) = \frac{-1}{\log^2 a} (H'G)(1/\log a) \sim \frac{-1}{\log^3 a} \quad \text{as } t \rightarrow \infty$$

and hence (4.36) (after a further time shift). \square

4.3. Proofs of Theorem 2.1 and Corollary 2.2. Let u be the solution of (2.1)–(2.4) and let \underline{u}, \bar{u} be the sub-/supersolutions provided by Lemmas 4.1 and 4.2. The asymptotic expansion (2.8)–(2.9) in Theorem 2.1 will be an immediate consequence of the following two Lemmas. The first one guarantees that u lies between suitable time-shifts of \underline{u} and \bar{u} . The second one shows that (shifted versions of) \underline{u} and \bar{u} satisfy the required asymptotic behavior.

LEMMA 4.3.

(i) *There exists $T_1 > 0$ such that*

$$(4.43) \quad u(\cdot, t) \geq \underline{u}(\cdot, t - T_1), \quad t \geq T_1.$$

(ii) *Let τ be as in Lemma 3.3. Then there exists $T_2 > 0$ such that*

$$u(\cdot, t) \leq \bar{u}(\cdot, t + T_2), \quad t \geq \tau.$$

Proof. (i) Since $\underline{u}(\cdot, 0) \in C^1([0, 1])$ with $\underline{u}(0, 0) = 0$, $\underline{u}(1, 0) < 1$, and $\underline{u}_x(\cdot, 0) > 0$, it follows from Lemma 3.4 that $u(\cdot, T_1) \geq \underline{u}(\cdot, 0)$ for some $T_1 > 0$. The assertion then follows from the comparison principle.

(ii) Due to (4.38), (4.14), and $h(0) = 0$, we have

$$f(y) \geq c_1 \log(y + 1) \quad \text{and} \quad h(y) \leq c_2 (y + 1) \log(y + 1), \quad y \geq 0,$$

for some $c_1, c_2 > 0$. This along with (4.29)–(4.31) implies that, for all $t \geq t_2$ large enough and all $x \in [0, 1]$,

$$\begin{aligned} \bar{u}(x, t) &= \frac{ax}{ax + 1} + bf(ax) - (1 + \varepsilon(t))b^2h(ax) \\ &\geq \frac{ax}{ax + 1} + b \log(ax + 1)(c_1 - 2c_2(ax + 1)b) \\ &\geq \frac{ax}{ax + 1} + b \log(ax + 1) \left(c_1 - 3c_2 \frac{a + 1}{a \log a} \right) \geq \frac{ax}{ax + 1}. \end{aligned}$$

Take τ, η as in Lemma 3.3 and set $x_0 = 1/4K$. For all $x \in [0, x_0]$ and $t \geq t_2$, with $t_2 \geq \tau$ possibly larger, we have $a/2K \geq ax + 1$; hence

$$\bar{u}(x, t) \geq \frac{ax}{ax + 1} \geq 2Kx \geq u(x, \tau), \quad 0 \leq x \leq x_0, \quad t \geq t_2.$$

On the other hand, (4.38) implies that

$$f'(y) \leq c_3 \quad \text{and} \quad |h'(y)| \leq c_3 + \log(y + 1), \quad y \geq 0,$$

for some $c_3 > 0$. For all $x \in (x_0, 1]$ and $t \geq t_2$ (possibly larger), we thus have

$$\begin{aligned} \bar{u}_x &= \frac{a}{(1 + ax)^2} + baf'(ax) - (1 + \varepsilon(t))b^2ah'(ax) \\ &\leq \frac{a}{(1 + ax_0)^2} + c_3ba + 2(c_3 + \log(a + 1))b^2a \leq \eta, \end{aligned}$$

due to (4.29) and (4.30); hence

$$\bar{u}(x, t) \geq 1 - \eta(1 - x) \geq u(x, \tau), \quad x_0 < x \leq 1, \quad t \geq t_2.$$

Therefore, $\bar{u}(\cdot, t_2) \geq u(\cdot, \tau)$ in $[0, 1]$. The assertion, with $T_2 = t_2 - \tau$, thus follows from the comparison principle. \square

LEMMA 4.4. *Let*

$$A(t) = \exp \left[\frac{5}{2} + \sqrt{2t} \right].$$

Then, for any $T \in \mathbb{R}$, each of the functions $w = \underline{w}$ and $w = \bar{w}$ satisfies

$$1 - w(x, t + T) = \frac{1}{1 + A(t)x} \left[1 - x + O(t^{-1/2} \log t) \right]$$

as $t \rightarrow \infty$, uniformly in $[0, 1]$.

Proof. We shall give the proof for $w = \underline{w}$, the other case being completely similar.

Set $\tilde{a}(t) = a(t + T)$ and $\tilde{b}(t) = b(t + T)$. We first note that, by (4.3), we have

$$\tilde{a}(t) = \left(1 + O(t^{-1/2} \log t) \right) A(t)$$

and

$$(4.44) \quad \log \tilde{a}(t) = \log A(t) + O(t^{-1/2} \log t) = \left(1 + O(t^{-1} \log t) \right) \log A(t).$$

Also, by (4.4), we have

$$(4.45) \quad \tilde{b}(t) = \frac{1 + O(t^{-1/2})}{\tilde{a} \log \tilde{a}}.$$

Moreover, due to Lemma 3.5, there exists $C > 0$ such that

$$(4.46) \quad |f(y) - \log(1 + y)| \leq C, \quad y \geq 0.$$

Now, using (4.45), (4.46), and (4.44), we compute

$$\begin{aligned} \frac{1}{1 + \tilde{a}x} - \tilde{b}f(\tilde{a}x) &= \frac{1}{1 + \tilde{a}x} \left[1 - \frac{(1 + \tilde{a}x)f(\tilde{a}x)}{\tilde{a} \log \tilde{a}} \left(1 + O(t^{-1/2}) \right) \right] \\ &= \frac{1}{1 + \tilde{a}x} \left[1 - \frac{(1 + \tilde{a}x) \log(1 + \tilde{a}x)}{\tilde{a} \log \tilde{a}} \left(1 + O(t^{-1/2}) \right) \right] \\ &= \frac{1}{1 + \tilde{a}x} \left[1 - \frac{x \log(1 + \tilde{a}x)}{\log \tilde{a}} + O(t^{-1/2}) \right] \\ &= \frac{1}{1 + \tilde{a}x} \left[1 - x + R(x, t) + O(t^{-1/2}) \right], \end{aligned}$$

where

$$R(x, t) := \frac{x(\log \tilde{a} - \log(1 + \tilde{a}x))}{\log \tilde{a}}.$$

Here and in what follows, the O 's are uniform in $[0, 1]$. To control R we note that, if $\tilde{a}x \geq 1$, then

$$\log(\tilde{a}x) \leq \log(1 + \tilde{a}x) \leq \log(\tilde{a}x) + \log 2,$$

and hence

$$|R(x, t)| \leq \frac{x(|\log x| + \log 2)}{\log \tilde{a}} \leq \frac{C}{\log \tilde{a}},$$

whereas, if $\tilde{a}x < 1$, then $|R(x, t)| \leq x \leq 1/\tilde{a}$. It follows that $\sup_{x \in [0,1]} |R(x, t)| = O(t^{-1/2})$ as $t \rightarrow \infty$. Since, by (4.3), we have

$$\begin{aligned} 1 + \tilde{a}x &= (1 + Ax) \left[1 + \frac{(\tilde{a} - A)x}{1 + Ax} \right] = (1 + Ax) \left[1 + \frac{O(t^{-1/2} \log t)Ax}{1 + Ax} \right] \\ &= (1 + Ax)(1 + O(t^{-1/2} \log t)), \end{aligned}$$

we deduce that

$$\frac{1}{1 + \tilde{a}x} - \tilde{b}f(\tilde{a}x) = \frac{1 + O(t^{-1/2} \log t)}{1 + Ax} (1 - x + O(t^{-1/2})),$$

and hence

$$(4.47) \quad \frac{1}{1 + \tilde{a}x} - \tilde{b}f(\tilde{a}x) = \frac{1 - x + O(t^{-1/2} \log t)}{1 + Ax}.$$

On the other hand, due to Lemma 3.6 and $g(0) = 0$, there exists $C_1 > 0$ such that

$$|g(y)| \leq C_1(1 + y) \log(1 + y), \quad y \geq 0.$$

Consequently, using also (4.45), we obtain, for $t \rightarrow \infty$,

$$(4.48) \quad (1 + Ax)\tilde{b}^2|g(\tilde{a}x)| \leq 2C_1 \frac{(1 + \tilde{a}x)^2 \log(1 + \tilde{a}x)}{\tilde{a}^2 \log^2 \tilde{a}} \leq \frac{3C_1}{\log \tilde{a}} = O(t^{-1/2}).$$

Combining (4.47) and (4.48) finally yields

$$1 - \underline{u}(x, t + T) = \frac{1}{1 + \tilde{a}x} - \tilde{b}f(\tilde{a}x) + \tilde{b}^2g(\tilde{a}x) = \frac{1 - x + O(t^{-1/2} \log t)}{1 + Ax}. \quad \square$$

Proof of (2.10) and (2.11). First notice that estimate (3.27) in Lemma 3.3 and Lemma 4.3(ii) guarantee the control of the slope at $x = 1$, namely, (3.22), for any finite $T > 0$. The C^1 regularity property (2.10) is then a consequence of Lemma 3.2. To show (2.11), note that $\underline{u}_x(0, t) = a(t)$ due to (4.7) and (4.2). The lower estimate corresponding to (2.11) is then a consequence of (4.43) and (4.3). The proof of the upper part is similar by using \bar{u} . \square

Proof of Corollary 2.2. Assertion (i) is an immediate consequence of (2.8) and (2.9). To show (ii), it suffices to observe that, due to (2.9),

$$\begin{aligned} \int_0^1 \frac{dx}{1 + A(t)x} &= \left[\frac{\log(1 + A(t)x)}{A(t)} \right]_0^1 = \frac{\log(1 + A(t))}{A(t)} \\ &= (1 + O(t^{-1/2}))\sqrt{2t} \exp \left[-\frac{5}{2} - \sqrt{2t} \right] \end{aligned}$$

as $t \rightarrow \infty$ and to use

$$\int_0^1 \frac{x dx}{1 + A(t)x} \leq \frac{1}{A(t)}. \quad \square$$

Acknowledgments. This research was performed while Nikos I. Kavallaris was a visitor at the Laboratoire Analyse Géométrie et Applications in Université Paris-Nord. He is grateful to this institution for its hospitality and stimulating atmosphere. He would like to express his sincere thanks to Piotr Biler, who introduced him to the chemotaxis problem, and to Andrew A. Lacey for stimulating discussions.

REFERENCES

- [1] F. BAVAUD, *Equilibrium properties of the Vlasov functional: The generalized Poisson-Boltzmann-Emden equation*, Rev. Modern Phys., 63 (1991), pp. 129–149.
- [2] P. BILER, *Existence and nonexistence of solutions for a model of gravitational interaction of particles III*, Colloq. Math., 68 (1995), pp. 229–239.
- [3] P. BILER, *Local and global solvability of some parabolic systems modeling chemotaxis*, Adv. Math. Sci. Appl., 8 (1998), pp. 715–743.
- [4] P. BILER, D. HILHORST, AND T. NADZIEJA, *Existence and nonexistence of solutions for a model of gravitational interaction of particles II*, Colloq. Math., 67 (1994), pp. 297–308.
- [5] P. BILER, G. KARCH, P. LAURENÇOT, AND T. NADZIEJA, *The 8π -problem for radially symmetric solutions of a chemotaxis model in a disc*, Topol. Methods Nonlinear Anal., 27 (2006), pp. 133–144.
- [6] P. BILER, G. KARCH, P. LAURENÇOT, AND T. NADZIEJA, *The 8π -problem for radially symmetric solutions of a chemotaxis model in the plane*, Math. Methods Appl. Sci., 29 (2006), pp. 1563–1583.
- [7] P. BILER AND T. NADZIEJA, *Existence and nonexistence of solutions for a model of gravitational interaction of particles I*, Colloq. Math., 66 (1994), pp. 319–334.
- [8] P. BILER AND T. NADZIEJA, *Growth and accretion of mass in an astrophysical model II*, Appl. Math., 23 (1995), pp. 351–361.
- [9] P. BILER AND T. NADZIEJA, *A nonlocal singular parabolic problem modelling gravitational interaction of particles*, Adv. Differential Equations, 3 (1998), pp. 177–197.
- [10] A. BLANCHET, J.A. CARRILLO, AND N. MASMOUDI, *Infinite time aggregation for the critical two-dimensional Patlak-Keller-Segel model*, Comm. Pure Appl. Math., 61 (2008), pp. 1449–1481.
- [11] A. BLANCHET, J. DOLBEAULT, AND B. PERTHAME, *Two dimensional Keller-Segel model: Optimal critical mass and qualitative properties of solutions*, Electron. J. Differential Equations, 44 (2006), pp. 1–32.
- [12] M.P. BRENNER, P. CONSTANTIN, L.P. KADANOFF, A. SCHENKEL, AND S.C. VENKATARAMANI, *Diffusion, attraction and collapse*, Nonlinearity, 12 (1999), pp. 1071–1098.
- [13] S. CHANDRASEKHAR, *An Introduction to the Study of Stellar Structure*, Dover, New York, 1967.
- [14] S. CHILDRRESS AND J.K. PERCUS, *Nonlinear aspects of chemotaxis*, Math. Biosci., 56 (1981), pp. 217–237.
- [15] M. CHIPOT AND F.B. WEISSLER, *Some blowup results for a nonlinear parabolic equation with a gradient term*, SIAM J. Math. Anal., 20 (1989), pp. 886–907.
- [16] M. DOI AND S.F. EDWARDS, *The Theory of Polymer Dynamics*, Clarendon Press, Oxford, UK, 1986.
- [17] J. DOLBEAULT AND B. PERTHAME, *Optimal critical mass in the two-dimensional Keller-Segel model in \mathbb{R}^2* , C. R. Math. Acad. Sci. Paris, 339 (2004), pp. 611–616.
- [18] J.W. DOLD, V.A. GALAKTIONOV, A.A. LACEY, AND J.L. VÁZQUEZ, *Rate of approach to a singular steady state in quasilinear reaction-diffusion equations*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 26 (1998), pp. 663–687.
- [19] E. FEIREISL, H. PETZELTOVA, AND P. LAURENÇOT, *On convergence to equilibria for the Keller-Segel chemotaxis model*, J. Differential Equations, 236 (2007), pp. 551–569.
- [20] M. FILA, M. WINKLER, AND E. YANAGIDA, *Grow-up rate of solutions for a supercritical semilinear diffusion equation*, J. Differential Equations, 205 (2004), pp. 365–389.
- [21] M. FILA, J.R. KING, M. WINKLER, AND E. YANAGIDA, *Optimal lower bound of the grow-up rate for a supercritical parabolic equation*, J. Differential Equations, 228 (2006), pp. 339–356.
- [22] V.A. GALAKTIONOV AND J.R. KING, *Composite structure of global unbounded solutions of nonlinear heat equations with critical Sobolev exponents*, J. Differential Equations, 189 (2003), pp. 199–233.
- [23] I.A. GUERRA AND M.A. PELETIER, *Self-similar blow-up for a diffusion-attraction problem*, Nonlinearity, 17 (2004), pp. 2137–2162.

- [24] M.A. HERRERO, *The mathematics of chemotaxis*, in Handbook of Differential Equations: Evolutionary Equations, Vol. 3, Elsevier, New York, 2007, pp. 137–193.
- [25] M.A. HERRERO AND J.J.L. VELÁZQUEZ, *Singularity patterns in a chemotaxis model*, Math. Ann., 306 (1996), pp. 583–623.
- [26] M.A. HERRERO AND J.J.L. VELÁZQUEZ, *Chemotactic collapse for the Keller-Segel model*, J. Math. Biol., 35 (1996), pp. 177–194.
- [27] M.A. HERRERO AND J.J.L. VELÁZQUEZ, *A blow-up mechanism for a chemotaxis model*, Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4), 24 (1997), pp. 633–683.
- [28] D. HORSTMANN, *From 1970 until present: The Keller-Segel model in chemotaxis and its consequences I*, Jahresber. Deutsch. Math.-Verein., 105 (2003), pp. 103–165.
- [29] D. HORSTMANN, *From 1970 until present: The Keller-Segel model in chemotaxis and its consequences II*, Jahresber. Deutsch. Math.-Verein., 106 (2004), pp. 51–69.
- [30] W. JÄGER AND S. LUCKHAUS, *On explosions of solutions to a system of partial differential equations modelling chemotaxis*, Trans. Amer. Math. Soc., 329 (1992), pp. 819–824.
- [31] N.I. KAVALLARIS AND T. SUZUKI, *On the finite-time blow-up of a non-local parabolic equation describing chemotaxis*, Differential Integral Equations, 20 (2007), pp. 293–308.
- [32] E.F. KELLER AND L.A. SEGEL, *Initiation of slime mold aggregation viewed as an instability*, J. Theor. Biol., 26 (1970), pp. 399–415.
- [33] A.A. LACEY AND D. TZANETIS, *Global existence and convergence to a singular steady state for a semilinear heat equation*, Proc. Roy. Soc. Edinburgh Sect. A, 105 (1987), pp. 289–305.
- [34] N. MIZOGUCHI, *Growup of solutions for a semilinear heat equation with supercritical nonlinearity*, J. Differential Equations, 227 (2006), pp. 652–669.
- [35] T. NAGAI, *Blow-up of radially symmetric solutions to a chemotaxis system*, Adv. Math. Sci. Appl., 5 (1995), pp. 1–21.
- [36] T. NAGAI, *Blow-up of nonradial solutions to parabolic-elliptic systems modeling chemotaxis in two-dimensional domains*, J. Inequal. Appl., 6 (2001), pp. 37–55.
- [37] T. NAGAI, T. SENBA, AND K. YOSHIDA, *Application of the Trudinger-Moser inequality to a parabolic system of chemotaxis*, Funkcial. Ekvac., 40 (1997), pp. 411–433.
- [38] V. NANJUNDIAH, *Chemotaxis, signal relaying, and aggregation morphology*, J. Theor. Biol., 42 (1973), pp. 63–105.
- [39] K. OHTSUKA, T. SENBA, AND T. SUZUKI, *Blowup in infinite time in the simplified system of chemotaxis*, Adv. Math. Sci. Appl., 17 (2007), pp. 445–472.
- [40] C.S. PATLAK, *Random walk with persistence and external bias*, Bull. Math. Biol. Biophys., 15 (1953), pp. 311–338.
- [41] B. PERTHAME, *PDE models for chemotactic movements: Parabolic, hyperbolic and kinetic*, Appl. Math., 49 (2004), pp. 539–564.
- [42] P. POLAČIK AND E. YANAGIDA, *On bounded and unbounded global solutions of a supercritical semilinear heat equation*, Math. Ann., 327 (2003), pp. 745–771.
- [43] P. QUITTNER AND PH. SOUPLLET, *Superlinear Parabolic Problems. Blow-up, Global Existence and Steady States*, Birkhäuser Advanced Texts, Birkhäuser Verlag, Basel, 2007.
- [44] T. SENBA AND T. SUZUKI, *Chemotactic collapse in a parabolic-elliptic system of mathematical biology*, Adv. Differential Equations, 6 (2001), pp. 21–50.
- [45] T. SENBA, *Blowup in infinite time of radial solutions for a parabolic-elliptic system in high dimensional Euclidean spaces*, Nonlinear Anal., to appear.
- [46] C. SIRE AND P.-H. CHAVANIS, *Thermodynamics and collapse of self-gravitating Brownian particles in D dimensions*, Phys. Rev. E (3), 66 (2002), 046133.
- [47] PH. SOUPLLET AND J.L. VÁZQUEZ, *Stabilization towards a singular steady state with gradient blow-up for a diffusion-convection problem*, Discrete Contin. Dyn. Syst., 14 (2006), pp. 221–234.
- [48] T. SUZUKI, *Free Energy and Self-Interacting Particles*, Birkhäuser Boston, Boston, 2005.
- [49] G. WOLANSKY, *On steady distributions of self-attracting clusters under friction and fluctuations*, Arch. Ration. Mech. Anal., 119 (1992), pp. 355–391.
- [50] G. WOLANSKY, *On the evolution of self-interacting clusters and applications to semilinear equations with exponential nonlinearity*, J. Anal. Math., 59 (1992), 251–272.
- [51] G. WOLANSKY, *A critical parabolic estimate and application to nonlocal equations arising in chemotaxis*, Appl. Anal., 66 (1997), 291–321.

MATHEMATICAL RESULTS ON EXISTENCE FOR VISCOELASTODYNAMIC PROBLEMS WITH UNILATERAL CONSTRAINTS*

ADRIEN PETROV[†] AND MICHELLE SCHATZMAN[‡]

Abstract. This paper focuses on a damped wave equation and the evolution of a Kelvin–Voigt viscoelastic material, both problems being subject to unilateral boundary conditions. Under appropriate regularity assumptions on the initial data, both problems possess a weak solution which is obtained as the limit of a sequence of solutions of penalized problems; the functional properties of all the traces are precisely identified through Fourier analysis, and this enables us to infer the existence of a strong solution, i.e., a solution satisfying almost everywhere the unilateral conditions.

Key words. viscoelasticity, Signorini conditions, penalty method, traces, variational inequality, convolution

AMS subject classifications. 35L85, 49J40, 73D99, 73V25

DOI. 10.1137/070695101

1. Introduction and notation. This paper aims to give some new mathematical results on existence for a damped wave equation with an obstacle and for full viscoelasticity in the particular case of a Kelvin–Voigt material with unilateral boundary conditions.

We consider in section 2 a damped wave equation taking place in a half-space, with an obstacle at the boundary. Let $u(x, t)$ be the displacement at time t of the material point of spatial coordinate $x = (x_1, x') \in (-\infty, 0] \times \mathbb{R}^{d-1}$ at rest with $d \geq 2$. We will agree that if we write a function of space and time as a function of three variables, then the first variable is the normal space variable x_1 , the second variable is the tangential space variable x' , and the last variable is time. Let $f(x_1, x', t)$ denote a density of external forces, depending on space and time. Define $\Omega \stackrel{\text{def}}{=} (-\infty, 0] \times \mathbb{R}^{d-1}$, and let α be a positive number. The mathematical problem is formulated as follows:

$$(1.1) \quad u_{tt} - \Delta u - \alpha \Delta u_t = f, \quad x \in \Omega, \quad t > 0,$$

with Cauchy initial data

$$(1.2) \quad u(\cdot, 0) = u_0 \quad \text{and} \quad u_t(\cdot, 0) = u_1$$

and Signorini boundary conditions at $x_1 = 0$, $t > 0$,

$$(1.3) \quad 0 \leq u \perp u_{x_1} + \alpha u_{x_1 t} \geq 0,$$

where $(\cdot)_t \stackrel{\text{def}}{=} \frac{\partial}{\partial t}(\cdot)$ and $(\cdot)_{x_1} \stackrel{\text{def}}{=} \frac{\partial}{\partial x_1}(\cdot)$. The orthogonality has the natural meaning: if we have enough regularity, it means that the product $u(u_{x_1} + \alpha u_{x_1 t})$ vanishes almost everywhere on the boundary. If we do not have enough regularity, the above inequality

*Received by the editors June 22, 2007; accepted for publication (in revised form) September 23, 2008; published electronically January 7, 2009.

<http://www.siam.org/journals/sima/40-5/69510.html>

[†]Weierstraß-Institut für Angewandte Analysis und Stochastik, Mohrenstraße 39, 10117 Berlin, Germany (petrov@wias-berlin.de).

[‡]CNRS, Université de Lyon, Institut Camille Jordan, 21 Avenue Claude Bernard, F-69622 Villeurbanne Cedex, France (schatz@math.univ-lyon1.fr).

is integrated on an appropriate set of test functions, yielding a weak formulation for the unilateral condition. The main result of section 2 is that, indeed, (1.3) holds almost everywhere on the boundary; i.e., we have a strong solution to our problem.

We suppose that the initial position u_0 belongs to the Sobolev space $H^2(\Omega)$ and satisfies the compatibility condition $u_0(0, \cdot, \cdot) \geq 0$, the initial velocity u_1 belongs to $H^1(\Omega)$, and the density of forces f belongs to $L^2_{loc}([0, \infty); L^2(\Omega))$. The choice of a function f defined for all nonnegative time is justified by the use of a Fourier transform in the latter part of the article. This is not a significant restriction as we can always extend f by 0 if it is defined only for finite times.

Let us describe the weak formulation of the problem. Denote by K the convex set

$$K \stackrel{\text{def}}{=} \{v \in H^1_{loc}(\Omega \times [0, \infty)) : \nabla v_t \in L^2_{loc}([0, \infty); L^2(\Omega)), v|_{\{0\} \times \mathbb{R}^{d-1}} \geq 0\}.$$

This unusual convex set has been devised in order to write a weak formulation of our problem. Since we expect to find a scalar product $(\nabla u_t, \nabla w)$, we require ∇u_t to be square integrable. Thus, the weak formulation associated to (1.1)–(1.3) is obtained by multiplying (1.1) by $v - u$, $v \in K$, and by integrating formally over $\Omega \times (0, \tau)$. Then, we get

$$(1.4) \quad \begin{cases} \text{find } u \in K \text{ such that for all } v \in K \text{ and for all } \tau \in (0, \infty), \\ \int_{\Omega} (u_t(v - u))|_0^\tau dx - \int_0^\tau \int_{\Omega} u_t(v_t - u_t) dx dt \\ + \int_0^\tau \int_{\Omega} (\nabla u + \alpha \nabla u_t)(\nabla v - \nabla u) dx dt \geq \int_0^\tau \int_{\Omega} f(v - u) dx dt. \end{cases}$$

We approximate this problem by a regularized problem, and the bulk of our work is to get enough regularity to obtain traces of the tangential first derivatives of the solution on the boundary. This is done essentially through an analysis of the Dirichlet to Neumann operator for the damped wave equation. For this purpose, one must go to tangential Fourier variables and perform an estimate on a pseudodifferential term.

In section 3, we treat the evolution of a Kelvin–Voigt material (see [DaL90]) occupying a three dimensional half-space, satisfying Signorini conditions at the boundary and Cauchy data at $t = 0$. We make the assumptions of small deformations. Let $\varepsilon_{ij}(u) \stackrel{\text{def}}{=} (u_{j,x_i} + u_{i,x_j})/2$ be the strain tensor, and let there be given two Hooke tensors, a^n_{ijkl} , $n = 0, 1$. We define the two stress tensors σ^n_{ij} corresponding, respectively, to the elastic and the viscous parts of the stress:

$$(1.5) \quad \sigma^n_{ij}(u) \stackrel{\text{def}}{=} a^n_{ijkl} \varepsilon_{kl}(u);$$

here, we have used the summation convention on repeated indices. The displacement field u satisfies the system

$$(1.6) \quad \rho u_{i,tt} = \sigma^0_{ij,x_j}(u) + \sigma^1_{ij,x_j}(u_t) + f_i, \quad x \in \Omega, \quad t > 0.$$

The initial data are given by

$$(1.7) \quad u(\cdot, 0) = v_0 \quad \text{and} \quad u_t(\cdot, 0) = v_1.$$

The components of the unit external normal are δ_{1j} (δ is the Kronecker symbol), and a basis of tangential vectors can be taken as $\tau_j = \delta_{2j}$ and $\tau'_j = \delta_{3j}$. Denote by

$\Sigma = \{0\} \times \mathbb{R}^{d-1}$ the boundary of Ω . Then, the boundary conditions on $\Sigma \times [0, \infty)$ are

$$(1.8a) \quad 0 \geq u_1 \perp \sigma_{11}^0(u) + \sigma_{11}^1(u_t) \leq 0,$$

$$(1.8b) \quad \sigma_{12}^0(u) + \sigma_{12}^1(u_t) = 0, \quad \text{and} \quad \sigma_{13}^0(u) + \sigma_{13}^1(u_t) = 0.$$

One of the main results of section 3 is to show that (1.8a) holds almost everywhere, because the relevant traces exist.

In order to simplify the problem, we have considered a homogeneous and isotropic material; then, the Hooke tensors a_{ijkl}^n are defined with the help of Lamé constants λ^n and μ^n :

$$a_{ijkl}^n \stackrel{\text{def}}{=} \lambda^n \delta_{ij} \delta_{kl} + 2\mu^n \delta_{ik} \delta_{jl}, \quad n = 0, 1.$$

We define two elasticity operators A^n by

$$A^n u \stackrel{\text{def}}{=} a_{ijkl}^n \partial_j \varepsilon_{kl}(u), \quad n = 0, 1.$$

Then, the problem (1.6)–(1.8) can be rewritten as follows:

$$(1.9a) \quad \rho u_{tt} - A^0 u - A^1 u_t = f, \quad x \in \Omega, \quad t > 0,$$

$$(1.9b) \quad 0 \geq u_1 \perp (\sigma_{11}^0(u) + \sigma_{11}^1(u_t)) \leq 0 \quad \text{on} \quad \Sigma \times [0, \infty),$$

$$(1.9c) \quad \sigma_{12}^0(u) + \sigma_{12}^1(u_t) = 0 \quad \text{and} \quad \sigma_{13}^0(u) + \sigma_{13}^1(u_t) = 0 \quad \text{on} \quad \Sigma \times [0, \infty),$$

$$(1.9d) \quad u(\cdot, 0) = v_0 \quad \text{and} \quad u_t(\cdot, 0) = v_1.$$

Let us describe now the functional hypotheses on the data; if X is a space of scalar functions, the bold-face notation \mathbf{X} denotes the space X^d . For the final result, we require v_0 to belong to $\mathbf{H}^{5/2}(\Omega)$, v_1 to $\mathbf{H}^{3/2}(\Omega)$, and f to $\mathbf{H}_{\text{loc}}^1([0, \infty); L^2(\Omega))$. The initial data must satisfy the compatibility condition $(v_0)_1(0, x') \leq 0$ for all $x' \in \Sigma$. Let K be the convex set defined by

$$K \stackrel{\text{def}}{=} \{v \in \mathbf{H}^1(\Omega \times (0, \tau)) : \nabla v_t \in \mathbf{L}^2(\Omega \times (0, \tau)), v(0, \cdot) \leq 0\}.$$

Define two bilinear forms by

$$a^0(u, v) \stackrel{\text{def}}{=} \int_{\Omega} a_{ijkl}^0 \varepsilon_{ij}(u) \varepsilon_{kl}(v) \, dx \quad \text{and} \quad a^1(u, v) \stackrel{\text{def}}{=} \int_{\Omega} a_{ijkl}^1 \varepsilon_{ij}(u) \varepsilon_{kl}(v) \, dx.$$

We obtain a weak formulation of the problem (1.9) as follows: we multiply (1.9a) by $v - u$, $v \in K$, and we formally integrate the result over $\Omega \times (0, \tau)$; we obtain then the variational inequality

$$(1.10) \quad \left\{ \begin{array}{l} \text{find } u \in K \text{ such that for all } v \in K \text{ and for all } \tau \in (0, \infty), \\ \int_0^\tau \int_{\Omega} \rho u_{tt} \cdot (v - u) \, dx \, dt + \int_0^\tau a^0(u, u - v) \, dt \\ + \int_0^\tau a^1(u_t, v - u) \, dt \geq \int_0^\tau \int_{\Omega} f \cdot (v - u) \, dx \, dt. \end{array} \right.$$

The existence result for (1.1)–(1.3) is easily established by the penalty method and was already proved by Jarušek et al. [JM*92] in the case of distributed constraints.

Jarušek has also proved in [Jar96] an existence result for (1.9), in a much more general and complicated case, since it allows for contact, a given friction at the boundary, a nonlinear constitutive law for viscoelasticity, and a general geometry. However,

the boundary conditions must be understood in the sense of duality, since the traces themselves are defined by duality. The reader can also refer to [KuS04a]. Moreover, notice that [KuS04b] contains some regularity results for a Signorini problem with normal compliance.

The reader may wonder why our proofs of existence for weak solutions are longer than other known proofs; the reason is that we prove more estimates in order to get more information. Hence, we cannot reuse previous proofs of existence of weak solutions. In particular, we need stronger results on the regularity of solutions, and these cannot be deduced from former results on regularity, which use stronger assumptions.

In the present paper, for both problems, we penalize the obstacle constraint, we construct a solution of the penalized problem, and we show the existence of a weak solution by passing to the limit with respect to the penalty parameter. Then, under appropriate regularity conditions on the data, we prove that the penalized solution has traces, which can be estimated, and therefore the limiting weak solution that we obtained is a strong solution.

We have chosen to limit ourselves to homogeneous and, in the case of viscoelasticity, isotropic media, since we use a Fourier transform and some calculus in the dual variables for obtaining the traces. We believe that more general cases could be considered, namely, variable coefficients and curved boundary. The method is probably straightforward, but it is exacting, since one would probably have to use the full technology of pseudodifferential operators.

Observe that nothing is known about uniqueness.

These two problems are treated in the same article, because they are quite close. Proofs for the second problem are shortened when very close to proofs for the first one. Nevertheless, there are substantial differences in detail, since the second problem is much more complicated than the first. In particular, the bulk of the proof in section 3 consists in obtaining a solution of a linear system through Fourier–Laplace transform and then in estimating this solution in anisotropic Sobolev spaces.

2. The damped wave equation with Signorini boundary conditions.

2.1. The penalized problem. We approximate (1.1)–(1.3) by the penalty method. This means that we replace the rigid constraint (1.3) by a very stiff response. When the constraint is active, the response is linear, and it vanishes when the constraint is not active. More precisely, letting $r^- \stackrel{\text{def}}{=} -\min(r, 0)$, we replace u by u^ϵ , which satisfies

$$(2.1) \quad u_{tt}^\epsilon - \Delta u^\epsilon - \alpha \Delta u_t^\epsilon = f, \quad x \in \Omega, \quad t > 0,$$

with initial data

$$(2.2) \quad u^\epsilon(\cdot, 0) = u_0 \quad \text{and} \quad u_t^\epsilon(\cdot, 0) = u_1$$

and boundary condition

$$(2.3) \quad (u_{x_1}^\epsilon + \alpha u_{x_1 t}^\epsilon)(0, \cdot, \cdot) = (u^\epsilon(0, \cdot, \cdot))^- / \epsilon.$$

Notice that (2.3) is the so-called *normal compliance condition* introduced by Martins and Oden [MaO88].

Define the following sets:

$$(2.4) \quad \forall \tau \in (0, \infty), \quad Q_\tau \stackrel{\text{def}}{=} \Omega \times (0, \tau) \quad \text{and} \quad I_\tau \stackrel{\text{def}}{=} \Sigma \times (0, \tau).$$

THEOREM 2.1. *Let $W_{\text{loc}} \stackrel{\text{def}}{=} \{u \in H^1_{\text{loc}}([0, \infty) \times \Omega) : \nabla u_t \in L^2_{\text{loc}}([0, \infty); L^2(\Omega))\}$. Assume that u_0 and u_1 belong to $H^1(\Omega)$, and f belongs to $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$; then for every $\epsilon > 0$ there exists a unique weak solution $u^\epsilon \in W_{\text{loc}}$ of the problem (2.1)–(2.3) such that*

$$(2.5a) \quad u^\epsilon \in L^\infty_{\text{loc}}([0, \infty); H^1(\Omega)),$$

$$(2.5b) \quad u_t^\epsilon \in L^2_{\text{loc}}([0, \infty); H^1(\Omega)),$$

$$(2.5c) \quad u_{tt}^\epsilon \in L^2_{\text{loc}}([0, \infty); L^2(\Omega)),$$

and for every $\tau \in (0, T)$ and for all $v \in W_{\text{loc}}$, the following variational equality is satisfied:

$$(2.6) \quad \int_{\Omega} ((u_t^\epsilon v)(\cdot, \tau) - u_1 v(\cdot, 0)) \, dx - \int_{Q_\tau} u_t^\epsilon v_t \, dx \, dt + \int_{Q_\tau} \nabla u^\epsilon \nabla v \, dx \, dt + \alpha \int_{Q_\tau} \nabla u_t^\epsilon \nabla v \, dx \, dt - \frac{1}{\epsilon} \int_{I_\tau} (u^\epsilon)^- v \, dx' \, dt = \int_{Q_\tau} f v \, dx \, dt.$$

Proof. The theorem is proved by the standard Galerkin method, and the reader can refer to [GGZ74] or the appendix of [JM*92]. \square

2.2. A priori estimates. We establish here estimates up to the boundary and interior estimates which will enable us later to infer the existence of a weak solution to (1.1)–(1.3).

LEMMA 2.2. *Assume that f belongs to $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, u_0 to $H^1(\Omega)$, and u_1 to $L^2(\Omega)$. Then u_t^ϵ and ∇u^ϵ are bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\Omega))$, ∇u_t^ϵ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, and $(u^\epsilon(0, \cdot, \cdot))^- / \sqrt{\epsilon}$ is bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\mathbb{R}^-))$ independently of $\epsilon > 0$. If, moreover, u_0 belongs to $H^2(\Omega)$, Δu^ϵ is bounded in $L^\infty_{\text{loc}}(0, \infty; L^2(\Omega))$ independently of ϵ .*

Proof. These estimates result from an application of the Gronwall lemma to the energy identity. We multiply (2.1) by u_t^ϵ , and we integrate this expression over Q_τ to get

$$\int_{Q_\tau} u_{tt}^\epsilon u_t^\epsilon \, dx \, dt - \int_{Q_\tau} \Delta u^\epsilon u_t^\epsilon \, dx \, dt - \alpha \int_{Q_\tau} \Delta u_t^\epsilon u_t^\epsilon \, dx \, dt = \int_{Q_\tau} f u_t^\epsilon \, dx \, dt.$$

We integrate the first integral in time in the above relation, we use Green’s formula for the second and third, and, with the help of the boundary conditions (2.3), we obtain

$$\begin{aligned} & \frac{1}{2} \int_{\Omega} (|u_t^\epsilon(\cdot, \tau)|^2 + |\nabla u^\epsilon(\cdot, \tau)|^2) \, dx + \alpha \int_{Q_\tau} |\nabla u_t^\epsilon|^2 \, dx \, dt + \frac{1}{2\epsilon} \int_{\Sigma} ((u^\epsilon(0, \cdot, \cdot))^-)^2 \Big|_0^\tau \, dx' \\ & = \int_{Q_\tau} f u_t^\epsilon \, dx \, dt + \frac{1}{2} \int_{\Omega} (|\nabla u_0|^2 + |u_1|^2) \, dx. \end{aligned}$$

We may deduce from the classical Gronwall lemma that u_t^ϵ and ∇u^ϵ are bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\Omega))$, ∇u_t^ϵ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, and $(u^\epsilon(0, \cdot, \cdot))^- / \sqrt{\epsilon}$ is bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\mathbb{R}^{d-1}))$ independently of $\epsilon > 0$.

For the bound on Δu^ϵ , we multiply (2.1) by $e^{t/\alpha}$, and we integrate from 0 to τ . After an integration by parts on the term

$$\int_0^\tau u_{tt}^\epsilon e^{t/\alpha} \, dt,$$

we find the identity

$$(2.7) \quad \begin{aligned} \Delta u^\epsilon(\cdot, \tau) &= \Delta u_0 e^{-\tau/\alpha} + \alpha^{-1} (u_t^\epsilon(\cdot, \tau) - u_1 e^{-\tau/\alpha}) \\ &\quad - \alpha^{-2} \int_0^\tau u_t^\epsilon(\cdot, t) e^{(t-\tau)/\alpha} dt - \alpha^{-1} \int_0^\tau f e^{(t-\tau)/\alpha} dt. \end{aligned}$$

The end of the proof is straightforward. \square

Remark 2.3. If we suppose that f vanishes for t large, then, independently of $\epsilon > 0$, u_t^ϵ , ∇u^ϵ , and Δu^ϵ are bounded in $L^\infty([0, \infty); L^2(\Omega))$ and ∇u_t^ϵ is bounded in $L^2([0, \infty); L^2(\Omega))$. These properties can be proved using the arguments given in the proof of Lemma 2.2, with the origin of time moved to T if $f(\cdot, t)$ vanishes for $t \geq T$; since the integral involving f vanishes, the conclusion is clear for u_t^ϵ , ∇u^ϵ , and ∇u_t^ϵ . For Δu^ϵ , we use the identity (2.7), with 0 replaced by $T \leq \tau$, and the conclusion is clear.

LEMMA 2.4. *Assume the hypotheses of Lemma 2.2. Then for all nonnegative, continuously differentiable, and compactly supported ψ on \mathbb{R}^{d-1} and for all $\tau \in [0, T]$,*

$$\int_{I_\tau} \frac{(u^\epsilon)^-}{\epsilon} \psi \, dx' \, dt$$

is bounded independently of $\epsilon > 0$. In particular, $(u^\epsilon(0, \cdot, \cdot))^-/\epsilon$ is a bounded measure on I_τ .

Proof. Let ϕ be a continuous function with compact support in \mathbb{R}^d ; we multiply (2.1) by ϕ and integrate over Q_τ ; thanks to the boundary conditions (2.3) and Green's formula, we obtain

$$\int_\Omega \phi u_t^\epsilon(\cdot, t) \Big|_0^\tau dx + \int_{Q_\tau} \nabla \phi (\nabla u^\epsilon + \alpha \nabla u_t^\epsilon) \, dx \, dt - \frac{1}{\epsilon} \int_{I_\tau} (u^\epsilon)^- \phi \, dx' \, dt = \int_{Q_\tau} \phi f \, dx \, dt.$$

Since the product $|zy|$ can be estimated by $|z|^2/2 + |y|^2/2$, we get the following inequality:

$$(2.8) \quad \begin{aligned} \frac{1}{\epsilon} \int_{I_\tau} (u^\epsilon)^- \phi \, dx' \, dt &\leq \frac{1}{2} \int_\Omega (|u_t^\epsilon(\cdot, \tau)|^2 + |u_1|^2) \, dx + \int_\Omega |\phi|^2 \, dx \\ &\quad + \int_{Q_\tau} |\nabla \phi \nabla u^\epsilon| \, dx \, dt + \alpha \int_{Q_\tau} |\nabla \phi \nabla u_t^\epsilon| \, dx \, dt + \int_{Q_\tau} |\phi f| \, dx \, dt. \end{aligned}$$

The right-hand side of (2.8) is bounded since f belongs to $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, u_1 to $L^2(\Omega)$, and u_t^ϵ , ∇u^ϵ , and ∇u_t^ϵ to $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$. Moreover, $(u^\epsilon(0, \cdot, \cdot))^-$ is nonnegative; if the trace ψ of ϕ over Σ is nonnegative, the inequality is clear. The last statement of the theorem is obtained by a classical approximation argument. Write $\mu^\epsilon = (u^\epsilon(0, \cdot, \cdot))^-/\epsilon$. Let ψ_n be an increasing sequence of nonnegative, continuously differentiable, and compactly supported functions on Σ , which are at most equal to ψ . Then the integrals of ψ_n against μ^ϵ converge to the integral of $\lim_n \psi_n$ against μ^ϵ , so that the integral of any nonnegative, continuous, and compactly supported function against μ^ϵ is nonnegative, and this is precisely the definition of a nonnegative measure on Σ . \square

Let us turn now to interior estimates.

LEMMA 2.5. *Assume the hypotheses of Lemma 2.2. Then for all $\beta > 0$, u_{tt}^ϵ and Δu_t^ϵ are bounded in $L^2_{\text{loc}}([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$, independently of $\epsilon > 0$.*

Proof. The idea of the proof is twofold: we multiply u^ϵ by a truncation function $\varphi \in C_0^\infty(\mathbb{R})$, and we define $v^\epsilon \stackrel{\text{def}}{=} \varphi u^\epsilon$; we will observe that $w^\epsilon \stackrel{\text{def}}{=} v_t^\epsilon$ satisfies a heat equation, whose right-hand side will be estimated thanks to the previous lemmas. Let us now go into the details.

Let φ be a truncation function which is equal to 1 if $x_1 \leq -\beta$ and to 0 if $x_1 \geq -\beta/2$ ($\beta > 0$). Then, we multiply u^ϵ by φ , which enables us to forget about the strongly nonlinear boundary conditions. Define

$$(2.9) \quad v^\epsilon(x_1, \cdot, \cdot) \stackrel{\text{def}}{=} \varphi(x_1)u^\epsilon(x_1, \cdot, \cdot).$$

The derivatives of v^ϵ are given by

$$(2.10a) \quad v_{tt}^\epsilon = \varphi u_{tt}^\epsilon,$$

$$(2.10b) \quad \Delta v^\epsilon = \varphi \Delta u^\epsilon + 2\varphi_{x_1} \nabla u^\epsilon + \varphi_{x_1 x_1} u^\epsilon,$$

$$(2.10c) \quad \Delta v_t^\epsilon = \varphi \Delta u_t^\epsilon + 2\varphi_{x_1} \nabla u_t^\epsilon + \varphi_{x_1 x_1} u_t^\epsilon.$$

Observe that, thanks to relations (2.1) and (2.10), we have

$$(2.11) \quad v_{tt}^\epsilon - \Delta v^\epsilon - \alpha \Delta v_t^\epsilon = \tilde{g}^\epsilon,$$

where $\tilde{g}^\epsilon \stackrel{\text{def}}{=} \varphi f - 2\varphi_{x_1}(\nabla u^\epsilon + \alpha \nabla u_t^\epsilon) - \varphi_{x_1 x_1}(u^\epsilon + \alpha u_t^\epsilon)$. Since $f, u_t^\epsilon, \nabla u^\epsilon, \nabla u_t^\epsilon$, and u^ϵ are bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, \tilde{g}^ϵ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$. Let us define

$$(2.12) \quad w^\epsilon \stackrel{\text{def}}{=} v_t^\epsilon \quad \text{and} \quad g^\epsilon \stackrel{\text{def}}{=} \tilde{g}^\epsilon + \Delta v^\epsilon.$$

Substituting (2.12) into (2.11), we obtain

$$(2.13) \quad w_t^\epsilon - \alpha \Delta w^\epsilon = g^\epsilon.$$

Let us prove now that w_t^ϵ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$. For this purpose, we multiply (2.13) by w_t^ϵ ; we integrate this expression over Ω :

$$\int_{\Omega} |w_t^\epsilon|^2 dx - \alpha \int_{\Omega} \Delta w^\epsilon w_t^\epsilon dx = \int_{\Omega} g^\epsilon w_t^\epsilon dx.$$

We use Green's formula in the second term on the left-hand side of the above expression, thus getting the following equality:

$$(2.14) \quad \int_{\Omega} |w_t^\epsilon|^2 dx + \alpha \int_{\Omega} \nabla w_t^\epsilon \nabla w^\epsilon dx = \int_{\Omega} g^\epsilon w_t^\epsilon dx.$$

We integrate (2.14) over $(0, \tau)$, we observe that the product $|g^\epsilon w_t^\epsilon|$ can be estimated by $|g^\epsilon|^2/2 + |w_t^\epsilon|^2/2$, and we obtain

$$(2.15) \quad \begin{aligned} & \frac{1}{2} \int_{Q_\tau} |w_t^\epsilon|^2 dx dt + \alpha \int_{\Omega} |\nabla w^\epsilon(\cdot, \tau)|^2 dx \\ & \leq \alpha \int_{\Omega} |\nabla w^\epsilon(\cdot, 0)|^2 dx + \frac{1}{2} \int_{Q_\tau} |g^\epsilon|^2 dx dt. \end{aligned}$$

Since u_1 belongs to $H^1(\Omega)$ and φ belongs to $C_0^\infty(\mathbb{R})$, $\nabla w^\epsilon(\cdot, 0) = \varphi_{x_1} u_1 + \varphi \nabla u_1$ is bounded in $L^2(\Omega)$. Moreover, g^ϵ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$ because Δv^ϵ and \tilde{g}^ϵ are bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$. Therefore, (2.9), (2.12), and (2.15) enable us to

deduce that u_{tt}^ϵ is bounded in $L^2_{loc}([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$. We use analogous arguments to show that Δu_t^ϵ is bounded in $L^2_{loc}([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$. We multiply (2.13) by Δw^ϵ , we integrate over Q_τ , and, thanks to Green's formula, we obtain

$$(2.16) \quad -\frac{1}{2} \int_{\Omega} |\nabla w^\epsilon|^2|_0^\tau dx - \alpha \int_{Q_\tau} |\Delta w^\epsilon|^2 dx dt = \int_{Q_\tau} g^\epsilon \Delta w^\epsilon dx dt.$$

The product $|g^\epsilon \Delta w^\epsilon|$ can be estimated by $|g^\epsilon|^2/(2\gamma) + \gamma|\Delta w^\epsilon|^2/2$, and if we choose $\gamma \in (0, 2\alpha)$, we obtain the following inequality:

$$(2.17) \quad \left(\alpha - \frac{\gamma}{2}\right) \int_{Q_\tau} |\Delta w^\epsilon|^2 dx dt \leq \frac{1}{2\gamma} \int_{Q_\tau} |g^\epsilon|^2 dx dt + \frac{1}{2} \int_{\Omega} |\nabla w^\epsilon(\cdot, 0)|^2 dx.$$

Since g^ϵ and $\nabla w^\epsilon(\cdot, 0)$ are, respectively, bounded in $L^2_{loc}([0, \infty); L^2(\Omega))$ and $L^2(\Omega)$, according to (2.9), (2.12), and (2.17), we infer that Δu_t^ϵ is bounded in the space $L^2_{loc}([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$. \square

2.3. Existence of a weak solution. In this section, we show that it is possible to pass to the limit in the variational formulation of the penalized problem to obtain a weak solution of (1.1)–(1.3). There is a minor subtlety due to the unboundedness of Ω .

THEOREM 2.6. *Assume the hypotheses of Lemma 2.2. Then there exists a solution of the variational inequality (1.4); this solution can be obtained as a limit of a subsequence of the penalty approximation defined by (2.1)–(2.3).*

Proof. Let v belong to K and φ be a function belonging to $C^\infty_0(\bar{\Omega} \times [0, \infty))$, which takes its values in $[0, 1]$. Multiplying (2.1) by $(v - u^\epsilon)\varphi$ and integrating over Q_τ and then observing that

$$\int_{I_\tau} ((u^\epsilon)^- \varphi(v - u^\epsilon)) dx' dt = \int_{I_\tau} (((u^\epsilon)^-)^2 \varphi) dx' dt + \int_{I_\tau} ((u^\epsilon)^- \varphi v) dx' dt$$

is nonnegative, we may deduce the following inequality:

$$(2.18) \quad \int_{\Omega} u_t^\epsilon \varphi(v - u^\epsilon)|_0^\tau dx - \int_{Q_\tau} u_t^\epsilon (\varphi(v - u^\epsilon))_t dx dt + \int_{Q_\tau} (\nabla u^\epsilon + \alpha \nabla u_t^\epsilon) \nabla (\varphi(v - u^\epsilon)) dx dt \geq \int_{Q_\tau} f \varphi(v - u^\epsilon) dx dt.$$

We infer from Lemmas 2.2 and 2.4 that it is possible to extract a subsequence, still denoted by u^ϵ , such that

$$(2.19a) \quad u^\epsilon \rightharpoonup u \text{ in } L^\infty_{loc}([0, \infty); L^2(\Omega)) \text{ weak } *$$

$$(2.19b) \quad u_t^\epsilon \rightharpoonup u_t \text{ in } L^\infty_{loc}([0, \infty); L^2(\Omega)) \text{ weak } *$$

$$(2.19c) \quad \nabla u^\epsilon \rightharpoonup \nabla u \text{ in } L^\infty_{loc}([0, \infty); L^2(\Omega)) \text{ weak } *$$

$$(2.19d) \quad \Delta u^\epsilon \rightharpoonup \Delta u \text{ in } L^\infty_{loc}([0, \infty); L^2(\Omega)) \text{ weak } *$$

$$(2.19e) \quad \nabla u_t^\epsilon \rightharpoonup \nabla u_t \text{ in } L^2_{loc}([0, \infty); L^2(\Omega)) \text{ weak.}$$

Define the set $Q_R \stackrel{\text{def}}{=} \{x : x_1 < 0, |x'| \leq R\} \times [0, R]$. Thanks to the classical compactness properties of injections of Sobolev spaces on bounded open sets, we see that for all $R > 0$, the restrictions of u^ϵ and ∇u^ϵ to Q_R converge strongly to their

respective limits in $L^2(Q_R)$; therefore, we can pass to the limit in all the terms of (2.18) except possibly the first two terms.

Let us prove that u_t is continuous from $[0, \infty)$ to $L^2(\Omega)$ equipped with the weak topology: we infer from the estimates of Lemma 2.5 that for all $\beta > 0$, u_{tt}^ϵ restricted to $x_1 < -\beta$ is bounded in $L^2_{loc}([0, \infty); L^2((-\infty, -\beta] \times \Sigma))$; therefore, u_t^ϵ converges to a function u_t whose restriction to $x_1 < -\beta$ is continuous from $[0, \infty)$ to $L^2((-\infty, -\beta] \times \Sigma)$. Let $t_j \in [0, \infty)$ be a sequence converging to $t_\infty < \infty$; as u_t belongs to $L^2_{loc}([0, \infty); L^2(\Omega))$, we may extract a subsequence, still denoted by t_j , such that

$$u_t^\epsilon(\cdot, t_j) \rightharpoonup z \quad \text{in } L^2(\Omega) \text{ weak.}$$

But since, for all $\beta > 0$,

$$u_t^\epsilon(\cdot, t_j)1_{\{x_1 < -\beta\}} \rightarrow u_t(\cdot, t_\infty)1_{\{x_1 < -\beta\}} \quad \text{in } L^2(\Omega),$$

we see that z must coincide with $u_t(\cdot, t_\infty)$ and that the whole sequence converges strongly to $u_t(\cdot, t_\infty)$; this proves that u_t is continuous from $[0, \infty)$ to $L^2(\Omega)$ weak.

Let us prove now that $u_t^\epsilon(\cdot, t)$ converges weakly to $u_t(\cdot, t)$ for all $t > 0$: let γ be an arbitrary positive number; let z belong to $L^2(\Omega)$; and denote by C_1 an upper bound for $|u_t^\epsilon|_{L^\infty([0, T]; L^2(\Omega))}$ with T fixed. We choose β so small that

$$\left(\int_{-\beta < x_1 < 0} |z|^2 dx \right)^{1/2} \leq \frac{\gamma}{4C_1};$$

then, for $t \in [0, T]$,

$$(2.20) \quad \left| \int_{\Omega} (u_t^\epsilon(\cdot, t) - u_t(\cdot, t))z dx \right| \leq \left| \int_{x_1 < -\beta} (u_t^\epsilon(\cdot, t) - u_t(\cdot, t))z dx \right| + \left(\int_{-\beta < x_1 < 0} |z|^2 dx \right)^{1/2} \left(\int_{-\beta < x_1 < 0} |u_t^\epsilon(\cdot, t) - u_t(\cdot, t)|^2 dx \right)^{1/2}.$$

By definition of C_1 , the second term on the right-hand side of (2.20) is estimated by $C_1\gamma/(2C_1) = \gamma/2$. As $u_t^\epsilon|_{(-\infty, -\beta) \times I_T}$ is bounded in $H^1((-\infty, -\beta) \times I_T)$, we see that

$$\int_{-\infty}^{-\beta} \int_{\Sigma} u_t^\epsilon z dx \quad \text{converges to} \quad \int_{-\infty}^{-\beta} \int_{\Sigma} u_t z dx$$

uniformly with respect to $t \in [0, T]$. It suffices therefore to choose ϵ so small that the first term on the right-hand side of (2.20) is estimated by $\gamma/2$. This proves that the convergence of $\int_{\Omega} u_t^\epsilon z dx$ to $\int_{\Omega} u_t z dx$ is uniform on compact sets in time. In particular, as ϵ tends to 0, it is plain that for all $\tau > 0$,

$$\int_{\Omega} u_t^\epsilon \varphi(v - u^\epsilon) dx \rightarrow \int_{\Omega} u_t \varphi(v - u) dx.$$

Let us turn now to the term

$$\int_{Q_\tau} u_t^\epsilon (\varphi_t(v - u^\epsilon) + \varphi(v_t - u_t^\epsilon)) dx dt.$$

It is clear that

$$\int_{Q_\tau} u_t^\epsilon (\varphi_t(v - u^\epsilon) + \varphi v_t) dx dt \rightarrow \int_{Q_\tau} u_t (\varphi_t(v - u) + \varphi v_t) dx dt.$$

There remains to prove the convergence

$$\int_{Q_\tau} |u_t^\epsilon|^2 \varphi \, dx \, dt \rightarrow \int_{Q_\tau} |u_t|^2 \varphi \, dx \, dt.$$

We observe that

$$\begin{aligned} \int_{Q_\tau} |u_t^\epsilon - u_t|^2 \varphi \, dx \, dt &\leq \int_0^\tau \int_\Sigma \int_{x_1 \leq -\beta} |u_t^\epsilon - u_t|^2 \varphi \, dx \, dt \\ &\quad + \int_0^\tau \int_\Sigma \int_{-\beta \leq x_1 \leq 0} |u_t^\epsilon - u_t|^2 \varphi \, dx \, dt. \end{aligned}$$

Let γ be any positive number. One can deduce from the estimates of $|u_t^\epsilon|_{L^2(I_\tau)}$ and $|\nabla u_t^\epsilon|_{L^2(Q_\tau)}$ that there exists a constant C_2 independent from ϵ such that

$$|u_t^\epsilon(x_1, \cdot, \cdot)|_{L^2(\Sigma \times (0, \tau))} \leq C_2.$$

Therefore,

$$\int_0^\tau \int_\Sigma \int_{-\beta \leq x_1 \leq 0} |u_t^\epsilon - u_t|^2 \varphi \, dx \, dt \leq C_2^2 \beta.$$

We choose β so small that $C_2^2 \beta \leq \gamma/2$; then we know from the estimates of Lemmas 2.2 and 2.5 that the restriction of u^ϵ to $\{x_1 < -\beta\}$ intersected with a ball containing the support of φ is bounded in H^2 of that set; therefore, for ϵ small enough,

$$\int_{Q_\tau} |u_t^\epsilon - u_t|^2 \varphi \, dx \, dt \leq \frac{\gamma}{2},$$

and the convergence of the first two terms of (2.18) is proved. We observe now that since $u, u_t, \nabla u,$ and ∇u_t belong to $L^2_{loc}([0, \infty); L^2(\Omega))$, we may replace φ by φ_R in the variational inequality where φ_R is equal to 1 over the set Q_R and vanishes outside of Q_{R+1} . It is plain that as $R \rightarrow \infty$ all the terms in (2.18) converge to their limit; thus we have proved the existence of the desired weak solution. \square

Remark 2.7. Nothing is known about uniqueness.

2.4. Auxiliary results on the damped wave equation with Dirichlet boundary conditions. We establish a priori estimates on the damped wave equation with Dirichlet boundary conditions. These estimates will enable us to give some properties on the trace spaces which we use in the next subsection.

LEMMA 2.8. *Assume u_0 belongs to $H^{5/2}(\Omega)$; then, there exists a function $z \in H^3(\Omega \times [0, \infty))$ with compact support in t such that the trace of z on Ω is equal to u_0 .*

Proof. We extend u_0 into a function belonging to $H^{5/2}(\mathbb{R}^d)$: as the boundary of Ω is smooth, this extension is a consequence of classical results on Sobolev spaces. Then there exists a function Z belonging to $H^3(\mathbb{R}^d \times [0, \infty))$ whose trace is u_0 . It suffices now to select a cut-off function $\varphi \in C^\infty([0, \infty))$ which is equal to 1 on $[0, 1]$ and to 0 on $[2, \infty)$ and to define z as the restriction of φZ to $\Omega \times [0, \infty)$. \square

LEMMA 2.9. *Assume u_0 belongs to $H^{5/2}(\Omega)$, u_1 belongs to $H^1(\Omega)$, and f belongs to $L^2_{loc}([0, \infty); L^2(\Omega))$. Define z as in Lemma 2.8, and let \bar{u} be the solution of (1.1) with the initial data (1.2) and boundary condition $\bar{u}(0, \cdot, \cdot) = z(0, \cdot, \cdot)$. Then the trace $\bar{g} \stackrel{\text{def}}{=} -(\bar{u}_{x_1} + \alpha \bar{u}_{x_1 t})(0, \cdot, \cdot)$ is well defined and belongs to $L^2_{loc}([0, \infty); L^2(\mathbb{R}^{d-1}))$. Moreover, if f is compactly supported in time,*

$$\int_0^\tau |\bar{g}(\cdot, t)|^2_{L^2(\Sigma)} \, dt$$

increases at most polynomially with respect to τ .

Proof. The function $\zeta \stackrel{\text{def}}{=} \bar{u} - z$ satisfies the equation

$$(2.21) \quad \zeta_{tt} - \Delta\zeta - \alpha\Delta\zeta_t = F, \quad x \in \Omega, \quad t > 0,$$

where $F \stackrel{\text{def}}{=}} f - z_{tt} + \Delta z + \alpha\Delta z_t$, with initial data

$$\zeta(\cdot, 0) = 0 \quad \text{and} \quad \zeta_t(\cdot, 0) = u_1$$

and the Dirichlet boundary condition $\zeta(0, \cdot, \cdot) = 0$. If we multiply (2.21) by ζ_t and integrate, and if we suppose that f and F are compactly supported in time, we may easily deduce that $\zeta_t, \nabla\zeta$ are bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\Omega))$ and $\nabla\zeta_t$ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$. In order to get more information, we multiply (2.21) by $\Delta\zeta_t$; since the boundary term vanishes, we get immediately the identity

$$\begin{aligned} & \alpha \int_{Q_\tau} |\Delta\zeta_t|^2 \, dx \, dt + \frac{1}{2} \int_\Omega |\Delta\zeta(\cdot, \tau)|^2 \, dx + \int_\Omega |\nabla\zeta_t(\cdot, \tau)|^2 \, dx \\ &= \int_\Omega |\nabla\zeta_t(\cdot, 0)|^2 \, dx - \int_{Q_\tau} F\Delta\zeta_t \, dx \, dt. \end{aligned}$$

The product $F\Delta\zeta_t$ can be estimated by $\alpha|\Delta\zeta_t|^2/2 + |F|^2/(2\alpha)$; then $\Delta\zeta_t$ is bounded in $L^2_{\text{loc}}([0, \infty); L^2(\Omega))$, and $\Delta\zeta$ and $\nabla\zeta_t$ are bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\Omega))$. In particular, if the support in time of F is bounded, $\Delta\zeta_t$ is bounded in $L^2([0, \infty); L^2(\Omega))$ and $\Delta\zeta$ and $\nabla\zeta_t$ are bounded in $L^\infty([0, \infty); L^2(\Omega))$. Hence $\zeta_{x_1 t}(0, \cdot, \cdot)$ and $\zeta_{x_1}(0, \cdot, \cdot)$ belong, respectively, to $L^2_{\text{loc}}([0, \infty); H^{1/2}(\mathbb{R}^{d-1}))$ and to $L^\infty_{\text{loc}}([0, \infty); H^{1/2}(\mathbb{R}^{d-1}))$, and if the support in time of f is bounded, the local character of these spaces may be removed. \square

2.5. Regularity of the trace. We characterize the trace spaces using Fourier analysis, and we prove that u is a strong solution of (1.1)–(1.3). Here, we mean by strong solution that all the traces can be defined.

Let ν be a positive number. Denote by $v^\epsilon \stackrel{\text{def}}{=} e^{-\nu t}(u^\epsilon - \bar{u})$ a solution of

$$(2.22a) \quad (\nu + \partial_t)^2 v^\epsilon - (1 + \alpha(\nu + \partial_t))\Delta v^\epsilon = 0, \quad x \in \Omega, \quad t > 0,$$

$$(2.22b) \quad (1 + \alpha(\nu + \partial_t))v^\epsilon_{x_1}(0, \cdot, \cdot) = e^{-\nu t}\bar{g} - (v^\epsilon(0, \cdot, \cdot) + e^{-\nu t}\bar{u}(0, \cdot, \cdot))^-/\epsilon,$$

$$(2.22c) \quad v^\epsilon(\cdot, 0) = 0 \quad \text{and} \quad v^\epsilon_t(\cdot, 0) = 0.$$

We denote by $\xi \stackrel{\text{def}}{=} (\xi_2, \dots, \xi_d)^\top$ and ω , respectively, the dual variables to $x' \stackrel{\text{def}}{=} (x_2, \dots, x_d)^\top$ and t . The Fourier transform of $u(0, x', t)$ is $\widehat{u}(0, \xi, \omega)$, where the convention for the Fourier transform is

$$\widehat{u}(0, \xi, \omega) = \int_{\mathbb{R}^d} e^{-i(\xi \cdot x' + \omega t)} u(0, x', t) \, dx' \, dt.$$

Then $u(0, x', t)$ belongs to the Sobolev space $H^{a,b}_{\text{loc}}(\mathbb{R}^{d-1} \times [0, \infty))$, $(a, b) \in \mathbb{R}^2$, iff $|\xi|^a \widehat{u}(0, \xi, \omega)$ and $|\omega|^b \widehat{u}(0, \xi, \omega)$ belong to $L^2(\mathbb{R}^d)$.

We apply a partial Fourier transform in the tangential variable to (2.22a), and we get the following differential equation:

$$(2.23) \quad \widehat{v}^\epsilon_{x_1 x_1} = \left(|\xi|^2 + \frac{(\nu + i\omega)^2}{1 + \alpha(\nu + i\omega)} \right) \widehat{v}^\epsilon.$$

Define $\widehat{\lambda}$ to be

$$\widehat{\lambda}(\xi, \omega) \stackrel{\text{def}}{=} \sqrt{|\xi|^2 + \frac{(\nu + i\omega)^2}{1 + \alpha(\nu + i\omega)}};$$

thus $\widehat{\lambda}$ is holomorphic in the lower half-plane $\Im(\omega) < 0$ and $\Re\widehat{\lambda} \geq 0$ for $\Im(\omega) = 0$. The general solution of (2.23) is given by $\widehat{a}^\epsilon e^{\widehat{\lambda}x_1} + \widehat{b}^\epsilon e^{-\widehat{\lambda}x_1}$; since we performed a Fourier transform on v^ϵ , we assumed implicitly that v^ϵ and \widehat{v}^ϵ are tempered, respectively, in (x', t) and (ξ, ω) . We remark that the term $\widehat{b}^\epsilon e^{-\widehat{\lambda}x_1}$ can be tempered only if \widehat{b}^ϵ decays at infinity very quickly, and since this must be true for all x_1 , it implies that \widehat{b}^ϵ vanishes; the proof is similar to that given in [PeS02]. We deduce that the solution of (2.23) is $\widehat{a}^\epsilon e^{\widehat{\lambda}x_1}$. In particular,

$$(2.24) \quad ((1 + \alpha(\nu + \partial_t))v_{x_1}^\epsilon) \mathcal{Y}(0, \xi, \omega) = \widehat{\lambda}_1 \widehat{v}^\epsilon(0, \xi, \omega),$$

where $\widehat{\lambda}_1 \stackrel{\text{def}}{=} (1 + \alpha(\nu + i\omega))\widehat{\lambda}$. Define

$$g(x', t) \stackrel{\text{def}}{=} e^{-\nu t} \widehat{g}(x', t) \quad \text{and} \quad h(x', t) \stackrel{\text{def}}{=} e^{-\nu t} \widehat{u}(0, x', t).$$

If we let $w^\epsilon(x', t)$ be the trace $v^\epsilon(0, x', t)$, (2.22) can now be written as

$$(2.25) \quad \lambda_1 * w^\epsilon = g + (w^\epsilon + h)^- / \epsilon,$$

where w^ϵ vanishes for all $t \leq 0$.

Remark 2.10. It is clear that $\widehat{\lambda}$ is a holomorphic function in $\Im(\omega) < 0$, and thus we may deduce that λ_1 is a causal distribution.

LEMMA 2.11. *Let u^ϵ be the solution of (2.1)–(2.3). Then we may extract a subsequence, still denoted by u^ϵ , such that*

$$u^\epsilon(0, \cdot, \cdot) \rightharpoonup u(0, \cdot, \cdot) \quad \text{weakly in} \quad H_{\text{loc}}^{1/2, 5/4}(\mathbb{R}^{d-1} \times [0, \infty)).$$

Moreover, u is a strong solution of (1.1)–(1.3).

Proof. Formally, we multiply (2.25) by $\alpha(\nu w^\epsilon + w_t^\epsilon) + w^\epsilon$, and we estimate the pseudodifferential term in the Fourier variable; we obtain

$$(2.26) \quad \begin{aligned} & \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{\lambda}_1 \widehat{w}^\epsilon \overline{(1 + \alpha(\nu + i\omega))\widehat{w}^\epsilon} \, d\omega \, d\xi \\ &= \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{g} \overline{(1 + \alpha(\nu + i\omega))\widehat{w}^\epsilon} \, d\omega \, d\xi \\ &+ \frac{1}{\epsilon} \int_0^\infty \int_{\mathbb{R}^{d-1}} (w^\epsilon + h)^- (1 + \alpha(\nu + \partial_t))w^\epsilon \, dx' \, dt. \end{aligned}$$

Since $(u^\epsilon(0, \cdot, \cdot))^- / \sqrt{\epsilon}$ is bounded in $L_{\text{loc}}^\infty([0, \infty); L^2(\mathbb{R}^{d-1}))$, the absolute value of the second integral in the right-hand side of (2.26) is bounded, and we infer that there exists $C_1 > 0$ such that

$$(2.27) \quad \Re \int_{\mathbb{R}^d} \widehat{\lambda}_1 |\widehat{w}^\epsilon|^2 \overline{(1 + \alpha(\nu + i\omega))} \, d\omega \, d\xi \leq C_1 + \Re \int_{\mathbb{R}^d} \widehat{g} \overline{(1 + \alpha(\nu + i\omega))\widehat{w}^\epsilon} \, d\omega \, d\xi.$$

On the other hand, we have

$$\Re \widehat{\lambda}^2 = |\xi|^2 + \frac{\nu^2(1 + \alpha\nu) + (-1 + \alpha\nu)\omega^2}{|1 + \alpha(\nu + i\omega)|^2} \quad \text{and} \quad \Im \widehat{\lambda}^2 = \frac{2\nu\omega + \alpha\omega(\nu^2 + \omega^2)}{|1 + \alpha(\nu + i\omega)|^2}.$$

We may choose ν such that $\nu\alpha = 1$; we get then

$$(2.28) \quad \Re\widehat{\lambda}^2 = |\xi|^2 + \frac{2}{\alpha^2|2 + i\alpha\omega|^2} \quad \text{and} \quad \Im\widehat{\lambda}^2 = \frac{\omega(3 + \alpha^2\omega^2)}{\alpha|2 + i\alpha\omega|^2}.$$

Therefore, we infer that

$$\arg \widehat{\lambda} = \frac{1}{2} \arctan \left(\frac{|\xi|^2|2 + i\alpha\omega|^2 + 2}{\alpha\omega(3 + \alpha^2\omega^2)} \right).$$

According to (2.28), $\arg \widehat{\lambda}$ belongs to $[0, \pi/4]$, and since $\widehat{\lambda}$ is never equal to zero, we get for $|\xi| + |\omega| \gg 1$ the following inequality:

$$(2.29) \quad \Re\widehat{\lambda} \geq C(1 + |\xi| + \sqrt{|\omega|}).$$

Therefore, we obtain

$$C \int_{\mathbb{R}^d} |2 + i\alpha\omega|^2 (1 + |\xi| + \sqrt{|\omega|}) |\widehat{w}^\epsilon|^2 \, d\omega \, d\xi \leq C_1 + \int_{\mathbb{R}^d} |2 + i\alpha\omega| |\widehat{g}| |\widehat{w}^\epsilon| \, d\omega \, d\xi.$$

We estimate the product $|zy|$ by $|z|^2/(2\gamma) + \gamma|y|^2/2$, $\gamma > 0$; we see that

$$(2.30) \quad \begin{aligned} & \left(C - \frac{\gamma}{2}\right) \int_{\mathbb{R}^d} |2 + i\alpha\omega|^2 (1 + |\xi| + \sqrt{|\omega|}) |\widehat{w}^\epsilon|^2 \, d\omega \, d\xi \\ & \leq C_1 + \frac{1}{2\gamma} \int_{\mathbb{R}^d} \frac{|\widehat{g}|^2}{1 + |\xi| + \sqrt{|\omega|}} \, d\omega \, d\xi. \end{aligned}$$

We choose γ such that $\gamma < 2C$; since g belongs to $L^2([0, \infty); H^{1/2}(\mathbb{R}^{d-1}))$, then it is easy to deduce from (2.30) that $u^\epsilon(0, \cdot, \cdot)$ is bounded in $H_{\text{loc}}^{1/2, 5/4}(\mathbb{R}^{d-1} \times [0, \infty))$. Moreover, it is clear that $(u_{x_1} + \alpha u_{x_1 t})(0, \cdot, \cdot)$ is bounded in $H_{\text{loc}}^{-1/2, -1/4}(\mathbb{R}^{d-1} \times [0, \infty))$. Therefore, all the traces are defined, and we may deduce that u is a strong solution of (2.1)–(2.3). \square

Remark 2.12. We have been unable to establish that the energy loss is purely viscous as in the case of the one dimensional viscously damped wave equation on the half-line and with unilateral boundary conditions [PeS02, PeS08].

3. The evolution of a Kelvin–Voigt material with Signorini boundary conditions. As for the damped wave equation with unilateral boundary conditions, a priori estimates on the penalized problem and care relative due to the unboundedness of Ω enable us to pass to the limit in the penalized variational formulation and to deduce the existence of a solution to (1.9). Korn’s inequality plays an important role here. If we denote by \bar{u} the solution of (1.9a) with initial data (1.9d) and Dirichlet boundary data at $x_1 = 0$, then we establish that the trace $-(a_{11kl}^0 \varepsilon_{kl}(\bar{u}) + a_{11kl}^1 \varepsilon_{kl}(\bar{u}_t))|_{\Sigma \times [0, \infty)}$ increases exponentially with time in $L_{\text{loc}}^2(\Sigma \times [0, \infty))$ and not polynomially as in the case of the damped wave equation with Dirichlet boundary conditions studied in subsection 2.4. We determine the trace spaces using analogous techniques already developed in section 2.5, but here we perform a Fourier transform in the tangential variables (x_2, x_3, t) and a Laplace transform in x_1 .

3.1. The penalized problem. We approximate (1.9) as in section 2.1. More precisely, letting $r^+ \stackrel{\text{def}}{=} \max(r, 0)$, we replace u by u^ϵ which is a solution of the following penalized problem:

$$(3.1) \quad \rho u_{tt}^\epsilon - A^0 u^\epsilon - A^1 u_t^\epsilon = f, \quad x \in \Omega, \quad t > 0,$$

with initial data

$$(3.2) \quad u^\epsilon(\cdot, 0) = v_0 \quad \text{and} \quad u_t^\epsilon(\cdot, 0) = v_1$$

and boundary conditions

$$(3.3a) \quad a_{11kl}^0 \varepsilon_{kl}(u^\epsilon) + a_{11kl}^1 \varepsilon_{kl}(u_t^\epsilon) = -(u_1^\epsilon)^+ / \epsilon,$$

$$(3.3b) \quad a_{12kl}^0 \varepsilon_{kl}(u^\epsilon) + a_{12kl}^1 \varepsilon_{kl}(u_t^\epsilon) = 0, \quad \text{and} \quad a_{13kl}^0 \varepsilon_{kl}(u^\epsilon) + a_{13kl}^1 \varepsilon_{kl}(u_t^\epsilon) = 0.$$

Recall that Q_τ and I_τ were defined by (2.4).

THEOREM 3.1. *Let $W \stackrel{\text{def}}{=} \{u \in \mathbf{H}_{\text{loc}}^1([0, \infty) \times \Omega) : \nabla u_t \in \mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))\}$. Then for each $\epsilon > 0$ there exists a unique weak solution $u^\epsilon \in W$ of the problem (3.1)–(3.3) such that*

$$\begin{aligned} u^\epsilon &\in \mathbf{L}_{\text{loc}}^\infty([0, \infty); H^1(\Omega)), \\ u_t^\epsilon &\in \mathbf{L}_{\text{loc}}^2([0, \infty); H^1(\Omega)), \\ u_{tt}^\epsilon &\in \mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega)), \end{aligned}$$

and for every $\tau \in (0, T)$ and for all $v \in W$, the following variational equality is satisfied:

$$(3.4) \quad \begin{aligned} &\int_{Q_\tau} \rho u_{tt}^\epsilon \cdot v \, dx \, dt + \int_0^\tau (a^0(u^\epsilon, v) + a^1(u_t^\epsilon, v)) \, dt \\ &+ \int_{I_\tau} \frac{(u_1^\epsilon)^+}{\epsilon} v_1 \, dx' \, dt \geq \int_{Q_\tau} f \cdot v \, dx \, dt. \end{aligned}$$

Proof. We leave the verification of the proof to the reader as it is analogous to the one developed in [Jar96]. \square

3.2. Estimates on the penalized solution. We establish a priori estimates which are essential to prove the existence of a weak solution to (3.1)–(3.3). These estimates are obtained thanks to the techniques already developed in section 2.2 for the damped wave equation and to Korn’s inequality.

LEMMA 3.2. *Assume that f belongs to $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$, v_0 to $\mathbf{H}^1(\Omega)$, and v_1 to $\mathbf{L}^2(\Omega)$. Then u_t^ϵ and ∇u^ϵ are bounded in $\mathbf{L}_{\text{loc}}^\infty([0, \infty); L^2(\Omega))$, ∇u_t^ϵ is bounded in $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$, and $(u_1^\epsilon(0, \cdot, \cdot))^+ / \sqrt{\epsilon}$ is bounded in $L_{\text{loc}}^\infty([0, \infty); L^2(\mathbb{R}^{d-1}))$, independently of $\epsilon > 0$.*

Proof. These estimates are a simple application of the Gronwall lemma to the energy estimate. We multiply (3.1) by u_t^ϵ and integrate this expression over Q_τ to get

$$(3.5) \quad \begin{aligned} &\frac{1}{2} \int_\Omega (\rho |u_t^\epsilon|^2 + a_{ijkl}^0 \varepsilon_{ij}(u^\epsilon) \varepsilon_{kl}(u^\epsilon))|_0^\tau \, dx + \int_{Q_\tau} a_{ijkl}^1 \varepsilon_{ij}(u_t^\epsilon) \varepsilon_{kl}(u_t^\epsilon) \, dx \, dt \\ &+ \frac{1}{2\epsilon} \int_\Sigma ((u_1^\epsilon)^+)^2|_0^\tau \, dx' = \int_{Q_\tau} f \cdot u^\epsilon \, dx \, dt. \end{aligned}$$

According to Korn’s inequality, it is possible to infer that there exist two positive constants C_1 and C_2 such that

$$\int_\Omega a_{ijkl}^n \varepsilon_{kl}(z) \varepsilon_{ij}(z) \, dz \geq C_1 \int_\Omega |\nabla z|^2 \, dz - C_2 \int_\Omega |z|^2 \, dz, \quad n = 0, 1.$$

As $f u_t^\epsilon$ can be estimated by $|f|^2/(2\gamma) + \gamma|u_t^\epsilon|^2/2$, $\gamma > 0$, and using the above inequality, we deduce from (3.5) that

$$\begin{aligned} & \frac{1}{2} \int_{\Omega} (\rho|u_t^\epsilon|^2 + C_1|\nabla u^\epsilon|^2)(\cdot, \tau) \, dx + C_1 \int_{Q_\tau} |\nabla u_t^\epsilon|^2 \, dx \, dt + \frac{1}{2\epsilon} \int_{\Sigma} ((u_1^\epsilon)^+)|_0^\tau \, dx' \\ & \leq \frac{C_2}{2} \int_{\Omega} |u^\epsilon(\cdot, \tau)|^2 \, dx + \left(C_2 + \frac{\gamma}{2}\right) \int_{Q_\tau} |u_t^\epsilon|^2 \, dx \, dt + \frac{1}{2\gamma} \int_{Q_\tau} |f|^2 \, dx \, dt \\ & \quad + \frac{1}{2} \int_{\Omega} (\rho|v_1|^2 + a_{ijkl}^0 \varepsilon_{ij}(v_0) \varepsilon_{kl}(v_0)) \, dx. \end{aligned}$$

A classical Gronwall lemma enables us to deduce that u_t^ϵ and ∇u^ϵ are bounded in the space $\mathbf{L}_{loc}^\infty([0, \infty); L^2(\Omega))$, ∇u_t^ϵ is bounded in $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$, and $(u_1^\epsilon(0, \cdot, \cdot))^+/\sqrt{\epsilon}$ is bounded in $\mathbf{L}_{loc}^\infty([0, \infty); L^2(\mathbb{R}^{d-1}))$. \square

Remark 3.3. If we suppose that f vanishes for large t , then, independently of $\epsilon > 0$,

$$\operatorname{ess\,sup}_{0 \leq t \leq T} |u^\epsilon(\cdot, t)|_{\mathbf{H}^1} \leq C(1 + T)$$

and

$$\left(\int_0^T |u_t^\epsilon(\cdot, t)|_{\mathbf{H}^1}^2 \, dt \right)^{1/2} \leq C(1 + T).$$

These properties can be proved using the arguments given in the proof of Lemma 3.2, with the origin of time moved to T if f vanishes for $t \geq T$; since the integral involving f vanishes, the conclusion is clear.

LEMMA 3.4. *Assume that f belongs to $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$, v_0 to $\mathbf{H}^1(\Omega)$, and v_1 to $\mathbf{L}^2(\Omega)$. Then, independently of $\epsilon > 0$, the trace $(u_1^\epsilon(0, \cdot, \cdot))^+/\epsilon$ is bounded in the space of measures on I_T .*

Proof. Let φ be a cut-off function which belongs to $\mathbf{C}^1(\mathbb{R}^{d-1})$, is equal to 1 in the sphere of center 0 and radius $R > 0$, and vanishes outside of a sphere of radius $R + 1$. We multiply (3.1) by φ and we integrate over Q_τ ; due to the boundary conditions (3.3), we obtain

$$\begin{aligned} & \int_{\Omega} \rho u_t^\epsilon \cdot \varphi|_0^\tau \, dx + \frac{1}{\epsilon} \int_{I_\tau} (u_1^\epsilon)^+ \varphi_1 \, dx' \, dt + \int_{Q_\tau} \sigma_{ij}^0(u^\epsilon) \varepsilon_{ij}(\varphi) \, dx \, dt \\ & \quad + \int_{Q_\tau} \sigma_{ij}^1(u_t^\epsilon) \varepsilon_{ij}(\varphi) \, dx \, dt = \int_{Q_\tau} f \cdot \varphi \, dx \, dt. \end{aligned}$$

As the product $|zy|$ can be estimated by $|z|^2/2 + |y|^2/2$, we get the following inequality:

$$\begin{aligned} (3.6) \quad & \frac{1}{\epsilon} \int_{I_\tau} (u_1^\epsilon)^+ \varphi_1 \, dx' \, dt \leq \frac{\rho}{2} \int_{\Omega} (|u_t^\epsilon(\cdot, \tau)|^2 + |v_1|^2) \, dx + \rho \int_{\Omega} |\varphi|^2 \, dx \\ & \quad + \int_{Q_\tau} |(\sigma_{ij}^0(u^\epsilon) + \sigma_{ij}^1(u_t^\epsilon)) \varepsilon_{ij}(\varphi)| \, dx \, dt + \int_{Q_\tau} |f \cdot \varphi| \, dx \, dt. \end{aligned}$$

We may deduce that the right-hand side of (3.6) is bounded using the Lemma 3.2. Since $(u_1^\epsilon(0, \cdot, \cdot))^+$ is nonnegative, the conclusion is clear. \square

LEMMA 3.5. *Assume that f , v_0 , and v_1 belong, respectively, to $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$, $\mathbf{H}^2(\Omega)$, and $\mathbf{L}^2(\Omega)$. Then $A^0 u^\epsilon$ and $A^1 u^\epsilon$ are bounded in $\mathbf{L}_{loc}^\infty([0, \infty), L^2(\Omega))$, independently of $\epsilon > 0$.*

Proof. Once again, we use energy techniques, but now we multiply relation (3.1) by $A^1 u^\epsilon$ and we integrate over Q_τ to obtain

$$(3.7) \quad \begin{aligned} \frac{1}{2} \int_{\Omega} |A^1 u^\epsilon(\cdot, \tau)|^2 dx &= \frac{1}{2} \int_{\Omega} |A^1 v_0|^2 dx + \int_{Q_\tau} \rho u_{tt}^\epsilon \cdot (A^1 u^\epsilon) dx dt \\ &- \int_{Q_\tau} (A^0 u^\epsilon) \cdot (A^1 u^\epsilon) dx dt - \int_{Q_\tau} f \cdot (A^1 u^\epsilon) dx dt. \end{aligned}$$

Notice that

$$(3.8) \quad \begin{aligned} \int_{Q_\tau} \rho u_{tt}^\epsilon \cdot (A^1 u^\epsilon) dx dt &= \rho \int_{\Omega} u_t^\epsilon \cdot (A^1 u^\epsilon)|_0^\tau dx \\ &- \rho \int_{I_\tau} u_{1,t}^\epsilon \sigma_{1j}^1(u^\epsilon) dx' dt + \int_{Q_\tau} a_{ijkl}^1 \varepsilon_{ij}(u_t^\epsilon) \varepsilon_{kl}(u_t^\epsilon) dx dt. \end{aligned}$$

Carrying (3.8) into (3.7) and using the boundary conditions (3.3), we obtain

$$(3.9) \quad \begin{aligned} \frac{1}{2} \int_{\Omega} |A^1 u^\epsilon(\cdot, \tau)|^2 dx &= \frac{1}{2} \int_{\Omega} |A^1 v_0|^2 dx - \int_{Q_\tau} (A^0 u^\epsilon) \cdot (A^1 u^\epsilon) dx dt \\ &- \int_{Q_\tau} f \cdot (A^1 u^\epsilon) dx dt + \rho \int_{\Omega} u_t^\epsilon \cdot (A^1 u^\epsilon)|_0^\tau dx + \frac{\rho}{\epsilon} \int_{I_\tau} u_{1,t}^\epsilon (u_1^\epsilon)^+ dx' dt \\ &+ \rho \int_{I_\tau} u_{1,t}^\epsilon \sigma_{1j}^0(u^\epsilon) dx' dt + \int_{Q_\tau} a_{ijkl}^1 \varepsilon_{ij}(u_t^\epsilon) \varepsilon_{kl}(u_t^\epsilon) dx dt. \end{aligned}$$

On the other hand, we observe that

$$(3.10) \quad \int_{I_\tau} |\sigma_{1j}^0(u^\epsilon)|^2 dx' dt \leq C \left(\int_{Q_\tau} |u^\epsilon|^2 dx dt + \int_{Q_\tau} |A^1 u^\epsilon|^2 dx dt \right),$$

and for all v belonging to $\mathbf{H}^1(\Omega)$ and $A^1 v$ belonging to $\mathbf{L}^2(\Omega)$, we get

$$(3.11) \quad |A^0 v|_{L^2(\Omega)} \leq C |v|_{L^2(\Omega)} + |A^1 v|_{L^2(\Omega)}.$$

Define

$$(3.12) \quad F(t) \stackrel{\text{def}}{=} \int_{\Omega} |A^1 u^\epsilon(\cdot, t)|^2 dx.$$

According to (3.10)–(3.12) and since $u_t^\epsilon \cdot (A^1 u^\epsilon)$ can be estimated by $|u_t^\epsilon|^2/(2\gamma) + \gamma|A^1 u^\epsilon|^2/2$, $\gamma > 0$, it is possible to infer from (3.9) the following inequality:

$$\begin{aligned} \left(\frac{1}{2} - \frac{\rho\gamma}{2}\right) F(\tau) &\leq \frac{1}{2} F(0) + (2 + C) \int_0^\tau F(t) dt + \frac{1}{2} \int_{Q_\tau} |f|^2 dx dt \\ &+ \frac{\rho}{2\gamma} \int_{\Omega} |u_t^\epsilon(\cdot, \tau)|^2 dx + \rho \int_{\Omega} |v_1 \cdot (A^1 v_0)| dx + \int_{Q_\tau} a_{ijkl}^1 \varepsilon_{ij}(u_t^\epsilon) \varepsilon_{kl}(u_t^\epsilon) dx dt \\ &+ \frac{\rho}{2\epsilon} \int_{\Sigma} (u_1^\epsilon(\cdot, \tau))^+ dx' + (1 + C) \int_{Q_\tau} |u^\epsilon|^2 dx dt + \int_{I_\tau} |u_t^\epsilon|^2 dx' dt. \end{aligned}$$

If we choose γ such that $\rho\gamma < 1$, we may infer using Lemma 3.2 and a classical Gronwall inequality that F is bounded in $L_{\text{loc}}^\infty([0, \infty))$. This proves the lemma. \square

Remark 3.6. If we suppose that f vanishes for t large, then, independently of $\epsilon > 0$, $A^0 u^\epsilon$ and $A^1 u^\epsilon$ increase polynomially. These properties can be proved using the arguments given in Remark 3.3.

Let us turn now to interior estimates.

LEMMA 3.7. *Assume that f belongs to $\mathbf{L}^2_{\text{loc}}([0, \infty); \mathbf{L}^2(\Omega))$, v_0 to $\mathbf{H}^2(\Omega)$, and v_1 to $\mathbf{L}^2(\Omega)$. Then for all $\beta > 0$, u^ϵ_{tt} and $A^1 u^\epsilon_t$ are bounded in $\mathbf{L}^2([0, \infty); \mathbf{L}^2((-\infty, -\beta) \times \Sigma))$, independently of $\epsilon > 0$.*

Proof. As for the proof of Lemma 2.5, we use a truncation function which enables us to forget about the strongly nonlinear boundary conditions. More precisely, we multiply u^ϵ by a cut-off function $\varphi(x_1) \in C^\infty([0, \infty))$ which is equal to 0 on $x_1 \leq -\beta$ and to 1 on $x_1 \geq -\beta/2$, $\beta > 0$. Define

$$(3.13) \quad v^\epsilon(x_1, \cdot, \cdot) \stackrel{\text{def}}{=} \varphi(x_1)u^\epsilon(x_1, \cdot, \cdot).$$

The derivatives of v^ϵ are given by

$$(3.14a) \quad v^\epsilon_{tt} = \varphi u^\epsilon_{tt},$$

$$(3.14b) \quad \varepsilon_{kl,x_j}(v^\epsilon) = \varphi \varepsilon_{kl,x_j}(u^\epsilon) + 2\varphi_{x_1} \varepsilon_{kl}(u^\epsilon) + \varphi_{x_1 x_1} u^\epsilon_k,$$

$$(3.14c) \quad \varepsilon_{kl,x_j}(v^\epsilon_t) = \varphi \varepsilon_{kl,x_j}(u^\epsilon_t) + 2\varphi_{x_1} \varepsilon_{kl}(u^\epsilon_t) + \varphi_{x_1 x_1} u^\epsilon_{k,t}.$$

Notice that thanks to relations (3.1) and (3.14), we have

$$(3.15) \quad v^\epsilon_{tt} - A^0 v^\epsilon - A^1 v^\epsilon = \tilde{g}^\epsilon,$$

where $\tilde{g}^\epsilon \stackrel{\text{def}}{=} \varphi f_i - 2\varphi_{x_1}(a^0_{ijkl} \varepsilon_{kl}(u^\epsilon) + a^1_{ijkl} \varepsilon_{kl}(u^\epsilon_t)) - \varphi_{x_1 x_1}(a^0_{ijkl} u^\epsilon_k + a^1_{ijkl} u^\epsilon_{k,t})$. Thanks to Lemma 3.2, we deduce that \tilde{g}^ϵ is bounded in $\mathbf{L}^2_{\text{loc}}([0, \infty); \mathbf{L}^2(\Omega))$. Define

$$(3.16) \quad w^\epsilon \stackrel{\text{def}}{=} v^\epsilon_t \quad \text{and} \quad g^\epsilon \stackrel{\text{def}}{=} \tilde{g}^\epsilon + A^0 v^\epsilon.$$

We substitute (3.16) into (3.15) and obtain

$$(3.17) \quad w^\epsilon_t - A^1 w^\epsilon = g^\epsilon.$$

We will prove that w^ϵ_t is bounded in $\mathbf{L}^2_{\text{loc}}([0, \infty); \mathbf{L}^2(\Omega))$. For this purpose, we multiply (3.17) by w^ϵ_t ; we integrate this expression over Q_τ to obtain

$$\int_{Q_\tau} |w^\epsilon_t|^2 \, dx \, dt - \int_{Q_\tau} (A^1 w^\epsilon) \cdot w^\epsilon_t \, dx \, dt = \int_{Q_\tau} g^\epsilon \cdot w^\epsilon_t \, dx \, dt.$$

As

$$(3.18) \quad \int_{Q_\tau} (A^1 w^\epsilon) \cdot w^\epsilon_t \, dx \, dt = -\frac{1}{2} \int_\Omega a^1_{ijkl} \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon) \Big|_0^\tau \, dx,$$

we infer that

$$(3.19) \quad \begin{aligned} & \int_{Q_\tau} |w^\epsilon_t|^2 \, dx \, dt + \frac{1}{2} \int_\Omega a^1_{ijkl} \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon) \Big|_{t=\tau} \, dx \\ &= \frac{1}{2} \int_\Omega a^1_{ijkl} \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon) \Big|_{t=0} \, dx + \int_{Q_\tau} g^\epsilon \cdot w^\epsilon_t \, dx \, dt. \end{aligned}$$

According to Korn's inequality, we infer that there exist C_1 and C_2 such that

$$(3.20) \quad \int_\Omega a^1_{ijkl} \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon) \, dx \geq C_1 \int_\Omega |\nabla w^\epsilon|^2 \, dx - C_2 \int_\Omega |w^\epsilon|^2 \, dx.$$

Carrying the above inequality into (3.19) and observing that $g^\epsilon \cdot w_t^\epsilon$ can be estimated by $|g^\epsilon|^2/2 + |w_t^\epsilon|^2/2$, we get

$$(3.21) \quad \int_{Q_\tau} |w_t^\epsilon|^2 \, dx \, dt + C_1 \int_\Omega |\nabla w^\epsilon(\cdot, \tau)|^2 \, dx \leq \int_\Omega a_{ijkl}^1 \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon)|_{t=0} \, dx + C_2 \int_\Omega |w^\epsilon(\cdot, \tau)|^2 \, dx + \int_{Q_\tau} |g^\epsilon|^2 \, dx \, dt.$$

As v_0 belongs to $\mathbf{H}^2(\Omega)$, v_1 belongs to $\mathbf{H}^1(\Omega)$, φ belongs to $C_0^\infty(\mathbb{R})$, and g^ϵ is bounded in $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$, we infer that the right-hand side of (3.21) is bounded. Therefore, using identities (3.13) and (3.16), it is possible to deduce that u_{tt}^ϵ is bounded in $\mathbf{L}_{loc}^2([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$.

We will show that $A^1 w^\epsilon$ is bounded in $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$ using an analogous method. We multiply (3.17) by $A^1 w^\epsilon$, we integrate over Q_τ , and we obtain

$$(3.22) \quad \int_{Q_\tau} w_t^\epsilon \cdot (A^1 w^\epsilon) \, dx \, dt - \int_{Q_\tau} |A^1 w^\epsilon|^2 \, dx \, dt = \int_{Q_\tau} g^\epsilon \cdot (A^1 w^\epsilon) \, dx \, dt.$$

Carrying (3.18) and (3.20) into (3.22), $g^\epsilon \cdot (A^1 w^\epsilon)$ being estimated by $|g^\epsilon|^2/2 + |A^1 w^\epsilon|^2/2$, we obtain

$$(3.23) \quad \int_{Q_\tau} |A^1 w^\epsilon|^2 \, dx \, dt + C_1 \int_\Omega |\nabla w^\epsilon(\cdot, \tau)|^2 \, dx \leq \int_\Omega a_{ijkl}^1 \varepsilon_{ij}(w^\epsilon) \varepsilon_{kl}(w^\epsilon)|_{t=0} \, dx + C_2 \int_\Omega |w^\epsilon(\cdot, \tau)|^2 \, dx + \int_{Q_\tau} |g^\epsilon|^2 \, dx \, dt.$$

Thanks to (3.13) and (3.16), we may deduce from (3.23) that $A^1 u_t^\epsilon$ is bounded in $\mathbf{L}_{loc}^2([0, \infty); L^2((-\infty, -\beta) \times \Sigma))$. \square

3.3. Existence of a weak solution. Thanks to the estimates obtained in section 3.2, we are able to pass to the limit in the variational formulation associated to the penalized problem (3.1)–(3.3). Therefore, it is routine to deduce that there exists a solution to (1.9).

Because Ω is an unbounded set, the proof will be technical but similar to the one developed in section 2.3.

THEOREM 3.8. *Assume that f belongs to $\mathbf{L}_{loc}^2([0, \infty); L^2(\Omega))$, v_0 to $\mathbf{H}^2(\Omega)$, and v_1 to $\mathbf{L}^2(\Omega)$. Then there exists a solution to the variational inequality (1.10); this solution is the limit of a subsequence of the penalty approximation defined by (3.1)–(3.3).*

Proof. Let $\varphi \in C_0^\infty(\bar{\Omega} \times [0, \infty))$ be a function which takes its values between 0 and 1. We suppose here that v belongs to K . Multiplying (3.1) by $(v - u^\epsilon)\varphi$ and integrating over Q_τ ,

$$\int_{I_\tau} (u_1^\epsilon)^+ (\varphi(v_1 - u_1^\epsilon)) \, dx' \, dt = - \int_{I_\tau} ((u_1^\epsilon)^+)^2 \varphi \, dx' \, dt + \int_{I_\tau} (u_1^\epsilon)^+ \varphi v_1 \, dx' \, dt$$

being negative, we get the following inequality:

$$(3.24) \quad \int_\Omega \rho u_t^\epsilon \cdot (\varphi(v - u^\epsilon))|_0^\tau \, dx - \int_{Q_\tau} \rho u_t^\epsilon \cdot (\varphi(v - u^\epsilon))_t \, dx \, dt + \int_{Q_\tau} (a_{ijkl}^0 \varepsilon_{kl}(u^\epsilon) \varepsilon_{ij}(u^\epsilon) + a_{ijkl}^1 \varepsilon_{kl}(u_t^\epsilon) \varepsilon_{ij}(u^\epsilon)) (\varphi(v_i - u_i^\epsilon)) \, dx \, dt \geq \int_{Q_\tau} f \cdot (\varphi(v - u^\epsilon)) \, dx \, dt.$$

We may deduce from Lemmas 3.2 and 3.5 that there exists a subsequence, still denoted by u^ϵ , such that

$$\begin{aligned}
 (3.25a) \quad & u^\epsilon \rightharpoonup u \text{ in } \mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega)) \text{ weak } *, \\
 (3.25b) \quad & u_t^\epsilon \rightharpoonup u_t \text{ in } \mathbf{L}_{\text{loc}}^\infty([0, \infty); L^2(\Omega)) \text{ weak } *, \\
 (3.25c) \quad & \nabla u^\epsilon \rightharpoonup \nabla u \text{ in } \mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega)) \text{ weak } *, \\
 (3.25d) \quad & A^n u^\epsilon \rightharpoonup A^n u \text{ in } \mathbf{L}_{\text{loc}}^\infty([0, \infty); L^2(\Omega)) \text{ weak } *, \quad n = 0, 1, \\
 (3.25e) \quad & \nabla u_t^\epsilon \rightharpoonup \nabla u_t \text{ in } \mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega)) \text{ weak } *.
 \end{aligned}$$

Thanks to the classical compactness properties of Sobolev space injections on bounded open sets, we see that for all $R > 0$, the restrictions of u^ϵ and $a_{ijkl}^n \varepsilon_{kl}(u^\epsilon)$, $n = 0, 1$, to $Q_R = \{x : x_1 < 0, |x'| \leq R\} \times [0, R]$ (a set which has already been defined in section 2.3) converge strongly to their respective limits in $\mathbf{L}^2(Q_R)$. On the other hand, using the same techniques as those of section 2.3, we may prove that u^ϵ converges strongly to u in $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$. The complete proof can be found in [Pet02, pp. 113–115].

We observe now that since u and u_t belong to $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$, we may replace φ by φ_R in the variational inequality where φ_R is equal to 1 over the set Q_R and vanishes outside of Q_{R+1} . When R tends to infinity all the terms in (3.24) converge to their limit; thus we have proved the existence of a weak solution. \square

Remark 3.9. As for the damped wave equation with Signorini boundary conditions, the uniqueness is still an open problem.

3.4. Preliminary results. In this section, we establish estimates on the problem (1.9a) with initial data (1.9d) and the Dirichlet boundary condition which enable us to characterize the trace spaces in the next section.

LEMMA 3.10. *Assume v_0 and v_1 belong, respectively, to $\mathbf{H}^{5/2}(\Omega)$ and $\mathbf{H}^{3/2}(\Omega)$; then, there exists a function with compact support in t such that the traces of z and z_t on Σ are, respectively, v_0 and v_1 .*

Proof. We extend v_0 and v_1 into functions belonging, respectively, to $\mathbf{H}^{5/2}(\mathbb{R}^d)$ and $\mathbf{H}^{3/2}(\mathbb{R}^d)$. Then there exists a function Z belonging to $\mathbf{H}^3(\mathbb{R}^d \times [0, \infty))$ such that $Z|_{\mathbb{R}^d \times \{0\}} = v_0$ and $Z_t|_{\mathbb{R}^d \times \{0\}} = v_1$. We select a cut-off function $\varphi \in C^\infty([0, \infty))$ which is equal to 1 on $[0, 1]$ and to 0 on $[2, \infty)$, and we define z as the restriction of $\varphi(x)Z(x, t)$ to $\Omega \times [0, \infty)$. \square

LEMMA 3.11. *Assume v_0 belongs to $\mathbf{H}^{5/2}(\Omega)$, v_1 belongs to $\mathbf{H}^1(\Omega)$, and f belongs to $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$. Define z as in Lemma 3.10, and let \bar{u} be the solution of (1.9a) with initial data (1.9d) and boundary condition $\bar{u}(0, \cdot, \cdot) = z(0, \cdot, \cdot)$. Then the trace $\bar{g} \stackrel{\text{def}}{=} -(a_{11kl}^0 \varepsilon_{kl}(\bar{u}) + a_{11kl}^1 \varepsilon_{kl}(\bar{u}_t))|_{\Sigma \times [0, \infty)}$ is well defined and belongs to the space $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Sigma))$. Moreover, there exists $K > 0$ such that $e^{-Kt} \bar{g} \in L^2(\Sigma \times [0, \infty))$.*

Proof. Let $\zeta \stackrel{\text{def}}{=} \bar{u} - z$ be the solution of the following problem:

$$(3.26) \quad \rho \zeta_{tt} - A^0 \zeta - A^1 \zeta_t = F, \quad x \in \Omega, \quad t > 0,$$

where $F \stackrel{\text{def}}{=} f - \rho z_{tt} + A^0 z + A^1 z_t$ with initial data $\zeta(\cdot, 0) = \zeta_t(\cdot, 0) = 0$ and boundary condition $\zeta(0, \cdot, \cdot) = 0$. Multiplying (3.26) by ζ_t and integrating over Q_τ , Korn's inequality enables us to deduce that ζ_t and $\nabla \zeta$ are bounded in $\mathbf{L}_{\text{loc}}^\infty([0, \infty); L^2(\Omega))$ and $\nabla \zeta_t$ is bounded in $\mathbf{L}_{\text{loc}}^2([0, \infty); L^2(\Omega))$. If we multiply (3.26) by $A^1 \zeta$, we may deduce that $A^0 \zeta$ and $A^1 \zeta$ are bounded in $\mathbf{L}_{\text{loc}}^\infty([0, \infty); L^2(\Omega))$, arguing as in the proof of Lemma 3.5. On the other hand, we have

$$\int_{Q_\tau} \zeta_{tt} \cdot (A^1 \zeta_t) \, dx \, dt = -\frac{1}{2} \int_{\Omega} a^1_{ijkl} \varepsilon_{kl}(\zeta_t) \varepsilon_{ij}(\zeta_t) \Big|_0^\tau \, dx.$$

Therefore, we multiply (3.26) by $A^1 \zeta_t$, we integrate over Q_τ , and, thanks to the above identity, we get

$$\begin{aligned} & \frac{\rho}{2} \int_{\Omega} a^1_{ijkl} \varepsilon_{kl}(\zeta_{tt}) \varepsilon_{ij}(\zeta_t) \Big|_{t=\tau} \, dx + \int_{Q_\tau} |A^1 \zeta_t|^2 \, dx \, dt + \frac{1}{2} \int_{\Omega} (A^0 \zeta) \cdot (A^1 \zeta) \Big|_0^\tau \, dx \\ &= \frac{\rho}{2} \int_{\Omega} a^1_{ijkl} \varepsilon_{kl}(\zeta_t) \varepsilon_{ij}(\zeta_t) \Big|_{t=0} \, dx - \int_{Q_\tau} F \cdot (A^1 \zeta_t) \, dx \, dt. \end{aligned}$$

According to Gronwall’s lemma, there exists $K > 0$ such that

$$\int_{Q_\tau} |A^1 \zeta_t|^2 \, dx \, dt \leq C e^{K\tau} \left(|F|_{L^2(0,\tau;L^2(\Omega))}^2 + |\xi_t(\cdot, 0)|_{H^1(\Omega)}^2 + |\xi|_{L^2(0,\tau;L^2(\Omega))}^2 \right).$$

The lemma is now clear. \square

3.5. The trace spaces. We proceed as in section 2.5. A Fourier–Laplace transform and Lemma 3.11 enable us to infer that all the traces can be defined. Therefore, it is plain that a weak solution of (1.9) is also a strong one.

Let us remark first that the problem (3.3) can be written under an equivalent form: let us extend by 0 for $t \leq 0$ the difference $v^\epsilon \stackrel{\text{def}}{=} e^{-\nu t}(u^\epsilon - \bar{u})$; then it satisfies

$$(3.27) \quad \begin{aligned} & \rho(\nu + \partial_t)^2 v_i^\epsilon - ((\lambda^0 + \mu^0) + (\lambda^1 + \nu^1)(\nu + \partial_t)) \operatorname{div} v^\epsilon \\ & - (\mu^0 + \mu^1(\nu + \partial_t)) \Delta v_i^\epsilon = 0, \quad x \in \Omega, \quad t > 0, \end{aligned}$$

with boundary conditions at $\{x_1 = 0\}$

$$(3.28a) \quad (\mu^0 + \nu \mu^1)(v_{j,x_1}^\epsilon + v_{1,x_j}^\epsilon) + \mu^1(v_{j,x_1 t}^\epsilon + v_{1,x_j t}^\epsilon) = 0, \quad j = 2, 3,$$

$$(3.28b) \quad (\lambda^0 + \lambda^1(\nu + \partial_t)) \operatorname{div} v^\epsilon + 2(\mu^0 + \mu^1(\nu + \partial_t))v_{x_1}^\epsilon = e^{-\nu t} \bar{g} - \frac{(v_1^\epsilon - e^{-\nu t} \bar{u})^+}{\epsilon}$$

and with initial data

$$(3.29) \quad v^\epsilon(\cdot, 0) = 0 \quad \text{and} \quad v_t^\epsilon(\cdot, 0) = 0.$$

If v^ϵ is a tempered distribution, we may perform a Fourier transform in the tangential variable (x', t) and a Laplace transform in x_1 . Denoting by ξ and ω the dual variables of x' and t and by η the dual variable of x_1 , we are led to the system

$$(3.30) \quad \begin{aligned} & \rho(\nu + i\omega)^2 \widehat{v}^\epsilon - ((\lambda^0 + \mu^0) + (\lambda^1 + \mu^1)(\nu + i\omega)) \begin{pmatrix} \eta \\ i\xi \end{pmatrix} (\eta, i\xi^\top) \widehat{v}^\epsilon \\ & + (\mu^0 + \mu^1(\nu + i\omega))(|\xi|^2 - \eta^2) \widehat{v}^\epsilon = 0. \end{aligned}$$

Equation (3.30) is a linear system of equations; we seek its eigenvalues η_i and its eigenvectors ϕ_i :

$$(3.31a) \quad \eta_1^2 = |\xi|^2 + \frac{\rho(\nu + i\omega)^2}{\mu^0 + \mu^1(\nu + i\omega)} \quad \text{and} \quad \phi_1 = \begin{pmatrix} 0 \\ i\xi^\perp \end{pmatrix},$$

$$(3.31b) \quad \eta_2^2 = |\xi|^2 + \frac{\rho(\nu + i\omega)^2}{\mu^0 + \mu^1(\nu + i\omega)} \quad \text{and} \quad \phi_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix},$$

$$(3.31c) \quad \eta_3^2 = |\xi|^2 + \frac{\rho(\nu + i\omega)^2}{\lambda^0 + 2\mu^0 + (\lambda^1 + 2\mu^1)(\nu + i\omega)} \quad \text{and} \quad \phi_3 = \begin{pmatrix} \eta_3 \\ i\xi \end{pmatrix},$$

where ξ^\perp is obtained from ξ by a rotation of $\pi/2$. We choose η_i to be the causal determination of the square root of η_i^2 ; therefore, η_i is holomorphic in the lower half-plane $\Im(\omega) < 0$. Let us denote by \widehat{v}^ϵ the partial Fourier transform of v^ϵ with respect to the tangential variables. As v^ϵ and \widetilde{v}^ϵ are tempered distributions, \widehat{v}^ϵ is also tempered; therefore, it can include only factors of the form $e^{\eta_i x_1}$, and thus, it must be of the form

$$(3.32) \quad \widehat{v}^\epsilon(x_1, \xi, \omega) = \sum_{i=1}^3 \theta_i(\xi, \omega) \phi_i e^{\eta_i x_1}.$$

Our goal now is to determine the θ_i 's. Define $v^\epsilon \stackrel{\text{def}}{=} (v_1^\epsilon, (v^\epsilon)')$. If we apply a partial Fourier transform in the tangential variable to the boundary condition (3.28a), we obtain

$$(3.33) \quad (\widehat{v}^\epsilon)'_{x_1}(0, \xi, \omega) = -i\xi \widehat{v}_1^\epsilon(0, \xi, \omega).$$

Carrying (3.32) into (3.33), we infer that at $x_1 = 0$,

$$i\xi^\perp \eta_2 \theta_1 + i\xi \eta_3 \theta_3 = -i\xi(\theta_2 + \eta_3 \theta_3);$$

thus it is clear that $\theta_1 = 0$ and $\theta_2 = -2\eta_3 \theta_3$. Furthermore, relation (3.32) taken at $x_1 = 0$ enables us to deduce that $\theta_3 = -\widehat{v}_1^\epsilon(0, \xi, \omega)/\eta_3$. Finally, we obtain

$$(3.34) \quad \widehat{v}^\epsilon(x_1, \xi, \omega) = 2\widehat{v}_1^\epsilon(0, \xi, \omega) \phi_2 e^{\eta_2 x_1} - \widehat{v}_1^\epsilon(0, \xi, \omega) \phi_3 e^{\eta_3 x_1} / \eta_3.$$

At last, using (3.34), we can write the left-hand side of (3.28b) as a product of convolution: if we perform a Fourier transform of the left-hand side of (3.28b) and since

$$\widehat{v}_{1,x_1}^\epsilon(0, \xi, \omega) = (2\eta_2 - \eta_3) \widehat{v}_1^\epsilon(0, \xi, \omega) \quad \text{and} \quad (\widehat{v}^\epsilon)'(0, \xi, \omega) = -i\xi \widehat{v}_1^\epsilon(0, \xi, \omega) / \eta_3,$$

we obtain

$$((\lambda^0 + \lambda^1(\nu + \partial_t)) \operatorname{div} v^\epsilon + 2(\mu^0 + \mu^1(\nu + \partial_t)) v_{1,x_1}^\epsilon) \check{\gamma}(0, \xi, \omega) = \widehat{b} \widehat{v}_1^\epsilon(0, \xi, \omega),$$

where

$$\widehat{b} \stackrel{\text{def}}{=} (\lambda^0 + 2\mu^0 + (\lambda^1 + 2\mu^1)(\nu + i\omega))(2\eta_2 - \eta_3) + (\lambda^0 + \lambda^1(\nu + i\omega))|\xi|^2 / \eta_3.$$

Let $w^\epsilon(x', t)$ be the trace $v^\epsilon(0, x', t)$; then (3.28b) can now be written as

$$(3.35) \quad b * w_1^\epsilon = e^{-\nu t} \bar{g} - \frac{(w_1^\epsilon - e^{-\nu t} \bar{u}_1(0, \cdot, \cdot))^+}{\epsilon}.$$

LEMMA 3.12. *Let $u^\epsilon = (u_1^\epsilon, u_2^\epsilon, u_3^\epsilon)^\top$ be the solution of (3.1)–(3.3a). Then we may extract a subsequence, still denoted by u_1^ϵ , such that*

$$u_1^\epsilon(0, \cdot, \cdot) \rightharpoonup u_1(0, \cdot, \cdot) \quad \text{weakly in } H_{\text{loc}}^{1/2, 5/4}(\mathbb{R}^{d-1} \times [0, \infty)).$$

Moreover, u is a strong solution of (1.9).

Proof. We denote by $\widehat{\psi}$ and \widehat{g} the respective Fourier transforms of $\psi \stackrel{\text{def}}{=} \lambda^0 + 2\mu^0 + (\lambda^1 + 2\mu^1)(\nu + \partial_t)$ and $g \stackrel{\text{def}}{=} e^{-\nu t} \bar{g}$. Multiplying (3.35) by ψw_1^ϵ and using the Plancherel identity, we obtain

$$\begin{aligned} & \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{\psi} \widehat{b} |\widehat{w}_1^\epsilon|^2 \, d\xi \, d\omega \\ &= \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{g} \overline{\widehat{\psi} \widehat{w}_1^\epsilon} \, d\xi \, d\omega - \int_0^\infty \int_{\mathbb{R}^{d-1}} \frac{(w_1^\epsilon - e^{-\nu t} \bar{u}_1(0, \cdot, \cdot))^+}{\epsilon} \psi w_1^\epsilon \, dx' \, dt. \end{aligned}$$

According to the Cauchy-Schwarz inequality and since $(u_1^\epsilon(0, \cdot, \cdot))^+ / \sqrt{\epsilon}$ is bounded in $L^\infty_{\text{loc}}([0, \infty); L^2(\mathbb{R}^{d-1}))$, the absolute value of the second integral on the right-hand side of the above inequality is bounded by C_1 ; therefore, we get

$$(3.36) \quad \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{\psi} \widehat{b} |\widehat{w}_1^\epsilon|^2 \, d\xi \, d\omega \leq C_1 + \frac{1}{(2\pi)^d} \Re \int_{\mathbb{R}^d} \widehat{g} \overline{\widehat{\psi} \widehat{w}_1^\epsilon} \, d\xi \, d\omega.$$

Define

$$\begin{aligned} \kappa &\stackrel{\text{def}}{=} \frac{-\rho(\nu + i\omega)^2 - 2|\xi|^2(\mu^0 + (\nu + i\omega)\mu^1)}{-\rho(\nu + i\omega)^2 - |\xi|^2(\lambda^0 + 2\mu^0 + (\nu + i\omega)(\lambda^1 + 2\mu^1))}, \\ x_0 &\stackrel{\text{def}}{=} \sqrt{\frac{2\rho(\lambda^0 + \nu\lambda^1)}{4\lambda^1\mu^1 + (\lambda^1)^2}}. \end{aligned}$$

Then $\widehat{b} = \widehat{\psi}(2\eta_2 - \kappa\eta_3)$, and we remark also that it is sufficient to find a function h which depends on ξ and ω such that $\Re(2\eta_2 - \kappa\eta_3) \geq |h|$. If we assume that $|\xi| + |\omega| \gg 1$, we have two cases to consider according to the values taken by $|\xi|$. We suppose first that $|\xi|^2 + 2\rho(\nu\mu^1 - \mu^0)/(\mu^1)^2 \geq 0$; then η_2 can be approximated by $\widetilde{\eta}_2$ defined as follows:

$$|\widetilde{\eta}_2|^2 = \left(|\xi|^2 + \frac{\rho(\nu\mu^1 - \mu^0)}{(\mu^1)^2} \right)^2 + \left(\frac{\rho\mu^1\omega}{(\mu^1)^2} \right)^2.$$

Therefore, it is easy to deduce that $|\widetilde{\eta}_2|^2 \geq |\xi|^4/4 + \rho^2\omega^2/(\mu^1)^2$, and then, in the case $|\xi|^2 + 2\rho(\nu\mu^1 - \mu^0)/(\mu^1)^2 \geq 0$, we obtain the following estimate:

$$(3.37) \quad \Re\eta_2 \geq \cos(\pi/4)|\eta_2| \geq \frac{1}{\sqrt{2}} \left(\frac{|\xi|^4}{4} + \frac{\rho^2\omega^2}{(\mu^1)^2} \right)^{1/4}.$$

In the other case, we suppose $|\xi|^2 + 2\rho(\nu\mu^1 - \mu^0)/(\mu^1)^2 \leq 0$; then, it is plain that

$$|\Re\eta_2^2| \leq \frac{3\rho(\nu\mu^1 + \mu^0)}{(\mu^1)^2} \quad \text{and} \quad |\Im\eta_2^2| \geq \frac{\rho\mu^1|\omega|}{(\mu^1 + \nu\mu^0)^2 + (\mu^0)^2},$$

which implies that there exists $C > 0$ such that

$$|\text{arc cotan } \eta_2^2| \leq \frac{3(\nu\mu^1 + \mu^0)((\mu^1 + \nu\mu^0)^2 + (\mu^0)^2)}{(\mu^1)^3|\omega|} \leq \frac{C}{|\omega|}.$$

We deduce from the above inequality and from $|\eta_2|^2 \geq C|\omega|$ that $|\arg \eta_2^2| \leq \pi/2 + C/|\omega|$ and thus $\cos(\arg \eta_2) \geq 1/2$. In the case $|\xi|^2 + 2\rho(\nu\mu^1 - \mu^0)/(\mu^1)^2 \leq 0$, we get

$$(3.38) \quad \Re\eta_2 \geq C\sqrt{|\omega|}/2.$$

Therefore, in both cases, we infer from (3.37) and (3.38) that there exists $M > 0$ such that

$$(3.39) \quad \Re \eta_2 \geq M (\omega^2 + |\xi|^4)^{1/4}.$$

Furthermore, there exists $C > 0$ such that $|\kappa|^2 \leq 1 + C 1_{\{|\xi| \leq x_0\}} / |\omega|^2$ and for $|\xi|$ large enough, $|\eta_2| \geq |\eta_3|$. Then (3.39) enables us to deduce

$$(3.40) \quad \Re(2\eta_2 - \kappa\eta_3) \geq M (\omega^2 + |\xi|^4)^{1/4}.$$

Carrying (3.40) into (3.36), we obtain

$$(3.41) \quad M \int_{\mathbb{R}^d} |\widehat{\psi}|^2 (\omega^2 + |\xi|^4)^{1/4} |\widehat{w}_1^\varepsilon|^2 d\xi d\omega \leq C_1 + \int_{\mathbb{R}^d} |\widehat{g}| |\widehat{\psi}| |\widehat{w}_1^\varepsilon| d\xi d\omega.$$

We estimate the product zy by $|z|^2/(2\gamma) + \gamma|y|^2/2$, $\gamma > 0$, and we see that

$$\left(M - \frac{\gamma}{2}\right) \int_{\mathbb{R}^d} |\omega|^2 (\omega^2 + |\xi|^4)^{1/4} |\widehat{w}_1^\varepsilon|^2 d\xi d\omega \leq C_1 + \frac{1}{2\gamma} \int_{\mathbb{R}^d} \frac{|\widehat{g}|^2}{(\omega^2 + |\xi|^4)^{1/4}} d\xi d\omega.$$

We choose γ such that $\gamma < 2M$. On the other hand, $e^{-Kt}\bar{g}(\cdot, t)$ is bounded in $L^2(\Sigma \times [0, \infty))$, so that $g(\cdot, t)$ is bounded in $L^2(\Sigma \times [0, \infty))$ if we choose $\nu > K$. Therefore, u_1^ε is bounded in $H_{\text{loc}}^{1/2, 5/4}(\Sigma \times [0, \infty))$. In particular, $(\lambda^0 + \lambda^1(\nu + \partial_t)) \operatorname{div} v^\varepsilon + 2(\mu^0 + \mu^1(\nu + \partial_t))v_{x_1}^\varepsilon$ is bounded in $H_{\text{loc}}^{-1/2, -1/4}(\Sigma \times [0, \infty))$. We conclude that u is a strong solution of (1.9) because all the traces can be defined. \square

Acknowledgment. We would like to thank M. Chambat for reading this work and for providing helpful comments.

REFERENCES

- [DaL90] R. DAUTRAY AND J.-L. LIONS, *Mathematical Analysis and Numerical Methods for Science and Technology*, Vol. 3, Springer-Verlag, Berlin, 1990.
- [GGZ74] H. GAJEWSKI, K. GRÖGER, AND K. ZACHARIAS, *Nichtlineare Operatorgleichungen und Operatordifferentialgleichungen*, Mathematische Lehrbücher und Monographien, II. Abteilung, Mathematische Monographien, Band 38, Akademie-Verlag, Berlin, 1974.
- [Jar96] J. JARUŠEK, *Dynamic contact problems with given friction for viscoelastic bodies*, Czechoslovak Math. J., 46 (1996), pp. 475–487.
- [JM*92] J. JARUŠEK, J. MÁLEK, J. NEČAS, AND V. ŠVERÁK, *Variational inequality for a viscous drum vibrating in the presence of an obstacle*, Rend. Mat. Appl. (7), 12 (1992), pp. 943–958.
- [KuS04a] K. L. KUTTLER AND M. SHILLOR, *Dynamic contact with Signorini's condition and slip rate dependent friction*, Electron. J. Differential Equations, 83 (2004).
- [KuS04b] K. L. KUTTLER AND M. SHILLOR, *Regularity of solutions to a dynamic frictionless contact problem with normal compliance*, Nonlinear Anal., 59 (2004), pp. 1063–1075.
- [MaO88] J. A. C. MARTINS AND J. T. ODEN, *Corrigendum: Existence and uniqueness results for dynamic contact problems with nonlinear normal and friction interface laws*, Nonlinear Anal., 12 (1988), p. 747.
- [Pet02] A. PETROV, *Modélisation mathématique de procédés d'usinage: abrasion et mouillage [Mathematical Modelization of Maching Process: Abrasion and Wetting]*, Ph.D. thesis, University Claude Bernard (Lyon 1), pp. 199–2002; available online at <http://www.wias-berlin.de/people/petrov/pub.html>, 2002.
- [PeS02] A. PETROV AND M. SCHATZMAN, *Viscoélastodynamique monodimensionnelle avec conditions de Signorini*, C. R. Math. Acad. Sci. Paris, 334 (2002), pp. 983–988.
- [PeS08] A. PETROV AND M. SCHATZMAN, *A pseudodifferential linear complementarity problem related to one dimensional viscoelastic model with Signorini conditions*, Arch. Ration. Mech. Anal., to appear.

L^2 DECAY OF SOLUTIONS TO A MICRO-MACRO MODEL FOR POLYMERIC FLUIDS NEAR EQUILIBRIUM*

LINGBING HE[†] AND PING ZHANG[‡]

Abstract. In this paper, we consider the long time decay of the L^2 norm to the global solutions (u, ψ) constructed in [F.-H. Lin, C. Liu, and P. Zhang, *Comm. Pure Appl. Math.*, 60 (2007), pp. 838–866] for a micro-macro model of polymeric fluids near equilibrium $(0, M)$. Under the additional assumption that $u_0 \in H^{-\kappa}(\mathbf{R}^3)$, $(\psi_0 - M)/\sqrt{M} \in L^2(\mathbf{R}_q^3; H^{-\kappa}(\mathbf{R}_x^3))$, we prove that $(u(t), \psi(t))$ tends to $(0, M)$ as t goes to infinity with decaying rate $\|u(t)\|_{L^2} \leq C(1+t)^{-\frac{b}{2}}$ and $\|\frac{\psi(t)-M}{\sqrt{M}}\|_{L^2} \leq C(1+t)^{-\frac{b+1}{2}}$ for $b = \min(\kappa, \frac{3}{2})$. In general, without this additional assumption, we shall present an explicit long time decaying formula for $\|u(t)\|_{L^2}$ and $\|\frac{\psi(t)-M}{\sqrt{M}}\|_{L^2}$.

Key words. L^2 decay, micro-macro model, Fourier transform

AMS subject classifications. 74A25, 76D99

DOI. 10.1137/07712031

1. Introduction. In this paper, we consider the long time behavior to the global solutions constructed in [10] for a coupled microscopic-macroscopic model for polymeric fluid near equilibrium. The micro-mechanical models for polymeric liquids usually consist of beads joined by springs or rods [1, 6]. In the simplest case, a molecule configuration can be described by its end-to-end vector q . Taking into account the elastic effect together with the thermo-fluctuation, the distribution function $\psi(t, x, q)$ of molecule orientations q satisfies a Fokker–Planck equation. The convection velocity u satisfies the Navier–Stokes equations with an elastic stress which reflects the microscopic contribution of the polymer molecules to the overall macroscopic flow fields. Mathematically, this system reads (one may check [10] for a formal energetic variational derivation) as

$$(1.1) \quad \begin{cases} u_t + u \cdot \nabla u + \nabla p = \mu \Delta u + \nabla \cdot \tau, & x \in \mathbf{R}^3, \\ \nabla \cdot u = 0, & x \in \mathbf{R}^3, \\ \psi_t + u \cdot \nabla \psi = \sigma \Delta_q \psi - \nabla_q \cdot (\nabla u q \psi - \sigma \nabla_q U \psi), & (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3, \end{cases}$$

where the polymer stress τ is given by

$$(1.2) \quad \tau = \int_{\mathbf{R}^3} \nabla_q U \otimes q \psi dq,$$

with $U(q) = U(|q|^2)$ being the potential function.

Existence results for micro-macro models of polymeric fluids are usually limited to small time existence [7, 17, 13] and uniqueness of strong solutions or global existence

*Received by the editors December 29, 2007; accepted for publication (in revised form) August 15, 2008; published electronically January 14, 2009.

<http://www.siam.org/journals/sima/40-5/71203.html>

[†]Department of Mathematical Sciences, Tsinghua University, Beijing 100084, China (lbhe@math.tsinghua.edu.cn).

[‡]Academy of Mathematics & Systems Science, CAS, Beijing 100080, China (zp@amss.ac.cn). This author was partially supported by NSF of China under grant 10525101 and 10421101, National 973 project, and the innovation grant from Chinese Academy of Sciences.

of weak solutions [12]. In the setting when the last equation of (1.1) is formulated as a stochastic PDE, we refer to [7] (see also [17] for a polynomial force). Concerning the general coupled PDE system, some preliminary studies were made in the earlier work [13].

In a recent work [10], Lin, Liu, and Zhang studied the global existence of smooth solutions to (1.1) near equilibrium, which is a sort of extension to a related result of the Oldroyd model [9, 3], and which corresponds to the Hooke dumbbell model. In two space dimensions, Constantin et al. [4] proved the global existence of smooth solutions to a coupled nonlinear Fokker–Planck and Navier–Stokes system when the convection velocity u in the Fokker–Planck equation is replaced by a sort of time averaged one. Later this assumption was removed by Constantin and Masmoudi [5]. Meanwhile Lin, Zhang, and Zhang [11] independently proved the global regularity for the two-dimensional corotational FENE model with smooth initial data.

On the other hand, in [8], Jourdain et al. investigated the long time behavior of both Hookean models and FENE models for various special flows in a bounded domain with suitable boundary conditions. The main aim of this paper is to consider the long time decay rate for the global solutions to (1.1) constructed in [10].

Before we proceed, let us introduce some basic notation that will be used throughout this paper. Similarly to [10], after renormalization, we assume that $\int_{\mathbf{R}^3} \exp\{-U\} dq = 1$; furthermore,

$$(1.3) \quad \begin{aligned} &|q| \leq C(|\nabla_q U| + 1), \quad \Delta_q U \leq C + \delta|\nabla_q U|^2, \text{ with } \delta < 1, \\ &\int_{\mathbf{R}^3} |\nabla_q U|^2 e^{-U} dq \leq C, \quad \int_{\mathbf{R}^3} |q|^4 e^{-U} dq \leq C, \end{aligned}$$

and

$$(1.4) \quad \begin{aligned} &|\nabla_q^k(q\nabla_q U)| \leq C(|q||\nabla_q U| + 1), \quad \int_{\mathbf{R}^3} |\nabla_q^k(q\nabla_q U e^{-\frac{U}{2}})|^2 dq \leq C, \\ &\left| \nabla_q^k \left(\Delta_q U - \frac{|\nabla_q U|^2}{2} \right) \right| \leq C(1 + |\nabla_q U|^2), \end{aligned}$$

where the positive integer $0 \leq k \leq s$, which will be fixed in the later sections.

In what follows, for any given function $\phi(x, q)$, we denote

$$|\phi|_{H^s} = \left\{ \int_{\mathbf{R}^3} \int_{\mathbf{R}^3} \sum_{|\alpha| \leq s} |\nabla_x^\alpha \phi|^2 dq dx \right\}^{\frac{1}{2}},$$

and $\|\phi\|_{H^s}$ the standard Sobolev norm of ϕ . We shall use the convention (f, g) to stand for both the inner product on \mathbf{R}^3 , $\int_{\mathbf{R}^3} fg dx$, and on $\mathbf{R}^3 \times \mathbf{R}^3$, $\int_{\mathbf{R}^3} \int_{\mathbf{R}^3} fg dq dx$. And we will denote ∇^s to be any of ∇_x^α , where α is any multiple indices with $|\alpha| = s$. And we use similar notation for ∇_q^s .

It is easy to check that $(0, e^{-U})$ is a stationary solution to the system (1.1). In [10], Lin, Liu, and Zhang considered solutions of (1.1) with ψ being of the form

$$(1.5) \quad \psi = e^{-U} + e^{-U/2} f \stackrel{\text{def}}{=} M + \sqrt{M} f.$$

As we assume that $\int_{\mathbf{R}^3} \psi_0(x, q) dq = 1$, there holds $\int_{\mathbf{R}^3} \psi(t, x, q) dq = 1$, which together with (1.5) yields

$$(1.6) \quad \int_{\mathbf{R}^3} \sqrt{M} f dq = 0.$$

Moreover, plugging (1.5) into (1.1), we obtain the following system for (u, f) :

$$(1.7) \quad \begin{cases} u_t + u \cdot \nabla u + \nabla p = \mu \Delta u + \nabla \cdot \left(\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes q f \, dq \right), & x \in \mathbf{R}^3, \\ f_t + u \cdot \nabla f + \nabla u q \cdot \nabla_q f - \sigma \left(\Delta_q f + \frac{\Delta_q U}{2} f - \frac{|\nabla_q U|^2}{4} f \right) \\ \quad = \sqrt{M} \nabla u q \cdot \nabla_q U + \frac{1}{2} \nabla u q \cdot \nabla_q U f, & (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3, \\ \operatorname{div} u = 0, & x \in \mathbf{R}^3, \end{cases}$$

together with the initial conditions

$$(1.8) \quad u|_{t=0} = u_0, \quad f|_{t=0} = f_0, \quad \text{for } (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3,$$

where f_0 satisfies (1.6) and $\operatorname{div} u_0 = 0$.

For the reader's convenience, we first recall the following result from [10].

THEOREM 1.1. *Let $s \geq 7$ be an integer, f_0 satisfy (1.6), and $\psi_0 = M + \sqrt{M} f_0 > 0$. Then, there exists a sufficiently small constant ε such that if*

$$\begin{aligned} \frac{1}{2} \|u_0\|_{L^2}^2 + \int_{\mathbf{R}^3} \int_{\mathbf{R}^3} [\psi_0 \ln \psi_0 + U \psi_0] \, dq \, dx &\leq \varepsilon^{1+a} \mu \sigma \min(\mu, \sigma), \\ \|u_0\|_{H^s}^2 + \|f_0\|_{H^s}^2 + |q f_0|_{H^4}^2 &\leq \varepsilon \min(\mu, \sigma), \end{aligned}$$

where $a > 0$ is a small positive constant, then (1.7) and (1.8) has a unique global classical solution (u, ψ) with $\psi = M + \sqrt{M} f > 0$, and

$$(1.9) \quad \begin{aligned} &\sup_{0 \leq t < \infty} \left(\|u(t)\|_{H^s}^2 + |q f|_{H^4}^2 + \|f(t)\|_{H^s}^2 \right) \\ &+ \int_0^t \left[\mu \|\nabla u\|_{H^s}^2 + \sigma \left(\left\| |q| \left(\nabla_q f + \frac{1}{2} \nabla_q U f \right) \right\|_{H^4}^2 + \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{H^s}^2 \right) \right] \, d\tau \\ &\leq C \varepsilon \min(\mu, \sigma). \end{aligned}$$

The main aim of this paper is to prove that the solution (u, f) constructed in Theorem 1.1 decays to $(0, 0)$ as t goes to ∞ . More precisely, motivated by [14, 15] and with the additional assumption on the low frequency part of (u_0, f_0) , we obtain the following decay rates for the L^2 norm of (u, f) .

THEOREM 1.2. *Under the assumptions of Theorem 1.1, we assume further that there exists $\kappa, C > 0$ and a small enough number $c > 0$ such that*

$$(1.10) \quad \int_{|\xi| \leq c} |\xi|^{-2\kappa} |\hat{u}_0|^2 \, d\xi, \quad \int_{\mathbf{R}^3} \int_{|\xi| \leq c} |\xi|^{-2\kappa} |\mathcal{F}_x(f_0)(\xi, q)|^2 \, d\xi \, dq \leq C.$$

Then the L^2 norm of (u, f) decays to $(0, 0)$ according to

$$(1.11) \quad \|u(t)\|_{L^2} \leq C(1+t)^{-\frac{b}{2}}, \quad \|f(t)\|_{L^2} \leq C(1+t)^{-\frac{b+1}{2}},$$

with $b = \min(\kappa, \frac{3}{2})$.

In general, without the additional assumptions (1.10), to study the precise long time decay rates for (u, f) , we need some basic facts from the Littlewood–Paley theory. One may check [2] and [16] for more details.

Let $\mathcal{C} \stackrel{\text{def}}{=} \{\xi \in \mathbf{R}^3, \frac{3}{4} \leq |\xi| \leq \frac{8}{3}\}$. Let $\varphi \in C_c^\infty(\mathcal{C})$, which satisfies

$$\sum_{j \in \mathbf{Z}} \varphi(2^{-j}\xi) = 1 \quad \forall \xi \in \mathbf{R}^3 \setminus \{0\}.$$

We denote $h \stackrel{\text{def}}{=} \mathcal{F}^{-1}\varphi$; then the Littlewood–Paley operators $\dot{\Delta}_j$ and \dot{S}_j can be defined as follows:

$$\begin{aligned} \dot{\Delta}_j f &\stackrel{\text{def}}{=} \varphi(2^{-j}D)f = 2^{3j} \int_{\mathbf{R}^3} h(2^j y) f(x-y) dy \quad \text{for } j \in \mathbf{Z}, \\ (1.12) \quad \dot{S}_j f &\stackrel{\text{def}}{=} \sum_{j' \leq j-1} \dot{\Delta}_{j'} f. \end{aligned}$$

With the introduction of $\dot{\Delta}_j$, we present the following decay estimates for the solution (u, f) constructed in Theorem 1.1.

THEOREM 1.3. *Under the assumptions of Theorem 1.1, the global solutions (u, f) constructed in Theorem 1.1 decay to $(0, 0)$ by*

$$(1.13) \quad \|u(t)\|_{L^2}^2 \leq C \left[(1+t)^{-\frac{1}{2}} + \sum_{j < 0} (\|\dot{\Delta}_j u_0\|_{L^2} + \|\dot{\Delta}_j f_0\|_{L^2}) e^{-c2^{2j}t} \right],$$

$$(1.14) \quad \|f(t)\|_{L^2}^2 \leq C \left[(1+t)^{-\frac{3}{2}} + \sum_{j < 0} (\|\dot{\Delta}_j u_0\|_{L^2} + \|\dot{\Delta}_j f_0\|_{L^2}) 2^{2j} e^{-c2^{2j}t} \right]$$

for some $c > 0$. Moreover, for any $j \in \mathbf{Z}$,

$$(1.15) \quad \|(1 - \dot{S}_j)u(t)\|_{L^2}^2 \leq C_j (1+t)^{-\frac{1}{2}}, \quad \|(1 - \dot{S}_j)f(t)\|_{L^2}^2 \leq C_j (1+t)^{-\frac{3}{2}}.$$

Remark 1.1. (i) The estimates (1.13) and (1.14) can be improved when $\|u(t)\|_{L^2} \leq C(1+t)^{-\frac{1}{4}}$. The main reason lies in the fact that in this case $\|u(t)\|_{L^2}$ belongs to L^p for $p > 4$. One may check the proof of Proposition 2.1 for details. Then under the assumption (1.10), we can recover Theorem 1.2 from Theorem 1.3 by using the fact that

$$\sup_{t \geq 0} \left\{ \sum_{j \in \mathbf{Z}} (2^{2j}t)^\kappa e^{-c2^{2j}t} \right\} \leq C.$$

(ii) The estimate (1.13) implies that $\|u(t)\|_{L^2}$ tends to 0 as t goes to ∞ . While thanks to [2, 16], we have

$$\|f(t)\|_{L^2}^2 \leq \frac{\varrho(t)}{1+t},$$

where $\varrho(t) \rightarrow 0$ as $t \rightarrow \infty$.

(iii) Without (1.10), Theorem 1.3 provides a more precise decay estimate for the L^2 norm of $u(t)$ and $f(t)$. For instance, let ϕ be a smooth function with low frequency part satisfying

$$(1.16) \quad (\mathcal{F}\phi)(\xi) \Big|_{|\xi| \leq \frac{1}{2}} = |\log |\xi||^{-1} |\xi|^{-\frac{3}{2}}.$$

It is easy to check that $\phi \notin \dot{H}^{-\kappa}(\mathbf{R}^3)$ for any positive κ . While for any $j < 0$, we have

$$\|\dot{\Delta}_j \phi\|_{L^2}^2 \approx j^{-2}.$$

Let $\omega(t) \stackrel{\text{def}}{=} \sum_{j>0} j^{-2} e^{-c2^{-2j}t}$ and $\varpi(t) \stackrel{\text{def}}{=} \sum_{j>0} j^{-2} 2^{-2j} e^{-c2^{-2j}t}$; then there exist $N, C_N > 0$ such that

$$\begin{aligned} \omega(2^{2k}) &= \sum_{j>0} j^{-2} e^{-c2^{2(k-j)}} \\ &= \sum_{j>k} j^{-2} e^{-c2^{2(k-j)}} + \sum_{0<j\leq k} 2^{-2\log_2 j} e^{-c2^{2(k-j)}} \\ &\leq \int_k^\infty x^{-2} dx + k^{-2} C_N \sum_{0<j\leq k} 2^{2(\log_2 k - \log_2 j)} 2^{-2N(k-j)} \\ &\leq C_N k^{-1}. \end{aligned}$$

By the monotone property of $\omega(t)$ and $\varpi(t)$, we conclude that

$$\omega(t) \leq C(\log t)^{-1}, \quad \varpi(t) \leq C(1+t)^{-1}(\log t)^{-\frac{3}{2}}.$$

On the other hand, it is easy to construct examples so that the low frequency part of u_0 coincides with $\varepsilon\phi$, and the low frequency part of f_0 coincides with $\varepsilon\phi(x)\psi(q)$ with $\psi \in \mathcal{S}(\mathbf{R}^3)$. Then for ε sufficiently small, let (u, f) be the unique solution of (1.7) and (1.8) obtained in Theorem 1.1. The above calculations show that $u(t)$ has the logarithm decay, while $f(t)$ has a better decay rate than $(1+t)^{-1}(\log t)^{-1}$, which can be obtained by directly using Schonbek’s devices [14, 15].

(iv) Under the assumptions of Theorem 1.1, we can provide an explicit decay rate for $\|u(t)\|_{L^p}$ and $\|f(t)\|_{L^p}$ for $p > 2$ (see Corollary 2.1).

Remark 1.2. We can study the decay rates for the global smooth solutions to the Oldroyd model in [10, 3] by using exactly the same procedure as that in this paper. We omit the details here.

2. The proof of Theorem 1.2. Motivated by [14, 15], we shall present the proof of Theorem 1.2 in this section. However, compared with [14, 15], the source term in the Navier–Stokes equations of (1.7) does not decay as fast as what is required in [14, 15]. To deal with this term, we need to use the coupling effect between u and f in (1.7).

Let us first recall the following lemma from [10], which shall be of constant use in what follows.

LEMMA 2.1. *Let f satisfy $\int_{\mathbf{R}^3} f \sqrt{M} dq = 0$, and let U satisfy (1.3). Then there hold*

$$\begin{aligned} \|f\|_{L^2} &\leq C \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}, \\ \|\nabla_q U f\|_{L^2} &\leq C \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}, \\ \|\nabla_q U q f\|_{L^2} &\leq C \|(1 + |q|^2)^{\frac{1}{2}} \left(\nabla_q f + \frac{1}{2} \nabla_q U f \right)\|_{L^2}. \end{aligned}$$

Moreover, for any integer $s_1 \geq 0$ and $s_2 \geq 1$, there holds

$$\begin{aligned} \|\nabla_q U \nabla^{s_1} \nabla_q^{s_2} f\|_{L^2} &\leq C \sum_{k=0}^{s_2} \left\| \nabla_q \nabla^{s_1} \nabla_q^k f + \frac{1}{2} \nabla_q U \nabla^{s_1} \nabla_q^k f \right\|_{L^2}, \\ \|\nabla^{s_1} \nabla_q^{s_2} f\|_{L^2} &\leq C \sum_{k=0}^{s_2-1} \left\| \nabla_q \nabla^{s_1} \nabla_q^k f + \frac{1}{2} \nabla_q U \nabla^{s_1} \nabla_q^k f \right\|_{L^2}. \end{aligned}$$

As a first step in the proof of Theorem 1.2, we shall prove the following L^2 decay estimate for (u, f) .

PROPOSITION 2.1. *Under the assumptions of Theorem 1.2, the L^2 norm of the global classical solutions (u, f) to (1.7) constructed in Theorem 1.1 decay to $(0, 0)$ as t goes to ∞ according to*

$$(2.1) \quad \|u(t)\|_{L^2}^2 + \|f(t)\|_{L^2}^2 \leq C(1+t)^{-b},$$

with $b = \min(\kappa, \frac{3}{2})$.

Proof. We first get by using standard energy estimation to (1.7) that

$$\begin{aligned} \frac{d}{dt} (\|u\|_{L^2}^2 + \|f\|_{L^2}^2) + 2\mu \|\nabla u\|_{L^2}^2 + 2\sigma \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}^2 \\ = (\nabla u q \cdot \nabla_q U f, f) \leq |qf|_{H^4} \|\nabla u\|_{L^2} \|\nabla_q U f\|_{L^2}, \end{aligned}$$

which together with Lemma 2.1 and Theorem 1.1 applied gives

$$(2.2) \quad \frac{d}{dt} (\|u\|_{L^2}^2 + \|f\|_{L^2}^2) + \mu \|\nabla u\|_{L^2}^2 + \sigma \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}^2 \leq 0.$$

Thanks to (2.2), to use Schonbek’s strategy in [14, 15], we need to split the phase space into two time-dependent domains. More precisely, we decompose $\|\nabla u\|_{L^2}^2$ as

$$\|\nabla u\|_{L^2}^2 = \int_{S(t)} |\xi|^2 |\hat{u}(\xi)|^2 d\xi + \int_{S(t)^c} |\xi|^2 |\hat{u}(\xi)|^2 d\xi,$$

where $S(t) \stackrel{\text{def}}{=} \{\xi : |\xi| \leq C^{\frac{1}{2}}(1+t)^{-\frac{1}{2}}\}$ and the constant C will be chosen later on. Then, thanks to (2.2), we have

$$(2.3) \quad \frac{d}{dt} (\|u\|_{L^2}^2 + \|f\|_{L^2}^2) + \frac{C\mu}{1+t} \|u\|_{L^2}^2 + \sigma \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}^2 \leq \frac{C\mu}{1+t} \int_{S(t)} |\hat{u}(\xi)|^2 d\xi.$$

In what follows, we shall focus on the L^2 estimate to the low frequency part of u . In order to do so, we take Fourier transform with respect to x variables in (1.7) to get

$$(2.4) \quad \begin{cases} \hat{u}_t + \mathcal{F}_x(u \cdot \nabla u) + i\xi \hat{p} = -\mu |\xi|^2 \hat{u} + i\xi \cdot (\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes q \hat{f} dq), \\ \hat{f}_t + \mathcal{F}_x(u \cdot \nabla f) + \mathcal{F}_x(\nabla u q \cdot \nabla_q f) - \sigma (\Delta_q \hat{f} + \frac{\Delta_q U}{2} \hat{f} - \frac{|\nabla_q U|^2}{4} \hat{f}) \\ \quad = i(\xi_j \hat{u}_i) (\nabla_q U \otimes q \sqrt{M}) + \frac{1}{2} \mathcal{F}_x(\nabla u q \cdot \nabla_q U f), \\ \xi \cdot \hat{u} = \xi \cdot \hat{\hat{u}} = 0, \end{cases}$$

where $\hat{f}(\xi, q, t) \stackrel{\text{def}}{=} \mathcal{F}_x(f)(\xi, q, t)$. Notice that $\xi \cdot \hat{u} = \xi \cdot \tilde{u} = 0$. By multiplying the first equation of (2.4) by $\tilde{u}(t, \xi)$ and taking the resulting real part we get that

$$(2.5) \quad \frac{1}{2} \frac{d}{dt} |\hat{u}(\xi, t)|^2 + \mathcal{R}e[\mathcal{F}_x(u \cdot \nabla u) \cdot \tilde{u}(\xi, t)] \\ = -\mu|\xi|^2 |\hat{u}(\xi, t)|^2 + \mathcal{R}e \left[i\xi \otimes \tilde{u}(\xi, t) : \left(\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes q \hat{f} \, dq \right) \right].$$

A similar process applied to the microscopic equation of (2.4) gives

$$(2.6) \quad \frac{1}{2} \frac{d}{dt} \int_{\mathbf{R}^3} |\hat{f}(\xi, q, t)|^2 dq + \mathcal{R}e \left[\int_{\mathbf{R}^3} \left(\mathcal{F}_x(u \cdot \nabla f) \tilde{f}(\xi, q, t) + \mathcal{F}_x(\nabla u q \cdot \nabla_q f) \tilde{f}(\xi, q, t) \right) dq \right] \\ + \sigma \left\| \nabla_q \hat{f} + \frac{1}{2} \nabla_q U \hat{f} \right\|_{L_q^2}^2 \\ = \mathcal{R}e \left[i\xi \otimes \hat{u}(\xi, t) : \left(\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes q \tilde{f}(\xi, q, t) \, dq \right) \right] \\ + \frac{1}{2} \mathcal{R}e \left[\int_{\mathbf{R}^3} \mathcal{F}_x(\nabla u q \cdot \nabla_q U f) \tilde{f}(\xi, q, t) dq \right].$$

Thanks to (2.5)–(2.6), we obtain that for any $\delta > 0$, there is a C_δ such that

$$\frac{d}{dt} \left(|\hat{u}(\xi, t)|^2 + \int_{\mathbf{R}^3} |\hat{f}(\xi, q, t)|^2 dq \right) + \mu|\xi|^2 |\hat{u}(\xi, t)|^2 + 2\sigma \left\| \nabla_q \hat{f} + \frac{1}{2} \nabla_q U \hat{f} \right\|_{L_q^2}^2 \\ \leq C |\mathcal{F}_x(u \otimes u)|^2 + C_\delta \int_{\mathbf{R}^3} (|\mathcal{F}_x(u \cdot \nabla f)|^2 + |\mathcal{F}_x(\nabla u q \cdot \nabla_q f)|^2 \\ + |\mathcal{F}_x(\nabla u q \cdot \nabla_q U f)|^2) dq + \delta \int_{\mathbf{R}^3} |\tilde{f}(\xi, q, t)|^2 dq,$$

from which we deduce that

$$(2.7) \quad |\hat{u}(\xi, t)|^2 + \int_{\mathbf{R}^3} |\hat{f}(\xi, q, t)|^2 dq + 2\sigma \int_0^t \int_{\mathbf{R}^3} e^{-\mu|\xi|^2(t-s)} \left| \nabla_q \hat{f} + \frac{1}{2} \nabla_q U \hat{f} \right|^2 dq ds \\ \leq e^{-\mu|\xi|^2 t} \left(|\hat{u}_0(\xi)|^2 + \int_{\mathbf{R}^3} |\hat{f}_0(\xi, q)|^2 dq \right) \\ + \int_0^t \int_{\mathbf{R}^3} e^{-\mu|\xi|^2(t-s)} (\mu|\xi|^2 + \delta) |\hat{f}(\xi, q, s)|^2 dq ds \\ + C \int_0^t |\mathcal{F}_x(u \otimes u)|^2 ds + C_\delta \int_0^t \int_{\mathbf{R}^3} (|\mathcal{F}_x(u \cdot \nabla f)(\xi, q, s)|^2 \\ + |\mathcal{F}_x(\nabla u q \cdot \nabla_q f)(\xi, q, s)|^2 + |\mathcal{F}_x(\nabla u q \cdot \nabla_q U f)(\xi, q, s)|^2) dq ds.$$

Now we are in a position to estimate the right-hand side of (2.3). First, it is easy to observe that

$$\begin{aligned} & \int_0^t \int_{\mathbf{R}^3} \int_{S(t)} e^{-|\xi|^2(t-s)} (\mu|\xi|^2 + \delta) |\hat{f}(\xi, q, s)|^2 d\xi dq ds \\ & \leq \left(\frac{C\mu}{1+t} + \delta \right) \int_0^t \int_{\mathbf{R}^3} \int_{S(t)} e^{-\mu|\xi|^2(t-s)} |\hat{f}(\xi, q, s)|^2 d\xi dq ds. \end{aligned}$$

While noticing that $\int_{\mathbf{R}^3} \sqrt{M} \hat{f} dq = 0$, applying Lemma 2.1 gives

$$\begin{aligned} \int_{\mathbf{R}^3} \int_{S(t)} e^{-\mu|\xi|^2(t-s)} |\hat{f}(\xi, q, s)|^2 d\xi dq &= \int_{S(t)} e^{-\mu|\xi|^2(t-s)} \int_{\mathbf{R}^3} |\hat{f}|^2 dq d\xi \\ &\leq C \int_{\mathbf{R}^3} \int_{S(t)} e^{-\mu|\xi|^2(t-s)} \left| \nabla_q \hat{f} + \frac{1}{2} \nabla_q U \hat{f} \right|^2 d\xi dq. \end{aligned}$$

To deal with the remaining terms (2.7), we introduce $\phi \in \mathcal{S}(\mathbf{R}^3)$, the Fourier transform of which satisfies

$$\hat{\phi}(\xi) = \begin{cases} 1, & |\xi| \leq 2C^{\frac{1}{2}}, \\ 0, & |\xi| \geq 3C^{\frac{1}{2}}, \end{cases}$$

and we denote $\phi_t(x) \stackrel{\text{def}}{=} (1+t)^{-\frac{3}{2}} \phi((1+t)^{-\frac{1}{2}}x)$. Then, by using Young’s inequality and (1.9) we get that

$$\begin{aligned} \int_0^t \int_{S(t)} |\mathcal{F}_x(u \otimes u)|^2 d\xi ds &\leq \int_0^t \int_{\mathbf{R}^3} |(u \otimes u) * \phi_t|^2 dx ds \\ &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \|u \otimes u\|_{L^1}^2 ds \leq C(1+t)^{-\frac{3}{2}} \int_0^t \|u\|_{L^2}^4 ds \\ (2.8) \qquad \qquad \qquad &\leq C(1+t)^{-\frac{1}{2}}. \end{aligned}$$

While thanks to Minkowski’s inequality and (1.9), we obtain

$$\begin{aligned} \int_0^t \int_{\mathbf{R}^3} \int_{S(t)} |\mathcal{F}_x(u \cdot \nabla f)|^2 d\xi dq ds &\leq \int_0^t \int_{\mathbf{R}^3} \int_{\mathbf{R}^3} |(u \cdot \nabla f) * \phi_t|^2 dx dq ds \\ &\leq \int_0^t \int_{\mathbf{R}^3} \left| (\|u\|_{L^2_q} \|\nabla f\|_{L^2_q}) * |\phi_t| \right|^2 dx ds \leq C(1+t)^{-\frac{3}{2}} \int_0^t \left\| \|u\|_{L^2_q} \|\nabla f\|_{L^2_q} \right\|_{L^1_x}^2 ds \\ &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \|u\|_{L^2}^2 \|\nabla f\|_{L^2}^2 ds \leq C(1+t)^{-\frac{3}{2}} \int_0^t \left| \nabla_q f + \frac{1}{2} \nabla_q U f \right|_{H^4}^2 ds \\ &\leq C(1+t)^{-\frac{3}{2}}. \end{aligned}$$

A similar argument for the last two terms in (2.7) yields

$$\begin{aligned} \int_0^t \int_{\mathbf{R}^3} \int_{S(t)} |\mathcal{F}_x(\nabla u q \cdot \nabla_q f)|^2 d\xi dq ds &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \left\| |\nabla u| \|q \nabla_q f\|_{L^2_q} \right\|_{L^1_x}^2 ds \\ &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \|q \nabla_q f\|_{L^2}^2 ds \leq C(1+t)^{-\frac{3}{2}} \int_0^t \left\| (1+|q|^2)^{\frac{1}{2}} \left(\nabla_q f + \frac{1}{2} \nabla_q U f \right) \right\|_{L^2}^2 ds \\ &\leq C(1+t)^{-\frac{3}{2}} \end{aligned}$$

and

$$\int_0^t \int_{\mathbf{R}^3} \int_{S(t)} |\mathcal{F}_x(\nabla u q \cdot \nabla_q U f)|^2 d\xi dq ds \leq C(1+t)^{-\frac{3}{2}}.$$

Therefore, thanks to (2.7), we obtain

$$\begin{aligned} \int_{S(t)} |\hat{u}(\xi, t)|^2 d\xi &\leq \int_{S(t)} e^{-\mu|\xi|^2 t} |\hat{u}(\xi, 0)|^2 d\xi + \int_{S(t)} \int_{\mathbf{R}^3} e^{-\mu|\xi|^2 t} |\hat{f}(\xi, q, 0)|^2 d\xi dq \\ &\quad + C(1+t)^{-\frac{1}{2}} \end{aligned}$$

for t large enough. Then with the additional assumption (1.10), we conclude

$$(2.9) \quad \int_{S(t)} |\hat{u}(\xi, t)|^2 d\xi \leq C(1+t)^{-\kappa} + C(1+t)^{-\frac{1}{2}},$$

for t large enough. This together with (2.3) ensures that

$$\frac{d}{dt} (\|u\|_{L^2}^2 + \|f\|_{L^2}^2) + \frac{C\mu}{1+t} \|u\|_{L^2}^2 + \sigma \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}^2 \leq \frac{C\mu}{1+t} (1+t)^{-b},$$

where $b = \min(\kappa, \frac{1}{2})$, from which we deduce that

$$(2.10) \quad \|u(t)\|_{L^2}^2 + \|f(t)\|_{L^2}^2 \leq C(1+t)^{-b}.$$

With the above decay estimate for $\|u(t)\|_{L^2}$, we can improve the estimate in (2.8). Suppose that $\kappa > \frac{1}{2}$; then $b = \frac{1}{2}$ in (2.10) and

$$\begin{aligned} \int_0^t \int_{S(t)} |\mathcal{F}_x(u \otimes u)|^2 d\xi ds &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \|u\|_{L^2}^4 ds \\ &\leq C(1+t)^{-\frac{7}{6}} \left(\int_0^t \|u\|_{L^2}^6 ds \right)^{\frac{2}{3}} \\ &\leq C(1+t)^{-\frac{7}{6}}. \end{aligned}$$

Then the proof of (2.9) implies that

$$\int_{S(t)} |\hat{u}(\xi, t)|^2 d\xi \leq C(1+t)^{-\kappa} + C(1+t)^{-\frac{7}{6}},$$

for t large enough, from which, with the proof of (2.10), we deduce that

$$(2.11) \quad \|u(t)\|_{L^2}^2 + \|f(t)\|_{L^2}^2 \leq C(1+t)^{-b},$$

with $b = \min(\kappa, \frac{7}{6})$. Now if $\kappa > \frac{7}{6}$, then by using (2.11) we get that

$$\begin{aligned} \int_0^t \int_{S(t)} |\mathcal{F}_x(u \otimes u)|^2 d\xi ds &\leq C(1+t)^{-\frac{3}{2}} \int_0^t \|u\|_{L^2}^4 ds \\ &\leq C(1+t)^{-\frac{3}{2}}. \end{aligned}$$

Then the proof of (2.9) yields

$$\int_{S(t)} |\hat{u}(\xi, t)|^2 d\xi \leq C(1+t)^{-\kappa} + C(1+t)^{-\frac{3}{2}},$$

from which, with the proof of (2.10), we obtain

$$\|u(t)\|_{L^2}^2 + \|f(t)\|_{L^2}^2 \leq C(1+t)^{-b},$$

with $b = \min(\kappa, \frac{3}{2})$. This completes the proof of the proposition. \square

With Proposition 2.1, to complete the proof of Theorem 1.2, we still need to improve the decay rate for the L^2 norm of f . In order to do so, we shall first prove the L^2 decays for the derivatives of u and f , which, in particular, implies a sort of L^p decay estimate for u, f without any additional assumption on the initial data.

PROPOSITION 2.2. *Under the assumptions of Theorem 1.1, we have the following decay estimate for $\nabla u, \nabla f$, and f :*

$$(2.12) \quad \|\nabla u(t)\|_{L^2}^2 + \|\nabla f(t)\|_{L^2}^2 + \|f(t)\|_{L^2}^2 \leq C(1+t)^{-1}.$$

Proof. Taking ∇ to (1.7), by a standard energy estimation we get that

$$\begin{aligned} \frac{d}{dt} (\|\nabla u\|_{L^2}^2 + \|\nabla f\|_{L^2}^2) + 2\mu \|\nabla^2 u\|_{L^2}^2 + 2\sigma \left\| \nabla_q \nabla f + \frac{1}{2} \nabla_q U \nabla f \right\|_{L^2}^2 \\ = -2(\nabla(u \cdot \nabla u), \nabla u) - 2(\nabla(u \cdot \nabla f), \nabla f) \\ (2.13) \quad -2(\nabla(\nabla u q \cdot \nabla_q f), \nabla f) + (\nabla(\nabla u q \cdot \nabla_q U f), \nabla f) \end{aligned}$$

Thanks to (1.9), we have

$$\begin{aligned} |(\nabla(u \cdot \nabla u), \nabla u)| &= |(\nabla(u \otimes u), \nabla \nabla u)| \\ &\leq C \|u\|_{L^3} \|\nabla u\|_{L^6} \|\nabla^2 u\|_{L^2} \leq C\epsilon \|\nabla^2 u\|_{L^2}^2. \end{aligned}$$

Similarly, as $\operatorname{div} u = 0$ and $\operatorname{div}_q(\nabla u q) = 0$, we have

$$\begin{aligned} |(\nabla(u \cdot \nabla f), \nabla f)| &= |((\nabla u) \cdot \nabla f, \nabla f)| \leq \|\nabla u\|_{H^2} \|\nabla f\|_{L^2}^2, \\ |(\nabla(\nabla u q \cdot \nabla_q f), \nabla f)| &= |((\nabla \nabla u) q \cdot \nabla_q f, \nabla f)| \\ &\leq \|\nabla_q f\|_{H^4} \|\nabla^2 u\|_{L^2} \|q \nabla f\|_{L^2}, \end{aligned}$$

and

$$\begin{aligned} |(\nabla(\nabla u q \cdot \nabla_q U f), \nabla f)| \\ = |(\nabla(\nabla u) q \cdot \nabla_q U f, \nabla f) + (\nabla u q \cdot \nabla_q U(\nabla f), \nabla f)| \\ \leq \|q f\|_{H^4} \|\nabla^2 u\|_{L^2} \|\nabla_q U \nabla f\|_{L^2} + \|\nabla u\|_{H^2} \|q \nabla f\|_{L^2} \|\nabla_q U \nabla f\|_{L^2}. \end{aligned}$$

Then we deduce from (1.9), (2.13) that

$$(2.14) \quad \frac{d}{dt} (\|\nabla u\|_{L^2}^2 + \|\nabla f\|_{L^2}^2) + \mu \|\nabla^2 u\|_{L^2}^2 + \sigma \left\| \nabla_q \nabla f + \frac{1}{2} \nabla_q U \nabla f \right\|_{L^2}^2 \leq 0.$$

Then Schonbek’s strategy [14, 15] applied to (2.14) gives

$$(2.15) \quad \frac{d}{dt}(\|\nabla u\|_{L^2}^2 + \|\nabla f\|_{L^2}^2) + \frac{C\mu}{1+t}\|\nabla u\|_{L^2}^2 + \sigma \left\| \nabla_q(\nabla f) + \frac{1}{2}\nabla_q U(\nabla f) \right\|_{L^2}^2 \leq \frac{C\mu}{1+t} \int_{S(t)} |\xi|^2 |\hat{u}(\xi, t)|^2 d\xi,$$

from which we deduce that

$$(2.16) \quad \|\nabla u(t)\|_{L^2}^2 + \|\nabla f(t)\|_{L^2}^2 \leq C(1+t)^{-1}.$$

With (2.16), we now turn to the proof of the L^2 decay of f . First by a standard energy estimation we get that

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|f\|_{L^2}^2 + \sigma \left\| \nabla_q f + \frac{1}{2} \nabla_q U f \right\|_{L^2}^2 \\ & \leq \|\nabla u\|_{L^2} \|\sqrt{M} \nabla_q U\|_{L^2} \|qf\|_{L^2} + |qf|_{H^4} \|\nabla u\|_{L^2} \|\nabla_q U f\|_{L^2}. \end{aligned}$$

Then thanks to Lemma 2.1 and (1.9), we obtain

$$(2.17) \quad \frac{d}{dt} \|f\|_{L^2}^2 + \sigma \|f\|_{L^2}^2 \leq C \|\nabla u\|_{L^2}^2 \leq C(1+t)^{-1},$$

from which we deduce that

$$\begin{aligned} \|f(t)\|_{L^2}^2 & \leq C \|f_0\|_{L^2}^2 e^{-\sigma t} + \int_0^t e^{-\sigma(t-s)} (1+s)^{-1} ds. \\ & \leq C(1+t)^{-1}, \end{aligned}$$

which together with (2.16) gives (2.12). This completes the proof of Proposition 2.2. \square

Remark 2.1. Assume that $\|u(t)\|_{L^2}^2 \leq C(1+t)^{-b}$; then, thanks to (2.15), we can improve (2.16) to $\|\nabla u(t)\|_{L^2}^2 + \|\nabla f(t)\|_{L^2}^2 \leq C(1+t)^{-b-1}$, which, in turn, implies that $\|f(t)\|_{L^2}^2 \leq C(1+t)^{-b-1}$.

An immediate corollary of Proposition 2.2 and Sobolev imbedding is the following corollary.

COROLLARY 2.1. *Under the assumption of Theorem 1.1, we have*

$$\|u(t)\|_{L^p}^2 + \|f(t)\|_{L^2_q(L^p_x)}^2 \leq \begin{cases} C_p(1+t)^{-\frac{6}{p}} & \text{if } 6 \leq p \leq \infty, \\ C_p(1+t)^{-\frac{3(p-2)}{2p}} & \text{if } 2 < p \leq 6. \end{cases}$$

Now we are in a position to complete the proof of Theorem 1.2.

Proof of Theorem 1.2. Combining Proposition 2.1 and Remark 2.1, we complete the proof of the theorem. \square

3. The proof of Theorem 1.3. In this section, we consider the L^2 decay of u without the additional assumption (1.10). It is easy to observe from the proof of Theorem 1.2 that the low frequency part of initial data makes this problem difficult. To get rid of this difficulty, we are going to study the decay of solution to the linearized

problem of (1.7) first, then we estimate the difference between the solution of the linearized equation (3.1) and that of (1.7).

The linearized equation of (1.7) reads as

$$(3.1) \quad \begin{cases} v_t + \nabla h = \mu \Delta v + \nabla \cdot \left(\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes qg \, dq \right), & x \in \mathbf{R}^3, \\ g_t - \sigma \left(\Delta_q g + \frac{\Delta_q U}{2} g - \frac{|\nabla_q U|^2}{4} g \right) = \sqrt{M} \nabla v q \cdot \nabla_q U, & (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3, \\ \operatorname{div} v = 0, & x \in \mathbf{R}^3, \end{cases}$$

together with the initial condition

$$(3.2) \quad v|_{t=0} = u_0, \quad g|_{t=0} = f_0 \quad \text{for } (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3.$$

Notice from (3.1) that

$$\frac{d}{dt} \int_{\mathbf{R}^3} g \sqrt{M} dq = 0,$$

which together with the fact that $\int_{\mathbf{R}^3} f_0 \sqrt{M} dq = 0$ ensures that

$$\int_{\mathbf{R}^3} g \sqrt{M} dq = 0,$$

so that Lemma 2.1 can still be applied for g .

PROPOSITION 3.1. *Under the assumptions of Theorem 1.1, (3.1)–(3.2) has a unique solution (v, g) so that*

$$(3.3) \quad \begin{aligned} & \sup_{0 \leq t < \infty} \left(\|v(t)\|_{H^s}^2 + \|qg\|_{H^4}^2 + \|g(t)\|_{H^s}^2 \right) \\ & + \int_0^t \left[\mu \|\nabla v\|_{H^s}^2 + \sigma \left(\|q\| \left(\|\nabla_q g + \frac{1}{2} \nabla_q U g \right)\|_{H^4}^2 + \left\| \nabla_q g + \frac{1}{2} \nabla_q U g \right\|_{H^s}^2 \right) \right] d\tau \\ & \leq C\varepsilon \min(\mu, \sigma). \end{aligned}$$

Furthermore, there holds

$$(3.4) \quad \begin{aligned} & \| |D|^\lambda v \|_{L^2}^2 \leq C_\lambda (1+t)^{-\lambda}, \\ & \| |D|^\lambda g \|_{L^2}^2 + \| |D|^\lambda \nabla_q g \|_{L^2}^2 + \| |q| |D|^\lambda g \|_{L^2}^2 \leq C_\lambda (1+t)^{-\lambda-1}, \end{aligned}$$

for $\lambda \in [0, s]$, where $|D|^\lambda$ is the Fourier multiplier with the symbol $|\xi|^\lambda$.

Proof. A standard functional analysis method can be applied to prove the global existence of a smooth solution to (3.1). Furthermore, a similar proof of (1.9) in [10] ensures that (v, g) satisfies (3.3). We omit the details here.

Now let us turn to the long time behavior of (v, g) . First, let $\lambda \in \mathbf{N}$. We get by using standard energy estimation that

$$\frac{d}{dt} \left(\|\nabla^\lambda v\|_{L^2}^2 + \|\nabla^\lambda g\|_{L^2}^2 \right) + 2\mu \|\nabla \nabla^\lambda v\|_{L^2}^2 + 2\sigma \left\| \nabla_q (\nabla^\lambda g) + \frac{1}{2} \nabla_q U (\nabla^\lambda g) \right\|_{L^2}^2 = 0.$$

Schonbek’s strategy applied gives

$$\begin{aligned} & \frac{d}{dt} \left(\|\nabla^\lambda v\|_{L^2}^2 + \|\nabla^\lambda g\|_{L^2}^2 \right) + \frac{\lambda+1}{1+t} \|\nabla^\lambda v\|_{L^2}^2 + 2\sigma \left\| \nabla_q (\nabla^\lambda g) + \frac{1}{2} \nabla_q U (\nabla^\lambda g) \right\|_{L^2}^2 \\ & \leq \frac{\lambda+1}{1+t} \int_{|\xi| \leq (\lambda+1)^{\frac{1}{2}} (1+t)^{-\frac{1}{2}}} |\mathcal{F}(\nabla^\lambda v)(\xi)|^2 d\xi \leq C_\lambda (1+t)^{-\lambda-1}, \end{aligned}$$

from which we deduce that

$$(3.5) \quad \|\nabla^\lambda v\|_{L^2}^2 + \|\nabla^\lambda g\|_{L^2}^2 \leq C_\lambda(1+t)^{-\lambda}.$$

On the other hand, by applying ∇^λ and $\nabla^\lambda \nabla_q$ to the microscopic equation (3.1) and taking the L^2 inner product of the resulting equations with $\nabla^\lambda g$ and $\nabla^\lambda \nabla_q g$, respectively, we obtain

$$\frac{d}{dt} \|\nabla^\lambda g\|_{L^2}^2 + \sigma \left\| \nabla_q(\nabla^\lambda g) + \frac{1}{2} \nabla_q U(\nabla^\lambda g) \right\|_{L^2}^2 \leq C \|\nabla \nabla^\lambda v\|_{L^2}^2$$

and

$$\begin{aligned} & \frac{d}{dt} \|\nabla^\lambda \nabla_q g\|_{L^2}^2 + \sigma \left\| \nabla_q(\nabla^\lambda \nabla_q g) + \frac{1}{2} \nabla_q U(\nabla^\lambda \nabla_q g) \right\|_{L^2}^2 \\ & \leq C(\|\nabla_q U \nabla^\lambda g\|_{L^2} \|\nabla_q U \nabla^\lambda \nabla_q g\|_{L^2} + \|\nabla^\lambda g\|_{L^2} \|\nabla^\lambda \nabla_q g\|_{L^2}) \\ & \quad + \|\nabla \nabla^\lambda v\|_{L^2} \|\nabla^\lambda \nabla_q g\|_{L^2}, \end{aligned}$$

where we used (1.3) and (1.4) in the estimate of the second inequality.

In order to get the estimate for $\| |q| \nabla^\lambda g \|_{L^2}$, we apply ∇^λ to the microscopic equation (3.1) and then take the L^2 inner product of the resulting equation with $|q|^2 \nabla^\lambda g$ to get

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \| |q| \nabla^\lambda g \|_{L^2}^2 + \sigma \left\| |q| \left(\nabla_q(\nabla^\lambda g) + \frac{1}{2} \nabla_q U(\nabla^\lambda g) \right) \right\|_{L^2}^2 \\ & \leq \|\nabla^\lambda g\|_{L^2}^2 + \int_{\mathbf{R}^3 \times \mathbf{R}^3} |\nabla_q U q| |\nabla^\lambda g|^2 dq dx \\ & \quad + \|\nabla^{1+\lambda} v\|_{L^2} \left(\int_{\mathbf{R}^3} (M|q|^2) dq \right)^{\frac{1}{2}} \| |q| \nabla_q U \nabla^\lambda g \|_{L^2}. \end{aligned}$$

Therefore, by using Lemma 2.1 we get that

$$\begin{aligned} & \frac{d}{dt} \left(\|\nabla^\lambda g\|_{L^2}^2 + c \|\nabla^\lambda \nabla_q g\|_{L^2}^2 + c \| |q| \nabla^\lambda g \|_{L^2}^2 \right) + \sigma \left(\left\| \nabla_q(\nabla^\lambda g) + \frac{1}{2} \nabla_q U(\nabla^\lambda g) \right\|_{L^2}^2 \right. \\ & \quad \left. + c \left\| \nabla_q(\nabla^\lambda \nabla_q g) + \frac{1}{2} \nabla_q U(\nabla^\lambda \nabla_q g) \right\|_{L^2}^2 + c \left\| |q| \left(\nabla_q(\nabla^\lambda g) + \frac{1}{2} \nabla_q U(\nabla^\lambda g) \right) \right\|_{L^2}^2 \right) \\ & \leq C \|\nabla \nabla^\lambda v\|_{L^2}^2, \end{aligned}$$

by taking c sufficiently small, which implies

$$(3.6) \quad \|\nabla^\lambda g(t)\|_{L^2}^2 + \|\nabla^\lambda \nabla_q g(t)\|_{L^2}^2 + \| |q| \nabla^\lambda g(t) \|_{L^2}^2 \leq C_\lambda(1+t)^{-\lambda-1}.$$

With (3.5) and (3.6), we complete the proof of (3.4) via a trivial interpolation argument. \square

Remark 3.1. Plugging the second inequality of (3.4) into the v equation of (3.1) and using Schonbek’s strategy in [14, 15], we cannot obtain an improved decay estimate for v . In particular, the decay rate of $\|v(t)\|_{L^2}$ is not as good as that in [14, 15].

The main reason lies in the fact that elastic stress τ defined in (3.1) does not decay as fast as the source term in [14, 15]. A similar remark applies to (1.7) as well.

PROPOSITION 3.2. *Let (u, f) be the unique solution of (1.7) and (1.8) obtained in Theorem 1.1, and let (v, g) be the unique solution of (3.1) and (3.2) obtained in Proposition 3.1. Then there hold*

$$(3.7) \quad \begin{aligned} \|(u - v)(t)\|_{L^2}^2 &\leq C(1 + t)^{-\frac{1}{2}}, \\ \|(f - g)(t)\|_{L^2}^2 &\leq C(1 + t)^{-\frac{3}{2}}. \end{aligned}$$

Proof. Thanks to (1.7) and (3.1), we obtain

$$(3.8) \quad \left\{ \begin{aligned} &(u - v)_t + u \cdot \nabla u + \nabla(p - h) \\ &= \mu \Delta(u - v) + \nabla \cdot \left(\int_{\mathbf{R}^3} \sqrt{M} \nabla_q U \otimes q(f - g) \, dq \right), \quad x \in \mathbf{R}^3, \\ &(f - g)_t + u \cdot \nabla f + \nabla u q \cdot \nabla_q f - \sigma \left(\Delta_q(f - g) + \frac{\Delta_q U}{2}(f - g) - \frac{|\nabla_q U|^2}{4}(f - g) \right) \\ &= \sqrt{M} \nabla(u - v) q \cdot \nabla_q U + \frac{1}{2} \nabla u q \cdot \nabla_q U f, \quad (x, q) \in \mathbf{R}^3 \times \mathbf{R}^3, \\ &\operatorname{div}(u - v) = 0, \quad x \in \mathbf{R}^3. \end{aligned} \right.$$

Note that as $\operatorname{div} u = 0$ and $\operatorname{div} v = 0$, by standard energy estimates we get that

$$\begin{aligned} &\frac{d}{dt} (\|u - v\|_{L^2}^2 + \|f - g\|_{L^2}^2) + 2\mu \|\nabla(u - v)\|_{L^2}^2 + 2\sigma \left\| \nabla_q(f - g) + \frac{1}{2} \nabla_q U(f - g) \right\|_{L^2}^2 \\ &= -2(u \cdot \nabla u, u - v) - 2(u \cdot \nabla f, f - g) - 2(\nabla u q \cdot \nabla_q f, f - g) \\ &\quad + (\nabla u q \cdot \nabla_q U f, f - g) \\ &= 2(u \cdot \nabla v, u - v) + 2(u \cdot \nabla g, f - g) + 2(\nabla u q \cdot \nabla_q g, f - g) \\ &\quad + (\nabla u q \cdot \nabla_q U(f - g), f - g) + (\nabla u q \cdot \nabla_q U g, f - g). \end{aligned}$$

Thanks to (1.9) and Lemma 2.1, we have

$$\begin{aligned} |(\nabla u q \cdot \nabla_q U(f - g), f - g)| &\leq \|\nabla u\|_{H^2} \|q(f - g)\|_{L^2} \|\nabla_q U(f - g)\|_{L^2} \\ &\leq \epsilon \min\{\mu, \sigma\} \left\| \nabla_q(f - g) + \frac{1}{2} \nabla_q U(f - g) \right\|_{L^2}^2. \end{aligned}$$

Then by taking ϵ small enough we get

$$(3.9) \quad \begin{aligned} &\frac{d}{dt} (\|u - v\|_{L^2}^2 + \|f - g\|_{L^2}^2) + \mu \|\nabla(u - v)\|_{L^2}^2 + \sigma \left\| \nabla_q(f - g) + \frac{1}{2} \nabla_q U(f - g) \right\|_{L^2}^2 \\ &\leq 2|(u \cdot \nabla v, u - v)| + 2|(u \cdot \nabla g, f - g)| + 2|(\nabla u q \cdot \nabla_q g, f - g)| \\ &\quad + |(\nabla u q \cdot \nabla_q U g, f - g)|. \end{aligned}$$

In what follows, we shall deal with the right-hand side of (3.9) term by term. Our main idea is to make full use of the decay estimate to the solutions of the linear system

(3.4). First of all, thanks to (1.9), (2.12), and (3.4), we have

$$\begin{aligned} |(u \cdot \nabla v, u - v)| &\leq C \|u\|_{L^4} \|v\|_{L^4} \|\nabla(u - v)\|_{L^2} \\ &\leq C_\delta \|u\|_{\dot{H}^{\frac{3}{4}}}^2 \|v\|_{\dot{H}^{\frac{3}{4}}}^2 + \delta \|\nabla(u - v)\|_{L^2}^2 \\ &\leq C_\delta (1 + t)^{-\frac{3}{2}} + \delta \|\nabla(u - v)\|_{L^2}^2, \end{aligned}$$

for any $\delta > 0$. A similar procedure applied to the remaining terms in (3.9) gives

$$\begin{aligned} |(u \cdot \nabla g, f - g)| &\leq C \int_{\mathbf{R}^3} \|u\|_{L^4} \|f - g\|_{L_x^2} \|\nabla g\|_{L_x^4} dq \\ &\leq C \|u\|_{\dot{H}^{\frac{3}{4}}} \|f - g\|_{L^2} \|g\|_{\dot{H}^{\frac{7}{4}}} \leq C_\delta (1 + t)^{-\frac{5}{2}} + \delta \|f - g\|_{L^2}^2, \\ |(\nabla u q \cdot \nabla_q g, f - g)| &\leq C \|\nabla u\|_{\dot{H}^{\frac{3}{4}}} \|\nabla_q g\|_{\dot{H}^{\frac{3}{4}}} \|q(f - g)\|_{L^2} \\ &\leq C_\delta (1 + t)^{-2} + \delta \|q(f - g)\|_{L^2}^2, \end{aligned}$$

and

$$\begin{aligned} |(\nabla u q \cdot \nabla_q U g, f - g)| &\leq C \|\nabla u\|_{\dot{H}^{\frac{3}{4}}} \|qg\|_{\dot{H}^{\frac{3}{4}}} \|\nabla_q U(f - g)\|_{L^2} \\ &\leq C_\delta (1 + t)^{-2} + \delta \|\nabla_q U(f - g)\|_{L^2}^2. \end{aligned}$$

Then taking δ small enough, we deduce from (3.9) and Schonbek’s strategy in [14, 15] that

$$\begin{aligned} &\frac{d}{dt} (\|u - v\|_{L^2}^2 + \|f - g\|_{L^2}^2) + \frac{C\mu}{t+1} \|u - v\|_{L^2}^2 + \sigma \left\| \nabla_q(f - g) + \frac{1}{2} \nabla_q U(f - g) \right\|_{L^2}^2 \\ (3.10) \quad &\leq \frac{C\mu}{t+1} \int_{S(t)} |\mathcal{F}(u - v)(\xi)|^2 d\xi + C(1 + t)^{-\frac{3}{2}}. \end{aligned}$$

Now we repeat the procedure in section 2 to estimate the low frequency of $u - v$. First, thanks to (3.8), we have

$$\begin{aligned} &\frac{d}{dt} |\mathcal{F}(u - v)(\xi, t)|^2 + \frac{d}{dt} \int_{\mathbf{R}^3} |\mathcal{F}_x(f - g)(\xi, q, t)|^2 dq \\ &\quad + \mu |\xi|^2 |\mathcal{F}(u - v)(\xi, t)|^2 + 2\sigma \left| \nabla_q \mathcal{F}_x(f - g) + \frac{1}{2} \nabla_q U \mathcal{F}_x(f - g) \right|_{L_q^2}^2 \\ &\leq C |\mathcal{F}(u \otimes u)|^2 + C_\delta \int_{\mathbf{R}^3} (|\mathcal{F}_x(u \cdot \nabla f)|^2 + |\mathcal{F}_x(\nabla u q \cdot \nabla_q f)|^2 \\ &\quad + |\mathcal{F}_x(\nabla u q \cdot \nabla_q U f)|^2) dq + \delta \int_{\mathbf{R}^3} |\mathcal{F}(f - g)(\xi, q, t)|^2 dq, \end{aligned}$$

which implies

$$\begin{aligned}
 & |\mathcal{F}(u - v)(\xi, t)|^2 + \int_{\mathbf{R}^3} |\mathcal{F}_x(f - g)(\xi, q, t)|^2 dq \\
 & + 2\sigma \int_0^t \int_{\mathbf{R}^3} e^{-\mu|\xi|^2(t-s)} \left| \nabla_q \mathcal{F}_x(f - g) + \frac{1}{2} \nabla_q U \mathcal{F}_x(f - g) \right|^2 dq ds \\
 & \leq \int_0^t \int_{\mathbf{R}^3} e^{-|\xi|^2(t-s)} (\delta + \mu|\xi|^2) |\mathcal{F}(f - g)(\xi, q, s)|^2 dq ds + C \int_0^t |\mathcal{F}(u \otimes u)|^2 ds \\
 & + C \int_0^t \int_{\mathbf{R}^3} (|\mathcal{F}_x(u \cdot \nabla f)(\xi, q, s)|^2 + |\mathcal{F}_x(\nabla u q \cdot \nabla_q f)(\xi, q, s)|^2 \\
 & + |\mathcal{F}_x(\nabla u q \cdot \nabla_q U f)(\xi, q, s)|^2) dq ds,
 \end{aligned}$$

from which, with the proof of Proposition 2.1, we deduce that

$$(3.11) \quad \int_{S(t)} |\mathcal{F}(u - v)(\xi, t)|^2 d\xi \leq C(1 + t)^{-\frac{1}{2}}.$$

Plugging (3.11) into (3.10), we get for large enough t that

$$\frac{d}{dt} ((1 + t)^{C\mu} \|u - v\|_{L^2}^2 + (1 + t)^{C\mu} \|f - g\|_{L^2}^2) \leq C(1 + t)^{C\mu - \frac{3}{2}}.$$

Let us choose C in the definition of $S(t)$ to be large enough so that $C\mu \geq 2$. Then, we obtain

$$(3.12) \quad \|(u - v)(t)\|_{L^2}^2 + \|(f - g)(t)\|_{L^2}^2 \leq C(1 + t)^{-\frac{1}{2}}.$$

A similar procedure applied to $\nabla(u - v)$ and $\nabla(f - g)$ gives

$$(3.13) \quad \|\nabla(u - v)(t)\|_{L^2}^2 + \|\nabla(f - g)(t)\|_{L^2}^2 \leq C(1 + t)^{-\frac{3}{2}}.$$

With (3.12) and (3.13), we are in a position to prove (3.7). In fact, standard energy estimation applied to the microscopic equation of (3.8) yields

$$\begin{aligned}
 & \frac{d}{dt} \|f - g\|_{L^2}^2 + 2\sigma \left\| \nabla_q(f - g) + \frac{1}{2} \nabla_q U(f - g) \right\|_{L^2}^2 \\
 & = -2(u \cdot \nabla f, f - g) + 2(\sqrt{M} \nabla(u - v) q \cdot \nabla_q U, f - g) \\
 & \quad - 2(\nabla u q \cdot \nabla_q f, f - g) + (\nabla u q \cdot \nabla_q U f, f - g) \\
 & = 2(u \cdot \nabla g, f - g) + 2(\sqrt{M} \nabla(u - v) q \cdot \nabla_q U, f - g) + (\nabla u q \cdot \nabla_q g, f - g) \\
 & \quad + (\nabla u q \cdot \nabla_q U(f - g), f - g) + (\nabla u q \cdot \nabla_q U g, f - g),
 \end{aligned}$$

from which, with the arguments from (3.9) to (3.10), we obtain

$$\|(f - g)(t)\|_{L^2}^2 \leq C(1 + t)^{-\frac{3}{2}}.$$

This together with (3.12) completes the proof of (3.7). \square

Now we present the proof of Theorem 1.3.

Proof of Theorem 1.3. Thanks to Proposition 3.2, we only need to take care of the decay estimate of the L^2 norm of (v, g) . In order to do so, we first act $\dot{\Delta}_j$ to (3.1), then use the micro-local energy estimate to get

$$\frac{d}{dt}(\|\dot{\Delta}_j v\|_{L^2}^2 + \|\dot{\Delta}_j g\|_{L^2}^2) + 2\mu\|\nabla \dot{\Delta}_j v\|_{L^2}^2 + 2\sigma\left\|\nabla_q(\dot{\Delta}_j g) + \frac{1}{2}\nabla_q U(\dot{\Delta}_j g)\right\|_{L^2}^2 = 0,$$

from which we deduce that there exists some small positive constant c such that

$$\frac{d}{dt}(\|\dot{\Delta}_j v\|_{L^2}^2 + \|\dot{\Delta}_j g\|_{L^2}^2) + c\|\dot{\Delta}_j v\|_{L^2}^2 + c\|\dot{\Delta}_j g\|_{L^2}^2 \leq 0 \quad \text{for } j \geq 0,$$

and

$$\frac{d}{dt}(\|\dot{\Delta}_j v\|_{L^2}^2 + \|\dot{\Delta}_j g\|_{L^2}^2) + c2^{2j}\|\dot{\Delta}_j v\|_{L^2}^2 + c2^{2j}\|\dot{\Delta}_j g\|_{L^2}^2 \leq 0 \quad \text{for } j < 0.$$

Then applying Gronwall's inequality gives

$$\begin{aligned} &\|\dot{\Delta}_j v(t)\|_{L^2}^2 + \|\dot{\Delta}_j g(t)\|_{L^2}^2 \leq (\|\dot{\Delta}_j u_0\|_{L^2}^2 + \|\dot{\Delta}_j f_0\|_{L^2}^2)e^{-ct} \quad \text{for } j \geq 0, \\ (3.14) \quad &\|\dot{\Delta}_j v(t)\|_{L^2}^2 + \|\dot{\Delta}_j g(t)\|_{L^2}^2 \leq (\|\dot{\Delta}_j u_0\|_{L^2}^2 + \|\dot{\Delta}_j f_0\|_{L^2}^2)e^{-c2^{2j}t} \quad \text{for } j < 0. \end{aligned}$$

Next, let us turn to the decay estimate of $\|\Delta_j g(t)\|_{L^2}$. First by using a micro-local energy estimate to the microscopic equation of (3.1) we get that

$$\frac{d}{dt}\|\dot{\Delta}_j g\|_{L^2}^2 + 2\sigma\left\|\nabla_q(\dot{\Delta}_j g) + \frac{1}{2}\nabla_q U(\dot{\Delta}_j g)\right\|_{L^2}^2 \leq \|\nabla(\dot{\Delta}_j v)\|_{L^2}\|q\sqrt{M}\nabla_q U\|_{L^2_q}\|\dot{\Delta}_j g\|_{L^2},$$

from which, with (3.14) and (2.1), we deduce that

$$(3.15) \quad \|\dot{\Delta}_j g(t)\|_{L^2}^2 \leq C(\|\dot{\Delta}_j u_0\|_{L^2}^2 + \|\dot{\Delta}_j f_0\|_{L^2}^2)2^{2j}e^{-c2^{2j}t}.$$

On the other hand, thanks to [2, 16], we have $\|u\|_{L^2}^2 \equiv \sum_{j \in \mathbf{Z}} \|\Delta_j u\|_{L^2}^2$. Therefore, thanks to (3.7) and (3.14), we obtain

$$\begin{aligned} \|u(t)\|_{L^2}^2 &\leq C\left(\|u(t) - v(t)\|_{L^2}^2 + \sum_{j \in \mathbf{Z}} \|\Delta_j v(t)\|_{L^2}^2\right) \\ &\leq C\left[(1+t)^{-\frac{1}{2}} + \sum_{j \leq 0} (\|\dot{\Delta}_j u_0\|_{L^2}^2 + \|\dot{\Delta}_j f_0\|_{L^2}^2)e^{-c2^{2j}t}\right]. \end{aligned}$$

This proves (1.13) and the first part of (1.15). A similar argument gives (1.14) and the second part of (1.15). This completes the proof of Theorem 1.3. \square

Acknowledgments. Both authors would like to thank the support of Morning-side Center of Mathematics of the Chinese Academy of Sciences. The authors also thank the anonymous referees for valuable suggestions.

REFERENCES

[1] R. B. BIRD, C. F. CURTIS, R. C. ARMSTRONG, AND O. HASSAGER, *Dynamics of Polymeric Liquids, Volume 2: Kinetic Theory*, 2nd ed., Wiley Interscience, New York, 1987.

- [2] J.-Y. CHEMIN, *Localization in Fourier space and Navier-Stokes system*, in Phase Space Analysis of Partial Differential Equations, Vol. 1, Pubbl. Cent. Ric. Mat. Ennio Giorgi, Scuola Norm. Sup., Pisa, 2004, pp. 53–135.
- [3] Y. CHEN AND P. ZHANG, *The global existence of small solutions to the incompressible viscoelastic fluid system in general space dimensions*, Comm. Partial Differential Equations, 31 (2006), pp. 1793–1810.
- [4] P. CONSTANTIN, C. FEFFERMAN, E. TITI, AND A. ZARNESCU, *Regularity for coupled two-dimensional nonlinear Fokker-Planck and Navier-Stokes systems*, Comm. Math. Phys., 270 (2007), pp. 789–811.
- [5] P. CONSTANTIN AND N. MASMOUDI, *Global well-posedness for a Smoluchowski equation coupled with Navier-Stokes equations in 2D*, Comm. Math. Phys., 278 (2008), pp. 179–191.
- [6] M. DOI AND S. F. EDWARDS, *The Theory of Polymer Dynamics*, Oxford Science Publications, Oxford, UK, 1986.
- [7] B. JOURDAIN, T. LELIÈVRE, AND C. LE BRIS, *Existence of solution for a micro-macro model of polymeric fluid: The FENE model*, J. Funct. Anal., 209 (2004), pp. 162–193.
- [8] B. JOURDAIN, T. LELIÈVRE, C. LE BRIS, AND F. OTTO, *Long-time asymptotics of a multiscale model for polymeric fluid flows*, Arch. Ration. Mech. Anal., 181 (2006), pp. 97–148.
- [9] F.-H. LIN, C. LIU, AND P. ZHANG, *On hydrodynamics of viscoelastic fluids*, Comm. Pure Appl. Math., 58 (2005), pp. 1437–1471.
- [10] F.-H. LIN, C. LIU, AND P. ZHANG, *On a micro-macro model for polymeric fluids near equilibrium*, Comm. Pure Appl. Math., 60 (2007), pp. 838–866.
- [11] F.-H. LIN, P. ZHANG, AND Z. ZHANG, *On the global existence of smooth solution to the 2-D FENE dumbbell model*, Comm. Math. Phys., 277 (2008), pp. 531–553.
- [12] P.-L. LIONS AND N. MASMOUDI, *Global existence of weak solutions to micro-macro models*, C. R. Math. Acad. Sci. Paris, 345 (2007), pp. 15–20.
- [13] M. RENARDY, *An existence theorem for model equations resulting from kinetic theories of polymer solutions*, SIAM J. Math. Anal., 22 (1991), pp. 313–327.
- [14] M. E. SCHONBEK, *L^2 decay for weak solutions of the Navier-Stokes equations*, Arch. Rational Mech. Anal., 88 (1985), pp. 209–222.
- [15] M. E. SCHONBEK, *Large time behavior of solutions to the Navier-Stokes equations*, Comm. Partial Differential Equations, 11 (1986), pp. 733–763.
- [16] H. TRIEBEL, *Theory of Function Spaces*, Monogr. Math. 78, Birkhäuser Verlag, Basel, 1983.
- [17] E. WEINAN, T. J. LI, AND P. W. ZHANG, *Well-posedness for the dumbbell model of polymeric fluids*, Comm. Math. Phys., 248 (2004), pp. 409–427.

ANISOTROPIC INHOMOGENEOUS RECTANGULAR THIN-WALLED BEAMS*

LORENZO FREDDI[†], FRANÇOIS MURAT[‡], AND ROBERTO PARONI[§]

Abstract. This paper is devoted to the asymptotic analysis of the problem of linear elasticity for an anisotropic and inhomogeneous body occupying, in its reference configuration, a cylindrical domain with a rectangular cross section with sides proportional to ε and ε^2 and clamped on one of its bases. The sequence of solutions u^ε of the equilibrium problem is shown to converge in an appropriate topology, as ε goes to zero, to the solution of a problem for a beam in which the extensional, flexural, and torsional effects are all coupled together.

Key words. asymptotic analysis, calculus of variations, thin-walled beams, dimension reduction, variational convergence, linear elasticity

AMS subject classifications. 49J45, 74K10, 74B05

DOI. 10.1137/080720279

1. Introduction. Geometrically, a thin-walled beam is a slender structural element whose length is much larger than the diameter of the cross section which, on its hand, is larger than the thickness of the thin wall. These kinds of beams have been used for a long time in civil and mechanical engineering and, most of all, in flight vehicle structures because of their high ratio between maximum strength and weight. More recently, their importance has increased because of the introduction of fiber-reinforced composite materials in structural components. These materials are finding more and more applications for their high resistance to corrosion and high strength. Composite beams are usually made up by fiber-reinforced laminates and, hence, are anisotropic and inhomogeneous, even in cross-section planes. These peculiarities make classical thin-walled beam theories not applicable. The problem though has attracted the interest of several researchers and by now a huge number of articles can be found on the subject; see, for instance, [11] and the references therein. This is strongly remarked also in the first sentences of the abstract of [12]: “There is no lack of composite beam theories. Quite to the contrary, there might be too many of them. Different approaches, notations, etc., are used by the authors of those theories, so it is not always straightforward to compare the assumptions made and to assess the quantitative consequences of those assumptions.”

The problem under study has a huge technological interest. One very suggestive, mentioned in [2], concerns the rotor blades of helicopters. The blades are composite beams and, hence, anisotropic and inhomogeneous. The anisotropy and the inhomogeneity introduce, as we shall also deduce, structural couplings between bending, extension, and twisting behaviors. It has been observed experimentally that these

*Received by the editors April 4, 2008; accepted for publication (in revised form) August 15, 2008; published electronically January 14, 2009. The first and third authors were partially supported by the Italian Ministry of University and Research within the Project PRIN 2005.

<http://www.siam.org/journals/sima/40-5/72027.html>

[†]Dipartimento di Matematica e Informatica, Università di Udine, via delle Scienze 206, 33100 Udine, Italy (freddi@dimi.uniud.it). This author was supported by INDAM with a GNAMPA research project 2008.

[‡]Laboratoire Jacques-Louis Lions, Université Paris VI, Boîte courrier 187, 75252 Paris Cedex 05, France (murat@ann.jussieu.fr).

[§]Dipartimento di Architettura e Pianificazione, Università degli Studi di Sassari, Palazzo del Pou Salit, Piazza Duomo, 07041 Alghero, Italy (paroni@uniss.it).

couplings have a powerful influence on blade dynamics including vibrations and the aeroelastic stability; see [2]. If a model of composite beams that accurately describes the structural couplings was at our disposal, then we could try to vary the anisotropy and inhomogeneity so as to minimize undesired effects like, for instance, vibrations. Through the control of lamination parameters (ply orientation and stacking sequence), it would then be possible for industry to minimize the undesired effects.

Our aim here is to deduce a “rigorous” model for a composite thin-walled beam, that is, a inhomogeneous and anisotropic beam. We shall achieve our goal by means of well-established asymptotic methods starting from the three-dimensional linear theory of elasticity.

This paper is devoted to the asymptotic analysis of the linearized system of equilibrium equations of a body which occupies, in its reference configuration, a cylindrical domain with a rectangular cross section with sides proportional to ε and ε^2 and clamped on one of its two bases. In particular, we study the compactness properties of the sequence of solutions u^ε of the equilibrium problems and, letting ε go to zero, we are concerned with the identification of the limit problem. The same problem has been studied from the point of view of Γ -convergence: in [4] in the simpler setting of homogeneous and isotropic material and in [3] in the case of an anisotropic material which is inhomogeneous only along the longitudinal axis and subject to residual stress. Trabucho and Viano [9] also studied the same problem by superimposing two asymptotic analyses where the lengths of the two sides of the cross section go to zero independently.

Besides the material properties of the body, our treatment differs from the preceding works also in the topology used in the passage from the three-dimensional problem to the one-dimensional: the one used in the present paper delivers much more information on the deformation of the beam. The approach is close to the one developed in a recent paper of Murat and Sili [8] for a thin cylinder of radius ε . The lack of isotropy or homogeneity assumptions leads to a limit problem where the extensional, flexural, and torsional effects are coupled together. In fact, we prove that the limit problem can be written as a system of five equations in a 5-tuple of unknowns (u, v, w, p, q) (see Theorem 6.1) and that $u^\varepsilon - (u + \varepsilon v + \varepsilon^2 w + \varepsilon^3 p + \varepsilon^4 q)$ converges strongly to zero in $H^1(\Omega)$, under some regularity assumptions on $v, w, p,$ and q (see Corollary 6.1). We also derive the set of Euler equations of the variational limit problem, that is, the system of equilibrium differential equations, in the fully general case. Then we show that a strong simplification and a partial decoupling occurs when the material is homogeneous, and a complete decoupling is obtained for a homogeneous orthotropic material.

Notation. Throughout this paper $\Omega_1, \Omega_2,$ and Ω_3 will denote the following three intervals:

$$\Omega_\alpha := (-a_\alpha/2, +a_\alpha/2) \quad \text{for } \alpha = 1, 2 \quad \text{and } \Omega_3 := (0, \ell),$$

where $a_1, a_2,$ and ℓ are three positive real numbers. Also, for $i, j = 1, 2, 3$ we set

$$\Omega_{ij} := \Omega_i \times \Omega_j$$

and

$$\Omega := \Omega_1 \times \Omega_2 \times \Omega_3.$$

Unless otherwise specified, we use the Einstein summation convention. Moreover, we use the following convention for indexing vector and tensor components: Greek indices α and β take their values in the set $\{1, 2\}$ and Latin indices i, j , and k in the set $\{1, 2, 3\}$. With a little abuse of notation, and because this is a common practice and does not give rise to any mistakes, we call “sequences” even those families indicized by a continuous parameter $\varepsilon \in (0, 1)$. The component k of a vector v will be denoted either with $(v)_k$ or v_k , and an analogous notation will be used to denote tensor components. $\mathcal{E}_{\alpha\beta}$ denotes the Ricci’s symbol, that is, $\mathcal{E}_{11} = \mathcal{E}_{22} = 0$, $\mathcal{E}_{12} = 1$, and $\mathcal{E}_{21} = -1$. Since usually $x = (x_1, x_2, x_3)$, we shall then denote by $x' := (x_1, x_2)$. A wide use will be made of vector valued distributions and Sobolev spaces; for a brief account of which and for the current notation we refer the reader to the book of Le Dret [5]. Throughout this paper C will denote a constant which may change line by line.

2. The three-dimensional problem. We consider a body which occupies, in its reference configuration, the region

$$\Omega_\varepsilon := \varepsilon^2\Omega_1 \times \varepsilon\Omega_2 \times \Omega_3 \subset \mathbb{R}^3.$$

We denote by $E(u)$ the strain of the displacement u , whose components are

$$E_{ij}(u) = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right).$$

The elasticity tensor, with respect to the reference configuration Ω_ε , of the material will be denoted by \mathbb{C}^ε . We assume it to be essentially bounded,

$$\mathbb{C}^\varepsilon \in L^\infty(\Omega_\varepsilon; \mathbb{R}^{3 \times 3 \times 3 \times 3});$$

to have the minor symmetries,

$$\mathbb{C}_{ijkl}^\varepsilon = \mathbb{C}_{jikl}^\varepsilon = \mathbb{C}_{ijlk}^\varepsilon;$$

and to be positive definite. That is, there exists a constant $c > 0$ such that

$$(1) \quad \mathbb{C}^\varepsilon A \cdot A \geq c|A|^2,$$

for all three by three symmetric matrices A and for all ε . We consider the body clamped on $\Gamma_b^\varepsilon := \partial\Omega_\varepsilon \cap \{x_3 = 0\}$ and we denote by

$$H_{dn}^1(\Omega_\varepsilon; \mathbb{R}^3) := \{\varphi \in H^1(\Omega_\varepsilon; \mathbb{R}^3) : \varphi = 0 \text{ on } \Gamma_b^\varepsilon\}.$$

The weak form of the equilibrium problem can be written as

$$(2) \quad \begin{cases} \tilde{u}^\varepsilon \in H_{dn}^1(\Omega_\varepsilon; \mathbb{R}^3), \\ \int_{\Omega_\varepsilon} \mathbb{C}^\varepsilon E(\tilde{u}^\varepsilon) \cdot E(\varphi) \, dv = \int_{\Omega_\varepsilon} \tilde{F}^\varepsilon \cdot E(\varphi) \, dv \quad \forall \varphi \in H_{dn}^1(\Omega_\varepsilon; \mathbb{R}^3), \end{cases}$$

where the matrix field \tilde{F}^ε , which takes into account the presence of external forces, is assumed to be an element of $L^2(\Omega_\varepsilon; \mathbb{R}_{\text{sym}}^{3 \times 3})$.

If \tilde{F}^ε is not just in $L^2(\Omega_\varepsilon; \mathbb{R}_{\text{sym}}^{3 \times 3})$ but in

$$H(\text{div}, \Omega_\varepsilon) := \{T \in L^2(\Omega_\varepsilon; \mathbb{R}_{\text{sym}}^{3 \times 3}) : \text{div } T \in L^2(\Omega_\varepsilon; \mathbb{R}^3)\},$$

then the previous problem can be seen as the weak form of the following problem:

$$(3) \quad \begin{cases} \operatorname{div} \mathbb{C}^\varepsilon E(\tilde{u}^\varepsilon) + \tilde{b}^\varepsilon = 0 & \text{in } \Omega_\varepsilon, \\ \mathbb{C}^\varepsilon E(\tilde{u}^\varepsilon)n^\varepsilon = \tilde{c}^\varepsilon & \text{on } \Gamma_c^\varepsilon, \\ \tilde{u}^\varepsilon = 0 & \text{on } \Gamma_b^\varepsilon, \end{cases}$$

where $\Gamma_c^\varepsilon := \partial\Omega_\varepsilon \setminus \Gamma_b^\varepsilon$, and n^ε denotes the unit outward normal vector to Ω_ε , while the *body loads* \tilde{b}^ε and the *contact loads* \tilde{c}^ε are simply given by

$$(4) \quad \tilde{b}^\varepsilon := -\operatorname{div} \tilde{F}^\varepsilon \quad \text{in } \Omega_\varepsilon, \quad \tilde{c}^\varepsilon := \tilde{F}^\varepsilon n^\varepsilon \quad \text{in } \Gamma_c^\varepsilon.$$

Note that given $\tilde{b}^\varepsilon \in L^2(\Omega_\varepsilon; \mathbb{R}^3)$ and $\tilde{c}^\varepsilon \in H^{-1/2}(\partial\Omega_\varepsilon; \mathbb{R}^3)$ it is always possible to find an $\tilde{F}^\varepsilon \in H(\operatorname{div}, \Omega_\varepsilon)$ which satisfies (4).

3. The rescaled problem. It is convenient to work with the domain Ω instead of the domain Ω_ε . We therefore rescale the problem by means of the scaling map $s_\varepsilon : \Omega \rightarrow \Omega_\varepsilon$,

$$s_\varepsilon(x_1, x_2, x_3) = (\varepsilon^2 x_1, \varepsilon x_2, x_3).$$

Let \tilde{u}^ε be the solution of (2); then we define the “rescaled solution” u^ε by

$$u_1^\varepsilon := \varepsilon^2 \tilde{u}_1^\varepsilon \circ s_\varepsilon, \quad u_2^\varepsilon := \varepsilon \tilde{u}_2^\varepsilon \circ s_\varepsilon, \quad u_3^\varepsilon := \tilde{u}_3^\varepsilon \circ s_\varepsilon.$$

Let E^ε be the “rescaled strain” defined by

$$(5) \quad E^\varepsilon(\varphi) := \begin{pmatrix} \frac{1}{\varepsilon^4} E_{11}(\varphi) & \frac{1}{\varepsilon^3} E_{12}(\varphi) & \frac{1}{\varepsilon^2} E_{13}(\varphi) \\ \frac{1}{\varepsilon^3} E_{21}(\varphi) & \frac{1}{\varepsilon^2} E_{22}(\varphi) & \frac{1}{\varepsilon} E_{23}(\varphi) \\ \frac{1}{\varepsilon^2} E_{31}(\varphi) & \frac{1}{\varepsilon} E_{32}(\varphi) & E_{33}(\varphi) \end{pmatrix}.$$

It follows that $E^\varepsilon(u^\varepsilon) = E(\tilde{u}^\varepsilon) \circ s_\varepsilon$.

We further assume that there exists a $\mathbb{C} \in L^\infty(\Omega; \mathbb{R}^{3 \times 3 \times 3 \times 3})$ such that

$$\mathbb{C}^\varepsilon = \mathbb{C} \circ s_\varepsilon,$$

and we denote with $F^\varepsilon = \tilde{F}^\varepsilon \circ s_\varepsilon^{-1} \in L^2(\Omega; \mathbb{R}_{\text{sym}}^{3 \times 3})$. With this notation u^ε turns out to be the unique solution of

$$(6) \quad \begin{cases} u^\varepsilon \in H_{dn}^1(\Omega; \mathbb{R}^3), \\ \int_\Omega \mathbb{C} E^\varepsilon(u^\varepsilon) \cdot E^\varepsilon(\varphi) \, dx = \int_\Omega F^\varepsilon \cdot E^\varepsilon(\varphi) \, dx \quad \forall \varphi \in H_{dn}^1(\Omega; \mathbb{R}^3), \end{cases}$$

where $\Gamma_b := \partial\Omega \cap \{x_3 = 0\}$ and

$$H_{dn}^1(\Omega; \mathbb{R}^3) := \left\{ \varphi \in H^1(\Omega; \mathbb{R}^3) : \varphi = 0 \text{ on } \Gamma_b \right\}.$$

By taking $\varphi = u^\varepsilon$ and using (1), we find

$$(7) \quad c \|E^\varepsilon(u^\varepsilon)\|_{L^2(\Omega)} \leq \|F^\varepsilon\|_{L^2(\Omega)}.$$

Thus a uniform bound on $\|F^\varepsilon\|_{L^2(\Omega)}$ would lead to rescaled strains uniformly bounded in ε . We augment this requirement by assuming

$$(8) \quad F^\varepsilon \rightharpoonup F \quad \text{in } L^2(\Omega; \mathbb{R}_{\text{sym}}^{3 \times 3}).$$

Remark 3.1. Let $S^\varepsilon := \text{diag}(1/\varepsilon^2, 1/\varepsilon, 1)$; then $E^\varepsilon(\varphi) = S^\varepsilon E(\varphi) S^\varepsilon$. If we assume $F^\varepsilon \in H(\text{div}, \Omega)$, then we can write

$$\begin{aligned} \int_{\Omega} F^\varepsilon \cdot E^\varepsilon(\varphi) \, dx &= \int_{\Omega} S^\varepsilon F^\varepsilon S^\varepsilon \cdot E(\varphi) \, dx \\ &= - \int_{\Omega} \text{div}(S^\varepsilon F^\varepsilon S^\varepsilon) \cdot \varphi \, dx + \langle S^\varepsilon F^\varepsilon S^\varepsilon n, \varphi \rangle_{\partial\Omega}, \end{aligned}$$

for all $\varphi \in H^1_{dn}(\Omega; \mathbb{R}^3)$, and conclude that instead of considering F^ε we could have been using the following body and contact forces

$$b^\varepsilon := -\text{div}(S^\varepsilon F^\varepsilon S^\varepsilon) \quad \text{in } \Omega \quad \text{and} \quad c^\varepsilon := S^\varepsilon F^\varepsilon S^\varepsilon n \quad \text{in } \Gamma_c.$$

Note that if $F^\varepsilon = S^{\varepsilon^{-1}} F^{(0)} S^{\varepsilon^{-1}}$ for some $F^{(0)} \in H(\text{div}, \Omega)$ the body and the contact forces would be independent of ε . Since $S^{\varepsilon^{-1}} = \text{diag}(\varepsilon^2, \varepsilon, 1)$, the sequence $\{S^{\varepsilon^{-1}} F^{(0)} S^{\varepsilon^{-1}}\}$, with $F^{(0)} \in H(\text{div}, \Omega)$, strongly converges in $L^2(\Omega; \mathbb{R}^{3 \times 3}_{\text{sym}})$ and hence satisfies assumption (8). Assumption (8) though allows us to consider “stronger” forces than $F^\varepsilon = S^{\varepsilon^{-1}} F^{(0)} S^{\varepsilon^{-1}}$, like $F^\varepsilon = F^{(0)}$ or, more generally,

$$F^\varepsilon = F^{(0)} + \varepsilon F^{(1)} + \varepsilon^2 F^{(2)} + \varepsilon^3 F^{(3)} + \varepsilon^4 F^{(4)},$$

with $F^{(i)} \in H(\text{div}, \Omega)$ for $i = 0, 1, \dots, 4$, where, for instance, the term $\varepsilon^4 F^{(4)}$ would lead to the definition of body and contact forces independent of ε .

4. Partial Korn’s inequalities. In this section we state and prove several Korn’s inequalities. The proofs of Theorems 4.2 and 4.3 follow some of the lines of that of Theorem 4.4, which is due to Monneau, Murat, and Sili [7].

THEOREM 4.1. *There exists a constant C such that*

$$\left\| u_1 - \int_{\Omega_1} u_1 \, dx_1 \right\|_{L^2(\Omega)} \leq C \left\| \frac{\partial u_1}{\partial x_1} \right\|_{L^2(\Omega)}$$

for every $u_1 \in H^1(\Omega_1; L^2(\Omega_{23}))$.

Proof. By density we may restrict ourselves to considering $u_1 \in C^1(\bar{\Omega})$. For every x_2 and x_3 there exists a $\xi = \xi(x_2, x_3)$ such that $u_1(\xi, x_2, x_3) = \int_{\Omega_1} u_1(s, x_2, x_3) \, ds$ and

$$u_1(x_1, \cdot, \cdot) - u_1(\xi, \cdot, \cdot) = \int_{\xi}^{x_1} \frac{\partial u_1}{\partial x_1}(s, \cdot, \cdot) \, ds.$$

Taking squares and applying Jensen’s inequality we conclude the proof. □

THEOREM 4.2. *There exists a constant C such that*

$$\begin{aligned} \left\| u_2 - \left(\int_{\Omega_1} u_2 \, dx_1 - x_1 \frac{\partial}{\partial x_2} \int_{\Omega_1} u_1 \, dx_1 \right) \right\|_{H^{-1}(\Omega_2; L^2(\Omega_{13}))} \\ \leq C \left(\|E_{11}(u)\|_{L^2(\Omega)} + \|E_{12}(u)\|_{L^2(\Omega)} \right) \end{aligned}$$

for every $u \in H^1(\Omega; \mathbb{R}^2)$.

Proof. Let

$$\bar{u}_1 := u_1 - \int_{\Omega_1} u_1 \, dx_1, \quad \bar{u}_2 := u_2 - \int_{\Omega_1} u_2 \, dx_1 + x_1 \frac{\partial}{\partial x_2} \int_{\Omega_1} u_1 \, dx_1,$$

and note that $\int_{\Omega_1} \bar{u}_2 dx_1 = 0$ and $\bar{u}_2 \in H^1(\Omega_1; L^2(\Omega_{23}))$. Let $\psi \in H_0^1(\Omega_2)$; then

$$\begin{aligned} \frac{\partial}{\partial x_1} \int_{\Omega_2} \psi \bar{u}_2 dx_2 &= \int_{\Omega_2} \psi \left(\frac{\partial u_2}{\partial x_1} + \frac{\partial}{\partial x_2} \int_{\Omega_1} u_1 dx_1 \right) dx_2 \\ &= \int_{\Omega_2} \psi \left(2E_{12}(u) - \frac{\partial \bar{u}_1}{\partial x_2} \right) dx_2 \\ &= \int_{\Omega_2} 2\psi E_{12}(u) + \bar{u}_1 \frac{d\psi}{dx_2} dx_2. \end{aligned}$$

Since $\int_{\Omega_2} \psi \bar{u}_2 dx_2 \in H^1(\Omega_1; L^2(\Omega_{23}))$ and $\int_{\Omega_1} \int_{\Omega_2} \psi \bar{u}_2 dx_2 dx_1 = 0$, by Theorem 4.1 and the above equation we deduce

$$\begin{aligned} \left\| \int_{\Omega_2} \psi \bar{u}_2 dx_2 \right\|_{L^2(\Omega)} &\leq C \left\| \frac{\partial}{\partial x_1} \int_{\Omega_2} \psi \bar{u}_2 dx_2 \right\|_{L^2(\Omega)} \\ &\leq C \|\psi\|_{H^1(\Omega_2)} (\|E_{12}(u)\|_{L^2(\Omega)} + \|\bar{u}_1\|_{L^2(\Omega)}) \\ &\leq C \|\psi\|_{H^1(\Omega_2)} (\|E_{12}(u)\|_{L^2(\Omega)} + \|E_{11}(u)\|_{L^2(\Omega)}). \end{aligned}$$

Let $\varphi \in L^2(\Omega_{13})$; then

$$\begin{aligned} \left| \int_{\Omega} \varphi \psi \bar{u}_2 dx \right| &\leq \|\varphi\|_{L^2(\Omega_{13})} \left\| \int_{\Omega_2} \psi \bar{u}_2 dx_2 \right\|_{L^2(\Omega_{13})} \\ &\leq C \|\varphi\|_{L^2(\Omega_{13})} \|\psi\|_{H^1(\Omega_2)} (\|E_{11}(u)\|_{L^2(\Omega_{13})} + \|E_{12}(u)\|_{L^2(\Omega)}). \end{aligned}$$

A density argument concludes the proof. Indeed, let $\{\varphi_n\}$ be an orthonormal basis of $L^2(\Omega_{13})$ and for any $v \in H_0^1(\Omega_2; L^2(\Omega_{13}))$, let

$$\psi_n(x_2) := \int_{\Omega_{13}} \varphi_n(x_1, x_3) v(x) dx_1 dx_3 \in H_0^1(\Omega_2).$$

Then for $v_N := \sum_{n=1}^N \psi_n \varphi_n$ we have

$$\left| \int_{\Omega} v_N \bar{u}_2 dx \right| \leq \|v_N\|_{H_0^1(\Omega_2; L^2(\Omega_{13}))} (\|E_{11}(u)\|_{L^2(\Omega)} + \|E_{12}(u)\|_{L^2(\Omega)}),$$

and letting N go to infinity we conclude the proof. \square

Remark 4.1. In spite of the fact that the left-hand side belongs to $L^2(\Omega)$, the inequality of Theorem 4.2 does not hold true if one replaces the norm H^{-1} with the norm of L^2 , because of the following counterexample, which is inspired by an example contained in [7] in a quite similar framework.

Consider two scalar smooth functions $\varphi_\alpha \in C^\infty(\bar{\Omega}_\alpha)$ with φ_1 satisfying

$$\int_{\Omega_1} \varphi_1 dx_1 = \int_{\Omega_1} \frac{\partial \varphi_1}{\partial x_1} dx_1 = 0.$$

Define

$$u_1 := -\varphi_2 \frac{\partial \varphi_1}{\partial x_1}, \quad u_2 := \varphi_1 \frac{\partial \varphi_2}{\partial x_2}, \quad u_3 := 0.$$

Then $u \in H^1(\Omega; \mathbb{R}^3)$, and the inequality of Theorem 4.2 reduces to

$$\left\| \varphi_1 \frac{\partial \varphi_2}{\partial x_2} \right\|_{H^{-1}(\Omega_2; L^2(\Omega_1))} \leq C \left\| \varphi_2 \frac{\partial^2 \varphi_1}{\partial x_1^2} \right\|_{L^2(\Omega)},$$

which cannot be true if we replace H^{-1} by L^2 because in such a case, taking $\|\varphi_1\|_{L^2(\Omega_1)} = \|\partial^2 \varphi_1 / \partial x_1^2\|_{L^2(\Omega_1)}$ would imply that

$$\left\| \frac{\partial \varphi_2}{\partial x_2} \right\|_{L^2(\Omega)} \leq C \|\varphi_2\|_{L^2(\Omega)}$$

for any $\varphi_2 \in C^\infty(\bar{\Omega}_2)$, which is clearly impossible.

Define (using the summation convention)

$$rd_2 = \{r \in L^2(\Omega_{12}; \mathbb{R}^2) : \exists c \in \mathbb{R}, d \in \mathbb{R}^2 \text{ such that } r_\alpha(y) = \mathcal{E}_{\beta\alpha} x_\beta d + c_\alpha\},$$

where \mathcal{E} denotes the Ricci's symbol. The elements of rd_2 are the infinitesimal rigid displacements on Ω_{12} . It is easy to see that $rd_2 \subset H^1(\Omega_{12}; \mathbb{R}^2)$; moreover, being finite-dimensional, it is closed in $L^2(\Omega_{12}; \mathbb{R}^2)$. Thus, the orthogonal projection operator of $L^2(\Omega_{12}; \mathbb{R}^2)$ on rd_2 , which will be denoted by \wp , is well defined. Given a vector function $v \in L^2(\Omega_{12}; \mathbb{R}^m)$ with $m \geq 2$, we define

$$(9) \quad \vartheta(v) := \frac{1}{I_O} \int_{\Omega_{12}} (x_1 v_2 - x_2 v_1) dx', \text{ where } I_O = \int_{\Omega_{12}} (x_1^2 + x_2^2) dx'.$$

If $v \in L^2(\Omega_{12}; \mathbb{R}^2)$, then the components of \wp turn out to be

$$(10) \quad \wp_\alpha(v) = \mathcal{E}_{\beta\alpha} x_\beta \vartheta(v) + \int_{\Omega_{12}} v_\alpha dx'.$$

Furthermore, the two-dimensional Korn's inequality can be written as

$$(11) \quad \|v - \wp(v)\|_{H^1(\Omega_{12}; \mathbb{R}^2)} \leq C \sum_{\alpha, \beta} \|E_{\alpha\beta}(v)\|_{L^2(\Omega_{12}; \mathbb{R}^{2 \times 2})}$$

for all $v \in H^1(\Omega_{12}; \mathbb{R}^2)$ with a constant C which is independent of v .

Similarly, given a vector function $u \in L^2(\Omega_3; L^2(\Omega_{12}, \mathbb{R}^m))$ with $m \geq 2$, we analogously define $\vartheta(u) \in L^2(\Omega_3)$ and $\wp_\alpha(u)$. The operator \wp associates to any $u \in L^2(\Omega_3; L^2(\Omega_{12}, \mathbb{R}^m))$ a function $\wp(u) \in L^2(\Omega_3; L^2(\Omega_{12}, \mathbb{R}^2))$ which is an infinitesimal rigid displacement on Ω_{12} for almost every $x_3 \in \Omega_3$.

Let us observe that the orthogonal complement, with respect to the $L^2(\Omega_{12}; \mathbb{R}^2)$ inner product, of rd_2 in $H^1(\Omega_{12}; \mathbb{R}^2)$ can be then characterized as

$$rd_2^\perp = \{v \in H^1(\Omega_{12}; \mathbb{R}^2) : \wp(v) = 0\} = \{v \in H_m^1(\Omega_{12}; \mathbb{R}^2) : \vartheta(v) = 0\}.$$

Moreover, we denote by

$$(12) \quad RD_2^\perp(\Omega) = \{v \in L^2(\Omega_3; H_m^1(\Omega_{12}; \mathbb{R}^2)) : \vartheta(v) = 0 \text{ a.e. } x_3 \in \Omega_3\}.$$

Hereafter, for any $u \in L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^m))$, $m \geq 2$, we set

$$(13) \quad \tilde{u}_\alpha := u_\alpha - \wp_\alpha(u).$$

Of course $\tilde{u} \in L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))$ and $\int_{\Omega_{12}} \tilde{u} dx_1 dx_2 = 0$ and $\vartheta(\tilde{u}) = 0$ a.e. in Ω_3 , where the latter follows from the linearity of ϑ and the fact that $\vartheta(u) = \vartheta(\wp(u))$. Thus $\tilde{u} \in RD_2^\perp(\Omega)$.

LEMMA 4.1. *There exists a constant C such that*

$$\|\tilde{u}\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} \leq C \sum_{\alpha, \beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega_3; L^2(\Omega_{12}))}$$

for every $u \in L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^m))$, $m \geq 2$.

Proof. Since $u_\alpha(\cdot, \cdot, x_3) \in L^2(\Omega_{12})$ for almost every $x_3 \in \Omega_3$, from (10) and (11) we have that the relations

$$(14) \quad \wp_\alpha(u) = \mathcal{E}_{\beta\alpha} x_\beta \vartheta(u) + \frac{1}{|\Omega_{12}|} \int_{\Omega_{12}} u_\alpha dx',$$

$$(15) \quad \|\tilde{u}\|_{H^1(\Omega_{12}; \mathbb{R}^2)} \leq C \sum_{\alpha, \beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega_{12})}$$

hold for almost every x_3 in Ω_3 , and the claimed inequality follows by integration. \square

A different proof of the lemma above can be found in Le Dret [6].

PROPOSITION 4.1. $RD_2^\perp(\Omega)$ is a Hilbert space with the norm

$$\|v\|_{RD_2^\perp(\Omega)} := \left(\sum_{\alpha, \beta} \|E_{\alpha\beta}(v)\|_{L^2(\Omega)}^2 \right)^{1/2}.$$

Proof. We have only to prove that $\|v\|_{RD_2^\perp(\Omega)}$ is equivalent to the norm induced on $RD_2^\perp(\Omega)$ by $L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))$ since the former space is a closed subspace of the latter. For any $v \in RD_2^\perp(\Omega)$, recalling that $\wp(v) = 0$, so that $v = \tilde{v}$, and using Lemma 4.1, we then have

$$\|v\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} = \|\tilde{v}\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} \leq C \|v\|_{RD_2^\perp(\Omega)}$$

while the opposite inequality is trivially satisfied. \square

THEOREM 4.3. There exists a constant C such that

$$\begin{aligned} \left\| u_3 - \left(x_1 x_2 \frac{d\vartheta(u)}{dx_3} - x_1 \frac{d}{dx_3} \int_{\Omega_{12}} u_1 dx' + \int_{\Omega_1} u_3 dx_1 \right) \right\|_{H^{-1}(\Omega_3; L^2(\Omega_{12}))} \\ \leq C \left(\sum_{\alpha, \beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega)} + \|E_{13}(u)\|_{L^2(\Omega)} \right) \end{aligned}$$

for every $u \in H^1(\Omega; \mathbb{R}^3)$.

Proof. Let $\tilde{u}_\alpha := u_\alpha - \wp_\alpha(u)$ as in (13) and

$$\tilde{u}_3 := u_3 - \left(x_1 x_2 \frac{d\vartheta(u)}{dx_3} - x_1 \frac{d}{dx_3} \int_{\Omega_{12}} u_1 dx' + \int_{\Omega_1} u_3 dx_1 \right).$$

Since $E_{13}(u) = E_{13}(\tilde{u})$,

$$\frac{\partial \tilde{u}_3}{\partial x_1} = 2E_{13}(u) - \frac{\partial \tilde{u}_1}{\partial x_3}.$$

Let $\psi \in H_0^1(\Omega_3)$; then

$$\begin{aligned} \frac{\partial}{\partial x_1} \int_{\Omega_3} \psi \tilde{u}_3 dx_3 &= \int_{\Omega_3} \psi \left(2E_{13}(u) - \frac{\partial \tilde{u}_1}{\partial x_3} \right) dx_3 \\ &= \int_{\Omega_3} 2 \left(\psi E_{13}(u) + \tilde{u}_1 \frac{d\psi}{dx_3} \right) dx_3. \end{aligned}$$

Since $\int_{\Omega_3} \psi \tilde{u}_3 dx_3 \in H^1(\Omega)$ and $\int_{\Omega_1} \int_{\Omega_3} \psi \tilde{u}_3 dx_3 dx_1 = 0$, by Theorem 4.1, Lemma 4.1, and the previous equality we deduce

$$\begin{aligned} \left\| \int_{\Omega_3} \psi \tilde{u}_3 dx_3 \right\|_{L^2(\Omega)} &\leq C \left\| \frac{\partial}{\partial x_1} \int_0^\ell \psi \tilde{u}_3 dx_3 \right\|_{L^2(\Omega)} \\ &\leq C \|\psi\|_{H^1(\Omega_3)} (\|E_{13}(u)\|_{L^2(\Omega)} + \|\tilde{u}_1\|_{L^2(\Omega)}) \\ &\leq C \|\psi\|_{H^1(\Omega_3)} (\|E_{13}(u)\|_{L^2(\Omega)} + \sum_{\alpha,\beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega)}). \end{aligned}$$

Let $\varphi \in L^2(\Omega_{12})$; then

$$\begin{aligned} \left| \int_{\Omega} \varphi \psi \tilde{u}_3 dx \right| &= \left| \int_{\Omega_{12}} \varphi \int_{\Omega_3} \psi \tilde{u}_3 dx_3 dx' \right| \leq \|\varphi\|_{L^2(\Omega_{12})} \left\| \int_{\Omega_3} \psi \tilde{u}_3 dx_3 \right\|_{L^2(\Omega_{12})} \\ &\leq C \|\varphi\|_{L^2(\Omega_{12})} \|\psi\|_{H^1(\Omega_3)} \left(\|E_{13}(u)\|_{L^2(\Omega)} + \sum_{\alpha,\beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega)} \right). \end{aligned}$$

Arguing as in the proof of Theorem 4.2, we conclude the proof. \square

Remark 4.2. As in Remark 4.1, in spite of the fact that the left-hand side belongs to $L^2(\Omega)$, the inequality of Theorem 4.3 does not hold true if one replaces the norm H^{-1} with the norm of L^2 , because of the following counterexample.

Consider three scalar smooth functions $\varphi_i \in C^\infty(\bar{\Omega}_i)$ with φ_α satisfying

$$\int_{\Omega_\alpha} \varphi_\alpha dx_\alpha = \int_{\Omega_\alpha} \frac{\partial \varphi_\alpha}{\partial x_\alpha} dx_\alpha = 0.$$

Define

$$u_1 := -\varphi_2 \frac{\partial \varphi_1}{\partial x_1} \varphi_3, \quad u_2 := -\varphi_1 \frac{\partial \varphi_2}{\partial x_2} \varphi_3, \quad u_3 := +\varphi_1 \varphi_2 \frac{\partial \varphi_3}{\partial x_3}.$$

Then $u \in H^1(\Omega; \mathbb{R}^3)$, and the inequality of Theorem 4.3 reduces to

$$\left\| \varphi_1 \varphi_2 \frac{\partial \varphi_3}{\partial x_3} \right\|_{H^{-1}(\Omega_2; L^2(\Omega_1))} \leq C \left\| \varphi_3 \left(\frac{\partial^2 \varphi_1}{\partial x_1^2} \varphi_2 + \varphi_1 \frac{\partial^2 \varphi_2}{\partial x_2^2} \right) \right\|_{L^2(\Omega)},$$

which cannot be true if we replace H^{-1} by L^2 , because in such a case taking $\|\varphi_\alpha\|_{L^2(\Omega_\alpha)} = \|\partial^2 \varphi_\alpha / \partial x_1^\alpha\|_{L^2(\Omega_\alpha)}$ would imply that

$$\left\| \frac{\partial \varphi_3}{\partial x_3} \right\|_{L^2(\Omega)} \leq C \|\varphi_3\|_{L^2(\Omega)}$$

for any $\varphi_3 \in C^\infty(\bar{\Omega}_3)$, which is clearly impossible.

The next partial Korn's inequality is proved in Monneau, Murat, and Sili [7].

THEOREM 4.4. *There exists a constant C such that*

$$\begin{aligned} \left\| u_3 - \left(\int_{\Omega_{12}} u_3 dx' - x_\alpha \frac{d}{dx_3} \int_{\Omega_{12}} u_\alpha dx' \right) \right\|_{H^{-1}(\Omega_3; L^2(\Omega_{12}))} \\ \leq C \left(\sum_{\alpha\beta} \|E_{\alpha\beta}(u)\|_{L^2(\Omega)} + \sum_{\alpha} \|E_{\alpha 3}(u)\|_{L^2(\Omega)} \right) \end{aligned}$$

for every $u \in H_{dn}^1(\Omega; \mathbb{R}^3)$.

5. Limit strain characterization. Let

$$\begin{aligned} H_m^1(\Omega_{12}) &:= \{z \in H^1(\Omega_{12}) : \int_{\Omega_{12}} z \, dx' = 0\}, \\ H_{m_1}^1(\Omega_1; L^2(\Omega_{23})) &:= \{z \in H^1(\Omega_1; L^2(\Omega_{23})) : \int_{\Omega_1} z \, dx_1 = 0 \text{ a.e. in } \Omega_{23}\}, \\ H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12})) &:= \{v \in H^{-1}(\Omega_3; L^2(\Omega_{12})) : \langle z_3, \varphi \rangle = 0 \\ &\quad \forall \varphi \in H_0^1(\Omega_3; L^2(\Omega_2))\}, \end{aligned}$$

where the bracket in the last definition has to be understood in the sense of the duality $H^{-1}(\Omega_3; L^2(\Omega_{12})) \times H_0^1(\Omega_3; L^2(\Omega_{12}))$.

Let (u^ε) be the sequence of solutions to problems (6). From (7) and assumption (8) it follows that

$$(16) \quad \sup_\varepsilon \|E^\varepsilon(u^\varepsilon)\|_{L^2(\Omega)} < +\infty.$$

Hence, possibly passing to a subsequence, we have that

$$E^\varepsilon(u^\varepsilon) \rightharpoonup E \quad \text{in } L^2(\Omega; \mathbb{R}_{\text{sym}}^{3 \times 3})$$

for some $E \in L^2(\Omega; \mathbb{R}_{\text{sym}}^{3 \times 3})$.

In this section we characterize the limit strain E . For clarity we state several lemmas.

LEMMA 5.1 (component 33). *There exists a function \bar{u} in the set of the so-called Bernoulli–Navier displacements*

$$\mathcal{U} := \{u \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{\alpha i}(u) = 0\}$$

such that

$$E_{33} = \frac{\partial \bar{u}_3}{\partial x_3} = E_{33}(\bar{u}).$$

Proof. From (16), the structure of E^ε , and Korn’s inequality, we have that

$$C \geq \|E^\varepsilon(u^\varepsilon)\|_{L^2(\Omega)} \geq \|E(u^\varepsilon)\|_{L^2(\Omega)} \geq C_K \|u^\varepsilon\|_{H^1(\Omega)},$$

where C is a constant independent of ε and C_K is Korn’s constant. Hence, up to a subsequence $u^\varepsilon \rightharpoonup \bar{u}$ in $H^1(\Omega; \mathbb{R}^3)$, for some $\bar{u} \in H_{dn}^1(\Omega; \mathbb{R}^3)$. The claim follows by noticing that $\|E_{\alpha i}(u^\varepsilon)\|_{L^2(\Omega)} \leq C\varepsilon$, and $E_{33}^\varepsilon(u^\varepsilon) = \partial u_3^\varepsilon / \partial x_3$. In fact it follows that $\bar{u} \in \{u \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{\alpha i}(u) = 0\}$. \square

Remark 5.1 (representation of the space \mathcal{U}). It is well known (see, for instance, Le Dret [5]) that the space of Bernoulli–Navier displacements admits the following representation:

$$\mathcal{U} := \left\{ u \in H_{dn}^2(\Omega)^2 \times H_{dn}^1(\Omega) : \text{exists } \zeta \in H_{dn}^2(\Omega_3)^2 \times H_{dn}^1(\Omega_3) \text{ such that} \right. \\ \left. u_1 = \zeta_1, u_2 = \zeta_2, u_3 = \zeta_3 - x_1 \frac{d\zeta_1}{dx_3} - x_2 \frac{d\zeta_2}{dx_3} \right\}.$$

Moreover, \mathcal{U} is a Hilbert space with the norm

$$\|u\|_{\mathcal{U}} := \|E_{33}(u)\|_{L^2(\Omega)},$$

which is equivalent to that induced by $H_{dn}^1(\Omega; \mathbb{R}^3)$ (see [8]).

LEMMA 5.2 (component 23). *There exists a function \bar{v} in the space*

$$\mathcal{V} := \left\{ v \in H_{dn}^1(\Omega)^2 \times L^2(\Omega_3; H_m^1(\Omega_{12})) : \text{exist } \vartheta \in H^1(\Omega_3) \text{ such that } \vartheta(0) = 0 \right. \\ \left. \text{and } \varrho \in L^2(\Omega_3; H^1(\Omega_2)) \text{ such that } v_1(x) = -x_2\vartheta(x_3), v_2(x) = x_1\vartheta(x_3), \right. \\ \left. v_3(x) = x_1x_2 \frac{d\vartheta}{dx_3}(x_3) + \varrho(x_2, x_3) \right\}$$

such that

$$E_{23} = E_{23}(\bar{v}).$$

Proof. Let $v_i^\varepsilon := \frac{1}{\varepsilon}u_i^\varepsilon$ and $\vartheta^\varepsilon := \vartheta(v^\varepsilon)$ (see (9)).

First of all, by adapting an argument of [4] and Lemmas 4.4 and 4.5, we prove that there exists $\vartheta \in H^1(\Omega_3)$ with $\vartheta(0) = 0$, such that, up to subsequences,

$$(17) \quad \vartheta^\varepsilon \rightarrow \vartheta \text{ in } L^2(\Omega_3).$$

Applying Lemma 4.1, there exists a constant C such that

$$(18) \quad \|\tilde{v}^\varepsilon\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} \leq C \sum_{\alpha, \beta} \|E_{\alpha\beta}(v^\varepsilon)\|_{L^2(\Omega_3; L^2(\Omega_{12}))}.$$

Since, furthermore,

$$E_{11}(v^\varepsilon)_{11} = \varepsilon^3 E_{11}^\varepsilon(u^\varepsilon), \quad E_{12}(v^\varepsilon) = \varepsilon^2 E_{12}^\varepsilon(u^\varepsilon), \quad E_{22}(v^\varepsilon) = \varepsilon E_{22}^\varepsilon(u^\varepsilon),$$

and using (16) we have

$$(19) \quad \|E_{\alpha\beta}(v^\varepsilon)\|_{L^2(\Omega)} \leq \varepsilon \|E_{\alpha\beta}^\varepsilon(u^\varepsilon)\|_{L^2(\Omega)} \leq C \varepsilon.$$

From (18) we get

$$(20) \quad \|\tilde{v}^\varepsilon\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} \leq C\varepsilon;$$

hence

$$(21) \quad \tilde{v}^\varepsilon \rightarrow 0 \text{ in } L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2)).$$

Now let $\eta \in C_c^\infty(\Omega_{12})$ be such that

$$\int_{\Omega_{12}} \eta \, dx' = -\frac{I_O}{2}.$$

Then, taking into account (14), we have

$$\begin{aligned}
 I_O \vartheta^\varepsilon &= -2\vartheta^\varepsilon \int_{\Omega_{12}} \eta \, dx' = -\vartheta^\varepsilon \int_{\Omega_{12}} \eta D_\alpha x_\alpha \, dx' \\
 &= \vartheta^\varepsilon \int_{\Omega_{12}} D_\alpha \eta \, x_\alpha \, dx' = \vartheta^\varepsilon \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} \mathcal{E}_{\beta\gamma} D_\alpha \eta \, x_\beta \, dx' \\
 &= \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta \mathcal{E}_{\beta\gamma} x_\beta \vartheta^\varepsilon \, dx' \\
 &= \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta \left(\wp_\gamma(v^\varepsilon) - \frac{1}{|\Omega_{12}|} \int_{\Omega_{12}} v_\gamma^\varepsilon \, dx' \right) dx' \\
 &= \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta \wp_\gamma(v^\varepsilon) \, dx' \\
 &= \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta v_\gamma^\varepsilon \, dx' - \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta (v^\varepsilon - \wp_\gamma(v^\varepsilon)) \, dx'.
 \end{aligned}$$

Hence, denoting by

$$\tilde{\vartheta}^\varepsilon = \frac{1}{I_O} \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta w_\gamma^\varepsilon \, dx'$$

and recalling (21), we find

$$(22) \quad \vartheta^\varepsilon - \tilde{\vartheta}^\varepsilon \rightarrow 0 \text{ in } L^2(\Omega_3).$$

We now show that $D_3 \tilde{\vartheta}^\varepsilon$ is bounded in L^2 . Since $\mathcal{E}_{\alpha\gamma} D_\alpha D_\gamma \eta = 0$ everywhere in Ω_3 and $D_\alpha \eta = 0$ on $\partial\Omega_3$, we have

$$I_O D_3 \tilde{\vartheta}^\varepsilon = \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta D_3 v_\gamma^\varepsilon \, dx' = 2 \int_{\Omega_{12}} \mathcal{E}_{\alpha\gamma} D_\alpha \eta E_{\gamma 3}(v^\varepsilon) \, dx',$$

but $E_{13}(v^\varepsilon) = \varepsilon E_{13}^\varepsilon(u^\varepsilon)$ and $E_{23}(v^\varepsilon) = E_{23}^\varepsilon(u^\varepsilon)$, and therefore $D_3 \tilde{\vartheta}^\varepsilon$ is bounded in $L^2(\Omega_3)$. Since $\tilde{\vartheta}^\varepsilon(0) = 0$, $\tilde{\vartheta}^\varepsilon$ is then bounded in $H^1(\Omega_3)$ so that there exists $\vartheta \in H^1(\Omega_3)$ with $\vartheta(0) = 0$, such that, up to subsequences,

$$\tilde{\vartheta}^\varepsilon \rightharpoonup \vartheta \text{ in } H^1(\Omega_3).$$

Thus, from (22) we obtain (17).

Let us now set

$$\bar{v}_1^\varepsilon := v_1^\varepsilon - \int_{\Omega_{12}} v_1^\varepsilon \, dx', \quad \bar{v}_2^\varepsilon := v_2^\varepsilon - \int_{\Omega_{12}} v_2^\varepsilon \, dx',$$

and

$$\bar{v}_3^\varepsilon := v_3^\varepsilon - \left(\int_{\Omega_{12}} v_3^\varepsilon \, dx' - x_\alpha \frac{d}{dx_3} \int_{\Omega_{12}} v_\alpha^\varepsilon \, dx' \right).$$

Observing that, by the definitions,

$$\bar{v}_1^\varepsilon = \tilde{v}_1^\varepsilon - x_2 \vartheta^\varepsilon, \quad \bar{v}_2^\varepsilon = \tilde{v}_2^\varepsilon + x_1 \vartheta^\varepsilon,$$

from (17) and (21) we have that

$$\bar{v}_1^\varepsilon \rightarrow -x_2\vartheta, \quad \bar{v}_2^\varepsilon \rightarrow +x_1\vartheta \quad \text{in } L^2(\Omega_3; H^1(\Omega_{12})).$$

By Theorem 4.4, we have that

$$\|\bar{v}_3^\varepsilon\|_{H^{-1}(\Omega_3; L^2(\Omega_{12}))} \leq C,$$

and hence there exists $\bar{v}_3 \in H^{-1}(\Omega_3; L^2(\Omega_{12}))$ such that, up to subsequences,

$$(23) \quad \bar{v}_3^\varepsilon \rightharpoonup \bar{v}_3 \text{ in } H^{-1}(\Omega_3; L^2(\Omega_{12})).$$

Moreover, let us set

$$\bar{v}_1 = -x_2\vartheta, \quad \bar{v}_2 = +x_1\vartheta,$$

and check that the vector field \bar{v} so defined satisfies the properties claimed in the statement of Lemma 4.1. A simple computation shows that

$$E_{13}(\bar{v}^\varepsilon) = E_{13}(v^\varepsilon) \quad \text{and} \quad E_{23}(\bar{v}^\varepsilon) = E_{23}(v^\varepsilon).$$

Noticing that $E_{13}(v^\varepsilon) = \varepsilon E_{13}^\varepsilon(u^\varepsilon) \rightarrow 0$ in $L^2(\Omega)$ and $E_{13}(\bar{v}^\varepsilon) \rightarrow E_{13}(\bar{v})$ in $\mathcal{D}'(\Omega)$, we obtain that $E_{13}(\bar{v}) = 0$. Hence,

$$\frac{\partial \bar{v}_3}{\partial x_1} = -\frac{\partial \bar{v}_1}{\partial x_3} = x_2 \frac{d\vartheta}{dx_3}$$

and, integrating with respect to x_1 ,

$$\bar{v}_3 = x_1 x_2 \frac{d\vartheta}{dx_3} + \varrho(x_2, x_3)$$

for some function $\varrho \in L^2(\Omega_3; H^1(\Omega_2))$. Moreover, $\bar{v}_3 \in L^2(\Omega_3; H^1(\Omega_{12}))$ and, from (23) and the fact that $\bar{v}_3, \bar{v}_3^\varepsilon \in L^2(\Omega)$, we then obtain easily that $\int_{\Omega_{12}} \bar{v}_3 dx' = 0$, which concludes the proof. \square

LEMMA 5.3 (characterization of the space \mathcal{V}). *The space \mathcal{V} admits the following characterization:*

$$(24) \quad \mathcal{V} = \{v \in H_{dn}^1(\Omega)^2 \times L^2(\Omega_3; H_m^1(\Omega_{12})) : E_{\alpha\beta}(v) = 0, E_{13}(v) = 0, \\ E_{23}(v) \in L^2(\Omega) \text{ and } \int_{\Omega_{12}} v_\alpha dx' = 0 \text{ a.e.}\}.$$

Moreover, it is a Hilbert space with the norm

$$\|v\|_{\mathcal{V}} := \|W_{13}(v)\|_{L^2(\Omega)} + \|E_{23}(v)\|_{L^2(\Omega)},$$

where

$$W_{13}(v) = \frac{1}{2} \left(\frac{\partial v_1}{\partial x_3} - \frac{\partial v_3}{\partial x_1} \right).$$

Proof. Let us call V the space at the right-hand side of equality (24) and let \mathcal{V} be as in the statement of Lemma 5.2. It is trivial to check that $\mathcal{V} \subseteq V$. Let us prove the opposite inclusion. Let $v \in V$. Since $v_\alpha \in H_{dn}^1(\Omega)$ and $E_{\alpha\beta}(v) = 0$, by integration there exists $\vartheta \in L^2(\Omega_3)$ such that

$$v_1 = -x_2\vartheta(x_3) + a_1(x_3), \quad v_2 = x_1\vartheta(x_3) + a_2(x_3),$$

and since $\int_{\Omega_{12}} v_\alpha dx' = 0$, we have $a_1 = a_2 = 0$ a.e.. From the resulting expression of v_α it follows that $\vartheta \in H_{dn}^1(\Omega)$ and $\vartheta(0) = 0$.

Since $E_{13}(v) = 0$ we obtain that $\partial v_3 / \partial x_1 = x_2 \vartheta'(x_3)$. Then there exists $\varrho \in L^2(\Omega_{23})$ such that

$$v_3 = x_1 x_2 \vartheta'(x_3) + \varrho(x_2, x_3),$$

from which it follows also that $\varrho \in L^2(\Omega_3; H_m^1(\Omega_2))$ and that $E_{23}(v) \in L^2(\Omega)$.

The last part of the claim follows from the fact that $L^2(\Omega_3; H_m^1(\Omega_{12}))$ is a Hilbert space with the scalar product

$$\langle u_3, v_3 \rangle_{L^2(\Omega_3; H_m^1(\Omega_{12}))} = \int_{\Omega_3} \langle u_3(x_3), v_3(x_3) \rangle_{H_m^1(\Omega_{12})} dx_3,$$

and \mathcal{V} is a closed subspace of $H_{dn}^1(\Omega)^2 \times L^2(\Omega_3; H_m^1(\Omega_{12}))$ which is Hilbert with the product norm

$$\|v_1\|_{H_{dn}^1(\Omega)} + \|v_2\|_{H_{dn}^1(\Omega)} + \|v_3\|_{L^2(\Omega_3; H_m^1(\Omega_{12}))}$$

induced by $H_{dn}^1(\Omega)^2 \times L^2(\Omega_3; H_m^1(\Omega_{12}))$. The proof that this norm is equivalent to $\|v\|_{\mathcal{V}}$ is an easy consequence of the Poincaré inequality, the representation lemma, Lemma 5.2, and the characterization of the space \mathcal{V} proved above. \square

LEMMA 5.4 (components 22 and 13). *There exists a function \bar{w} in the space*

$$\begin{aligned} \mathcal{W} = \{w \in RD_2^1(\Omega) \times H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12})) : E_{11}(w) = E_{12}(w) = 0, \\ E_{13}(w) \in L^2(\Omega)\} \end{aligned}$$

such that

$$E_{22} = \frac{\partial \bar{w}_2}{\partial x_2} = E_{22}(\bar{w}), \quad E_{13} = \frac{1}{2} \left(\frac{\partial \bar{w}_1}{\partial x_3} + \frac{\partial \bar{w}_3}{\partial x_1} \right) = E_{13}(\bar{w}).$$

Moreover, \mathcal{W} is a Hilbert space with the norm

$$\|w\|_{\mathcal{W}} := \|E_{22}(w)\|_{L^2(\Omega)} + \|E_{13}(w)\|_{L^2(\Omega)} + \|w_3\|_{H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12}))}.$$

Proof. Let

$$w^\varepsilon := \frac{1}{\varepsilon^2} u^\varepsilon;$$

then

$$(25) \quad \begin{aligned} \|E_{11}(w^\varepsilon)\|_{L^2(\Omega)} &\leq C\varepsilon^2, & \|E_{12}(w^\varepsilon)\|_{L^2(\Omega)} &\leq C\varepsilon, \\ \|E_{22}(w^\varepsilon)\|_{L^2(\Omega)} &\leq C, & \|E_{13}(w^\varepsilon)\|_{L^2(\Omega)} &\leq C. \end{aligned}$$

Let us recall that $\tilde{w}_\alpha^\varepsilon := w_\alpha^\varepsilon - \wp_\alpha(w^\varepsilon)$. By Lemma 4.1 and using (25), we have

$$\|\tilde{w}^\varepsilon\|_{L^2(\Omega_3; H^1(\Omega_{12}; \mathbb{R}^2))} \leq C \sum_{\alpha, \beta} \|E_{\alpha\beta}(w^\varepsilon)\|_{L^2(\Omega)} \leq C,$$

$\int_{\Omega_{12}} \tilde{w}^\varepsilon dx' = 0$, and $\vartheta(\tilde{w}^\varepsilon) = 0$. Thus, up to a subsequence,

$$(26) \quad \tilde{w}_\alpha^\varepsilon \rightharpoonup \bar{w}_\alpha \quad \text{in } L^2(\Omega_3; H^1(\Omega_{12})) \text{ for } \alpha = 1, 2;$$

moreover, a.e.,

$$(27) \quad \int_{\Omega_{12}} \bar{w}_\alpha dx' = 0, \quad \vartheta(\bar{w}) = 0.$$

Hence $(\bar{w}_1, \bar{w}_2) \in RD_2^\perp(\Omega)$ (see (12)). Using (25) we also obtain

$$E_{11}(\bar{w}) = E_{12}(\bar{w}) = 0.$$

Moreover,

$$E_{22}^\varepsilon(u^\varepsilon) = E_{22}(w^\varepsilon) = E_{22}(\tilde{w}^\varepsilon) \rightharpoonup E_{22}(\bar{w}) \text{ in } L^2(\Omega).$$

Let

$$\tilde{w}_3^\varepsilon := w_3^\varepsilon - \left(x_1 x_2 \frac{d\vartheta(w^\varepsilon)}{dx_3} - x_1 \frac{d}{dx_3} \int_{\Omega_{12}} w_1^\varepsilon dx' + \int_{\Omega_1} w_3^\varepsilon dx_1 \right).$$

By Theorem 4.3 we have that

$$\|\tilde{w}_3^\varepsilon\|_{H^{-1}(\Omega_3; L^2(\Omega_{12}))} \leq C \left(\sum_{\alpha, \beta} \|E_{\alpha\beta}(w^\varepsilon)\|_{L^2(\Omega)} + \|E_{13}(w^\varepsilon)\|_{L^2(\Omega)} \right) \leq C.$$

Hence, up to a subsequence,

$$\tilde{w}_3^\varepsilon \rightharpoonup w_3 \quad \text{in } H^{-1}(\Omega_3; L^2(\Omega_{12})).$$

Taking into account that (see (16))

$$\wp_1(w^\varepsilon) = \int_{\Omega_{12}} w_1^\varepsilon dx' - x_2 \vartheta(w^\varepsilon), \quad \wp_2(w^\varepsilon) = \int_{\Omega_{12}} w_2^\varepsilon dx' + x_1 \vartheta(w^\varepsilon),$$

we easily deduce that $E_{13}^\varepsilon(u^\varepsilon) = E_{13}(w^\varepsilon) = E_{13}(\tilde{w}^\varepsilon)$ and hence that

$$E_{13} = \frac{1}{2} \left(\frac{\partial \bar{w}_1}{\partial x_3} + \frac{\partial \bar{w}_3}{\partial x_1} \right).$$

Finally, since $\int_{\Omega_1} \tilde{w}_3^\varepsilon dx_1 = 0$, we have that

$$\langle \bar{w}_3, \varphi \rangle_{H^{-1}(\Omega_3; L^2(\Omega_{12})) \times H_0^1(\Omega_3; L^2(\Omega_{12}))} = 0$$

for all $\varphi \in H_0^1(\Omega_3; L^2(\Omega_2))$.

The last part of the claim follows from the fact that \mathscr{W} is a closed subspace of

$$\{z \in RD_2^\perp(\Omega) \times H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12})) : E_{13}(z) \in L^2(\Omega)\}$$

which, in turn, is a Hilbert space under the scalar product

$$\langle z, \zeta \rangle := \langle z, \zeta \rangle_{RD_2^\perp(\Omega) \times H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12}))} + \int_{\Omega} E_{13}(z) E_{13}(\zeta) dx$$

and from an application of Proposition 4.1. \square

LEMMA 5.5 (representation of the space \mathscr{W}). *The space \mathscr{W} admits the following representation:*

$$(28) \quad \mathscr{W} = \left\{ w \in L^2(\Omega_3; H^1(\Omega_{12}))^2 \times H_{m_1}^{-1}(\Omega_3; L^2(\Omega_{12})) : E_{13}(w) \in L^2(\Omega), \right. \\ \left. \text{there exists } \eta_1 \in L^2(\Omega_3; H^2(\Omega_2)) \text{ and } \eta_2 \in L^2(\Omega_3; H^1(\Omega_2)) \text{ such that} \right. \\ \left. w_1(x) = \eta_1(x_2, x_3), \quad w_2(x) = -x_1 \frac{\partial \eta_1}{\partial x_2}(x_2, x_3) + \eta_2(x_2, x_3), \right. \\ \left. \int_{\Omega_2} \eta_\alpha dx_2 = 0, \quad \int_{\Omega_2} \left(\frac{a_1^2}{12} \frac{\partial \eta_1}{\partial x_2} + x_2 \eta_1 \right) dx_2 = 0 \right\}.$$

Proof. Let us call W the set on the right-hand side of equality (28), and let \mathscr{W} be as in the statement of Lemma 5.4. Then it is trivial to check that $W \subseteq \mathscr{W}$. Let us prove the converse inequality. Let $w \in \mathscr{W}$ as in Lemma 5.4. Since $E_{11}(w) = E_{12}(w) = 0$, by integration we deduce that there exist $\eta_1 \in L^2(\Omega_3; H^2(\Omega_2))$ and $\eta_2 \in L^2(\Omega_3; H^1(\Omega_2))$ such that

$$w_1(x) = \eta_1(x_2, x_3), \quad w_2(x) = -x_1 \frac{\partial \eta_1}{\partial x_2}(x_2, x_3) + \eta_2(x_2, x_3).$$

Since

$$(29) \quad \int_{\Omega_{12}} w_\alpha dx' = 0, \quad \vartheta(w) = 0,$$

we have that

$$\int_{\Omega_2} \eta_\alpha dx_2 = 0 \quad \text{a.e. for } \alpha = 1, 2,$$

and

$$\int_{\Omega_2} \left(\frac{a_1^2}{12} \frac{\partial \eta_1}{\partial x_2} + x_2 \eta_1 \right) dx_2 = 0,$$

a.e.. \square

LEMMA 5.6 (component 12). *There exists a vector function \bar{p} in the set*

$$\mathscr{P} = \{0\} \times H_{m_1}^1(\Omega_1; L^2(\Omega_{23})) \times \{0\}$$

such that

$$(30) \quad E_{12} = \frac{1}{2} \frac{\partial \bar{p}_2}{\partial x_1} = E_{12}(\bar{p}).$$

Moreover, \mathscr{P} is a Hilbert space with the norm

$$\|p\|_{\mathscr{P}} := \|E_{12}(p)\|.$$

Proof. Let

$$p_\alpha^\varepsilon := \frac{1}{\varepsilon^3} u_\alpha^\varepsilon \quad \text{for } \alpha = 1, 2,$$

and

$$\bar{p}_1^\varepsilon := p_1^\varepsilon - \int_{\Omega_1} p_1^\varepsilon dx_1, \quad \bar{p}_2^\varepsilon := p_2^\varepsilon - \int_{\Omega_1} p_2^\varepsilon dx_1 + x_1 \frac{\partial}{\partial x_2} \int_{\Omega_1} p_1^\varepsilon dx_1.$$

Then $E_{11}(\bar{p}^\varepsilon) = E_{11}(p^\varepsilon) = \varepsilon E_{11}^\varepsilon(u^\varepsilon)$; hence, by Theorem 4.1 we have

$$\bar{p}_1^\varepsilon \rightarrow 0, \quad \frac{\partial \bar{p}_1^\varepsilon}{\partial x_1} \rightarrow 0 \quad \text{in } L^2(\Omega).$$

Since $E_{12}(\bar{p}^\varepsilon) = E_{12}(p^\varepsilon) = E_{12}^\varepsilon(u^\varepsilon)$, by Theorem 4.2 we also have that, up to subsequences,

$$\bar{p}_2^\varepsilon \rightharpoonup \bar{p}_2 \text{ in } H^{-1}(\Omega_2; L^2(\Omega_{13})).$$

Setting $\bar{p} := (0, \bar{p}_2, 0)$, we have

$$E_{12} = E_{12}(\bar{p}) = \frac{1}{2} \frac{\partial \bar{p}_2}{\partial x_1},$$

that is, (30). It remains then to prove that $\bar{p}_2 \in H^1(\Omega_1; L^2(\Omega_{23}))$ and that $\int_{\Omega_1} \bar{p}_2 dx_1 = 0$.

Since $\bar{p}_2 \in H^{-1}(\Omega_2; L^2(\Omega_{13}))$, for any $\varphi \in H_0^1(\Omega_2)$ and any $\psi \in L^2(\Omega_{13})$ the product $\varphi \otimes \psi$ belongs to $H_0^1(\Omega_2; L^2(\Omega_{13}))$, and the linear map

$$P_\varphi : L^2(\Omega_{13}) \rightarrow \mathbb{R}, \quad \langle P_\varphi, \psi \rangle := \langle \bar{p}_2, \varphi \otimes \psi \rangle$$

satisfies the estimate

$$|\langle P_\varphi, \psi \rangle| \leq \|\bar{p}_2\|_{H^{-1}} \|\varphi\|_{H_0^1} \|\psi\|_{L^2}.$$

Thus $P_\varphi \in L^2(\Omega_{13})$. Moreover, from the definition of P_φ and the fact that $\frac{\partial \bar{p}_2}{\partial x_1} \in L^2(\Omega)$, we obtain that $\frac{\partial P_\varphi}{\partial x_1} \in L^2(\Omega_{13})$ and also that

$$(31) \quad \frac{\partial P_\varphi}{\partial x_1} = \int_{\Omega_2} \frac{\partial \bar{p}_2}{\partial x_1} \varphi dx_2.$$

Since $\int_{\Omega_1} \bar{p}_2^\varepsilon dx_1 = 0$ it follows from the definitions that

$$(32) \quad \int_{\Omega_1} P_\varphi(s, x_3) ds = \langle \bar{p}_2, \varphi \rangle = 0$$

for almost every $x_3 \in \Omega_3$; using this fact, the following Poincaré inequality holds:

$$\|P_\varphi(\cdot, x_3)\|_{L^2(\Omega_1)} \leq C \left\| \frac{\partial P_\varphi}{\partial x_1}(\cdot, x_3) \right\|_{L^2(\Omega_1)},$$

where the constant C depends only on the domain Ω_1 and is therefore independent of x_3 . By substituting (31) inside the Poincaré inequality, we obtain that

$$(33) \quad \|P_\varphi\|_{L^2(\Omega_{13})} \leq C \left\| \frac{\partial \bar{p}_2}{\partial x_1} \right\|_{L^2(\Omega)} \|\varphi\|_{L^2(\Omega_2)}.$$

Using the density of $C_c^\infty(\Omega_2) \otimes C_c^\infty(\Omega_{13})$ in $L^2(\Omega)$ (see, for instance, Treves [10, Theorem 39.2 and subsequent Corollary 3]), the fact that $P_\varphi \in L^2(\Omega_{13})$, and inequality (33), we have that

$$\begin{aligned} \|\bar{p}_2\|_{L^2(\Omega)} &= \sup_{\varphi \in C_c^\infty(\Omega_2), \psi \in C_c^\infty(\Omega_{13})} \frac{|\langle \bar{p}_2, \varphi \psi \rangle|}{\|\varphi\|_{L^2(\Omega_2)} \|\psi\|_{L^2(\Omega_{13})}} \\ &= \sup_{\varphi \in C_c^\infty(\Omega_2), \psi \in C_c^\infty(\Omega_{13})} \frac{|\langle P_\varphi, \psi \rangle|}{\|\varphi\|_{L^2(\Omega_2)} \|\psi\|_{L^2(\Omega_{13})}} \\ &\leq \sup_{\varphi \in C_c^\infty(\Omega_2)} \frac{\|P_\varphi\|_{L^2(\Omega_{13})}}{\|\varphi\|_{L^2(\Omega_2)}} \leq C \left\| \frac{\partial \bar{p}_2}{\partial x_1} \right\|_{L^2(\Omega)}; \end{aligned}$$

hence $\bar{p}_2 \in L^2(\Omega)$. Thus $\bar{p}_2 \in H^1(\Omega_1; L^2(\Omega_{23}))$, and (32) implies $\int_{\Omega_1} \bar{p}_2 dx_1 = 0$.

The last part of the claim follows from the fact that $H^1_{m_1}(\Omega_1; L^2(\Omega_{23}))$ is a Hilbert space with the norm

$$\|p\|_{H^1_{m_1}(\Omega_1; L^2(\Omega_{23}))} := \left\| \frac{\partial p_2}{\partial x_1} \right\|_{L^2(\Omega)},$$

which, by Theorem 4.1, turns out to be equivalent to the canonical one. □

LEMMA 5.7 (component 11). *There exists a function \bar{q} in the space*

$$\mathcal{Q} := H^1_{m_1}(\Omega_1; L^2(\Omega_{23})) \times \{0\}^2$$

such that

$$E_{11} = \frac{\partial \bar{q}_1}{\partial x_1} = E_{11}(\bar{q}).$$

Moreover, \mathcal{Q} is a Hilbert space with the norm

$$\|q\|_{\mathcal{Q}} := \|E_{11}(q)\|.$$

Proof. Let

$$q_1^\varepsilon := \frac{1}{\varepsilon^4} \left(u_1^\varepsilon - \int_{\Omega_1} u_1^\varepsilon dx_1 \right);$$

then

$$\sup_\varepsilon \left\| \frac{\partial q_1^\varepsilon}{\partial x_1} \right\|_{L^2(\Omega)} = \sup_\varepsilon \|E_{11}^\varepsilon(u^\varepsilon)\|_{L^2(\Omega)} \leq C,$$

and by Theorem 4.1 we have $\sup_\varepsilon \|q_1^\varepsilon\|_{L^2(\Omega)} \leq C$. Then, up to a subsequence, $q_1^\varepsilon \rightharpoonup \bar{q}_1$ in $L^2(\Omega)$, and $E_{11}^\varepsilon(u^\varepsilon) = \partial \bar{q}_1^\varepsilon / \partial x_1 \rightharpoonup \partial \bar{q}_1 / \partial x_1$ in $L^2(\Omega)$ for some $\bar{q}_1 \in L^2(\Omega)$.

The last part of the claim follows from the fact that $H^1_{m_1}(\Omega_1; L^2(\Omega_{23}))$ is a Hilbert space with the norm

$$\|q\|_{H^1_{m_1}(\Omega_1; L^2(\Omega_{23}))} := \|q_{1,1}\|_{L^2(\Omega)},$$

which, by Theorem 4.1, turns out to be equivalent to the canonical one. □

6. The limit problem. Let us consider the space $\mathcal{A} := \mathcal{U} \times \mathcal{V} \times \mathcal{W} \times \mathcal{P} \times \mathcal{Q}$. According to the notation and the results proved in the previous section, \mathcal{A} is a Hilbert space when endowed with the product norm

$$\|(u, v, w, p, q)\|_{\mathcal{A}} := \|u\|_{\mathcal{U}} + \|v\|_{\mathcal{V}} + \|w\|_{\mathcal{W}} + \|p\|_{\mathcal{P}} + \|q\|_{\mathcal{Q}}.$$

Given a 5-tuple of vector valued distributions $(u, v, w, p, q) \in \mathcal{D}'(\Omega; \mathbb{R}^3)^5$, let us define

$$(34) \quad E(u, v, w, p, q) := \begin{pmatrix} E_{11}(q) & E_{12}(p) & E_{13}(w) \\ & E_{22}(w) & E_{23}(v) \\ \text{Sym.} & & E_{33}(u) \end{pmatrix}.$$

We are now in a position to state the main result of this paper.

THEOREM 6.1. Let \mathbb{C} be a positive definite fourth order tensor field on Ω with the minor symmetries, i.e., $\mathbb{C}_{ijkl} = \mathbb{C}_{jikl} = \mathbb{C}_{ijlk}$. Let F^ε be a second order symmetric tensor field which belongs to $L^2(\Omega; \mathbb{R}^{3 \times 3})$. Then problem (6), that is,

$$(35) \quad \begin{cases} u^\varepsilon \in H_{dn}^1(\Omega; \mathbb{R}^3), \\ \int_{\Omega} \mathbb{C}E^\varepsilon(u^\varepsilon) \cdot E^\varepsilon(\varphi) \, dx = \int_{\Omega} F^\varepsilon \cdot E^\varepsilon(\varphi) \, dx \quad \forall \varphi \in H_{dn}^1(\Omega; \mathbb{R}^3), \end{cases}$$

admits a unique solution u^ε . Moreover, if $F^\varepsilon \rightarrow F$ in $L^2(\Omega; \mathbb{R}^{3 \times 3})$, then we have the following:

1. the problem

$$(36) \quad \int_{\Omega} \mathbb{C}E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \cdot E(u, v, w, p, q) \, dx = \int_{\Omega} F \cdot E(u, v, w, p, q) \, dx$$

$$\forall (u, v, w, p, q) \in \mathcal{A}$$

admits a unique solution $(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \in \mathcal{A}$;

- 2. $u^\varepsilon \rightharpoonup \bar{u}$ in $H^1(\Omega; \mathbb{R}^3)$;
- 3. $E^\varepsilon(u^\varepsilon) \rightarrow E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q})$ in $L^2(\Omega; \mathbb{R}^{3 \times 3})$.

The following corollary can be seen as a corrector result.

COROLLARY 6.1. If the solution $(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q})$ of problem (36) is such that

$$(37) \quad \frac{\partial \bar{v}_3}{\partial x_3}, E_{23}(\bar{w}), \frac{\partial \bar{w}_3}{\partial x_3}, \frac{\partial \bar{q}_1}{\partial x_2}, \frac{\partial \bar{q}_1}{\partial x_3}, \frac{\partial \bar{p}_2}{\partial x_2}, \frac{\partial \bar{p}_2}{\partial x_3} \in L^2(\Omega),$$

then

$$\|E^\varepsilon(u^\varepsilon) - E^\varepsilon(\bar{u}^\varepsilon)\|_{L^2(\Omega; \mathbb{R}^{3 \times 3})} \rightarrow 0,$$

where

$$\bar{u}^\varepsilon := \bar{u} + \varepsilon \bar{v} + \varepsilon^2 \bar{w} + \varepsilon^3 \bar{p} + \varepsilon^4 \bar{q}.$$

Proof. Since

$$E^\varepsilon(\bar{u}^\varepsilon) = E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) + \varepsilon \begin{pmatrix} 0 & E_{12}(\bar{q}) & 0 \\ & E_{22}(\bar{p}) & E_{23}(\bar{w}) \\ \text{Sym.} & & E_{33}(\bar{v}) \end{pmatrix} + \varepsilon^2 \begin{pmatrix} 0 & 0 & E_{13}(\bar{q}) \\ & 0 & E_{23}(\bar{p}) \\ \text{Sym.} & & E_{33}(\bar{w}) \end{pmatrix}$$

the additional regularity assumptions imply that $E^\varepsilon(\bar{u}^\varepsilon) \in L^2(\Omega; \mathbb{R}^{3 \times 3})$ and

$$\|E^\varepsilon(\bar{u}^\varepsilon) - E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q})\|_{L^2(\Omega; \mathbb{R}^{3 \times 3})} \rightarrow 0.$$

Then the claim follows from step 3 of Theorem 6.1. \square

In order to prove Theorem 6.1, we introduce the subspaces of $H_{dn}^1(\Omega; \mathbb{R}^3)$,

$$\hat{\mathcal{U}} := \{u \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{\alpha\beta}(u) = E_{\alpha 3}(u) = 0\},$$

$$\hat{\mathcal{V}} := \{v \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{\alpha\beta}(v) = E_{13}(v) = 0\},$$

$$\hat{\mathcal{W}} := \{w \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{1\beta}(w) = 0\},$$

$$\hat{\mathcal{P}} := \{p \in H_{dn}^1(\Omega; \mathbb{R}^3) : E_{11}(p) = 0\},$$

$$\hat{\mathcal{Q}} := H_{dn}^1(\Omega; \mathbb{R}^3),$$

and define $\hat{\mathcal{A}} := \hat{\mathcal{U}} \times \hat{\mathcal{V}} \times \hat{\mathcal{W}} \times \hat{\mathcal{P}} \times \hat{\mathcal{Q}}$.

Let us note that $\mathcal{U} = \hat{\mathcal{U}}$, but similar equalities are not true for the spaces \mathcal{V} , \mathcal{W} , \mathcal{P} , and \mathcal{Q} . Nevertheless, for such spaces we can prove the following approximation lemma.

LEMMA 6.1. *For every $v \in \mathcal{V}$, $w \in \mathcal{W}$, $p \in \mathcal{P}$, and $q \in \mathcal{Q}$ there exist sequences (\hat{v}^n) in $\hat{\mathcal{V}}$, (\hat{w}^n) in $\hat{\mathcal{W}}$, (\hat{p}^n) in $\hat{\mathcal{P}}$, and (\hat{q}^n) in $\hat{\mathcal{Q}}$ such that the following convergences hold in the norm of $L^2(\Omega)$:*

- (i) $E_{23}(\hat{v}^n) \rightarrow E_{23}(v)$,
- (ii) $E_{13}(\hat{w}^n) \rightarrow E_{13}(w)$ and $E_{22}(\hat{w}^n) \rightarrow E_{22}(w)$,
- (iii) $E_{12}(\hat{p}^n) \rightarrow E_{12}(p)$,
- (iv) $E_{11}(\hat{q}^n) \rightarrow E_{11}(q)$.

Proof. Let us prove (i). Since $v \in \mathcal{V}$ (see Lemma 5.2), there exist $\vartheta \in H^1(\Omega_3)$ with $\vartheta(0) = 0$ and $\varrho \in L^2(\Omega_3; H^1(\Omega_2))$ such that

$$\begin{aligned} v_1(x) &= -x_2\vartheta(x_3), & v_2(x) &= x_1\vartheta(x_3), \\ v_3(x) &= x_1x_2\frac{d\vartheta}{dx_3}(x_3) + \varrho(x_2, x_3). \end{aligned}$$

Then there exist sequences $\vartheta_n \in H^2_{dn}(\Omega_3)$ and $\varrho_n \in H^1_{dn}(\Omega_{23})$ such that

$$\vartheta_n \rightarrow \vartheta \text{ in } H^1(\Omega_3), \quad \varrho_n \rightarrow \varrho \text{ in } L^2(\Omega_3; H^1(\Omega_2)).$$

Setting

$$\begin{aligned} \hat{v}_1^n(x) &= -x_2\vartheta_n(x_3), & \hat{v}_2^n(x) &= x_1\vartheta_n(x_3), \\ \hat{v}_3^n(x) &= x_1x_2\frac{d\vartheta_n}{dx_3}(x_3) + \varrho_n(x_2, x_3), \end{aligned}$$

we obtain the claim.

Let us prove (ii). As $w \in \mathcal{W}$ (see Lemma 5.4), there exist $\eta_1 \in L^2(\Omega_3; H^2(\Omega_2))$ and $\eta_2 \in L^2(\Omega_3; H^1(\Omega_2))$ such that

$$w_1(x) = \eta_1(x_2, x_3), \quad w_2(x) = -x_1\frac{\partial\eta_1}{\partial x_2}(x_2, x_3) + \eta_2(x_2, x_3).$$

Then, there exist a sequence $\eta_1^n \in H^2_{dn}(\Omega)$ with $\partial\eta_1^n/\partial x_1 = 0$ such that

$$\eta_1^n \rightarrow \eta_1 \quad \text{in } L^2(\Omega_3; H^2(\Omega_2))$$

and a sequence $\eta_2^n \in H^1_{dn}(\Omega)$ with $\partial\eta_2^n/\partial x_1 = 0$ such that

$$\eta_2^n \rightarrow \eta_2 \quad \text{in } L^2(\Omega_3; H^1(\Omega_2)).$$

Since $\partial w_3/\partial x_1 = 2E_{13}(w) - \partial\eta_1/\partial x_3$ and $E_{13}(w) \in L^2(\Omega)$, by integration, there exists $G_{13} \in H^1(\Omega_1; L^2(\Omega_{23}))$ such that

$$\frac{\partial G_{13}}{\partial x_1} = E_{13}(w) \quad \text{and} \quad w_3 = 2G_{13} - x_1\frac{\partial\eta_1}{\partial x_3},$$

where we have also used the fact that η_1 does not depend on x_1 .

We may also find a sequence $G_{13}^n \in H^1_{dn}(\Omega)$ such that

$$G_{13}^n \rightarrow G_{13} \quad \text{in } H^1(\Omega_1; L^2(\Omega_{23}))$$

and define

$$\hat{w}_1^n := \eta_1^n, \quad \hat{w}_2^n := -x_1 \frac{\partial \eta_1^n}{\partial x_2} + \eta_2^n, \quad \hat{w}_3^n := G_{13}^n - x_1 \frac{\partial \eta_1^n}{\partial x_3}.$$

Then $\hat{w}^n \in H_{dn}^1(\Omega)$, $E_{11}(\hat{w}^n) = \partial \eta_1^n / \partial x_1 = 0$, and $E_{12}(\hat{w}^n) = 0$ so that $\hat{w}^n \in \hat{\mathcal{W}}$. Moreover,

$$E_{13}(\hat{w}^n) = \frac{\partial G_{13}^n}{\partial x_1} \rightarrow \frac{\partial G_{13}}{\partial x_1} = E_{13}(w) \quad \text{in } L^2(\Omega),$$

and

$$E_{22}(\hat{w}^n) = \frac{\partial \hat{w}_2^n}{\partial x_2} = -x_1 \frac{\partial^2 \eta_1^n}{\partial x_2^2} + \frac{\partial \eta_2^n}{\partial x_2} \rightarrow \frac{\partial \hat{w}_2}{\partial x_2} = E_{22}(\hat{w}) \quad \text{in } L^2(\Omega).$$

To prove (iii) it is enough to consider, for a given $p = (0, p_2, 0) \in \mathcal{P}$, a sequence $\hat{p}_2^n \in H_{dn}^1(\Omega)$, which converges to p_2 in the norm of $H^1(\Omega_1; L^2(\Omega_{23}))$, and set $\hat{p}^n = (0, \hat{p}_2^n, 0)$. Finally, claim (iv) simply follows from the density of $H_{dn}^1(\Omega)$ in $H_{m_1}^1(\Omega_1; L^2(\Omega_{23}))$. \square

Proof of Theorem 6.1. The existence and uniqueness of the solution of problem (36) follows from an application of the Lax–Milgram lemma to the symmetric bilinear form defined on \mathcal{A} by

$$a[(u, v, w, p, q), (\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q})] := \int_{\Omega} \mathbb{C}E(u, v, w, p, q) \cdot E(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \, dx,$$

which is continuous and coercive with respect to the Hilbertian norm on \mathcal{A} defined at the beginning of this section.

Part 2 of the statement of Theorem 6.1 is actually a consequence of step 3. Let us now prove part 3. According to the results proved in the previous section, we have

$$(38) \quad E^\varepsilon(u^\varepsilon) \rightharpoonup E \text{ in } L^2(\Omega; \mathbb{R}^{3 \times 3}),$$

and there exists a $(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \in \mathcal{A}$ such that

$$E = \begin{pmatrix} E_{11}(\bar{q}) & E_{12}(\bar{p}) & E_{13}(\bar{w}) \\ & E_{22}(\bar{w}) & E_{23}(\bar{v}) \\ \text{Sym.} & & E_{33}(\bar{u}) \end{pmatrix} = E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}).$$

The result will be achieved in two steps: (i) we prove that $(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q})$ satisfies equality (36) and therefore coincides with the unique solution of the variational problem, and (ii) we show that the convergence in (38) is indeed strong.

Let $(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \in \hat{\mathcal{A}}$ and set

$$\hat{\varphi}^\varepsilon := \hat{u} + \varepsilon \hat{v} + \varepsilon^2 \hat{w} + \varepsilon^3 \hat{p} + \varepsilon^4 \hat{q};$$

then $\hat{\varphi}^\varepsilon \in H_{dn}^1(\Omega; \mathbb{R}^3)$ and an easy computation shows that, as $\varepsilon \rightarrow 0$,

$$(39) \quad E^\varepsilon(\hat{\varphi}^\varepsilon) \rightarrow \begin{pmatrix} E_{11}(\hat{q}) & E_{12}(\hat{p}) & E_{13}(\hat{w}) \\ & E_{22}(\hat{w}) & E_{23}(\hat{v}) \\ \text{Sym.} & & E_{33}(\hat{u}) \end{pmatrix} = E(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q})$$

in the norm convergence of $L^2(\Omega; \mathbb{R}^{3 \times 3})$. Taking $\varphi = \hat{\varphi}^\varepsilon$ in (35) and passing to the limit we find

$$\int_{\Omega} \mathbb{C}E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \cdot E(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \, dx = \int_{\Omega} F \cdot E(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \, dx,$$

for every $(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \in \hat{\mathcal{A}}$. This equality holds in fact for any $(u, v, w, p, q) \in \mathcal{A}$ in place of $(\hat{u}, \hat{v}, \hat{w}, \hat{p}, \hat{q}) \in \hat{\mathcal{A}}$ because of the approximation Lemma 6.1 which ensures that there exists a sequence $(u, \hat{v}^n, \hat{w}^n, \hat{p}^n, \hat{q}^n) \in \hat{\mathcal{A}}$ such that

$$\|(u, \hat{v}^n, \hat{w}^n, \hat{p}^n, \hat{q}^n) - (u, v, w, p, q)\|_{\mathcal{A}} \rightarrow 0.$$

To show that the convergence in (38) is indeed strong, it suffices to prove that $\lim_{\varepsilon \rightarrow 0} \|E^\varepsilon(u^\varepsilon)\|_{L^2(\Omega; \mathbb{R}^{3 \times 3})} = \|E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q})\|_{L^2(\Omega; \mathbb{R}^{3 \times 3})}$ or, equivalently, that

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \mathbb{C}E^\varepsilon(u^\varepsilon) \cdot E^\varepsilon(u^\varepsilon) \, dx &= \lim_{\varepsilon \rightarrow 0} \int_{\Omega} F^\varepsilon \cdot E^\varepsilon(u^\varepsilon) \, dx \\ &= \int_{\Omega} F \cdot E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \, dx \\ &= \int_{\Omega} \mathbb{C}E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \cdot E(\bar{u}, \bar{v}, \bar{w}, \bar{p}, \bar{q}) \, dx, \end{aligned}$$

where we passed to the limit thanks to the strong convergence of F^ε . □

7. Equilibrium differential equations. In this last section we derive the differential formulation of the limit problem. For simplicity we assume here that the elasticity tensor also satisfies the major symmetries; that is, $\mathbb{C}_{ijkl} = \mathbb{C}_{klij}$ for any i, j, k, l . Nevertheless, the same computation can be performed also in the general case.

To make (36) more explicit and to keep the notation compact, in writing the elasticity tensor components \mathbb{C}_{ijkl} we associate to a pair of components ij a single component s following the rule $11 \mapsto 1, 22 \mapsto 2, 33 \mapsto 3, 23 \mapsto 4, 13 \mapsto 5, 12 \mapsto 6$, and we write, for instance, c_{14} for \mathbb{C}_{1123} ; see Auld [1] for more details on the notation used. Clearly $c_{ij} = c_{ji}$. Still, for brevity, define $\bar{e}_1 = E_{11}(\bar{q})$, $\bar{e}_2 = E_{22}(\bar{w})$, $\bar{e}_3 = E_{33}(\bar{u})$, $\bar{e}_4 = 2E_{23}(\bar{v})$, $\bar{e}_5 = 2E_{13}(\bar{w})$, $\bar{e}_6 = 2E_{12}(\bar{p})$. Letting

$$\mathcal{L}(u, v, w, p, q) := \int_{\Omega} F \cdot E(u, v, w, p, q) \, dx,$$

we can then rewrite (36) as

$$\begin{aligned} \int_{\Omega} \sum_{j=1}^6 [c_{1j} \bar{e}_j E_{11}(q) + c_{2j} \bar{e}_j E_{22}(w) + c_{3j} \bar{e}_j E_{33}(u) + 2c_{4j} \bar{e}_j E_{23}(v) \\ + 2c_{5j} \bar{e}_j E_{13}(w) + 2c_{6j} \bar{e}_j E_{12}(p)] \, dx = \mathcal{L}(u, v, w, p, q), \end{aligned}$$

for every $(u, v, w, p, q) \in \mathcal{A}$. Thus

$$(40) \quad \int_{\Omega} \sum_{j=1}^6 c_{1j} \bar{e}_j E_{11}(q) \, dx = \mathcal{L}(0, 0, 0, 0, q),$$

$$(41) \quad \int_{\Omega} \sum_{j=1}^6 2c_{6j} \bar{e}_j E_{12}(p) \, dx = \mathcal{L}(0, 0, 0, p, 0),$$

$$(42) \quad \int_{\Omega} \sum_{j=1}^6 [c_{2j} \bar{e}_j E_{22}(w) + 2c_{5j} \bar{e}_j E_{13}(w)] \, dx = \mathcal{L}(0, 0, w, 0, 0),$$

$$(43) \quad \int_{\Omega} \sum_{j=1}^6 2c_{4j} \bar{e}_j E_{23}(v) \, dx = \mathcal{L}(0, v, 0, 0, 0),$$

$$(44) \quad \int_{\Omega} \sum_{j=1}^6 c_{3j} \bar{e}_j E_{33}(u) \, dx = \mathcal{L}(u, 0, 0, 0, 0).$$

In this section, for simplicity, we assume

$$(45) \quad \begin{aligned} \mathcal{L}(0, 0, \cdot, \cdot, \cdot) &= 0, \quad \mathcal{L}(0, v, 0, 0, 0) = \int_0^\ell m(x_3) \vartheta(x_3) \, dx_3, \\ \mathcal{L}(u, 0, 0, 0, 0) &= \int_{\Omega} b \cdot u \, dx + \int_{\Gamma_\ell} s \cdot u \, d\mathcal{H}^2, \end{aligned}$$

where $\Gamma_\ell = \partial\Omega \cap \{x_3 = \ell\}$; $u \in \mathcal{U}$; $b \in L^2(\Omega)$; $s \in L^2(\Gamma_\ell)$; $v \in \mathcal{V}$; $\vartheta \in H^1(\Omega_3)$, $\vartheta(0) = 0$, is related to v as in Lemma 5.2; and $m \in L^2(\Omega_3)$. Such assumptions are quite often satisfied in engineering applications.

We now derive the equilibrium equations in differential form. Let $\psi \in L^2(\Omega)$ and define

$$q_1 := \int_{-a_1/2}^{x_1} \psi(s, \cdot, \cdot) \, ds - \int_{\Omega_1} \int_{-a_1/2}^{x_1} \psi(s, \cdot, \cdot) \, ds \, dx_1.$$

Then $q := (q_1, 0, 0) \in \mathcal{Q}$ and $E_{11}(q) = \psi$; hence, from (40) and (45) we deduce

$$(46) \quad \sum_{j=1}^6 c_{1j} \bar{e}_j = 0 \quad \text{a.e..}$$

With the same argument it follows from (41) that

$$(47) \quad \sum_{j=1}^6 c_{6j} \bar{e}_j = 0 \quad \text{a.e..}$$

From (46) and (47) we deduce, since $c_{11}c_{66} - c_{16}^2 > 0$, that

$$(48) \quad \bar{e}_1 = E_{11}(\bar{q}) = - \sum_{j=2}^5 \frac{c_{66}c_{1j} - c_{16}c_{6j}}{c_{11}c_{66} - c_{16}^2} \bar{e}_j \quad \text{a.e.,}$$

$$(49) \quad \bar{e}_6 = 2E_{12}(\bar{p}) = - \sum_{j=2}^5 \frac{c_{11}c_{6j} - c_{16}c_{1j}}{c_{11}c_{66} - c_{16}^2} \bar{e}_j \quad \text{a.e..}$$

Using (45), (48), and (49) we can rewrite (42), after setting

$$\tilde{c}_{ij} = c_{ij} - c_{i1} \frac{c_{66}c_{1j} - c_{16}c_{6j}}{c_{11}c_{66} - c_{16}^2} - c_{i6} \frac{c_{11}c_{6j} - c_{16}c_{1j}}{c_{11}c_{66} - c_{16}^2},$$

for $i, j = 2, \dots, 5$, as

$$(50) \quad \int_{\Omega} \sum_{j=2}^5 [\tilde{c}_{2j}\bar{e}_j E_{22}(w) + 2\tilde{c}_{5j}\bar{e}_j E_{13}(w)] dx = \mathcal{L}(0, 0, w, 0, 0) = 0.$$

Since $w \in \mathcal{W}$, it then admits the representation given in Lemma 5.5 in terms of functions η_1 and η_2 . Choosing $\eta_1 = \eta_2 = 0$, so that $E_{22}(w) = 0$, and w_3 like it has been chosen q_1 previously, we find from (50) that

$$(51) \quad \sum_{j=2}^5 \tilde{c}_{5j}\bar{e}_j = 0 \quad \text{a.e..}$$

Let $\psi \in L^2(\Omega_{23})$. Taking $\eta_1 = w_3 = 0$ and

$$\eta_2 := \int_{-a_2/2}^{x_2} \psi(s, \cdot) ds - \int_{\Omega_2} \int_{-a_2/2}^{x_2} \psi(s, \cdot) ds dx_2$$

so that $E_{22}(w) = \psi$, we find from (50) that

$$(52) \quad \sum_{j=2}^5 \int_{\Omega_1} \tilde{c}_{2j}\bar{e}_j dx_1 = 0 \quad \text{a.e..}$$

Taking instead $\eta_2 = w_3 = 0$ and

$$\eta_1 := \int_{-a_2/2}^{x_2} \int_{-a_2/2}^t \psi(s, \cdot) ds dt - K_1 x_2 - K_2,$$

where the constants K_1 and K_2 are chosen in order to satisfy the mean integral conditions required on η_1 by (28), we have $E_{22}(w) = -x_1\psi$, and hence, from (50) we deduce

$$(53) \quad \sum_{j=2}^5 \int_{\Omega_1} x_1 \tilde{c}_{2j}\bar{e}_j dx_1 = 0 \quad \text{a.e..}$$

From (51), and observing that the positive definiteness of the elastic tensor implies $\tilde{c}_{55} > 0$, we find

$$(54) \quad \bar{e}_5 = 2E_{13}(\bar{w}) = - \sum_{j=2}^4 \frac{\tilde{c}_{5j}}{\tilde{c}_{55}} \bar{e}_j \quad \text{a.e..}$$

To solve (52) and (53) we need to write explicitly $\bar{e}_2 = E_{22}(\bar{w})$. Since $\bar{w} \in \mathcal{W}$, by Lemma 5.5, we can write

$$\bar{w}_1(x) = \bar{\eta}_1(x_2, x_3), \quad \bar{w}_2(x) = -x_1 \frac{\partial \bar{\eta}_1}{\partial x_2}(x_2, x_3) + \bar{\eta}_2(x_2, x_3),$$

where $\bar{\eta}_1$ and $\bar{\eta}_2$ belong to the appropriate spaces. Since

$$(55) \quad \bar{e}_2 = -x_1 \frac{\partial^2 \bar{\eta}_1}{\partial x_2^2} + \frac{\partial \bar{\eta}_2}{\partial x_2}$$

and using (54), we can rewrite (52) and (53) as

$$s_{1;22} \frac{\partial^2 \bar{\eta}_1}{\partial x_2^2} - s_{0;22} \frac{\partial \bar{\eta}_2}{\partial x_2} = \sum_{j=3}^4 \int_{\Omega_1} \hat{c}_{2j} \bar{e}_j \, dx_1,$$

$$s_{2;22} \frac{\partial^2 \bar{\eta}_1}{\partial x_2^2} - s_{1;22} \frac{\partial \bar{\eta}_2}{\partial x_2} = \sum_{j=3}^4 \int_{\Omega_1} x_1 \hat{c}_{2j} \bar{e}_j \, dx_1,$$

where we have set

$$(56) \quad \hat{c}_{ij} := \frac{\tilde{c}_{55} \tilde{c}_{ij} - \tilde{c}_{i5} \tilde{c}_{j5}}{\tilde{c}_{55}}, \quad s_{k;ij} := \int_{\Omega_1} x_1^k \hat{c}_{ij} \, dx_1,$$

for $i, j = 2, 3, 4$ and $k = 0, 1, 2$. From these equations we find

$$\frac{\partial^2 \bar{\eta}_1}{\partial x_2^2} = \frac{1}{s_{0;22} s_{2;22} - s_{1;22}^2} \left(s_{0;22} \sum_{j=3}^4 \int_{\Omega_1} x_1 \hat{c}_{2j} \bar{e}_j \, dx_1 - s_{1;22} \sum_{j=3}^4 \int_{\Omega_1} \hat{c}_{2j} \bar{e}_j \, dx_1 \right),$$

$$\frac{\partial \bar{\eta}_2}{\partial x_2} = \frac{1}{s_{0;22} s_{2;22} - s_{1;22}^2} \left(s_{1;22} \sum_{j=3}^4 \int_{\Omega_1} x_1 \hat{c}_{2j} \bar{e}_j \, dx_1 - s_{2;22} \sum_{j=3}^4 \int_{\Omega_1} \hat{c}_{2j} \bar{e}_j \, dx_1 \right),$$

and then by integration $\bar{\eta}_1$ and $\bar{\eta}_2$ (the fact that $s_{0;22} s_{2;22} - s_{1;22}^2 > 0$ can be checked, for instance, by using Hölder's inequality; see Wheeden and Zygmund [13, Chapter 8, Exercise 4]).

According to Remark 5.1 and Lemma 5.2, we let

$$(57) \quad \bar{e}_3 = \bar{\zeta}'_3 - x_1 \bar{\zeta}''_1 - x_2 \bar{\zeta}''_2, \quad \bar{e}_4 = 2x_1 \bar{\vartheta}' + \frac{\partial \bar{\varrho}}{\partial x_2}.$$

Setting

$$S_{mpqr}^{ijkl} := \frac{s_{i;2j} s_{k;2l} - s_{m;2p} s_{q;2r}}{s_{0;22} s_{2;22} - s_{1;22}^2},$$

we then have

$$\frac{\partial^2 \bar{\eta}_1}{\partial x_2^2} = S_{0312}^{0213} (\bar{\zeta}'_3 - x_2 \bar{\zeta}''_2) - S_{1213}^{0223} \bar{\zeta}''_1 + 2S_{1214}^{0224} \bar{\vartheta}' + S_{0412}^{0214} \frac{\partial \bar{\varrho}}{\partial x_2},$$

$$\frac{\partial \bar{\eta}_2}{\partial x_2} = S_{0322}^{1213} (\bar{\zeta}'_3 - x_2 \bar{\zeta}''_2) - S_{2213}^{1223} \bar{\zeta}''_1 + 2S_{1422}^{1224} \bar{\vartheta}' + S_{0422}^{1214} \frac{\partial \bar{\varrho}}{\partial x_2},$$

and taking into account the relations (48), (49), (54), (55), and (57), we find

$$(58) \quad \sum_{j=1}^6 c_{ij} \bar{e}_j = (\hat{c}_{i3} - x_1 \hat{c}_{i2} S_{0312}^{0213} + \hat{c}_{i2} S_{0322}^{1213}) (\bar{\zeta}'_3 - x_2 \bar{\zeta}''_2)$$

$$- (x_1 \hat{c}_{i3} - x_1 \hat{c}_{i2} S_{1213}^{0223} + \hat{c}_{i2} S_{2213}^{1223}) \bar{\zeta}''_1$$

$$+ 2(x_1 \hat{c}_{i4} - x_1 \hat{c}_{i2} S_{1214}^{0224} + \hat{c}_{i2} S_{1422}^{1224}) \bar{\vartheta}'$$

$$+ (\hat{c}_{i4} - x_1 \hat{c}_{i2} S_{0412}^{0214} + \hat{c}_{i2} S_{0422}^{1214}) \frac{\partial \bar{\varrho}}{\partial x_2},$$

for $i = 3, 4$. Now let $v \in \mathcal{V}$ and ϑ and ϱ be as in Lemma 5.2. With $\psi \in L^2(\Omega_{23})$ and $\vartheta = 0$ and

$$\varrho := \int_{-a_2/2}^{x_2} \psi(s, \cdot) ds - \int_{\Omega_2} \int_{-a_2/2}^{x_2} \psi(s, \cdot) ds dx_2,$$

we find from (43) that

$$\int_{\Omega_1} \sum_{j=1}^6 c_{4j} \bar{e}_j dx_1 = 0.$$

It then follows that

$$\begin{aligned} 0 = & (s_{0;43} - s_{1;42} S_{0312}^{0213} + s_{0;42} S_{0322}^{1213}) (\bar{\zeta}'_3 - x_2 \bar{\zeta}''_2) \\ & - (s_{1;43} - s_{1;42} S_{1213}^{0223} + s_{0;42} S_{2213}^{1223}) \bar{\zeta}''_1 \\ & + 2(s_{1;44} - s_{1;42} S_{1214}^{0224} + s_{0;42} S_{1422}^{1224}) \bar{\vartheta}' \\ & + (s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214}) \frac{\partial \bar{\varrho}}{\partial x_2}, \end{aligned}$$

and provided that the last coefficient $s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214} \neq 0$, one finds

$$\begin{aligned} \frac{\partial \bar{\varrho}}{\partial x_2} = & \frac{-1}{s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214}} \\ & \cdot \left[(s_{0;43} - s_{1;42} S_{0312}^{0213} + s_{0;42} S_{0322}^{1213}) (\bar{\zeta}'_3 - x_2 \bar{\zeta}''_2) \right. \\ & - (s_{1;43} - s_{1;42} S_{1213}^{0223} + s_{0;42} S_{2213}^{1223}) \bar{\zeta}''_1 \\ & \left. + (s_{1;44} - s_{1;42} S_{1214}^{0224} + s_{0;42} S_{1422}^{1224}) 2\bar{\vartheta}' \right]. \end{aligned}$$

We may rewrite (58) as

$$(59) \quad \sum_{j=1}^6 c_{ij} \bar{e}_j = F_{i3} \bar{\zeta}'_3 - F_{i2} \bar{\zeta}''_2 - F_{i1} \bar{\zeta}''_1 + F_{i4} \bar{\vartheta}'$$

for $i = 3, 4$, where

$$\begin{aligned} F_{i1} = & (x_1 \hat{c}_{i3} - x_1 \hat{c}_{i2} S_{1213}^{0223} + \hat{c}_{i2} S_{2213}^{1223}) \\ & - \frac{(\hat{c}_{i4} - x_1 \hat{c}_{i2} S_{0412}^{0214} + \hat{c}_{i2} S_{0422}^{1214}) (s_{1;43} - s_{1;42} S_{1213}^{0223} + s_{0;42} S_{2213}^{1223})}{s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214}}, \\ F_{i2} = & x_2 F_{i3}, \\ F_{i3} = & (\hat{c}_{i3} - x_1 \hat{c}_{i2} S_{0312}^{0213} + \hat{c}_{i2} S_{0322}^{1213}) \\ & - \frac{(\hat{c}_{i4} - x_1 \hat{c}_{i2} S_{0412}^{0214} + \hat{c}_{i2} S_{0422}^{1214}) (s_{0;43} - s_{1;42} S_{0312}^{0213} + s_{0;42} S_{0322}^{1213})}{s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214}}, \\ F_{i4} = & 2(x_1 \hat{c}_{i4} - x_1 \hat{c}_{i2} S_{1214}^{0224} + \hat{c}_{i2} S_{1422}^{1224}) \\ & - 2 \frac{(\hat{c}_{i4} - x_1 \hat{c}_{i2} S_{0412}^{0214} + \hat{c}_{i2} S_{0422}^{1214}) (s_{1;44} - s_{1;42} S_{1214}^{0224} + s_{0;42} S_{1422}^{1224})}{s_{0;44} - s_{1;42} S_{0412}^{0214} + s_{0;42} S_{0422}^{1214}}. \end{aligned}$$

Let

$$A_{ij}(x_3) := \int_{\Omega_{12}} F_{ij}(\cdot, \cdot, x_3) dx_1 dx_2$$

and

$$K_{ij}(x_3) := \int_{\Omega_{12}} x_1 F_{ij}(\cdot, \cdot, x_3) dx_1 dx_2, \quad L_{ij}(x_3) := \int_{\Omega_{12}} x_2 F_{ij}(\cdot, \cdot, x_3) dx_1 dx_2.$$

Then, from (43), (44), and (59) we finally deduce the following system of equilibrium differential equations:

$$(60) \quad \begin{cases} (A_{33}\bar{\zeta}_3' - A_{31}\bar{\zeta}_1'' - A_{32}\bar{\zeta}_2'' + A_{34}\bar{\vartheta}')' = -p_3, \\ (K_{33}\bar{\zeta}_3' - K_{31}\bar{\zeta}_1'' - K_{32}\bar{\zeta}_2'' + K_{34}\bar{\vartheta}')'' = -p_1, \\ (L_{33}\bar{\zeta}_3' - L_{31}\bar{\zeta}_1'' - L_{32}\bar{\zeta}_2'' + L_{34}\bar{\vartheta}')'' = -p_2, \\ 2(K_{43}\bar{\zeta}_3' - K_{41}\bar{\zeta}_1'' - K_{42}\bar{\zeta}_2'' + K_{44}\bar{\vartheta}')' = -m, \end{cases}$$

where

$$\begin{aligned} p_1 &= \left(\int_{\Omega_{12}} x_1 b_3 dx_1 dx_2 \right)' + \int_{\Omega_{12}} b_1 dx_1 dx_2, \\ p_2 &= \left(\int_{\Omega_{12}} x_2 b_3 dx_1 dx_2 \right)' + \int_{\Omega_{12}} b_2 dx_1 dx_2, \\ p_3 &= \int_{\Omega_{12}} b_3 dx_1 dx_2. \end{aligned}$$

The system (60) should then be completed with the suitable boundary conditions.

7.1. The homogeneous beam. In the general inhomogeneous and anisotropic case, the torsional, flexional, and extensional problems are all coupled together in the equilibrium differential system (60). A partial decoupling occurs already in the homogeneous fully anisotropic case. Indeed, in this case, from (56), we have

$$s_{0;ij} = \hat{c}_{ij}, \quad s_{1;ij} = 0, \quad s_{2;ij} = \frac{a_1^3}{12} \hat{c}_{ij},$$

and therefore

$$S_{0312}^{0213} = S_{2213}^{1223} = S_{0412}^{0214} = S_{1422}^{1224} = 0, \quad S_{1213}^{0223} = \frac{\hat{c}_{23}}{\hat{c}_{22}}, \quad S_{0322}^{1213} = -\frac{\hat{c}_{23}}{\hat{c}_{22}}, \quad S_{1422}^{1214} = \frac{\hat{c}_{24}}{\hat{c}_{22}},$$

which causes many of the coefficients of the system A_{ij} , K_{ij} , and L_{ij} to be zero. In this case, the system (60) simply rewrites as

$$(61) \quad \begin{cases} A_{33}\bar{\zeta}_3'' = -p_3, \\ (-K_{31}\bar{\zeta}_1'' + K_{34}\bar{\vartheta}')'' = -p_1, \\ -L_{32}\bar{\zeta}_2^{(iv)} = -p_2, \\ 2(-K_{41}\bar{\zeta}_1'' + K_{44}\bar{\vartheta}')' = -m, \end{cases}$$

where

$$\begin{aligned}
 A_{33} &= a_1 a_2 \left[\frac{\hat{c}_{22} \hat{c}_{33} - \hat{c}_{23}^2}{\hat{c}_{22}} + \frac{(\hat{c}_{22} \hat{c}_{34} - \hat{c}_{23} \hat{c}_{24})^2}{\hat{c}_{22}(\hat{c}_{22} \hat{c}_{44} - \hat{c}_{24}^2)} \right], \\
 K_{31} &= \frac{a_1^3 a_2}{12} \frac{\hat{c}_{22} \hat{c}_{33} - \hat{c}_{23}^2}{\hat{c}_{22}}, \quad K_{34} = \frac{a_1^3 a_2}{6} \frac{\hat{c}_{22} \hat{c}_{34} - \hat{c}_{23} \hat{c}_{24}}{\hat{c}_{22}}, \\
 K_{41} &= \frac{a_1^3 a_2}{12} \frac{\hat{c}_{22} \hat{c}_{34} - \hat{c}_{23} \hat{c}_{24}}{\hat{c}_{22}}, \quad K_{44} = \frac{a_1^3 a_2}{6} \frac{\hat{c}_{22} \hat{c}_{44} - \hat{c}_{24}^2}{\hat{c}_{22}}, \\
 L_{32} &= \frac{a_1 a_2^3}{12} \left[\frac{\hat{c}_{22} \hat{c}_{33} - \hat{c}_{23}^2}{\hat{c}_{22}} + \frac{(\hat{c}_{22} \hat{c}_{34} - \hat{c}_{23} \hat{c}_{24})^2}{\hat{c}_{22}(\hat{c}_{22} \hat{c}_{44} - \hat{c}_{24}^2)} \right].
 \end{aligned}$$

Thus for a fully anisotropic but homogeneous beam there is only coupling between twisting and bending in direction 1.

7.2. The homogeneous orthotropic/isotropic beam. When the material is orthotropic a complete decoupling occurs. Indeed, for orthotropic material we have that $c_{ki} = 0$ for $k = 1, 2, 3$ and $i = 4, 5, 6$, and $c_{45} = c_{46} = c_{56} = 0$. It then follows that $K_{34} = K_{41} = 0$ and the system (61) reduces to

$$(62) \quad \begin{cases} A_{33} \bar{\zeta}_3'' = -p_3, \\ -K_{31} \bar{\zeta}_1^{(iv)} = -p_1, \\ -L_{32} \bar{\zeta}_2^{(iv)} = -p_2, \\ 2K_{44} \bar{\vartheta}'' = -m. \end{cases}$$

Finally, if the material is isotropic, that is, if it is orthotropic and $c_{12} = c_{23} = c_{13} =: \lambda$, $c_{44} = c_{55} = c_{66} =: \mu$, and $c_{11} = c_{22} = c_{33} = \lambda + 2\mu$, where λ and μ are the Lamé moduli, then we have that

$$A_{33} = a_1 a_2 E, \quad K_{31} = \frac{a_1^3 a_2}{12} E, \quad 2K_{44} = \frac{a_1^3 a_2}{3} \mu, \quad L_{32} = \frac{a_1 a_2^3}{12} E,$$

where $E := \mu \frac{3\lambda + 2\mu}{\lambda + \mu}$ is the Young modulus of the material. Hence, in the isotropic and homogeneous case we recover the usual form of the differential system of equilibrium equations.

REFERENCES

[1] B. A. AULD, *Acoustic Fields and Waves in Solids*, Vol. 1, John Wiley and Sons, New York, 1973.
 [2] R. CHANDRA, A. D. STEMPLE, AND I. CHOPRA, *Thin-walled composite beams under bending, torsional, and extensional effects*, J. Aircraft, 27 (1990), pp. 619–626.
 [3] L. DELLA LONGA AND A. LONDERO, *Thin Walled Beams with Residual Stress*, submitted.
 [4] L. FREDDI, A. MORASSI, AND R. PARONI, *Thin-walled beams: The case of the rectangular cross-section*, J. Elasticity, 76 (2004), pp. 45–66.
 [5] H. LE DRET, *Problèmes Variationnels dans les Multi-Domaines*, Recherches en Mathématiques Appliquées 19, Masson, Paris, 1991.
 [6] H. LE DRET, *Convergence of displacements and stresses in linearly elastic slender rods as the thickness goes to zero*, Asymptot. Anal., 10 (1995), pp. 367–402.
 [7] R. MONNEAU, F. MURAT, AND A. SILI, *A Partial Korn's Inequality and Error Estimates for the 3d-1d Dimension Reduction in Anisotropic Heterogeneous Linearized Elasticity*, in preparation.

- [8] F. MURAT AND A. SILI, *Anisotropic, Heterogeneous, Linearized Elasticity Problems in Thin Cylinders*, in preparation.
- [9] L. TRABUCHO AND J. M. VIANO, *Mathematical modelling of rods*, in Handbook of Numerical Analysis, Vol. 4, North-Holland, Amsterdam, 1996, pp. 487–974.
- [10] F. TREVES, *Topological Vector Spaces, Distributions and Kernels*, Academic Press, New York, London, 1967.
- [11] V. V. VOLOVOI AND D. H. HODGES, *Theory of anisotropic thin-walled beams*, J. Appl. Mech., 67 (2000), pp. 453–459.
- [12] V. V. VOLOVOI, D. H. HODGES, C. E. S. CESNIK, AND B. POPESCU, *Assessment of beam modeling methods for rotor blade applications*, Math. Comput. Modelling, 33 (2001), pp. 1099–1112.
- [13] R. L. WHEEDEN AND A. ZYGMUND, *Measure and Integral. An Introduction to Real Analysis*, Monogr. Textbooks Pure Appl. Math. 43, Marcel Dekker, New York, Basel, 1977.

CORRECTORS FOR THE HOMOGENIZATION OF A CLASS OF HYPERBOLIC EQUATIONS WITH IMPERFECT INTERFACES*

PATRIZIA DONATO[†], LUISA FAELLA[‡], AND SARA MONSURRO[§]

Abstract. We present here some corrector results for the homogenization of the wave equation in a two-component composite with ε -periodic connected inclusions which complete the homogenization results proved in [P. Donato, L. Faella, and S. Monsurro, *J. Math. Pures Appl.*, 87 (2007), pp. 119–143] by the authors. On the interface separating the two components we prescribe a jump of the solution proportional to the conormal derivatives via a function of order ε^γ , with $-1 < \gamma \leq 1$. Due to different expressions of the energies of the limit problems, the cases $-1 < \gamma < 1$ and $\gamma = 1$ need to be treated separately. The second one, where a memory effect appears in the homogenized problem, is the most interesting. For this critical case, displaying lack of compactness, we in particular establish the central upper semicontinuity type inequalities by splitting a related energy term into a compact part and a part vanishing in appropriate norms.

Key words. homogenization, correctors, hyperbolic equations

AMS subject classifications. 35B27, 35L05, 82B24

DOI. 10.1137/080712684

1. Introduction. In this paper we prove some corrector results for the homogenization of a linear hyperbolic problem in a domain Ω of \mathbb{R}^n made up of two components, a connected one $\Omega_{1\varepsilon}$ and a disconnected one $\Omega_{2\varepsilon}$, which is a union of ε -periodic connected inclusions of size ε . The conditions prescribed on the interface $\Gamma^\varepsilon := \partial\Omega_{2\varepsilon}$ between the two components are the continuity of the conormal derivatives and a jump of the solution proportional to the conormal derivatives via a function of order ε^γ . We suppose here that $-1 < \gamma \leq 1$. This work completes the corresponding homogenization results proved by the authors in [9].

Let A be a periodic, symmetric, bounded, and elliptic matrix field and h a bounded and periodic function. We consider the following problem:

$$(1.1) \quad \begin{cases} u_\varepsilon'' - \operatorname{div}(A^\varepsilon \nabla u_\varepsilon) = f_\varepsilon & \text{in } (\Omega_{1\varepsilon} \times \Omega_{2\varepsilon}) \times]0, T[, \\ [A^\varepsilon \nabla u_\varepsilon] \cdot n_{1\varepsilon} = 0 & \text{on } \Gamma^\varepsilon \times]0, T[, \\ A^\varepsilon \nabla u_{1\varepsilon} \cdot n_{1\varepsilon} = -\varepsilon^\gamma h^\varepsilon[u_\varepsilon] & \text{on } \Gamma^\varepsilon \times]0, T[, \\ u_\varepsilon = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_\varepsilon(0) = U_\varepsilon^0 & \text{in } \Omega, \\ u_\varepsilon'(0) = U_\varepsilon^1 & \text{in } \Omega, \end{cases}$$

where $A^\varepsilon(x) := A(x/\varepsilon)$, $h^\varepsilon(x) := h(x/\varepsilon)$, $u_\varepsilon = (u_{1\varepsilon}, u_{2\varepsilon})$ is defined in $\Omega_{1\varepsilon} \times \Omega_{2\varepsilon}$, $[\cdot]$ denotes the jump trough Γ^ε , and $n_{i\varepsilon}$ is the unitary outward normal to $\Omega_{i\varepsilon}$, $i = 1, 2$.

*Received by the editors January 7, 2008; accepted for publication (in revised form) September 2, 2008; published electronically January 21, 2009.

<http://www.siam.org/journals/sima/40-5/71268.html>

[†]Laboratoire de Mathématiques Raphaël Salem, Université de Rouen, Tecnopôle du Madrillet, Avenue de l'Université, B.P. 12, 76801 Saint Etienne du Rouvray, France (Patrizia.Donato@univ-rouen.fr).

[‡]Dipartimento di Automazione, Elettromagnetismo, Ingegneria, dell'Informazione e Matematica Industriale, Università di Cassino, via G. Di Biasio n.43, 03043 Cassino (FR), Italy (l.faella@unicas.it).

[§]Dipartimento di Matematica ed Informatica, Università degli Studi di Salerno, via Ponte don Melillo, 84084, Fisciano (SA), Italy (smonsurro@unisa.it).

This problem models the wave propagation in a medium made up of two materials with very different coefficients of propagation, which gives rise to the jump in the boundary condition on the interface. For the physical model we refer to [4], where these kinds of interface conditions are derived.

The homogenization results proved in [9] (recalled in section 2) show two different asymptotic behaviors for $-1 < \gamma < 1$ and $\gamma = 1$, the last case being the most interesting one, since a memory effect appears at the limit (see also Remark 2.7).

When $-1 < \gamma < 1$, under suitable assumptions on the data (see Theorem 2.3), we have, as $\varepsilon \rightarrow 0$, the following convergences:

$$\begin{cases} P_1^\varepsilon u_{1\varepsilon} \rightharpoonup u_1 & \text{weakly* in } L^\infty(0, T; H_0^1(\Omega)), \\ P_1^\varepsilon u'_{1\varepsilon} \rightharpoonup u'_1 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

while

$$\begin{cases} \widetilde{u_{2\varepsilon}} \rightharpoonup \theta_2 u_1 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \\ \widetilde{u'_{2\varepsilon}} \rightharpoonup \theta_2 u'_1 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

and

$$(1.2) \quad \begin{cases} A^\varepsilon \widetilde{\nabla u_{1\varepsilon}} \rightharpoonup A^0 \nabla u_1 & \text{weakly* in } L^\infty(0, T; [L^2(\Omega)]^n), \\ A^\varepsilon \widetilde{\nabla u_{2\varepsilon}} \rightharpoonup 0 & \text{weakly* in } L^\infty(0, T; [L^2(\Omega)]^n), \end{cases}$$

where P_1^ε is a suitable extension operator, θ_2 is the proportion of material occupying the inclusions, $\widetilde{}$ denotes the zero extension to the whole of Ω , and u_1 is the solution of the homogenized problem

$$\begin{cases} u_1'' - \operatorname{div} (A^0 \nabla u_1) = f_1 + f_2 & \text{in } \Omega \times]0, T[, \\ u_1 = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_1(0) = U_0 & \text{in } \Omega, \\ u'_1(0) = U_1 & \text{in } \Omega. \end{cases}$$

The matrix field A^0 is the same constant positive definite matrix obtained by Cioranescu and Saint Jean Paulin (see [7] and also [6]) for the homogenization of the elliptic problem in the perforated domain $\Omega_{1\varepsilon}$, with a Neumann condition on the boundary of the holes.

All of these convergences are weak, but sufficient, to obtain the homogenized problem. At this point, the question is how to improve them, in particular give more information on the gradient and on the first time derivative of the solution. Such results are known in standard homogenization as corrector results. The aim of this paper is precisely to answer the above questions.

The corrector result given in Theorem 2.4 states that, under some additional assumptions,

$$\begin{cases} \widetilde{u'_{1\varepsilon}} + \widetilde{u'_{2\varepsilon}} \rightarrow u'_1 & \text{strongly in } C^0(0, T; L^2(\Omega)), \\ \lim_{\varepsilon \rightarrow 0} \|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0(0, T; [L^1(\Omega_{1\varepsilon})]^n)} = 0, \\ \lim_{\varepsilon \rightarrow 0} \|\nabla u_{2\varepsilon}\|_{C^0(0, T; [L^2(\Omega_{2\varepsilon})]^n)} = 0, \end{cases}$$

where C^ε is the same corrector matrix of the problem studied in [7].

The additional assumptions concern the data. We assume strong convergence for both f_ε and U_ε^1 , and we suppose that U_ε^0 solves the particular elliptic equation given by (2.10). These assumptions are needed for the convergence of the energy of problem (1.1) to that of the homogenized one. This convergence is necessary to prove the corrector results, as already evidenced in the classical homogenization of the wave equation (see [2]) or in a perforated domain with a Neumann condition on the boundaries of the holes (see [18]).

If $\gamma = 1$, the situation is more complicated since the limit problem consists in a p.d.e. coupled with an o.d.e., as proved in [9]. One still has the convergences

$$\begin{cases} P_1^\varepsilon u_{1\varepsilon} \rightharpoonup u_1 & \text{weakly* in } L^\infty(0, T; H_0^1(\Omega)), \\ P_1^\varepsilon u'_{1\varepsilon} \rightharpoonup u'_1 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

and

$$\begin{cases} \widetilde{u_{2\varepsilon}} \rightharpoonup u_2 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \\ \widetilde{u'_{2\varepsilon}} \rightharpoonup u'_2 & \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

but u_2 is not $\theta_2 u'_1$ anymore. In this case (see Theorem 2.6), (u_1, u_2) is the unique solution of the coupled problem

$$\begin{cases} \theta_1 u''_1 - \operatorname{div}(A^0 \nabla u_1) + c_h(\theta_2 u_1 - u_2) = f_1 & \text{in } \Omega \times]0, T[, \\ u''_2 - c_h(\theta_2 u_1 - u_2) = f_2 & \text{in } \Omega \times]0, T[, \\ u_1 = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_1(x, 0) = U^0, \quad u_2(x, 0) = \theta_2 U^0 & \text{in } \Omega, \\ u'_1(x, 0) = U^1, \quad u'_2(x, 0) = \theta_2 U^1 & \text{in } \Omega, \end{cases}$$

where $\theta_1 = 1 - \theta_2$ and c_h is a constant depending on the function h . Furthermore, convergences (1.2) still hold. Let us mention that in [9, Remark 2.11] the limit function u_2 is explicitly computed in terms of u_1 and it is shown that u_1 is the solution of a linear hyperbolic problem with a linear memory effect.

Under the same assumptions as in the previous case, the corrector result for $\gamma = 1$ (Theorem 2.8) states the convergences

$$\begin{cases} \lim_{\varepsilon \rightarrow 0} \|u'_{1\varepsilon} - u'_1\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))} = 0, \\ \lim_{\varepsilon \rightarrow 0} \|u'_{2\varepsilon} - \theta_2^{-1} u'_2\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))} = 0, \\ \lim_{\varepsilon \rightarrow 0} \|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0(0, T; [L^1(\Omega_{1\varepsilon})]^n)} = 0, \\ \lim_{\varepsilon \rightarrow 0} \|\nabla u_{2\varepsilon}\|_{C^0(0, T; [L^2(\Omega_{2\varepsilon})]^n)} = 0. \end{cases}$$

The corrector results are proved in section 5 for $-1 < \gamma < 1$ and in section 6 for $\gamma = 1$.

Let us point out the main difficulties of our situation due to the presence of the interface. The first concerns the study of the convergence of the energy associated to problem (1.1). This is in general straightforward, but here the expression of the energy of problem (1.1) is more complicated than usual, since it contains a boundary term. Moreover, for $\gamma = 1$, the energy of the homogenized coupled system contains also a zero order term which is a linear combination of u_1 and u_2 . This can be seen comparing formulas (4.5) (case $-1 < \gamma < 1$) and (4.8) (case $\gamma = 1$). We can prove in section 4 (Theorems 4.3 and 4.4) that in both cases the energy converges to those of the homogenized problems.

The second and more difficult point consists in proving two upper semicontinuity type inequalities given for the two cases in Propositions 5.1 and 6.1, respectively. They require specific arguments and are crucial, since they allow us to end, by a density argument, the proof of the corrector result. Let us mention that usually in the literature one obtains a convergence, which gives rise to an equality for the limit (see Remark 5.2), but here we can only obtain an inequality on an upper limit. Nevertheless, this is enough to conclude.

Let us briefly describe Proposition 6.1 concerning the case $\gamma = 1$, which is the most interesting one. It states that if we set, for $\Phi, \Psi \in C^\infty([0, T], \mathcal{D}(\Omega))$,

$$\hat{X}_\varepsilon = \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon} - \Phi'|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon} - \Psi'|^2 dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon (\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi) (\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi) dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon} \nabla u_{2\varepsilon} dx \right],$$

then

$$(1.3) \quad \limsup_{\varepsilon \rightarrow 0} \|\hat{X}_\varepsilon\|_{C^0([0, T])} \leq \|\hat{X}\|_{C^0([0, T])},$$

where

$$\hat{X} = \frac{1}{2} \left[\theta_1 \|u'_1 - \Phi'\|_{L^2(\Omega)}^2 + \theta_2^{-1} \|u'_2 - \theta_2 \Psi'\|_{L^2(\Omega)}^2 + \int_{\Omega} A^0 (\nabla u_1 - \nabla \Phi) (\nabla u_1 - \nabla \Phi) dx \right].$$

The main difficulty when proving this result is due to the boundary term in the energy associated to problem (1.1). Indeed, we cannot have estimates for the time derivative of \hat{X}_ε so that we cannot derive any compactness of \hat{X}_ε in $C^0([0, T])$. To overcome this difficulty, in Step 4 of the proof of Proposition 6.1 we decompose \hat{X}_ε as a sum of two terms and we show that the first one is compact and the second one goes to zero in $C^0([0, T])$ (see also Remarks 6.2 and 6.3 for details). This allows us to prove (1.3).

The correctors for the corresponding elliptic case have been studied by the first author in [8]. The first homogenization result for these kinds of boundary conditions was done in the elliptic case, for some values of the parameter γ , by Auriault and Ene [1] by the multiple scales method. We refer to Lipton [15] for the study of the limit problem when $\gamma = 0$, to S. Monsurrò [17] for the case $\gamma \leq -1$, and to [10] for the case $\gamma > -1$. For similar elliptic homogenization problems we refer also to [3], [12], [13], [16], [19].

2. Formulation of the problem and main results. Along this paper, Ω will denote an open bounded subset of \mathbb{R}^n and $\{\varepsilon\}$ a sequence of positive real numbers converging to zero.

Let $Y =]0, l_1[\times \dots \times]0, l_n[$ be a reference cell in \mathbb{R}^n and Y_1 and Y_2 be two nonempty open and disjoint subsets of Y such that

$$Y = Y_1 \cup \overline{Y_2},$$

with Y_1 connected and $\Gamma := \partial Y_2$ of class C^2 .

For any $k \in \mathbb{Z}^n$ we denote

$$Y_i^k := k_l + Y_i, \quad \Gamma_k := k_l + \Gamma,$$

where $k_l = (k_1 l_1, \dots, k_n l_n)$ and $i = 1, 2$, and we suppose that

$$(2.1) \quad \partial\Omega \cap \left(\bigcup_{k \in \mathbb{Z}^n} (\varepsilon \Gamma_k) \right) := \emptyset.$$

We introduce then, for any ε , the two components of Ω and the interface, respectively, by

$$\Omega_{i\varepsilon} := \Omega \cap \left\{ \bigcup_{k \in K_\varepsilon} \varepsilon Y_i^k \right\}, \quad i = 1, 2, \quad \text{and} \quad \Gamma^\varepsilon := \partial\Omega_{2\varepsilon},$$

where K_ε is the set of the n -tuples such that $\varepsilon\Gamma_k$ is included in Ω , namely,

$$K_\varepsilon := \{k \in Z^n \mid \varepsilon\Gamma_k \cap \Omega \neq \emptyset\},$$

so that $\partial\Omega \cap \Gamma^\varepsilon = \emptyset$.

Due to (2.1), the set $\Omega_{1\varepsilon}$ is connected while $\Omega_{2\varepsilon}$ is a union of $c\varepsilon^{-n}$ disjoint translated sets of εY_2 , c being a constant independent of ε .

In what follows, we will denote by χ_ω the characteristic function of any open set $\omega \subset \mathbb{R}^n$. We know that (see, for instance, [6])

$$(2.2) \quad \chi_{\Omega_{i\varepsilon}} \rightharpoonup \theta_i := \frac{|Y_i|}{|Y|} \quad \text{weakly in } L^2(\Omega), \quad i = 1, 2.$$

We define the normed space V^ε by

$$V^\varepsilon := \{v_1 \in H^1(\Omega_{1\varepsilon}) \mid v_1 = 0\}$$

endowed with the norm

$$\|v_1\|_{V^\varepsilon} := \|\nabla v_1\|_{L^2(\Omega_{1\varepsilon})}.$$

Remark 2.1. Since we do not assume any regularity on $\partial\Omega$, the condition on $\partial\Omega$ in the definition of V^ε has to be understood in a density sense. More precisely, V^ε is the closure with respect to the $H^1(\Omega_{1\varepsilon})$ -norm of the set of the functions of $C^\infty(\Omega_{1\varepsilon})$ with a compact support contained in Ω . This make sense because of (2.1).

As in [9], we also define, for every $\gamma \in \mathbb{R}$,

$$H_\gamma^\varepsilon := \{v = (v_1, v_2) \mid v_1 \in V^\varepsilon \quad \text{and} \quad v_2 \in H^1(\Omega_{2\varepsilon})\}$$

equipped with the norm

$$\|v\|_{H_\gamma^\varepsilon}^2 := \|\nabla v_1\|_{L^2(\Omega_{1\varepsilon})}^2 + \|\nabla v_2\|_{L^2(\Omega_{2\varepsilon})}^2 + \varepsilon^\gamma \|v_1 - v_2\|_{L^2(\Gamma^\varepsilon)}^2.$$

Obviously, if $0 < \varepsilon < 1$ and $\gamma_1 \leq \gamma_2$, then

$$\|v\|_{H_{\gamma_2}^\varepsilon} \leq \|v\|_{H_{\gamma_1}^\varepsilon}.$$

Moreover, as shown in [17, Lemmas 2.7 and 2.8 and their proofs], for every fixed ε the norms of H_γ^ε and $V^\varepsilon \times H^1(\Omega_{2\varepsilon})$ are equivalent.

To introduce the coefficient matrix, we define, for $\alpha, \beta \in \mathbb{R}$ with $0 < \alpha < \beta$, the set $M(\alpha, \beta, Y)$ of the $n \times n$ Y -periodic matrix-valued functions in $L^\infty(Y)$ such that

$$(A(x)\lambda, \lambda) \geq \alpha|\lambda|^2, \quad |A(x)\lambda| \leq \beta\lambda$$

for any $\lambda \in \mathbb{R}^n$ and a.e. in Y .

We assume that

$$(2.3) \quad A \in M(\alpha, \beta, Y), \quad A \text{ symmetric,}$$

and we set, for any $\varepsilon > 0$,

$$(2.4) \quad A^\varepsilon(x) := A(x/\varepsilon).$$

Moreover, we consider a Y -periodic function h such that

$$h \in L^\infty(\Gamma) \text{ and } \exists h_0 \in \mathbb{R} \text{ such that } 0 < h_0 < h(y) \text{ a.e. in } \Gamma$$

and set

$$(2.5) \quad h^\varepsilon(x) := h\left(\frac{x}{\varepsilon}\right).$$

Throughout this work, we suppose $T > 0$ be given and $-1 < \gamma \leq 1$. We consider the following problem:

$$(2.6) \quad \begin{cases} u''_{1\varepsilon} - \operatorname{div}(A^\varepsilon \nabla u_{1\varepsilon}) = f_{1\varepsilon} & \text{in } \Omega_{1\varepsilon} \times]0, T[, \\ u''_{2\varepsilon} - \operatorname{div}(A^\varepsilon \nabla u_{2\varepsilon}) = f_{2\varepsilon} & \text{in } \Omega_{2\varepsilon} \times]0, T[, \\ A^\varepsilon \nabla u_{1\varepsilon} \cdot n_{1\varepsilon} = -A^\varepsilon \nabla u_{2\varepsilon} \cdot n_{2\varepsilon} & \text{on } \Gamma^\varepsilon \times]0, T[, \\ A^\varepsilon \nabla u_{1\varepsilon} \cdot n_{1\varepsilon} = -\varepsilon^\gamma h^\varepsilon(u_{1\varepsilon} - u_{2\varepsilon}) & \text{on } \Gamma^\varepsilon \times]0, T[, \\ u_{1\varepsilon} = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_{1\varepsilon}(0) = U_{1\varepsilon}^0 & \text{in } \Omega_{1\varepsilon}, \quad u_{2\varepsilon}(0) = U_{2\varepsilon}^0 & \text{in } \Omega_{2\varepsilon}, \\ u'_{1\varepsilon}(0) = U_{1\varepsilon}^1 & \text{in } \Omega_{1\varepsilon}, \quad u'_{2\varepsilon}(0) = U_{2\varepsilon}^1 & \text{in } \Omega_{2\varepsilon}, \end{cases}$$

where $n_{i\varepsilon}$ is the unitary outward normal to $\Omega_{i\varepsilon}$, $i = 1, 2$, and

$$(2.7) \quad \begin{cases} \text{(i) } f_\varepsilon := (f_{1\varepsilon}|_{\Omega_{1\varepsilon}}, f_{2\varepsilon}|_{\Omega_{2\varepsilon}}) \text{ with } f_{i\varepsilon} \in L^2(0, T; L^2(\Omega)), \quad i = 1, 2, \\ \text{(ii) } U_\varepsilon^0 := (U_{1\varepsilon}^0, U_{2\varepsilon}^0) \in V^\varepsilon \times H^1(\Omega_{2\varepsilon}), \\ \text{(iii) } U_\varepsilon^1 \in L^2(\Omega). \end{cases}$$

Its variational (weak) formulation is

$$(2.8) \quad \begin{cases} \text{find } u_\varepsilon = (u_{1\varepsilon}, u_{2\varepsilon}) \text{ in } L^2(0, T; V^\varepsilon) \times L^2(0, T; H^1(\Omega_{2\varepsilon})) \\ \text{such that } u'_{1\varepsilon} \in L^2(0, T; L^2(\Omega_{1\varepsilon})), u'_{2\varepsilon} \in L^2(0, T; L^2(\Omega_{2\varepsilon})) \text{ and} \\ \langle u''_{1\varepsilon}, v_1 \rangle_{(V^\varepsilon)', V^\varepsilon} + \langle u''_{2\varepsilon}, v_2 \rangle_{(H^1(\Omega_{2\varepsilon}))', H^1(\Omega_{2\varepsilon})} \\ + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon} \cdot \nabla v_1 \, dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon} \cdot \nabla v_2 \, dx \\ + \varepsilon^\gamma \int_{\Gamma^\varepsilon} h^\varepsilon(u_{1\varepsilon} - u_{2\varepsilon})(v_1 - v_2) \, d\sigma_x = \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} v_1 \, dx + \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} v_2 \, dx \\ \text{for every } (v_1, v_2) \in V^\varepsilon \times H^1(\Omega_{2\varepsilon}) \text{ in } \mathcal{D}'(0, T), \\ u_{1\varepsilon}(0) = U_{1\varepsilon}^0 \text{ in } \Omega_{1\varepsilon}, \quad u_{2\varepsilon}(0) = U_{2\varepsilon}^0 \text{ in } \Omega_{2\varepsilon}, \\ u'_{1\varepsilon}(0) = U_{1\varepsilon}^1 & \text{in } \Omega_{1\varepsilon}, \quad u'_{2\varepsilon}(0) = U_{2\varepsilon}^1 & \text{in } \Omega_{2\varepsilon}. \end{cases}$$

The homogenization of this problem has been studied in [9] under more general assumptions than (2.7), obtaining different limit problems according to different values of γ . We prove here some corrector results for the two cases $-1 < \gamma < 1$ and $\gamma = 1$.

It is well known that (see, for instance, [2] and [18]) corrector results for the wave equation need stronger assumptions than that of the convergence results. In particular, we suppose that

$$(2.9) \quad \begin{cases} \text{(i) } (f_{1\varepsilon}, f_{2\varepsilon}) \rightarrow (\theta_1^{-1} f_1, \theta_2^{-1} f_2) \text{ strongly in } [L^2(0, T; L^2(\Omega))]^2, \\ \text{(ii) } U_\varepsilon^1 \rightarrow U^1 \text{ strongly in } L^2(\Omega) \end{cases}$$

and, concerning the initial condition U_ε^0 , we assume that it is the unique solution of the problem

$$(2.10) \quad \begin{cases} -\operatorname{div} (A^\varepsilon \nabla U_{1\varepsilon}^0) = Q_1^{\varepsilon*} (-\operatorname{div} (A^0 \nabla U^0)) & \text{in } \Omega_{1\varepsilon}, \\ -\operatorname{div} (A^\varepsilon \nabla U_{2\varepsilon}^0) = 0 & \text{in } \Omega_{2\varepsilon}, \\ A^\varepsilon \nabla U_{1\varepsilon}^0 \cdot n_{1\varepsilon} = -A^\varepsilon \nabla U_{2\varepsilon}^0 \cdot n_{2\varepsilon} & \text{on } \Gamma^\varepsilon, \\ A^\varepsilon \nabla U_{1\varepsilon}^0 \cdot n_{1\varepsilon} = -\varepsilon^\gamma h^\varepsilon (U_{1\varepsilon}^0 - U_{2\varepsilon}^0) & \text{on } \Gamma^\varepsilon, \\ U_{1\varepsilon}^0 = 0 & \text{on } \partial\Omega, \end{cases}$$

where U^0 is given in $H_0^1(\Omega)$ and $Q_1^{\varepsilon*}$ is the adjoint of a suitable extension operator Q_1^ε , defined in section 3.

Remark 2.2. The homogenization result proved in [8] for the elliptic case (recalled in section 3 for the reader’s convenience) applies to (2.10) and gives (see the proof of Lemma 4.2 for details) the following convergences:

$$(2.11) \quad \begin{cases} \text{(i)} \quad Q_1^\varepsilon U_{1\varepsilon}^0 \rightharpoonup U^0 & \text{weakly in } H_0^1(\Omega), \\ \text{(ii)} \quad A^\varepsilon \widetilde{\nabla U_{1\varepsilon}^0} \rightharpoonup A^0 \nabla U^0 & \text{weakly in } [L^2(\Omega)]^n, \\ \text{(iii)} \quad A^\varepsilon \widetilde{\nabla U_{2\varepsilon}^0} \rightharpoonup 0 & \text{weakly in } [L^2(\Omega)]^n, \\ \text{(iv)} \quad \|U_\varepsilon^0\|_{H_\gamma^\varepsilon} \leq C & \text{with } C \text{ independent of } \varepsilon, \\ \text{(v)} \quad \widetilde{U_{2\varepsilon}^0} \rightharpoonup \theta_2 U^0 & \text{weakly in } H_0^1(\Omega), \end{cases}$$

where $\widetilde{}$ denotes the zero extension to the whole of Ω and the homogenized matrix A^0 is defined by

$$(2.12) \quad A^0 \lambda = \frac{1}{|Y_1|} \int_{Y_1} A \nabla w_\lambda \, dy$$

with $w_\lambda \in H^1(Y_1)$ solution, for any $\lambda \in \mathbb{R}^n$, of

$$(2.13) \quad \begin{cases} -\operatorname{div} (A \nabla w_\lambda) = 0 & \text{in } Y_1, \\ (A \nabla w_\lambda) \cdot n_1 = 0 & \text{on } \Gamma, \\ w_\lambda - \lambda \cdot y & Y\text{-periodic}, \\ \frac{1}{|Y_1|} \int_{Y_1} (w_\lambda - \lambda \cdot y) \, dy = 0. \end{cases}$$

Moreover, from (2.9) and (2.11)(i) and (v) one has

$$\begin{cases} \text{(i)} \quad (\chi_{\Omega_{1\varepsilon}} f_{1\varepsilon}, \chi_{\Omega_{2\varepsilon}} f_{2\varepsilon}) \rightharpoonup (f_1, f_2) & \text{weakly in } [L^2(0, T; L^2(\Omega))]^2, \\ \text{(ii)} \quad (\chi_{\Omega_{1\varepsilon}} U_\varepsilon^1, \chi_{\Omega_{2\varepsilon}} U_\varepsilon^1) \rightharpoonup (\theta_1 U^1, \theta_2 U^1) & \text{weakly in } L^2(\Omega) \times L^2(\Omega), \\ \text{(iii)} \quad (\widetilde{U_{1\varepsilon}^0}, \widetilde{U_{2\varepsilon}^0}) \rightharpoonup (\theta_1 U^0, \theta_2 U^0) & \text{weakly in } L^2(\Omega) \times L^2(\Omega). \end{cases}$$

Let us describe first the case $-1 < \gamma < 1$, for which the homogenization result proved in [9] applies and, under the above assumptions, reads as follows.

THEOREM 2.3 (see [9]). *For $-1 < \gamma < 1$, let A^ε and h^ε be defined by (2.4) and (2.5), respectively, suppose that (2.7), (2.9), and (2.10) hold, and let u_ε be the solution of problem (2.6). Then, there exist a constant $C > 0$ (independent of ε) and an*

extension operator $P_1^\varepsilon \in \mathcal{L}(L^\infty(0, T; H^k(\Omega_{1\varepsilon})); L^\infty(0, T; H^k(\Omega)))$, for $k = 1, 2$, such that

$$(2.14) \quad \begin{cases} \text{(i)} & P_1^\varepsilon u_{1\varepsilon} \rightharpoonup u_1 \quad \text{weakly}^* \text{ in } L^\infty(0, T; H_0^1(\Omega)), \\ \text{(ii)} & \widetilde{u_{2\varepsilon}} \rightharpoonup u_2 := \theta_2 u_1 \quad \text{weakly}^* \text{ in } L^\infty(0, T; L^2(\Omega)), \\ \text{(iii)} & \varepsilon^{\gamma/2} \|u_{1\varepsilon} - u_{2\varepsilon}\|_{L^\infty(0, T; L^2(\Gamma^\varepsilon))} < C, \end{cases}$$

$$(2.15) \quad \begin{cases} \text{(i)} & P_1^\varepsilon u'_{1\varepsilon} \rightharpoonup u'_1 \quad \text{weakly}^* \text{ in } L^\infty(0, T; L^2(\Omega)), \\ \text{(ii)} & \widetilde{u'_{2\varepsilon}} \rightharpoonup u'_2 = \theta_2 u'_1 \quad \text{weakly}^* \text{ in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

$$(2.16) \quad \begin{cases} \text{(i)} & A^\varepsilon \widetilde{\nabla u_{1\varepsilon}} \rightharpoonup A^0 \nabla u_1 \quad \text{weakly}^* \text{ in } L^\infty(0, T; [L^2(\Omega)]^n), \\ \text{(ii)} & A^\varepsilon \widetilde{\nabla u_{2\varepsilon}} \rightharpoonup 0 \quad \text{weakly}^* \text{ in } L^\infty(0, T; [L^2(\Omega)]^n), \end{cases}$$

where θ_2 is given by (2.2) and u_1 is the unique solution in $L^2(0, T; H_0^1(\Omega))$, with u'_1 in $L^2(0, T; L^2(\Omega))$, of the problem

$$(2.17) \quad \begin{cases} u''_1 - \operatorname{div}(A^0 \nabla u_1) = f_1 + f_2 & \text{in } \Omega \times]0, T[, \\ u_1 = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_1(0) = U^0 & \text{in } \Omega, \\ u'_1(0) = U^1 & \text{in } \Omega, \end{cases}$$

with A^0 defined in (2.12).

Let $(e_j)_{j=1, \dots, n}$ be the canonical basis and $w_j \in H^1(Y_1)$ be the solution of problem (2.13), written for $\lambda = e_j$, for $j = 1, \dots, n$. The corrector matrix C^ε is defined, for any ε , by

$$(2.18) \quad C^\varepsilon(x) = \widetilde{C}\left(\frac{x}{\varepsilon}\right),$$

where

$$C_{ij}(y) := \frac{\partial w_j}{\partial y_i}(y), \quad \text{for } i, j = 1, \dots, n,$$

and $\widetilde{}$ denotes the zero extension to the whole reference cell Y . Let us mention that, as in the elliptic case (see [8]), C^ε is the same corrector obtained in the case of a periodic perforated domain with an homogeneous Neumann condition on the holes studied in [7]. Here, due to the symmetry of A , the matrix C^ε is symmetric too.

Our first corrector result is as follows.

THEOREM 2.4 (correctors for $-1 < \gamma < 1$). *Let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7), (2.9), and (2.10) hold. If u_ε is the solution of problem (2.6) with $-1 < \gamma < 1$ and u_1 is the solution of problem (2.17), then*

$$(2.19) \quad \begin{cases} \text{(i)} & \widetilde{u'_{1\varepsilon}} + \widetilde{u'_{2\varepsilon}} \rightarrow u'_1 \quad \text{strongly in } C^0(0, T; L^2(\Omega)), \\ \text{(ii)} & \lim_{\varepsilon \rightarrow 0} \|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0(0, T; [L^1(\Omega_{1\varepsilon})]^n)} = 0, \\ \text{(iii)} & \lim_{\varepsilon \rightarrow 0} \|\nabla u_{2\varepsilon}\|_{C^0(0, T; [L^2(\Omega_{2\varepsilon})]^n)} = 0. \end{cases}$$

This theorem will be proved in section 5.

Remark 2.5. Convergence (2.19)(i) is equivalent to

$$\lim_{\varepsilon \rightarrow 0} \|u'_{i\varepsilon} - u'_1\|_{C^0(0, T; L^2(\Omega_{i\varepsilon}))} = 0, \quad i = 1, 2,$$

since

$$\int_{\Omega} |\widetilde{u}'_{1\varepsilon} + \widetilde{u}'_{2\varepsilon} - u'_1|^2 dx = \int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon} - u'_1|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon} - u'_1|^2 dx,$$

$\Omega_{1\varepsilon}$ and $\Omega_{2\varepsilon}$ being disjoint sets.

We present now the more complicated case $\gamma = 1$, which is also the more interesting, since the limit function u_2 is not $\theta_2 u_1$ anymore and the limit problem satisfied by (u_1, u_2) is a coupled system consisting in a p.d.e. and an o.d.e. Indeed, we have the following theorem.

THEOREM 2.6 (see [9]). *For $\gamma = 1$, let A^ε and h^ε be defined by (2.4) and (2.5), respectively, suppose that (2.7), (2.10), and (2.9) hold, and let u_ε be the solution of the problem (2.6). Then, there exist a constant $C > 0$ (independent of ε) and an extension operator $P_1^\varepsilon \in \mathcal{L}(L^\infty(0, T; H^k(\Omega_{1\varepsilon}); L^\infty(0, T; H^k(\Omega)))$, for $k = 1, 2$, such that*

$$(2.20) \quad \begin{cases} \text{(i)} & P_1^\varepsilon u_{1\varepsilon} \rightharpoonup u_1 \quad \text{weakly* in } L^\infty(0, T; H_0^1(\Omega)), \\ \text{(ii)} & \widetilde{u}_{2\varepsilon} \rightharpoonup u_2 \quad \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \\ \text{(iii)} & \varepsilon^{1/2} \|u_{1\varepsilon} - u_{2\varepsilon}\|_{L^\infty(0, T; L^2(\Gamma^\varepsilon))} < C, \end{cases}$$

$$(2.21) \quad \begin{cases} \text{(i)} & P_1^\varepsilon u'_{1\varepsilon} \rightharpoonup u'_1 \quad \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \\ \text{(ii)} & \widetilde{u}'_{2\varepsilon} \rightharpoonup u'_2 \quad \text{weakly* in } L^\infty(0, T; L^2(\Omega)), \end{cases}$$

$$(2.22) \quad \begin{cases} \text{(i)} & A^\varepsilon \widetilde{\nabla} u_{1\varepsilon} \rightharpoonup A^0 \nabla u_1 \quad \text{weakly* in } L^\infty(0, T; [L^2(\Omega)]^n), \\ \text{(ii)} & A^\varepsilon \widetilde{\nabla} u_{2\varepsilon} \rightharpoonup 0 \quad \text{weakly* in } L^\infty(0, T; [L^2(\Omega)]^n), \end{cases}$$

where the couple (u_1, u_2) is the unique solution $L^2(0, T; H_0^1(\Omega)) \times L^2(0, T; L^2(\Omega))$, with (u'_1, u'_2) in $L^2(0, T; L^2(\Omega)) \times L^2(0, T; L^2(\Omega))$, of the problem

$$(2.23) \quad \begin{cases} \theta_1 u''_1 - \operatorname{div}(A^0 \nabla u_1) + c_h(\theta_2 u_1 - u_2) = f_1 & \text{in } \Omega \times]0, T[, \\ u''_2 - c_h(\theta_2 u_1 - u_2) = f_2 & \text{in }]0, T[\text{ for a.e. } x \in \Omega, \\ u_1 = 0 & \text{on } \partial\Omega \times]0, T[, \\ u_1(0) = U^0, \quad u_2(0) = \theta_2 U^0 & \text{in } \Omega, \\ u'_1(0) = U^1, \quad u'_2(0) = \theta_2 U^1 & \text{in } \Omega, \end{cases}$$

where θ_i , for $i = 1, 2$, are given by (2.2), $c_h = \frac{1}{|Y_2|} \int_{\Gamma} h(y) d\sigma_y > 0$, and the homogenized matrix A^0 is defined by (2.12).

Remark 2.7. (i) In [9], it is shown that u_1 is the solution of the equation

$$\theta_1 u''_1 - \operatorname{div}(A^0 \nabla u_1) + c_h \theta_2 u_1 - c_h^2 \theta_2 \int_0^t \mathcal{K}(t, s) u_1(s) ds = F \text{ in } \Omega \times]0, T[,$$

containing a periodic memory kernel \mathcal{K} ; moreover, the limit solution u_2 and the function F are explicitly computed. In this paper we will need only to use the limit problem under the form (2.23), which is well adapted to the study of the energies.

(ii) Observe that convergences (2.14)–(2.15) and (2.20)–(2.21) imply that for $-1 < \gamma \leq 1$

$$\chi_{\Omega_{1\varepsilon}} P_1^\varepsilon u'_{1\varepsilon} \rightharpoonup \theta_1 u'_1 \quad \text{weakly* in } L^\infty(0, T; L^2(\Omega)).$$

Moreover, from classical compactness results (see [14])

$$P_1^\varepsilon u_{1\varepsilon} \rightarrow u_1 \text{ strongly in } C^0(0, T; L^2(\Omega)).$$

Let us state the corrector result associated to the case $\gamma = 1$.

THEOREM 2.8 (correctors for $\gamma = 1$). *Let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7), (2.9), and (2.10) hold. If u_ε is the solution of problem (2.6) with $\gamma = 1$ and (u_1, u_2) is the solution of problem (2.23), then*

$$(2.24) \quad \begin{cases} \text{(i)} \lim_{\varepsilon \rightarrow 0} \|u'_{1\varepsilon} - u'_1\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))} = 0, \\ \text{(ii)} \lim_{\varepsilon \rightarrow 0} \|u'_{2\varepsilon} - \theta_2^{-1} u'_2\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))} = 0, \\ \text{(iii)} \lim_{\varepsilon \rightarrow 0} \|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0(0, T; [L^1(\Omega_{1\varepsilon})]^n)} = 0, \\ \text{(iv)} \lim_{\varepsilon \rightarrow 0} \|\nabla u_{2\varepsilon}\|_{C^0(0, T; [L^2(\Omega_{2\varepsilon})]^n)} = 0. \end{cases}$$

Remark 2.9. Observe that in both cases, due to convergences (2.19)(iii) and (2.24)(iv), respectively, the additional assumptions (2.9) and (2.10) imply that convergences (2.16)(ii) and (2.22)(ii) are actually strong.

Theorem 2.4 will be proved in section 6. Its proof is more delicate than that of the case $-1 < \gamma < 1$, due to the particular form of the limit problem, and requires specific arguments (see Remark 6.3 for details).

3. Some preliminary results. In this section we will recall some technical results needed in the proofs of Theorems 2.4 and 2.8 as well as the homogenization result proved in [8] and [10] for the elliptic problem.

LEMMA 3.1 (see [7]). (i) *There exists a linear continuous extension operator Q_1 belonging to $\mathcal{L}(H^1(Y_1); H^1(Y)) \cap \mathcal{L}(L^2(Y_1); L^2(Y))$ such that, for some positive constant C ,*

$$\|\nabla Q_1 v_1\|_{L^2(Y)} \leq C \|\nabla v_1\|_{L^2(Y_1)}$$

for every $v_1 \in H^1(Y_1)$.

(ii) *There exists an extension operator Q_1^ε belonging to $\mathcal{L}(L^2(\Omega_{1\varepsilon}); L^2(\Omega)) \cap \mathcal{L}(V^\varepsilon; H_0^1(\Omega))$ such that, for some positive constant C (independent of ε),*

$$\|Q_1^\varepsilon v_1\|_{L^2(\Omega)} \leq C \|v_1\|_{L^2(\Omega_{1\varepsilon})}$$

and

$$\|\nabla Q_1^\varepsilon v_1\|_{L^2(\Omega)} \leq C \|\nabla v_1\|_{L^2(\Omega_{1\varepsilon})}$$

for every $v_1 \in V^\varepsilon$.

LEMMA 3.2 (see [5]). *There exists a linear continuous extension operator P_1^ε belonging to $\mathcal{L}(L^\infty(0, T; H^k(\Omega_{1\varepsilon})); L^\infty(0, T; H^k(\Omega)))$, $k = 1, 2$, such that, for some positive constant C (independent of ε),*

$$\begin{cases} \text{(i)} P_1^\varepsilon \varphi = \varphi \text{ in } \Omega_{1\varepsilon} \times]0, T[, \\ \text{(ii)} P_1^\varepsilon \varphi' = (P^\varepsilon \varphi)' \text{ in } \Omega \times]0, T[, \\ \text{(iii)} \|P_1^\varepsilon \varphi\|_{L^\infty(0, T; L^2(\Omega))} \leq C \|\varphi\|_{L^\infty(0, T; L^2(\Omega_{1\varepsilon}))}, \\ \text{(iv)} \|P_1^\varepsilon \varphi'\|_{L^\infty(0, T; L^2(\Omega))} \leq C \|\varphi'\|_{L^\infty(0, T; L^2(\Omega_{1\varepsilon}))}, \\ \text{(v)} \|\nabla(P_1^\varepsilon \varphi)\|_{L^\infty(0, T; [L^2(\Omega)]^n)} \leq C \|\nabla \varphi\|_{L^\infty(0, T; [L^2(\Omega_{1\varepsilon})]^n)} \end{cases}$$

for any $\varphi \in L^\infty(0, T; H^k(\Omega_{1\varepsilon}))$.

The following homogenization result concerns the elliptic problem associated to problem (2.6).

THEOREM 3.3 (see [6] and [8]). *Let $-1 < \gamma \leq 1$, A^ε and h^ε be defined by (2.4) and (2.5), respectively, and u_ε be the solution of problem*

$$\begin{cases} -\operatorname{div}(A^\varepsilon \nabla u_{1\varepsilon}) = b_1^\varepsilon + Q_1^{\varepsilon*}(g) & \text{in } \Omega_{1\varepsilon}, \\ -\operatorname{div}(A^\varepsilon \nabla u_{2\varepsilon}) = b_2^\varepsilon & \text{in } \Omega_{2\varepsilon}, \\ A^\varepsilon \nabla u_{1\varepsilon} \cdot n_{1\varepsilon} = -A^\varepsilon \nabla u_{2\varepsilon} \cdot n_{2\varepsilon} & \text{on } \Gamma^\varepsilon, \\ A^\varepsilon \nabla u_{1\varepsilon} \cdot n_{1\varepsilon} = -\varepsilon^\gamma h^\varepsilon(u_{1\varepsilon} - u_{2\varepsilon}) & \text{on } \Gamma^\varepsilon, \\ u_{1\varepsilon} = 0 & \text{on } \partial\Omega, \end{cases}$$

where $Q_1^\varepsilon \in \mathcal{L}(H^{-1}(\Omega); V^\varepsilon)$ is the adjoint of the extension operator given in Lemma 3.1(ii). Suppose that $b_1^\varepsilon \in L^2(\Omega_{1\varepsilon})$ and $b_2^\varepsilon \in L^2(\Omega_{2\varepsilon})$ satisfy

$$\begin{cases} \widetilde{b}_1^\varepsilon \rightharpoonup \theta_1 b_1 & \text{weakly in } L^2(\Omega), \\ \widetilde{b}_2^\varepsilon \rightharpoonup \theta_2 b_2 & \text{weakly in } L^2(\Omega), \end{cases}$$

and let g be given in $L^2(\Omega)$. Then, there exists a positive constant C (independent of ε) such that

$$\begin{cases} Q_1^\varepsilon u_{1\varepsilon} \rightharpoonup u_1 & \text{weakly in } H_0^1(\Omega), \\ A^\varepsilon \nabla u_{1\varepsilon} \rightharpoonup A^0 \nabla u_1 & \text{weakly in } [L^2(\Omega)]^n, \\ \|u_{1\varepsilon} - u_{2\varepsilon}\|_{L^2(\Gamma^\varepsilon)} < C\varepsilon^{-\gamma/2}, \end{cases}$$

and the following convergences hold:

$$\begin{cases} \widetilde{u_{2\varepsilon}} \rightharpoonup u_2 & \text{weakly in } L^2(\Omega), \\ A^\varepsilon \nabla \widetilde{u_{2\varepsilon}} \rightharpoonup 0 & \text{weakly in } [L^2(\Omega)]^n. \end{cases}$$

The function u_1 is the unique solution in $H_0^1(\Omega)$ of the problem

$$\begin{cases} -\operatorname{div}(A^0 \nabla u_1) = \theta_1 b_1 + \theta_2 b_2 + g & \text{in } \Omega, \\ u_1 = 0 & \text{on } \partial\Omega, \end{cases}$$

with θ_i , $i = 1, 2$, given by (2.2) and A^0 by (2.12).

Moreover, for $-1 < \gamma < 1$, one has

$$\begin{cases} u_2 = \theta_2 u_1, \\ \|Q_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon}\|_{L^2(\Omega_{2\varepsilon})}^2 \rightarrow 0, \end{cases}$$

while, for $\gamma = 1$,

$$u_2 = \theta_2 (u_1 + c_h^{-1} b_2),$$

where $c_h = \frac{1}{|Y_2|} \int_\Gamma h(y) d\sigma_y$.

The following lemma is a straightforward adaptation of Lemma 3.3 of [8] to the time-dependent case.

LEMMA 3.4. *Suppose that Γ is of class C^2 . Let $g \in L^\infty(\Gamma)$, and set $c_g := \frac{1}{|Y_2|} \int_\Gamma g(y) d\sigma_y$.*

(i) *There exists $\psi_g \in W^{1,\infty}(Y_2)$ such that for every $v_\varepsilon \in L^2(0, T; W^{1,1}(\Omega_{2\varepsilon}))$ one has*

$$\varepsilon \int_{\Gamma^\varepsilon} g(x/\varepsilon)v_\varepsilon(x, t) d\sigma_x = c_g \int_{\Omega_{2\varepsilon}} v_\varepsilon(x, t) dx + \varepsilon \int_{\Omega_{2\varepsilon}} \nabla_y \psi_g(x/\varepsilon) \nabla_x v_\varepsilon(x, t) dx$$

for all $t \in [0, T]$.

(ii) *If for some positive constant C (independent of ε) one has*

$$\|v_\varepsilon\|_{L^2(0, T; W^{1,1}(\Omega_{2\varepsilon}))} \leq C,$$

then

$$\liminf_{\varepsilon \rightarrow 0} \varepsilon \int_{\Gamma^\varepsilon} g(x/\varepsilon)v_\varepsilon(x, t) d\sigma_x = \liminf_{\varepsilon \rightarrow 0} c_g \int_{\Omega_{2\varepsilon}} v_\varepsilon(x, t) dx .$$

We will also need the following result.

LEMMA 3.5 (see [11]). *Let \mathcal{O} be an open set of \mathbb{R}^n and $\{\mathcal{O}_\varepsilon\}_\varepsilon \subset \mathcal{O}$ a sequence of open subsets of \mathcal{O} . Suppose that $\{v_\varepsilon\}_\varepsilon \subset L^p(0, T; L^p(\mathcal{O}_\varepsilon))$, $p > 1$, is such that*

$$\begin{cases} \chi_{\mathcal{O}_\varepsilon} \rightharpoonup \chi_0 & \text{in } L^\infty(\mathcal{O}) \text{ weakly } *, \\ \tilde{v}_\varepsilon \rightharpoonup \chi_0 v & \text{weakly in } L^p(0, T; L^p(\mathcal{O})) . \end{cases}$$

Then,

$$\liminf_{\varepsilon \rightarrow 0} \int_{\mathcal{O}_\varepsilon} |v_\varepsilon(t)|^p dx \geq \int_{\mathcal{O}} \chi_0 |v(t)|^p dx .$$

4. Asymptotic behavior of the energy. In this section we study the convergence of the energy associated to problem (2.6) which, for every ε , is defined by (4.1)

$$\begin{aligned} E^\varepsilon(t) := & \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon}(t)|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon}(t)|^2 dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \nabla u_{1\varepsilon}(t) dx \right. \\ & \left. + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) dx + \varepsilon^\gamma \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right] . \end{aligned}$$

LEMMA 4.1. *For $-1 < \gamma \leq 1$, let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7) holds. If u_ε is the solution of problem (2.6), then*

$$\begin{aligned} E^\varepsilon(t) = & \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |U_\varepsilon^1|^2 dx + \int_{\Omega_{2\varepsilon}} |U_\varepsilon^1|^2 dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla U_{1\varepsilon}^0 \nabla U_{1\varepsilon}^0 dx \right. \\ (4.2) \quad & \left. + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla U_{2\varepsilon}^0 \nabla U_{2\varepsilon}^0 dx + \varepsilon^\gamma \int_{\Gamma^\varepsilon} h^\varepsilon |U_{1\varepsilon}^0 - U_{2\varepsilon}^0|^2 d\sigma_x \right] \\ & + \int_0^t \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} u'_{1\varepsilon} dx d\tau + \int_0^t \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} u'_{2\varepsilon} dx d\tau \quad \forall t \in [0, T] . \end{aligned}$$

Proof. If we take $(u'_{1\varepsilon}, u'_{2\varepsilon})$ as test functions in the variational formulation of problem (2.6) (actually a standard density argument has to be used, see, for instance,

[4] and [20]) we get

$$\begin{aligned}
 & \int_0^t \langle u''_{1\varepsilon}, u'_{1\varepsilon} \rangle_{(V^\varepsilon)', V^\varepsilon} d\tau + \int_0^t \langle u''_{2\varepsilon}, u'_{2\varepsilon} \rangle_{(H^1(\Omega_{2\varepsilon}))', H^1(\Omega_{2\varepsilon})} d\tau \\
 (4.3) \quad & + \int_0^t \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon} \nabla u'_{1\varepsilon} dx d\tau + \int_0^t \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon} \nabla u'_{2\varepsilon} dx d\tau \\
 & + \varepsilon^\gamma \int_0^t \int_{\Gamma^\varepsilon} h^\varepsilon(u_{1\varepsilon} - u_{2\varepsilon})(u'_{1\varepsilon} - u'_{2\varepsilon}) d\sigma_x d\tau \\
 & = \int_0^t \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} u'_{1\varepsilon} dx d\tau + \int_0^t \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} u'_{2\varepsilon} dx d\tau.
 \end{aligned}$$

Moreover,

$$\int_0^t \langle u''_{i\varepsilon}, u'_{i\varepsilon} \rangle_{(V^\varepsilon)', V^\varepsilon} d\tau = \frac{1}{2} \int_0^t \frac{d}{d\tau} \int_{\Omega_{i\varepsilon}} |u'_{i\varepsilon}|^2 dx d\tau, \quad i = 1, 2.$$

Similarly, since the matrix A is symmetric we get

$$\int_0^t \int_{\Omega_{i\varepsilon}} A^\varepsilon \nabla u_{i\varepsilon} \cdot \nabla u'_{i\varepsilon} dx d\tau = \frac{1}{2} \int_0^t \frac{d}{d\tau} \int_{\Omega_{i\varepsilon}} A^\varepsilon \nabla u_{i\varepsilon} \cdot \nabla u_{i\varepsilon} dx d\tau, \quad i = 1, 2.$$

On the other hand,

$$\varepsilon^\gamma \int_0^t \int_{\Gamma^\varepsilon} h^\varepsilon(u_{1\varepsilon} - u_{2\varepsilon})(u'_{1\varepsilon} - u'_{2\varepsilon}) d\sigma_x d\tau = \frac{\varepsilon^\gamma}{2} \int_0^t \frac{d}{d\tau} \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon} - u_{2\varepsilon}|^2 d\sigma_x d\tau.$$

Therefore, identity (4.3) can be rewritten, using (4.1), as

$$E^\varepsilon(t) - E^\varepsilon(0) = \int_0^t \frac{dE^\varepsilon}{d\tau} d\tau = \int_0^t \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} u'_{1\varepsilon} dx d\tau + \int_0^t \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} u'_{2\varepsilon} dx d\tau.$$

Hence, again by (4.1) and (2.6), we have (4.2). \square

LEMMA 4.2. For $-1 < \gamma \leq 1$, let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7), (2.9), and (2.10) hold. If u_ε is the solution of problem (2.6), then

$$\begin{aligned}
 \lim_{\varepsilon \rightarrow 0} E^\varepsilon(t) &= \frac{1}{2} \left[\int_\Omega |U^1|^2 dx + \int_\Omega A^0 \nabla U^0 \nabla U^0 dx \right] \\
 &+ \int_0^t \int_\Omega f_1 u'_1 dx d\tau + \theta_2^{-1} \int_0^t \int_\Omega f_2 u'_2 dx d\tau \quad \forall t \in [0, T],
 \end{aligned}$$

where we have $u'_2 = \theta_2 u'_1$, if $-1 < \gamma < 1$.

Proof. We make use of Lemma 4.1, passing to the limit in (4.2) for any fixed $t \in [0, T]$.

For the first two integrals, using (2.9), we have

$$\lim_{\varepsilon \rightarrow 0} \left(\int_{\Omega_{1\varepsilon}} |U^\varepsilon_1|^2 dx + \int_{\Omega_{2\varepsilon}} |U^\varepsilon_2|^2 dx \right) = \int_\Omega |U^1|^2 dx.$$

Observe now that (2.11) holds in both cases $-1 < \gamma < 1$ and $\gamma = 1$. Indeed, Theorem 3.3 written for $u_{1\varepsilon} = U^0_{1\varepsilon}$, $u_{2\varepsilon} = U^0_{2\varepsilon}$, $b_{1\varepsilon} = b_{2\varepsilon} = 0$, and $g = -\operatorname{div} A^0 \nabla U^0$

gives, by the uniqueness of the solution of the limit problem, $u_1 = U^0$. Hence, if we take $(U_{1\varepsilon}^0, U_{2\varepsilon}^0)$ as test functions in (2.10), in view of (2.11)(i) we have

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \left[\int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla U_{1\varepsilon}^0 \nabla U_{1\varepsilon}^0 \, dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla U_{2\varepsilon}^0 \nabla U_{2\varepsilon}^0 \, dx + \varepsilon^\gamma \int_{\Gamma^\varepsilon} h^\varepsilon (U_{1\varepsilon}^0 - U_{2\varepsilon}^0)^2 \, d\sigma_x \right] \\ = \lim_{\varepsilon \rightarrow 0} \langle -\operatorname{div} (A^0 \nabla U^0), Q_1^\varepsilon U_{1\varepsilon}^0 \rangle = \int_{\Omega} A^0 \nabla U^0 \nabla U^0 \, dx. \end{aligned}$$

On the other hand, (2.9) and Remark 2.7(ii) imply that

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \int_0^t \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} u'_{1\varepsilon} \, dx \, d\tau + \int_0^t \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} u'_{2\varepsilon} \, dx \, d\tau \\ = \int_0^t \int_{\Omega} f_1 u'_1 \, dx \, d\tau + \theta_2^{-1} \int_0^t \int_{\Omega} f_2 u'_2 \, dx \, d\tau, \end{aligned}$$

where for $-1 < \gamma < 1$ and $\gamma = 1$ we used convergences (2.15) and (2.21), respectively. \square

From now on, the convergence of the energy in the two cases $-1 < \gamma < 1$ and $\gamma = 1$ need to be studied separately.

THEOREM 4.3 (convergence of energy for $-1 < \gamma < 1$). *Let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7), (2.9), and (2.10) hold. If u_ε is the solution of problem (2.6) with $-1 < \gamma < 1$, then*

$$(4.4) \quad E^\varepsilon \rightarrow E \text{ in } C^0([0, T]),$$

where

$$(4.5) \quad E(t) := \frac{1}{2} \left[\int_{\Omega} |u'_1(t)|^2 \, dx + \int_{\Omega} A^0 \nabla u_1(t) \nabla u_1(t) \, dx \right]$$

is the energy associated to problem (2.17).

Proof. Taking u'_1 as test function in problem (2.17) and arguing as in the proof of Lemma 4.1, we obtain

$$E(t) = \frac{1}{2} \left[\int_{\Omega} |U^1|^2 \, dx + \int_{\Omega} A^0 \nabla U^0 \nabla U^0 \, dx \right] + \int_0^t \int_{\Omega} (f_1 + f_2) u'_1 \, dx \, d\tau.$$

Hence, by Lemma 4.2 we obtain

$$(4.6) \quad \lim_{\varepsilon \rightarrow 0} E^\varepsilon(t) = E(t) \quad \forall t \in [0, T],$$

since in this case $u'_2 = \theta_2 u'_1$.

Using Lemma 4.1, Theorem 2.3, (2.7), (2.9), and (2.10) we deduce that there exists a constant C (independent of ε) such that

$$|E^\varepsilon(t)| \leq C \quad \forall t \in [0, T]$$

and, from the Hölder inequality,

$$\begin{aligned} |E^\varepsilon(t + \sigma) - E^\varepsilon(t)| \\ = \left| \int_t^{t+\sigma} \int_{\Omega_{1\varepsilon}} f_{1\varepsilon} u'_{1\varepsilon} \, dx \, d\tau + \int_t^{t+\sigma} \int_{\Omega_{2\varepsilon}} f_{2\varepsilon} u'_{2\varepsilon} \, dx \, d\tau \right| \leq \sigma^{1/2} C. \end{aligned}$$

Consequently, by the Ascoli–Arzelà theorem, there exist a subsequence, still denoted by E^ε , and a function ζ in $C^0([0, T])$ such that

$$(4.7) \quad E^\varepsilon \rightarrow \zeta \text{ in } C^0([0, T]).$$

Hence, from (4.6) and (4.7), by uniqueness, one has $\zeta = E$, which concludes the proof. \square

THEOREM 4.4 (convergence of energy for $\gamma = 1$). *Let A^ε and h^ε be defined by (2.4) and (2.5), respectively, and suppose that (2.7), (2.9), and (2.10) hold. If u_ε is the solution of problem (2.6) with $\gamma = 1$, then*

$$E^\varepsilon \rightarrow \hat{E} \text{ in } C^0([0, T]),$$

where

$$(4.8) \quad \begin{aligned} \hat{E}(t) := & \frac{1}{2} \left[\theta_1 \int_\Omega |u'_1(t)|^2 dx + \theta_2^{-1} \int_\Omega |u'_2(t)|^2 dx \right. \\ & \left. + \int_\Omega A^0 \nabla u_1(t) \nabla u_1(t) dx + c_h \theta_2^{-1} \int_\Omega |\theta_2 u_1(t) - u_2(t)|^2 dx \right] \end{aligned}$$

is the energy associated to problem (2.23).

Proof. Let us prove first that

$$(4.9) \quad \begin{aligned} \hat{E}(t) = & \frac{1}{2} \left[\int_\Omega |U^1|^2 dx + \int_\Omega A^0 \nabla U^0 \nabla U^0 dx \right] \\ & + \int_0^t \int_\Omega f_1 u'_1 dx d\tau + \theta_2^{-1} \int_0^t \int_\Omega f_2 u'_2 dx d\tau \quad \forall t \in [0, T]. \end{aligned}$$

Let us take u'_1 as test function in the first equation of (2.23). Multiplying the second equation by $\theta_2^{-1} u'_2$ and integrating over $\Omega \times [0, t]$ we have by summation

$$(4.10) \quad \begin{aligned} & \int_0^t \langle \theta_1 u''_1, u'_1 \rangle_{L^2(\Omega), L^2(\Omega)} d\tau + \int_0^t \langle u''_2, \theta_2^{-1} u'_2 \rangle_{L^2(\Omega), L^2(\Omega)} d\tau \\ & \quad + \int_0^t \int_\Omega A^0 \nabla u_1 \nabla u'_1 dx d\tau + c_h \int_0^t \int_\Omega (\theta_2 u_1 - u_2) u'_1 dx d\tau \\ & \quad - c_h \int_0^t \int_\Omega (\theta_2 u_1 - u_2) \theta_2^{-1} u'_2 dx d\tau \\ & = \int_0^t \int_\Omega f_1 u'_1 dx d\tau + \theta_2^{-1} \int_0^t \int_\Omega f_2 u'_2 dx d\tau \quad \forall t \in [0, T]. \end{aligned}$$

Now, observe that

$$\begin{aligned} \int_0^t \langle \theta_1 u''_1, u'_1 \rangle_{L^2(\Omega), L^2(\Omega)} d\tau &= \frac{1}{2} \theta_1 \int_0^t \frac{d}{d\tau} \int_\Omega |u'_1|^2 dx d\tau, \\ \int_0^t \langle u''_2, \theta_2^{-1} u'_2 \rangle_{L^2(\Omega), L^2(\Omega)} d\tau &= \frac{1}{2} \theta_2^{-1} \int_0^t \frac{d}{d\tau} \int_\Omega |u'_2|^2 dx d\tau. \end{aligned}$$

Similarly, taking into account that the matrix A^0 is symmetric we get

$$\int_0^t \int_\Omega A^0 \nabla u_1 \nabla u'_1 dx d\tau = \frac{1}{2} \int_0^t \frac{d}{d\tau} \int_\Omega A^0 \nabla u_1 \cdot \nabla u_1 dx d\tau.$$

On the other hand,

$$\begin{aligned} & c_h \int_0^t \int_{\Omega} (\theta_2 u_1 - u_2) u_1' \, dx \, d\tau - c_h \int_0^t \int_{\Omega} (\theta_2 u_1 - u_2) \theta_2^{-1} u_2' \, dx \, d\tau \\ &= c_h \theta_2^{-1} \int_0^t \int_{\Omega} (\theta_2 u_1 - u_2) (\theta_2 u_1' - u_2') \, dx \, d\tau = \frac{c_h \theta_2^{-1}}{2} \int_0^t \frac{d}{d\tau} \int_{\Omega} (\theta_2 u_1 - u_2)^2 \, dx \, d\tau. \end{aligned}$$

Therefore, taking into account (4.8) and identity (4.10) we obtain

$$\hat{E}(t) - \hat{E}(0) = \int_0^t \frac{d\hat{E}}{d\tau} \, d\tau = \int_0^t \int_{\Omega} f_1 u_1' \, dx \, d\tau + \theta_2^{-1} \int_0^t \int_{\Omega} f_2 u_2' \, dx \, d\tau \quad \forall t \in [0, T].$$

Using (2.23) and again (2.8), we deduce that

$$\begin{aligned} \hat{E}(0) &= \frac{1}{2} \left[\theta_1 \int_{\Omega} |U^1|^2 \, dx + \theta_2^{-1} \int_{\Omega} |\theta_2 U^1|^2 \, dx + \int_{\Omega} A^0 \nabla U^0 \nabla U^0 \, dx \right] \\ &\quad + c_h \theta_2^{-1} \int_{\Omega} |\theta_2 U^0 - \theta_2 U^1|^2 \, dx \\ &= \frac{1}{2} \left[\int_{\Omega} |U^1|^2 \, dx + \int_{\Omega} A^0 \nabla U^0 \nabla U^0 \, dx \right]. \end{aligned}$$

This gives (4.9) which, together with Lemma 4.2, implies that

$$\lim_{\varepsilon \rightarrow 0} E^\varepsilon(t) \rightarrow \hat{E}(t) \quad \forall t \in [0, T].$$

Now, arguing as in the proof of Theorem 4.3, but using here Theorem 2.6 instead of Theorem 2.3, we conclude the proof. \square

5. Proof of the corrector result in the case $-1 < \gamma < 1$. Let Q_1 be defined by Lemma 3.1 and introduce the test functions associated to the solution w_λ of problem (2.13) by

$$(5.1) \quad \chi_{1\lambda} = \lambda \cdot y - w_\lambda(y), \quad w_\lambda^\varepsilon(x) := \lambda \cdot x - \varepsilon(Q_1(\chi_{1\lambda})(x/\varepsilon)).$$

A simple change of scale gives that

$$(5.2) \quad \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_\lambda^\varepsilon \cdot \nabla v_1 \, dx = 0$$

for every $v_1 \in V^\varepsilon$ and (see [7]) the following convergences hold true:

$$(5.3) \quad \begin{cases} w_\lambda^\varepsilon \rightharpoonup \lambda \cdot x & \text{weakly in } H^1(\Omega), \\ w_\lambda^\varepsilon \rightarrow \lambda \cdot x & \text{strongly in } L^2(\Omega), \\ \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_\lambda^\varepsilon \rightharpoonup A^0 \lambda & \text{weakly in } [L^2(\Omega)]^n, \end{cases}$$

where A^0 is given by (2.12).

From now on, we will adopt the Einstein summation convention. The proof of the corrector result is based on the following proposition.

PROPOSITION 5.1. *Suppose that the assumptions of Theorem 2.4 are fulfilled. Set for any $\Phi \in C^\infty([0, T], \mathcal{D}(\Omega))$*

$$\begin{aligned} X_\varepsilon(t) &= \frac{1}{2} \left[\int_{\Omega} \left| \widetilde{u_{1\varepsilon}'}(t) + \widetilde{u_{2\varepsilon}'}(t) - \Phi'(t) \right|^2 \, dx \right. \\ &\quad \left. + \int_{\Omega_{1\varepsilon}} A^\varepsilon (\nabla u_{1\varepsilon}(t) - C^\varepsilon \nabla \Phi(t)) (\nabla u_{1\varepsilon}(t) - C^\varepsilon \nabla \Phi(t)) \, dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) \, dx \right], \end{aligned}$$

where C^ε is given in (2.18). Then

$$(5.4) \quad \limsup_{\varepsilon \rightarrow 0} \|X_\varepsilon\|_{C^0([0,T])} \leq \|X\|_{C^0([0,T])},$$

where

$$X(t) = \frac{1}{2} \left[\|u'_1(t) - \Phi'(t)\|_{L^2(\Omega)}^2 + \int_{\Omega} A^0(\nabla u_1(t) - \nabla \Phi(t))(\nabla u_1(t) - \nabla \Phi(t)) \, dx \right].$$

Remark 5.2. Observe that in general in literature (see, for instance, [5], [17], [18]) one can introduce suitable quantities X_ε and X such that

$$X_\varepsilon \rightarrow X \quad \text{in } C^0([0, T]).$$

Here, we can prove only (5.4), which is just an inequality on the upper limit of $\|X_\varepsilon\|_{C^0([0,T])}$. Nevertheless, this will be sufficient to prove the corrector result.

Proof of Proposition 5.1. Using the symmetry of A^ε and C^ε , one has

$$\begin{aligned} X_\varepsilon(t) &= \frac{1}{2} \left[\int_{\Omega} \left| \widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) - \Phi'(t) \right|^2 \right. \\ &\quad + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \nabla u_{1\varepsilon}(t) \, dx - 2 \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) \nabla u_{1\varepsilon}(t) \, dx \\ &\quad \left. + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) C^\varepsilon \nabla \Phi(t) \, dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) \, dx \right]. \end{aligned}$$

Set, for $t \in [0, T]$,

$$\begin{aligned} X_{1\varepsilon}(t) &= \frac{1}{2} \left[\int_{\Omega} \left| \widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right|^2 + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \nabla u_{1\varepsilon}(t) \, dx \right. \\ &\quad \left. + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) \, dx \right], \\ X_{2\varepsilon}(t) &= \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi'(t) \, dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) \nabla u_{1\varepsilon}(t) \, dx, \end{aligned}$$

and

$$X_{3\varepsilon}(t) = \frac{1}{2} \left[\int_{\Omega} |\Phi'(t)|^2 + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) C^\varepsilon \nabla \Phi(t) \, dx \right]$$

so that

$$X_\varepsilon = X_{1\varepsilon} - X_{2\varepsilon} + X_{3\varepsilon}.$$

For $X_{1\varepsilon}$, recalling (4.1), we have

$$(5.5) \quad X_{1\varepsilon}(t) \leq E^\varepsilon(t) \quad \forall t \in [0, T].$$

For $X_{2\varepsilon}$, taking into account (2.18), (5.1) written for $\lambda = e_i$, and (5.2) we obtain

$$\begin{aligned}
 X_{2\varepsilon}(t) &= \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi'(t) \, dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \frac{\partial \Phi}{\partial x_i}(t) \nabla u_{1\varepsilon}(t) \, dx \\
 &= \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi'(t) \, dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) u_{1\varepsilon}(t) \right] \, dx \\
 (5.6) \quad &\quad - \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] u_{1\varepsilon}(t) \, dx \\
 &= \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi'(t) \, dx \\
 &\quad - \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) \, dx \quad \forall t \in [0, T].
 \end{aligned}$$

Hence, by (2.14), (2.15), and (5.3) we conclude that

$$\begin{aligned}
 X_{2\varepsilon} &\rightarrow \int_{\Omega} (\theta_1 u'_1 + \theta_2 u'_1) \Phi' \, dx - \int_{\Omega} A^0 e_i \nabla \left[\frac{\partial \Phi}{\partial x_i} \right] u_1 \, dx \\
 (5.7) \quad &= \int_{\Omega} u'_1 \Phi' \, dx + \int_{\Omega} A^0 \nabla \Phi \nabla u_1 \, dx \text{ in } \mathcal{D}'(0, T).
 \end{aligned}$$

Moreover, taking into account (5.6), by Theorem 2.3, (5.2), and (5.3) we obtain that there exists a constant C (independent of ε) such that

$$(5.8) \quad \|X_{2\varepsilon}\|_{L^\infty(0,T)} \leq C.$$

Consequently, to prove that convergence (5.7) takes place in $C^0([0, T])$ it is enough to show that

$$(5.9) \quad \left\| \frac{\partial X_{2\varepsilon}}{\partial t} \right\|_{L^2(0,T)} \leq C.$$

Using the variational formulation (2.8), we have

$$\begin{aligned}
 \frac{\partial}{\partial t} \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi'(t) \, dx &= \int_{\Omega} \left(\widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) \right) \Phi''(t) \\
 &\quad - \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \cdot \nabla \Phi'(t) \, dx - \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \cdot \nabla \Phi'(t) \, dx \\
 &\quad + \int_{\Omega_{1\varepsilon}} f_{1\varepsilon}(t) \Phi'(t) \, dx + \int_{\Omega_{2\varepsilon}} f_{2\varepsilon}(t) \Phi'(t) \, dx,
 \end{aligned}$$

which is bounded in $L^2(0, T)$ due to (2.15), (2.16), and (2.7).

Moreover, we have

$$\begin{aligned}
 \frac{\partial}{\partial t} \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) \, dx &= \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi'}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) \, dx \\
 &\quad + \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u'_{1\varepsilon}(t) \, dx,
 \end{aligned}$$

which is bounded in $L^2(0, T)$ as a consequence of (5.3) and Theorem 2.3.

Thus, from (5.6) we obtain (5.9). Due to classical compactness results, from this inequality and (5.8) it turns out that $X_{2\varepsilon}$ is relatively compact in $C^0([0, T])$. This, together with (5.7), gives

$$(5.10) \quad X_{2\varepsilon} \rightarrow X_2 := \int_{\Omega} u'_1 \Phi' \, dx + \int_{\Omega} A^0 \nabla \Phi \nabla u_1 \, dx \text{ strongly in } C^0([0, T]).$$

For $X_{3\varepsilon}$, by the smoothness of Φ we have that it is bounded in $L^\infty(0, T)$ as well as its time derivative. Moreover, by (5.2) we deduce that

$$\begin{aligned} \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) C^\varepsilon \nabla \Phi(t) \, dx &= \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \frac{\partial \Phi}{\partial x_i}(t) \nabla w_j^\varepsilon \frac{\partial \Phi}{\partial x_j}(t) \, dx \\ &= - \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon w_j^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \frac{\partial \Phi}{\partial x_j}(t) \right] \, dx. \end{aligned}$$

Thus, from (5.3)

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} X_{3\varepsilon}(t) &= \lim_{\varepsilon \rightarrow 0} \frac{1}{2} \left[\|\Phi'(t)\|_{L^2(\Omega)}^2 + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) C^\varepsilon \nabla \Phi(t) \, dx \right] \\ &= \frac{1}{2} \left[\|\Phi'(t)\|_{L^2(\Omega)}^2 - \int_{\Omega} A^0 e_{ij} x_j \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \frac{\partial \Phi}{\partial x_j}(t) \right] \, dx \right] \\ &= \frac{1}{2} \left[\|\Phi'(t)\|_{L^2(\Omega)}^2 + \int_{\Omega} A^0 \nabla \Phi(t) \nabla \Phi(t) \, dx \right] \quad \forall t \in [0, T], \end{aligned}$$

which gives

$$(5.11) \quad X_{3\varepsilon} \rightarrow X_3 := \frac{1}{2} \left[\|\Phi'\|_{L^2(\Omega)}^2 + \int_{\Omega} A^0 \nabla \Phi \nabla \Phi \, dx \right] \text{ strongly in } C^0([0, T]).$$

Convergence (5.11), together with (5.10), (4.4), and (4.5), gives

$$(5.12) \quad E^\varepsilon - X_{2\varepsilon} + X_{3\varepsilon} \rightarrow E - X_2 + X_3 = X \quad \text{strongly in } C^0([0, T]).$$

Since from (5.5) one has

$$0 \leq X_\varepsilon(t) \leq E^\varepsilon(t) - X_{2\varepsilon}(t) + X_{3\varepsilon}(t) \quad \forall t \in [0, T],$$

by (5.12) we conclude that

$$\limsup_{\varepsilon \rightarrow 0} \|X_\varepsilon\|_{C^0([0, T])} \leq \lim_{\varepsilon \rightarrow 0} \|E^\varepsilon - X_{2\varepsilon} + X_{3\varepsilon}\|_{C^0([0, T])} = \|X\|_{C^0([0, T])}. \quad \square$$

Remark 5.3. Let us point out that the main difficulty when proving Proposition 5.1 is due to the fact that, as a consequence of the presence of the boundary term in the energy, we have only $X_{1\varepsilon} \leq E^\varepsilon$ and not $X_{1\varepsilon} = E^\varepsilon$ as in the classical cases.

Nevertheless, the fact that $E^\varepsilon - X_{2\varepsilon} + X_{3\varepsilon} \rightarrow X$ in $C^0([0, T])$ allows us to prove the result.

Proof of Theorem 2.4. The function $u_1 \in L^2(0, T; H_0^1(\Omega)) \cap C^0([0, T]; L^2(\Omega))$ and its derivative $u'_1 \in C^0([0, T]; L^2(\Omega))$; thus, from classical density results, for any $\delta > 0$ there exists $\Phi \in C^\infty([0, T]; \mathcal{D}(\Omega))$ such that

$$(5.13) \quad \begin{cases} \|u'_1 - \Phi'\|_{C^0([0, T]; L^2(\Omega))} \leq \delta, \\ \|\nabla u_1 - \nabla \Phi\|_{L^2(0, T; L^2(\Omega))} \leq \delta. \end{cases}$$

Therefore,

$$\begin{aligned}
 & \left\| \widetilde{u}'_{1\varepsilon} + \widetilde{u}'_{2\varepsilon} - u'_1 \right\|_{C^0(0,T;L^2(\Omega))}^2 \\
 (5.14) \quad & \leq 2 \left(\left\| \widetilde{u}'_{1\varepsilon} + \widetilde{u}'_{2\varepsilon} - \Phi' \right\|_{C^0(0,T;L^2(\Omega))}^2 + \|\Phi' - u'_1\|_{C^0(0,T;L^2(\Omega))}^2 \right) \\
 & \leq 2 \left\| \widetilde{u}'_{1\varepsilon} + \widetilde{u}'_{2\varepsilon} - \Phi' \right\|_{C^0(0,T;L^2(\Omega))}^2 + 2\delta^2.
 \end{aligned}$$

On the other hand, due to the Cauchy–Schwarz inequality, definition (2.18) of C^ε , and (5.13) there exist two constants C_1, C_2 such that

$$\begin{aligned}
 (5.15) \quad & \|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0([0,T];[L^1(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})]^n)}^2 \\
 & \leq 2\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{C^0([0,T];[L^1(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})]^n)}^2 \\
 & \quad + 2\|C^\varepsilon(\nabla \Phi - \nabla u_1)\|_{C^0([0,T];[L^1(\Omega_{1\varepsilon})]^n)}^2 \\
 & \leq 2C_1\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{C^0([0,T];[L^2(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})]^n)}^2 + 2C_2\delta^2.
 \end{aligned}$$

From properties (2.3) and definition (2.4) of A^ε , we have

$$\begin{aligned}
 & \frac{1}{2} \left\| \widetilde{u}'_{1\varepsilon}(t) + \widetilde{u}'_{2\varepsilon}(t) - \Phi'(t) \right\|_{L^2(\Omega)}^2 \\
 & \quad + \frac{\alpha}{2} \left(\|\nabla u_{1\varepsilon}(t) - C^\varepsilon \nabla \Phi(t)\|_{[L^2(\Omega_{1\varepsilon})]^n}^2 + \|\nabla u_{2\varepsilon}(t)\|_{[L^2(\Omega_{2\varepsilon})]^n}^2 \right) \leq X_\varepsilon(t) \quad \forall t \in [0, T];
 \end{aligned}$$

moreover, there exists a constant C_3

$$\|X\|_{C^0([0,T])} \leq C_3\delta^2,$$

where X_ε and X are defined in Proposition 5.1 and are applied to the function Φ , given by (5.13).

From this last inequality and by Proposition 5.1, we conclude that

$$\begin{aligned}
 & \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{1}{2} \left\| \widetilde{u}'_{1\varepsilon} + \widetilde{u}'_{2\varepsilon} - \Phi' \right\|_{C^0([0,T];L^2(\Omega))}^2 \right. \\
 & \quad \left. + \frac{\alpha}{2} \left(\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{C^0([0,T];[L^2(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})]^n)}^2 \right) \right\} \\
 & \leq \limsup_{\varepsilon \rightarrow 0} \|X_\varepsilon\|_{C^0([0,T])} \leq \|X\|_{C^0([0,T])} \leq C_3\delta^2.
 \end{aligned}$$

This, together with (5.14) and (5.15), gives convergences (2.19), δ being arbitrary. \square

6. Proof of the corrector result in the case $\gamma = 1$. The proof is based on the result below, which is analogous to that shown for $-1 < \gamma < 1$ in Proposition 5.1. Let us emphasize that here the proof is more delicate and contains some technically specific parts.

PROPOSITION 6.1. *Suppose that the assumptions of Theorem 2.8 are fulfilled. Set for any $\Phi, \Psi \in C^\infty([0, T], \mathcal{D}(\Omega))$*

$$\begin{aligned} \hat{X}_\varepsilon(t) &= \frac{1}{2} \int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon}(t) - \Phi'(t)|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon}(t) - \Psi'(t)|^2 dx \\ &\quad + \int_{\Omega_{1\varepsilon}} A^\varepsilon (\nabla u_{1\varepsilon}(t) - C^\varepsilon \nabla \Phi(t)) (\nabla u_{1\varepsilon}(t) - C^\varepsilon \nabla \Phi(t)) dx \\ &\quad + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) dx \Big]. \end{aligned}$$

Then,

$$\limsup_{\varepsilon \rightarrow 0} \|\hat{X}_\varepsilon\|_{C^0([0, T])} \leq \|\hat{X}\|_{C^0([0, T])},$$

where

$$(6.1) \quad \begin{aligned} \hat{X}(t) &= \frac{1}{2} \left[\theta_1 \|u'_1(t) - \Phi'(t)\|_{L^2(\Omega)}^2 + \theta_2^{-1} \|u'_2(t) - \theta_2 \Psi'(t)\|_{L^2(\Omega)}^2 \right. \\ &\quad \left. + \int_{\Omega} A^0 (\nabla u_1(t) - \nabla \Phi(t)) (\nabla u_1(t) - \nabla \Phi(t)) dx \right]. \end{aligned}$$

Remark 6.2. Here also (see Remark 5.2) we have only an inequality on the upper limit of $\|\hat{X}_\varepsilon\|_{C^0([0, T])}$. Due to the special form of the energy, to prove Proposition 6.1 we need to introduce more complex arguments than those used in the proof of Proposition 5.1.

Proof of Proposition 6.1. By the symmetry of A^ε and C^ε , one has

$$\begin{aligned} \hat{X}_\varepsilon(t) &= \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon}(t) - \Phi'(t)|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon}(t) - \Psi'(t)|^2 dx \right. \\ &\quad + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \nabla u_{1\varepsilon}(t) dx - 2 \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) \nabla u_{1\varepsilon}(t) dx \\ &\quad \left. + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) C^\varepsilon \nabla \Phi(t) dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) dx \right]. \end{aligned}$$

Decompose \hat{X}_ε as

$$\hat{X}_\varepsilon = X_{1\varepsilon} - \hat{X}_{2\varepsilon} + \hat{X}_{3\varepsilon},$$

where, as in the previous case,

$$\begin{aligned} X_{1\varepsilon}(t) &= \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |u'_{1\varepsilon}(t)|^2 dx + \int_{\Omega_{2\varepsilon}} |u'_{2\varepsilon}(t)|^2 dx \right. \\ &\quad \left. + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \nabla u_{1\varepsilon}(t) dx + \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla u_{2\varepsilon}(t) dx \right], \end{aligned}$$

while

$$\begin{aligned} \hat{X}_{2\varepsilon}(t) &= \int_{\Omega_{1\varepsilon}} u'_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega_{2\varepsilon}} u'_{2\varepsilon}(t) \Psi'(t) dx \\ &\quad + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi(t) \nabla u_{1\varepsilon}(t) dx \end{aligned}$$

and

$$\hat{X}_{3\varepsilon}(t) = \frac{1}{2} \left[\int_{\Omega_{1\varepsilon}} |\Phi'(t)|^2 dx + \int_{\Omega_{2\varepsilon}} |\Psi'(t)|^2 dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon C^\varepsilon \nabla \Phi C^\varepsilon \nabla \Phi dx \right].$$

We proceed in several steps.

Step 1. Let us prove that

$$(6.2) \quad \limsup_{\varepsilon \rightarrow 0} X_{1\varepsilon}(t) \leq \hat{X}_1 := \frac{1}{2} \left[\theta_1 \int_{\Omega} |u'_1(t)|^2 dx + \theta_2^{-1} \int_{\Omega} |u'_2(t)|^2 dx + \int_{\Omega} A^0 \nabla u_1(t) \nabla u_1(t) dx \right] \quad \forall t \in [0, T].$$

Since from (4.1) we have

$$(6.3) \quad X_{1\varepsilon} = E^\varepsilon - \frac{1}{2} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon} - u_{2\varepsilon}|^2 d\sigma_x \right),$$

following the same ideas of [8] we show first that

$$(6.4) \quad \liminf_{\varepsilon \rightarrow 0} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right) \geq c_h \theta_2^{-1} \int_{\Omega} |\theta_2 u_1(t) - u_2(t)|^2 dx$$

for any $t \in [0, T]$.

To do that, we apply Lemma 3.4(ii) with $g = h$ and $v_\varepsilon = (P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon})^2$ getting

$$(6.5) \quad \begin{aligned} & \liminf_{\varepsilon \rightarrow 0} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right) \\ &= \liminf_{\varepsilon \rightarrow 0} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |P_1^\varepsilon u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right) \\ &= \liminf_{\varepsilon \rightarrow 0} c_h \int_{\Omega_{2\varepsilon}} |P_1^\varepsilon u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 dx \quad \forall t \in [0, T]. \end{aligned}$$

Observe now that, thanks to (2.2), (2.20), and Remark 2.7(ii), we have

$$\widetilde{P_1^\varepsilon u_{1\varepsilon}}_{|\Omega_{2\varepsilon}}(t) - \widetilde{u_{2\varepsilon}}(t) = \chi_{|\Omega_{2\varepsilon}} P_1^\varepsilon u_{1\varepsilon}(t) - \widetilde{u_{2\varepsilon}}(t) \rightharpoonup \theta_2 (u_1(t) - \theta_2^{-1} u_2(t)) \text{ weakly in } L^2(\Omega)$$

for any $t \in [0, T]$.

Hence, we can apply Lemma 3.5 with $p = 2$, $\mathcal{O}_\varepsilon = \Omega_{2\varepsilon}$, $\chi_0 = \theta_2$, $v_\varepsilon = P_1^\varepsilon u_{1\varepsilon}|_{\Omega_{2\varepsilon}} - u_{2\varepsilon}$, and $v = (u_1 - \theta_2^{-1} u_2)$. We have

$$\liminf_{\varepsilon \rightarrow 0} c_h \int_{\Omega_{2\varepsilon}} |P_1^\varepsilon u_{1\varepsilon}|_{\Omega_{2\varepsilon}}(t) - u_{2\varepsilon}(t)|^2 dx \geq c_h \int_{\Omega} \theta_2 |u_1(t) - \theta_2^{-1} u_2(t)|^2 dx \quad \forall t \in [0, T],$$

which, together with (6.5), gives (6.4).

By (6.3), (6.4), and Theorem 4.4 we conclude that

$$\begin{aligned} \limsup_{\varepsilon \rightarrow 0} X_{1\varepsilon}(t) &= \limsup_{\varepsilon \rightarrow 0} \left[E^\varepsilon(t) - \frac{1}{2} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right) \right] \\ &= \lim_{\varepsilon \rightarrow 0} E^\varepsilon(t) - \liminf_{\varepsilon \rightarrow 0} \left[\frac{1}{2} \left(\varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon |u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 d\sigma_x \right) \right] \\ &\leq \frac{1}{2} \left[\theta_1 \int_{\Omega} |u'_1(t)|^2 dx + \theta_2^{-1} \int_{\Omega} |u'_2(t)|^2 dx \right. \\ &\quad \left. + \int_{\Omega} A^0 \nabla u_1(t) \nabla u_1(t) dx + c_h \theta_2^{-1} \int_{\Omega} |\theta_2 u_1(t) - u_2(t)|^2 dx \right. \\ &\quad \left. - c_h \theta_2^{-1} \int_{\Omega} |\theta_2 u_1(t) - u_2(t)|^2 dx \right] \quad \forall t \in [0, T], \end{aligned}$$

which proves (6.2).

Step 2. Observe now that, due to the presence of the boundary term in (6.3), we cannot have a priori estimates on the time derivative of $X_{1\varepsilon}$. Hence, we cannot deduce any compactness of $X_{1\varepsilon}$ in $C^0([0, T])$. To overcome this difficulty, we decompose $X_{1\varepsilon}$ as a sum of a compact part $E^\varepsilon - Y_\varepsilon$ and a rest \hat{Y}_ε , which goes to zero in $C^0([0, T])$ as $\varepsilon \rightarrow 0$. To do that we use Lemma 3.4(i) in (6.3) to obtain

$$X_{1\varepsilon} = E^\varepsilon - Y_\varepsilon - \hat{Y}_\varepsilon,$$

where

$$Y_\varepsilon(t) := c_h \int_{\Omega_{2\varepsilon}} |P_1^\varepsilon u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 dx$$

and

$$(6.6) \quad \hat{Y}_\varepsilon(t) := \varepsilon \int_{\Omega_{2\varepsilon}} \nabla_y \psi_h(x/\varepsilon) \nabla_x |P_1^\varepsilon u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 dx \quad \forall t \in [0, T].$$

Let us show first that

$$(6.7) \quad E^\varepsilon - Y_\varepsilon \text{ is compact in } C^0([0, T]).$$

Indeed, from Theorem 4.4 E^ε converges in $C^0([0, T])$. Moreover, Y_ε is bounded in $L^\infty(0, T)$ and

$$\left\| \frac{\partial Y_\varepsilon}{\partial t} \right\|_{L^\infty(0, T)} \leq C \|P_1^\varepsilon u'_{1\varepsilon} - u'_{2\varepsilon}\|_{L^\infty(0, T; L^2(\Omega_{2\varepsilon}))} \|P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon}\|_{L^\infty(0, T; L^2(\Omega_{2\varepsilon}))} \leq C.$$

This gives (6.7).

To conclude this step, let us prove that

$$(6.8) \quad \hat{Y}_\varepsilon \rightarrow 0 \text{ in } C^0([0, T]).$$

Indeed,

$$\begin{aligned} \varepsilon \left| \int_{\Omega_{2\varepsilon}} \nabla_y \psi_h(x/\varepsilon) \nabla_x |P_1^\varepsilon u_{1\varepsilon}(t) - u_{2\varepsilon}(t)|^2 dx \right| \\ \leq 2\varepsilon \|\nabla_y \psi_h\|_{L^\infty(\mathbb{R}^n)} \|\nabla(P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon})(P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon})\|_{L^\infty(0, T; L^1(\Omega_{2\varepsilon}))} \\ \leq C_3 \varepsilon \|P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon}\|_{C^0(0, T; H^1(\Omega_{2\varepsilon}))} \leq C\varepsilon \quad \forall t \in [0, T] \end{aligned}$$

with C independent of ε and t . Therefore,

$$\varepsilon \left| \int_{\Omega_{2\varepsilon}} \nabla_y \psi_h(x/\varepsilon) \nabla_x (P_1^\varepsilon u_{1\varepsilon} - u_{2\varepsilon})^2 dx \right| \rightarrow 0 \text{ in } C^0([0, T]),$$

which gives (6.8).

Step 3. We prove here the convergences of $\hat{X}_{2\varepsilon}$ and $\hat{X}_{3\varepsilon}$ in $C^0([0, T])$.

Consider first $\hat{X}_{2\varepsilon}$. From (2.18), (5.1) written for $\lambda = e_i$, and (5.2) we have

$$\begin{aligned} (6.9) \quad \hat{X}_{2\varepsilon}(t) &= \int_{\Omega_{1\varepsilon}} u'_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega_{2\varepsilon}} u'_{2\varepsilon}(t) \Psi'(t) dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla u_{1\varepsilon}(t) \frac{\partial \Phi}{\partial x_i}(t) dx \\ &= \int_{\Omega} \widetilde{u}'_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega} \widetilde{u}'_{2\varepsilon}(t) \Psi'(t) dx + \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[u_{1\varepsilon}(t) \frac{\partial \Phi}{\partial x_i}(t) \right] dx \\ &\quad - \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] u_{1\varepsilon}(t) dx = \int_{\Omega} \widetilde{u}'_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega} \widetilde{u}'_{2\varepsilon}(t) \Psi'(t) dx \\ &\quad - \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) dx. \end{aligned}$$

Due to Theorem 2.6, Remark 2.7(ii), and (5.3) one concludes that

$$\begin{aligned} (6.10) \quad \hat{X}_{2\varepsilon} &\rightarrow \int_{\Omega} \theta_1 u'_1 \Phi' dx + \int_{\Omega} u'_2 \Psi' dx - \int_{\Omega} A^0 e_i \nabla \left[\frac{\partial \Phi}{\partial x_i} \right] u_1 dx \\ &= \int_{\Omega} \theta_1 u'_1 \Phi' dx + \int_{\Omega} u'_2 \Psi' dx + \int_{\Omega} A^0 \nabla \Phi \nabla u_1 dx \text{ in } \mathcal{D}'(0, T). \end{aligned}$$

We have to prove that this convergence takes place in $C^0([0, T])$. From (6.9), Theorem 2.6, (5.2), and (5.3) there exists a constant C (independent of ε) such that

$$(6.11) \quad \|\hat{X}_{2\varepsilon}\|_{L^\infty(0, T)} \leq C.$$

Let us show that

$$(6.12) \quad \left\| \frac{\partial \hat{X}_{2\varepsilon}}{\partial t} \right\|_{L^2(0, T)} \leq C.$$

The variational formulation (2.8) written for $v_1 = \Phi$ and $v_2 = \Psi$ reads

$$\begin{aligned} &\frac{\partial}{\partial t} \left(\int_{\Omega} u'_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega} u'_{2\varepsilon}(t) \Psi'(t) dx \right) \\ &= \int_{\Omega} \widetilde{u}'_{1\varepsilon}(t) \Phi''(t) + \int_{\Omega} \widetilde{u}'_{2\varepsilon}(t) \Psi''(t) - \int_{\Omega_{1\varepsilon}} A^\varepsilon \nabla u_{1\varepsilon}(t) \cdot \nabla \Phi'(t) dx \\ &\quad - \int_{\Omega_{2\varepsilon}} A^\varepsilon \nabla u_{2\varepsilon}(t) \nabla \Psi'(t) dx - \varepsilon \int_{\Gamma^\varepsilon} h^\varepsilon(u_{1\varepsilon}(t) - u_{2\varepsilon}(t)) (\Phi'(t) - \Psi'(t)) d\sigma_x \\ &\quad + \int_{\Omega_{1\varepsilon}} f_{1\varepsilon}(t) \Phi'(t) dx + \int_{\Omega_{2\varepsilon}} f_{2\varepsilon}(t) \Psi'(t) dx, \end{aligned}$$

which is bounded in $L^2(0, T)$ due to Theorems 2.6 and (2.9)(i).

Moreover, we have

$$\begin{aligned} \frac{\partial}{\partial t} \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) \, dx &= \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi'}{\partial x_i}(t) \right] P_1^\varepsilon u_{1\varepsilon}(t) \, dx \\ &+ \int_{\Omega} \chi_{\Omega_{1\varepsilon}} A^\varepsilon \nabla w_i^\varepsilon \nabla \left[\frac{\partial \Phi}{\partial x_i}(t) \right] P_1^\varepsilon u'_{1\varepsilon}(t) \, dx, \end{aligned}$$

and the right-hand side of this equality is bounded in $L^2(0, T)$ in view of (5.3) and Theorem 2.6. Hence (6.12) holds; this, together with (6.11), shows that $\hat{X}_{2\varepsilon}$ is relatively compact in $C^0([0, T])$ so that from (6.10) we obtain

$$(6.13) \quad \hat{X}_{2\varepsilon} \rightarrow \hat{X}_2 := \int_{\Omega} \theta_1 u'_1 \Phi' \, dx + \int_{\Omega} u'_2 \Psi' \, dx + \int_{\Omega} A^0 \nabla \Phi \nabla u_1 \, dx \text{ strongly in } C^0([0, T]).$$

Concerning $\hat{X}_{3\varepsilon}$, a similar argument to that used in the proof of Proposition 5.1 to show (5.11) gives

$$(6.14) \quad \hat{X}_{3\varepsilon} \rightarrow \hat{X}_3 := \frac{1}{2} \left[\theta_1 \|\Phi'\|_{L^2(\Omega)}^2 + \theta_2 \|\Psi'\|_{L^2(\Omega)}^2 + \int_{\Omega} A^0 \nabla \Phi \nabla \Phi \, dx \right] \text{ strongly in } C^0([0, T]).$$

Step 4. Observe first that from (6.2), (6.13), and (6.14) we have

$$(6.15) \quad 0 \leq \limsup_{\varepsilon \rightarrow 0} \hat{X}_\varepsilon(t) \leq \hat{X}_1(t) - \hat{X}_2(t) + \hat{X}_3(t) = \hat{X}(t) \quad \forall t \in [0, T],$$

where \hat{X} is given by (6.1).

On the other hand, if we write \hat{X}_ε as

$$(6.16) \quad \hat{X}_\varepsilon = Z_\varepsilon - \hat{Y}_\varepsilon,$$

where

$$(6.17) \quad Z_\varepsilon := E^\varepsilon - Y_\varepsilon - \hat{X}_{2\varepsilon} + \hat{X}_{3\varepsilon}$$

and \hat{Y}_ε is given by (6.6), due to (6.8) we have

$$(6.18) \quad \limsup_{\varepsilon \rightarrow 0} \|\hat{X}_\varepsilon\|_{C^0([0, T])} \leq \limsup_{\varepsilon \rightarrow 0} \|Z_\varepsilon\|_{C^0([0, T])} + \lim_{\varepsilon \rightarrow 0} |\hat{Y}_\varepsilon| = \limsup_{\varepsilon \rightarrow 0} \|Z_\varepsilon\|_{C^0([0, T])}.$$

Consequently, to conclude, it is enough to prove that

$$(6.19) \quad \limsup_{\varepsilon \rightarrow 0} \|Z_\varepsilon\|_{C^0([0, T])} \leq \|\hat{X}\|_{C^0([0, T])},$$

which together with (6.18) will give the claimed result.

Let $\{\varepsilon'\}$ be a subsequence such that

$$(6.20) \quad \limsup_{\varepsilon \rightarrow 0} \|Z_\varepsilon\|_{C^0([0, T])} = \lim_{\varepsilon' \rightarrow 0} \|Z_{\varepsilon'}\|_{C^0([0, T])}.$$

From (6.7), (6.13), and (6.14) it follows that Z_ε is compact in $C^0([0, T])$. Therefore, there exists a subsequence (still denoted $\{\varepsilon'\}$) such that

$$(6.21) \quad Z_{\varepsilon'} \rightarrow Z \text{ strongly in } C^0([0, T]).$$

But, from (6.16), (6.15), and (6.8) we have

$$0 \leq Z(t) = \lim_{\varepsilon' \rightarrow 0} [\hat{X}_{\varepsilon'}(t) + \hat{Y}_{\varepsilon'}(t)] = \lim_{\varepsilon' \rightarrow 0} \hat{X}_{\varepsilon'}(t) \leq X(t) \quad \forall t \in [0, T]$$

so that from (6.21) one has

$$\lim_{\varepsilon' \rightarrow 0} \|Z_{\varepsilon'}\|_{C^0([0, T])} = \|Z\|_{C^0([0, T])} \leq \|\hat{X}\|_{C^0([0, T])},$$

which together with (6.20) gives (6.19) and ends the proof. \square

Remark 6.3. Let us emphasize the main difficulties in the proof of Proposition 6.1. First, as in the previous case (see Remark 5.3), due to the presence of the boundary term in the energy, we have only $X_{1\varepsilon} \leq E^\varepsilon$ and not $X_{1\varepsilon} = E^\varepsilon$, as usual. Moreover, here the limit of $E^\varepsilon - \hat{X}_{2\varepsilon} + \hat{X}_{3\varepsilon}$ is bigger than \hat{X} so that we cannot conclude as in the previous case. On the other hand, since we have no estimates for the time derivative of the boundary term of $X_{1\varepsilon}$, we cannot derive any compactness of $X_{1\varepsilon}$ in $C^0([0, T])$. This is why we introduced in Step 4 of the proof a more technical argument, which relies on the decomposition (6.17) of \hat{X}_ε as a sum of a compact term $Z_\varepsilon - Y_\varepsilon$ and a rest \hat{Y}_ε , which goes to zero in $C^0([0, T])$ as $\varepsilon \rightarrow 0$.

Proof of Theorem 2.8. Since $u_1 \in L^2(0, T; H_0^1(\Omega)) \cap C^0([0, T]; L^2(\Omega))$, $u'_1 \in C^0([0, T]; L^2(\Omega))$, and $u'_2 \in C^0([0, T]; L^2(\Omega))$, from classical density results, we have that for any $\delta > 0$ there exist $\Phi, \Psi \in C^\infty([0, T]; \mathcal{D}(\Omega))$ such that

$$(6.22) \quad \begin{cases} \text{(i)} \ \|u'_1 - \Phi'\|_{C^0([0, T]; L^2(\Omega))} \leq \delta, \\ \text{(ii)} \ \|\theta_2^{-1}u'_2 - \Psi'\|_{C^0([0, T]; L^2(\Omega))} \leq \delta, \\ \text{(iii)} \ \|\nabla u_1 - \nabla \Phi\|_{L^2(0, T; L^2(\Omega))} \leq \delta. \end{cases}$$

Therefore,

$$(6.23) \quad \begin{aligned} & \|u'_{1\varepsilon} - u'_1\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))}^2 + \|u'_{2\varepsilon} - \theta_2^{-1}u'_2\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))}^2 \\ & \leq 2 \left(\|u'_{1\varepsilon} - \Phi'\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))}^2 + \|\Phi' - u'_1\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))}^2 \right. \\ & \quad \left. + \|u'_{2\varepsilon} - \Psi'\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))}^2 + \|\Psi' - \theta_2^{-1}u'_2\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))}^2 \right) \\ & \leq 2 \left(\|u'_{1\varepsilon} - \Phi'\|_{C^0(0, T; L^2(\Omega_{1\varepsilon}))}^2 + \|u'_{2\varepsilon} - \Psi'\|_{C^0(0, T; L^2(\Omega_{2\varepsilon}))}^2 \right) + 2\delta^2. \end{aligned}$$

The same argument used to prove (5.15) gives that there exists a constant C_1 such that

$$(6.24) \quad \begin{aligned} & \left(\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla u_1\|_{C^0([0, T]; [L^1(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0, T]; [L^2(\Omega_{2\varepsilon})]^n)}^2 \right) \\ & \leq C_1 \left(\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{C^0([0, T]; [L^2(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0, T]; [L^2(\Omega_{2\varepsilon})]^n)}^2 + \delta^2 \right). \end{aligned}$$

On the other hand, from the properties of A^ε ,

$$\begin{aligned} & \frac{1}{2} \left[\|u'_{1\varepsilon}(t) - \Phi'(t)\|_{L^2(\Omega_{1\varepsilon})}^2 + \|u'_{2\varepsilon}(t) - \Psi'(t)\|_{L^2(\Omega_{2\varepsilon})}^2 \right. \\ & \quad \left. + \alpha \left(\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{[L^2(\Omega_{1\varepsilon})]^n}^2 + \|\nabla u_{2\varepsilon}\|_{[L^2(\Omega_{2\varepsilon})]^n}^2 \right) \right] \leq \hat{X}_\varepsilon \quad \forall t \in [0, T] \end{aligned}$$

and there exists a constant C_2 such that

$$\|\hat{X}\|_{C^0([0, T])} \leq C_2 \delta^2,$$

where \hat{X}_ε and \hat{X} , defined in Proposition 6.1, are applied to the function Φ and Ψ , given by (6.22).

Therefore, from Proposition 6.1 we conclude

$$\begin{aligned} & \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{1}{2} \left[\|u'_{1\varepsilon} - \Phi'\|_{C^0([0,T];[L^2(\Omega_{1\varepsilon})])}^2 + \|u'_{2\varepsilon} - \Psi'\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})])}^2 \right] \right. \\ & \quad \left. + \frac{\alpha}{2} \left(\|\nabla u_{1\varepsilon} - C^\varepsilon \nabla \Phi\|_{C^0([0,T];[L^2(\Omega_{1\varepsilon})]^n)}^2 + \|\nabla u_{2\varepsilon}\|_{C^0([0,T];[L^2(\Omega_{2\varepsilon})]^n)}^2 \right) \right\} \\ & \leq \limsup_{\varepsilon \rightarrow 0} \|\hat{X}_\varepsilon\|_{C^0([0,T])} \leq \|\hat{X}\|_{C^0([0,T])} \leq C_2 \delta^2, \end{aligned}$$

which, together with (6.23) and (6.24), gives (2.24), since δ is arbitrary. \square

REFERENCES

- [1] J.L. AURIAULT AND H. ENE, *Macroscopic modelling of heat transfer in composites with interfacial thermal barrier*, Internat. J. Heat Mass Transfer, 37 (1994), pp. 2885–2892.
- [2] S. BRAHIM-OTSMANE, G.A. FRANCFORT, AND F. MURAT, *Correctors for the homogenization of the wave and heat equations*, J. Math. Pure. Appl., 71 (1992), pp. 197–231.
- [3] E. CANON AND J.N. PERNIN, *Homogenization of diffusion in composite media with interfacial barrier*, Rev. Roumaine Math. Pures Appl., 44 (1999), pp. 23–36.
- [4] H.S. CARSLAW AND J.C. JAEGER, *Conduction of Heat in Solids*, Clarendon Press, Oxford, 1947.
- [5] D. CIORANESCU AND P. DONATO, *Exact internal controllability in perforated domains*, J. Math. Pures Appl., 68 (1989), pp. 185–213.
- [6] D. CIORANESCU AND P. DONATO, *An Introduction to Homogenization*, Oxford Lecture Ser. Math., Appl. 17, Oxford University Press, New York, 1999.
- [7] D. CIORANESCU AND J. SAINT JEAN PAULIN, *Homogenization in open sets with holes*, J. Math. Anal. Appl., 71 (1979), pp. 590–607.
- [8] P. DONATO, *Some corrector results for composites with imperfect interface*, Rend. Mat. Ser. VII, 26 (2006), pp. 189–209.
- [9] P. DONATO, L. FAELLA, AND S. MONSURRÒ, *Homogenization of the wave equation in composites with imperfect interface: A memory effect*, J. Math. Pures Appl., 87 (2007), pp. 119–143.
- [10] P. DONATO AND S. MONSURRÒ, *Homogenization of two heat conductors with interfacial contact resistance*, Anal. Appl., 2 (2004), pp. 247–273.
- [11] P. DONATO AND A. NABIL, *Approximate controllability of linear parabolic equations in perforated domains*, ESAIM Control Optim. Calc. Var., 6 (2001), pp. 21–38.
- [12] H. ENE AND D. POLISEVSKI, *Model of diffusion in partially fissured media*, Z. Angew. Math. Phys., 53 (2002), pp. 1052–1059.
- [13] H.C. HUMMEL, *Homogenization for heat transfer in polycrystals with interfacial resistances*, Appl. Anal., 75 (2000), pp. 403–424.
- [14] J.L. LIONS AND E. MAGENES, *Problèmes aux Limites Non Homogènes et Applications*, Dunod, Paris, 1968.
- [15] R. LIPTON, *Heat conduction in fine scale mixtures with interfacial contact resistance*, SIAM J. Appl. Math., 58 (1998), pp. 55–72.
- [16] R. LIPTON AND B. VERNESCU, *Composite with imperfect interface*, Proc. R. Soc. Lond. Ser. A, 452 (1996), pp. 329–358.
- [17] S. MONSURRÒ, *Homogenization of a two-component composite with interfacial thermal barrier*, Adv. Math. Sci. Appl., 13 (2003), pp. 43–63.
- [18] A. NABIL, *A Corrector result for the wave equations in perforated domains*, GAKUTO Internat. Ser., Math. Sci. Appl., 9 (1997), pp. 309–321.
- [19] J.N. PERNIN, *Homogénéisation d'un problème de diffusion en milieu composite à deux composantes*, C.R. Acad. Sci. Paris Sér. I, 321 (1995), pp. 949–952.
- [20] E. ZEIDLER, *Nonlinear Functional Analysis and Its Applications*, II/B. *Nonlinear Monotone Operators*, Springer-Verlag, New York, 1990.

RELAXATION-TIME LIMIT IN THE ISOTHERMAL HYDRODYNAMIC MODEL FOR SEMICONDUCTORS*

JIANG XU[†]

Abstract. This work is concerned with the relaxation-time limit in the multidimensional isothermal hydrodynamic model for semiconductors in the critical Besov space. As the initial data are sufficiently close to equilibrium, the uniform (global) classical solutions are constructed by the high- and low-frequency decomposition methods. Furthermore, it is shown that the scaled classical solutions strongly converge towards that of a drift-diffusion model, as the relaxation time tends to zero.

Key words. relaxation-time limit, isothermal hydrodynamic model, semiconductors

AMS subject classifications. 35L45, 35B25, 35M20

DOI. 10.1137/080721893

1. Introduction and main results. In this work, we are interested in the isothermal hydrodynamic model for semiconductors, which is of the form

$$(1.1) \quad \begin{cases} n_t + \nabla \cdot (n\mathbf{u}) = 0, \\ (n\mathbf{u})_t + \nabla \cdot (n\mathbf{u} \otimes \mathbf{u}) + a^2 \nabla n = n \nabla \Phi - \frac{n\mathbf{u}}{\tau}, \\ \Delta \Phi = n - \bar{n}, \quad \Phi \rightarrow 0 \text{ as } |x| \rightarrow +\infty \end{cases}$$

for $(t, x) \in [0, +\infty) \times \mathbf{R}^N$ ($N \geq 2$). Here $n, \mathbf{u} = (u^1, u^2, \dots, u^N)^\top$ (\top transpose), and Φ stand for the electron density, the electron velocity, and the electrostatic potential, respectively; $a > 0$ is the sound speed; and $0 < \tau \leq 1$ is the (small) momentum relaxation time for electrons. The symbols ∇, Δ , and \otimes are the gradient operator, Laplacian operator, and the tensor products of two vectors, respectively. The positive constant \bar{n} stands for the density of positively charged background ions.

The system (1.1) is supplemented with the initial data

$$(1.2) \quad (n, \mathbf{u})(x, 0) = (n_0, \mathbf{u}_0)(x), \quad x \in \mathbf{R}^N.$$

Concerning the small relaxation-time analysis, we define the scaled variables as in [12]:

$$(1.3) \quad (n^\tau, \mathbf{u}^\tau, \mathbf{e}^\tau)(x, s) = \left(n, \frac{1}{\tau} \mathbf{u}, \mathbf{e} \right) \left(x, \frac{s}{\tau} \right), \quad (\mathbf{e} = \nabla \Phi).$$

Then the new variables satisfy the equations

$$(1.4) \quad \begin{cases} n_s^\tau + \nabla \cdot (n^\tau \mathbf{u}^\tau) = 0, \\ \tau^2 \mathbf{u}_s^\tau + \tau^2 (\mathbf{u}^\tau \cdot \nabla) \mathbf{u}^\tau + \mathbf{u}^\tau = \mathbf{e}^\tau - a^2 \frac{\nabla n^\tau}{n^\tau}, \\ \nabla \cdot \mathbf{e}^\tau = n^\tau - \bar{n}, \end{cases}$$

with the initial data

$$(1.5) \quad (n^\tau, \tau \mathbf{u}^\tau)(x, 0) = (n_0, \mathbf{u}_0).$$

*Received by the editors April 22, 2008; accepted for publication September 2, 2008; published electronically January 21, 2009. This work was supported by NUAA's Scientific Fund for the Introduction of Qualified Personnel.

<http://www.siam.org/journals/sima/40-5/72189.html>

[†]Department of Mathematics, Nanjing University of Aeronautics and Astronautics, Nanjing 211100, People's Republic of China (jiangxu_79@nuaa.edu.cn, jiangxu_79@yahoo.com.cn).

Let $\tau \rightarrow 0$; then, formally, we arrive at the classical drift-diffusion model

$$(1.6) \quad \begin{cases} \mathcal{N}_s + \nabla \cdot (\mathcal{N}\mathcal{E} - a^2\nabla\mathcal{N}) = 0, \\ \nabla \cdot \mathcal{E} = \mathcal{N} - \bar{n}, \\ \mathcal{N}(x, 0) = n_0. \end{cases}$$

This relaxation-time limit for the system (1.1)–(1.2) was first studied by Marcati and Natalini [12]. They obtained the uniform weak solutions with respect to τ and proved that the scaled weak solutions converged towards that of the drift-diffusion model (1.6). Subsequently, some rigorous results related to the relaxation limit have appeared (we refer the reader to [1, 3, 7, 8, 9, 10]); however, these results are restricted to one space dimension. Physically, it is more important and more interesting to study this asymptotic limit in several space dimensions; unfortunately, due to the serious difficulties in establishing the global existence of weak or smooth solutions to (1.1)–(1.2), up to now, only partial relaxation limit results are available in several space dimensions. In [11], Lattanzio and Marcati considered the three-dimensional isentropic hydrodynamic model, assuming the existence of L^∞ -solutions in an τ -independent time interval, and justified the relaxation limit. Inspired by the Maxwell iteration, Yong [15] dealt with the periodic initial-value problem for the multidimensional isentropic hydrodynamic model in the Sobolev space $H^\ell(\mathbf{T}^N)$; however, the regularity index is required to be high ($\ell > 1 + N/2$, an integer).

Recently, in [6], we first established the global existence and exponential stability of classical solutions to the (isentropic and isothermal) hydrodynamic model for semiconductors in the critical Besov space $B_{2,1}^{1+N/2}(\mathbf{R}^N)$, but it is not clear whether the results are independent of the relaxation time τ . In this paper, we construct the uniform classical solutions (close to equilibrium) to the *isothermal* hydrodynamic model (1.1)–(1.2) and justify the above formal limit rigorously. For the proof of the global result, different from that in [6], we add some “new” ideas. More concretely speaking, we take full advantage of the special structure (skew-symmetry) of (1.1)–(1.2), which can help us avoid differentiating the system with respect to variable t once, and develop some “new” frequency-localization estimates; for details, see Lemmas 4.3, 4.4, and 4.5. In this sense, the proof of the global result is shortened heavily. Here, we first state the main results as follows.

THEOREM 1.1. *Let $\bar{n} > 0$ be a constant reference density. Suppose that $n_0 - \bar{n}, \mathbf{u}_0$, and $\mathbf{e}_0 \in B_{2,1}^\sigma(\mathbf{R}^N)$ ($\sigma = 1 + N/2$). There exists a positive constant δ_0 independent of τ , such that if*

$$\|(n_0 - \bar{n}, \mathbf{u}_0, \mathbf{e}_0)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} \leq \delta_0,$$

then the system (1.1)–(1.2) admits a unique global solution $(n, \mathbf{u}, \mathbf{e})$ satisfying

$$(n - \bar{n}, \mathbf{u}, \mathbf{e}) \in \mathcal{C}([0, \infty), B_{2,1}^\sigma(\mathbf{R}^N)).$$

Moreover, the uniform energy estimate holds:

$$(1.7) \quad \begin{aligned} & \sup_{t \geq 0} \left(\|(n - \bar{n}, \mathbf{u}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} \right) \\ & + \mu_0 \int_0^\infty \left(\tau \|(n - \bar{n}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} + \|\mathbf{u}(\cdot, t)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} \right) dt \\ & \leq C_0 \|(n_0 - \bar{n}, \mathbf{u}_0, \mathbf{e}_0)\|_{B_{2,1}^\sigma(\mathbf{R}^N)}, \end{aligned}$$

where μ_0, C_0 are the positive constants independent of τ , and $\mathbf{e}_0 := \nabla\Delta^{-1}(n_0 - \bar{n})$.

Remark 1.1. From the energy estimate (1.7) and the smallness of τ ($0 < \tau \leq 1$), we can obtain the uniform exponential decay of classical solution $(n, \mathbf{u}, \mathbf{e})$ near to equilibrium $(\bar{n}, 0, 0)$ in [6]:

$$\|(n - \bar{n}, \mathbf{u}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} \leq C_0 \|(n_0 - \bar{n}, \mathbf{u}_0, \mathbf{e}_0)\|_{B_{2,1}^\sigma(\mathbf{R}^N)} \exp(-\mu_0 \tau t), \quad t \geq 0.$$

Remark 1.2. To our knowledge, Theorem 1.1 cannot directly be applied to the generally isentropic hydrodynamic model for semiconductors. As a matter of fact, for the isentropic model, we can establish a similar relaxation-limit result in stronger Besov space $B_{2,2}^{\sigma+\epsilon}(\mathbf{R}^N)$ ($\epsilon > 0$).

Then, by considering an “ $\mathcal{O}(1/\tau)$ time scale” in (1.3), we justify the following convergence to the drift-diffusion model (1.6) by the standard weak convergence methods and the application of the compactness theorem in [14].

THEOREM 1.2. *Let $(n, \mathbf{u}, \mathbf{e})$ be the global solution of (1.1)–(1.2) given by Theorem 1.1. Then there exists some function $(\mathcal{N}, \mathcal{E})$ which is a global weak solution to (1.6) satisfying $(\mathcal{N} - \bar{n}, \mathcal{E}) \in \mathcal{C}([0, \infty), B_{2,1}^\sigma(\mathbf{R}^N))$ such that as $\tau \rightarrow 0$, it yields the following ($1 \leq p < +\infty$, $\sigma' < \sigma$):*

$$(n, \mathbf{e})\left(x, \frac{s}{\tau}\right) \rightarrow (\mathcal{N}, \mathcal{E})(x, s) \quad \text{strongly in } L^p(0, T; (B_{2,1}^{\sigma'}(\mathbf{R}^N))_{\text{loc}}) \quad \text{for any } T > 0.$$

Remark 1.3. Let us mention that this strong convergence result is *weaker* than that obtained in [1, 4], which can be regarded as a supplement to the theory of diffusive limit for hyperbolic problems.

The rest of this paper is organized as follows. In section 2, we present some definitions and basic facts on the Littlewood–Paley decomposition theory and Besov space. In section 3, we rewrite the isothermal hydrodynamic model as a symmetric hyperbolic system and recall the local existence result of classical solutions. Then in section 4, we establish the uniform a priori estimate, which is used to derive the global existence of uniform classical solutions. Finally, in section 5, we perform the relaxation-time limit.

Throughout this paper, $f \approx g$ means that $f \leq Cg$ and $g \leq Cf$, where $C > 0$ is a uniform constant with respect to τ . All functional spaces are considered in \mathbf{R}^N , so we may omit the space dependence for simplicity.

2. Littlewood–Paley analysis. In this section, we review briefly the Littlewood–Paley decomposition theory and the characterization of Besov space; see also, e.g., [2] or [6].

Let \mathcal{S} be the Schwarz class. (φ, χ) is a couple of smooth functions valued in $[0, 1]$ such that φ is supported in the shell $\mathbf{C}(0, \frac{3}{4}, \frac{8}{3}) = \{\xi \in \mathbf{R}^N \mid \frac{3}{4} \leq |\xi| \leq \frac{8}{3}\}$, χ is supported in the ball $\mathbf{B}(0, \frac{4}{3}) = \{\xi \in \mathbf{R}^N \mid |\xi| \leq \frac{4}{3}\}$, and

$$\chi(\xi) + \sum_{q=0}^{\infty} \varphi(2^{-q}\xi) = 1, \quad q \in \mathbf{Z}, \quad \xi \in \mathbf{R}^N.$$

For $f \in \mathcal{S}'$ (denote the set of temperate distributions which is the dual of \mathcal{S}), one can define the nonhomogeneous dyadic blocks as follows:

$$\Delta_{-1}f := \chi(D)f = \tilde{h} * f \quad \text{with } \tilde{h} = \mathcal{F}^{-1}\chi,$$

$$\Delta_q f := \varphi(2^{-q}D)f = 2^{qN} \int h(2^q y) f(x - y) dy \quad \text{with } h = \mathcal{F}^{-1}\varphi \text{ if } q \geq 0,$$

where $*$, \mathcal{F}^{-1} represent the convolution operator and the inverse Fourier transform, respectively. The nonhomogeneous Littlewood–Paley decomposition is

$$f = \sum_{q \geq -1} \Delta_q f \quad \text{in } \mathcal{S}'.$$

Define the low-frequency cut-off by

$$S_q f := \sum_{p \leq q-1} \Delta_p f.$$

Of course, $S_0 f = \Delta_{-1} f$. The above Littlewood–Paley decomposition is almost orthogonal in L^2 .

PROPOSITION 2.1. *For any $f, g \in \mathcal{S}'$, the following properties hold:*

$$\Delta_p \Delta_q f \equiv 0 \quad \text{if } |p - q| \geq 2,$$

$$\Delta_q (S_{p-1} f \Delta_p g) \equiv 0 \quad \text{if } |p - q| \geq 5.$$

Besov space can be characterized in virtue of the Littlewood–Paley decomposition.

DEFINITION 2.2. *Let $1 \leq p \leq \infty$ and $s \in \mathbf{R}$. For $1 \leq r < \infty$, the Besov spaces $B_{p,r}^s$ are defined by*

$$f \in B_{p,r}^s \Leftrightarrow \left(\sum_{q \geq -1} (2^{qs} \|\Delta_q f\|_{L^p})^r \right)^{\frac{1}{r}} < \infty$$

and $B_{p,\infty}^s$ are defined by

$$f \in B_{p,\infty}^s \Leftrightarrow \sup_{q \geq -1} 2^{qs} \|\Delta_q f\|_{L^p} < \infty.$$

LEMMA 2.3 (Bernstein’s inequality). *Let $k \in \mathbf{N}$ and $0 < R_1 < R_2$. There exists a constant C depending only on R_1, R_2 , and N such that for all $1 \leq a \leq b \leq \infty$ and $f \in L^a$, we have*

$$\text{Supp } \mathcal{F}f \subset \mathbf{B}(0, R_1 \lambda) \Rightarrow \sup_{|\alpha|=k} \|\partial^\alpha f\|_{L^b} \leq C^{k+1} \lambda^{k+N(\frac{1}{a}-\frac{1}{b})} \|f\|_{L^a},$$

$$\text{Supp } \mathcal{F}f \subset \mathbf{C}(0, R_1 \lambda, R_2 \lambda) \Rightarrow C^{-k-1} \lambda^k \|f\|_{L^a} \leq \sup_{|\alpha|=k} \|\partial^\alpha f\|_{L^a} \leq C^{k+1} \lambda^k \|f\|_{L^a}.$$

Here, $\mathcal{F}f$ (or $\widehat{f} = \int_{\mathbf{R}^N} f(x) \exp(-ix \cdot \xi) dx$) represents the Fourier transform on f .

There is a compactness result in Besov space, which we show in the following proposition.

PROPOSITION 2.4. *Let $1 \leq p, r \leq \infty$, $s \in \mathbf{R}$, and $\varepsilon > 0$. For all $\phi \in C_c^\infty$, the map $f \mapsto \phi f$ is compact from $B_{p,r}^{s+\varepsilon}$ to $B_{p,r}^s$.*

Finally, we state a result of continuity for the composition to end this section.

PROPOSITION 2.5. *Let $1 \leq p, r \leq \infty$, and I be an open interval of \mathbf{R} . Let $s > 0$ and let n be the smallest integer such that $n \geq s$. Let $F : I \rightarrow \mathbf{R}$ satisfy $F(0) = 0$ and $F' \in W^{n,\infty}(I; \mathbf{R})$. Assume that $v \in B_{p,r}^s$ takes values in $J \subset\subset I$. Then $F(v) \in B_{p,r}^s$ and there exists a constant C depending only on s, I, J , and N such that*

$$\|F(v)\|_{B_{p,r}^s} \leq C(1 + \|v\|_{L^\infty})^n \|F'\|_{W^{n,\infty}(I)} \|v\|_{B_{p,r}^s}.$$

3. Reformulation and local existence. Let us introduce the enthalpy $\mathcal{H}(\varrho) = a^2 \ln \varrho$ ($\varrho > 0$), and set

$$(3.1) \quad m(t, x) = (\mathcal{H}(n(t, x)) - \mathcal{H}(\bar{n}))/a.$$

Then (1.1) can be transformed into the symmetric hyperbolic form

$$(3.2) \quad \partial_t U + \sum_{j=1}^N A_j(\mathbf{u}) \partial_{x_j} U = \begin{pmatrix} 0 \\ -\frac{1}{\tau} \mathbf{u} + \mathbf{e} \end{pmatrix},$$

coupled with the dynamic electron field equation

$$(3.3) \quad \text{dive} = h(m),$$

where

$$U = \begin{pmatrix} m \\ \mathbf{u} \end{pmatrix}, \quad A_j(\mathbf{u}) = \begin{pmatrix} u^j & ae_j^\top \\ ae_j & u^j I_N \end{pmatrix}$$

(I_N denotes the unit matrix of order N),

and $h(m) = \bar{n}\{\exp(m/a) - 1\}$ is a smooth function on the domain $\{m \mid -\infty < m < +\infty\}$ satisfying $h(0) = 0$. The initial data (1.2) become

$$(3.4) \quad (m, \mathbf{u}, \mathbf{e})|_{t=0} = (a(\ln n_0 - \ln \bar{n}), \mathbf{u}_0, \mathbf{e}_0).$$

Remark 3.1. The variable change is from the open set $\{(n, \mathbf{u}, \mathbf{e}) \in (0, +\infty) \times \mathbf{R}^N \times \mathbf{R}^N\}$ to the whole space $\{(m, \mathbf{u}, \mathbf{e}) \in \mathbf{R} \times \mathbf{R}^N \times \mathbf{R}^N\}$. It is easy to show that for classical solutions $(n, \mathbf{u}, \mathbf{e})$ away from the vacuum, (1.1)–(1.2) is equivalent to (3.2)–(3.4).

Now, we recall a local existence and uniqueness result of the classical solutions to (3.2)–(3.4) which has been established in [6].

PROPOSITION 3.1. *For any fixed relaxation time $\tau > 0$, assume that $(m_0, \mathbf{u}_0, \mathbf{e}_0) \in B_{2,1}^\sigma$; then, there exist a time $T_0 > 0$ (depending only on the initial data) and a unique solution $(m, \mathbf{u}, \mathbf{e})$ to (3.2)–(3.4) such that $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}^1([0, T_0] \times \mathbf{R}^N)$ and $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}([0, T_0], B_{2,1}^\sigma) \cap \mathcal{C}^1([0, T_0], B_{2,1}^{\sigma-1})$.*

4. A uniform estimate and global existence. In this section, we shall establish a uniform a priori estimate, which is used to derive the global existence of classical solutions to (3.2)–(3.4).

PROPOSITION 4.1. *There exist three positive constants δ_1, C_1 , and μ_1 independent of τ such that for any $T > 0$, if*

$$(4.1) \quad \sup_{0 \leq t \leq T} \|(m, \mathbf{u}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma} \leq \delta_1,$$

then

$$(4.2) \quad \begin{aligned} & \|(m, \mathbf{u}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma} + \mu_1 \int_0^t \left(\tau \|(m, \mathbf{e})(\cdot, \varsigma)\|_{B_{2,1}^\sigma} + \|\mathbf{u}(\cdot, \varsigma)\|_{B_{2,1}^\sigma} \right) d\varsigma \\ & \leq C_1 \|(m, \mathbf{u}, \mathbf{e})(\cdot, 0)\|_{B_{2,1}^\sigma}, \quad t \in [0, T]. \end{aligned}$$

The main ingredients in the proof of Proposition 4.1 are the high-frequency ($q \geq 0$) estimates and low-frequency ($q = -1$) estimates on $(m, \mathbf{u}, \mathbf{e})$. To do this, we need to establish some lemmas.

LEMMA 4.2. *If $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}([0, T], B_{2,1}^\sigma) \cap \mathcal{C}^1([0, T], B_{2,1}^{\sigma-1})$ is a solution of (3.2)–(3.4) for any given $T > 0$, then the following estimate holds ($q \geq -1$):*

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) + \frac{1}{\tau} \|\Delta_q \mathbf{u}\|_{L^2}^2 \\
 & \leq \frac{1}{2} \|\nabla \mathbf{u}\|_{L^\infty} (\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2) + \|[\mathbf{u}, \Delta_q] \cdot \nabla m\|_{L^2} \|\Delta_q m\|_{L^2} \\
 (4.3) \quad & + \|[\mathbf{u}, \Delta_q] \cdot \nabla \mathbf{u}\|_{L^2} \|\Delta_q \mathbf{u}\|_{L^2} + \frac{1}{\bar{n}} \|\Delta_q(h(m)\mathbf{u})\|_{L^2} \|\Delta_q \mathbf{e}\|_{L^2},
 \end{aligned}$$

where the commutator $[f, g] = fg - gf$.

Remark 4.1. By applying the operator Δ_q to (3.2), then multiplying the resulting equations by the conjugator of $\Delta_q m(\overline{\Delta_q m})$ and $\overline{\Delta_q \mathbf{u}}$, respectively, we can directly achieve (4.3) without extra trouble; see [6].

Below, we formulate an important skew-symmetry lemma which was developed in [4, 13, 16].

LEMMA 4.3 (Shizuta–Kawashima). *For all $\xi \in \mathbf{R}^N$, $\xi \neq 0$, the system (3.2) admits a real skew-symmetric smooth matrix $K(\xi)$ which is defined in the unit sphere \mathbf{S}^{N-1} :*

$$(4.4) \quad K(\xi) = \begin{pmatrix} 0 & \frac{\xi^\top}{|\xi|} \\ -\frac{\xi}{|\xi|} & 0 \end{pmatrix},$$

and then

$$(4.5) \quad K(\xi) \sum_{j=1}^N \xi_j A_j(0) = \begin{pmatrix} a|\xi| & 0 \\ 0 & -a\frac{\xi \otimes \xi}{|\xi|} \end{pmatrix}.$$

Thanks to the skew-symmetry of the system (3.2), we can develop some “new” frequency-localization estimates and avoid performing the t -derivative to (3.2), which is different from the estimates in [6].

LEMMA 4.4. *If $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}([0, T], B_{2,1}^\sigma) \cap \mathcal{C}^1([0, T], B_{2,1}^{\sigma-1})$ is a solution of (3.2)–(3.4) for any given $T > 0$, then the following estimates hold:*

$$\begin{aligned}
 & \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi + \tau \frac{a}{2} 2^{2q} \|\Delta_q m\|_{L^2}^2 \\
 & \leq \frac{C}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2 + C\tau 2^q \|\Delta_q U\|_{L^2} \|\Delta_q \mathcal{G}\|_{L^2} \\
 (4.6) \quad & + C\tau \|\Delta_q(\tilde{h}(m)m)\|_{L^2} \|\Delta_q m\|_{L^2} \quad (q \geq 0);
 \end{aligned}$$

$$\begin{aligned}
 & \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_{-1} U})^* K(\xi) \widehat{\Delta_{-1} U} \right) d\xi + \tau \frac{\bar{n}}{a} \|\Delta_{-1} m\|_{L^2}^2 \\
 & \leq \frac{C}{\tau} \|\Delta_{-1} \mathbf{u}\|_{L^2}^2 + C\tau \|\Delta_{-1} U\|_{L^2} \|\Delta_{-1} \mathcal{G}\|_{L^2} \\
 (4.7) \quad & + C\tau \|\Delta_{-1}(\tilde{h}(m)m)\|_{L^2} \|\Delta_{-1} m\|_{L^2},
 \end{aligned}$$

where the function \mathcal{G} is given in (4.9), $\tilde{h}(m) = \int_0^1 h'(\varsigma m) d\varsigma - \frac{\bar{n}}{a}$ is a smooth function on $\{|m| - \infty < \varsigma m < \infty, \varsigma \in [0, 1]\}$ satisfying $\tilde{h}(0) = 0$, and $C > 0$ is a uniform constant independent of τ .

Proof. The system (3.2) can be written as the linearized form

$$(4.8) \quad \partial_t U + \sum_{j=1}^N A_j(0) \partial_{x_j} U = \mathcal{G} + \begin{pmatrix} 0 \\ -\frac{1}{\tau} \mathbf{u} + \mathbf{e} \end{pmatrix},$$

where

$$(4.9) \quad \mathcal{G} = \sum_{j=1}^N \left\{ A_j(0) - A_j(\mathbf{u}) \right\} \partial_{x_j} U.$$

Applying the operator Δ_q to the system (4.8) gives

$$(4.10) \quad \partial_t \Delta_q U + \sum_{j=1}^N A_j(0) \partial_{x_j} \Delta_q U = \Delta_q \mathcal{G} + \begin{pmatrix} 0 \\ -\frac{1}{\tau} \Delta_q \mathbf{u} + \Delta_q \mathbf{e} \end{pmatrix}.$$

By performing the Fourier transform with respect to the space variable x for (4.10) and multiplying the resulting equation by $-i\tau(\widehat{\Delta_q U})^* K(\xi)$ ($*$ represents transpose and conjugator), then taking the real part of each term in the equality, we can obtain

$$(4.11) \quad \begin{aligned} & \tau \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \frac{d}{dt} \widehat{\Delta_q U} \right) + \tau (\widehat{\Delta_q U})^* K(\xi) \left(\sum_{j=1}^N \xi_j A_j(0) \right) \widehat{\Delta_q U} \\ & = \tau \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) (\widehat{\Delta_q \mathcal{G}}) \right) - \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{u}} \right) + \tau \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{e}} \right). \end{aligned}$$

Using the skew-symmetry of $K(\xi)$, we have

$$(4.12) \quad \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \frac{d}{dt} \widehat{\Delta_q U} \right) = \frac{1}{2} \frac{d}{dt} \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right).$$

Substituting (4.5) into the second term on the left-hand side of (4.11), it is not difficult to get

$$(4.13) \quad \begin{aligned} & \tau \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \frac{d}{dt} \widehat{\Delta_q U} \right) + \tau (\widehat{\Delta_q U})^* K(\xi) \left(\sum_{j=1}^N \xi_j A_j(0) \right) \widehat{\Delta_q U} \\ & \geq \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) + a\tau |\xi| |\widehat{\Delta_q U}|^2 - 2a\tau |\xi| |\widehat{\Delta_q \mathbf{u}}|^2. \end{aligned}$$

With the help of Young’s inequality, the right-hand side of (4.11) can be estimated as

$$(4.14) \quad \begin{aligned} & \tau \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) (\widehat{\Delta_q \mathcal{G}}) \right) - \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{u}} \right) + \tau \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{e}} \right) \\ & \leq \tau \frac{a}{2} |\xi| |\widehat{\Delta_q U}|^2 + \frac{C}{\tau |\xi|} |\widehat{\Delta_q \mathbf{u}}|^2 + \tau |\widehat{\Delta_q U}| |\widehat{\Delta_q \mathcal{G}}| + \tau \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{e}} \right), \end{aligned}$$

where we have used the uniform boundedness of the matrix $K(\xi)$ ($\xi \neq 0$); the positive constant C is independent of τ . Combining this with the equality (4.11) and the inequality (4.13)–(4.14), we deduce that

$$(4.15) \quad \begin{aligned} & \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) + \tau \frac{a}{2} |\xi| |\widehat{\Delta_q U}|^2 \\ & \leq \frac{C}{\tau} \left(|\xi| + \frac{1}{|\xi|} \right) |\widehat{\Delta_q \mathbf{u}}|^2 + \tau |\widehat{\Delta_q U}| |\widehat{\Delta_q \mathcal{G}}| + \tau \operatorname{Im} \left(\overline{(\widehat{\Delta_q m})} \frac{\xi^\top}{|\xi|} \widehat{\Delta_q \mathbf{e}} \right). \end{aligned}$$

Multiplying (4.15) by $|\xi|$ and integrating it over \mathbf{R}^N , from Plancherel’s theorem, we obtain

$$\begin{aligned}
 & \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi + \tau \frac{a}{2} \|\Delta_q \nabla U\|_{L^2}^2 \\
 (4.16) \quad & \leq \frac{C}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2 + C\tau 2^q \|\Delta_q U\|_{L^2} \|\Delta_q \mathcal{G}\|_{L^2} + \tau \operatorname{Im} \int \left((\overline{\widehat{\Delta_q m}}) \xi^\top \widehat{\Delta_q \mathbf{e}} \right) d\xi.
 \end{aligned}$$

For the third term on the right-hand side of (4.16), we have

$$(4.17) \quad \tau \operatorname{Im} \int \left((\overline{\widehat{\Delta_q m}}) \xi^\top \widehat{\Delta_q \mathbf{e}} \right) d\xi = \tau (J + \bar{J}),$$

where

$$\begin{aligned}
 J + \bar{J} &= -\frac{i}{2} \int \left((\overline{\widehat{\Delta_q m}}) \xi^\top \widehat{\Delta_q \mathbf{e}} \right) d\xi + \frac{i}{2} \int \left((\widehat{\Delta_q m}) \xi^\top \overline{\widehat{\Delta_q \mathbf{e}}} \right) d\xi \\
 &= \frac{1}{2} \int (\overline{\Delta_q \nabla m}) \cdot \widehat{\Delta_q \mathbf{e}} d\xi + \frac{1}{2} \int (\Delta_q \nabla m) \cdot \overline{\widehat{\Delta_q \mathbf{e}}} d\xi \\
 &= \frac{1}{2} (2\pi)^N \left(\int \overline{\Delta_q \nabla m} \cdot \Delta_q \mathbf{e} dx + \int \Delta_q \nabla m \cdot \overline{\Delta_q \mathbf{e}} dx \right) \\
 &= -\frac{1}{2} (2\pi)^N \left(\int \overline{\Delta_q m} \Delta_q \operatorname{div} dx + \int \Delta_q m \overline{\Delta_q \operatorname{div}} dx \right) \\
 &= -\frac{1}{2} (2\pi)^N \left(\int \overline{\Delta_q m} \Delta_q (h(m) - h(0)) dx + \int \Delta_q m \overline{\Delta_q (h(m) - h(0))} dx \right) \\
 &= -\frac{(2\pi)^N \bar{n}}{a} \|\Delta_q m\|_{L^2}^2 - \frac{(2\pi)^N}{2} \left(\int \overline{\Delta_q m} \Delta_q (\tilde{h}(m)m) dx \right. \\
 (4.18) \quad & \left. + \int \Delta_q m \overline{\Delta_q (\tilde{h}(m)m)} dx \right).
 \end{aligned}$$

Here, $\tilde{h}(m) = \int_0^1 h'(\zeta m) d\zeta - \bar{n}/a$ is a smooth function on $\{|m| - \infty < \zeta m < \infty, \zeta \in [0, 1]\}$ satisfying $\tilde{h}(0) = 0$. Therefore, from (4.16)–(4.18), we arrive at

$$\begin{aligned}
 & \frac{\tau}{2} \frac{d}{dt} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi + \tau \frac{a}{2} \|\Delta_q \nabla U\|_{L^2}^2 + \tau \frac{(2\pi)^N \bar{n}}{a} \|\Delta_q m\|_{L^2}^2 \\
 (4.19) \quad & \leq \frac{C}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2 + C\tau 2^q \|\Delta_q U\|_{L^2} \|\Delta_q \mathcal{G}\|_{L^2} + C\tau \|\Delta_q (\tilde{h}(m)m)\|_{L^2} \|\Delta_q m\|_{L^2}.
 \end{aligned}$$

In view of Lemma 2.3,

$$\|\Delta_q \nabla m\|_{L^2} \approx 2^q \|\Delta_q m\|_{L^2} \quad (q \geq 0)$$

follows, so we get the estimates (4.6) and (4.7) immediately. \square

From the proof of Lemma 4.4, we see that the Poisson equation remedies the low-frequency estimate on $\|\Delta_{-1} m\|_{L^2}^2$, which plays a key role in the uniform exponential decay of classical solutions (see Remark 1.1). This is essentially different from the Euler equations studied in [5]. On the electron field \mathbf{e} , we also have some “new” a priori estimates.

LEMMA 4.5. *If $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}([0, T], B_{2,1}^\sigma) \cap \mathcal{C}^1([0, T], B_{2,1}^{\sigma-1})$ is a solution of (3.2)–(3.4) for any given $T > 0$, then*

$$(4.20) \quad 2^q \|\Delta_q \mathbf{e}\|_{L^2}^2 \leq C \left(\frac{\bar{n}}{a} \|\Delta_q m\|_{L^2} + \|\Delta_q (\tilde{h}(m)m)\|_{L^2} \right) \|\Delta_q \mathbf{e}\|_{L^2} \quad (q \geq 0);$$

$$\begin{aligned}
 & -\frac{d}{dt} \int \Delta_{-1} \mathbf{e} \cdot \overline{\Delta_{-1} \mathbf{u}} dx + \|\Delta_{-1} \mathbf{e}\|_{L^2}^2 \\
 & \leq C(\bar{n} \|\Delta_{-1} \mathbf{u}\|_{L^2} + \|\Delta_{-1}(h(m)\mathbf{u})\|_{L^2}) \|\Delta_{-1} \mathbf{u}\|_{L^2} + C \left(a \|\Delta_{-1} \nabla m\|_{L^2} \right. \\
 (4.21) \quad & \left. + \frac{1}{\tau} \|\Delta_{-1} \mathbf{u}\|_{L^2} + \|\mathbf{u}\|_{L^\infty} \|\Delta_{-1} \nabla \mathbf{u}\|_{L^2} + \|[\mathbf{u}, \Delta_{-1}] \cdot \nabla \mathbf{u}\|_{L^2} \right) \|\Delta_{-1} \mathbf{e}\|_{L^2},
 \end{aligned}$$

where $C > 0$ is a uniform constant independent of τ .

Proof. (I) By applying the operator Δ_q to both sides of $\text{dive} = h(m)$ ($q \geq 0$), integrating it over \mathbf{R}^N after multiplying $\Delta_q \text{dive}$, and noticing the irrotationality of \mathbf{e} , we can obtain (4.20) in virtue of Hölder's inequality.

(II) From (1.1) and (3.1), we get

$$(4.22) \quad \mathbf{e}_t = -\nabla \Delta^{-1} \nabla \cdot \{h(m)\mathbf{u} + \bar{n}\mathbf{u}\},$$

where the nonlocal term $\nabla \Delta^{-1} \nabla \cdot f$ is the product of Riesz transforms on f . From (3.2) and (4.22), we have

$$\begin{aligned}
 & -\frac{d}{dt} \int \Delta_q \mathbf{e} \cdot \overline{\Delta_q \mathbf{u}} dx \\
 & = -\int \Delta_q \mathbf{e}_t \overline{\Delta_q \mathbf{u}} dx - \int \Delta_q \mathbf{e} \overline{\Delta_q \mathbf{u}_t} dx \\
 & = \int \nabla \Delta^{-1} \nabla \cdot \Delta_q \{h(m)\mathbf{u} + \bar{n}\mathbf{u}\} \overline{\Delta_q \mathbf{u}} dx \\
 (4.23) \quad & - \int \Delta_q \mathbf{e} \cdot \left(-a \overline{\Delta_q \nabla m} - \frac{1}{\tau} \overline{\Delta_q \mathbf{u}} - \mathbf{u} \overline{\Delta_q \nabla \mathbf{u}} + \overline{[\mathbf{u}, \Delta_q] \cdot \nabla \mathbf{u}} + \overline{\Delta_q \mathbf{e}} \right) dx.
 \end{aligned}$$

Using the L^2 -boundedness of the Riesz transform and Hölder's inequality, we derive (4.21) immediately. \square

For those estimates of the commutators in (4.3) and (4.21), we have the following fact.

LEMMA 4.6 (see [6]). *Let $s > 0$ and $1 < p < \infty$; then the following inequalities are true:*

$$\begin{aligned}
 & 2^{qs} \| [f, \Delta_q] \mathcal{A}g \|_{L^p} \\
 (4.24) \quad & \leq \begin{cases} Cc_q \|f\|_{B_{p,1}^s} \|g\|_{B_{p,1}^s}, & f, g \in B_{p,1}^s, \quad s = 1 + N/p, \\ Cc_q \|f\|_{B_{p,1}^{s+1}} \|g\|_{B_{p,1}^{s+1}}, & f \in B_{p,1}^s, \quad g \in B_{p,1}^{s+1}, \quad s = N/p, \\ Cc_q \|f\|_{B_{p,1}^{s+1}} \|g\|_{B_{p,1}^s}, & f \in B_{p,1}^{s+1}, \quad g \in B_{p,1}^s, \quad s = N/p. \end{cases}
 \end{aligned}$$

In particular, if $f = g$, then

$$(4.25) \quad 2^{qs} \| [f, \Delta_q] \mathcal{A}g \|_{L^p} \leq Cc_q \|\nabla f\|_{L^\infty} \|g\|_{B_{p,1}^s}, \quad s > 0,$$

where the operator $\mathcal{A} = \text{div}$ or ∇ , C is a harmless constant, and c_q denotes a sequence such that $\|(c_q)\|_{l^1} \leq 1$.

In addition, in the proof of Proposition 4.1, an elementary algebra inequality will be used. For clarity, we put it into a lemma.

LEMMA 4.7. *Assume that $f, g, h > 0$ and $0 < \tau \leq 1$; then the following inequality is true:*

$$(4.26) \quad f + \tau g + \tau h \leq C \frac{\frac{1}{\tau} f^2 + \tau g^2 + \tau h^2}{(f^2 + g^2 + h^2)^{1/2}},$$

where $C > 0$ is a uniform constant independent of τ .

Proof. Note that

$$(4.27) \quad (f^2 + g^2 + h^2)^{1/2} \approx f + g + h.$$

We need only show

$$(4.28) \quad (f + \tau g + \tau h)(f + g + h) \leq C \left(\frac{1}{\tau} f^2 + \tau g^2 + \tau h^2 \right),$$

which is obvious in virtue of the smallness of τ ($0 < \tau \leq 1$) and Young's inequality. \square

Proof of Proposition 4.1. From the a priori estimate assumption (4.1), we have

$$(4.29) \quad \sup_{0 \leq t \leq T} (\| (m, \mathbf{u}, \mathbf{e})(\cdot, t) \|_{W^{1,\infty}}) \leq C \delta_1.$$

In what follows, we are going to divide the proof into some lemmas.

LEMMA 4.8 ($q \geq 0$). *There exist some positive constants K_1, K_2, μ_2 independent of τ such that the following estimate holds:*

$$(4.30) \quad \begin{aligned} & 2^{q(\sigma-1)} \frac{d}{dt} \left\{ \frac{K_1}{2} 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) \right. \\ & \left. + \frac{K_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \right\}^{1/2} \\ & + \mu_2 2^{q\sigma} \left(\tau \|\Delta_q m\|_{L^2} + \|\Delta_q \mathbf{u}\|_{L^2} + \tau \|\Delta_q \mathbf{e}\|_{L^2} \right) \\ & \leq C \{ 2^{q\sigma} \|\nabla \mathbf{u}\|_{L^\infty} (\|\Delta_q m\|_{L^2} + \|\Delta_q \mathbf{u}\|_{L^2}) + c_q \|\mathbf{u}\|_{B_{2,1}^{\sigma-1}} (\|m\|_{B_{2,1}^{\sigma-1}} + \|\mathbf{u}\|_{B_{2,1}^{\sigma-1}}) \\ & + \tau 2^{q(\sigma-1)} \|\Delta_q \mathcal{G}\|_{L^2} + 2^{q\sigma} (\|\Delta_q (h(m)\mathbf{u})\|_{L^2} + \tau \|\Delta_q (\tilde{h}(m)m)\|_{L^2}) \}, \end{aligned}$$

where K_1, K_2 are given in (4.32), and $C > 0$ is a uniform constant independent of τ .

Proof. Combining (4.3), (4.6), and (4.20), we have

$$(4.31) \quad \begin{aligned} & \frac{d}{dt} \left\{ \frac{K_1}{2} 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) \right. \\ & \left. + \frac{K_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \right\} \\ & + \frac{K_1}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2 + \tau \frac{a}{2} K_2 2^{2q} \|\Delta_q m\|_{L^2}^2 + K_3 \tau 2^{2q} \|\Delta_q \mathbf{e}\|_{L^2}^2 \\ & \leq K_1 2^{2q} \left\{ \frac{1}{2} \|\nabla \mathbf{u}\|_{L^\infty} (\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2) \right. \\ & \quad + \|[\mathbf{u}, \Delta_q] \cdot \nabla m\|_{L^2} \|\Delta_q m\|_{L^2} + \|[\mathbf{u}, \Delta_q] \cdot \nabla \mathbf{u}\|_{L^2} \|\Delta_q \mathbf{u}\|_{L^2} \\ & \quad \left. + \frac{1}{\bar{n}} \|\Delta_q (h(m)\mathbf{u})\|_{L^2} \|\Delta_q \mathbf{e}\|_{L^2} \right\} + K_2 \left\{ \frac{C}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2 \right. \\ & \quad \left. + C \tau 2^{2q} \|\Delta_q U\|_{L^2} \|\Delta_q \mathcal{G}\|_{L^2} + C \tau \|\Delta_q (\tilde{h}(m)m)\|_{L^2} \|\Delta_q m\|_{L^2} \right\} \\ & \quad + K_3 2^q \tau C \left\{ \left(\frac{\bar{n}}{a} \|\Delta_q m\|_{L^2} + \|\Delta_q (\tilde{h}(m)m)\|_{L^2} \right) \|\Delta_q \mathbf{e}\|_{L^2} \right\}, \end{aligned}$$

where these positive constants K_1, K_2 , and K_3 (independent of τ) satisfy

$$(4.32) \quad K_2 = \frac{K_1}{4C}, \quad K_3 = \frac{a^3}{2C^2 \bar{n}^2} K_2.$$

We introduce them in order to ensure that

$$\begin{aligned}
 & \frac{K_1}{2} 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) \\
 & + \frac{K_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \\
 (4.33) \quad & \approx 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \|\Delta_q \mathbf{e}\|_{L^2}^2 \right),
 \end{aligned}$$

since

$$(4.34) \quad \left| \frac{K_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \right| \leq \frac{CK_2}{2} 2^{2q} (\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2),$$

and we eliminate the quadratic term $\frac{C}{\tau} 2^{2q} \|\Delta_q \mathbf{u}\|_{L^2}^2$, $K_3 2^q \tau C \frac{\bar{n}}{a} \|\Delta_q m\|_{L^2} \|\Delta_q \mathbf{e}\|_{L^2}$ in the right-hand side of (4.31) with the aid of Young’s inequality; for similar details, see [6]. Then, dividing the resulting inequality by

$$\left\{ \frac{K_1}{2} 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) + \frac{K_2}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \right\}^{1/2}$$

and after eliminating the quadratic terms and multiplying the factor $2^{q(\sigma-1)}$ on both sides of the inequality, we arrive at (4.30) immediately with the help of Lemmas 4.6 and 4.7. \square

Similarly, for the case of low frequency ($q = -1$), we have the following a priori estimate.

LEMMA 4.9 ($q = -1$). *There exist some positive constants $\bar{K}_1, \bar{K}_2, \bar{K}_3$, and μ_3 independent of τ such that the following estimate holds:*

$$\begin{aligned}
 & 2^{-(\sigma-1)} \frac{d}{dt} \left\{ \frac{\bar{K}_1}{2} 2^{-2} \left(\|\Delta_{-1} m\|_{L^2}^2 + \|\Delta_{-1} \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_{-1} \mathbf{e}\|_{L^2}^2 \right) \right. \\
 & + \frac{\bar{K}_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_{-1} U})^* K(\xi) \widehat{\Delta_{-1} U} \right) d\xi - \bar{K}_3 \tau \int \Delta_{-1} \mathbf{e} \cdot \overline{\Delta_{-1} \mathbf{u}} dx \left. \right\}^{1/2} \\
 & + \mu_3 2^{-\sigma} \left(\tau \|\Delta_{-1} m\|_{L^2} + \|\Delta_{-1} \mathbf{u}\|_{L^2} + \tau \|\Delta_{-1} \mathbf{e}\|_{L^2} \right) \\
 & \leq C \{ 2^{-\sigma} (\|\nabla \mathbf{u}\|_{L^\infty} + \|\mathbf{u}\|_{L^\infty}) (\|\Delta_{-1} m\|_{L^2} + \|\Delta_{-1} \mathbf{u}\|_{L^2}) \\
 & + c_{-1} \|\mathbf{u}\|_{B_{2,1}^\sigma} (\|m\|_{B_{2,1}^\sigma} + \|\mathbf{u}\|_{B_{2,1}^\sigma}) + \tau 2^{-(\sigma-1)} \|\Delta_{-1} \mathcal{G}\|_{L^2} \\
 (4.35) \quad & + 2^{-\sigma} (\|\Delta_{-1} (h(m) \mathbf{u})\|_{L^2} + \tau \|\Delta_{-1} (\tilde{h}(m) m)\|_{L^2}) \},
 \end{aligned}$$

where $C > 0$ is a uniform constant independent of τ .

Remark 4.2. Similar to the proof of Lemma 4.8, the constants $\bar{K}_1, \bar{K}_2, \bar{K}_3$ are used to ensure that

$$\begin{aligned}
 & \frac{\bar{K}_1}{2} 2^{-2} \left(\|\Delta_{-1} m\|_{L^2}^2 + \|\Delta_{-1} \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_{-1} \mathbf{e}\|_{L^2}^2 \right) \\
 & + \frac{\bar{K}_2 \tau}{2} \operatorname{Im} \int |\xi| \left((\widehat{\Delta_{-1} U})^* K(\xi) \widehat{\Delta_{-1} U} \right) d\xi - \bar{K}_3 \tau \int \Delta_{-1} \mathbf{e} \cdot \overline{\Delta_{-1} \mathbf{u}} dx \\
 (4.36) \quad & \approx 2^{-2} \left(\|\Delta_{-1} m\|_{L^2}^2 + \|\Delta_{-1} \mathbf{u}\|_{L^2}^2 + \|\Delta_{-1} \mathbf{e}\|_{L^2}^2 \right)
 \end{aligned}$$

and to eliminate some quadratic terms in the right-hand side of inequality (4.35).

Summing (4.30) on $q \in \mathbf{N} \cup \{0\}$ and adding with (4.35), then, according to the smallness of τ ($0 < \tau \leq 1$), a priori assumption (4.1), (4.29), and Moser's estimates (Proposition 2.5), we obtain the following differential inequality:

$$(4.37) \quad \begin{aligned} & \frac{d}{dt}Q + \mu_4 \left(\tau \|m\|_{B_{2,1}^\sigma} + \|\mathbf{u}\|_{B_{2,1}^\sigma} + \tau \|\mathbf{e}\|_{B_{2,1}^\sigma} \right) \\ & \leq C\delta_1 \left(\tau \|m\|_{B_{2,1}^\sigma} + \|\mathbf{u}\|_{B_{2,1}^\sigma} + \tau \|\mathbf{e}\|_{B_{2,1}^\sigma} \right), \end{aligned}$$

where

$$(4.38) \quad \begin{aligned} Q = & \sum_{q \geq 0} 2^{q(\sigma-1)} \left\{ \frac{K_1}{2} 2^{2q} \left(\|\Delta_q m\|_{L^2}^2 + \|\Delta_q \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_q \mathbf{e}\|_{L^2}^2 \right) \right. \\ & + \frac{K_2 \tau}{2} \text{Im} \int |\xi| \left((\widehat{\Delta_q U})^* K(\xi) \widehat{\Delta_q U} \right) d\xi \Big\}^{1/2} \\ & + \left\{ \frac{\bar{K}_1}{2} 2^{-2} \left(\|\Delta_{-1} m\|_{L^2}^2 + \|\Delta_{-1} \mathbf{u}\|_{L^2}^2 + \frac{1}{\bar{n}} \|\Delta_{-1} \mathbf{e}\|_{L^2}^2 \right) \right. \\ & \left. + \frac{\bar{K}_2 \tau}{2} \text{Im} \int |\xi| \left((\widehat{\Delta_{-1} U})^* K(\xi) \widehat{\Delta_{-1} U} \right) d\xi - \bar{K}_3 \tau \int \Delta_{-1} \mathbf{e} \cdot \overline{\Delta_{-1} \mathbf{u}} dx \right\}^{1/2} \end{aligned}$$

and $\mu_4 > 0$ is a uniform constant independent of τ . Note that

$$(4.39) \quad Q \approx \|(m, \mathbf{u}, \mathbf{e})(\cdot, t)\|_{B_{2,1}^\sigma}, \quad t \geq 0,$$

and by choosing $\delta_1 = \frac{\mu_4}{2C}, \mu_1 = \frac{\mu_4}{2}$, we conclude the proof of Proposition 4.1.

Based on Propositions 3.1 and 4.1, we establish the global existence of classical solutions to the system (3.2)–(3.4) by virtue of the standard continuation argument. Using the imbedding property in Besov space $B_{2,1}^\sigma$, $(m, \mathbf{u}, \mathbf{e}) \in \mathcal{C}^1([0, \infty) \times \mathbf{R}^N)$ solves (3.2)–(3.4). From Remark 3.1, we know $(n, \mathbf{u}, \mathbf{e}) \in \mathcal{C}^1([0, \infty) \times \mathbf{R}^N)$ is a solution of (1.1)–(1.2) with $n > 0$. Furthermore, we attain Theorem 1.1.

5. Relaxation-time limit. In this section, we give the proof of Theorem 1.2.

Proof. From the uniform energy estimate (1.7) in Theorem 1.1 and the scaled variables (1.3), we have

$$(5.1) \quad \begin{aligned} & \sup_{s \geq 0} \left(\|(n^\tau - \bar{n}, \tau \mathbf{u}^\tau, \mathbf{e}^\tau)(\cdot, s)\|_{B_{2,1}^\sigma} \right) + \mu_0 \int_0^\infty \left(\|(n^\tau - \bar{n}, \mathbf{u}^\tau, \mathbf{e}^\tau)(\cdot, s)\|_{B_{2,1}^\sigma} \right) ds \\ & \leq C_0 \|(n_0 - \bar{n}, \mathbf{u}_0, \mathbf{e}_0)\|_{B_{2,1}^\sigma}. \end{aligned}$$

According to Proposition 2.4 and the compactness theorem in [14], there exists a function $(\mathcal{N}, \mathcal{E}) \in \mathcal{C}([0, \infty), \bar{n} + B_{2,1}^\sigma) \times \mathcal{C}([0, \infty), B_{2,1}^\sigma)$ such that, for any $T > 0$, the sequence (up to a subsequence)

$$\{n^\tau\} \rightarrow \mathcal{N} \quad \text{strongly in } L^p(0, T; (B_{2,1}^{\sigma'})_{\text{loc}}) (1 \leq p < +\infty, \sigma' < \sigma),$$

$$\{\mathbf{e}^\tau\} \rightarrow \mathcal{E} \quad \text{strongly in } L^p(0, T; (B_{2,1}^{\sigma'})_{\text{loc}}) (1 \leq p < +\infty, \sigma' < \sigma), \text{ as } \tau \rightarrow 0.$$

From (5.1), we can derive that

$$\mathbf{u}^\tau \text{ is uniformly bounded in } L^1(0, T; B_{2,1}^\sigma)$$

$$\text{and } \tau(\mathbf{u}^\tau \cdot \nabla)\mathbf{u}^\tau \text{ is uniformly bounded in } L^1(0, T; B_{2,1}^{\sigma-1}).$$

Hence, in the system (1.4)–(1.5), the above convergence properties allow us to pass to the limit $\tau \rightarrow 0$ in the sense of distributions, which implies that $(\mathcal{N}, \mathcal{E})$ is a global weak solution to the drift-diffusion model (1.6). \square

Acknowledgments. The author thanks his adviser, Professor Daoyuan Fang, for his constant support and encouragement. He also thanks the anonymous referees for their useful comments and suggestions.

REFERENCES

- [1] G. ALÌ, D. BINI AND R. RIONERO, *Global existence and relaxation limit for smooth solutions to the Euler–Poisson model for semiconductors*, SIAM J. Math. Anal., 32 (2000), pp. 572–587.
- [2] J. Y. CHEMIN, *Perfect Incompressible Fluids*, Oxford Lecture Ser. Math. Appl. 14, Oxford University Press, New York, 1998 (in English).
- [3] G. Q. CHEN, J. JEROME, AND B. ZHANG, *Particle hydrodynamic models in biology and microelectronics: Singular relaxation limits*, Nonlinear Anal., 30 (1997), pp. 233–244.
- [4] J. F. COULOMBEL AND T. GOUDON, *The strong relaxation limit of the multidimensional isothermal Euler equations*, Trans. Amer. Math. Soc., 359 (2007), pp. 637–648.
- [5] D. Y. FANG AND J. XU, *Existence and asymptotic behavior of C^1 solutions to the multidimensional compressible Euler equations with damping*, Nonlinear Anal., 70 (2009), pp. 244–261.
- [6] D. Y. FANG, J. XU, AND T. ZHANG, *Global exponential stability of classical solutions to the hydrodynamic model for semiconductors*, Math. Models Methods Appl. Sci., 17 (2007), pp. 1507–1530.
- [7] L. HSIAO AND K. ZHANG, *The relaxation of the hydrodynamic model for semiconductors to the drift-diffusion equations*, J. Differential Equations, 165 (2000), pp. 315–354.
- [8] A. JÜNGEL AND Y. J. PENG, *A hierarchy of hydrodynamic models for plasmas: Zero-relaxation-time limits*, Comm. Partial Differential Equations, 24 (1999), pp. 1007–1033.
- [9] A. JÜNGEL AND Y. J. PENG, *Zero-relaxation-time limits in the hydrodynamic models for plasmas revisited*, Z. Angew. Math. Phys., 51 (2000), pp. 385–396.
- [10] S. JUNCA AND M. RASCLE, *Relaxation of the isothermal Euler–Poisson system to the drift-diffusion equations*, Quart. Appl. Math., 58 (2000), pp. 511–521.
- [11] C. LATTANZIO AND P. MARCATI, *The relaxation to the drift-diffusion system for the 3-D isentropic Euler–Poisson model for semiconductors*, Discrete Contin. Dyn. Syst., 5 (1999), pp. 449–455.
- [12] P. MARCATI AND R. NATALINI, *Weak solutions to a hydrodynamic model for semiconductors and relaxation to the drift-diffusion equations*, Arch. Ration. Mech. Anal., 129 (1995), pp. 129–145.
- [13] Y. SHIZUTA AND S. KAWASHIMA, *Systems of equations of hyperbolic-parabolic type with applications to the discrete Boltzmann equation*, Hokkaido Math. J., 14 (1985), pp. 249–275.
- [14] J. SIMON, *Compact sets in the space $L^p(0, T; B)$* , Ann. Mat. Pura Appl. (4), 146 (1987), pp. 65–96.
- [15] W. A. YONG, *Diffusive relaxation limit of multidimensional isentropic hydrodynamical models for semiconductors*, SIAM J. Appl. Math., 64 (2004), pp. 1737–1748.
- [16] W. A. YONG, *Entropy and global existence for hyperbolic balance laws*, Arch. Ration. Mech. Anal., 172 (2004), pp. 247–266.

EXISTENCE AND STABILITY RESULTS FOR PERIODIC STOKESIAN HELE–SHAW FLOWS*

JOACHIM ESCHER[†] AND BOGDAN-VASILE MATIOC[†]

Abstract. We consider here a 2π -periodic and two-dimensional Hele–Shaw flow modelling the motion of a viscous and incompressible fluid. The free surface is moving under the influence of gravity and is modelled by a modified Darcy law for Stokesian fluids. The bottom of the cell is assumed to be impermeable. We prove the existence of a unique classical solution if the initial data is near a constant, identify the equilibria of the flow, and study their stability.

Key words. quasi-linear elliptic equation, nonlinear parabolic equation, Hele–Shaw flow, non-Newtonian fluid, Oldroyd-B fluid, power law fluid

AMS subject classifications. 35J65, 35K55, 35R35, 42A45

DOI. 10.1137/070707671

1. Introduction.

1.1. The mathematical model. Using a fluid model with a shear-rate dependent viscosity, the authors in [12] derive a Darcy law coupling the pressure p and the velocity field v of the bulk fluid situated between the parallel plates of a Hele–Shaw cell:

$$(1) \quad v = -\frac{Dp}{\bar{\mu}(|Dp|^2)}.$$

Here $\bar{\mu}$ is the so-called *effective viscosity*. In the situation studied in [12] the cell is horizontal and the gap d between the plates is small in relation to the other two dimensions of the cell. Having averaged over the gap d they consider (1) for a two-dimensional fluid. From the conservation of mass one concludes that the pressure p of a finite patch Ω of fluid with boundary Γ must solve the Dirichlet problem

$$(2) \quad \operatorname{div} \left(\frac{Dp}{\bar{\mu}(|Dp|^2)} \right) = 0 \quad \text{in } \Omega, \quad p = -\sigma\kappa \quad \text{on } \Gamma,$$

where σ is the surface tension parameter and κ the curvature of Γ . The solvability of (2) is discussed in terms of the monotonicity of $\bar{\mu}$. It should be emphasized that in [12] only problems on temporal fixed domains are rigorously investigated. However, the full problem describing the more complex situation in which the boundary Γ moves and which is obtained by coupling (2) with an evolution equation describing the motion of the fluid molecules situated on Γ is not addressed by the authors in [12].

In this paper we study the general Cauchy problem of the full system describing the flow of a non-Newtonian fluid in a vertical Hele–Shaw cell with an impermeable bottom and laterals (see Figure 1). Hence, effects caused by gravity must be incorporated into our setting. As in [12] the gap d is again small compared to the height and length of the cell. In order to avoid the contact angle problem a two-dimensional,

*Received by the editors November 8, 2007; accepted for publication (in revised form) September 13, 2008; published electronically January 21, 2009.

<http://www.siam.org/journals/sima/40-5/70767.html>

[†]Institute of Applied Mathematics, Leibniz University of Hanover, Welfengarten 1, Hanover, Germany (escher@ifam.uni-hannover.de, matioc@ifam.uni-hannover.de).

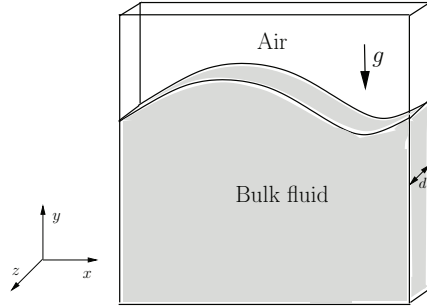


FIG. 1.

strip-like geometry will be considered. Given a positive function $f \in C^1(\mathbb{R})$, bounded away from 0, we define the set

$$\tilde{\Omega}_f := \{(x, y) \in \mathbb{R}^2 : 0 < y < f(x)\},$$

identified as the fluid domain, and denote the components of its boundary $\partial\tilde{\Omega}_f$ by

$$\tilde{\Gamma}_f := \{(x, f(x)) : x \in \mathbb{R}\}, \quad \tilde{\Gamma}_0 := \mathbb{R} \times \{0\}.$$

We denote by ν the outward-pointing normal of $\partial\tilde{\Omega}_f$. The motion of the fluid is modelled using the Darcy law (1):

$$v = -\frac{Du}{\bar{\mu}(|Du|^2)},$$

with p replaced by the *velocity potential* u defined by

$$u(x, y) = \frac{p(x, y)}{g \cdot \rho} + y, \quad (x, y) \in \tilde{\Omega}_f.$$

Here g is the gravity acceleration constant, ρ is the density constant, and $Du = (\partial_1 u, \partial_2 u)$ is the gradient of u . Also called *piezometric head*, u expresses the mechanical energy due to gravity and pressure of the fluid, per unit weight of fluid (see [3]).

The viscosity $\mu \in C^\infty([0, \infty), (0, \infty))$ is assumed to satisfy $\mu(r) + 2r\mu'(r) > 0$ for $r \geq 0$. Particularly, this implies that the mapping $[0, \infty) \ni r \mapsto h(r) := r\mu^2(r)$ is invertible. The effective viscosity $\bar{\mu}$ (see [12]) is defined by

$$\frac{1}{\bar{\mu}(r)} := c \int_{-1}^1 \frac{s^2}{\mu(rs^2)} ds,$$

where c is a positive constant and $\tilde{\mu} := \mu \circ h^{-1}$. The free surface $\tilde{\Gamma}_f$ separating the fluid from air, at pressure normalized to be zero, is moving under the influence of gravity. Surface tension effects are neglected. Consequently, $p = 0$ on $\tilde{\Gamma}_f$ and

$$u = f \quad \text{on} \quad \tilde{\Gamma}_f.$$

On the fixed boundary $\tilde{\Gamma}_0$ we have no flux, i.e.,

$$\partial_\nu u = 0 \quad \text{on} \quad \tilde{\Gamma}_0.$$

Furthermore, the interface $\tilde{\Gamma}_f$ is implicitly given by $F(t, z) = 0$, where $z = (x, y)$ and $F(t, z) = y - f(t, x)$. Assuming that a particle located on the interface remains there, we obtain by differentiating this equation with respect to t the relation

$$0 = \frac{d}{dt}F(t, z) = -f_t(t, x) + (-f_x, 1)z'.$$

Finally, we shall make the following periodicity requirement on f and u :

$$\begin{aligned} f(t, x + 2\pi) &= f(t, x) & \forall x \in \mathbb{R}, t \geq 0, \\ u(x + 2\pi, y) &= u(x, y) & \forall (x, y) \in \tilde{\Omega}_{f(t)}, t \geq 0. \end{aligned}$$

Summarizing, we arrive at the following system:

$$\begin{aligned} (3) \quad \operatorname{div} \left(\frac{Du}{\bar{\mu}(|Du|^2)} \right) &= 0 & \text{in } \Omega_{f(t)}, \quad t \geq 0, \\ \partial_\nu u &= 0 & \text{on } \Gamma_0, \quad t \geq 0, \\ u &= f & \text{on } \Gamma_{f(t)}, \quad t \geq 0, \\ \partial_t f(t, \cdot) + \frac{\sqrt{1 + \partial_x f^2(t, \cdot)}}{\bar{\mu}(|Du(\cdot, f(t, \cdot))|^2)} \partial_\nu u(\cdot, f(t, \cdot)) &= 0 & \text{on } \mathbb{S}^1, \quad t > 0, \\ f(0, \cdot) &= f_0 & \text{on } \mathbb{S}^1, \end{aligned}$$

where, given $t \geq 0$, we use the notation

$$\begin{aligned} \Omega_{f(t)} &:= \{(x, y) \in \mathbb{S}^1 \times \mathbb{R} : 0 < y < f(t, x)\}, \\ \Gamma_{f(t)} &:= \{(x, f(t, x)) : x \in \mathbb{S}^1\}, \quad \Gamma_0 = \mathbb{S}^1 \times \{0\}, \end{aligned}$$

and where \mathbb{S}^1 stands for the unit circle. For the sake of simplicity, we identify periodic functions on \mathbb{R} with functions on \mathbb{S}^1 , and periodic functions in the x variable in $\tilde{\Omega}_f$ with functions in Ω_f , for positive functions f on \mathbb{S}^1 , respectively. Given $k \in \mathbb{N}$ and $\alpha \in (0, 1)$, we define the so-called little Hölder spaces $h^{k+\alpha}(\mathbb{S}^1)$ as the closure of $C^\infty(\mathbb{S}^1)$ in $C^{k+\alpha}(\mathbb{S}^1)$.

For our analysis we fix $\alpha \in (0, 1)$ and define

$$\mathcal{U} := \left\{ f \in C^{2+\alpha}(\mathbb{S}^1) : \min_{x \in \mathbb{S}^1} f(x) > 0 \right\},$$

respectively, $\mathcal{V} := \mathcal{U} \cap h^{2+\alpha}(\mathbb{S}^1)$. For $f \in \mathcal{U}$ we denote by $buc^{k+\alpha}(\Omega_f)$ the closure of $BUC^\infty(\Omega_f)$ in the Hölder space $BUC^{k+\alpha}(\Omega_f)$.

A pair (u, f) is called a *classical Hölder solution* of (3) on $[0, T]$, with $T > 0$, if

$$\begin{aligned} f &\in C([0, T], \mathcal{V}) \cap C^1([0, T], h^{1+\alpha}(\mathbb{S}^1)), \\ u(\cdot, t) &\in buc^{2+\alpha}(\Omega_{f(t)}), \quad t \in [0, T], \end{aligned}$$

and if (u, f) satisfies the equations in (3) pointwise.

For simplicity we denote

$$\mathcal{Q}u := \operatorname{div} \left(\frac{Du}{\bar{\mu}(|Du|^2)} \right).$$

Then $\mathcal{Q}u = a_{ij}(Du)u_{ij}$, with

$$a_{ij}(q_1, q_2) = \frac{\delta_{ij}}{\bar{\mu}(|q|^2)} - \frac{2q_i q_j \bar{\mu}'(|q|^2)}{\bar{\mu}^2(|q|^2)}, \quad 1 \leq i, j \leq 2, q = (q_1, q_2) \in \mathbb{R}^2.$$

The eigenvalues of the matrix $[a_{ij}(q)]_{1 \leq i, j \leq 2}$, $q \in \mathbb{R}^2$, are

$$\lambda_1(q) = \frac{1}{\bar{\mu}(|q|^2)}, \quad \lambda_2(q) = \frac{1}{\bar{\mu}(|q|^2)} - \frac{2|q|^2 \bar{\mu}'(|q|^2)}{\bar{\mu}^2(|q|^2)},$$

and we have

$$c|\xi|^2 \leq a_{ij}(q_1, q_2)\xi_i \xi_j \leq C|\xi|^2 \quad \forall \xi = (\xi_1, \xi_2) \in \mathbb{R}^2, (q_1, q_2) \in \mathbb{R}^2,$$

if we assume that there exist positive constants c and C such that

$$(A_1) \quad c \leq \frac{1}{\bar{\mu}(r)} \leq C \quad \forall r \geq 0,$$

$$(A_2) \quad c \leq \frac{1}{\bar{\mu}(r)} - \frac{2r\bar{\mu}'(r)}{\bar{\mu}^2(r)} \leq C \quad \forall r \geq 0.$$

It can be shown (cf. [6]) that (A₁) and (A₂) hold true, provided the viscosity μ satisfies

$$(V_1) \quad m \leq \mu(r) \leq M,$$

$$(V_2) \quad m \leq \mu(r) + 2r\mu'(r) \leq M,$$

for all $r \geq 0$, where m and M are positive constants. If the viscosity μ is constant, then the fluid is called Newtonian. Results regarding this situation can be found in [4] and [7], [8], [9], [10]. The class of Stokesian, sometimes also called non-Newtonian, fluids satisfying (V₁) and (V₂) is quite large; see [6]. This class includes particularly numerous Oldroyd-B and power law fluids.

Under the above assumptions \mathcal{Q} is a uniformly elliptic operator in \mathbb{R}^2 . Consider now the problem

$$\begin{aligned} \operatorname{div} \left(\frac{Du}{\bar{\mu}(|Du|^2)} \right) &= 0 \quad \text{in } G_{f(t)}, \quad t \geq 0, \\ u &= f \quad \text{on } \partial G_{f(t)}, \quad t \geq 0, \\ \partial_t f(t, \cdot) + \frac{\sqrt{1 + \partial_x f^2(t, \cdot)}}{\bar{\mu}(|Du(\cdot, f(t, \cdot))|^2)} \partial_\nu u(\cdot, f(t, \cdot)) &= 0 \quad \text{on } \mathbb{S}^1, \quad t > 0, \\ f(0, \cdot) &= f_0 \quad \text{on } \mathbb{S}^1, \end{aligned} \tag{4}$$

where

$$G_{f(t)} := \{(x, y) \in \mathbb{S}^1 \times \mathbb{R} : -f(t, x) < y < f(t, x)\},$$

for $t \geq 0$. A pair (u, f) is called a *classical Hölder solution* of (4) on $[0, T]$ if

$$f \in C([0, T], \mathcal{V}) \cap C^1([0, T], h^{1+\alpha}(\mathbb{S}^1)),$$

$$u(\cdot, t) \in buc^{2+\alpha}(G_{f(t)}), \quad t \in [0, T],$$

and if (u, f) satisfies the equations in (4) pointwise.

Given $f \in \mathcal{V}$, let $u \in buc^{2+\alpha}(G_f)$ be the unique solution (existence will be discussed later) of the quasi-linear Dirichlet problem

$$(5) \quad \begin{aligned} \mathcal{Q}u &= 0 \quad \text{in } G_f, \\ u &= f \quad \text{on } \partial G_f. \end{aligned}$$

Because of the symmetry of the domain G_f and of the boundary conditions we obtain

$$u(x, -y) = u(x, y)$$

for all $(x, y) \in G_f$. Consequently, $\partial_\nu u = 0$ on Γ_0 ; thus the restriction of u to Ω_f is the unique solution of

$$(6) \quad \begin{aligned} \mathcal{Q}u &= 0 \quad \text{in } \Omega_f, \\ \partial_\nu u &= 0 \quad \text{on } \Gamma_0, \\ u &= f \quad \text{on } \Gamma_f. \end{aligned}$$

We deduce that there exists a one-to-one correspondence between solutions to (3) and (4). Namely, if (u, f) is a solution to (3) on the interval $[0, T]$, $T > 0$, then for each $t \in [0, T]$ the velocity potential $u(\cdot, t) \in buc^{2+\alpha}(\Omega_{f(t)})$ can be extended by reflection on the entire domain $G_{f(t)}$. We obtain in this manner a solution to (4). Moreover, if (\tilde{u}, f) is a solution to (4) on the interval $[0, T]$, $T > 0$, then we can restrict for $t \in [0, T]$ the function $\tilde{u}(\cdot, t)$ to $\Omega_{f(t)}$, and so obtain a solution to (3).

It is more convenient to treat (4), because we have in this case a Dirichlet problem for the velocity potential, which can be solved using the techniques presented in [11].

We state now the first main result of this work.

THEOREM 1.1. *Assume that (A₁) and (A₂) hold true.*

(a) *Let c be a positive constant. Then we find an open neighborhood \mathcal{O} of c in \mathcal{V} such that for each $f_0 \in \mathcal{O}$, problem (3) has a classical Hölder solution (u, f) on an interval $[0, T]$ with $T > 0$. Moreover, there exists a constant $\gamma \in (0, 1)$ such that $f \in C_\gamma^\gamma((0, T], h^{2+\alpha}(\mathbb{S}^1))$.*

(b) *Let (u_1, f_1) and (u_2, f_2) be solutions of (3) with $f_1 \in C_\gamma^\gamma((0, T], h^{2+\alpha}(\mathbb{S}^1))$, $\gamma \in (0, 1)$, and $f_2 \in C_\delta^\delta((0, T], h^{2+\alpha}(\mathbb{S}^1))$, $\delta \in (0, 1)$. If $f_1([0, T]) \subset \mathcal{O}$ and $f_2([0, T]) \subset \mathcal{O}$, then $(u_1, f_1) = (u_2, f_2)$.*

For a definition of the weighted Hölder spaces $C_\gamma^\gamma((0, T], h^{2+\alpha}(\mathbb{S}^1))$, $\gamma \in (0, 1)$, see [13].

In section 3 we prove that each pair $(c, c) \in buc^{2+\alpha}(\Omega_c) \times \mathcal{V}$, with $c > 0$, is a stationary solution of the flow (3). Furthermore, we prove that if the initial data $f_0 \in \mathcal{O}$ is close enough to c in $h^{2+\alpha}(\mathbb{S}^1)$ and if Ω_{f_0} and Ω_c enclose the same fluid volume, then the solution corresponding to f_0 exists globally in time and is attracted at an exponential rate by c (see Theorem 3.3).

1.2. The transformed problem. In order to solve the problem we transform it on a fixed reference manifold $G := \mathbb{R} \times (0, 2)$. In order to do so we define for each $f \in \mathcal{U}$ a diffeomorphism $\phi_f \in Diff^{2+\alpha}(G, G_f)$ by

$$\phi_f(x, y) := (x, (1 - y)f(x)), \quad (x, y) \in G,$$

and the push-forward and pull-back operators induced by ϕ_f :

$$\begin{aligned} \phi_f^* : BUC(G_f) &\rightarrow BUC(G), \quad u \mapsto u \circ \phi_f, \\ \phi_f^\# : BUC(G) &\rightarrow BUC(G_f), \quad v \mapsto v \circ \phi_f^{-1}. \end{aligned}$$

The transformed operators $\mathcal{A}(f)$ and \mathcal{B} , acting on $BUC^2(G)$, respectively, $\mathcal{U} \times BUC^{2+\alpha}(G)$, are defined by

$$\begin{aligned} \mathcal{A}(f) &:= \phi_f^* \circ \mathcal{Q} \circ \phi_*^f, \\ \mathcal{B}(f, v)(x) &:= \frac{D(\phi_*^f v)}{\bar{\mu}(|D(\phi_*^f v)|^2)}(x, f(x)) \cdot n(x), \quad x \in \mathbb{S}^1, \end{aligned}$$

with $n(x) := (-f'(x), 1)$, $x \in \mathbb{S}^1$. Transformation of (4) to G yields

$$(7) \quad \begin{aligned} \mathcal{A}(f)v &= 0 && \text{in } G \times [0, \infty), \\ v &= f && \text{on } \partial G \times [0, \infty), \\ \partial_t f + \mathcal{B}(f, v) &= 0 && \text{on } \mathbb{S}^1 \times (0, \infty), \\ f(0) &= f_0, \end{aligned}$$

where $v := \phi_f^* u$. A pair (v, f) is called a *classical Hölder solution* of (7) on $[0, T]$, with $T > 0$, if

$$\begin{aligned} f &\in C([0, T], \mathcal{V}) \cap C^1([0, T], h^{1+\alpha}(\mathbb{S}^1)), \\ v(\cdot, t) &\in buc^{2+\alpha}(G), \quad t \in [0, T], \end{aligned}$$

and if (v, f) satisfies the equations in (7) pointwise. It is obvious that problems (4) and (7) are equivalent in the following sense.

LEMMA 1.2. *Let $f_0 \in \mathcal{V}$ be given.*

(a) *If (u, f) is a classical Hölder solution of (4), then $(\phi_f^* u, f)$ is a classical Hölder solution of (7).*

(b) *If (v, f) is a classical Hölder solution of (7), then $(\phi_*^f v, f)$ is a classical Hölder solution of (4).*

Proof. See, for example, [6]. □

LEMMA 1.3. *Given $f \in \mathcal{U}$, we have*

$$\mathcal{A}(f)v = b_{ij}(y, f, Dv)v_{ij} + b(y, f, Dv)v_2 \quad \text{for } v \in BUC^2(G),$$

where, using the notation $D_f v := (v_1 + \frac{(1-y)f'}{f}v_2, -\frac{1}{f}v_2)$ for $f \in \mathcal{U}$, $v \in BUC^2(G)$, and $y \in [0, 2]$, we have

$$\begin{aligned} b_{11}(y, f, Dv) &= a_{11}(D_f v), \\ b_{12}(y, f, Dv) &= b_{21}(y, f, Dv) = \frac{(1-y)f'}{f}a_{11}(D_f v) - \frac{1}{f}a_{12}(D_f v), \\ b_{22}(y, f, Dv) &= \frac{(1-y)^2 f'^2}{f^2}a_{11}(D_f v) - \frac{2(1-y)f'}{f^2}a_{12}(D_f v) + \frac{1}{f^2}a_{22}(D_f v), \\ b(y, f, Dv) &= (1-y) \left(\frac{f''}{f} - \frac{2f'^2}{f^2} \right) a_{11}(D_f v) + \frac{2f'}{f^2}a_{12}(D_f v). \end{aligned}$$

Proof. This follows by direct computation. □

We have the following existence and uniqueness theorem. The proof is similar to that of Lemma 2.2 in [6].

THEOREM 1.4. *Let $f \in \mathcal{V}$ be given. There exists a unique solution $v \in buc^{2+\alpha}(G)$ of the quasi-linear Dirichlet problem*

$$(8) \quad \begin{aligned} \mathcal{A}(f)v &= 0 \quad \text{in } G, \\ v &= f \quad \text{on } \partial G. \end{aligned}$$

Given $f \in \mathcal{V}$, we denote by $\mathcal{T}(f) \in buc^{2+\alpha}(G)$ the solution to (8). The mapping $[\mathcal{V} \ni f \mapsto \mathcal{T}(f) \in buc^{2+\alpha}(G)]$ is smooth.

2. The nonlinear Cauchy problem. Due to Theorem 1.4 we can reduce problem (4) to an abstract, fully nonlinear Cauchy problem on \mathbb{S}^1 :

$$(9) \quad \partial_t f + \Phi(f) = 0, \quad f(0) = f_0,$$

where

$$\Phi(f) := \mathcal{B}(f, \mathcal{T}(f)).$$

Note that the mapping $f \mapsto \Phi(f)$ is a nonlinear and nonlocal operator of first order. It carries also a fully nonlinear structure in the sense that there is no leading linear part in \mathcal{Q} . We solve the Cauchy problem (9) by applying the theory of maximal regularity; cf. [13]. Further on we verify that the assumption in Theorem 8.1.1 in [13] are satisfied.

We first restrict our attention to the operator \mathcal{B} . Let $(f, v) \in \mathcal{V} \times buc^{2+\alpha}(G)$ be given. The function $\mathcal{B}(f, v)$ is given by

$$\mathcal{B}(f, v) = -\frac{1}{\bar{\mu}(|\gamma_0 D_f v|^2)} \left(f' \gamma_0 v_1 + \frac{1}{f} (1 + f'^2) \gamma_0 v_2 \right),$$

with γ_0 denoting the trace operator on Γ_0 . Together with the relation

$$|\gamma_0 D_f v|^2 = \gamma_0 v_1^2 + 2 \frac{f'}{f} \gamma_0 v_1 v_2 + \frac{1 + f'^2}{f^2} \gamma_0 v_2^2$$

we obtain that the operator \mathcal{B} defined above is smooth.

LEMMA 2.1. *The mapping $\mathcal{B} : \mathcal{V} \times buc^{2+\alpha}(G) \rightarrow h^{1+\alpha}(\mathbb{S}^1)$ is smooth. Given $(f, v) \in \mathcal{V} \times buc^{2+\alpha}(G)$, we have*

$$\begin{aligned} \partial \mathcal{B}(f, v)[h, u] &= -\frac{1}{\bar{\mu}} (|\gamma_0 D_f v|^2) \left[f' \gamma_0 u_1 + h' \gamma_0 v_1 + \frac{1}{f} (1 + f'^2) \gamma_0 u_2 \right. \\ &\quad \left. - \left(\frac{h}{f^2} - \frac{2f'h'}{f} + \frac{hf'^2}{f^2} \right) \gamma_0 v_2 \right] \\ &\quad - 2 \left(\frac{1}{\bar{\mu}} \right)' (|\gamma_0 D_f v|^2) \left(f' \gamma_0 v_1 + \frac{1}{f} (1 + f'^2) \gamma_0 v_2 \right) \left[\gamma_0 v_1 u_1 \right. \\ &\quad \left. + \frac{h'}{f} \gamma_0 v_1 v_2 + \frac{f'}{f} \gamma_0 u_1 v_2 + \frac{f'}{f} \gamma_0 v_1 u_2 - \frac{f'h}{f^2} \gamma_0 v_1 v_2 + \frac{f'h'}{f^2} \gamma_0 v_2^2 \right. \\ &\quad \left. + \frac{f'^2}{f^2} \gamma_0 v_2 u_2 - \frac{hf'^2}{f^3} \gamma_0 v_2^2 + \frac{1}{f^2} \gamma_0 v_2 u_2 - \frac{h}{f^3} \gamma_0 v_2^2 \right] \end{aligned}$$

for all $[h, u] \in h^{2+\alpha}(\mathbb{S}^1) \times buc^{2+\alpha}(G)$.

LEMMA 2.2. *Given $f \in \mathcal{V}$ and $h \in h^{2+\alpha}(\mathbb{S}^1)$, the function $\partial\mathcal{T}(f)[h]$ is the unique solution of the linear Dirichlet problem*

$$\begin{aligned}
 & b_{ij}w_{ij} + bw_2 + D_f w \left[u_{11}\partial a_{11}(D_f u) + 2u_{12} \left(\frac{(1-y)f'}{f} \partial a_{11}(D_f u) - \frac{1}{f} \partial a_{12}(D_f u) \right) \right. \\
 & + u_{22} \left(\frac{(1-y)^2 f'^2}{f^2} \partial a_{11}(D_f u) - 2 \frac{(1-y)f'}{f^2} \partial a_{12}(D_f u) + \frac{1}{f^2} \partial a_{22}(D_f u) \right) \\
 & \left. + u_2 \left((1-y) \left(\frac{f''}{f} - 2 \frac{f'^2}{f^2} \right) \partial a_{11}(D_f u) + 2 \frac{f'}{f^2} \partial a_{12}(D_f u) \right) \right] \\
 & = -u_2 \left((1-y) \frac{fh' - f'h}{f^2}, \frac{h}{f^2} \right) \cdot \left[u_{11}\partial a_{11}(D_f u) \right. \\
 & + 2u_{12} \left(\frac{(1-y)f'}{f} \partial a_{11}(D_f u) - \frac{1}{f} \partial a_{12}(D_f u) \right) + u_{22} \left(\frac{(1-y)^2 f'^2}{f^2} \partial a_{11}(D_f u) \right. \\
 & \left. - 2 \frac{(1-y)f'}{f^2} \partial a_{12}(D_f u) + \frac{1}{f^2} \partial a_{22}(D_f u) \right) + u_2 \left((1-y) \left(\frac{f''}{f} - 2 \frac{f'^2}{f^2} \right) \partial a_{11}(D_f u) \right. \\
 & \left. + 2 \frac{f'}{f^2} \partial a_{12}(D_f u) \right) \left. \right] - 2u_{12} \left((1-y) \frac{fh' - f'h}{f^2} a_{11}(D_f u) + \frac{h}{f^2} a_{12}(D_f u) \right) \\
 & - 2u_{22} \left(\frac{(1-y)^2 (ff'h' - f'^2 h)}{f^3} a_{11}(D_f u) - (1-y) \frac{fh' - 2f'h}{f^3} a_{12}(D_f u) \right. \\
 & \left. - \frac{h}{f^3} a_{22}(D_f u) \right) - u_2 \left((1-y) \left(\frac{fh'' - f''h}{f^2} - 4 \frac{ff'h' - f'^2 h}{f^3} \right) a_{11}(D_f u) \right. \\
 & \left. + 2 \frac{fh' - 2f'h}{f^3} a_{12}(D_f u) \right) \quad \text{in } G, \\
 & w = h \quad \text{on } \partial G,
 \end{aligned}$$

where $u := \mathcal{T}(f)$ and $b_{ij} = b_{ij}(y, f, Du)$, $b = b(y, f, Du)$ are the coefficients computed in Lemma 1.3.

Proof. The proof follows using standard arguments and we omit it. □

Consequently, we obtain from Lemma 2.1 and Theorem 1.4 that the operator Φ is smooth. We can use the chain rule to compute its derivative. More precisely, we have the following result.

THEOREM 2.3. *The map Φ belongs to $C^\infty(\mathcal{V}, h^{1+\alpha}(\mathbb{S}^1))$ and*

$$\partial\Phi(f) = \partial\mathcal{B}(f, \mathcal{T}(f)) \circ (\text{id}_{h^{2+\alpha}(\mathbb{S}^1)}, \partial\mathcal{T}(f)) \quad \forall f \in \mathcal{V}.$$

Our goal is to show that for $c \in \mathbb{R}_{>0}$ the Fréchet derivative $-\partial\Phi(c)$ generates a strongly continuous analytic semigroup in $\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1))$. Then, using general results from the theory of maximal regularity (cf. [13]), we can prove Theorem 1.1. It is clear that we consider here $-\partial\Phi(c)$ as an unbounded operator in $h^{1+\alpha}(\mathbb{S}^1)$ with domain $h^{2+\alpha}(\mathbb{S}^1)$.

Let $c \in \mathbb{R}_{>0}$ be given. The solution $\mathcal{T}(c)$ to (8) is the constant function c on G . Given $h \in h^{2+\alpha}(\mathbb{S}^1)$, the map $\partial\mathcal{T}(c)[h]$ is the solution of the Dirichlet problem

$$(10) \quad \begin{aligned} c^2 w_{11} + w_{22} &= 0 && \text{in } G, \\ w &= h && \text{on } \partial G. \end{aligned}$$

Consequently,

$$-\partial\Phi(c)[h](x) = \frac{\zeta}{c} w_2(x, 0), \quad x \in \mathbb{S}^1,$$

where $\zeta := 1/\overline{\mu}(0)$.

We expand h and w in the following way:

$$h(x) = \sum_{k \in \mathbb{Z}} c_k e^{ikx}, \quad w(x, y) = \sum_{k \in \mathbb{Z}} C_k(y) e^{ikx},$$

and we substitute these expressions into (10). Comparing the coefficients of e^{ikx} , $k \in \mathbb{Z}$, we get the following equations for C_k :

$$(11) \quad \begin{aligned} C_k'' - c^2 k^2 C_k &= 0, && 0 < y < 2, \\ C_k(0) &= c_k, \\ C_k(2) &= c_k, \end{aligned}$$

for $k \in \mathbb{Z} \setminus \{0\}$, and

$$(12) \quad \begin{aligned} C_0'' &= 0, && 0 < y < 2, \\ C_0(0) &= c_0, \\ C_0(2) &= c_0. \end{aligned}$$

The solution to (12) is the constant function $C_0 = c_0$, and for $k \in \mathbb{Z} \setminus \{0\}$ we have

$$C_k(y) = c_k d_k(y), \quad 0 \leq y \leq 2,$$

where

$$d_k(y) := \frac{e^{2ck} - 1}{e^{4ck} - 1} e^{cky} + \frac{e^{4ck} - e^{2ck}}{e^{4ck} - 1} e^{-cky}, \quad 0 \leq y \leq 2.$$

Summarizing, we obtain

$$(13) \quad -\partial\Phi(c) \left[\sum_{k \in \mathbb{Z}} c_k e^{ikx} \right] = \sum_{k \in \mathbb{Z}} \lambda_k c_k e^{ikx},$$

for all $h = \sum_{k \in \mathbb{Z}} c_k e^{ikx} \in h^{2+\alpha}(\mathbb{S}^1)$, with

$$(14) \quad \lambda_k := -\zeta \frac{e^{2ck} - 1}{e^{2ck} + 1} k, \quad k \in \mathbb{Z}.$$

The multiplying coefficients λ_k satisfy $\lambda_k = \lambda_{-k}$ for all $k \in \mathbb{Z}$. Furthermore,

$$0 = \lambda_0 > \lambda_{\pm 1} > \lambda_{\pm 2} > \cdots > \lambda_{\pm k} > \lambda_{\pm(k+1)} > \cdots.$$

For the purpose of proving that $-\partial\Phi(c)$ generates a strongly continuous analytic semigroup in $\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1))$, i.e., $\partial\Phi(c) \in \mathcal{H}(h^{2+\alpha}(\mathbb{S}^1), h^{1+\alpha}(\mathbb{S}^1))$, it is enough (cf. [1]) to find constants $\omega > 0$ and $\kappa \geq 1$ such that

$$(15) \quad \lambda + \partial\Phi(c) \in \mathcal{L}is(h^{2+\alpha}(\mathbb{S}^1), h^{1+\alpha}(\mathbb{S}^1)),$$

$$(16) \quad |\lambda| \cdot \|R(\lambda, -\partial\Phi(c))\|_{\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1))} \leq \kappa$$

for all $\text{Re } \lambda \geq \omega$.

For $r \geq 0$ we introduce the Sobolev space

$$H^r(\mathbb{S}^1) := \left\{ f \in L^2(\mathbb{S}^1) : \sum_{n \in \mathbb{Z}} (1+n^2)^r |\widehat{f}(n)|^2 < \infty \right\},$$

endowed with the scalar product

$$\langle f, g \rangle := \sum_{n \in \mathbb{Z}} (1+n^2)^r \widehat{f}(n) \overline{\widehat{g}(n)}.$$

The smooth functions are dense in $H^r(\mathbb{S}^1)$ and the Sobolev embedding $H^{k+r}(\mathbb{S}^1) \hookrightarrow C^k(\mathbb{S}^1)$ holds for all $k \in \mathbb{N}$, provided $r > 1/2$. Therefore we have

$$H^{k+s}(\mathbb{S}^1) \xrightarrow{d} h^{k+\beta}(\mathbb{S}^1)$$

for all $k \in \mathbb{N}$, $\beta \in [0, 1]$, and $s > 3/2$.

We choose $\omega = 1$ and prove that relations (15) and (16) hold.

LEMMA 2.4. *Given $r \geq 0$ and $\text{Re } \lambda \geq \omega$, we have*

$$\lambda + \partial\Phi(c) \in \mathcal{L}is(H^{r+1}(\mathbb{S}^1), H^r(\mathbb{S}^1)).$$

Proof. The proof is just a consequence of the fact that

$$\lim_{k \rightarrow \infty} \frac{\lambda_k}{k} = -\zeta.$$

Given $r \geq 0$ and $\text{Re } \lambda \geq \omega$, the inverse of $\lambda + \partial\Phi(c)$ is the operator $R(\lambda, -\partial\Phi(c)) \in \mathcal{L}is(H^r(\mathbb{S}^1), H^{r+1}(\mathbb{S}^1))$ defined by

$$(17) \quad R(\lambda, -\partial\Phi(c)) \left[\sum_{k \in \mathbb{Z}} c_k e^{ikx} \right] = \sum_{k \in \mathbb{Z}} \frac{1}{\lambda - \lambda_k} c_k e^{ikx}$$

for all $h := \sum_{k \in \mathbb{Z}} c_k e^{ikx} \in H^r(\mathbb{S}^1)$. \square

In view of the above-mentioned facts, the multiplying operator $R(\lambda, -\partial\Phi(c))$ belongs to $\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1), h^{k+\alpha}(\mathbb{S}^1))$, $k \in \{1, 2\}$, if it belongs to $\mathcal{L}(C^{1+\alpha}(\mathbb{S}^1), C^{k+\alpha}(\mathbb{S}^1))$, $k \in \{1, 2\}$. In the next proposition we show that, with our choice of ω , relation (15) holds; i.e., for $\text{Re } \lambda \geq \omega$ it follows that λ belongs to the resolvent set $\rho(-\partial\Phi(c))$ of the operator $-\partial\Phi(c)$.

PROPOSITION 2.5.

$$\{\lambda \in \mathbb{C} : \text{Re } \lambda \geq \omega\} \subset \rho(-\partial\Phi(c)).$$

Proof. It is sufficient to prove that $R(\lambda, -\partial\Phi(c)) \in \mathcal{L}(C^{1+\alpha}(\mathbb{S}^1), C^{2+\alpha}(\mathbb{S}^1))$ for all $\text{Re } \lambda \geq \omega$. Let $\text{Re } \lambda \geq \omega$. Then $R(\lambda, -\partial\Phi(c)) \in \mathcal{L}(C^{1+\alpha}(\mathbb{S}^1), C^{2+\alpha}(\mathbb{S}^1))$, if we have

- (i) $\sup_{k \in \mathbb{Z}} |k| |M_k| < \infty$,
- (ii) $\sup_{k \in \mathbb{Z}} |k|^2 |M_{k+1} - M_k| < \infty$,
- (iii) $\sup_{k \in \mathbb{Z}} |k|^3 |M_{k+2} - 2M_{k+1} + M_k| < \infty$,

where $M_k := 1/(\lambda - \lambda_k)$ (see, e.g., [2] and [14]). The relation

$$\lim_{k \rightarrow \infty} \frac{k}{\lambda - \lambda_k} = \frac{1}{\zeta}$$

implies (i). For $k \geq 1$ we have

$$k^2 |M_{k+1} - M_k| = \frac{k}{\lambda - \lambda_{k+1}} \frac{k}{\lambda - \lambda_k} |\lambda_{k+1} - \lambda_k| \xrightarrow{k \rightarrow \infty} \frac{1}{\zeta},$$

because of $\lambda_k - \lambda_{k+1} \rightarrow \zeta$. Also

$$\begin{aligned} k^3 |M_{k+2} - 2M_{k+1} + M_k| &= \frac{k}{\lambda - \lambda_{k+2}} \frac{k}{\lambda - \lambda_{k+1}} \frac{k}{\lambda - \lambda_k} |\lambda(\lambda_{k+2} - 2\lambda_{k+1} + \lambda_k) \\ &\quad + \lambda_k(\lambda_{k+1} - \lambda_{k+2}) + \lambda_{k+2}(\lambda_{k+1} - \lambda_k)|, \end{aligned}$$

and $\lambda_{k+2} - 2\lambda_{k+1} + \lambda_k \rightarrow 0$, respectively, $\lambda_k(\lambda_{k+1} - \lambda_{k+2}) + \lambda_{k+2}(\lambda_{k+1} - \lambda_k) \rightarrow 2\zeta^2$. This completes the proof. \square

PROPOSITION 2.6. *There exists $\kappa \geq 1$ such that*

$$|\lambda| \cdot \|R(\lambda, -\partial\Phi(c))\|_{\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1))} \leq \kappa$$

for all $\text{Re } \lambda \geq \omega$.

Proof. Given $h = \sum_{k \in \mathbb{Z}} \widehat{h}(k) e^{ikx} \in h^{1+\alpha}(\mathbb{S}^1)$, we have

$$|\lambda| R(\lambda, -\partial\Phi(c)) \left[\sum_{k \in \mathbb{Z}} \widehat{h}(k) e^{ikx} \right] = \sum_{k \in \mathbb{Z}} M_k^\lambda \widehat{h}(k) e^{ikx},$$

where

$$M_k^\lambda = \frac{|\lambda|}{\lambda - \lambda_k} \quad \forall k \in \mathbb{Z}, \text{Re } \lambda \geq \omega.$$

It suffices to prove (see [2], [14]) that there exist positive constants s_1, s_2 , and s_3 such that

- (i) $\sup_{k \in \mathbb{Z}} |M_k^\lambda| \leq s_1$,
- (ii) $\sup_{k \in \mathbb{Z}} |k| |M_{k+1}^\lambda - M_k^\lambda| \leq s_2$,
- (iii) $\sup_{k \in \mathbb{Z}} |k|^2 |M_{k+2}^\lambda - 2M_{k+1}^\lambda + M_k^\lambda| \leq s_3$

for all $\operatorname{Re} \lambda \geq \omega$. For $\operatorname{Re} \lambda \geq \omega$ and $k \in \mathbb{Z}$ we have $\operatorname{Re}(\lambda - \lambda_k) \geq 1$ and $|\lambda - \lambda_k| \geq |\lambda|$. Thus, (i) holds with $s_1 = 1$. Moreover, for $k \neq 0$, we have $|\lambda - \lambda_k| \geq |\lambda_k|$ for all $\operatorname{Re} \lambda \geq \omega$; hence

$$|k| |M_{k+1}^\lambda - M_k^\lambda| = \frac{|\lambda|}{|\lambda - \lambda_{k+1}|} \frac{|k|}{|\lambda - \lambda_k|} |\lambda_{k+1} - \lambda_k| \leq \frac{|k|}{|\lambda_k|} |\lambda_{k+1} - \lambda_k|,$$

and (ii) holds. For $k \geq 1$ we further have

$$\begin{aligned} k^2 |M_{k+2} - 2M_{k+1} + M_k| &= \frac{|\lambda|}{|\lambda - \lambda_{k+2}|} \frac{k}{|\lambda - \lambda_{k+1}|} \frac{k}{|\lambda - \lambda_k|} |\lambda(\lambda_{k+2} - 2\lambda_{k+1} + \lambda_k) \\ &\quad + \lambda_k(\lambda_{k+1} - \lambda_{k+2}) + \lambda_{k+2}(\lambda_{k+1} - \lambda_k)| \\ &\leq \frac{k}{|\lambda_k|} |k(\lambda_{k+2} - 2\lambda_{k+1} + \lambda_k)| \\ &\quad + \frac{k}{|\lambda_k|} \frac{k}{|\lambda_{k+1}|} |\lambda_k(\lambda_{k+1} - \lambda_{k+2}) + \lambda_{k+2}(\lambda_{k+1} - \lambda_k)|. \end{aligned}$$

Taking into account

$$k(\lambda_{k+2} - 2\lambda_{k+1} + \lambda_k) \xrightarrow[k \rightarrow \infty]{} 0,$$

and using the symmetry of the coefficients λ_k , we have (iii). \square

In conclusion, given $c \in \mathbb{R}_{>0}$, we have proved that $-\partial\Phi(c)$ generates a strongly continuous analytic semigroup in $\mathcal{L}(h^{1+\alpha}(\mathbb{S}^1))$. The proof of Theorem 1.1 is similar to that of Theorem 8.1.1 in [13], and the assumptions of this theorem are all satisfied (see also Theorem 2.3).

3. Stability of the equilibria. We want to identify the steady states of this moving boundary problem, i.e., solutions which do not depend on time, and to study their stability. A pair $(u, f) \in buc^{2+\alpha}(\Omega_f) \times \mathcal{V}$ is a stationary solution of the flow (3) iff it is a solution of the free boundary problem

$$\begin{aligned} \operatorname{div} \left(\frac{Du}{\mu(|Du|^2)} \right) &= 0 \quad \text{in } \Omega_f, \\ \partial_\nu u &= 0 \quad \text{on } \Gamma_0, \\ u &= f \quad \text{on } \Gamma_f, \\ \partial_\nu u &= 0 \quad \text{on } \Gamma_f. \end{aligned} \tag{18}$$

If $f \in \mathcal{V}$ is known, then $u \in buc^{2+\alpha}(\Omega_f)$ is a solution of the Neumann problem

$$\begin{aligned} \operatorname{div} \left(\frac{Du}{\mu(|Du|^2)} \right) &= 0 \quad \text{in } \Omega_f, \\ \partial_\nu u &= 0 \quad \text{on } \partial\Omega_f. \end{aligned} \tag{19}$$

Thus there exists a constant $c \in \mathbb{R}$ such that $u = c$. The condition $u = f$ on Γ_f , together with $f \in \mathcal{V}$, implies $c \in \mathbb{R}_{>0}$. In conclusion, $(c, c) \in buc^{2+\alpha}(\Omega_c) \times \mathcal{V}$, $c > 0$, are the stationary solutions of the flow (3), or, equivalently, $(c, c) \in buc^{2+\alpha}(G) \times \mathcal{V}$, $c > 0$, are the stationary solutions of the flow (7).

Let $c \in \mathbb{R}_{>0}$ be given. We want to identify the spectrum of the operator $-\partial\Phi(c)$. Of course, we have

$$\{\lambda_k : k \in \mathbb{N}\} \subset \sigma_p(-\partial\Phi(c)) \subset \sigma(-\partial\Phi(c)).$$

Using a similar argument as in Proposition 2.5, one can in fact show that

$$\{\lambda_k : k \in \mathbb{N}\} = \sigma_p(-\partial\Phi(c)) = \sigma(-\partial\Phi(c)).$$

Thus, the spectrum of $-\partial\Phi(c)$ contains just eigenvalues. The largest eigenvalue is $\lambda_0 = 0$. This eigenvalue is simple and has a one-dimensional eigenspace consisting of constant functions. All other eigenvalues $\lambda_k, k \geq 1$, have a two-dimensional eigenspace spanned by $\{e^{ikx}, e^{-ikx}\}$. Compared to [5], the analysis is more involved here, since we have to take the eigenvalue $\lambda_0 = 0$ into consideration.

In order to prove the stability we transfer problem (9) into a neighborhood of the origin in $h^{2+\alpha}(\mathbb{S}^1)$. Let $\mathcal{V}_c := \mathcal{V} - c$. Then \mathcal{V}_c is an open neighborhood of the origin in $h^{2+\alpha}(\mathbb{S}^1)$. We define $\psi : \mathcal{V}_c \rightarrow h^{1+\alpha}(\mathbb{S}^1)$ by $\psi(f) := -\Phi(f + c)$ and consider the following Cauchy problem:

$$(20) \quad \partial_t f = \psi(f), \quad f(0) = f_0.$$

We have $\psi(0) = -\Phi(c) = 0$ and $\partial\psi(0) = -\partial\Phi(c)$. Thus, 0 is a stationary solution of (20), and there exists a one-to-one correspondence between solutions to (9) and (20). Namely, if f is a solution to (20) corresponding to the initial data $f_0 \in \mathcal{V}_c$, then $f + c$ is the solution to (9) corresponding to $f_0 + c \in \mathcal{V}$. Conversely, if f is a solution to (9) corresponding to the initial data $f_0 \in \mathcal{V}$, then $f - c$ is the solution to (20) corresponding to $f_0 - c \in \mathcal{V}_c$.

Hence, instead of studying the stability of the steady state c for (9), we study the stability of 0 for (20). Note that 0 is an eigenvalue of the derivative $\partial\psi(0)$. Moreover, the null solution is not asymptotically stable because each neighborhood of 0 in \mathcal{V}_c contains constants, which are also steady states of this flow.

However, it is inappropriate to pose the problem in this way because, as the next lemma shows, we have conservation of the fluid volume.

We introduce first some notation. Define by $\tilde{h}^{k+\alpha}(\mathbb{S}^1), k \in \{1, 2\}$, the closure of the set of all smooth functions on \mathbb{S}^1 with mean 0 in $C^{k+\alpha}(\mathbb{S}^1)$. We have defined in this way closed subspaces of the little Hölder spaces. It is obvious that

$$\tilde{h}^{2+\alpha}(\mathbb{S}^1) \stackrel{d}{\hookrightarrow} \tilde{h}^{1+\alpha}(\mathbb{S}^1) \quad \text{and} \quad \tilde{h}^{k+\alpha}(\mathbb{S}^1) = \left\{ f \in h^{k+\alpha}(\mathbb{S}^1) : \int_{\mathbb{S}^1} f \, dx = 0 \right\}, k \in \{1, 2\}.$$

Further, let $\tilde{\mathcal{V}} := \mathcal{V}_c \cap \tilde{h}^{2+\alpha}(\mathbb{S}^1)$.

LEMMA 3.1 (conservation of volume). *Given $g \in \tilde{\mathcal{V}}$, we have $\psi(g) \in \tilde{h}^{1+\alpha}(\mathbb{S}^1)$.*

Proof. Let $g + c = f \in \mathcal{V}$ and denote by u the solution of (6). We have

$$\begin{aligned} \int_{\mathbb{S}^1} \psi(g) dx &= - \int_{\mathbb{S}^1} \Phi(f) dx = - \int_{\mathbb{S}^1} \mathcal{B}(f, \mathcal{I}(f)) dx \\ &= - \int_{\mathbb{S}^1} \frac{Du}{\bar{\mu}(|Du|^2)}(x, f(x)) \cdot n(x) \, dx = - \int_{\Gamma_f} \frac{\partial_\nu u}{\bar{\mu}(|Du|^2)} \, d\sigma \\ &= - \int_{\Omega_f} \operatorname{div} \left(\frac{Du}{\bar{\mu}(|Du|^2)} \right) \, dx + \int_{\Gamma_0} \frac{\partial_\nu u}{\bar{\mu}(|Du|^2)} \, d\sigma = 0. \quad \square \end{aligned}$$

We consider now the restriction $\tilde{\psi}$ of ψ to $\tilde{\mathcal{V}}$. From Lemma 3.1 and Theorem 2.3 we conclude that $\tilde{\psi} \in C^\infty(\tilde{\mathcal{V}}, \tilde{h}^{1+\alpha}(\mathbb{S}^1))$. Given $h = \sum_{k \in \mathbb{Z}} \hat{h}(k)e^{ikx} \in \tilde{h}^{2+\alpha}(\mathbb{S}^1)$ we have

$$\partial\tilde{\psi}(0) \left[\sum_{k \in \mathbb{Z}} \hat{h}(k)e^{ikx} \right] = \sum_{k \in \mathbb{Z}} \lambda_k \hat{h}(k)e^{ikx},$$

where the coefficients (λ_k) are given by (14).

THEOREM 3.2. *Consider $\partial\tilde{\psi}(0)$ as an unbounded operator in $\tilde{h}^{1+\alpha}(\mathbb{S}^1)$ with dense domain $\tilde{h}^{2+\alpha}(\mathbb{S}^1)$. Then*

$$(21) \quad \sigma(\partial\tilde{\psi}(0)) = \{\lambda_k : k \geq 1\},$$

$$(22) \quad |\lambda| \cdot \|R(\lambda, -\partial\tilde{\psi}(0))\|_{\mathcal{L}(\tilde{h}^{1+\alpha}(\mathbb{S}^1))} \leq \kappa \quad \forall \operatorname{Re} \lambda \geq \omega,$$

where κ and ω are the constants from Proposition 2.6.

Proof. Given $\lambda \in \rho(\partial\psi(0))$ and $h = \sum_{k \in \mathbb{Z}} \hat{h}(k)e^{ikx} \in \tilde{h}^{1+\alpha}(\mathbb{S}^1)$ there exists a unique

$$u = \sum_{k \in \mathbb{Z}} \frac{1}{\lambda - \lambda_k} \hat{h}(k)e^{ikx} \in h^{2+\alpha}(\mathbb{S}^1),$$

such that $(\lambda - \partial\psi(0))u = h$. Since $\hat{h}(0) = 0$ we conclude that $\hat{u}(0) = 0$; thus $u \in \tilde{h}^{2+\alpha}(\mathbb{S}^1)$ and $\lambda \in \rho(\partial\tilde{\psi}(0))$. Moreover, $\partial\tilde{\psi}(0) \in \mathcal{L}is(\tilde{h}^{2+\alpha}(\mathbb{S}^1), \tilde{h}^{1+\alpha}(\mathbb{S}^1))$ and (21) holds. Relation (22) follows from Proposition 2.6, using the fact that

$$(\lambda - \partial\psi(0))[h] = (\lambda - \partial\tilde{\psi}(0))[h],$$

for all $h \in \tilde{h}^{2+\alpha}(\mathbb{S}^1)$ and $\operatorname{Re} \lambda \geq \omega$. □

In conclusion $-\partial\tilde{\psi}(0)$ belongs to $\mathcal{H}(\tilde{h}^{2+\alpha}(\mathbb{S}^1), \tilde{h}^{1+\alpha}(\mathbb{S}^1))$, and therefore, if the initial value f_0 in (20) belongs to $\tilde{h}^{2+\alpha}(\mathbb{S}^1)$, then the evolution takes place in $\tilde{h}^{2+\alpha}(\mathbb{S}^1)$. We have the following stability result for the flow (3).

THEOREM 3.3 (exponential stability). *For any $\omega < -\lambda_1 = \tanh c/\bar{\mu}(0)$, there are positive constants r and C such that for any $f_0 \in c + \tilde{h}^{2+\alpha}(\mathbb{S}^1)$ with $\|f_0 - c\|_{C^{2+\alpha}(\mathbb{S}^1)} \leq r$ the solution to (3) exists globally in time and the estimate*

$$\|f(t) - c\|_{C^{2+\alpha}(\mathbb{S}^1)} + \|f'(t)\|_{C^{1+\alpha}(\mathbb{S}^1)} \leq Ce^{-\omega t} \|f_0 - c\|_{C^{2+\alpha}(\mathbb{S}^1)} \quad \forall t \geq 0$$

holds.

Proof. Transferring problem (9) into a neighborhood of 0 in $\tilde{h}^{2+\alpha}(\mathbb{S}^1)$ we find that all of the assumptions of Theorem 9.1.2 in [13] are satisfied (cf. Theorem 3.2). The result follows directly from this theorem. □

REFERENCES

[1] H. AMANN, *Linear and Quasilinear Parabolic Problems*, Vol. I, Birkhäuser, Basel, 1995.
 [2] W. ARENDT AND S. BU, *Operator-valued Fourier multipliers on periodic Besov spaces and applications*, Proc. Edinb. Math. Soc. (2), 47 (2004), pp. 15–33.
 [3] J. BEAR AND Y. BACHMAT, *Introduction to Modeling of Transport Phenomena in Porous Media*, Kluwer Academic Publishers, Boston, 1990.
 [4] J. ESCHER, *On moving boundaries in deformable media*, Adv. Math. Sci. Appl., 7 (1997), pp. 275–316.

- [5] J. ESCHER AND B.-V. MATIOC, *Stability of the equilibria for periodic Stokesian Hele-Shaw flows*, J. Evol. Equ., 8 (2008), pp. 513–522.
- [6] J. ESCHER AND B.-V. MATIOC, *A moving boundary problem for periodic Stokesian Hele-Shaw flows*, Interfaces Free Bound., to appear.
- [7] J. ESCHER AND G. PROKERT, *Stability of the equilibria for spatially periodic flows in porous media*, Nonlinear Anal., 45 (2001), pp. 1061–1080.
- [8] J. ESCHER AND G. SIMONETT, *Maximal regularity for a free boundary problem*, NoDEA Nonlinear Differential Equations Appl., 2 (1995), pp. 463–510.
- [9] J. ESCHER AND G. SIMONETT, *Analyticity of the interface in a free boundary problem*, Math. Ann., 305 (1996), pp. 435–459.
- [10] J. ESCHER AND G. SIMONETT, *Classical solutions of multidimensional Hele–Shaw models*, SIAM J. Math. Anal., 28 (1997), pp. 1028–1047.
- [11] D. GILBARG AND T. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, Springer-Verlag, New York, 1997.
- [12] L. KONDIC, P. PALFFY-MAHORNÝ, AND M. J. SHELLEY, *Models of non-Newtonian Hele-Shaw flow*, Phys. Rev. E (3), 54 (1996), pp. R4536–R4539.
- [13] A. LUNARDI, *Analytic Semigroups and Optimal Regularity in Parabolic Problems*, Birkhäuser, Basel, 1995.
- [14] H.-J. SCHMEISSER AND H. TRIEBEL, *Topics in Fourier Analysis and Function Spaces*, John Wiley and Sons, New York, 1987.

CONVERGENCE TO EQUILIBRIUM FOR PARABOLIC-HYPERBOLIC TIME-DEPENDENT GINZBURG–LANDAU–MAXWELL EQUATIONS*

MAURIZIO GRASSELLI[†], HAO WU[‡], AND SONGMU ZHENG[§]

Abstract. We consider a Ginzburg–Landau–Maxwell model which describes the behavior of a two-dimensional superconducting material. The state variables are the complex-valued order parameter ψ , the magnetic potential \mathbf{A} , and the electric potential Φ . Under the choice of Coulomb (i.e., London) gauge, the resulting system is a parabolic-hyperbolic coupled system of nonlinear partial differential equations subject to suitable boundary and initial conditions. Global well-posedness results were proved in [M. Tsutsumi and H. Kasai, *Nonlinear Anal.*, 37 (1999), pp. 187–216], while the existence of global attractor and exponential attractors was proved in [V. Berti and S. Gatti, *Quart. Appl. Math.*, 64 (2006), pp. 617–639]. In this paper we use an extended Lojasiewicz–Simon approach to show that for any initial datum in certain phase space, the corresponding global solution converges to an equilibrium as time goes to infinity. Besides, we also provide an estimate on the convergence rate with respect to the phase space metric.

Key words. Ginzburg–Landau–Maxwell equations, superconductivity, Coulomb gauge, convergence to equilibrium, Lojasiewicz–Simon inequality

AMS subject classifications. Primary, 35B40; Secondary, 82D55

DOI. 10.1137/080717833

1. Introduction. This paper is concerned with the asymptotic behavior of solutions to a two-dimensional Ginzburg–Landau–Maxwell model of superconductivity proposed in [27] (see also [2, 13, 32]). When a superconducting material is kept close to a certain critical temperature, its behavior can be macroscopically described, within the Ginzburg–Landau phase transition theory (see [18] and its references, cf. also [8, Chapter 11]), by the state variables (ψ, \mathbf{A}, Φ) . Here ψ is the complex order parameter, whose squared modulus represents the concentration of the superconducting electrons, while \mathbf{A} and Φ are the magnetic and the electric potentials, respectively.

On account of [32], we introduce the equations governing the evolution of (ψ, \mathbf{A}, Φ) . Suppose that the material occupies a bounded domain $\Omega \subset \mathbb{R}^2$ with smooth boundary Γ for any time $t \geq 0$. Then the state variables $\psi : \Omega \times (0, \infty) \rightarrow \mathbb{C}$, $\mathbf{A} : \Omega \times (0, \infty) \rightarrow \mathbb{R}^2$ and $\Phi : \Omega \times (0, \infty) \rightarrow \mathbb{R}$ satisfy the following equations:

$$(1.1) \quad \psi_t - i\Phi\psi - D_{\mathbf{A}}^2 \psi - \lambda^2 (1 - |\psi|^2) \psi = 0,$$

$$(1.2) \quad \varepsilon \left(\tilde{\mathbf{A}}_t - \nabla \Phi \right)_t + \sigma \left(\tilde{\mathbf{A}}_t - \nabla \Phi \right) + \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} \left(\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi} \right) + \operatorname{curl} \mathbf{H}_{ext} = 0,$$

in $\Omega \times (0, \infty)$, subject to the boundary conditions

$$(1.3) \quad D_{\tilde{\mathbf{A}}} \psi \cdot \mathbf{n} = 0, \quad \left(\operatorname{curl} \tilde{\mathbf{A}} + \mathbf{H}_{ext} \right) \times \mathbf{n} = 0, \quad \left(\tilde{\mathbf{A}}_t - \nabla \Phi \right) \cdot \mathbf{n} = 0, \quad \text{on } \Gamma \times (0, \infty),$$

*Received by the editors March 6, 2008; accepted for publication (in revised form) September 23, 2008; published electronically January 21, 2009.

<http://www.siam.org/journals/sima/40-5/71783.html>

[†]Dipartimento di Matematica, Politecnico di Milano, 20133 Milano, Italy (maurizio.grasselli@polimi.it). This author was partially supported by the Italian PRIN Research Project 2006 *Problemi a frontiera libera, transizioni di fase e modelli di isteresi*.

[‡]School of Mathematical Sciences, Fudan University, Shanghai 200433, China (haowufd@yahoo.com). This author was supported by the China Postdoctoral Science Foundation.

[§]Corresponding author. Institute of Mathematics, Fudan University, Shanghai 200433, China (songmuzheng@yahoo.com). This author was supported by NSF of China under the grant 10631020 and by the Chinese Ministry of Education under grant 20050246002.

where

$$(1.4) \quad D_{\tilde{\mathbf{A}}} \psi = \nabla \psi - i\tilde{\mathbf{A}}\psi, \quad D_{\tilde{\mathbf{A}}}^2 \psi = \Delta \psi - 2i\tilde{\mathbf{A}} \cdot \nabla \psi - i\psi \nabla \cdot \tilde{\mathbf{A}} - |\tilde{\mathbf{A}}|^2 \psi.$$

Here \mathbf{H}_{ext} is the external magnetic field, which is assumed to be time independent for the sake of simplicity throughout this paper; the positive constants λ , ε , and σ represent the Ginzburg–Landau parameter, the dielectric coefficient, and the electric conductivity, respectively, and \mathbf{n} denotes the unit outward normal to the boundary Γ .

We recall that (1.2) is derived from Maxwell’s equations. In the so-called quasi-steady approximation, namely, when the displacement current $\varepsilon(\tilde{\mathbf{A}}_t - \nabla \Phi)_t$ is negligible, we obtain the well-known system proposed by Gor’kov and Éliashberg [13], which has been widely studied in the literature (see, e.g., [7, 9, 11, 18, 19, 22, 23, 24, 26, 27, 29, 30, 33] and reference therein). However, if surface charge must be taken into account, then the displacement current cannot be neglected (see [2]).

Similarly to the classical time-dependent Ginzburg–Landau equations, the evolution equations (1.1)–(1.2) and the boundary conditions (1.3) are invariant under the gauge transformation:

$$(1.5) \quad \mathcal{T}_\chi : (\psi, \tilde{\mathbf{A}}, \Phi) \longmapsto (\psi e^{i\chi}, \tilde{\mathbf{A}} + \nabla \chi, \Phi + \chi_t).$$

The gauge χ can be any (smooth) real scalar-valued function depending on x and t . In [3, 32], the authors choose the Coulomb gauge (also known as the London gauge). This choice entails that the following identities hold:

$$(1.6) \quad \operatorname{div} \tilde{\mathbf{A}} = 0 \quad \text{in } \Omega, \quad \int_{\Omega} \Phi dx = 0,$$

$$(1.7) \quad \tilde{\mathbf{A}} \cdot \mathbf{n} = 0, \quad \partial_{\mathbf{n}} \Phi = 0, \quad \text{on } \Gamma.$$

For the existence of the gauge χ corresponding to (1.6) and (1.7), one may refer, e.g., to [23, Lemma 2.2, Remark 2.4] where even more general cases have been discussed (see also [18, section 1.4.4]). Using the Coulomb gauge, $\tilde{\mathbf{A}}$ is a solenoidal vector field so that

$$(1.8) \quad D_{\tilde{\mathbf{A}}}^2 \psi = \Delta \psi - 2i\tilde{\mathbf{A}} \cdot \nabla \psi - |\tilde{\mathbf{A}}|^2 \psi.$$

Moreover, system (1.1)–(1.2) is subject to the boundary conditions (cf. (1.3), (1.7))

$$(1.9) \quad \partial_{\mathbf{n}} \psi = 0, \quad \tilde{\mathbf{A}} \cdot \mathbf{n} = 0, \quad \left(\operatorname{curl} \tilde{\mathbf{A}} + \mathbf{H}_{ext} \right) \times \mathbf{n} = 0, \quad \text{on } \Gamma \times (0, \infty).$$

It will be more convenient to treat problem (1.1), (1.2), (1.9) in a homogeneous form. For this purpose, we consider the following boundary value problem:

$$(1.10) \quad \begin{cases} \operatorname{curl}^2 \mathbf{A}_{ext} = \operatorname{curl} \mathbf{H}_{ext}, & \operatorname{div} \mathbf{A}_{ext} = 0, & \text{in } \Omega, \\ \mathbf{A}_{ext} \cdot \mathbf{n} = 0, & (\operatorname{curl} \mathbf{A}_{ext} - \mathbf{H}_{ext}) \times \mathbf{n} = 0, & \text{on } \Gamma. \end{cases}$$

Throughout this paper we always assume that $\mathbf{H}_{ext} \in H^1$. Then, it is easy to see that the convex quadratic functional

$$(1.11) \quad J(\mathbf{A}) = \int_{\Omega} |\operatorname{curl} \mathbf{A} - \mathbf{H}_{ext}|^2 dx$$

on the domain

$$(1.12) \quad \mathcal{D} := \{ \mathbf{A} \in \mathbf{H}^1(\Omega) : \operatorname{div} \mathbf{A} = 0 \text{ in } \Omega, \mathbf{A} \cdot \mathbf{n}|_{\Gamma} = 0 \}$$

admits a unique minimizer $\mathbf{A}_{ext} \in \mathcal{D}$, which is the (unique) solution to (1.10). The linear mapping $\mathbf{H}_{ext} \rightarrow \mathbf{A}_{ext}$ is continuous from $\mathbf{H}^\alpha(\Omega)$ to $\mathbf{H}^{\alpha+1}(\Omega)$ for $0 \leq \alpha \leq 1$ (see [11, Lemma 3, Lemma 4]; cf. also [9]).

Set now

$$(1.13) \quad \mathbf{A} = \tilde{\mathbf{A}} + \mathbf{A}_{ext}.$$

Then the Ginzburg–Landau–Maxwell equations (1.1) and (1.2) can be reduced to the following homogeneous form in terms of ψ and \mathbf{A} (cf. [3, 32]):

$$(1.14) \quad \psi_t - i\Phi\psi - D_{\tilde{\mathbf{A}}}^2\psi - \lambda^2(1 - |\psi|^2)\psi = 0,$$

$$(1.15) \quad \varepsilon(\mathbf{A}_t - \nabla\Phi)_t + \sigma(\mathbf{A}_t - \nabla\Phi) + \operatorname{curl}^2\mathbf{A} + \frac{i}{2}(\overline{\psi}D_{\tilde{\mathbf{A}}}\psi - \psi\overline{D_{\tilde{\mathbf{A}}}\psi}) = 0.$$

By applying the spatial divergence operator to (1.15) and taking (1.14) into account, we obtain the evolution equation for Φ , namely,

$$(1.16) \quad -\varepsilon\Delta\Phi_t - \sigma\Delta\Phi + \frac{i}{2}(\overline{\psi}\psi_t - \psi\overline{\psi}_t) + |\psi|^2\Phi = 0.$$

Moreover, we have (see (1.5))

$$(1.17) \quad \operatorname{div}\mathbf{A} = 0 \quad \text{in } \Omega, \quad \int_{\Omega} \Phi dx = 0,$$

and (1.14)–(1.17) are subject to the homogeneous boundary conditions (cf. (1.7)–(1.10))

$$(1.18) \quad \partial_{\mathbf{n}}\psi = 0, \quad \mathbf{A} \cdot \mathbf{n} = 0, \quad \operatorname{curl}\mathbf{A} \times \mathbf{n} = 0, \quad \partial_{\mathbf{n}}\Phi = 0, \quad \text{on } \Gamma \times (0, \infty),$$

as well as to the initial conditions

$$(1.19) \quad \psi|_{t=0} = \psi_0, \quad \mathbf{A}|_{t=0} = \mathbf{A}_0, \quad \mathbf{A}_t|_{t=0} = \mathbf{A}_1, \quad \Phi|_{t=0} = \Phi_0, \quad \text{in } \Omega$$

for some given initial data.

The main concern of this paper is the convergence to equilibrium of the solution to problem (1.14)–(1.19) as time goes to infinity. In what follows we recall related results in the literature on problem (1.14)–(1.19). The global well-posedness of problem (1.14)–(1.19) has been carefully investigated in [32] in both two and three dimensions, even though, physically speaking, the present model is essentially two-dimensional (see [18, sections 1.5, 1.6, 1.8]; cf. also [5, 23]). The results of [32] are also concerned with the case $\varepsilon = 0$, which has been studied by many other authors (see, e.g., [18, 23, 26, 29, 30, 33] and references therein). Regarding the large time behavior of solutions, problem (1.14)–(1.19) has recently been analyzed in [3] within the theory of infinite-dimensional dissipative dynamical systems. In that paper the authors show that the quoted problem generates a strongly continuous semigroup which possesses the global attractor and an exponential attractor. However, the issue of convergence to equilibrium for single solutions to problem (1.14)–(1.19) has not been studied in [3]. This is the goal of the present contribution. More precisely, we will prove that, for any initial datum in the phase space introduced in [3], the corresponding global solution to (1.14)–(1.19) converges to a single equilibrium as time goes to infinity. Besides, an estimate of the convergence rate in the phase space metric will also be given. The key ingredients of the present paper can be summarized as follows:

(I) Since (1.15) is hyperbolic, we can no longer take advantage of the smoothing effects for solutions to parabolic equations like, e.g., in [9] where the Lorentz gauge is used. Instead, we will use a result proved in [3] to show the precompactness of solutions.

(II) We need to develop an extended Lojasiewicz–Simon approach for which it is necessary to establish a suitable Lojasiewicz–Simon-type inequality. Such an inequality is a suitable adaptation of the one proved in [9] for a problem similar to ours with $\varepsilon = 0$ but subject to the Lorentz gauge $\Phi = -\omega(\nabla \cdot \mathbf{A})$, with ω being a nonnegative constant (for $\omega = 0$, see also [24]).

(III) Due to the hyperbolic nature of (1.15), the standard Lojasiewicz–Simon approach used in the parabolic case must be modified by introducing an appropriate auxiliary functional \mathfrak{F} (see section 5), which usually depends on the problem under consideration (see, e.g., [15, 19, 35] and references therein for similar problems).

(IV) In the existing literature the convergence rate in a (lower order) norm is usually obtained directly by using the Lojasiewicz–Simon approach (see, e.g., [9, 17, 36]). Then, estimates in higher order norms can be deduced by means of interpolation inequalities (cf. [17]) and, consequently, the decay exponent deteriorates. On the contrary, in this paper and some earlier papers by the authors, using the original equations and suitable energy estimates and constructing proper differential inequalities, we are able to preserve the original convergence rate when we estimate the convergence with respect to the phase space metric. This technique has been successfully applied to other equations as well (see, e.g., [14, 34, 35, 36]).

In [9] (see also [19]) the convergence to single equilibria for the case $\varepsilon = 0$ has been investigated in detail in the hypothesis of a time-dependent applied magnetic field and using the Lorentz gauge. The present contribution is a first step towards a more comprehensive analysis of the case $\varepsilon > 0$. The corresponding results for other widely used gauges (e.g., the Lorentz gauge) as well as for the case of a time-dependent applied magnetic field will be considered in the near future. We recall that the choice of the gauge plays an important role in the mathematical treatment of these models (see [10] and references therein).

This paper is organized as follows. In section 2 we introduce the notation and state the main results of this paper. Section 3 contains some preliminary results on the time-dependent Ginzburg–Landau–Maxwell model. In section 4 we state a suitable version of a Lojasiewicz–Simon-type inequality. The last two sections are devoted to the proof of our main results. More precisely, in section 5 we prove the convergence to a single equilibrium, while in section 6, we establish the convergence rate estimate.

2. Notation and main results. We use the notation introduced in [32] (see also [3]). As usual, $L^p(\Omega)$ and $W^{k,p}(\Omega)$ stand for the Lebesgue and the Sobolev spaces of real-valued functions, with the convention that $H^k(\Omega) = W^{k,2}(\Omega)$. We denote by bold letters the spaces of vector-valued functions, whereas a subscript \mathbb{C} characterizes those of complex-valued functions. Without further specifications, $\|\cdot\|$ stands for the norm in $L^2_{\mathbb{C}}(\Omega)$, $\mathbf{L}^2(\Omega)$, or $L^2(\Omega)$, according to the context. This norm is always induced by the scalar inner product

$$\langle u, v \rangle = \int_{\Omega} u(x)\bar{v}(x)dx,$$

where for a vector-valued function, the product $u\bar{v}$ is replaced by the Euclidean inner product $\mathbf{u} \cdot \bar{\mathbf{v}}$. Next, we introduce the function spaces

$$(2.1) \quad \mathbf{X} = \{ \mathbf{u} \in \mathbf{H}^2(\Omega) : \operatorname{div} \mathbf{u} = 0 \text{ in } \Omega, \quad \mathbf{u} \cdot \mathbf{n}|_{\Gamma} = 0, \quad \operatorname{curl} \mathbf{u} \times \mathbf{n}|_{\Gamma} = \mathbf{0} \}$$

equipped with the usual $\mathbf{H}^2(\Omega)$ norm. Let \mathbf{X}_0 be the completion of \mathbf{X} in $\mathbf{H}^1(\Omega)$ and

$$H_{0m}^1(\Omega) = \left\{ u \in H^1(\Omega) : \int_{\Omega} u dx = 0 \right\},$$

$$H_{0m}^2(\Omega) = \left\{ u \in H^2(\Omega) : \partial_{\mathbf{n}} u|_{\Gamma} = 0, \int_{\Omega} u dx = 0 \right\}.$$

We can thus define the phase space we shall work with:

$$(2.2) \quad \mathbb{X}_{\infty}^0 = \left\{ \psi \in H_{\mathbb{C}}^2(\Omega) : \partial_{\mathbf{n}} \psi|_{\Gamma} = 0, \|\psi\|_{L^{\infty}(\Omega)} \leq 1 \right\} \times \mathbf{X} \times \mathbf{X}_0 \times H_{0m}^2(\Omega),$$

which is a closed subset of the Banach space

$$(2.3) \quad \mathbb{X}^0 = \left\{ \psi \in H_{\mathbb{C}}^2(\Omega) : \partial_{\mathbf{n}} \psi|_{\Gamma} = 0 \right\} \times \mathbf{X} \times \mathbf{X}_0 \times H_{0m}^2(\Omega).$$

The norm on \mathbb{X}^0 is defined by

$$(2.4) \quad \|z\|_0^2 = \|z_1\|_{H_{\mathbb{C}}^2(\Omega)}^2 + \|z_2\|_{\mathbf{H}^2(\Omega)}^2 + \|z_3\|_{\mathbf{H}^1(\Omega)}^2 + \|z_4\|_{H^2(\Omega)}^2$$

for any $z = (z_1, z_2, z_3, z_4) \in \mathbb{X}^0$. We also introduce the Banach space

$$(2.5) \quad \mathbb{X}^{-1} = H_{\mathbb{C}}^1(\Omega) \times \mathbf{X}_0 \times \mathbf{L}^2(\Omega) \times H_{0m}^1(\Omega) \supset \mathbb{X}^0,$$

endowed with the norm

$$(2.6) \quad \|z\|_{-1}^2 = \|z_1\|_{H_{\mathbb{C}}^1(\Omega)}^2 + \|z_2\|_{\mathbf{H}^1(\Omega)}^2 + \|z_3\|_{\mathbf{L}^2(\Omega)}^2 + \|z_4\|_{H^1(\Omega)}^2$$

for any $z = (z_1, z_2, z_3, z_4) \in \mathbb{X}^{-1}$.

In what follows, we will make use of the Gagliardo–Nirenberg inequalities in dimension two, namely,

$$(2.7) \quad \|u\|_{L^4(\Omega)}^2 \leq C \|u\| \|u\|_{H^1(\Omega)} \quad \forall u \in H^1(\Omega),$$

$$(2.8) \quad \|u\|_{L^{\infty}(\Omega)}^2 \leq C \|u\| \|u\|_{H^2(\Omega)} \quad \forall u \in H^2(\Omega),$$

where $C > 0$ depends only on Ω .

We also recall that there exist constants $\kappa > 0$ and $c > 0$ depending only on Ω such that (see, e.g., [25, Corollary 3.51])

$$(2.9) \quad \|\mathbf{A}\| \leq \kappa \|\operatorname{curl} \mathbf{A}\| \quad \forall \mathbf{A} \in \mathbf{X}_0,$$

and (cf. [12, Chapter I, (5.45), p. 92] also [32, section 2])

$$(2.10) \quad \|\mathbf{A}\|_{\mathbf{H}^1(\Omega)} \leq c \|\operatorname{curl} \mathbf{A}\| \quad \forall \mathbf{A} \in \mathbf{X}_0.$$

In the remaining part of the paper, we always assume that

$$(2.11) \quad \mathbf{A}_{ext} \in \mathbf{X}_0 \cap \mathbf{H}^2(\Omega).$$

We now recall the following result.

PROPOSITION 2.1 ([3, Theorem 2.1, Proposition 3.1]). *Denote $z_0 = (\psi_0, \mathbf{A}_0, \mathbf{A}_1, \Phi_0) \in \mathbb{X}^0$. Problem (1.14)–(1.19) generates a strongly continuous semigroup $\{S(t)\}_{t \geq 0}$*

on the space \mathbb{X}^0 . If $z_0 = (\psi_0, \mathbf{A}_0, \mathbf{A}_1, \Phi_0) \in \mathbb{X}_\infty^0$, then $S(t)z_0 \in \mathbb{X}_\infty^0$. Moreover, there exists a positive constant C depending only on $\|z_0\|_0$ such that

$$(2.12) \quad \|S(t)z_0\|_0 \leq C \quad \forall t \geq 0.$$

The stationary problem corresponding to (1.14)–(1.19) is

$$(2.13) \quad -(\nabla - i(\mathbf{A}_\infty - \mathbf{A}_{ext}))^2 \psi_\infty - \lambda^2 (1 - |\psi_\infty|^2) \psi_\infty = 0,$$

$$(2.14) \quad \operatorname{curl}^2 \mathbf{A}_\infty + \frac{i}{2} (\bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty) + |\psi|^2 (\mathbf{A}_\infty - \mathbf{A}_{ext}) = 0,$$

$$(2.15) \quad \operatorname{div} \mathbf{A}_\infty = 0,$$

$$(2.16) \quad \partial_{\mathbf{n}} \psi_\infty = 0, \quad \mathbf{A}_\infty \cdot \mathbf{n} = 0, \quad \operatorname{curl} \mathbf{A}_\infty \times \mathbf{n} = 0, \quad \text{on } \Gamma,$$

$$(2.17) \quad \Phi_\infty = 0.$$

Set now

$$(2.18) \quad \mathbf{A}_\infty = \tilde{\mathbf{A}}_\infty + \mathbf{A}_{ext}.$$

Then problem (2.13)–(2.17) becomes

$$(2.19) \quad -(\nabla - i\tilde{\mathbf{A}}_\infty)^2 \psi_\infty - \lambda^2 (1 - |\psi_\infty|^2) \psi_\infty = 0,$$

$$(2.20) \quad \operatorname{curl}^2 \tilde{\mathbf{A}}_\infty + \frac{i}{2} (\bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty) + |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty + \operatorname{curl} \mathbf{H}_{ext} = 0,$$

$$(2.21) \quad \operatorname{div} \tilde{\mathbf{A}}_\infty = 0,$$

$$(2.22) \quad \partial_{\mathbf{n}} \psi_\infty = 0, \quad \tilde{\mathbf{A}}_\infty \cdot \mathbf{n} = 0, \quad (\operatorname{curl} \tilde{\mathbf{A}}_\infty + \mathbf{H}_{ext}) \times \mathbf{n} = 0, \quad \text{on } \Gamma,$$

and

$$(2.23) \quad \Phi_\infty = 0.$$

Define the functional

$$(2.24) \quad \mathcal{E}_0(\psi_\infty, \tilde{\mathbf{A}}_\infty) = \int_\Omega \left[\frac{\lambda^2}{4} (1 - |\psi_\infty|^2)^2 + \frac{1}{2} \left| (\nabla - i\tilde{\mathbf{A}}_\infty) \psi_\infty \right|^2 + \frac{1}{2} \left| \operatorname{curl} \tilde{\mathbf{A}}_\infty + \mathbf{H}_{ext} \right|^2 \right] dx,$$

on $H_{\mathbb{C}}^1 \times \mathbf{X}_0$. It is easy to see that a critical point $(\psi_\infty, \tilde{\mathbf{A}}_\infty)$ of \mathcal{E}_0 is a weak solution of (2.19)–(2.22). This stationary problem embodies the macroscopic quantum-mechanical nature of the superconducting state. There are many works on the stationary Ginzburg–Landau equations. For instance, the existence of vortex-like solutions to the time-independent Ginzburg–Landau equations was proved in [24]. An existence result for the stationary problem of the Ginzburg–Landau–Maxwell equations with Coulomb gauge in the L^p framework was recently given in [1]. An equivalence relation between solutions of the time-independent and time-dependent Ginzburg–Landau equations that describe the same physical state of a superconductor was established in [21].

We are now in a position to state our main theorem.

THEOREM 2.1. *For any initial datum $z_0 = (\psi_0, \mathbf{A}_0, \mathbf{A}_1, \Phi_0) \in \mathbb{X}_\infty^0$, the global solution $z(t) = (\psi(t), \mathbf{A}(t), \mathbf{A}_t(t), \Phi(t))$ to problem (1.14)–(1.19) converges in \mathbb{X}^0 to a unique equilibrium $z_\infty = (\psi_\infty, \mathbf{A}_\infty, \mathbf{0}, 0)$, where $(\psi_\infty, \mathbf{A}_\infty)$ satisfies (2.13)–(2.16). Moreover, the following estimate holds:*

$$(2.25) \quad \|z(t) - z_\infty\|_0 \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0,$$

where $C > 0$ is a constant depending only on $\|z_0\|_0$ and $\theta \in (0, \frac{1}{2})$ depends on z_∞ .

Remark 2.1. We recall that θ is the so-called Łojasiewicz exponent (see Lemma 4.2 below). Note that the convergence rate is estimated in the same norm where convergence takes place. This requires some work (see section 6).

3. Preliminaries. In this section we first recall the basic results proven in [3], which entail the precompactness of the trajectory originated from an initial datum in \mathbb{X}_∞^0 .

Let

$$(3.1) \quad \mathbb{X}^1 = H_{\mathbb{C}}^3(\Omega) \times [\mathbf{X} \cap \mathbf{H}^3(\Omega)] \times \mathbf{X} \times [H_{0m}^2(\Omega) \cap H^3(\Omega)],$$

with the standard norm defined by

$$(3.2) \quad \|z\|_1^2 = \|z_1\|_{H_{\mathbb{C}}^3(\Omega)}^2 + \|\mathbf{z}_2\|_{\mathbf{H}^3(\Omega)}^2 + \|\mathbf{z}_3\|_{\mathbf{H}^2(\Omega)}^2 + \|z_4\|_{H^3(\Omega)}^2$$

for any $z = (z_1, \mathbf{z}_2, \mathbf{z}_3, z_4) \in \mathbb{X}^1$.

For any given $z_0 \in \mathbb{X}_\infty^0$, let $z(t) = S(t)z_0$ be the solution to problem (1.14)–(1.19). Following [3, section 4], in order to prove the precompactness of $z(t)$, we make a decomposition of the solution z into a uniformly stable part z_d which decays to zero and a compact part z_c , that is,

$$(3.3) \quad z(t) = z_d(t) + z_c(t).$$

Here $z_d(t) = (\psi^d(t), \mathbf{A}^d(t), \mathbf{A}_t^d(t), \Phi^d(t))$ solves

$$(3.4) \quad \begin{cases} \psi_t^d - \Delta \psi^d + \psi^d = 0, \\ \varepsilon (\mathbf{A}_t^d - \nabla \Phi^d)_t + \sigma (\mathbf{A}_t^d - \nabla \Phi^d) + \operatorname{curl}^2 \mathbf{A}^d = 0, \\ -\varepsilon \Delta \Phi_t^d - \sigma \Delta \Phi^d = 0, \\ \operatorname{div} \mathbf{A}^d = 0, \quad \int_{\Omega} \Phi^d = 0, \\ \partial_{\mathbf{n}} \psi^d|_{\Gamma} = 0, \quad \mathbf{A}^d \cdot \mathbf{n}|_{\Gamma} = 0, \quad \operatorname{curl} \mathbf{A}^d \times \mathbf{n}|_{\Gamma} = 0, \quad \partial_{\mathbf{n}} \Phi^d|_{\Gamma} = 0, \\ \psi^d(0) = \psi_0, \quad \mathbf{A}^d(0) = \mathbf{A}_0, \quad \mathbf{A}_t^d(0) = \mathbf{A}_1, \quad \Phi^d(0) = \Phi_0, \end{cases}$$

and $z_c(t) = (\psi^c(t), \mathbf{A}^c(t), \mathbf{A}_t^c(t), \Phi^c(t))$ is the solution to

$$(3.5) \quad \begin{cases} \psi_t^c - \Delta \psi^c + \psi^c = F(\psi, \tilde{\mathbf{A}}, \Phi), \\ \varepsilon (\mathbf{A}_t^c - \nabla \Phi^c)_t + \sigma (\mathbf{A}_t^c - \nabla \Phi^c) + \operatorname{curl}^2 \mathbf{A}^c = \mathbf{G}(\psi, \tilde{\mathbf{A}}), \\ -\varepsilon \Delta \Phi_t^c - \sigma \Delta \Phi^c = H(\psi, \Phi), \\ \operatorname{div} \mathbf{A}^c = 0, \quad \int_{\Omega} \Phi^c = 0, \\ \partial_{\mathbf{n}} \psi^c|_{\Gamma} = 0, \quad \mathbf{A}^c \cdot \mathbf{n}|_{\Gamma} = 0, \quad \operatorname{curl} \mathbf{A}^c \times \mathbf{n}|_{\Gamma} = 0, \quad \partial_{\mathbf{n}} \Phi^c|_{\Gamma} = 0, \\ \psi^c(0) = 0, \quad \mathbf{A}^c(0) = \mathbf{0}, \quad \mathbf{A}_t^c(0) = \mathbf{0}, \quad \Phi^c(0) = 0, \end{cases}$$

where, in (3.5), we have set

$$(3.6) \quad F(\psi, \tilde{\mathbf{A}}, \Phi) = i\Phi\psi - |\tilde{\mathbf{A}}|^2\psi - 2i\tilde{\mathbf{A}} \cdot \nabla\psi + \lambda^2(1 - |\psi|^2)\psi + \psi,$$

$$(3.7) \quad \mathbf{G}(\psi, \tilde{\mathbf{A}}) = \frac{i}{2}(\psi\nabla\bar{\psi} - \bar{\psi}\nabla\psi) - \tilde{\mathbf{A}}|\psi|^2,$$

$$(3.8) \quad H(\psi, \Phi) = -\frac{i}{2}(\bar{\psi}\psi_t - \psi\bar{\psi}_t) - |\psi|^2\Phi.$$

Then the following lemmas show that the \mathbb{X}^0 -norm of $z_d(t)$ exponentially decays to zero as time goes to infinity, while $z_c(t)$ remains in a bounded set of \mathbb{X}^1 which is compact in \mathbb{X}^0 .

LEMMA 3.1 (cf. [3, Lemma 4.2]). *There exist constants $\gamma > 0$ and $\nu > 0$ independent of $\|z_0\|_0$ such that*

$$(3.9) \quad \|z_d(t)\|_0^2 \leq \gamma^2 e^{-2\nu t} \|z_0\|_0^2 \quad \forall t \geq 0.$$

LEMMA 3.2 (cf. [3, Lemma 4.3]). *There exists a constant $C > 0$ depending on $\|z_0\|_0$ such that*

$$(3.10) \quad \|z_c(t)\|_1 \leq C \quad \forall t \geq 0.$$

Let us define

$$(3.11) \quad \begin{aligned} \mathcal{E}(\psi, \tilde{\mathbf{A}}, \tilde{\mathbf{A}}_t, \Phi) &= \frac{1}{2} \|D_{\tilde{\mathbf{A}}}\psi\|^2 + \frac{\lambda^2}{4} \|1 - |\psi|^2\|^2 + \frac{1}{2} \|\operatorname{curl}\tilde{\mathbf{A}} + \mathbf{H}_{ext}\|^2 \\ &+ \frac{\varepsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 + \frac{\varepsilon}{2} \|\nabla\Phi\|^2. \end{aligned}$$

It is easy to see that \mathcal{E} is well defined and continuous on \mathbb{X}^{-1} (cf. (2.5)). In addition, we have what follows.

LEMMA 3.3. *\mathcal{E} is a global Lyapunov functional for the dynamical system $(\mathbb{X}_\infty^0, S(t))$.*

Proof. We only need to prove that \mathcal{E} decreases along any given trajectory $z(t) = S(t)z_0$. Here and in what follows, we will sometimes use formal arguments which can be rigorously justified, e.g., by using a Galerkin approximation scheme (see [32]).

Multiplying (1.1) by $\bar{\psi}_t$ and its conjugate by ψ_t , respectively, integrating on Ω and adding the results together, we get

$$(3.12) \quad \begin{aligned} &\frac{d}{dt} \left(\frac{1}{2} \|D_{\tilde{\mathbf{A}}}\psi\|^2 + \frac{\lambda^2}{4} \|1 - |\psi|^2\|^2 \right) + \|\psi_t\|^2 + \frac{i}{2} \langle \Phi, \bar{\psi}\psi_t - \psi\bar{\psi}_t \rangle \\ &- \langle \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle - \frac{i}{2} \langle \tilde{\mathbf{A}}_t, \bar{\psi}\nabla\psi - \psi\nabla\bar{\psi} \rangle = 0. \end{aligned}$$

On the other hand, multiplying (1.2) by $\tilde{\mathbf{A}}_t$ in $L^2(\Omega)$, we have

$$(3.13) \quad \begin{aligned} &\frac{d}{dt} \left(\frac{1}{2} \|\operatorname{curl}\tilde{\mathbf{A}} + \mathbf{H}_{ext}\|^2 + \frac{\varepsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 \right) + \sigma \|\tilde{\mathbf{A}}_t\|^2 + \langle \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle \\ &+ \frac{i}{2} \langle \bar{\psi}\nabla\psi - \psi\nabla\bar{\psi}, \tilde{\mathbf{A}}_t \rangle = 0. \end{aligned}$$

Moreover, taking the inner product of (1.16) with Φ in $L^2(\Omega)$ yields

$$(3.14) \quad \frac{\varepsilon}{2} \frac{d}{dt} \|\nabla\Phi\|^2 + \sigma \|\nabla\Phi\|^2 + \langle \Phi^2, |\psi|^2 \rangle + \frac{i}{2} \langle \bar{\psi}\psi_t - \psi\bar{\psi}_t, \Phi \rangle = 0.$$

Hence, we infer from (3.12)–(3.14) that, for all $t \geq 0$,

$$\begin{aligned}
 & \frac{d}{dt} \mathcal{E} \left(\psi, \tilde{\mathbf{A}}, \tilde{\mathbf{A}}_t, \Phi \right) \\
 &= -\|\psi_t\|^2 - i \langle \Phi, \overline{\psi} \psi_t - \psi \overline{\psi}_t \rangle - \sigma \|\tilde{\mathbf{A}}_t\|^2 - \sigma \|\nabla \Phi\|^2 - \langle \Phi^2, |\psi|^2 \rangle \\
 &= -\|\psi_t - i\Phi\psi\|_{L^2_{\mathbb{C}}(\Omega)}^2 - \sigma \|\tilde{\mathbf{A}}_t\|^2 - \sigma \|\nabla \Phi\|^2 \\
 (3.15) \quad & \leq 0,
 \end{aligned}$$

and the assertion follows. \square

Due to the above lemmas, we can see that $(S(t), \mathbb{X}_{\infty}^0)$ is a gradient system. Furthermore, on account of Lemmas 3.1, 3.2, 3.3, and some well-known results in dynamical systems (see, for instance, [4, Chapter 9]), we deduce the following.

LEMMA 3.4. *For any $z_0 \in \mathbb{X}_{\infty}^0$, the ω -limit set of z_0 is a nonempty compact connected subset of \mathbb{X}^0 . Furthermore, we have*

- (i) $\text{dist}_{\mathbb{X}^0}(S(t)z_0, \omega(z_0)) \rightarrow 0$ as $t \rightarrow \infty$;
- (ii) $\omega(z_0)$ consists of equilibria of the form $z_{\infty} = (\psi_{\infty}, \mathbf{A}_{\infty}, \mathbf{0}, 0)$, where $(\psi_{\infty}, \mathbf{A}_{\infty})$ satisfies (2.13)–(2.16);
- (iii) $S(t)\omega(z_0) = \omega(z_0) \forall t \geq 0$;
- (iv) \mathcal{E} is constant on $\omega(z_0)$.

Remark 3.1. If the set of equilibria were discrete, then we could conclude immediately that each solution converges to a single equilibrium (see (i) of Lemma 3.4). However, as is well known, in more than one spatial dimension, the set of stationary solutions can be a continuum for physically reasonable nonlinearity. The reader is referred, for instance, to [16, Remark 2.3.13]), where the following two-dimensional equation $-\Delta u + u^3 - \lambda u = 0$, $\lambda > 0$, endowed with a standard Dirichlet homogeneous boundary condition, is considered.

Remark 3.2. Lemmas 3.1 and 3.2 play a basic role in [3]. In particular, they enable the authors to prove that $(S(t), \mathbb{X}_{\infty}^0)$ has the global attractor $\mathcal{A}_{\varepsilon}$ bounded in \mathbb{X}^1 . In addition, due to Lemma 3.3, we infer that $\mathcal{A}_{\varepsilon}$ coincides with the unstable manifold of the set of equilibria (see, e.g., [31, Chapter 7, section 4]).

4. Extended Łojasiewicz–Simon inequality. Let us set $\psi = \psi_1 + i\psi_2$, where $\psi_1, \psi_2 : \Omega \times \mathbb{R}^+ \rightarrow \mathbb{R}$, and define $\vec{\psi} = (\psi_1, \psi_2)$. Recalling (2.24), we rewrite the functional \mathcal{E}_0 in the real form E_0 :

$$\begin{aligned}
 E_0 \left[(\psi_1, \psi_2), \tilde{\mathbf{A}} \right] &= \mathcal{E}_0 \left(\psi, \tilde{\mathbf{A}} \right) \\
 &= \int_{\Omega} \left[\frac{\lambda^2}{4} (1 - |\psi|^2)^2 + \frac{1}{2} \left| (\nabla - i\tilde{\mathbf{A}}) \psi \right|^2 + \frac{1}{2} \left| \text{curl} \tilde{\mathbf{A}} + \mathbf{H}_{ext} \right|^2 \right] dx \\
 &= \int_{\Omega} \left[\frac{\lambda^2}{4} (1 - \psi_1^2 - \psi_2^2)^2 + \frac{1}{2} \left| \nabla \psi_1 + \tilde{\mathbf{A}} \psi_2 \right|^2 + \frac{1}{2} \left| \nabla \psi_2 - \tilde{\mathbf{A}} \psi_1 \right|^2 \right. \\
 (4.1) \quad & \left. + \frac{1}{2} \left| \text{curl} \tilde{\mathbf{A}} + \mathbf{H}_{ext} \right|^2 \right] dx.
 \end{aligned}$$

In analogy to [9, Lemma 5.1], we can see that the functional E_0 is analytic on

$$(4.2) \quad \mathcal{X} = [H^1(\Omega)]^2 \times \mathbf{X}_0$$

in the sense of Deimling (see [6, Definition 15.1, pp. 150]). Then a direct calculation shows what follows.

LEMMA 4.1. *The Fréchet derivative DE_0 of E_0 is an analytic mapping from \mathcal{X} to its dual \mathcal{X}' , and it is given by*

$$(4.3) \quad \left\langle DE_0[(\psi_1, \psi_2), \tilde{\mathbf{A}}], [(v_1, v_2), \mathbf{V}] \right\rangle = \langle \partial_{\psi_1} E_0, v_1 \rangle + \langle \partial_{\psi_2} E_0, v_2 \rangle + \langle \partial_{\tilde{\mathbf{A}}} E_0, \mathbf{V} \rangle,$$

for any $[(\psi_1, \psi_2), \tilde{\mathbf{A}}], [(v_1, v_2), \mathbf{V}] \in \mathcal{X}$, where

$$(4.4) \quad \left\langle \partial_{\psi_1} E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}], v_1 \right\rangle \\ = \int_{\Omega} \left[(\nabla \psi_1 + \tilde{\mathbf{A}} \psi_2) \cdot \nabla v_1 - (\nabla \psi_2 - \tilde{\mathbf{A}} \psi_1) \cdot \tilde{\mathbf{A}} v_1 \right] dx - \lambda^2 \int_{\Omega} (1 - |\psi|^2) \psi_1 v_1 dx,$$

$$(4.5) \quad \left\langle \partial_{\psi_2} E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}], v_2 \right\rangle \\ = \int_{\Omega} \left[(\nabla \psi_2 - \tilde{\mathbf{A}} \psi_1) \cdot \nabla v_2 + (\nabla \psi_1 + \tilde{\mathbf{A}} \psi_2) \cdot \tilde{\mathbf{A}} v_2 \right] dx - \lambda^2 \int_{\Omega} (1 - |\psi|^2) \psi_2 v_2 dx,$$

$$(4.6) \quad \left\langle \partial_{\tilde{\mathbf{A}}} E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}], \mathbf{V} \right\rangle \\ = \int_{\Omega} (\operatorname{curl} \tilde{\mathbf{A}} + \mathbf{H}_{ext}) \cdot (\operatorname{curl} \mathbf{V}) dx + \int_{\Omega} \left[-(\psi_1 \nabla \psi_2 - \psi_2 \nabla \psi_1) + \tilde{\mathbf{A}} |\psi|^2 \right] \cdot \mathbf{V} dx.$$

Let us introduce now the quadratic form $\mathcal{J} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ given by

$$(4.7) \quad \mathcal{J}([(v_1, v_2), \mathbf{V}], [(w_1, w_2), \mathbf{W}]) := \int_{\Omega} (\nabla \vec{v} \cdot \nabla \vec{w} + \vec{v} \cdot \vec{w}) dx + \int_{\Omega} (\operatorname{curl} \mathbf{V}) \cdot (\operatorname{curl} \mathbf{W}) dx,$$

$\forall [(v_1, v_2), \mathbf{V}], [(w_1, w_2), \mathbf{W}] \in \mathcal{X}$.

It follows from (2.10) that $(\mathcal{J}([(v_1, v_2), \mathbf{V}], [(v_1, v_2), \mathbf{V}]))^{1/2}$ defines an equivalent norm of $[(v_1, v_2), \mathbf{V}]$ on \mathcal{X} . Based on this fact, we can retrace the steps of [9, section 6] to prove the following version of the Lojasiewicz–Simon inequality, which is a slight modification of [9, Proposition 6.1]. We thus omit the details. This lemma can be viewed as an extended version of Simon’s result [28] for the scalar case under the use of L^2 norm.

LEMMA 4.2 (extended Lojasiewicz–Simon inequality). *Let $[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_{\infty}] \in \mathcal{X}$ be a critical point of E_0 . Then there exist constants $\beta > 0$ and $\theta \in (0, \frac{1}{2})$ depending on $[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_{\infty}]$ such that, for any $[(\psi_1, \psi_2), \tilde{\mathbf{A}}] \in \mathcal{X}$ satisfying*

$$(4.8) \quad \left\| [(\psi_1, \psi_2), \tilde{\mathbf{A}}] - [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_{\infty}] \right\|_{\mathcal{X}} < \beta,$$

we have

$$(4.9) \quad \left| E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_{\infty}] \right|^{1-\theta} \leq \left\| DE_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}] \right\|_{\mathcal{X}'}.$$

5. Convergence to equilibrium. In this section we prove the convergence to equilibrium for the global solutions to our Ginzburg–Landau–Maxwell equations. First, we show the convergence for $\tilde{\mathbf{A}}_t(t)$ and $\Phi(t)$ to 0, which is given by the following lemma.

LEMMA 5.1. *We have*

$$(5.1) \quad \lim_{t \rightarrow +\infty} \left\| \tilde{\mathbf{A}}_t(t) \right\|_{\mathbf{X}_0} = 0,$$

$$(5.2) \quad \lim_{t \rightarrow +\infty} \left\| \Phi(t) \right\|_{H^2(\Omega)} = 0.$$

Proof. Integrating (3.15) from 0 to t , we obtain

$$(5.3) \quad \begin{aligned} \mathcal{E}(t) + \int_0^t \|\psi_t(\tau) - i\Phi(\tau)\psi(\tau)\|^2 d\tau + \sigma \int_0^t \|\tilde{\mathbf{A}}_t(\tau)\|^2 d\tau + \sigma \int_0^t \|\nabla\Phi(\tau)\|^2 d\tau \\ \leq \mathcal{E}(0) < \infty. \end{aligned}$$

Since \mathcal{E} is nonnegative, we deduce

$$(5.4) \quad \int_0^\infty \|\tilde{\mathbf{A}}_t(\tau)\|^2 d\tau < \infty, \quad \int_0^\infty \|\nabla\Phi(\tau)\|^2 d\tau < \infty,$$

and

$$(5.5) \quad \int_0^\infty \|\psi_t(\tau) - i\Phi(\tau)\psi(\tau)\|^2 d\tau < \infty.$$

From (1.16), we infer that

$$(5.6) \quad \begin{aligned} \varepsilon \|\Delta\Phi_t\| &\leq \sigma \|\Delta\Phi\| + \left\| \frac{i}{2} (\overline{\psi}\psi_t - \psi\overline{\psi}_t) \right\| + \|\psi\|^2 \|\Phi\| \\ &\leq \sigma \|\Delta\Phi\| + C (\|\psi_t\| + \|\nabla\Phi\|) \end{aligned}$$

and, on account of (1.17), $\int_\Omega \Phi_t dx = 0$. Thus we have

$$(5.7) \quad \|\nabla\Phi_t\| \leq \|\Phi_t\|_{H^2(\Omega)} \leq C \|\Delta\Phi_t\| \leq C (\|\Delta\Phi\| + \|\psi_t\| + \|\nabla\Phi\|).$$

On the other hand, from (1.1), (1.8), and the Sobolev embedding theorem, we have

$$(5.8) \quad \begin{aligned} \|\psi_t\| &\leq C \left(\|\Phi\psi\| + \|\Delta\psi\| + \|\tilde{\mathbf{A}} \cdot \nabla\psi\| + \|\tilde{\mathbf{A}}\|^2 \|\psi\| + \|(1 - |\psi|^2)\psi\| \right) \\ &\leq C \left(\|\Phi\| \|\psi\|_{L^\infty(\Omega)} + \|\psi\|_{H^2_c(\Omega)} + \|\tilde{\mathbf{A}}\|_{\mathbf{H}^2(\Omega)} \|\nabla\psi\| + \|\tilde{\mathbf{A}}\|_{\mathbf{H}^2(\Omega)}^2 \|\psi\| \right. \\ &\quad \left. + \|\psi\|_{H^2_c(\Omega)}^3 \right). \end{aligned}$$

Then, it follows from (5.7), (5.8), and the uniform estimate (2.12) that

$$(5.9) \quad \|\nabla\Phi_t(t)\| \leq C \quad \forall t \geq 0.$$

Similarly, it follows from (1.15) and estimates (2.12) and (5.9) that

$$(5.10) \quad \begin{aligned} \|\tilde{\mathbf{A}}_{tt}\| &\leq C \left(\|\nabla\Phi_t\| + \|\tilde{\mathbf{A}}_t\| + \|\nabla\Phi\| + \|\operatorname{curl}^2 \tilde{\mathbf{A}}\| + \|\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}\| \right) \\ &\leq C \left(\|\nabla\Phi_t\| + \|\tilde{\mathbf{A}}_t\| + \|\nabla\Phi\| + \|\operatorname{curl}^2 \tilde{\mathbf{A}}\| + \|\psi\|_{L^\infty(\Omega)} \|\nabla\psi\| \right. \\ &\quad \left. + \|\psi\|_{L^\infty(\Omega)}^2 \|\tilde{\mathbf{A}}\| \right) \leq C. \end{aligned}$$

Let $h(t) = \|\tilde{\mathbf{A}}_t(t)\|^2$. Then we have

$$(5.11) \quad \left| \frac{dh}{dt} \right| = 2 \left| \langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}}_{tt} \rangle \right| \leq 2 \|\tilde{\mathbf{A}}_t\| \|\tilde{\mathbf{A}}_{tt}\| \leq C.$$

From (5.4) and (5.11), we infer that $h \in L^1(0, \infty)$, and it is globally Lipschitz on $(0, \infty)$. Hence we have

$$(5.12) \quad \lim_{t \rightarrow +\infty} \|\tilde{\mathbf{A}}_t(t)\| = 0.$$

Thus, on account of Lemmas 3.1 and 3.2, we deduce (5.1).

We recall that $\|\psi_0\|_{L^\infty(\Omega)} \leq 1$. Therefore, Proposition 2.1 implies

$$(5.13) \quad \|\psi(t)\|_{L^\infty(\Omega)} \leq 1 \quad \forall t \geq 0.$$

Thus, recalling (1.6), the Poincaré inequality for Φ combined with (5.4) and (5.13) implies that

$$(5.14) \quad \int_0^\infty \|-i\Phi(\tau)\psi(\tau)\|^2 d\tau \leq C \int_0^\infty \|\nabla\Phi(\tau)\|^2 d\tau < \infty.$$

This, together with (5.5), yields

$$(5.15) \quad \int_0^\infty \|\psi_t(\tau)\|^2 d\tau < \infty.$$

Observe now that

$$(5.16) \quad |i\langle \bar{\psi}\psi_t - \psi\bar{\psi}_t, \Phi \rangle| \leq 2\|\psi\Phi\|\|\psi_t\| \leq \langle \Phi^2, |\psi|^2 \rangle + \|\psi_t\|^2.$$

Hence, from (3.14), the Hölder inequality, and the Cauchy–Schwarz inequality, we get

$$(5.17) \quad \begin{aligned} & \frac{\varepsilon}{2} \frac{d}{dt} \|\nabla\Phi\|^2 + \sigma \|\nabla\Phi\|^2 + \langle \Phi^2, |\psi|^2 \rangle \\ & \leq \left| \frac{i}{2} \langle \bar{\psi}\psi_t - \psi\bar{\psi}_t, \Phi \rangle \right| \leq \|\psi\Phi\|\|\psi_t\| \leq \frac{1}{2} \langle \Phi^2, |\psi|^2 \rangle + \frac{1}{2} \|\psi_t\|^2. \end{aligned}$$

Setting $h_1(t) = \|\nabla\Phi(t)\|^2$, it follows from (5.4), (5.15), (5.17), and [36, Lemma 6.2.1] that

$$(5.18) \quad \lim_{t \rightarrow +\infty} h_1(t) = 0,$$

which entails

$$(5.19) \quad \lim_{t \rightarrow +\infty} \|\Phi(t)\|_{H^1(\Omega)} = 0.$$

Hence, (5.2) is a consequence of Lemma 3.1, Lemma 3.2, and (5.19). \square

To establish the convergence for ψ and $\tilde{\mathbf{A}}$ (or \mathbf{A}), it is convenient to rewrite \mathcal{E} in the real form. As before, we denote $\psi = \psi_1 + i\psi_2$, with ψ_1, ψ_2 being real functions. Then we have (cf. (3.11) and (4.1))

$$(5.20) \quad \begin{aligned} E \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] & := \mathcal{E} \left(\psi, \tilde{\mathbf{A}}, \tilde{\mathbf{A}}_t, \Phi \right) \\ & = E_0 \left[(\psi_1, \psi_2), \tilde{\mathbf{A}} \right] + \frac{\varepsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 + \frac{\varepsilon}{2} \|\nabla\Phi\|^2 \end{aligned}$$

and accordingly, (3.15) becomes

$$(5.21) \quad \frac{d}{dt} E \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] = -\|\psi_{1t} + \Phi\psi_2\|^2 - \|\psi_{2t} - \Phi\psi_1\|^2 - \sigma \|\tilde{\mathbf{A}}_t\|^2 - \sigma \|\nabla\Phi\|^2 \leq 0$$

$\forall t \geq 0$.

It follows from (2.12) and the Poincaré inequality that

$$(5.22) \quad \|\psi_{1t}\|^2 \leq 2 \left(\|\psi_{1t} + \Phi\psi_2\|^2 + \|\Phi\psi_2\|^2 \right) \leq C_1 \left(\|\psi_{1t} + \Phi\psi_2\|^2 + \|\nabla\Phi\|^2 \right),$$

$$(5.23) \quad \|\psi_{2t}\|^2 \leq 2 \left(\|\psi_{2t} - \Phi\psi_1\|^2 + \|\Phi\psi_1\|^2 \right) \leq C_1 \left(\|\psi_{2t} - \Phi\psi_1\|^2 + \|\nabla\Phi\|^2 \right),$$

where $C_1 > 0$ is a constant. By choosing

$$(5.24) \quad C_2 \in \left(0, \min \left\{ \frac{\sigma}{4C_1}, \frac{1}{2C_1} \right\} \right),$$

we infer from (5.21) that

$$(5.25) \quad \frac{d}{dt} E \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] \leq -C_2 \|\psi_{1t}\|^2 - C_2 \|\psi_{2t}\|^2 - \frac{\sigma}{2} \|\tilde{\mathbf{A}}_t\|^2 - \frac{\sigma}{2} \|\nabla\Phi\|^2.$$

Next, we introduce the following auxiliary functional:

$$(5.26) \quad \mathcal{G} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] = \left\langle \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext}, \tilde{\mathbf{A}}_t - \nabla\Phi \right\rangle_{\mathbf{X}'_0}.$$

Then we have what follows.

LEMMA 5.2. *The following inequality holds for all $t \geq 0$:*

$$(5.27) \quad \begin{aligned} & \frac{d}{dt} \mathcal{G} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] \\ & \leq -\frac{1}{2\varepsilon} \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0}^2 \\ & + C_3 \left(\|\psi_{1t}\|^2 + \|\psi_{2t}\|^2 + \|\nabla\psi_{1t}\|^2 + \|\nabla\psi_{2t}\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla\Phi\|^2 \right). \end{aligned}$$

Proof. By direct calculation and (1.2), we get

$$(5.28) \quad \begin{aligned} & \frac{d}{dt} \mathcal{G} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] \\ & = \left\langle \left(\operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right)_t, \tilde{\mathbf{A}}_t - \nabla\Phi \right\rangle_{\mathbf{X}'_0} \\ & \quad - \frac{1}{\varepsilon} \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0}^2 \\ & \quad - \frac{\sigma}{\varepsilon} \left\langle \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext}, \tilde{\mathbf{A}}_t - \nabla\Phi \right\rangle_{\mathbf{X}'_0} \\ & := I_1 + I_2 + I_3. \end{aligned}$$

Let us estimate the right-hand side term by term. Observe first that

$$(5.29) \quad \left| \left\langle \operatorname{curl}^2 \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}}_t - \nabla\Phi \right\rangle_{\mathbf{X}'_0} \right| = \left| \left\langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}}_t \right\rangle - \left\langle \tilde{\mathbf{A}}_t, \nabla\Phi \right\rangle \right| \leq \frac{3}{2} \|\tilde{\mathbf{A}}_t\|^2 + \frac{1}{2} \|\nabla\Phi\|^2.$$

Recalling that we are in dimension two, it follows from a well-known Sobolev embedding theorem that

$$(5.30) \quad \mathbf{X}_0 \subset \mathbf{H}^1(\Omega) \hookrightarrow \mathbf{L}^p(\Omega) \quad \forall p > 1,$$

$$(5.31) \quad \mathbf{L}^q(\Omega) \hookrightarrow (\mathbf{H}^1(\Omega))' \subset \mathbf{X}'_0 \quad \forall q > 1.$$

Therefore, using (2.12) and taking $q = \frac{4}{3}$ in (5.31), we have

$$\begin{aligned}
& \left| \left\langle \frac{i}{2} (\bar{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi})_t, \tilde{\mathbf{A}}_t - \nabla \Phi \right\rangle_{\mathbf{X}'_0} \right| \\
& \leq \left\| \left[-(\psi_1 \nabla \psi_2 - \psi_2 \nabla \psi_1) + \tilde{\mathbf{A}} |\psi|^2 \right]_t \right\|_{\mathbf{X}'_0} \left\| \tilde{\mathbf{A}}_t - \nabla \Phi \right\|_{\mathbf{X}'_0} \\
& \leq \left(\|\psi_{1t} \nabla \psi_2\|_{\mathbf{L}^{\frac{4}{3}}(\Omega)} + \|\psi_1 \nabla \psi_{2t}\| + \|\psi_{2t} \nabla \psi_1\|_{\mathbf{L}^{\frac{4}{3}}(\Omega)} + \|\psi_2 \nabla \psi_{1t}\| \right. \\
& \quad \left. + \|\tilde{\mathbf{A}}_t |\psi|^2\| + 2 \|\tilde{\mathbf{A}} (\psi_1 \psi_{1t} + \psi_2 \psi_{2t})\| \right) \left(\|\tilde{\mathbf{A}}_t\| + \|\nabla \Phi\| \right) \\
& \leq \left(\|\nabla \psi_2\|_{\mathbf{L}^4(\Omega)} \|\psi_{1t}\| + \|\nabla \psi_1\|_{\mathbf{L}^4(\Omega)} \|\psi_{2t}\| + C \|\nabla \psi_{1t}\| + C \|\nabla \psi_{2t}\| \right. \\
& \quad \left. + C \|\tilde{\mathbf{A}}_t\| + C \|\psi_{1t}\| + C \|\psi_{2t}\| \right) \left(\|\tilde{\mathbf{A}}_t\| + \|\nabla \Phi\| \right) \\
(5.32) \quad & \leq C \left(\|\psi_{1t}\|^2 + \|\psi_{2t}\|^2 + \|\nabla \psi_{1t}\|^2 + \|\nabla \psi_{2t}\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right).
\end{aligned}$$

As a result, we get

$$(5.33) \quad I_1 \leq C \left(\|\psi_{1t}\|^2 + \|\psi_{2t}\|^2 + \|\nabla \psi_{1t}\|^2 + \|\nabla \psi_{2t}\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right).$$

On the other hand, thanks to the Cauchy–Schwarz inequality, we obtain

$$(5.34) \quad I_3 \leq \frac{1}{2\varepsilon} \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\bar{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0}^2 + C \left(\|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right).$$

Thus, (5.27) follows from (5.28), (5.33), and (5.34). \square

Let

$$(5.35) \quad \mathcal{G}_1[(\psi_1, \psi_2), \Phi] = \frac{1}{2} \|\psi_{1t}\|^2 + \frac{1}{2} \|\psi_{2t}\|^2 + \frac{\varepsilon}{2} \|\Delta \Phi\|^2 + \frac{\sigma}{2} \|\nabla \Phi\|^2.$$

Then we have the following.

LEMMA 5.3. *The following inequality holds for all $t \geq 0$:*

$$\begin{aligned}
(5.36) \quad \frac{d}{dt} \mathcal{G}_1[(\psi_1, \psi_2), \Phi] & \leq -\frac{1}{2} \|\nabla \psi_{1t}\|^2 - \frac{1}{2} \|\nabla \psi_{2t}\|^2 - \frac{\sigma}{2} \|\Delta \Phi\|^2 - \frac{\varepsilon}{2} \|\nabla \Phi_t\|^2 \\
& + C_4 \left(\|\psi_{1t}\|^2 + \|\psi_{2t}\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right).
\end{aligned}$$

Proof. Differentiating (1.14) with respect to time, adding the resulting equation multiplied by $\bar{\psi}_t$ to its conjugate multiplied by ψ_t , then integrating on Ω yield

$$\begin{aligned}
(5.37) \quad \frac{1}{2} \frac{d}{dt} \|\psi_t\|^2 & = -\|\nabla \psi_t\|^2 + \frac{i}{2} \langle \Phi_t, \psi \bar{\psi}_t - \bar{\psi} \psi_t \rangle + \lambda^2 \langle (1 - 2|\psi|^2), |\psi_t|^2 \rangle \\
& - \frac{\lambda^2}{2} \left(\langle \psi^2, \bar{\psi}_t^2 \rangle + \langle \bar{\psi}^2, \psi_t^2 \rangle \right) - \langle |\tilde{\mathbf{A}}|^2, |\psi_t|^2 \rangle + i \langle \mathbf{A}_t, \psi \nabla \bar{\psi}_t - \bar{\psi} \nabla \psi_t \rangle \\
& - \langle \tilde{\mathbf{A}} \cdot \mathbf{A}_t, \psi \bar{\psi}_t + \bar{\psi} \psi_t \rangle - i \langle \tilde{\mathbf{A}}, \bar{\psi}_t \nabla \psi_t - \psi_t \nabla \bar{\psi}_t \rangle.
\end{aligned}$$

It follows from (2.12), standard Sobolev embeddings, and the Poincaré inequality that the terms on the right-hand side of (5.37) can be estimated as follows:

$$\begin{aligned}
(5.38) \quad \frac{i}{2} \langle \Phi_t, \psi \bar{\psi}_t - \bar{\psi} \psi_t \rangle & + \lambda^2 \langle (1 - 2|\psi|^2), |\psi_t|^2 \rangle - \frac{\lambda^2}{2} \left(\langle \psi^2, \bar{\psi}_t^2 \rangle + \langle \bar{\psi}^2, \psi_t^2 \rangle \right) \\
& \leq \frac{\varepsilon}{4} \|\nabla \Phi_t\|^2 + C \|\psi_t\|^2
\end{aligned}$$

$$\begin{aligned}
 (5.39) \quad & - \left\langle |\tilde{\mathbf{A}}|^2, |\psi_t|^2 \right\rangle + i \left\langle \mathbf{A}_t, \psi \nabla \bar{\psi}_t - \bar{\psi} \nabla \psi_t \right\rangle - \left\langle \tilde{\mathbf{A}} \cdot \mathbf{A}_t, \psi \bar{\psi}_t + \bar{\psi} \psi_t \right\rangle \\
 & - i \left\langle \tilde{\mathbf{A}}, \bar{\psi}_t \nabla \psi_t - \psi_t \nabla \bar{\psi}_t \right\rangle \\
 & \leq \frac{1}{2} \|\nabla \psi_t\|^2 + C \|\psi_t\|^2 + C \|\tilde{\mathbf{A}}_t\|^2.
 \end{aligned}$$

Multiplying now (1.16) by $-\Delta \Phi + \Phi_t$ and integrating on Ω yield

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} (\varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2) \\
 & = -\sigma \|\Delta \Phi\|^2 - \varepsilon \|\nabla \Phi_t\|^2 - \frac{i}{2} \langle \bar{\psi} \psi_t - \psi \bar{\psi}_t, -\Delta \Phi \rangle - \langle |\psi|^2 \Phi, -\Delta \Phi \rangle \\
 & \quad - \frac{i}{2} \langle \bar{\psi} \psi_t - \psi \bar{\psi}_t, \Phi_t \rangle - \langle |\psi|^2 \Phi, \Phi_t \rangle \\
 (5.40) \quad & \leq -\sigma \|\Delta \Phi\|^2 - \varepsilon \|\nabla \Phi_t\|^2 + \frac{\sigma}{2} \|\Delta \Phi\|^2 + \frac{\varepsilon}{4} \|\nabla \Phi_t\|^2 + C \|\psi_t\|^2 + C \|\nabla \Phi\|^2.
 \end{aligned}$$

Then from (5.37)–(5.40), it follows that

$$\begin{aligned}
 & \frac{1}{2} \frac{d}{dt} (\|\psi_t\|^2 + \varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2) \\
 (5.41) \quad & \leq -\frac{1}{2} \|\nabla \psi_t\|^2 - \frac{\sigma}{2} \|\Delta \Phi\|^2 - \frac{\varepsilon}{2} \|\nabla \Phi_t\|^2 + C \left(\|\psi_t\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right),
 \end{aligned}$$

which gives (5.36). \square

We are now able to prove that ψ_t vanishes at infinity.

LEMMA 5.4. *There holds*

$$(5.42) \quad \lim_{t \rightarrow +\infty} \|\psi_t(t)\| = 0.$$

Proof. It follows from (5.41) that

$$(5.43) \quad \frac{1}{2} \frac{d}{dt} (\|\psi_t\|^2 + \varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2) + \frac{\sigma}{2} \|\Delta \Phi\|^2 \leq C \left(\|\psi_t\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla \Phi\|^2 \right).$$

Since $z_0 = (\psi_0, \mathbf{A}_0, \mathbf{A}_1, \Phi_0) \in \mathbb{X}_\infty^0$, from (1.1) and the Sobolev embedding theorem we infer that $\|\psi_t|_{t=0}\| \leq C$, where C is a constant depending on $\|z_0\|_0$. Therefore, (5.4), (5.15), and (5.43) yield

$$(5.44) \quad \|\psi_t(t)\| \leq C \quad \forall t \geq 0,$$

$$(5.45) \quad \int_0^\infty \|\Delta \Phi(\tau)\|^2 d\tau < \infty.$$

Denote $h_2(t) = \|\psi_t\|^2 + \varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2$. Then (5.43), Lemma 5.1, and (5.44) imply

$$(5.46) \quad \frac{d}{dt} h_2 \leq C \quad \forall t \geq 0.$$

On the other hand, we deduce $h_2(t) \in L^1(0, \infty)$ from (5.4), (5.15), and (5.45). As a consequence of [36, Lemma 6.2.1],

$$(5.47) \quad \lim_{t \rightarrow +\infty} h_2(t) = 0.$$

The proof is complete. \square

As we mentioned in the Introduction, due to the hyperbolic nature of the equation for $\tilde{\mathbf{A}}$, a key point in the proof of convergence result is to construct an appropriate auxiliary functional in order to implement the Łojasiewicz–Simon approach. For this purpose, we introduce the following functional $\mathfrak{F} : \mathbb{X}^0 \rightarrow \mathbb{R}$:

$$(5.48) \quad \mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] := E [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] + \eta \mathcal{G} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] + \eta_1 \mathcal{G}_1 [(\psi_1, \psi_2), \Phi],$$

where η and η_1 are positive constants to be specified later.

Furthermore, it follows from (5.25), Lemma 5.2, and Lemma 5.3 that

$$(5.49) \quad \begin{aligned} & \frac{d}{dt} \mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] \\ & \leq -(C_2 - C_3\eta - C_4\eta_1) \left(\|\psi_{1t}\|^2 + \|\psi_{2t}\|^2 \right) \\ & \quad - \left(\frac{\sigma}{2} - C_3\eta - C_4\eta_1 \right) \left(\|\tilde{\mathbf{A}}_t\|^2 + \|\nabla\Phi\|^2 \right) \\ & \quad - \left(\frac{\eta_1}{2} - C_3\eta \right) \left(\|\nabla\psi_{1t}\|^2 + \|\nabla\psi_{2t}\|^2 \right) - \frac{\eta_1\sigma}{2} \|\Delta\Phi\|^2 - \frac{\eta_1\varepsilon}{2} \|\nabla\Phi_t\|^2 \\ & \quad - \frac{\eta}{2\varepsilon} \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0}^2. \end{aligned}$$

Taking

$$(5.50) \quad \eta = \frac{\min\{C_2, \sigma\}}{4C_3(1 + 4C_4)}, \quad \eta_1 = 4C_3\eta = \frac{\min\{C_2, \sigma\}}{(1 + 4C_4)},$$

we can deduce that

$$(5.51) \quad \begin{aligned} & \frac{d}{dt} \mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] \\ & \leq -C \left(\|\psi_{1t}\|_{H^1}^2 + \|\psi_{2t}\|_{H^1}^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla\Phi\|^2 + \|\Delta\Phi\|^2 + \|\nabla\Phi_t\|^2 \right. \\ & \quad \left. + \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0}^2 \right) \\ & \leq 0 \end{aligned}$$

for some $C > 0$.

Thanks to Lemma 3.1, Lemma 3.2, and Lemma 3.4, we can find an increasing unbounded sequence $\{t_n\}_{n=1}^\infty$ and a pair $(\psi_\infty, \tilde{\mathbf{A}}_\infty)$ satisfying (2.19)–(2.22) such that

$$(5.52) \quad \lim_{t_n \rightarrow +\infty} \|\tilde{\mathbf{A}}(t_n) - \tilde{\mathbf{A}}_\infty\|_{\mathbf{H}^2(\Omega)} = 0,$$

$$(5.53) \quad \lim_{t_n \rightarrow +\infty} \|\psi(t_n) - \psi_\infty\|_{H^2_c(\Omega)} = 0.$$

Thus, it follows from (5.52), (5.53), Lemma 5.1, and Lemma 5.4 that

$$(5.54) \quad \lim_{t_n \rightarrow +\infty} \mathfrak{F}(t_n) = E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty].$$

Since \mathfrak{F} is decreasing in time (see (5.51)), we have

$$(5.55) \quad \mathfrak{F}(t) \geq E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \quad \forall t \geq 0.$$

After these preparations, we proceed to prove the convergence result following a simple argument introduced in [20]. The key observation is that, after a certain time t_0 , the solution $[(\psi_1(t), \psi_2(t)), \tilde{\mathbf{A}}(t)]$ always satisfies the condition of Lemma 4.2, that is, it falls in a neighborhood of $[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty]$, where (4.9) holds.

We now consider all the possible cases.

(1) If there is a $t_0 > 0$ such that at this time $\mathfrak{F}[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] = E_0[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty]$, then $\forall t > t_0$, we deduce from (5.51) that $[(\psi_1(t), \psi_2(t)), \tilde{\mathbf{A}}(t)]$ is independent of t . We obtain the convergence for ψ and $\tilde{\mathbf{A}}$ from (5.52) and (5.53).

(2) If $\mathfrak{F}[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] > E_0[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \forall t \geq 0$ and there is $t_0 > 0$ such that $\forall t \geq t_0$, $[(\psi_1, \psi_2), \tilde{\mathbf{A}}]$ satisfies the condition of Lemma 4.2, i.e.,

$$(5.56) \quad \left\| [(\psi_1, \psi_2), \tilde{\mathbf{A}}] - [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right\|_{\mathcal{X}} < \beta,$$

then, for the Lojasiewicz exponent $\theta \in (0, \frac{1}{2})$ (see Lemma 4.2), we calculate

$$(5.57) \quad \begin{aligned} & \frac{d}{dt} \left(\mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right)^\theta \\ & = \theta \left(\mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right)^{\theta-1} \frac{d}{dt} \mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi]. \end{aligned}$$

Recalling (5.20) and the definitions of $\mathcal{G}[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi]$ and $\mathcal{G}_1[(\psi_1, \psi_2), \Phi]$, we have

$$(5.58) \quad \begin{aligned} & \left(\mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right)^{1-\theta} \\ & \leq C \left(\left| E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right|^{1-\theta} + \|\tilde{\mathbf{A}}_t\|^{2(1-\theta)} \right. \\ & \quad \left. + \|\nabla\Phi\|^{2(1-\theta)} + \|\Delta\Phi\|^{2(1-\theta)} + \|\psi_{1t}\|^{2(1-\theta)} + \|\psi_{2t}\|^{2(1-\theta)} \right. \\ & \quad \left. + \left\| \text{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \text{curl} \mathbf{H}_{ext} \right\|_{\mathbf{x}'_0} \right)^{2(1-\theta)}. \end{aligned}$$

On the other hand, on account of (4.9), Lemma 4.1 (cf. the explicit form of DE_0), Proposition 2.1, and the Poincaré inequality, we deduce that

$$(5.59) \quad \begin{aligned} & \left| E_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right|^{1-\theta} \leq \left\| DE_0 [(\psi_1, \psi_2), \tilde{\mathbf{A}}] \right\|_{\mathcal{X}'} \\ & \leq C \left(\|\psi_{1t} + \Phi\psi_2\| + \|\psi_{2t} - \Phi\psi_1\| + \left\| \text{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \text{curl} \mathbf{H}_{ext} \right\|_{\mathbf{x}'_0} \right) \\ & \leq C \left(\|\psi_{1t}\| + \|\psi_{2t}\| + \|\nabla\Phi\| + \left\| \text{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \text{curl} \mathbf{H}_{ext} \right\|_{\mathbf{x}'_0} \right). \end{aligned}$$

Since $\theta \in (0, \frac{1}{2})$, thanks to (5.1) and (5.2), it follows that

$$(5.60) \quad \begin{aligned} & \left(\mathfrak{F} [(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi] - E_0 [(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty] \right)^{1-\theta} \\ & \leq C \left(\|\psi_{1t}\| + \|\psi_{2t}\| + \|\tilde{\mathbf{A}}_t\| + \|\nabla\Phi\| + \|\Delta\Phi\| \right. \\ & \quad \left. + \left\| \text{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \text{curl} \mathbf{H}_{ext} \right\|_{\mathbf{x}'_0} \right). \end{aligned}$$

Hence, combining (5.51), (5.57), and (5.60), we infer that, for all $t \geq t_0$,

$$(5.61) \quad \begin{aligned} & \frac{d}{dt} \left(\mathfrak{F} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \right)^\theta + C_5 \left(\|\psi_{1t}\| + \|\psi_{2t}\| + \|\tilde{\mathbf{A}}_t\| \right. \\ & \left. + \|\nabla\Phi\| + \|\Delta\Phi\| + \|\nabla\Phi_t\| + \left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\overline{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{X}'_0} \right) \\ & \leq 0. \end{aligned}$$

As a result, we can deduce that

$$(5.62) \quad \int_{t_0}^t \|\tilde{\mathbf{A}}_t(\tau)\| d\tau < \infty,$$

$$(5.63) \quad \int_{t_0}^t (\|\psi_{1t}(\tau)\| + \|\psi_{2t}(\tau)\|) d\tau < \infty$$

$\forall t \geq t_0$. This yields

$$\lim_{t \rightarrow +\infty} \|\tilde{\mathbf{A}}(t) - \tilde{\mathbf{A}}_\infty\| = 0, \quad \lim_{t \rightarrow +\infty} \|\psi(t) - \psi_\infty\| = 0.$$

From Lemmas 3.1 and 3.2, owing to the uniqueness of the limit, we infer that

$$\lim_{t \rightarrow +\infty} \|\tilde{\mathbf{A}}(t) - \tilde{\mathbf{A}}_\infty\|_{\mathbf{H}^2(\Omega)} = 0, \quad \lim_{t \rightarrow +\infty} \|\psi(t) - \psi_\infty\|_{H^2_2(\Omega)} = 0.$$

Recalling Lemma 5.1, we conclude that $(\psi(t), \tilde{\mathbf{A}}(t), \tilde{\mathbf{A}}_t(t), \Phi(t))$ converges to $(\psi_\infty, \tilde{\mathbf{A}}_\infty, \mathbf{0}, 0)$ in \mathbb{X}^0 as t goes to infinity.

(3) Combining (5.52), (5.53), and Lemma 5.1, we obtain that, for any $\rho \in (0, \beta)$, there exists an integer $N = N(\rho)$ such that, when $n \geq N$,

$$(5.64) \quad \|\psi_1(t_n) - \psi_{1\infty}\|_{H^2(\Omega)} < \frac{\rho}{4}, \quad \|\psi_2(t_n) - \psi_{2\infty}\|_{H^2(\Omega)} < \frac{\rho}{4},$$

$$(5.65) \quad \|\tilde{\mathbf{A}}(t_n) - \tilde{\mathbf{A}}_\infty\|_{\mathbf{H}^2(\Omega)} < \frac{\rho}{4},$$

$$(5.66) \quad \frac{1}{C_5} \left(\mathfrak{F} \left[(\psi_1(t_n), \psi_2(t_n)), \tilde{\mathbf{A}}(t_n), \Phi(t_n) \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \right)^\theta < \frac{\rho}{4}.$$

We now define

$$(5.67) \quad \begin{aligned} \bar{t}_n = \sup_{t > t_n} \left\{ \right. & \|\psi_1(s) - \psi_{1\infty}\|_{H^2(\Omega)} + \|\psi_2(s) - \psi_{2\infty}\|_{H^2(\Omega)} \\ & \left. + \|\tilde{\mathbf{A}}(s) - \tilde{\mathbf{A}}_\infty\|_{\mathbf{H}^2(\Omega)} < \beta, \quad \forall s \in [t_n, t] \right\}, \end{aligned}$$

and we note that it follows from (5.64), (5.65), and the continuity of the orbit in \mathbb{X}^0 that $\bar{t}_n > t_n \forall n \geq N$. Then there are two possibilities, namely,

(i) There exists $n_0 \geq N$ such that $\bar{t}_{n_0} = +\infty$, then arguing as in (1) and (2), the convergence result can be obtained.

(ii) Otherwise, for all $n \geq N$, we have $t_n < \bar{t}_n < +\infty$, and, for all $t \in [t_n, \bar{t}_n]$, we have

$$\mathfrak{F} \left[(\psi_1(t), \psi_2(t)), \tilde{\mathbf{A}}(t), \Phi(t) \right] > E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right].$$

Hence, we deduce from (5.61) that

$$\begin{aligned}
 & \int_{t_n}^{\bar{t}_n} \left(\|\psi_{1t}(\tau)\| + \|\psi_{2t}(\tau)\| + \|\tilde{\mathbf{A}}_t(\tau)\| \right) d\tau \\
 & \leq \frac{1}{C_5} \left(\mathfrak{F} \left[(\psi_1(t_n), \psi_2(t_n)), \tilde{\mathbf{A}}(t_n), \Phi(t_n) \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \right)^\theta \\
 (5.68) \quad & < \frac{\rho}{4}.
 \end{aligned}$$

Noticing that $\rho > 0$ is arbitrary in the above argument, we are able to conclude that for any $\rho > 0$, there exist $N = N(\rho) \in \mathbb{N}$ such that, for all $n \geq N$, there holds

$$\begin{aligned}
 & \|\psi_1(\bar{t}_n) - \psi_{1\infty}\| + \|\psi_2(\bar{t}_n) - \psi_{2\infty}\| + \|\tilde{\mathbf{A}}(\bar{t}_n) - \tilde{\mathbf{A}}_\infty\| \\
 & \leq \|\psi_1(t_n) - \psi_{1\infty}\| + \|\psi_2(t_n) - \psi_{2\infty}\| + \|\tilde{\mathbf{A}}(t_n) - \tilde{\mathbf{A}}_\infty\| \\
 & \quad + \int_{t_n}^{\bar{t}_n} \left(\|\psi_{1t}(\tau)\| + \|\psi_{2t}(\tau)\| + \|\tilde{\mathbf{A}}_t(\tau)\| \right) d\tau \\
 (5.69) \quad & < \rho.
 \end{aligned}$$

This implies

$$(5.70) \quad \lim_{n \rightarrow +\infty} \left(\|\psi_1(\bar{t}_n) - \psi_{1\infty}\| + \|\psi_2(\bar{t}_n) - \psi_{2\infty}\| + \|\tilde{\mathbf{A}}(\bar{t}_n) - \tilde{\mathbf{A}}_\infty\| \right) = 0.$$

Therefore, due to Lemmas 3.1 and 3.2, there exists a subsequence of $\{\bar{t}_n\}$, denoted by $\{\bar{t}_{n_k}\}$ such that

$$(5.71) \quad \lim_{n_k \rightarrow +\infty} \left(\|\psi_1(\bar{t}_{n_k}) - \psi_{1\infty}\|_{H^2(\Omega)} + \|\psi_2(\bar{t}_{n_k}) - \psi_{2\infty}\|_{H^2(\Omega)} + \|\tilde{\mathbf{A}}(\bar{t}_{n_k}) - \tilde{\mathbf{A}}_\infty\|_{\mathbf{H}^2(\Omega)} \right) = 0,$$

which contradicts the definition of \bar{t}_n .

Summing up, we conclude that

$$(5.72) \quad \lim_{t \rightarrow +\infty} \|z(t) - z_\infty\|_0 = 0.$$

6. Convergence rate. In this section we demonstrate estimate (2.25) on convergence rate. Let us begin by noting that a combination of (5.60) and (5.61) implies, for all $t \geq t_0$, that

$$\begin{aligned}
 & \frac{d}{dt} \left(\mathfrak{F} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \right) \\
 (6.1) \quad & + C \left(\mathfrak{F} \left[(\psi_1, \psi_2), \tilde{\mathbf{A}}, \Phi \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \right)^{2(1-\theta)} \leq 0.
 \end{aligned}$$

Keeping in mind that $\theta \in (0, \frac{1}{2})$, we deduce

$$(6.2) \quad \mathfrak{F} \left[(\psi_1(t), \psi_2(t)), \tilde{\mathbf{A}}(t), \Phi(t) \right] - E_0 \left[(\psi_{1\infty}, \psi_{2\infty}), \tilde{\mathbf{A}}_\infty \right] \leq C(1+t)^{-1/(1-2\theta)},$$

for all $t \geq 0$ and for a suitable $C > 0$, because for $t \in [0, t_0]$, the term on the left-hand side is bounded by a constant depending only on initial data (cf. (2.12)). Hence, it

follows from (5.61) that

$$\begin{aligned}
& \int_t^\infty \left(\|\psi_{1t}\| + \|\psi_{2t}\| + \|\nabla\Phi\| + \|\nabla\Phi_t\| + \|\tilde{\mathbf{A}}_t\| \right) d\tau \\
& + \int_t^\infty \left(\left\| \operatorname{curl}^2 \tilde{\mathbf{A}} + \frac{i}{2} (\bar{\psi} D_{\tilde{\mathbf{A}}} \psi - \psi \overline{D_{\tilde{\mathbf{A}}} \psi}) + \operatorname{curl} \mathbf{H}_{ext} \right\|_{\mathbf{x}'_0} \right) d\tau \\
(6.3) \quad & \leq C(1+t)^{-\theta/(1-2\theta)}, \quad \forall t \geq 0,
\end{aligned}$$

and, as a consequence (cf. (1.2)),

$$(6.4) \quad \int_t^\infty \|\tilde{\mathbf{A}}_{tt}(\tau)\|_{\mathbf{x}'_0} d\tau \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.$$

Therefore, from (6.3) and (6.4), we deduce the following preliminary convergence rate estimate:

$$\begin{aligned}
& \|\psi(t) - \psi_\infty\| + \|\tilde{\mathbf{A}}(t) - \tilde{\mathbf{A}}_\infty\| + \|\tilde{\mathbf{A}}_t(t)\|_{\mathbf{x}'_0} + \|\nabla\Phi\| \\
& \leq \int_t^\infty \left(\|\psi_t(\tau)\| + \|\tilde{\mathbf{A}}_t(\tau)\| + \|\tilde{\mathbf{A}}_{tt}(\tau)\|_{\mathbf{x}'_0} + \|\nabla\Phi_t(\tau)\| \right) d\tau \\
(6.5) \quad & \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.
\end{aligned}$$

In order to get higher order estimates, we first subtract the stationary equations from the time-dependent equations. This gives

$$\begin{aligned}
& \psi_t - i\Phi\psi - \Delta(\psi - \psi_\infty) + 2i \left(\tilde{\mathbf{A}} \cdot \nabla\psi - \tilde{\mathbf{A}}_\infty \cdot \nabla\psi_\infty \right) + \tilde{\mathbf{A}}^2\psi - \tilde{\mathbf{A}}_\infty^2\psi_\infty \\
(6.6) \quad & -\lambda^2(1-|\psi|^2)\psi + \lambda^2(1-|\psi_\infty|^2)\psi_\infty = 0,
\end{aligned}$$

$$\begin{aligned}
& \varepsilon \left(\tilde{\mathbf{A}}_t - \nabla\Phi \right)_t + \sigma \left(\tilde{\mathbf{A}}_t - \nabla\Phi \right) + \operatorname{curl}^2 \tilde{\mathbf{A}} - \operatorname{curl}^2 \tilde{\mathbf{A}}_\infty + |\psi|^2 \tilde{\mathbf{A}} - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty \\
(6.7) \quad & + \frac{i}{2} (\bar{\psi} \nabla\psi - \psi \nabla\bar{\psi}) - \frac{i}{2} (\bar{\psi}_\infty \nabla\psi_\infty - \psi_\infty \nabla\bar{\psi}_\infty) = 0,
\end{aligned}$$

$$(6.8) \quad -\varepsilon \Delta\Phi_t - \sigma \Delta\Phi + \frac{i}{2} (\bar{\psi} \psi_t - \psi \bar{\psi}_t) + |\psi|^2 \Phi = 0,$$

in $\Omega \times (0, \infty)$, with boundary conditions

$$(6.9) \quad \partial_{\mathbf{n}}(\psi - \psi_\infty) = 0, \quad \left(\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \right) \cdot \mathbf{n} = 0, \quad \left(\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty \right) \times \mathbf{n} = 0, \quad \partial_{\mathbf{n}}\Phi = 0,$$

on $\Gamma \times (0, \infty)$.

Multiplying (6.6) by $\bar{\psi}_t$ and its conjugate by ψ_t , respectively, then integrating on Ω and adding the results together, we have

$$\begin{aligned}
& \frac{1}{2} \frac{d}{dt} \|\nabla\psi - \nabla\psi_\infty\|^2 + \|\psi_t\|^2 + \frac{i}{2} \langle \Phi, \bar{\psi} \psi_t - \psi \bar{\psi}_t \rangle \\
& + i \langle \tilde{\mathbf{A}} \cdot \nabla\psi - \tilde{\mathbf{A}}_\infty \cdot \nabla\psi_\infty, \bar{\psi}_t \rangle - i \langle \tilde{\mathbf{A}} \cdot \nabla\bar{\psi} - \tilde{\mathbf{A}}_\infty \cdot \nabla\bar{\psi}_\infty, \psi_t \rangle \\
& + \frac{1}{2} \langle \tilde{\mathbf{A}}^2\psi - \tilde{\mathbf{A}}_\infty^2\psi_\infty, \bar{\psi}_t \rangle + \frac{1}{2} \langle \tilde{\mathbf{A}}^2\bar{\psi} - \tilde{\mathbf{A}}_\infty^2\bar{\psi}_\infty, \psi_t \rangle \\
& - \frac{\lambda^2}{2} \int_\Omega [(1-|\psi|^2)\psi - (1-|\psi_\infty|^2)\psi_\infty] \bar{\psi}_t dx \\
(6.10) \quad & - \frac{\lambda^2}{2} \int_\Omega [(1-|\psi|^2)\bar{\psi} - (1-|\psi_\infty|^2)\bar{\psi}_\infty] \psi_t dx = 0.
\end{aligned}$$

It follows from inequalities (2.7), (2.8), and bound (2.12) that

$$\begin{aligned}
 I_4 &:= \left| i \left\langle \tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi_\infty, \bar{\psi}_t \right\rangle \right| \\
 &\leq \left| \left\langle \tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi, \bar{\psi}_t \right\rangle \right| + \left| \left\langle \tilde{\mathbf{A}}_\infty \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi_\infty, \bar{\psi}_t \right\rangle \right| \\
 &\leq \|\nabla \psi\|_{\mathbf{L}^4} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|_{\mathbf{L}^4} \|\psi_t\| + \|\tilde{\mathbf{A}}_\infty\|_{\mathbf{L}^\infty} \|\nabla \psi - \nabla \psi_\infty\| \|\psi_t\| \\
 &\leq C \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^{\frac{1}{2}} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\|^{\frac{1}{2}} \|\psi_t\| + C \|\nabla \psi - \nabla \psi_\infty\| \|\psi_t\| \\
 &\leq \frac{1}{16} \|\psi_t\|^2 + \frac{\alpha}{8} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\|^2 + C \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 \\
 (6.11) \quad &+ C_6 \|\nabla \psi - \nabla \psi_\infty\|^2.
 \end{aligned}$$

Similarly, we deduce

$$\begin{aligned}
 I_5 &:= \left| -i \left\langle \tilde{\mathbf{A}} \cdot \nabla \bar{\psi} - \tilde{\mathbf{A}}_\infty \cdot \nabla \bar{\psi}_\infty, \psi_t \right\rangle \right| \\
 &\leq \frac{1}{16} \|\psi_t\|^2 + \frac{\alpha}{8} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\|^2 + C \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 \\
 (6.12) \quad &+ C_6 \|\nabla \psi - \nabla \psi_\infty\|^2.
 \end{aligned}$$

Next, we have

$$\begin{aligned}
 I_6 &:= \frac{1}{2} \left| \left\langle \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_\infty^2 \psi_\infty, \bar{\psi}_t \right\rangle \right| \\
 &\leq \frac{1}{2} \left| \left\langle \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}^2 \psi_\infty, \bar{\psi}_t \right\rangle \right| + \frac{1}{2} \left| \left\langle \tilde{\mathbf{A}}^2 \psi_\infty - \tilde{\mathbf{A}}_\infty^2 \psi_\infty, \bar{\psi}_t \right\rangle \right| \\
 (6.13) \quad &\leq \frac{1}{16} \|\psi_t\|^2 + C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \|\psi - \psi_\infty\|^2 \right).
 \end{aligned}$$

In the same manner, we get

$$(6.14) \quad I_7 := \frac{1}{2} \left| \left\langle \tilde{\mathbf{A}}^2 \bar{\psi} - \tilde{\mathbf{A}}_\infty^2 \bar{\psi}_\infty, \psi_t \right\rangle \right| \leq \frac{1}{16} \|\psi_t\|^2 + C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \|\psi - \psi_\infty\|^2 \right).$$

Moreover, we have

$$\begin{aligned}
 I_8 &:= \frac{\lambda^2}{2} \left| \int_\Omega [(1 - |\psi|^2) \psi - (1 - |\psi_\infty|^2) \psi_\infty] \bar{\psi}_t dx \right| \\
 &\leq \frac{\lambda^2}{2} \left| \int_\Omega |\psi|^2 (\psi - \psi_\infty) \bar{\psi}_t dx \right| + \frac{\lambda^2}{2} \left| \int_\Omega \psi_\infty (|\psi|^2 - |\psi_\infty|^2) \bar{\psi}_t dx \right| \\
 &\quad + \frac{\lambda^2}{2} \left| \int_\Omega (\psi - \psi_\infty) \bar{\psi}_t dx \right| \\
 &\leq \frac{1}{32} \|\bar{\psi}_t\|^2 + C \|\psi - \psi_\infty\|^2 + C \|\psi_1^2 + \psi_2^2 - \psi_{1\infty}^2 - \psi_{2\infty}^2\| \|\bar{\psi}_t\| \\
 &\leq \frac{1}{32} \|\bar{\psi}_t\|^2 + C \|\psi - \psi_\infty\|^2 + C \|\bar{\psi}_t\| (\|\psi_1 - \psi_{1\infty}\| + \|\psi_2 - \psi_{2\infty}\|) \\
 &\leq \frac{1}{32} \|\bar{\psi}_t\|^2 + C \|\psi - \psi_\infty\|^2 + C \|\bar{\psi}_t\| \|\psi - \psi_\infty\| \\
 (6.15) \quad &\leq \frac{1}{16} \|\psi_t\|^2 + C \|\psi - \psi_\infty\|^2,
 \end{aligned}$$

(6.16)

$$I_9 := \frac{\lambda^2}{2} \left| \int_{\Omega} [(1 - |\psi|^2) \bar{\psi} - (1 - |\psi_{\infty}|^2) \bar{\psi}_{\infty}] \psi_t dx \right| \leq \frac{1}{16} \|\psi_t\|^2 + C \|\psi - \psi_{\infty}\|^2.$$

Multiplying (6.6) by $\bar{\psi} - \bar{\psi}_{\infty}$ and its conjugate by $\psi - \psi_{\infty}$, respectively, then integrating on Ω and adding the results together, we obtain

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \|\psi - \psi_{\infty}\|^2 + \|\nabla \psi - \nabla \psi_{\infty}\|^2 + \frac{i}{2} \langle \Phi, \psi \bar{\psi}_{\infty} - \bar{\psi} \psi_{\infty} \rangle \\ & + i \langle \tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle - i \langle \tilde{\mathbf{A}} \cdot \nabla \bar{\psi} - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \bar{\psi}_{\infty}, \psi - \psi_{\infty} \rangle \\ & + \frac{1}{2} \langle \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_{\infty}^2 \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle + \frac{1}{2} \langle \tilde{\mathbf{A}}^2 \bar{\psi} - \tilde{\mathbf{A}}_{\infty}^2 \bar{\psi}_{\infty}, \psi - \psi_{\infty} \rangle \\ & + \frac{1}{2} \langle -\lambda^2 (1 - |\psi|^2) \psi + \lambda^2 (1 - |\psi_{\infty}|^2) \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \\ (6.17) \quad & + \frac{1}{2} \langle -\lambda^2 (1 - |\psi|^2) \bar{\psi} + \lambda^2 (1 - |\psi_{\infty}|^2) \bar{\psi}_{\infty}, \psi - \psi_{\infty} \rangle = 0. \end{aligned}$$

Then, arguing as before, we infer the following estimates:

$$\begin{aligned} I_{10} & := \frac{1}{2} |i \langle \Phi, \psi \bar{\psi}_{\infty} - \bar{\psi} \psi_{\infty} \rangle| \\ & \leq \frac{1}{2} |\langle \Phi, \psi \bar{\psi}_{\infty} - \psi_{\infty} \bar{\psi}_{\infty} \rangle| + \frac{1}{2} |\langle \Phi, \psi_{\infty} \bar{\psi}_{\infty} - \psi_{\infty} \bar{\psi} \rangle| \\ & \leq \|\Phi\|_{L^4} \|\bar{\psi}_{\infty}\|_{L^4} \|\psi - \psi_{\infty}\| + \|\Phi\|_{L^4} \|\psi_{\infty}\|_{L^4} \|\bar{\psi} - \bar{\psi}_{\infty}\| \\ (6.18) \quad & \leq \|\nabla \Phi\|^2 + C \|\psi - \psi_{\infty}\|^2, \end{aligned}$$

$$\begin{aligned} I_{11} & := \left| i \langle \tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| \\ & \leq \left| \langle \tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \psi, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| + \left| \langle \tilde{\mathbf{A}}_{\infty} \cdot \nabla \psi - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| \\ & \leq \|\nabla \psi\|_{L^4} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_{\infty}\|_{L^4} \|\psi - \psi_{\infty}\| + \|\tilde{\mathbf{A}}_{\infty}\|_{L^{\infty}} \|\nabla \psi - \nabla \psi_{\infty}\| \|\psi - \psi_{\infty}\| \\ (6.19) \quad & \leq \frac{\epsilon \alpha}{16} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_{\infty}\|^2 + \frac{1}{16} \|\nabla \psi - \nabla \psi_{\infty}\|^2 + C \|\psi - \psi_{\infty}\|^2, \end{aligned}$$

$$\begin{aligned} I_{12} & := \left| i \langle \tilde{\mathbf{A}} \cdot \nabla \bar{\psi} - \tilde{\mathbf{A}}_{\infty} \cdot \nabla \bar{\psi}_{\infty}, \psi - \psi_{\infty} \rangle \right| \\ (6.20) \quad & \leq \frac{\epsilon \alpha}{16} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_{\infty}\|^2 + \frac{1}{16} \|\nabla \psi - \nabla \psi_{\infty}\|^2 + C \|\psi - \psi_{\infty}\|^2, \end{aligned}$$

$$\begin{aligned} I_{13} & := \left| \frac{1}{2} \langle \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_{\infty}^2 \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| \\ & \leq \left| \frac{1}{2} \langle \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_{\infty}^2 \psi, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| + \left| \frac{1}{2} \langle \tilde{\mathbf{A}}_{\infty}^2 \psi - \tilde{\mathbf{A}}_{\infty}^2 \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| \\ & \leq C \|\psi\|_{L^{\infty}} \|\tilde{\mathbf{A}} + \tilde{\mathbf{A}}_{\infty}\|_{L^{\infty}} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_{\infty}\| \|\psi - \psi_{\infty}\| + C \|\tilde{\mathbf{A}}_{\infty}^2\|_{L^{\infty}} \|\psi - \psi_{\infty}\|^2 \\ (6.21) \quad & \leq C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_{\infty}\|^2 + \|\psi - \psi_{\infty}\|^2 \right), \end{aligned}$$

$$(6.22) \quad I_{14} := \left| \frac{1}{2} \langle \tilde{\mathbf{A}}^2 \bar{\psi} - \tilde{\mathbf{A}}_{\infty}^2 \bar{\psi}_{\infty}, \psi - \psi_{\infty} \rangle \right| \leq C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_{\infty}\|^2 + \|\psi - \psi_{\infty}\|^2 \right),$$

$$(6.23) \quad I_{15} := \left| \frac{1}{2} \langle -\lambda^2 (1 - |\psi|^2) \psi + \lambda^2 (1 - |\psi_{\infty}|^2) \psi_{\infty}, \bar{\psi} - \bar{\psi}_{\infty} \rangle \right| \leq C \|\psi - \psi_{\infty}\|^2,$$

$$(6.24) \quad I_{16} := \left| \frac{1}{2} \langle -\lambda^2 (1 - |\psi|^2) \bar{\psi} + \lambda^2 (1 - |\psi_\infty|^2) \bar{\psi}_\infty, \psi - \psi_\infty \rangle \right| \leq C \|\psi - \psi_\infty\|^2.$$

In the above estimates, $\alpha > 0$ and $\epsilon > 0$ are constants to be chosen later.

We now consider (6.7) and take the inner product with $\tilde{\mathbf{A}}_t + \alpha(\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty)$ in $\mathbf{L}^2(\Omega)$. We thus have

$$(6.25) \quad \begin{aligned} & \frac{d}{dt} \left(\frac{1}{2} \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\|^2 + \frac{\epsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 + \alpha \epsilon \langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right. \\ & \quad \left. + \frac{\alpha \sigma}{2} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 \right) \\ & + \alpha \|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\|^2 + (\sigma - \alpha \epsilon) \|\tilde{\mathbf{A}}_t\|^2 \\ & + \langle \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle - \langle \tilde{\mathbf{A}}_\infty \cdot \tilde{\mathbf{A}}_t, |\psi_\infty|^2 \rangle \\ & + \frac{i}{2} \langle \bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}, \tilde{\mathbf{A}}_t \rangle - \frac{i}{2} \langle \bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}}_t \rangle \\ & + \alpha \langle |\psi|^2 \tilde{\mathbf{A}} - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \\ & + \frac{\alpha i}{2} \langle \bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle - \frac{\alpha i}{2} \langle \bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle = 0. \end{aligned}$$

Let us estimate some of the terms in (6.25). We have

$$(6.26) \quad \begin{aligned} I_{17} & := \left| \langle \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle - \langle \tilde{\mathbf{A}}_\infty \cdot \tilde{\mathbf{A}}_t, |\psi_\infty|^2 \rangle \right| \\ & \leq \left| \langle \tilde{\mathbf{A}} \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle - \langle \tilde{\mathbf{A}}_\infty \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle \right| \\ & \quad + \left| \langle \tilde{\mathbf{A}}_\infty \cdot \tilde{\mathbf{A}}_t, |\psi|^2 \rangle - \langle \tilde{\mathbf{A}}_\infty \cdot \tilde{\mathbf{A}}_t, |\psi_\infty|^2 \rangle \right| \\ & \leq \frac{\sigma}{16} \|\tilde{\mathbf{A}}_t\|^2 + C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \|\psi - \psi_\infty\|^2 \right), \end{aligned}$$

$$(6.27) \quad \begin{aligned} I_{18} & := \left| \frac{i}{2} \langle \bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}, \tilde{\mathbf{A}}_t \rangle - \frac{i}{2} \langle \bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}}_t \rangle \right| \\ & \leq \left| \langle \bar{\psi} \nabla \psi - \bar{\psi}_\infty \nabla \psi, \tilde{\mathbf{A}}_t \rangle \right| + \left| \langle \bar{\psi}_\infty \nabla \psi - \bar{\psi}_\infty \nabla \psi_\infty, \tilde{\mathbf{A}}_t \rangle \right| \\ & \quad + \left| \langle \psi \nabla \bar{\psi} - \psi_\infty \nabla \bar{\psi}, \tilde{\mathbf{A}}_t \rangle \right| + \left| \langle \psi_\infty \nabla \bar{\psi} - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}}_t \rangle \right| \\ & \leq \frac{\sigma}{16} \|\tilde{\mathbf{A}}_t\|^2 + C_7 \|\nabla \psi - \nabla \psi_\infty\|^2 + C \|\psi - \psi_\infty\|^2, \end{aligned}$$

$$(6.28) \quad \begin{aligned} I_{19} & := \left| \alpha \langle |\psi|^2 \tilde{\mathbf{A}} - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| \\ & \leq \alpha \left| \langle |\psi|^2 \tilde{\mathbf{A}} - |\psi|^2 \tilde{\mathbf{A}}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| + \alpha \left| \langle |\psi|^2 \tilde{\mathbf{A}}_\infty - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| \\ & \leq C \alpha \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \|\psi - \psi_\infty\|^2 \right), \end{aligned}$$

$$(6.29) \quad \begin{aligned} I_{20} & := \left| \frac{\alpha i}{2} \langle \bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle - \frac{\alpha i}{2} \langle \bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| \\ & \leq \alpha \left| \langle \bar{\psi} \nabla \psi - \bar{\psi}_\infty \nabla \psi, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| + \alpha \left| \langle \bar{\psi}_\infty \nabla \psi - \bar{\psi}_\infty \nabla \psi_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| \\ & \quad + \alpha \left| \langle \psi \nabla \bar{\psi} - \psi_\infty \nabla \bar{\psi}, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| + \alpha \left| \langle \psi_\infty \nabla \bar{\psi} - \psi_\infty \nabla \bar{\psi}_\infty, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle \right| \\ & \leq \frac{\alpha}{16} \|\nabla \psi - \nabla \psi_\infty\|^2 + C \alpha \left(\|\psi - \psi_\infty\|^2 + \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 \right). \end{aligned}$$

On the other hand, recalling (3.14) and (6.10), we observe that

$$(6.30) \quad I_{21} := |i \langle \Phi, \bar{\psi}\psi_t - \psi\bar{\psi}_t \rangle| \leq \frac{1}{16} \|\psi_t\|^2 + C \|\nabla\Phi\|^2.$$

Thus, adding (6.10), the product of (6.17) by constant $M > 0$, (6.25) and (3.14) together, using the above estimates of I_j ($j = 4, \dots, 21$), and taking

$$(6.31) \quad \alpha \in \left(0, \frac{\sigma}{2\varepsilon}\right), \quad M \geq \frac{\sigma}{\varepsilon} + 4C_6 + 2C_7, \quad \epsilon \in \left(0, \frac{3}{M}\right),$$

we can find two positive constants C and γ such that

$$(6.32) \quad \begin{aligned} & \frac{d}{dt} \left(\frac{1}{2} \|\nabla\psi - \nabla\psi_\infty\|^2 + \frac{M}{2} \|\psi - \psi_\infty\|^2 + \frac{1}{2} \|\operatorname{curl}\tilde{\mathbf{A}} - \operatorname{curl}\tilde{\mathbf{A}}_\infty\|^2 \right. \\ & \quad \left. + \frac{\varepsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 + \alpha\varepsilon \langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle + \frac{\alpha\sigma}{2} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \frac{\varepsilon}{2} \|\nabla\Phi\|^2 \right) \\ & \quad + \gamma \left(\|\psi_t\|^2 + \|\nabla\psi - \nabla\psi_\infty\|^2 + \|\operatorname{curl}\tilde{\mathbf{A}} - \operatorname{curl}\tilde{\mathbf{A}}_\infty\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla\Phi\|^2 \right) \\ & \leq C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \|\psi - \psi_\infty\|^2 + \|\nabla\Phi\|^2 \right). \end{aligned}$$

Let us now set

$$(6.33) \quad \begin{aligned} y := & \frac{1}{2} \|\nabla\psi - \nabla\psi_\infty\|^2 + \frac{M}{2} \|\psi - \psi_\infty\|^2 + \frac{1}{2} \|\operatorname{curl}\tilde{\mathbf{A}} - \operatorname{curl}\tilde{\mathbf{A}}_\infty\|^2 \\ & + \frac{\varepsilon}{2} \|\tilde{\mathbf{A}}_t\|^2 + \alpha\varepsilon \langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty \rangle + \frac{\alpha\sigma}{2} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|^2 + \frac{\varepsilon}{2} \|\nabla\Phi\|^2. \end{aligned}$$

Choosing α small enough, by the Cauchy–Schwarz inequality, we get

$$(6.34) \quad c_1 y \leq \|\psi_t\|^2 + \|\nabla\psi - \nabla\psi_\infty\|^2 + \|\operatorname{curl}\tilde{\mathbf{A}} - \operatorname{curl}\tilde{\mathbf{A}}_\infty\|^2 + \|\tilde{\mathbf{A}}_t\|^2 + \|\nabla\Phi\|^2 \leq c_2 y.$$

Therefore, recalling (6.5), we have

$$(6.35) \quad \frac{d}{dt} y(t) + \hat{\gamma} y(t) \leq C(1+t)^{-2\theta/(1-2\theta)} \quad \forall t \geq 0$$

for some constant $\hat{\gamma} > 0$. Thus, it follows that (see also [14, 34, 35])

$$(6.36) \quad y(t) \leq C(1+t)^{-2\theta/(1-2\theta)} \quad \forall t \geq 0.$$

Furthermore, on account of (6.34), we obtain the estimate of convergence rate in \mathbb{X}^{-1} -norm (cf. (2.5)), namely,

$$(6.37) \quad \|z(t) - z_\infty\|_{-1} \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.$$

To get (2.25), a further step is needed. We adapt a higher order estimate derived in [3, section 3]. Let us set

$$(6.38) \quad \begin{aligned} y_1 := & \alpha_1 \|\psi_t\|^2 + \|\operatorname{curl}\tilde{\mathbf{A}}_t\|^2 + \varepsilon \|\tilde{\mathbf{A}}_{tt}\|^2 + 2\alpha_2 \varepsilon \langle \tilde{\mathbf{A}}_t, \tilde{\mathbf{A}}_{tt} \rangle + \alpha_2 \sigma \|\tilde{\mathbf{A}}_t\|^2 \\ & + \varepsilon \|\Delta\Phi\|^2 + \sigma \|\nabla\Phi\|^2 \end{aligned}$$

and choose $\alpha_1, \alpha_2 > 0$ such that (cf. also (2.9))

$$(6.39) \quad \alpha_1 \geq \max \left\{ \frac{\sigma}{\varepsilon}, \frac{1}{\sigma\kappa} \right\} + \frac{24}{\sigma}, \quad 0 < \alpha_2 \leq \min \left\{ \frac{\sigma}{4\varepsilon}, \frac{1}{4\sigma\kappa} \right\}.$$

Then we observe that (see [3, equations (3.7) and (3.8)])

$$(6.40) \quad \frac{1}{2} \left(\alpha_1 \|\psi_t\|^2 + \|\operatorname{curl} \tilde{\mathbf{A}}_t\|^2 + \varepsilon \|\tilde{\mathbf{A}}_{tt}\|^2 + \varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2 \right) \leq y_1,$$

$$(6.41) \quad y_1 \leq \frac{3}{2} \left(\alpha_1 \|\psi_t\|^2 + \|\operatorname{curl} \tilde{\mathbf{A}}_t\|^2 + \varepsilon \|\tilde{\mathbf{A}}_{tt}\|^2 + \varepsilon \|\Delta \Phi\|^2 + \sigma \|\nabla \Phi\|^2 \right).$$

Furthermore, a refinement of the argument used in [3, section 3, p. 628] yields

$$(6.42) \quad \frac{d}{dt} y_1 + \gamma_1 y_1 + \frac{\alpha_1}{2} \|\nabla \psi_t\|^2 + \varepsilon \|\nabla \Phi_t\|^2 \leq C \left(\|\psi_t\|^2 + \|\nabla \Phi\|^2 + \|\tilde{\mathbf{A}}_t\|^2 \right),$$

for some $\gamma_1 > 0$, provided that α_1 and α_2 satisfy (6.39).

Hence, on account of (6.37), it follows that

$$(6.43) \quad \frac{d}{dt} y_1(t) + \gamma_1 y_1(t) \leq C(1+t)^{-2\theta/(1-2\theta)},$$

which gives, using (6.40),

$$(6.44) \quad \|\psi_t(t)\| + \|\operatorname{curl} \tilde{\mathbf{A}}_t(t)\| + \|\tilde{\mathbf{A}}_{tt}(t)\| + \|\Delta \Phi(t)\| \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.$$

On the other hand, we infer from (6.6) that

$$(6.45) \quad \begin{aligned} \|\Delta(\psi - \psi_\infty)\| &\leq \|\psi_t\| + C\|\Phi\| + C\|\tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi_\infty\| + \left\| \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_\infty^2 \psi_\infty \right\| \\ &\quad + \left\| -\lambda^2 (1 - |\psi|^2) \psi + \lambda^2 (1 - |\psi_\infty|^2) \psi_\infty \right\|, \end{aligned}$$

$$(6.46) \quad \begin{aligned} \|\tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi_\infty\| &\leq \|\tilde{\mathbf{A}} \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi\| + \|\tilde{\mathbf{A}}_\infty \cdot \nabla \psi - \tilde{\mathbf{A}}_\infty \cdot \nabla \psi_\infty\| \\ &\leq \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\|_{\mathbf{L}^4} \|\nabla \psi\|_{\mathbf{L}^4} + \|\tilde{\mathbf{A}}_\infty\|_{\mathbf{L}^\infty} \|\nabla \psi - \nabla \psi_\infty\| \\ &\leq C \left(\|\operatorname{curl} \tilde{\mathbf{A}} - \operatorname{curl} \tilde{\mathbf{A}}_\infty\| + \|\nabla \psi - \nabla \psi_\infty\| \right), \end{aligned}$$

$$(6.47) \quad \begin{aligned} \left\| \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}_\infty^2 \psi_\infty \right\| &\leq \left\| \tilde{\mathbf{A}}^2 \psi - \tilde{\mathbf{A}}^2 \psi_\infty \right\| + \left\| \tilde{\mathbf{A}}^2 \psi_\infty - \tilde{\mathbf{A}}_\infty^2 \psi_\infty \right\| \\ &\leq \left\| \tilde{\mathbf{A}}^2 \right\|_{\mathbf{L}^\infty} \|\psi - \psi_\infty\| + \|\tilde{\mathbf{A}} + \tilde{\mathbf{A}}_\infty\|_{\mathbf{L}^\infty} \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\| \|\psi_\infty\|_{L^\infty} \\ &\leq C \left(\|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\| + \|\psi - \psi_\infty\| \right), \end{aligned}$$

$$(6.48) \quad \left\| -\lambda^2 (1 - |\psi|^2) \psi + \lambda^2 (1 - |\psi_\infty|^2) \psi_\infty \right\| \leq C \|\psi - \psi_\infty\|.$$

Therefore, due to (6.37), we deduce

$$(6.49) \quad \|\Delta(\psi(t) - \psi_\infty)\| \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.$$

Similarly, from (6.7), we obtain

$$(6.50) \quad \begin{aligned} &\|\operatorname{curl}^2 \tilde{\mathbf{A}} - \operatorname{curl}^2 \tilde{\mathbf{A}}_\infty\| \\ &\leq \varepsilon \|\tilde{\mathbf{A}}_{tt}\| + \varepsilon \|\nabla \Phi_t\| + \sigma \|\tilde{\mathbf{A}}_t\| + \sigma \|\nabla \Phi\| + \left\| |\psi|^2 \tilde{\mathbf{A}} - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty \right\| \\ &\quad + \left\| \frac{i}{2} (\bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}) - \frac{i}{2} (\bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty) \right\|. \end{aligned}$$

Besides, the following estimates hold:

$$(6.51) \quad \left\| |\psi|^2 \tilde{\mathbf{A}} - |\psi_\infty|^2 \tilde{\mathbf{A}}_\infty \right\| \leq C \left(\|\psi - \psi_\infty\| + \|\tilde{\mathbf{A}} - \tilde{\mathbf{A}}_\infty\| \right),$$

$$(6.52) \quad \left\| \frac{i}{2} (\bar{\psi} \nabla \psi - \psi \nabla \bar{\psi}) - \frac{i}{2} (\bar{\psi}_\infty \nabla \psi_\infty - \psi_\infty \nabla \bar{\psi}_\infty) \right\| \leq C (\|\psi - \psi_\infty\| + \|\nabla \psi - \nabla \psi_\infty\|).$$

As a result, owing to (6.37), we infer from (5.7) and (6.50)–(6.52) that

$$(6.53) \quad \left\| \operatorname{curl}^2 \tilde{\mathbf{A}}(t) - \operatorname{curl}^2 \tilde{\mathbf{A}}_\infty \right\| \leq C(1+t)^{-\theta/(1-2\theta)} \quad \forall t \geq 0.$$

Finally, (2.25) follows from (6.37), (6.44), (6.49), and (6.53).

The proof of Theorem 2.1 is now complete.

Acknowledgments. This paper was conceived on occasion of a first author's visit to the Institute of Mathematics in Fudan University, whose hospitality is gratefully acknowledged.

REFERENCES

- [1] T. AKIYAMA AND Y. SHIBATA, *On an L^p approach to the stationary and nonstationary problems of the Ginzburg–Landau–Maxwell equations*, J. Differential Equations, 243 (2007), pp. 1–23.
- [2] J. BERGER, *Time-dependent Ginzburg–Landau equations with charged boundaries*, J. Math. Phys., 46 (2005), p. 095106.
- [3] V. BERTI AND S. GATTI, *Parabolic-hyperbolic time-dependent Ginzburg–Landau–Maxwell equations*, Quart. Appl. Math., 64 (2006), pp. 617–639.
- [4] T. CAZENAVE AND A. HARAUX, *An Introduction to Semilinear Evolution Equations*, Oxford University Press, New York, 1998.
- [5] S.J. CHAPMAN, Q. DU, M.D. GUNZBURGER, AND J.S. PETERSON, *Simplified G – L models for superconductivity valid for high kappa and high fields*, Adv. Math. Sci. Appl., 5 (1995), pp. 193–218.
- [6] K. DEIMLING, *Nonlinear functional analysis*, Springer, New York, 1985.
- [7] Q. DU, *Global existence and uniqueness of solutions of the time-dependent Ginzburg–Landau model for superconductivity*, Appl. Anal., 53 (1994), pp. 1–18.
- [8] M. FABRIZIO AND A. MORRO, *Electromagnetism of Continuous Media. Mathematical Modelling and Applications*, Oxford University Press, Oxford, 2003.
- [9] E. FEIREISL AND P. TAKÁČ, *Long-time stabilization of solutions to the Ginzburg–Landau equations of superconductivity*, Monatsh. Math., 133 (2001), pp. 197–221.
- [10] J. FLECKINGER–PELLÉ AND H.G. KAPER, *Gauges for the Ginzburg–Landau equations of superconductivity*, Z. Angew. Math. Mech., 96 (1996), pp. 345–348.
- [11] J. FLECKINGER–PELLÉ, H.G. KAPER, AND P. TAKÁČ, *Dynamics of the Ginzburg–Landau equations of superconductivity*, Nonlinear Anal., 32 (1998), pp. 647–665.
- [12] V. GIRAULT AND P.A. RAVIART, *Finite Element Methods for Navier–Stokes Equations*, Springer, New York, 1986.
- [13] L.P. GOR'KOV AND G.M. ÈLIASHBERG, *Generalizations of Ginzburg–Landau equations for nonstationary problems in the case of alloys with paramagnetic impurities*, Soviet. Phys. JETP, 27 (1968), pp. 328–334.
- [14] M. GRASSELLI, H. WU, AND S. ZHENG, *Asymptotic behavior of a non-isothermal Ginzburg–Landau model*, Quart. Appl. Math., to appear.
- [15] A. HARAUX AND M.A. JENDOUBI, *Convergence of bounded weak solutions of the wave equation with dissipation and analytic nonlinearity*, Calc. Var. Partial Differential Equations., 9 (1999), pp. 95–124.
- [16] A. HARAUX, *Systèmes Dynamiques Dissipatifs et Applications*, Masson, Paris, 1991.
- [17] A. HARAUX AND M.A. JENDOUBI, *Decay estimates to equilibrium for some evolution equations with an analytic nonlinearity*, Asymptot. Anal., 26 (2001), pp. 21–36.
- [18] K.-H. HOFFMANN AND Q. TANG, *Ginzburg–Landau phase-transition theory and superconductivity*, Internat. Ser. Numer. Math. 134, Birkhäuser, Basel, 2001.

- [19] S.-Z. HUANG AND P. TAKÁČ, *Convergence in gradient-like systems which are asymptotically autonomous and analytic*, *Nonlinear Anal.*, 46 (2001), pp. 675–698.
- [20] M.A. JENDOUBI, *A simple unified approach to some convergence theorem of L. Simon*, *J. Funct. Anal.*, 153 (1998), pp. 187–202.
- [21] H.G. KAPER AND P. TAKÁČ, *An equivalence relation for the Ginzburg–Landau equations of superconductivity*, *Z. Angew. Math. Phys.*, 48 (1997), pp. 665–675.
- [22] H.G. KAPER AND P. TAKÁČ, *Ginzburg–Landau dynamics with a time-dependent magnetic field*, *Nonlinearity*, 11 (1998), pp. 291–305.
- [23] J. LIANG AND Q. TANG, *Asymptotic behavior of the solutions of an evolutionary Ginzburg–Landau superconductivity model*, *J. Math. Anal. Appl.*, 195 (1995), pp. 92–107.
- [24] F.-H. LIN AND Q. DU, *Ginzburg–Landau vortices: Dynamics, pinning, and hysteresis*, *SIAM J. Math. Anal.*, 28 (1997), pp. 1265–1293.
- [25] P. MONK, *Finite Element Methods for Maxwell’s Equation*, Clarendon Press, Oxford, 2003.
- [26] A. RODRIGUEZ-BERNAL, B. WANG, AND R. WILLIE, *Asymptotic behavior of the time-dependent Ginzburg–Landau equations of superconductivity*, *Math. Methods Appl. Sci.*, 22 (1999), pp. 1647–1669.
- [27] A. SCHMID, *A time dependent Ginzburg–Landau equation and its application to the problem of resistivity in the mixed state*, *Phys. Kondens. Mater.*, 5 (1966), pp. 302–317.
- [28] L. SIMON, *Asymptotics for a class of nonlinear evolution equation with applications to geometric problems*, *Ann. of Math.*, 118 (1983), pp. 525–571.
- [29] Q. TANG, *On an evolutionary system of Ginzburg–Landau equations with fixed total magnetic flux*, *Comm. Partial Differential Equations*, 20 (1995), pp. 1–36.
- [30] Q. TANG AND S. WANG, *Time dependent Ginzburg–Landau equations of superconductivity*, *Phys. D*, 88 (1995), pp. 139–166.
- [31] R. TEMAM, *Infinite-dimensional dynamical systems in mechanics and physics*, *Appl. Math. Sci.* 68, Springer, New York, 1988.
- [32] M. TSUTSUMI AND H. KASAI, *The time-dependent Ginzburg–Landau–Maxwell equations*, *Nonlinear Anal.*, 37 (1999), pp. 187–216.
- [33] S. WANG AND M.Q. ZHAN, *L^p solutions to the time-dependent Ginzburg–Landau equations of superconductivity*, *Nonlinear Anal.*, 36 (1999), pp. 661–677.
- [34] H. WU, *Convergence to equilibrium for a Cahn–Hilliard model with the Wentzell boundary condition*, *Asymptot. Anal.*, 54 (2007), pp. 71–92.
- [35] H. WU, M. GRASSELLI, AND S. ZHENG, *Convergence to equilibrium for a parabolic–hyperbolic phase–field system with Neumann boundary conditions*, *Math. Models Methods Appl. Sci.*, 17 (2007), pp. 1–29.
- [36] S. ZHENG, *Nonlinear evolution equations*, Chapman & Hall–CRC Monogr. Surv. Pure Appl. Math. 133, Chapman & Hall/CRC, Boca Raton, FL, 2004.

ON A FOURTH ORDER NONLINEAR ELLIPTIC EQUATION WITH NEGATIVE EXPONENT*

ZONGMING GUO[†] AND JUNCHENG WEI[‡]

Abstract. We consider the following nonlinear fourth order equation: $T\Delta u - D\Delta^2 u = \frac{\lambda}{(L+u)^2}$, $-L < u < 0$, in Ω , $u = 0$, $\Delta u = 0$ on $\partial\Omega$, where $\lambda > 0$ is a parameter. This nonlinear equation models the deflection of charged plates in electrostatic actuators under the pinned boundary condition (Lin and Yang [*Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.*, 463 (2007), pp. 1323–1337]). Lin and Yang proved that there exists a $\lambda_c > 0$ such that for $\lambda > \lambda_c$ there is no solution, while for $\lambda < \lambda_c$ there is a branch of maximal solutions. In this paper, we show that in the physical domains (two or three dimensions) the maximal solution is unique and regular at $\lambda = \lambda_c$. In a two-dimensional (2D) convex smooth domain, we also establish the existence of a second mountain-pass solution for $\lambda \in (0, \lambda_c)$. The asymptotic behavior of the second solution is also studied. The main difficulty is the analysis of the touch-down behavior.

Key words. fourth order, electrostatic actuation, touch-down, pull-in threshold

AMS subject classifications. 35B45, 35J40

DOI. 10.1137/070703375

1. Introduction. We consider the structure of solutions to the problem

$$(P_\lambda) \quad \begin{cases} T\Delta u - D\Delta^2 u = \frac{\lambda}{(L+u)^2} & \text{in } \Omega, \\ -L < u \leq 0 & \text{in } \Omega, \\ u = 0, \quad \Delta u = 0 & \text{on } \partial\Omega, \end{cases}$$

where $\lambda > 0$ is a parameter, $T > 0$, $D > 0$, and $L > 0$ are fixed constants, and $\Omega \subset \mathbf{R}^N$ ($N \geq 2$) is a bounded smooth domain.

When $D = 0$, problem (P_λ) becomes

$$(Q_\lambda) \quad \begin{cases} T\Delta u = \frac{\lambda}{(L+u)^2} & \text{in } \Omega, \\ -L < u \leq 0 & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega, \end{cases}$$

which models a simple electrostatic microelectromechanical system (MEMS) device consisting of a thin dielectric elastic membrane with boundary supported at 0 above a rigid plate located at $-L$. Here $L + u$ represents the distance from the membrane to the plate. Recently there have been many studies on (Q_λ) . See, for example, [9], [14], [15], [16], [10], [11], [12], [7], [8], [17], [25], [24], and the references therein. These papers deal only with second order semilinear elliptic equations with singular nonlinearities. Equation (Q_λ) also appears in the study of thin film; see, for example, [2], [3], [5], [19], [20], [21], [18], and the references therein.

*Received by the editors September 21, 2007; accepted for publication (in revised form) October 2, 2008; published electronically January 21, 2009.

<http://www.siam.org/journals/sima/40-5/70337.html>

[†]Department of Mathematics, Henan Normal University, Xinxiang, 453007, People's Republic of China (gzm@henannu.edu.cn). The research of this author was supported by grants of NSFC (10571022 and 10871060).

[‡]Department of Mathematics, The Chinese University of Hong Kong, Shatin, Hong Kong (wei@math.cuhk.edu.hk). The research of this author was partially supported by earmarked grants from RGC of Hong Kong and a direct grant of CUHK.

In a recent paper [22], Lin and Yang derived the fourth order equation (P_λ) in the study of the deflection of charged plates in electrostatic actuators. Here $\lambda = aV^2$, where V is the electric voltage and a is positive constant. Associated with (P_λ) is the energy functional

$$(1.1) \quad E(u) = \int_{\Omega} \left\{ \frac{T}{2} |\nabla u|^2 + \frac{D}{2} |\Delta u|^2 - \frac{\lambda}{L+u} \right\},$$

where $P = \int_{\Omega} \frac{T}{2} |\nabla u|^2 dx$ is the stretching energy, $Q = \int_{\Omega} \frac{D}{2} |\Delta u|^2 dx$ corresponds to the bending energy, and $W = - \int_{\Omega} \frac{\lambda}{L+u(x)} dx$ is the electric potential energy.

Lin and Yang [22] considered two kinds of boundary conditions: the pinned boundary condition

$$u = \Delta u = 0 \text{ on } \partial\Omega$$

and the clamped boundary condition

$$u = \frac{\partial u}{\partial n} = 0 \text{ on } \partial\Omega.$$

For the pinned boundary condition problem (P_λ) , they found that there exists $0 < \lambda_c < \infty$ such that for $\lambda \in (0, \lambda_c)$, (P_λ) has a maximal regular solution u_λ , which can be obtained from an iterative scheme. (By a regular solution u_λ of (P_λ) , we mean that $u_\lambda \in C^4(\Omega) \cap C^3(\bar{\Omega})$ satisfies (P_λ) .) For $\lambda > \lambda_c$, (P_λ) does not have any regular solution. Moreover, if $\lambda', \lambda'' \in (0, \lambda_c)$ and $\lambda' < \lambda''$, then the corresponding maximal solutions $u_{\lambda'}$ and $u_{\lambda''}$ satisfy

$$u_{\lambda'} > u_{\lambda''} \text{ in } \Omega.$$

Physically, this is a natural relation because a higher supply voltage results in a greater elastic deformation or deflection.

The number λ_c , which determines the pull-in voltage, is called the pull-in threshold. It is known from [22] that, for $\lambda \in (0, \lambda_c)$, $\min_{\Omega}(L + u_\lambda) > 0$. Let $\Sigma_\lambda = \{x \in \Omega : L + u_\lambda(x) = 0\}$ be the singular set of (P_λ) . An interesting question is to study the limit of u_λ as $\lambda \nearrow \lambda_c$. The monotonicity of u_λ with respect to λ implies that there is a well-defined function U so that

$$U(x) = \lim_{\lambda \rightarrow \lambda_c^-} u_\lambda(x); \quad -L \leq U(x) < 0, \quad x \in \Omega.$$

However, $U(x)$ may touch down to $-L$ and cease to be a regular solution to (P_{λ_c}) . (By [22], $U \in W_{loc}^{2,2}(\Omega)$.) For the one-dimensional case, Lin and Yang showed that U is a regular solution; that is, the set $\Sigma_{\lambda_c} = \emptyset$.

In this paper, we will show that *for two dimensions and three dimensions, U is a regular solution.* Moreover, we also show that there is a *unique solution for (P_λ) at $\lambda = \lambda_c$.* To obtain our results, we first prove that the solutions u_λ for $\lambda \in (0, \lambda_c)$ obtained in [22] are stable in some sense. Furthermore, we also obtain the structure of solutions of (P_λ) in the two-dimensional (2D) case. Our main results of this paper are as follows.

THEOREM 1.1. *For dimension $N = 2$ or 3 , there exists a constant $0 < C := C(N, L)$ independent of λ such that for any $0 < \lambda < \lambda_c$, the maximal solution u_λ of (P_λ) satisfies $\min_{\Omega}(L + u_\lambda) \geq C$.*

Consequently, $u_{\lambda_c} = \lim_{\lambda \nearrow \lambda_c} u_\lambda$ exists in the topology of $C^4(\Omega)$. It is the unique regular solution to (P_{λ_c}) .

THEOREM 1.2. *Let $N = 2$ and Ω be a bounded, smooth, and convex domain in \mathbf{R}^2 . For $\lambda \in (0, \lambda_c]$, any solution of the problem (P_λ) is regular and the following hold.*

(i) *For $0 < \lambda < \lambda_c$, problem (P_λ) admits two solutions: the maximal solution and a mountain-pass solution.*

(ii) *For $\lambda = \lambda_c$, problem (P_λ) admits a unique regular solution.*

(iii) *For $\lambda > \lambda_c$, problem (P_λ) admits no regular solution.*

Furthermore, the mountain-pass solution V_λ has the following asymptotic behavior as $\lambda \rightarrow 0$:

$$(1.2) \quad \max_{\Omega} V_\lambda \rightarrow L \text{ as } \lambda \rightarrow 0, \quad \lim_{\lambda \rightarrow 0^+} \frac{[\min_{\Omega}(L - V_\lambda)]^3}{\lambda} = 0.$$

Remark. Theorem 1.2 shows that the bifurcation diagram of (P_λ) changes drastically when $D > 0$. In a nice 2D domain (see [16]), it has been proved in [16] that for λ small, the maximal solution is *unique*, and there exists $0 < \lambda_* < \lambda_c$ such that the solutions to (Q_λ) undergo infinitely many turning points. An interesting question is the asymptotic behavior as $D \rightarrow 0$. When $\Omega = B_1 \subset \mathbf{R}^2$, the complete bifurcation picture as well as the asymptotic behavior when $D \rightarrow 0$ has been considered in [23].

The organization of the paper is as follows: in section 2, we present some preliminary results on the first eigenvalue and the corresponding eigenfunction of the problem

$$-T\Delta\phi + D\Delta^2\phi = \sigma\phi \text{ in } \Omega, \quad \phi = \Delta\phi = 0 \text{ on } \partial\Omega.$$

In section 3, we derive a key L^1 bound for $\frac{1}{(L+u)^2}$. In section 4, we show the stability of the maximal solutions of (P_λ) . In section 5, we show that the solution at the pull-in threshold is regular for $N = 2$ or 3. In section 6, we show that any weak solution at the pull-in threshold is unique. In section 7, we present the structure of the solutions of (P_λ) for the 2D case. We show that for $0 < \lambda < \lambda_c$, (P_λ) admits at least two solutions: the maximal solution and a mountain-pass solution. Finally in section 8, we give some asymptotic behaviors of the mountain-pass solution as $\lambda \rightarrow 0^+$.

2. The first eigenfunction. In this section, we study the following eigenvalue problem:

$$(2.1) \quad -T\Delta\phi + D\Delta^2\phi = \sigma\phi \text{ in } \Omega, \quad \phi = \Delta\phi = 0 \text{ on } \partial\Omega,$$

where $T, D > 0$. We will show that (2.1) has the least eigenvalue σ_1 and the corresponding eigenfunction $\phi_1 > 0$ in Ω . Moreover, ϕ_1 is simple; i.e., all the eigenfunctions corresponding to σ_1 assume the forms of $C\phi_1$ with $C \in \mathbf{R}$.

PROPOSITION 2.1. *Problem (2.1) has the least eigenvalue σ_1 such that all the eigenfunctions corresponding to σ_1 assume the forms of $C\phi_1$, where $\phi_1 \in C^\infty(\Omega)$ and $\phi_1 > 0$ in Ω .*

Proof. This proposition may be known, but we cannot find the reference. We give a proof here for completeness.

Consider the following minimization problem:

$$(2.2) \quad \sigma_1 := \inf \left\{ \int_{\Omega} [T|\nabla\phi|^2 + D|\Delta\phi|^2] dx : \phi \in \mathcal{H}, \|\phi\|_{L^2(\Omega)} = 1 \right\},$$

where $\mathcal{H} = H^2(\Omega) \cap H_0^1(\Omega)$ is the function space obtained by taking the completion under the norm of $H^2(\Omega) \cap H_0^1(\Omega)$ (i.e., $\|\psi\| = (\int_{\Omega} [T|\nabla\psi|^2 + D|\Delta\psi|^2]dx)^{1/2}$) for the set of smooth functions that satisfy the boundary condition $\phi = \Delta\phi = 0$ on $\partial\Omega$. Since the Sobolev embedding $\mathcal{H} \hookrightarrow L^2(\Omega)$ is compact, by the standard direct method of calculus of variations, we have at least one minimizer ϕ_1 for the problem (2.2), where $\phi_1 \in \mathcal{H}$, $\|\phi_1\|_{L^2(\Omega)} = 1$. Furthermore, ϕ_1 is a weak solution to (2.1); namely,

$$(2.3) \quad \int_{\Omega} [T\nabla\phi_1\nabla\phi + D\Delta\phi_1\Delta\phi] = \sigma_1 \int_{\Omega} \phi_1\phi dx \quad \forall\phi \in \mathcal{H}.$$

Using the L^p -estimates due to Agmon, Douglis, and Nirenberg [1], we conclude that

$$\|\phi_1\|_{W^{4,p}(\Omega)} \leq C\|\phi_1\|_{L^p(\Omega)}$$

for any $p > 1$. Thus we have $\phi_1 \in C^4(\Omega) \cap C^3(\bar{\Omega})$, and hence $\Delta\phi_1 = 0$ on $\partial\Omega$, and ϕ_1 satisfies (2.1). (See a similar argument in Lemma B.3 of [27].)

It is clear that

$$\sigma_1 = \frac{\int_{\Omega} [T|\nabla\phi_1|^2 + D|\Delta\phi_1|^2]dx}{\int_{\Omega} \phi_1^2 dx} = \inf_{\phi \in \mathcal{H} \setminus \{0\}} \frac{\int_{\Omega} [T|\nabla\phi|^2 + D|\Delta\phi|^2]dx}{\int_{\Omega} \phi^2 dx}.$$

In order to show that ϕ_1 is of fixed sign, we consider the following new problem:

$$(2.4) \quad -T\Delta\psi_1 + D\Delta^2\psi_1 = \sigma_1|\phi_1| \quad \text{in } \Omega, \quad \psi_1 = \Delta\psi_1 = 0 \quad \text{on } \partial\Omega.$$

By the maximum principle, $\psi_1 > 0$, $-D\Delta\psi_1 + T\psi_1 > 0$ in Ω . Furthermore, we have $\psi_1 \geq \phi_1, \psi_1 \geq -\phi_1$, and hence $\psi_1 \geq |\phi_1|$ in Ω .

On the other hand, from (2.4) we obtain

$$(2.5) \quad \int_{\Omega} [T|\nabla\psi_1|^2 + D|\Delta\psi_1|^2]dx = \sigma_1 \int_{\Omega} \psi_1|\phi_1| dx \leq \sigma_1 \int_{\Omega} |\psi_1|^2 dx.$$

By the minimality of σ_1 , we have

$$(2.6) \quad \sigma_1 = \frac{\int_{\Omega} [T|\nabla\psi_1|^2 + D|\Delta\psi_1|^2]dx}{\int_{\Omega} |\psi_1|^2 dx}.$$

Thus ψ_1 also attains σ_1 , and hence the inequality of (2.5) is actually an equality. This implies that $\psi_1 = |\phi_1|$ in Ω . Since $\psi_1 > 0$ in Ω , we conclude that ϕ_1 is of fixed sign in Ω .

The above argument actually proves that any nonzero eigenfunction corresponding to σ_1 must be of fixed sign in Ω . So if ϕ_1 and ϕ_2 are two eigenfunctions corresponding to σ_1 , we may choose $\phi_1 > 0, \phi_2 > 0$. Let $x_0 \in \Omega$ and $C = \frac{\phi_1(x_0)}{\phi_2(x_0)}$. Then the function $\phi_1 - C\phi_2$ is again an eigenfunction corresponding to σ_1 . By the previous argument, we see that $\phi_1 \equiv C\phi_2$ in Ω . This completes the proof. \square

3. A uniform L^1 bound. In this section, we establish a key uniform L^1 bound for $\frac{1}{(L-v)^2}$, where v satisfies

$$(T\lambda) \quad \begin{cases} -T\Delta v + D\Delta^2 v = \frac{\lambda}{(L-v)^2} & \text{in } \Omega, \\ 0 < v < L & \text{in } \Omega, \\ v = 0, \quad \Delta v = 0 & \text{on } \partial\Omega \end{cases}$$

(which is equivalent to (P_λ) by taking $u = -v$). Note that $v \in C^4(\Omega) \cap C^2(\bar{\Omega})$ provided that v satisfies (T_λ) .

THEOREM 3.1. *Let Ω be a bounded, smooth, and convex domain. Then there exists a constant C (independent of λ) such that for any solution v to (T_λ) we have*

$$(3.1) \quad \int_{\Omega} \frac{1}{(L-v)^2} \leq \frac{C}{\lambda}.$$

As a consequence, we have

$$(3.2) \quad \int_{\Omega} (D|\Delta v|^2 + T|\nabla v|^2) \leq C.$$

Proof. Let ϕ_1 be as given in Proposition 2.1. Multiplying (T_λ) by ϕ_1 and integrating over Ω , we obtain

$$(3.3) \quad \lambda \int_{\Omega} \frac{1}{(L-v)^2} \phi_1 = \sigma_1 \int_{\Omega} v \phi_1 \leq C,$$

which implies that

$$(3.4) \quad \int_{\Omega'} \frac{1}{(L-v)^2} \leq \frac{C_{\Omega'}}{\lambda}$$

for any $\Omega' \subset \subset \Omega$, where $C_{\Omega'}$ is independent of λ .

We write (T_λ) as

$$\begin{cases} \Delta v + \frac{1}{D}w - \frac{T}{D}v = 0 & \text{in } \Omega, \\ \Delta w + \frac{\lambda}{(L-v)^2} = 0 & \text{in } \Omega, \\ v = w = 0 & \text{on } \partial\Omega. \end{cases}$$

If we denote $f_1(v, w) = -\frac{T}{D}v + \frac{1}{D}w$, $f_2(v, w) = \frac{\lambda}{(L-v)^2}$, we see that $\frac{\partial f_1}{\partial w} = \frac{1}{D} > 0$ and $\frac{\partial f_2}{\partial v} = \frac{2\lambda}{(L-v)^3} > 0$. Therefore, the convexity of Ω , Lemma 5.1 of [26], and the moving plane method near $\partial\Omega$ as in the appendix of [13] imply that there exist $t_0 > 0$ and $\alpha > 0$ depending only on the domain Ω , such that $v(x - t\nu)$ and $w(x - t\nu)$ are nondecreasing for $t \in [0, t_0]$, $\nu \in R^N$ satisfying $|\nu| = 1$ and $(\nu, n(x)) \geq \alpha$ and $x \in \partial\Omega$. Therefore, we can find $\gamma, \delta > 0$ such that for any $x \in \Omega_\delta := \{z \in \Omega : d(z, \partial\Omega) < \delta\}$ there exists a fixed-sized cone Γ_x (with x as its vertex) with

- (i) $meas(\Gamma_x) \geq \gamma$,
- (ii) $\Gamma_x \subset \{z \in \Omega : d(z, \partial\Omega) < \delta\}$, and
- (iii) $v(y) \geq v(x)$ for any $y \in \Gamma_x$.

Then, for any $x \in \Omega_\delta$, we have

$$\frac{1}{(L-v(x))^2} \leq \frac{1}{meas(\Gamma_x)} \int_{\Gamma_x} \frac{1}{(L-v)^2} \leq \frac{1}{\gamma} \int_{\Omega_\delta} \frac{1}{(L-v)^2} \leq \frac{C}{\lambda}.$$

This implies that $\frac{1}{(L-v)^2} \in L^\infty(\Omega_\delta)$ and there is $C > 0$ independent of λ such that

$$(3.5) \quad \sup_{x \in \Omega_\delta} v < L - C\sqrt{\lambda}.$$

Next, we derive estimates for w near $\partial\Omega$. Multiplying the second equation of the equivalent system of (T_λ) by φ_0 , the first eigenfunction of $-\Delta$, and integrating over Ω , we obtain

$$(3.6) \quad \lambda \int_{\Omega} \frac{1}{(L-v)^2} \varphi_0 = \lambda_1 \int_{\Omega} \varphi_0 w,$$

where λ_1 is the first eigenvalue of $-\Delta$. By (3.3), we have that

$$(3.7) \quad \lambda_1 \int_{\Omega} \varphi_0 w = \lambda \int_{\Omega} \frac{1}{(L-v)^2} \varphi_0 \leq \lambda C \int_{\Omega} \frac{1}{(L-v)^2} \phi_1 \leq C,$$

and hence

$$\int_{\Omega'} w \leq C_{\Omega'}$$

for any $\Omega' \subset\subset \Omega$. To see the second inequality of (3.7), we notice that there exist $\ell_i > 0$ ($i = 1, 2, 3, 4$) such that

$$\ell_1 d(x) \leq \varphi_0(x) \leq \ell_2 d(x), \quad \ell_3 d(x) \leq \phi_1(x) \leq \ell_4 d(x),$$

where $d(x) = \text{dist}(x, \partial\Omega)$. Hence $\varphi_0(x) \leq C\phi_1(x)$. The same reason as above shows that $w \leq C(\Omega_\delta)$.

By elliptic regularity applied to the system (T_λ) (noting that v, w , and $\frac{1}{(L-v)^2}$ are all bounded in Ω_δ), we have $v \in C^3(\Omega_\delta)$, and hence

$$\lambda \int_{\Omega} \frac{1}{(L-v)^2} = D \int_{\partial\Omega} \frac{\partial\Delta v}{\partial n} - T \int_{\Omega} \frac{\partial v}{\partial n} \leq C.$$

To prove the inequality (3.2), we multiply (T_λ) by v and integrate over Ω to obtain

$$T \int_{\Omega} |\nabla v|^2 + D \int_{\Omega} |\Delta v|^2 = \lambda \int_{\Omega} \frac{Lv}{(L-v)^2} \leq C. \quad \square$$

4. Stability of the maximal solutions of (P_λ) . In this section, we show that the maximal solutions u_λ to (P_λ) obtained in [22] for $\lambda \in (0, \lambda_c)$ are stable in some sense. Let $v_\lambda = -u_\lambda$. Then, from [22], for each $\lambda \in (0, \lambda_c)$, (T_λ) has a minimal positive solution v_λ .

We call v_λ stable if the first eigenvalue $\sigma_{1,\lambda}(v_\lambda)$ of the problem

$$(4.1) \quad -T\Delta h + D\Delta^2 h = \frac{2\lambda}{(L-v_\lambda)^3} h + \sigma h \text{ in } \Omega, \quad h = \Delta h = 0 \text{ on } \partial\Omega$$

is nonnegative. By arguments similar to those in the proof of Proposition 2.1, we see that the first eigenvalue $\sigma_{1,\lambda}(v_\lambda)$ exists and every eigenfunction corresponding to $\sigma_{1,\lambda}(v_\lambda)$ is of fixed sign if $\sigma_{1,\lambda}(v_\lambda) \geq 0$.

LEMMA 4.1. *Suppose that v is a regular solution of (T_λ) , and u is a regular supersolution of (T_λ) ; that is,*

$$\begin{cases} -T\Delta u + D\Delta^2 u \geq \frac{\lambda}{(L-u)^2} & \text{in } \Omega, \\ 0 < u < L & \text{in } \Omega, \\ u = 0, \Delta u = 0 & \text{on } \partial\Omega. \end{cases}$$

If $\sigma_{1,\lambda}(v) > 0$, then $u \geq v$ in Ω , and if $\sigma_{1,\lambda}(v) = 0$, then $u = v$ in Ω .

Proof. For a given λ and $x \in \Omega$, by the fact that $s \rightarrow (L - s)^{-2}$ is convex on $(0, L)$, we see that

$$(4.2) \quad -T\Delta(v + \tau(u - v)) + D\Delta^2(v + \tau(u - v)) - \frac{\lambda}{[L - (v + \tau(u - v))]^2} \geq 0 \text{ in } \Omega$$

for $\tau \in [0, 1]$. Note that (4.2) is an identity at $\tau = 0$, which means that the first derivative of the left-hand side of (4.2) with respect to τ is nonnegative at $\tau = 0$, i.e.,

$$(4.3) \quad \begin{cases} -T\Delta(u - v) + D\Delta^2(u - v) - \frac{2\lambda}{(L-v)^3}(u - v) \geq 0 & \text{in } \Omega, \\ u - v = 0, \quad \Delta(u - v) = 0 & \text{on } \partial\Omega. \end{cases}$$

Thus, the fact $\sigma_{1,\lambda}(v) > 0$ implies that $u \geq v$ in Ω . Indeed, on the contrary, we see that $0 \neq (u - v)^- \in H^2(\Omega) \cap H_0^1(\Omega)$. Multiplying $(u - v)^-$ on both the sides of (4.3) and integrating it on Ω , we see that

$$\begin{aligned} \sigma_{1,\lambda}(v) & \int_{\Omega} [(u - v)^-]^2 dx \\ & \leq T \int_{\Omega} |\nabla(u - v)^-|^2 dx + D \int_{\Omega} |\Delta(u - v)^-|^2 dx - \int_{\Omega} \frac{2\lambda}{(L - v)^3} [(u - v)^-]^2 dx \\ & \leq 0. \end{aligned}$$

This contradicts $\sigma_{1,\lambda}(v) > 0$.

If $\sigma_{1,\lambda}(v) = 0$, we have

$$-T\Delta(u - v) + D\Delta^2(u - v) - \frac{2\lambda}{(L - v)^3}(u - v) = 0 \text{ in } \Omega.$$

Moreover, the second derivative of the left-hand side of (4.2) with respect to τ at $\tau = 0$ is

$$-6\lambda(L - v)^{-4}(u - v)^2 \geq 0,$$

which implies that $u \equiv v$ in Ω . This completes the proof. \square

PROPOSITION 4.2. For each $\lambda \in (0, \lambda_c)$, the minimal positive solution v_λ of (T_λ) is stable.

Proof. Since $\sigma_1 > 0$, we easily see that the first eigenvalue $\sigma_{1,\lambda}(v_\lambda)$ of problem (4.1) is positive provided that λ is sufficiently small. Now we prove that $\sigma_{1,\lambda}(v_\lambda) > 0$ for $\lambda \in (0, \lambda_c)$.

We define

$$\lambda^* = \sup\{\rho : v_\lambda \text{ is a stable solution for } \lambda \in (0, \rho)\}.$$

It is clear that $\lambda^* \leq \lambda_c$. To show $\lambda^* = \lambda_c$, it suffices to prove that there is no regular minimal solution for (T_λ) with $\lambda > \lambda^*$. For that, suppose w is a regular minimal solution of $(T_{\lambda^*+\delta})$ with $\delta > 0$; then we would have for $\lambda \leq \lambda^*$

$$-T\Delta w + D\Delta^2 w = \frac{\lambda^* + \delta}{(L - w)^2} \geq \frac{\lambda}{(L - w)^2} \text{ in } \Omega.$$

Since for $0 < \lambda < \lambda^*$ the minimal solution v_λ is stable, it follows from Lemma 4.1 that $L > w \geq v_\lambda$. Consequently, $\bar{v} = \lim_{\lambda \nearrow \lambda^*} v_\lambda$ exists in $C^4(\Omega)$, and it is a regular solution to (T_{λ^*}) . Now, from the definition of λ^* and the implicit function theorem, we necessarily have $\sigma_{1,\lambda^*}(\bar{v}) = 0$. By Lemma 4.1 again, we obtain that $w \equiv \bar{v}$ in Ω and hence $\delta = 0$. This is a contradiction. Therefore, $\lambda^* = \lambda_c$. This completes the proof. \square

5. The regularity of the minimal solution of (T_λ) at $\lambda = \lambda_c$. In this section, we are concerned with the regularity of the minimal solutions of (T_λ) at $\lambda = \lambda_c$. Normally, the minimal solution v_λ at $\lambda = \lambda_c$ may have a singular set in Ω ; i.e., there exists a set $\Sigma_{\lambda_c} \subset \Omega$ such that $v_{\lambda_c}(x) = L$ for $x \in \Sigma_{\lambda_c}$. But we will see that for the lower dimensional case, $v_{\lambda_c} < L$ in Ω .

By a weak solution $v \in \mathcal{H}$ of (T_λ) we mean $0 < v \leq L$ in Ω and $(L - v)^{-2} \in L^1(\Omega)$ such that for any $\varphi \in \mathcal{H}$,

$$\int_{\Omega} [T \nabla v \cdot \nabla \varphi + D \Delta v \Delta \varphi] dx = \lambda \int_{\Omega} (L - v)^{-2} \varphi dx.$$

LEMMA 5.1. *If $v \in \mathcal{H}$ is a weak solution to (T_λ) , then there exists $C := C(\lambda) > 0$ such that*

$$\int_{\Omega} \frac{dx}{(L - v)^2} \leq C.$$

For $N \geq 2$, any solution v satisfying $(L - v)^{-2} \in L^p(\Omega)$ with $p = N/2$ is a classical solution.

Proof. For the first conclusion, we see that, since $v \in \mathcal{H}$ is a solution of (T_λ) ,

$$(5.1) \quad \int_{\Omega} \frac{v}{(L - v)^2} dx = \frac{1}{\lambda} \left[\int_{\Omega} (|\nabla v|^2 + |\Delta v|^2) dx \right] \leq C.$$

On the other hand, we see that

$$\frac{v}{(L - v)^2} = \frac{L}{(L - v)^2} - \frac{1}{(L - v)}.$$

Thus, (5.1) implies that

$$(5.2) \quad \int_{\Omega} \frac{L}{(L - v)^2} dx = \int_{\Omega} \frac{1}{(L - v)} dx + \int_{\Omega} \frac{v}{(L - v)^2} dx.$$

By Young's inequality, we have that

$$(5.3) \quad \int_{\Omega} \frac{1}{(L - v)} dx \leq \epsilon L \int_{\Omega} \frac{1}{(L - v)^2} dx + C(\epsilon, L) |\Omega|,$$

where $0 < \epsilon < 1/4$ and $C(\epsilon, L) > 0$ is a constant. Our first conclusion can be obtained from (5.1), (5.2), and (5.3).

For $N = 2$, suppose that v is a weak solution such that $\frac{1}{(L - v)^2} \in L^1(\Omega)$. Thus,

$$-D \Delta^2 v = \lambda(L - v)^{-2} - T \Delta v \in L^1(\Omega).$$

This and the Sobolev embedding imply that $\nabla^3 v \in L^q(\Omega)$ for any $1 < q < 2$. In particular, $\nabla v \in C^{2-\frac{2}{q}}(\Omega)$ for any $1 < q < 2$. This and the fact that $(L - v)^{-2} \in L^1(\Omega)$ clearly imply that $v < L$ in Ω . In fact, on the contrary, suppose that there exists $x_0 \in \Omega$ such that $v(x_0) = \max_{\Omega} v = L$. Then, $\nabla v(x_0) = 0$ and

$$v(x) - v(x_0) = \nabla v(\xi) \cdot (x - x_0) \text{ for } x \in \Omega \text{ near } x_0,$$

where $\xi = tx_0 + (1 - t)x$ with $t \in (0, 1)$. Moreover, since $\nabla v \in C^{2-\frac{2}{q}}(\Omega)$, we see that

$$|\nabla v(\xi) - \nabla v(x_0)| \leq M |\xi - x_0|^{2-\frac{2}{q}} \leq M |x - x_0|^{2-\frac{2}{q}},$$

and thus

$$|v(x) - v(x_0)| \leq M|x - x_0|^{3-\frac{2}{q}} \text{ for } x \in \Omega \text{ near } x_0.$$

This inequality shows that

$$\infty > \int_{\Omega} \frac{1}{(L-v)^2} dx \geq M^{-2} \int_{\Omega} |x - x_0|^{-(6-\frac{4}{q})} dx = \infty,$$

which is a contradiction, which implies that we must have $\|v\|_{C(\bar{\Omega})} < L$.

For $N \geq 3$, suppose that v is a weak solution such that $\frac{1}{(L-v)^2} \in L^p(\Omega)$ with $p = \frac{N}{2}$. By the regularity of Δ^2 , we see that $v \in W^{4,p}(\Omega)$. The Sobolev embedding theorem then implies that $v \in C^{1,\alpha}(\Omega)$ with $\alpha < 1$ since $4 - \frac{N}{p} = 2$. To show that v is a classical solution, it suffices to show that $v < L$ in Ω . Indeed, on the contrary, there exists $x_0 \in \Omega$ such that $v(x_0) = \max_{\Omega} v = L$. Then, $\nabla v(x_0) = 0$ and

$$(5.4) \quad v(x) - v(x_0) = \nabla v(\xi) \cdot (x - x_0) \text{ for } x \in \Omega \text{ near } x_0,$$

where $\xi = tx_0 + (1-t)x$ with $t \in (0, 1)$. Moreover, since $v \in C^{1,1}(\Omega)$, we see that $|\nabla v(\xi) - \nabla v(x_0)| \leq M|\xi - x_0| \leq M|x - x_0|$. This and (5.4) imply that

$$(5.5) \quad |v(x) - v(x_0)| \leq M|x - x_0|^{1+\alpha} \text{ for } x \in \Omega.$$

This inequality shows that

$$\infty > \int_{\Omega} \left(\frac{1}{(L-v)^2} \right)^p dx \geq M^{-2p} \int_{\Omega} |x - x_0|^{-2(1+\alpha)p} dx = \infty,$$

which is a contradiction, which implies that we must have $\|v\|_{C(\bar{\Omega})} < L$. This completes the proof. \square

PROPOSITION 5.2. *There exists a constant $C := C(L, \lambda) > 0$ such that for each $\lambda \in (0, \lambda_c)$, the minimal solution v_{λ} satisfies $\|(L - v_{\lambda})^{-2}\|_{L^{3/2}(\Omega)} \leq C$.*

Proof. Since the minimal solutions v_{λ} are stable, we have

$$(5.6) \quad \int_{\Omega} \frac{2\lambda}{(L - v_{\lambda})^3} w^2 dx \leq \int_{\Omega} [T|\nabla w|^2 + D|\Delta w|^2] dx$$

for all $0 < \lambda < \lambda_c$ and nonnegative $w \in \mathcal{H}$.

Let $w = v_{\lambda}$; we then have

$$(5.7) \quad \int_{\Omega} \frac{2\lambda}{(L - v_{\lambda})^3} v_{\lambda}^2 dx \leq \int_{\Omega} [T|\nabla v_{\lambda}|^2 + D|\Delta v_{\lambda}|^2] dx = \int_{\Omega} \frac{\lambda v_{\lambda}}{(L - v_{\lambda})^2} dx.$$

Since $v_{\lambda} < L$, this implies that

$$(5.8) \quad \int_{\Omega} \frac{v_{\lambda}^2}{(L - v_{\lambda})^3} dx \leq C$$

and

$$(5.9) \quad \int_{\Omega} \frac{L^2}{(L - v_{\lambda})^3} dx \leq \int_{\Omega} \frac{v_{\lambda}^2}{(L - v_{\lambda})^3} dx + \int_{\Omega} \frac{(L - v_{\lambda})^2}{(L - v_{\lambda})^3} dx \leq C + \int_{\Omega} \frac{1}{L - v_{\lambda}} dx.$$

Hence

$$(5.10) \quad \int_{\Omega} \frac{1}{(L - v_{\lambda})^3} dx \leq C.$$

This completes the proof. \square

Now we obtain the following theorem, from which our Theorem 1.1 can be obtained.

THEOREM 5.3. *For dimension $N = 2$ or 3 , there exists a constant $0 < C := C(N, L) < L$ independent of λ such that for any $0 < \lambda < \lambda_c$, the minimal solution v_{λ} of (T_{λ}) satisfies $\|v_{\lambda}\|_{C(\Omega)} \leq C$.*

Consequently, $v_{\lambda_c} = \lim_{\lambda \nearrow \lambda_c} v_{\lambda}$ exists in the topology of $C^4(\Omega)$. It is the unique classical solution to (T_{λ_c}) .

Proof. By Proposition 5.2 and (3.2), we see that there is $C > 0$ independent of λ such that

$$\|v_{\lambda}\|_{H^2(\Omega)} \leq C.$$

Since the mapping $\lambda \mapsto v_{\lambda}$ is increasing for $\lambda \in (0, \lambda_c)$, we see that there is a function $v_{\lambda_c} \in H^2(\Omega)$ such that

$$\lim_{\lambda \nearrow \lambda_c} v_{\lambda} = v_{\lambda_c} \text{ weakly in } H^2(\Omega).$$

Consequently, v_{λ_c} is a weak solution of the equation (T_{λ}) at the critical parameter λ_c ,

$$-T\Delta v_{\lambda_c} + D\Delta^2 v_{\lambda_c} = \frac{\lambda_c}{(L - v_{\lambda_c})^2} \text{ in } \Omega,$$

and in the sense of weak solutions, the critical value λ_c is attainable.

Now we show that v_{λ_c} is a classical solution. The implicit function theorem implies that the mapping $\lambda \mapsto v_{\lambda}$ from $(0, \lambda_c)$ to $C(\bar{\Omega})$ is continuous. Thus, we see that $\sigma_{1, \lambda_c} = 0$. (Otherwise, the implicit function theorem implies that v_{λ} will exist for $\lambda > \lambda_c$.) By arguments similar to those in the proof of Proposition 5.2, we see that

$$\|(L - v_{\lambda_c})^{-2}\|_{L^{3/2}(\Omega)} \leq C(L).$$

Note that (5.6) holds with the inequality replaced by an equality. Then Lemma 5.1 implies that for $N = 2$ and 3 , v_{λ_c} is a classical solution. Thus, there exists $C < L$ such that $\|v_{\lambda_c}\|_{C(\bar{\Omega})} \leq C$. Note that $\|v_{\lambda}\|_{C(\bar{\Omega})} \leq \|v_{\lambda_c}\|_{C(\bar{\Omega})} \leq C < L$ for $\lambda \in (0, \lambda_c)$. The uniqueness of v_{λ_c} of (T_{λ}) at $\lambda = \lambda_c$ follows from Lemma 4.1. This completes the proof. \square

6. Uniqueness of the solution of (T_{λ}) at $\lambda = \lambda_c$. We first note that the monotonicity with respect to λ and the uniform boundedness of the branch of the minimal solutions imply that the extremal function defined by $v_{\lambda_c} = \lim_{\lambda \nearrow \lambda_c} v_{\lambda}$ always exists and can always be considered as a solution for (T_{λ_c}) in a weak sense. On the other hand, if there is a $0 < C < L$ such that $\|v_{\lambda}\|_{C(\bar{\Omega})} \leq C$ for each $\lambda < \lambda_c$, just as in the case $N = 2$ or 3 , then we see from Theorem 5.3 that v_{λ_c} is the unique classical solution.

In the following, we consider only the case that v_{λ_c} is a weak solution (i.e., $v_{\lambda_c} \in W_{loc}^{2,2}(\Omega)$; note that we can obtain $v_{\lambda_c} \in \mathcal{H}$ provided that $v_{\lambda_c} \in W_{loc}^{2,2}(\Omega)$ by the moving plane argument) but with the possibility that $\|v_{\lambda_c}\|_{L^{\infty}(\Omega)} = L$.

THEOREM 6.1. For $\lambda > 0$, assume that $v \in \mathcal{H}$ is a weak solution to (T_λ) such that $\|v\|_{L^\infty(\Omega)} = L$. The following assertions are equivalent:

(i) $\sigma_{1,\lambda}(v) \geq 0$; that is, v satisfies

$$2\lambda \int_{\Omega} (L - v)^{-3} \phi^2 \leq \int_{\Omega} [T|\nabla\phi|^2 + D|\Delta\phi|^2] dx \quad \forall \phi \in \mathcal{H}.$$

(ii) $\lambda = \lambda_c$ and $v \equiv v_{\lambda_c}$ in Ω .

Theorem 6.1 can be easily obtained from the following proposition.

PROPOSITION 6.2. Let v_1, v_2 be two \mathcal{H} -weak solutions of (T_λ) so that $\sigma_{1,\lambda}(v_i) \geq 0$ for $i = 1, 2$. Then $v_1 = v_2$ a.e. in Ω .

Proof. For any $\theta \in [0, 1]$ and $\phi \in \mathcal{H}, \phi \geq 0$, we have that

$$\begin{aligned} I_{\theta,\phi} &:= T \int_{\Omega} \nabla(\theta v_1 + (1 - \theta)v_2) \nabla\phi dx + D \int_{\Omega} \Delta(\theta v_1 + (1 - \theta)v_2) \Delta\phi dx \\ &\quad - \lambda \int_{\Omega} [L - (\theta v_1 + (1 - \theta)v_2)]^{-2} \phi dx \\ &= \lambda \int_{\Omega} [\theta(L - v_1)^{-2} + (1 - \theta)(L - v_2)^{-2}] - [(L - (\theta v_1 + (1 - \theta)v_2))^{-2}] dx \\ &\geq 0 \end{aligned}$$

due to the convexity of $(L - s)^{-2}$ with respect to $s \in (0, L)$. Since $I_{0,\phi} = I_{1,\phi} = 0$, the derivative of $I_{\theta,\phi}$ at $\theta = 0, 1$ provides

$$\begin{aligned} \int_{\Omega} [T\nabla(v_1 - v_2) \nabla\phi + D\Delta(v_1 - v_2) \Delta\phi] - 2\lambda \int_{\Omega} (L - v_2)^{-3} (v_1 - v_2) \phi &\geq 0, \\ \int_{\Omega} [T\nabla(v_1 - v_2) \nabla\phi + D\Delta(v_1 - v_2) \Delta\phi] - 2\lambda \int_{\Omega} (L - v_1)^{-3} (v_1 - v_2) \phi &\leq 0 \end{aligned}$$

for any $\phi \in \mathcal{H}$ with $\phi \geq 0$. Testing the first inequality on $\phi = (v_1 - v_2)^-$ and the second one on $(v_1 - v_2)^+$, we obtain that

$$\begin{aligned} \int_{\Omega} [T|\nabla(v_1 - v_2)^-|^2 + D|\Delta(v_1 - v_2)^-|^2] - 2\lambda \int_{\Omega} (L - v_2)^{-3} ((v_1 - v_2)^-)^2 &\leq 0, \\ \int_{\Omega} [T|\nabla(v_1 - v_2)^+|^2 + D|\Delta(v_1 - v_2)^+|^2] - 2\lambda \int_{\Omega} (L - v_1)^{-3} ((v_1 - v_2)^+)^2 &\leq 0. \end{aligned}$$

Since $\sigma_{1,\lambda}(v_1) \geq 0$, we have the following:

- (1) If $\sigma_{1,\lambda}(v_1) > 0$, then $v_1 \leq v_2$ a.e. in Ω .
- (2) If $\sigma_{1,\lambda}(v_1) = 0$, then

$$\int_{\Omega} [T\nabla(v_1 - v_2) \nabla\bar{\varphi} + D\Delta(v_1 - v_2) \Delta\bar{\varphi}] - 2\lambda \int_{\Omega} (L - v_1)^{-3} (v_1 - v_2) \bar{\varphi} = 0,$$

where $\bar{\varphi} = (v_1 - v_2)^+$. Since $I_{\theta,\bar{\varphi}} \geq 0$ for any $\theta \in [0, 1]$ and $I_{1,\bar{\varphi}} = \partial I_{1,\bar{\varphi}} = 0$, we get that $\partial_{\theta\theta}^2 I_{1,\bar{\varphi}} = - \int_{\Omega} \frac{6\lambda}{(L - v_1)^4} ((v_1 - v_2)^+)^3 \geq 0$. Thus, $(v_1 - v_2)^+ = 0$ a.e. in Ω . Hence, $v_1 \leq v_2$ a.e. in Ω . The same argument applies to prove the reversed inequality, and the proof of the proposition is complete. \square

7. Structure of solutions of (T_λ) in the 2D case. In this section we obtain the structure of positive solutions of (T_λ) in the 2D case. The main theorem of this section is the following theorem. Our Theorem 1.2 can be obtained from this theorem and Theorem 8.2.

THEOREM 7.1. *Let Ω be a convex smooth domain in \mathbf{R}^2 . For $\lambda \in (0, \lambda_c]$, any solution of the problem (T_λ) is regular and the following hold.*

(i) *For $0 < \lambda < \lambda_c$, problem (T_λ) admits two solutions: the minimal solution and a mountain-pass solution.*

(ii) *For $\lambda = \lambda_c$, problem (T_λ) admits a unique regular solution.*

(iii) *For $\lambda > \lambda_c$, problem (T_λ) admits no regular solution.*

To prove this theorem, we first show the following lemma.

LEMMA 7.2. *For any fixed $\lambda > 0$, if $v_\lambda \in \mathcal{H}$ is a positive solution of (T_λ) , then there exists $0 < \tau_\lambda < L$ such that $v_\lambda \leq L - \tau_\lambda$ in Ω . This also implies that v_λ is regular.*

Proof. The embedding theorem implies that $v_\lambda \in C^\alpha(\bar{\Omega})$ for any $0 < \alpha < 1$, and thus the moving plane arguments as in the proof of Theorem 3.1 imply that if $v_\lambda(x_\lambda) = \max_\Omega v_\lambda$, then $x_\lambda \in \Omega_0$, where $\Omega_0 \subset\subset \Omega$. Moreover, by Theorem 3.1, we have

$$(7.1) \quad \int_\Omega (L - v_\lambda)^{-2} \leq \frac{C}{\lambda}$$

and

$$(7.2) \quad \int_\Omega [T|\nabla v_\lambda|^2 + D(\Delta v_\lambda)^2] dx \leq C.$$

Suppose that there is $\lambda_0 > 0$ and sequences $\{\lambda_i\}$ and $\{v_i\}$ with $\max_\Omega v_i = L - \epsilon_i$ such that $\lambda_i \rightarrow \lambda_0$, $\epsilon_i \rightarrow 0$ as $i \rightarrow \infty$. Making the transformation $w_i = L - v_i$, we see that w_i with $\min_\Omega w_i = \epsilon_i$ satisfies the problem

$$T\Delta w_i - D\Delta^2 w_i = \lambda_i w_i^{-2} \text{ in } \Omega, \quad w_i = L, \Delta w_i = 0 \text{ on } \partial\Omega.$$

Define $z_i = \Delta w_i$; then

$$(7.3) \quad -D\Delta z_i + Tz_i = \lambda_i w_i^{-2} \text{ in } \Omega, \quad z_i = 0 \text{ on } \partial\Omega.$$

It is known from (7.3) that $z_i(x) = \lambda_i \int_\Omega G_{T,D}(x, y) w_i^{-2}(y) dy$, where $G_{T,D}(x, y)$ is the Green's function of the operator $-D\Delta + TId$. Let $w_i(x_i) = \min_\Omega w_i$. Then $x_i \in \Omega_0 \subset\subset \Omega$. Setting $\tilde{w}_i(y) = \frac{w_i}{\epsilon_i}$ and $y = \lambda_i^{1/4} \epsilon_i^{-3/4} (x - x_i)$, we see that \tilde{w}_i with $\tilde{w}_i(0) = \min_{\Omega_i} \tilde{w}_i = 1$ and \tilde{w}_i satisfies the problem

$$(7.4) \quad \lambda_i^{-1/2} \epsilon_i^{3/2} T\Delta_y \tilde{w}_i - D\Delta_y^2 \tilde{w}_i = \tilde{w}_i^{-2} \text{ in } \Omega_i, \quad \tilde{w}_i = \frac{L_i}{\epsilon_i}, \Delta_y \tilde{w}_i = 0 \text{ on } \partial\Omega_i,$$

where $\Omega_i = \{y = \lambda_i^{1/4} \epsilon_i^{-3/4} (x - x_i) : x \in \Omega\}$. On the other hand,

$$\Delta_y \tilde{w}_i = \lambda_i^{-1/2} \epsilon_i^{1/2} \Delta_x w_i = \lambda_i^{1/2} \epsilon_i^{1/2} \int_\Omega G_{T,D}(x, \xi) w_i^{-2}(\xi) d\xi.$$

Note that $N = 2$ and $w_i \geq \epsilon_i$ in Ω . The Hölder inequality implies that

$$\begin{aligned} |\Delta_y \tilde{w}_i| &\leq C\epsilon_i^{1/2} \left(\int_{\Omega} [G_{T,D}(x, \xi)]^p d\xi \right)^{1/p} \left(\int_{\Omega} w_i^{-2} w_i^{2-2q}(\xi) d\xi \right)^{1/q} \\ &\leq C\epsilon_i^{1/2} \epsilon_i^{-2/p} \left(\int_{\Omega} [G_{T,D}(x, \xi)]^p d\xi \right)^{1/p} \left(\int_{\Omega} w_i^{-2}(\xi) d\xi \right)^{1/q} \\ &\leq C\epsilon_i^{\frac{1}{2} - \frac{2}{p}}, \end{aligned}$$

where we have applied (7.1). Choosing p sufficiently large, we see that

$$(7.5) \quad |\Delta_y \tilde{w}_i(y)| \rightarrow 0 \text{ for } y \in \Omega_i \text{ a.e. as } i \rightarrow \infty.$$

On the other hand, it follows from (7.4), (7.5), and the regularity of the operator $T\Delta - D\Delta^2$ that $\tilde{w}_i \rightarrow W$ in $C^4_{loc}(\mathbf{R}^2)$ as $i \rightarrow \infty$, where W with $W(0) = 1$ and $W \geq 1$ in \mathbf{R}^2 satisfies the equation

$$(7.6) \quad -D\Delta^2 W = W^{-2} \text{ in } \mathbf{R}^2, \quad W(0) = 1.$$

Meanwhile, (7.5) implies that $\Delta W = 0$ in \mathbf{R}^2 . This contradicts (7.6) and completes the proof of this lemma. \square

In the remainder of this section, we establish the existence of the second solution. Note that in the energy functional (1.1), the integral $\int_{\Omega} \frac{1}{L+u(x)} dx$ is not well defined for $u \in H^2(\Omega)$. Therefore, we do not have a good energy functional to work with. Our idea is to modify the nonlinearity so that the mountain-pass lemma works, and then show that the resulting solution has no singularity.

We first modify the nonlinearity. Since the nonlinearity $g(v) = \frac{1}{(L-v)^2}$ is singular at $v = L$, we need to consider a regularized C^1 nonlinearity $g_{\epsilon}(v)$, $0 < \epsilon < L$, of the following form:

$$g_{\epsilon}(v) = \begin{cases} \frac{1}{(L-v)^2}, & v \leq L - \epsilon, \\ \frac{1}{\epsilon^2} - \frac{(L-\epsilon)}{\epsilon^3} + \frac{1}{\epsilon^3(L-\epsilon)}v^2, & v > L - \epsilon. \end{cases}$$

For $\lambda \in (0, \lambda_c)$, we study the regularized semilinear elliptic problem:

$$(7.7) \quad -T\Delta v + D\Delta^2 v = \lambda g_{\epsilon}(v) \text{ in } \Omega, \quad v = \Delta v = 0 \text{ on } \partial\Omega.$$

From a variational viewpoint, the action functional associated to (7.7) is

$$J_{\epsilon, \lambda}(v) = \frac{1}{2} \int_{\Omega} [T|\nabla v|^2 + D(\Delta v)^2] dx - \lambda \int_{\Omega} G_{\epsilon}(v) dx, \quad v \in \mathcal{H},$$

where $G_{\epsilon}(v) = \int_{-\infty}^v g_{\epsilon}(s) ds$.

Fix now $0 < \epsilon < \tau_{\lambda}/4$, where τ_{λ} is as given in Lemma 7.2. The minimal solution \underline{v}_{λ} of (T_{λ}) is still a solution of (7.7) so that $\sigma_{1, \lambda}(\underline{v}_{\lambda}) > 0$. In order to motivate the choice of $g_{\epsilon}(v)$, we briefly sketch the proof of Theorem 7.1. First, we prove that \underline{v}_{λ} is a local minimum for $J_{\epsilon, \lambda}(v)$. Then, by the well-known mountain-pass theorem, we show the existence of a second solution $V_{\epsilon, \lambda}$ for (7.7). (A similar idea has been used in [6].) The subcritical growth

$$(7.8) \quad 0 \leq g_{\epsilon}(v) \leq C_{\epsilon}(1 + |v|^2)$$

and the inequality

$$(7.9) \quad 3G_\epsilon(v) \leq v g_\epsilon(v) \text{ for } v \geq L - \theta,$$

for some sufficiently small $\theta > 10\epsilon$ independent of ϵ , $C_\epsilon > 0$, will yield that $J_{\epsilon,\lambda}$ satisfies the Palais–Smale condition.

In order to complete the details of the proof of Theorem 7.1, we first need to show the following lemma.

LEMMA 7.3. *The minimal solution \underline{v}_λ of (T_λ) is a local minimum of $J_{\epsilon,\lambda}$ on \mathcal{H} .*

Proof. Since $\mathcal{H} \hookrightarrow C^\alpha(\bar{\Omega})$ for any $0 < \alpha < 1$, we need only to show that \underline{v}_λ is a local minimum of $J_{\epsilon,\lambda}$ in $C^\alpha(\bar{\Omega})$ for some $0 < \alpha < 1$. Indeed, since $\sigma_{1,\lambda}(\underline{v}_\lambda) > 0$, we have the following inequality:

$$(7.10) \quad \int_\Omega [T|\nabla\varphi|^2 + D(\Delta\varphi)^2]dx - 2\lambda \int_\Omega \frac{1}{(L - \underline{v}_\lambda)^3} \varphi^2 dx \geq \sigma_{1,\lambda} \int_\Omega \varphi^2 dx$$

for any $\varphi \in \mathcal{H}$, since $\underline{v}_\lambda \leq L - \tau_\lambda < L - \epsilon$ (see Lemma 7.2). Now, take any $\varphi \in \mathcal{H} \cap C^\alpha(\bar{\Omega})$ such that $\|\varphi\|_{C^\alpha} \leq \delta_\lambda$. Since $\underline{v}_\lambda \leq L - \tau_\lambda$, if $\delta_\lambda \leq \epsilon$, then $\underline{v}_\lambda + \varphi \leq L - \epsilon$, and we have that

$$\begin{aligned} & J_{\epsilon,\lambda}(\underline{v}_\lambda + \varphi) - J_{\epsilon,\lambda}(\underline{v}_\lambda) \\ &= \frac{1}{2} \int_\Omega [T|\nabla\varphi|^2 + D(\Delta\varphi)^2]dx + \int_\Omega [T\nabla\underline{v}_\lambda \cdot \nabla\varphi + D\Delta\underline{v}_\lambda \Delta\varphi]dx \\ &\quad - \lambda \int_\Omega \left(\frac{1}{L - \underline{v}_\lambda - \varphi} - \frac{1}{L - \underline{v}_\lambda} \right) \\ &\geq \frac{\sigma_{1,\lambda}}{2} \int_\Omega \varphi^2 - \lambda \int_\Omega \left(\frac{1}{L - \underline{v}_\lambda - \varphi} - \frac{1}{L - \underline{v}_\lambda} - \frac{\varphi}{(L - \underline{v}_\lambda)^2} - \frac{\varphi^2}{(L - \underline{v}_\lambda)^3} \right), \end{aligned}$$

where we have used (7.10). Since now

$$\left| \frac{1}{L - \underline{v}_\lambda - \varphi} - \frac{1}{L - \underline{v}_\lambda} - \frac{\varphi}{(L - \underline{v}_\lambda)^2} - \frac{\varphi^2}{(L - \underline{v}_\lambda)^3} \right| \leq C|\varphi|^3$$

for some $C > 0$, we have that

$$J_{\epsilon,\lambda}(\underline{v}_\lambda + \varphi) - J_{\epsilon,\lambda}(\underline{v}_\lambda) \geq \left(\frac{\sigma_{1,\lambda}}{2} - C\lambda\delta_\lambda \right) \int_\Omega \varphi^2 > 0$$

provided δ_λ is small enough. This proves that \underline{v}_λ is a local minimum of $J_{\epsilon,\lambda}$ in the C^α topology and completes the proof of this lemma. \square

Fix some ball $B_{2r} \subset \Omega$ of radius $2r$, $r > 0$. Take a cut-off function χ so that $\chi = 1$ on B_r and $\chi = 0$ outside B_{2r} . Let $w_\epsilon = (L - \epsilon)\chi \in \mathcal{H}$. We have that

$$J_{\epsilon,\lambda}(w_\epsilon) \leq \frac{(L - \epsilon)^2}{2} \int_\Omega [T|\nabla\chi|^2 + D(\Delta\chi)^2]dx - \frac{\lambda}{\epsilon^2} |B_r| \rightarrow -\infty$$

as $\epsilon \rightarrow 0$. Moreover, we can find for $\epsilon > 0$ small the inequality

$$(7.11) \quad J_{\epsilon,\lambda}(w_\epsilon) < J_{\epsilon,\lambda}(\underline{v}_\lambda).$$

Now fix $\epsilon > 0$ small enough in order that (7.11) holds, and define

$$c_{\epsilon,\lambda} = \inf_{\gamma \in \Gamma} \max_{v \in \gamma} J_{\epsilon,\lambda}(v),$$

where $\Gamma = \{\gamma : [0, 1] \rightarrow \mathcal{H}; \gamma \text{ continuous and } \gamma(0) = \underline{v}_\lambda, \gamma(1) = w_\epsilon\}$. We can then apply the mountain-pass theorem to get a solution $V_{\epsilon,\lambda}$ of (7.7), provided the Palais–Smale condition holds at level $c_{\epsilon,\lambda}$. The embedding theorem and the maximum principle imply that $V_{\epsilon,\lambda} > 0$ in Ω .

LEMMA 7.4. *Assume that $\{v_n\} \subset \mathcal{H}$ satisfies*

$$(7.12) \quad J_{\epsilon,\lambda_n}(v_n) \leq C, \quad J'_{\epsilon,\lambda_n}(v_n) \rightarrow 0 \text{ in } \mathcal{H}'$$

for $\lambda_n \rightarrow \lambda > 0$. Then the sequence $(v_n)_n$ is uniformly bounded in \mathcal{H} and therefore admits a convergent subsequence in \mathcal{H} .

Proof. By (7.12) we have that

$$\int_{\Omega} [T|\nabla v_n|^2 + D(\Delta v_n)^2] dx = \lambda_n \int_{\Omega} g_{\epsilon}(v_n)v_n dx + o(\|v_n\|_{\mathcal{H}})$$

as $n \rightarrow +\infty$, and then

$$\begin{aligned} C &\geq \frac{1}{2} \int_{\Omega} [T|\nabla v_n|^2 + D(\Delta v_n)^2] dx - \lambda_n \int_{\Omega} G_{\epsilon}(v_n) dx \\ &= \frac{1}{6} \int_{\Omega} [T|\nabla v_n|^2 + D(\Delta v_n)^2] dx + \lambda_n \int_{\Omega} \left(\frac{1}{3} v_n g_{\epsilon}(v_n) - G_{\epsilon}(v_n) \right) dx + o(\|v_n\|_{\mathcal{H}}) \\ &\geq \frac{1}{6} \int_{\Omega} [T|\nabla v_n|^2 + D(\Delta v_n)^2] dx + \lambda_n \int_{\{v_n \geq L-\theta\}} \left(\frac{1}{3} v_n g_{\epsilon}(v_n) - G_{\epsilon}(v_n) \right) dx \\ &\quad + o(\|v_n\|_{\mathcal{H}}) - C(\theta) \\ &\geq \frac{1}{6} \int_{\Omega} [T|\nabla v_n|^2 + D(\Delta v_n)^2] dx + o(\|v_n\|_{\mathcal{H}}) - C(\theta) \end{aligned}$$

in view of (7.9), where $C(\theta) > 0$ depends on θ but is independent of ϵ . Hence, $\sup_{n \in \mathbf{N}} \|v_n\|_{\mathcal{H}} < +\infty$.

The compactness of the embedding $\mathcal{H} \hookrightarrow C^{\alpha}(\bar{\Omega})$ for any $0 < \alpha < 1$ provides that, up to a subsequence, $v_n \rightarrow v$ weakly in \mathcal{H} and strongly in $C^{\alpha}(\bar{\Omega})$ for some $0 < \alpha < 1$ and some $v \in \mathcal{H}$. By (7.12) we get that $\int_{\Omega} [T|\nabla v|^2 + D(\Delta v)^2] dx = \lambda \int_{\Omega} g_{\epsilon}(v)v$, and then

$$\begin{aligned} &\int_{\Omega} [T|\nabla(v_n - v)|^2 + D(\Delta(v_n - v))^2] dx \\ &= T \left[\int_{\Omega} |\nabla v_n|^2 - \int_{\Omega} |\nabla v|^2 \right] + D \left[\int_{\Omega} (\Delta v_n)^2 - \int_{\Omega} (\Delta v)^2 \right] + o(1) \\ &= \lambda_n \int_{\Omega} g_{\epsilon}(v_n)v_n - \lambda \int_{\Omega} g_{\epsilon}(v)v + o(1) \rightarrow 0 \end{aligned}$$

as $n \rightarrow +\infty$. This completes the proof. \square

Proof of Theorem 7.1. We need only to show (i), and it is enough to show that for any fixed $\lambda > 0$, the mountain-pass solution $V_{\epsilon,\lambda}$ satisfies $V_{\epsilon,\lambda} \leq L - \epsilon$ in Ω .

Since $V_{\epsilon,\lambda} \in \mathcal{H}$, by the same argument as in Theorem 3.1, we easily see that

$$(7.13) \quad \int_{\Omega} g_{\epsilon}(V_{\epsilon,\lambda}) dx \leq C/\lambda,$$

where C is independent of ϵ . In fact, we see that

$$J_{\epsilon,\lambda}(V_{\epsilon,\lambda}) \leq \max_{v \in \gamma_0} J_{\epsilon,\lambda}(v),$$

where $\gamma_0 : [0, 1] \rightarrow \mathcal{H}$; $\gamma_0(v) = tv_\lambda + (1 - t)w_\epsilon$ for $t \in [0, 1]$. Thus,

$$J_{\epsilon,\lambda}(V_{\epsilon,\lambda}) \leq C,$$

where $C > 0$ is independent of ϵ . On the other hand, we see that

$$\begin{aligned} C &\geq \frac{1}{2} \int_{\Omega} [T|\nabla V_{\epsilon,\lambda}|^2 + D(\Delta V_{\epsilon,\lambda})^2] dx - \lambda_n \int_{\Omega} G_{\epsilon}(V_{\epsilon,\lambda}) dx \\ &= \frac{1}{6} \|V_{\epsilon,\lambda}\|_{\mathcal{H}}^2 + \lambda \int_{\Omega} \left(\frac{1}{3} V_{\epsilon,\lambda} g_{\epsilon}(V_{\epsilon,\lambda}) - G_{\epsilon}(V_{\epsilon,\lambda}) \right) dx \\ &\geq \frac{1}{6} \int_{\Omega} \|V_{\epsilon,\lambda}\|_{\mathcal{H}}^2 + \lambda \int_{\{V_{\epsilon,\lambda} \geq L - \theta\}} \left(\frac{1}{3} V_{\epsilon,\lambda} g_{\epsilon}(V_{\epsilon,\lambda}) - G_{\epsilon}(V_{\epsilon,\lambda}) \right) dx - C(\theta) \\ &\geq \frac{1}{6} \|V_{\epsilon,\lambda}\|_{\mathcal{H}}^2 dx - C(\theta). \end{aligned}$$

Thus,

$$\|V_{\epsilon,\lambda}\|_{\mathcal{H}} \leq C,$$

where $C > 0$ is independent of ϵ . The embedding $\mathcal{H} \hookrightarrow C^0(\overline{\Omega})$ implies $V_{\epsilon,\lambda} \leq C$ in Ω . By the moving plane argument, as in the proof of Theorem 3.1, we have that $V_{\epsilon,\lambda} \leq L - \theta$ in Ω_{δ} , where Ω_{δ} is as given in the proof of Theorem 3.1 and θ is as given by (7.9). This implies that (7.13) holds.

Let $W_{\epsilon} = \Delta V_{\epsilon,\lambda}$. Then W_{ϵ} satisfies the equation

$$-TW_{\epsilon} + D\Delta W_{\epsilon} = \lambda g_{\epsilon}(V_{\epsilon,\lambda}) \in L^1(\Omega).$$

Since $N = 2$, the Brezis–Merle inequality [4] implies that

$$\int_{\Omega} |W_{\epsilon}|^q dx \leq C \quad \forall q > 1,$$

where C is independent of ϵ . This also yields that

$$\|V_{\epsilon,\lambda}\|_{W^{2,q}(\Omega)} \leq C.$$

By choosing $q > 3$ sufficiently large, we see from the embedding $W_0^{2,q}(\Omega) \hookrightarrow C^{1+\frac{1}{2}}(\Omega)$ that $V_{\epsilon,\lambda} \leq C$ in Ω .

Now we show that $V_{\epsilon,\lambda} < L$ in Ω for ϵ sufficiently small. On the contrary, we suppose that there is a sequence $\{\epsilon_i\}$ with $\epsilon_i \rightarrow 0$ as $i \rightarrow \infty$ such that $\max_{\Omega} V_{\epsilon_i,\lambda} \geq L$. Denote $V_{\epsilon_i,\lambda}(x_i) = \max_{\Omega} V_{\epsilon_i,\lambda}$. By arguments similar to those in the proof of Lemma 5.1, we see that

$$V_{\epsilon_i,\lambda}(x_i) - V_{\epsilon_i,\lambda}(x) \leq C|x - x_i|^{3/2}.$$

Thus,

$$V_{\epsilon_i,\lambda}(x) \geq V_{\epsilon_i,\lambda}(x_i) - C|x - x_i|^{3/2} > L - \epsilon_i$$

provided that $|x - x_i| < (\epsilon_i/C)^{2/3}$. But

$$C \geq \int_{\Omega} g_{\epsilon_i}(V_{\epsilon_i,\lambda}) dx \geq \epsilon_i^{-2} \int_{\{|x-x_i| \leq (\epsilon_i/C)^{2/3}\}} dx = C\epsilon_i^{-2/3} \rightarrow \infty$$

as $i \rightarrow \infty$. This is a contradiction.

Now we claim that there exists $\delta > 0$ independent of ϵ such that

$$V_{\epsilon,\lambda} \leq L - \delta \text{ in } \Omega$$

for ϵ sufficiently small. On the contrary, there are sequences $\{\epsilon_i\}$ and $\{V_i\} \equiv \{V_{\epsilon_i,\lambda}\}$ with $\epsilon_i \rightarrow 0$ as $i \rightarrow \infty$ such that $\max_{\Omega} V_i = L - \xi_i$ and $\xi_i \rightarrow 0$ as $i \rightarrow \infty$. Set $Z_i = L - V_i$. Then $Z_i(x_i) := \min_{\Omega} Z_i = \xi_i$ and Z_i satisfies

$$T\Delta Z_i - D\Delta^2 Z_i = \lambda h_i(Z_i) \text{ in } \Omega, \quad Z_i = T, \quad \Delta Z_i = 0 \text{ on } \partial\Omega,$$

where

$$h_i(Z_i) = \begin{cases} \frac{1}{Z_i^2}, & Z_i \geq \epsilon_i, \\ \frac{1}{\epsilon_i^2} + \frac{2(\epsilon_i - Z_i)}{\epsilon_i^3} + \frac{(\epsilon_i - Z_i)^2}{\epsilon_i^3(L - \epsilon_i)}, & Z_i < \epsilon_i. \end{cases}$$

Making the transformations $\tilde{Z}_i(y) = Z_i/\xi_i$ and $y = \xi_i^{-3/4}(x - x_i)$, we see that $\tilde{Z}_i(0) = \min_{\Omega} \tilde{Z}_i = 1$ and \tilde{Z}_i satisfies the problem

$$(7.14) \quad \xi_i^{3/2} \Delta_y \tilde{Z}_i - D\Delta_y^2 \tilde{Z}_i = \lambda \tilde{h}_i(\tilde{Z}_i) \text{ in } \tilde{\Omega}_i, \quad \tilde{Z}_i = T/\xi_i, \quad \Delta_y \tilde{Z}_i = 0 \text{ on } \partial\tilde{\Omega}_i,$$

where $\tilde{\Omega}_i = \{y = \xi_i^{-3/4}(x - x_i) : x \in \Omega\}$ and

$$\tilde{h}_i(\tilde{Z}_i) = \begin{cases} \frac{1}{\tilde{Z}_i^2}, & \tilde{Z}_i \geq \frac{\epsilon_i}{\xi_i}, \\ 3\left(\frac{\xi_i}{\epsilon_i}\right)^2 - 2\left(\frac{\xi_i}{\epsilon_i}\right)^3 \tilde{Z}_i + \frac{\xi_i^2}{\epsilon_i(L - \epsilon_i)} - 2\left(\frac{\xi_i^3}{\epsilon_i^2(L - \epsilon_i)}\right) \tilde{Z}_i + \left(\frac{\xi_i^4}{\epsilon_i^3(L - \epsilon_i)}\right) \tilde{Z}_i^2, & \tilde{Z}_i < \frac{\epsilon_i}{\xi_i}. \end{cases}$$

We consider two cases for $\{\frac{\epsilon_i}{\xi_i}\}$ (we can choose subsequences if necessary):

- (i) There is $0 < A < \infty$ such that $\frac{\epsilon_i}{\xi_i} \leq A$ for all i ,
- (ii) $\frac{\epsilon_i}{\xi_i} \rightarrow \infty$ as $i \rightarrow \infty$.

For the first case, we have that there is $0 \leq A_1 \leq A$ such that $\lim_{i \rightarrow \infty} \frac{\epsilon_i}{\xi_i} = A_1$.

If $A_1 < 1$, since $\tilde{Z}_i \geq 1$, we have that $\tilde{h}_i = \tilde{Z}_i^{-2} \leq 1$ in $\tilde{\Omega}_i$ for i sufficiently large. If $1 \leq A_1 \leq A$, we also have that $\tilde{h}_i \leq C$ in $\tilde{\Omega}_i$ for i sufficiently large, where $0 < C < \infty$ is independent of i . Moreover, for $q > 3$,

$$\begin{aligned} \int_{\tilde{\Omega}_i} |\Delta_y \tilde{Z}_i|^q dy &= \xi_i^{q/2} \int_{\tilde{\Omega}_i} |\Delta_x Z_i|^q dy \\ &= \xi^{q-3} \int_{\Omega} |\Delta_x Z_i|^q dx \rightarrow 0 \end{aligned}$$

as $i \rightarrow \infty$. Thus, the regularity of Δ^2 implies that $\tilde{Z}_i \rightarrow \tilde{Z}$ in $C_{loc}^3(\mathbf{R}^2)$ with $\tilde{Z}(0) = \min_{\mathbf{R}^2} \tilde{Z} = 1$, and \tilde{Z} satisfies the equation

$$-D\Delta^2 \tilde{Z} = \lambda \tilde{Z}^{-2} \text{ in } \mathbf{R}^2$$

provided $A_1 \leq 1$ and the equation

$$D\Delta^2 \tilde{Z} = \lambda \tilde{h}(\tilde{Z}) \text{ in } \mathbf{R}^2$$

provided $1 < A_1 \leq A$, where

$$\tilde{h}(\tilde{Z}) = \begin{cases} \tilde{Z}^2, & \tilde{Z} \geq A_1, \\ \frac{3}{A_1^2} - \frac{2}{A_1^3} \tilde{Z}, & \tilde{Z} < A_1. \end{cases}$$

Moreover, for any large ball B_R of \mathbf{R}^2 , $\int_{B_R} |\Delta \tilde{Z}|^q(y) dy = 0$. This is impossible.

For the second case, we see that $\xi_i = o(\epsilon_i)$ for i sufficiently large. Thus, $Z_i(x_i) = \xi_i = o(\epsilon_i)$. Noting that $\int_{\Omega} |\Delta Z_i|^q dx \leq C$, we see that $Z_i \in W_0^{2,q}(\Omega)$. The embedding $W_0^{2,q}(\Omega) \hookrightarrow C^{1+\frac{1}{2}}(\Omega)$ gives

$$|Z_i(x)| \leq Z_i(x_i) + C|x - x_i|^{3/2} < \epsilon_i$$

provided $|x - x_i| \leq (\frac{\epsilon_i}{2C})^{2/3}$. Thus

$$C \geq \int_{\Omega} h_i(Z_i) dx \geq \frac{1}{\epsilon_i^2} \int_{Z_i < \epsilon_i} dx \geq C\epsilon_i^{-2/3} \rightarrow \infty$$

as $i \rightarrow \infty$. This is also a contradiction. Therefore,

$$V_{\epsilon,\lambda} \leq L - \delta \text{ in } \Omega,$$

where $\delta > 0$ is independent of ϵ . This also implies that $V_{\epsilon,\lambda}$ is a solution of (T_{λ}) . This completes the proof of (i) of Theorem 7.1. \square

8. The asymptotic behavior of the mountain-pass solution as $\lambda \rightarrow 0$.

In this section we will study the asymptotic behavior of the mountain-pass solution V_{λ} obtained in Theorem 7.1 as $\lambda \rightarrow 0$.

LEMMA 8.1.

$$\sigma_{1,\lambda}(V_{\lambda}) < 0 \text{ for } 0 < \lambda < \lambda_c.$$

Proof. Let \underline{v}_{λ} be the minimal solution of (T_{λ}) so that $V_{\lambda} \geq \underline{v}_{\lambda}$. If the linearization around V_{λ} had nonnegative first eigenvalue, then Lemma 4.1 would also yield $V_{\lambda} \leq \underline{v}_{\lambda}$ so that \bar{v}_{λ} and V_{λ} would necessarily coincide, which would be a contradiction. \square

THEOREM 8.2.

$$(8.1) \quad \max_{\Omega} V_{\lambda} \rightarrow L \text{ as } \lambda \rightarrow 0.$$

Moreover,

$$(8.2) \quad \lim_{\lambda \rightarrow 0^+} \frac{[\min_{\Omega}(L - V_{\lambda})]^3}{\lambda} = 0.$$

Proof. Suppose that there are sequences $\{\lambda_i\}$ and $\{V_i\} \equiv \{V_{\lambda_i}\}$ such that $\lambda_i \rightarrow 0$ as $i \rightarrow \infty$ and $\max_{\Omega} V_i \leq L - \delta$, where $0 < \delta < L$ is independent of i . Then it follows from the equation of V_i that $V_i \rightarrow 0$ in $C^0(\bar{\Omega})$ as $i \rightarrow \infty$ (we can choose subsequences if necessary). This contradicts the fact that $\sigma_{1,\lambda_i}(V_i) < 0$ for all i . Thus, (8.1) holds.

By Theorem 3.1, we see that

$$\lambda \int_{\Omega} (L - V_{\lambda})^{-2} dx + \int_{\Omega} |\Delta V_{\lambda}|^2 \leq C$$

for any λ sufficiently small, where C is independent of λ . Since ΔV_{λ} satisfies

$$-D\Delta(\Delta V_{\lambda}) = D\Delta V_{\lambda} + \frac{\lambda}{(L - V_{\lambda})^2} \in L^1(\Omega), \quad \Delta V_{\lambda} = 0 \text{ on } \partial\Omega,$$

by the Brezis–Merle inequality [4], we have, for any $q > 1$,

$$(8.3) \quad \int_{\Omega} |\Delta V_{\lambda}|^q dx \leq C$$

for any λ sufficiently small.

Let $V_\lambda(x_\lambda) = \max_\Omega V_\lambda$. Setting $W_\lambda = L - V_\lambda$, we see that $\xi_\lambda := W_\lambda(x_\lambda) = \min_\Omega W_\lambda$ and $\xi_\lambda \rightarrow 0$ as $\lambda \rightarrow 0$. Now we claim that

$$(8.4) \quad \lim_{\lambda \rightarrow 0} \frac{\xi_\lambda^3}{\lambda} = 0.$$

Suppose not; there are sequences $\{\lambda_i\}$ and $\{\xi_i\}$ with $\lambda_i \rightarrow 0$ as $i \rightarrow \infty$ such that $\frac{\xi_i^3}{\lambda_i} \rightarrow C > 0$ or $\frac{\xi_i^3}{\lambda_i} \rightarrow \infty$ as $i \rightarrow \infty$.

We first consider the case that $\frac{\xi_i^3}{\lambda_i} \rightarrow \infty$ as $i \rightarrow \infty$. Then, defining $\hat{W}_i = W_i/\xi_i$, we see that \hat{W}_i satisfies the problem

$$T\Delta\hat{W}_i - D\Delta^2\hat{W}_i = \frac{\lambda_i}{\xi_i^3}\hat{W}_i^{-2} \text{ in } \Omega, \quad \hat{W}_i = L/\xi_i, \quad \Delta\hat{W}_i = 0 \text{ on } \partial\Omega.$$

Since $\hat{W}_i \geq 1$, we see that $\hat{W}_i \rightarrow \hat{W}$ in $C_{loc}^3(\Omega)$ as $i \rightarrow \infty$ and \hat{W} with $\hat{W}(0) = \min_\Omega \hat{W} = 1$ satisfies the equation

$$(8.5) \quad T\Delta\hat{W} - D\Delta^2\hat{W} = 0 \text{ in } \Omega, \quad \hat{W} = \infty, \quad \Delta\hat{W} = 0 \text{ on } \partial\Omega.$$

Setting $Z = \Delta\hat{W}$, we see from (8.5) that

$$TZ - D\Delta Z = 0 \text{ in } \Omega, \quad Z = 0 \text{ on } \partial\Omega.$$

The strong maximum principle then implies that $Z \equiv 0$ in Ω and hence $\Delta\hat{W} \equiv 0$ in Ω . The maximum principle then implies that $\hat{W} \equiv 1$ in Ω , which is a contradiction.

Now we consider the case that $\lim_{i \rightarrow \infty} \frac{\xi_i^3}{\lambda_i} \rightarrow C > 0$. Defining $\hat{W}_i = W_i/\xi_i$ again, we see that \hat{W}_i satisfies the problem

$$(8.6) \quad T\Delta\hat{W}_i - D\Delta^2\hat{W}_i = \frac{\lambda_i}{\xi_i^3}\hat{W}_i^{-2} \text{ in } \Omega, \quad \hat{W}_i = L/\xi_i, \quad \Delta\hat{W}_i = 0 \text{ on } \partial\Omega.$$

Setting $\hat{Z}_i = \Delta\hat{W}_i$, we see that \hat{Z}_i satisfies the problem

$$(8.7) \quad T\hat{Z}_i - D\Delta\hat{Z}_i = \frac{\lambda_i}{\xi_i^3}\hat{W}_i^{-2} \text{ in } \Omega, \quad \hat{Z}_i = 0 \text{ on } \partial\Omega.$$

Therefore,

$$\hat{Z}_i = \frac{\lambda_i}{\xi_i^3} \int_\Omega G_{T,D}(x,y)\hat{W}_i^{-2}(y)dy,$$

and hence $|\hat{Z}_i| \leq C$, where $C > 0$ is independent of i . We now obtain from the regularity of Δ^2 and (8.6) that $\hat{W}_i \rightarrow \hat{W}$ in $C_{loc}^3(\Omega)$ and \hat{W} satisfies the equation

$$T\Delta\hat{W} - D\Delta^2\hat{W} = \frac{1}{C}\hat{W}^{-2} \text{ in } \Omega, \quad \hat{W} = \infty, \quad \Delta\hat{W} = 0 \text{ on } \partial\Omega.$$

On the other hand, we see from (8.7) that $\hat{Z}_i \rightarrow \hat{Z}$ in $C^1(\bar{\Omega})$ as $i \rightarrow \infty$ and $\hat{Z} \equiv \Delta\hat{W}$ satisfies the problem

$$T\hat{Z} - D\Delta\hat{Z} = \frac{1}{C}\hat{W}^{-2} \text{ in } \Omega, \quad \hat{Z} = 0 \text{ on } \partial\Omega.$$

Since we easily know that $\Delta\hat{W} \leq C$ on $\bar{\Omega}$ and hence $\Delta(\hat{W} - C\rho) \leq 0$ in Ω , where $-\Delta\rho = 1$ in Ω and $\rho = 0$ on $\partial\Omega$, the maximum principle implies that \hat{W} cannot be ∞ on $\partial\Omega$. Thus, (8.2) holds. This completes the proof of Theorem 8.2. \square

Acknowledgments. This paper was done while the first author was visiting the Department of Mathematics at the Chinese University of Hong Kong. He would like to thank the Department for its hospitality. We thank Professor Michael Ward and Professor Dong Ye for useful discussions. We also thank the referees for their valuable suggestions.

REFERENCES

- [1] S. AGMON, A. DOUGLIS, AND L. NIRENBERG, *Estimates near the boundary for solutions of elliptic partial differential equations, satisfying general boundary conditions, I*, Comm. Pure Appl. Math., 12 (1959), pp. 623–727.
- [2] A. L. BERTOZZI AND M. C. PUGH, *Long-wave instabilities and saturation in thin film equations*, Comm. Pure Appl. Math., 51 (1998), pp. 625–661.
- [3] A. L. BERTOZZI AND M. C. PUGH, *Finite-time blow-up of solutions of some long-wave unstable thin film equations*, Indiana Univ. Math. J., 49 (2000), pp. 1323–1366.
- [4] H. BREZIS AND F. MERLE, *Uniform estimates and blow-up behavior for solutions of $-\Delta u = V(x)e^u$ in two dimensions*, Comm. Partial Differential Equations, 16 (1991), pp. 1223–1254.
- [5] J. P. BURELBACH, S. G. BANKOFF, AND S. H. DAVIS, *Nonlinear stability of evaporating /condensing liquid films*, J. Fluid Mech., 195 (1998), pp. 463–494.
- [6] M. G. CRANDALL AND P. H. RABINOWITZ, *Some continuation and variational methods for positive solutions of nonlinear elliptic eigenvalue problems*, Arch. Ration. Mech. Anal., 58 (1975), pp. 207–218.
- [7] P. ESPOSITO, N. GHOUSSEB, AND Y. GUO, *Compactness along the branch of semi-stable and unstable solutions for an elliptic problem with a singular nonlinearity*, Comm. Pure Appl. Math., 60 (2008), pp. 1731–1768.
- [8] G. FLORES, G. A. MERCADO, AND J. A. PELESKO, *Dynamics and touchdown in electrostatic MEMS*, in Proceedings of ICMENS 2003, IEEE Computer Society, Washington, DC, 2003, pp. 182–187.
- [9] F. GAZZOLA AND H.-CH. GRUNAU, *Critical dimensions and higher order Sobolev inequalities with remainder terms*, NoDEA Nonlinear Differential Equations Appl., 8 (2001), pp. 35–44.
- [10] N. GHOUSSEB AND Y. GUO, *On the partial differential equations of electrostatic MEMS devices: Stationary case*, SIAM J. Math. Anal., 38 (2007), pp. 1423–1449.
- [11] N. GHOUSSEB AND Y. GUO, *On the partial differential equations of electrostatic MEMS devices II: Dynamic case*, NoDEA Nonlinear Differential Equations Appl., 15 (2008), pp. 115–145.
- [12] Y. GUO, *On the partial differential equations of electrostatic MEMS devices III: Refined touchdown behavior*, J. Differential Equations, 244 (2008), pp. 2277–2309.
- [13] Z. M. GUO AND J. R. L. WEBB, *Large and small solutions of a class of quasilinear elliptic eigenvalue problems*, J. Differential Equations, 180 (2002), pp. 1–50.
- [14] Z. M. GUO AND J. C. WEI, *Hausdorff dimension of ruptures for solutions of a semilinear elliptic equation with singular nonlinearity*, Manuscripta Math., 120 (2006), pp. 193–209.
- [15] Z. M. GUO AND J. C. WEI, *Symmetry of nonnegative solutions of a semilinear elliptic equation with singular nonlinearity*, Proc. Roy. Soc. Edinburgh Sect. A, 137 (2007), pp. 963–994.
- [16] Z. M. GUO AND J. C. WEI, *Infinitely many turning points for an elliptic problem with a singular nonlinearity*, J. London Math. Soc., 78 (2008), pp. 21–35.
- [17] Y. GUO, Z. PAN, AND M. J. WARD, *Touchdown and pull-in voltage behavior of a MEMS device with varying dielectric properties*, SIAM J. Appl. Math., 66 (2005), pp. 309–338.
- [18] C. C. HWANG, C. K. LIN, AND W. Y. UEN, *A nonlinear three-dimensional rupture theory of thin liquid films*, J. Colloid Interf. Sci., 190 (1997), pp. 250–252.
- [19] R. S. LAUGESEN AND M. C. PUGH, *Properties of steady states for thin film equations*, European J. Appl. Math., 11 (2000), pp. 293–351.
- [20] R. S. LAUGESEN AND M. C. PUGH, *Energy levels of steady-states for thin-film-type equations*, J. Differential Equations, 182 (2002), pp. 377–415.
- [21] R. S. LAUGESEN AND M. C. PUGH, *Linear stability of steady states for thin film and Cahn–Hilliard type equations*, Arch. Ration. Mech. Anal., 154 (2000), pp. 3–51.
- [22] F. H. LIN AND Y. S. YANG, *Nonlinear non-local elliptic equation modelling electrostatic actuation*, Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci., 463 (2007), pp. 1323–1337.
- [23] A. LINDSAY AND M. WARD, *Asymptotics of some nonlinear eigenvalue problems for a MEMS capacitor. Part I: Fold point asymptotics*, Methods Anal. Appl., to appear.

- [24] J. A. PELESKO, *Mathematical modeling of electrostatic MEMS with tailored dielectric properties*, SIAM J. Appl. Math., 62 (2002), pp. 888–908.
- [25] J. A. PELESKO AND D. H. BERNSTEIN, *Modeling MEMS and NEMS*, Chapman & Hall/CRC Press, Boca Raton, FL, 2002.
- [26] W. C. TROY, *Symmetry properties in systems of semilinear elliptic equations*, J. Differential Equations, 42 (1981), pp. 400–413.
- [27] R. C. A. M. VAN DER VORST, *Best constant for the embedding of the space $H^2 \cap H_0^1(\Omega)$ into $L^{2N/(N-4)}(\Omega)$* , Differential Integral Equations, 6 (1993), pp. 259–276.

STABILITY AND SYMMETRY IN THE NAVIER PROBLEM FOR THE ONE-DIMENSIONAL WILLMORE EQUATION*

KLAUS DECKELNICK[†] AND HANS-CHRISTOPH GRUNAU[†]

Abstract. We consider the one-dimensional Willmore equation subject to Navier boundary conditions; i.e., the position and the curvature are prescribed on the boundary. In a previous work, explicit symmetric solutions to symmetric data have been constructed. Within a certain range of boundary curvatures one has precisely two symmetric solutions, while for boundary curvatures outside the closure of this range there are none. The solutions are ordered; one is “small,” and the other is “large.” In the first part of this paper we address the stability problem and show that the small solution is (linearized) stable in the whole open range of admissible boundary curvatures, while the large one is unstable and has Morse index 1. A second goal is to investigate whether the small solution is minimal for the corresponding Willmore functional. It turns out that for a certain subrange of admissible boundary curvatures the small solution is the unique minimum, while for curvatures outside that range the minimum is not attained. As a byproduct of our argument we show that for any admissible function there exists a symmetric function with smaller Willmore energy.

Key words. Willmore equation, Navier boundary conditions, stability, Morse index, symmetry

AMS subject classifications. 53C21, 34B15, 35J65, 35B35

DOI. 10.1137/07069033X

1. Introduction. Recently, Willmore surfaces (see [19]) and the related flow have attracted quite some attention; see, e.g., [1, 8, 9, 10, 13, 17, 18], [3] for numerical studies, and [15, 5] for elastic curves, which are the one-dimensional analogues. The mentioned work is concerned with closed surfaces and curves, while only very few results concerning boundary value problems are available. Quite recently, Schätzle [16] considered Willmore surfaces with a boundary, which are subject to the constraint to be submanifolds of \mathbb{S}^n and which satisfy Dirichlet-type boundary conditions.

In order to gain some more insight into general boundary conditions for the “free” Willmore equation, in [4] we had a look at the one-dimensional case, where, in some situations, almost explicit solutions can be found for suitable boundary value problems. For further background information and references, see [4] and also [14]. In [4], we were interested in Willmore graphs and studied among others the *Navier boundary value problem* with symmetric data $\alpha \in \mathbb{R}$ for the *one-dimensional Willmore equation*:

$$(1.1) \quad \begin{cases} \frac{1}{\sqrt{1+u'(x)^2}} \frac{d}{dx} \left(\frac{\kappa'(x)}{\sqrt{1+u'(x)^2}} \right) + \frac{1}{2} \kappa^3(x) = 0, & x \in (0, 1), \\ u(0) = u(1) = 0, & \kappa(0) = \kappa(1) = -\alpha. \end{cases}$$

Here

$$(1.2) \quad \kappa(x) = \frac{d}{dx} \left(\frac{u'(x)}{\sqrt{1+u'(x)^2}} \right) = \frac{u''(x)}{(1+u'(x)^2)^{3/2}}$$

*Received by the editors May 2, 2007; accepted for publication (in revised form) October 2, 2008; published electronically January 21, 2009.

<http://www.siam.org/journals/sima/40-5/69033.html>

[†]Fakultät für Mathematik, Otto-von-Guericke-Universität Magdeburg, Postfach 4120, D-39016 Magdeburg, Germany (Klaus.Deckelnick@mathematik.uni-magdeburg.de, Hans-Christoph.Grunau@mathematik.uni-magdeburg.de).

denotes the curvature of the graph of u at the point $(x, u(x))$. Solutions of (1.1) are critical points of the modified one-dimensional Willmore functional

(1.3)

$$\tilde{W}_\alpha(u) = \int_{\text{graph}(u)} (\kappa(x)^2 + 2\alpha\kappa(x)) ds(x) = \int_0^1 (\kappa(x)^2 + 2\alpha\kappa(x)) \sqrt{1 + u'(x)^2} dx,$$

with $u \in H^2(0, 1) \cap H_0^1(0, 1)$; see [4, section 2]. The boundary conditions $u(0) = u(1) = 0$ are formulated by working in the space H_0^1 , while the curvature boundary conditions $\kappa(0) = \kappa(1) = -\alpha$ arise as natural boundary conditions since also the admissible testing functions have only to be in $H^2 \cap H_0^1$. By reflection it is sufficient to consider

$$\alpha \geq 0.$$

Our first results are concerned with solutions of (1.1) which are symmetric about $x = 1/2$. This class is particularly important for two reasons: on the one hand, it is possible to associate with each function in $H^2 \cap H_0^1$ a symmetric function with the same or smaller Willmore energy. The corresponding construction is nontrivial and will be described in the proof of Theorem 1.5. Moreover, exploiting an observation due to Euler [6, p. 234, line 13], we can give symmetric solutions more or less in closed form. This issue was discussed in detail in [4], where, among other things, the following result was proved.

PROPOSITION 1.1 (see [4, Theorem 1]). *There exists $\alpha_{\max} = 1.343799725\dots$ such that for $0 < \alpha < \alpha_{\max}$, the Navier boundary value problem (1.1) has precisely two smooth (graph) solutions u in the class of smooth functions that are symmetric around $x = \frac{1}{2}$. If $\alpha = \alpha_{\max}$, one has precisely one such solution, for $\alpha = 0$ one has only the trivial solution, and for $\alpha > \alpha_{\max}$ no such solutions exist.*

Both solutions are positive, and one of these solutions is larger than the other. The small solutions are ordered with respect to α , while the large ones become smaller for increasing α ; see Figure 1.1. For the bifurcation diagram, see Figure 1.2.

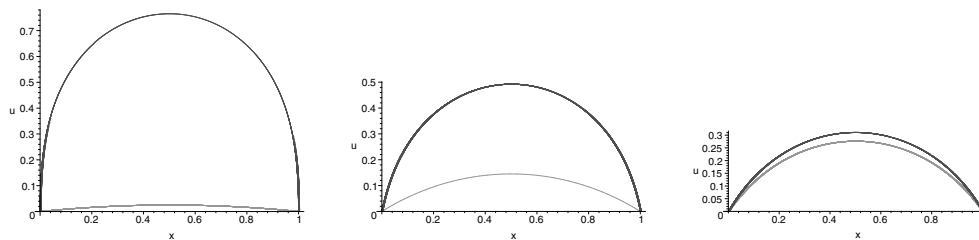


FIG. 1.1. Solutions of the Navier boundary value problem (1.1) for $\alpha = 0.2$, $\alpha = 1$, and $\alpha = 1.34$ (left to right). See [4, Figure 1].

It is an obvious conjecture that for $0 \leq \alpha < \alpha_{\max}$ the small solutions are (linearized) stable. This property was left open in [4], and to prove it is the first goal of this paper.

THEOREM 1.2. *Assume that $0 \leq \alpha < \alpha_{\max}$ and that u is the symmetric small solution of the Navier boundary value problem (1.1). Then, this solution is linearized stable; i.e., the spectrum of the (self-adjoint) linearization of (1.1) around u is contained in $(0, \infty)$.*

That these linearizations are the second variation of the functional \tilde{W}_α proves that the small solution is a local minimum of the functional \tilde{W}_α in $H^2 \cap H_0^1(0, 1)$.

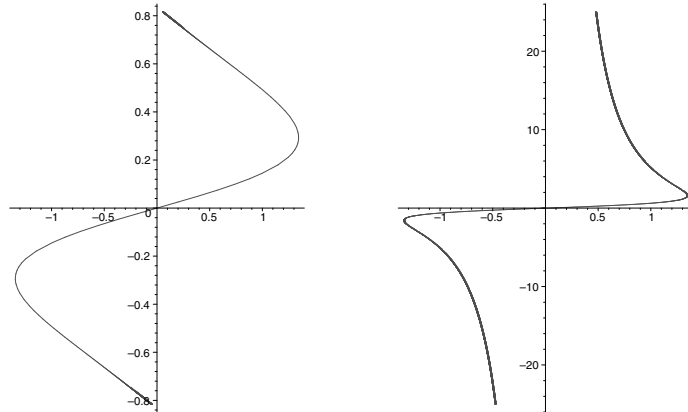


FIG. 1.2. Bifurcation diagram for (1.1): The extremals value of the solution $u(1/2)$ (left) and of the derivative $u'(0)$ (right) plotted over α . See [4, Figure 2].

Furthermore, we will show that on $0 < \alpha < \alpha_{\max}$, the large solutions are unstable. More precisely, we prove that the following holds.

THEOREM 1.3. *Assume that $0 < \alpha < \alpha_{\max}$ and that u is the symmetric large solution of the Navier boundary value problem (1.1). Then, this solution is unstable and has Morse index 1; i.e., one eigenvalue of the (self-adjoint) linearization of (1.1) is negative, while the remaining spectrum is contained in $(0, \infty)$.*

We emphasize that no symmetry assumptions are made in the discussion of the linearizations of (1.1).

A further important question is whether the small solutions are not only a local but also a global minimum of the functional \tilde{W}_α .

THEOREM 1.4. *There exists $\alpha^* = 1.132372323\dots \in (0, \alpha_{\max})$ such that for $0 \leq \alpha \leq \alpha^*$ the small solution u is the unique global minimum of the functional \tilde{W}_α in the class $H^2 \cap H_0^1(0, 1)$. If $\alpha^* < \alpha \leq \alpha_{\max}$, the infimum of \tilde{W}_α in $H^2 \cap H_0^1(0, 1)$ is not attained, and in that case*

$$\inf_{v \in H^2 \cap H_0^1(0,1)} \tilde{W}_\alpha(v) = \left(\int_{\mathbb{R}} \frac{1}{(1 + \tau^2)^{5/4}} d\tau \right)^2 - 2\alpha\pi.$$

The main idea of proving Theorem 1.4 consists in reducing the minimization of \tilde{W}_α over $H^2 \cap H_0^1(0, 1)$ to the minimization of a function of two variables. As a byproduct of this approach we shall see that the infimum of the Willmore energy in $H^2 \cap H_0^1(0, 1)$ coincides with the infimum in the subspace M of functions that are symmetric about $x = 1/2$; i.e., for every function in $H^2 \cap H_0^1(0, 1)$, there exists a symmetric function with the same or smaller Willmore energy. This is remarkable since we deal with a fourth order problem and the well-known symmetrization procedures do not apply.

THEOREM 1.5. *Let M be the class of functions in $H^2 \cap H_0^1(0, 1)$, which are symmetric about $x = 1/2$. Then we have*

$$\inf_{v \in H^2 \cap H_0^1(0,1)} \tilde{W}_\alpha(v) = \inf_{v \in M} \tilde{W}_\alpha(v).$$

2. Linearized stability. To prove Theorem 1.2 we describe in more detail how the symmetric solutions to (1.1) were obtained in [4].

In what follows, the function

$$(2.1) \quad G : \mathbb{R} \rightarrow \left(-\frac{c_0}{2}, \frac{c_0}{2}\right), \quad G(s) := \int_0^s \frac{1}{(1 + \tau^2)^{5/4}} d\tau,$$

$$c_0 = \int_{\mathbb{R}} \frac{1}{(1 + \tau^2)^{5/4}} d\tau = \mathcal{B}\left(\frac{1}{2}, \frac{3}{4}\right) = 2.396280469 \dots,$$

plays a crucial role. It is straightforward to see that G is strictly increasing and bijective with $G'(s) > 0$. So, also the inverse function

$$(2.2) \quad G^{-1} : \left(-\frac{c_0}{2}, \frac{c_0}{2}\right) \rightarrow \mathbb{R}$$

is strictly increasing, bijective, and smooth with $G^{-1}(0) = 0$.

LEMMA 2.1 (see [4, Lemma 4]). *Let $u \in C^4([0, 1])$ be a function symmetric about $x = 1/2$. Then u solves the Willmore equation in (1.1) iff there exists $c \in (-c_0, c_0)$ such that*

$$(2.3) \quad \forall x \in [0, 1] : \quad u'(x) = G^{-1}\left(\frac{c}{2} - cx\right).$$

For the curvature, one has that

$$(2.4) \quad \kappa(x) = -\frac{c}{\sqrt[4]{1 + G^{-1}\left(\frac{c}{2} - cx\right)^2}}.$$

Moreover, if we additionally assume that $u(0) = u(1) = 0$, then one has

$$(2.5) \quad u(x) = \frac{2}{c\sqrt[4]{1 + G^{-1}\left(\frac{c}{2} - cx\right)^2}} - \frac{2}{c\sqrt[4]{1 + G^{-1}\left(\frac{c}{2}\right)^2}} \quad (c \neq 0).$$

In order to solve the Navier boundary value problem (1.1), in [4], we had to study the function

$$(2.6) \quad h : (-c_0, c_0) \rightarrow \mathbb{R}, \quad h(c) = \frac{c}{\sqrt[4]{1 + G^{-1}\left(\frac{c}{2}\right)^2}}$$

and the equation $h(c) = \alpha$. See Figure 2.1. The range of h is precisely the set of α , for which the Navier boundary value problem (1.1) has a smooth symmetric graph solution. The number of solutions c of the equation $\alpha = h(c)$ is the number of such solutions of the boundary value problem.

LEMMA 2.2 (see [4, Lemma 6]). *We have $h > 0$ in $(0, c_0)$, $h < 0$ in $(-c_0, 0)$, $\lim_{c \nearrow c_0} h(c) = \lim_{c \searrow -c_0} h(c) = 0$. The function h is odd and has precisely one local maximum in $c_{\max} = 1.840428142 \dots$ and one local minimum in $c_{\min} = -c_{\max}$. The corresponding value is $\alpha_{\max} = h(c_{\max}) = 1.343799725 \dots$*

The small solutions correspond precisely to $c \in (0, c_{\max})$ and the large ones to $c \in (c_{\max}, c_0)$. Let us fix $c \in (0, c_0)$ with corresponding $\alpha = h(c)$ and solution u given by (2.5). First we have to calculate the linearization of (1.1) around u , i.e., the second variation of the modified Willmore functional \tilde{W}_α in u .

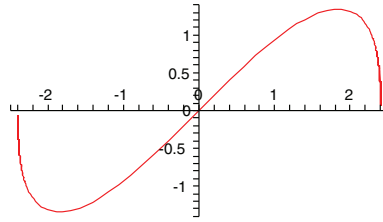


FIG. 2.1. The function $c \mapsto h(c)$. According to [4, Lemma 4], solutions to (1.1) are given by solving $h(c) = \alpha$.

LEMMA 2.3. We have

$$D^2\tilde{W}_\alpha(u)(\varphi, \eta) = 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx + 5 \int_0^1 \frac{1-u'(x)^2}{(1+u'(x)^2)^{3/2}} \kappa(x)^2 \varphi'(x)\eta'(x) dx + 6\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1, \quad \varphi, \eta \in H^2 \cap H_0^1(0, 1).$$

Proof. According to [4, Lemma 2 and Corollary 1], the first variation of $\tilde{W}_\alpha(u)$ is given by

$$D\tilde{W}_\alpha(u)(\varphi) = 2 \int_0^1 \frac{u''(x)\varphi''(x)}{(1+u'(x)^2)^{5/2}} dx - 5 \int_0^1 \frac{u'(x)u''(x)^2\varphi'(x)}{(1+u'(x)^2)^{7/2}} dx + 2\alpha \left[\frac{\varphi'(x)}{1+u'(x)^2} \right]_0^1, \quad \varphi \in H^2 \cap H_0^1(0, 1).$$

In order to obtain the second derivative, we also consider $\eta \in H^2 \cap H_0^1(0, 1)$ and differentiate the previous expression with respect to this direction:

$$\begin{aligned} D^2\tilde{W}_\alpha(u)(\varphi, \eta) &= \frac{d}{dt} D\tilde{W}_\alpha(u + t\eta)(\varphi)|_{t=0} \\ &= 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx - 10 \int_0^1 \frac{u'(x)u''(x)\varphi''(x)\eta'(x)}{(1+u'(x)^2)^{7/2}} dx \\ &\quad - 10 \int_0^1 \frac{u'(x)u''(x)\varphi'(x)\eta''(x)}{(1+u'(x)^2)^{7/2}} dx - 5 \int_0^1 \frac{u''(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{7/2}} dx \\ &\quad + 35 \int_0^1 \frac{u'(x)^2u''(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{9/2}} dx - 4\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1 \\ &= 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx - 5 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{\sqrt{1+u'(x)^2}} dx \\ &\quad - 10 \int_0^1 \kappa(x) \cdot \frac{u'(x)}{\sqrt{1+u'(x)^2}} \cdot \frac{1}{(1+u'(x)^2)^{3/2}} \cdot \frac{d}{dx} (\varphi'(x)\eta'(x)) dx \\ &\quad + 35 \int_0^1 \frac{u'(x)^2\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx - 4\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1. \end{aligned}$$

To proceed further we would like to integrate the third term by parts. Here we will exploit that u is a solution to (1.1). In particular, u is smooth and satisfies the Navier

boundary data $\kappa(x) = -\alpha, x \in \{0, 1\}$.

$$\begin{aligned}
 D^2\tilde{W}_\alpha(u)(\varphi, \eta) &= 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx - 5 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{\sqrt{1+u'(x)^2}} dx \\
 &\quad + 35 \int_0^1 \frac{u'(x)^2\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx - 4\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1 \\
 &\quad - 10 \left[\kappa(x) \frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1 + 10 \int_0^1 \frac{\kappa'(x)u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} dx \\
 &\quad + 10 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx - 30 \int_0^1 \frac{\kappa(x)u'(x)^2u''(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^3} dx \\
 &= 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx - 5 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{\sqrt{1+u'(x)^2}} dx \\
 &\quad + 5 \int_0^1 \frac{u'(x)^2\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx + 6\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1 \\
 &\quad + 10 \int_0^1 \frac{\kappa'(x)u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} dx + 10 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx.
 \end{aligned}$$

We infer from (2.3) and (2.4) that

$$\forall x \in [0, 1], \quad \kappa(x) (1 + u'(x)^2)^{1/4} = -c,$$

and hence

$$\forall x \in [0, 1], \quad \kappa'(x) (1 + u'(x)^2)^{1/4} + \frac{1}{2}u'(x)\kappa(x)^2 (1 + u'(x)^2)^{3/4} = 0.$$

Consequently,

$$\begin{aligned}
 D^2\tilde{W}_\alpha(u)(\varphi, \eta) &= 2 \int_0^1 \frac{\varphi''(x)\eta''(x)}{(1+u'(x)^2)^{5/2}} dx - 5 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{\sqrt{1+u'(x)^2}} dx \\
 &\quad + 6\alpha \left[\frac{u'(x)\varphi'(x)\eta'(x)}{(1+u'(x)^2)^2} \right]_0^1 + 10 \int_0^1 \frac{\kappa(x)^2\varphi'(x)\eta'(x)}{(1+u'(x)^2)^{3/2}} dx.
 \end{aligned}$$

This proves our claim. \square

Looking at $\eta \in H^2(0, 1) \cap H_0^1(0, 1)$ as a test function, we now obtain the linearization from the second variation with the help of integration by parts. Note that the first integral gives rise to a further boundary term $2 \left[\frac{\varphi''(x)\eta'(x)}{(1+u'(x)^2)^{5/2}} \right]_0^1$. Expressing u and α in terms of c according to Lemmas 2.1 and 2.2, the linearization of (1.1) around u reads as follows:

$$(2.7) \quad \left\{ \begin{aligned} &\left(\left(\frac{\varphi''(x)}{(1+G^{-1}(\frac{c}{2}-cx)^2)^{5/2}} \right)'' + \frac{5}{2}c^2 \left(\frac{G^{-1}(\frac{c}{2}-cx)^2-1}{(1+G^{-1}(\frac{c}{2}-cx)^2)^2} \varphi'(x) \right)' \right) = 0, \quad x \in (0, 1), \\ &\varphi(0) = \varphi(1) = 0, \\ &\frac{\varphi''(0)}{(1+G^{-1}(\frac{c}{2})^2)^{5/2}} + 3 \frac{cG^{-1}(\frac{c}{2})\varphi'(0)}{(1+G^{-1}(\frac{c}{2})^2)^{9/4}} = 0, \quad \frac{\varphi''(1)}{(1+G^{-1}(\frac{c}{2})^2)^{5/2}} - 3 \frac{cG^{-1}(\frac{c}{2})\varphi'(1)}{(1+G^{-1}(\frac{c}{2})^2)^{9/4}} = 0. \end{aligned} \right.$$

For $c = 0$, the small solution of (1.1) is $u(x) \equiv 0$, and $D^2\tilde{W}_0(u)(\varphi, \varphi) = \int_0^1 \varphi''(x)^2 dx$ is positive definite in $H^2 \cap H_0^1(0, 1)$ with respect to the $L^2(0, 1)$ -norm. The spectrum of the linearization as a regular elliptic operator on a bounded interval together with suitable boundary conditions consists only of eigenvalues. Since these eigenvalues depend smoothly on u and u depends smoothly on c , $D^2\tilde{W}_\alpha(u)(\varphi, \varphi)$ remains positive definite for c increasing from 0 as long as (2.7) has only the trivial solution $\varphi(x) \equiv 0$.

We assume that (2.7) has a solution φ and put

$$\chi(x) := \varphi'(x).$$

Then, there exists a constant $A \in \mathbb{R}$ such that χ solves the second order differential equation

$$\left(\frac{\chi'(x)}{\left(1 + G^{-1} \left(\frac{c}{2} - cx\right)^2\right)^{5/2}} \right)' + \frac{5}{2}c^2 \left(\frac{G^{-1} \left(\frac{c}{2} - cx\right)^2 - 1}{\left(1 + G^{-1} \left(\frac{c}{2} - cx\right)^2\right)^2} \chi \right) = c^2 A.$$

We introduce more suitable variables,

$$\begin{aligned} y &= G^{-1} \left(\frac{c}{2} - cx\right) \in \left[-G^{-1} \left(\frac{c}{2}\right), G^{-1} \left(\frac{c}{2}\right)\right], & x &= \frac{1}{2} - \frac{G(y)}{c}, \\ \psi(y) &:= \chi(x) = \chi \left(\frac{1}{2} - \frac{G(y)}{c}\right), & \chi(x) &= \psi \left(G^{-1} \left(\frac{c}{2} - cx\right)\right), \\ \chi'(x) &= -c \left(1 + G^{-1} \left(\frac{c}{2} - cx\right)^2\right)^{5/4} \psi' \left(G^{-1} \left(\frac{c}{2} - cx\right)\right), \end{aligned}$$

and conclude that ψ solves the following boundary value problem:

$$(2.8) \quad \begin{cases} \psi''(y) - \frac{5y}{2(1+y^2)}\psi'(y) + \frac{5(y^2-1)}{2(1+y^2)^2}\psi(y) = A, & y \in (-y_0, y_0), \\ \psi'(-y_0) + \frac{3y_0}{1+y_0^2}\psi(-y_0) = 0, & \psi'(y_0) - \frac{3y_0}{1+y_0^2}\psi(y_0) = 0. \end{cases}$$

Here, we denote

$$(2.9) \quad y_0 := G^{-1} \left(\frac{c}{2}\right).$$

To simplify the boundary conditions we make a last change of variables and put

$$(2.10) \quad \Phi(y) := \frac{\psi(y)}{(1+y^2)^{3/2}}, \quad y \in [-y_0, y_0],$$

and finally consider the following boundary value problem:

$$(2.11) \quad \begin{cases} (1+y^2)^{3/2}\Phi''(y) + \frac{7}{2}y(1+y^2)^{1/2}\Phi'(y) \\ \quad + (y^2 + \frac{1}{2})(1+y^2)^{-1/2}\Phi(y) = A, & y \in (-y_0, y_0), \\ \Phi'(-y_0) = \Phi'(y_0) = 0. \end{cases}$$

We recall the definition of $G(y) := \int_0^y \frac{1}{(1+\tau^2)^{5/4}} d\tau$ and put

$$(2.12) \quad \Phi_0(y) := -2\frac{1}{\sqrt{1+y^2}}, \quad \Phi_1(y) := \frac{1}{\sqrt[4]{1+y^2}}, \quad \Phi_2(y) := \frac{G(y)}{\sqrt[4]{1+y^2}}.$$

Then, one directly verifies that the general solution of the differential equation in (2.11) is given by

$$(2.13) \quad \Phi(y) := A \cdot \Phi_0(y) + \gamma_1 \cdot \Phi_1(y) + \gamma_2 \cdot \Phi_2(y)$$

with $\gamma_1, \gamma_2 \in \mathbb{R}$. Since $A \cdot \Phi_0(y) + \gamma_1 \cdot \Phi_1(y)$ is even and $\gamma_2 \cdot \Phi_2(y)$ is odd, the boundary conditions in (2.11) are equivalent to

$$(2.14) \quad A \cdot \Phi'_0(y_0) + \gamma_1 \cdot \Phi'_1(y_0) = 0 \text{ and } \gamma_2 \cdot \Phi'_2(y_0) = 0,$$

which in turn are equivalent to

$$(2.15) \quad \gamma_1 = \frac{4A}{\sqrt[4]{1+y_0^2}} \text{ and } (\gamma_2 = 0 \text{ or } \Phi'_2(y_0) = 0).$$

A beautiful coincidence between these solutions and the functions involved in the proof of Theorem 1.1 can be observed, namely,

$$(2.16) \quad \Phi_2(y) = \frac{1}{2}h(2G(y)), \quad \Phi'_2(y) = \frac{h'(2G(y))}{(1+y^2)^{5/4}}.$$

With the help of these observations we are now ready to conclude the following lemma.

LEMMA 2.4. *For $c \in [0, c_0) \setminus \{c_{\max}\}$, the boundary value problem (2.7) has only the trivial solution $\varphi(x) \equiv 0$. For $c = c_{\max}$, it has a one-dimensional null space which is spanned by*

$$\varphi(x) = \frac{1}{c} \int_{G^{-1}(\frac{c}{2}-cx)}^{G^{-1}(\frac{c}{2})} G(\eta) d\eta.$$

If $c = c_{\max}$, $\alpha = \alpha_{\max}$, instabilities will occur first from the corresponding solution u in direction of this function φ ; see Figure 2.2.

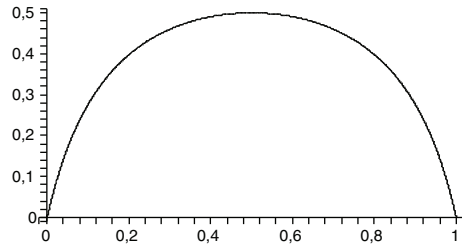


FIG. 2.2. Profile of the unstable direction in $c = c_{\max}$.

Proof. The case $c = 0$ is obvious, and we consider only $c \in (0, c_0)$. We denote

$$\tilde{\Phi}_1(y) := \Phi_0(y) + \frac{4}{\sqrt[4]{1+y_0^2}} \cdot \Phi_1(y) = \frac{2}{\sqrt[4]{1+y^2} \cdot \sqrt[4]{1+y_0^2}} \left(2 - \frac{\sqrt[4]{1+y_0^2}}{\sqrt[4]{1+y^2}} \right).$$

According to (2.13), we have to study

$$\Phi(y) = A\tilde{\Phi}_1(y) + \gamma_2\Phi_2(y)$$

with some suitable $A, \gamma_2 \in \mathbb{R}$. Let φ be the corresponding solution of (2.7) which is obtained from Φ by tracing back the changes of variables and integrating χ . We want to show first that necessarily $A = 0$ for any $c \in [0, c_0)$:

$$\begin{aligned} 0 &= \varphi(1) - \varphi(0) = \int_0^1 \chi(x) dx = \int_0^1 \psi \left(G^{-1} \left(\frac{c}{2} - cx \right) \right) dx \\ &= \frac{1}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \psi(y)(1+y^2)^{-5/4} dy = \frac{1}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \Phi(y)(1+y^2)^{1/4} dy \\ &= \frac{A}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \tilde{\Phi}_1(y)(1+y^2)^{1/4} dy + \frac{\gamma_2}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \Phi_2(y)(1+y^2)^{1/4} dy \\ &= \frac{A}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \tilde{\Phi}_1(y)(1+y^2)^{1/4} dy \end{aligned}$$

since Φ_2 is odd. Hence we may conclude that

$$\begin{aligned} 0 &= \frac{A}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \tilde{\Phi}_1(y)(1+y^2)^{1/4} dy \\ &= \frac{2A}{c} \int_{-G^{-1}(\frac{c}{2})}^{G^{-1}(\frac{c}{2})} \left(\frac{2}{(1+G^{-1}(\frac{c}{2}))^{1/4}} - \frac{1}{(1+y^2)^{1/4}} \right) dy \\ &= \frac{4A}{c} F \left(G^{-1} \left(\frac{c}{2} \right) \right), \end{aligned}$$

where F is defined by

$$F(\eta) := \frac{2\eta}{(1+\eta^2)^{1/4}} - \int_0^\eta \frac{1}{(1+s^2)^{1/4}} ds.$$

Since $F(0) = 0$ and

$$\begin{aligned} F'(\eta) &= \frac{2}{(1+\eta^2)^{1/4}} - \frac{\eta^2}{(1+\eta^2)^{5/4}} - \frac{1}{(1+\eta^2)^{1/4}} \\ &= \frac{1}{(1+\eta^2)^{1/4}} - \frac{\eta^2}{(1+\eta^2)^{5/4}} = \frac{1}{(1+\eta^2)^{5/4}} > 0, \end{aligned}$$

we have

$$F \left(G^{-1} \left(\frac{c}{2} \right) \right) > 0.$$

As a consequence, $A = 0$, and hence $\gamma_1 = 0$ by (2.15), and we are left with considering $\gamma_2\Phi_2$. We have that $h'(c) > 0$ for $c \in (0, c_{\max})$ and $h'(c) < 0$ for $c \in (c_{\max}, c_0)$. By making use of

$$\Phi'_2(y) = \frac{h'(2G(y))}{(1+y^2)^{5/4}}$$

and the boundary condition $\gamma_2\Phi'_2(G^{-1}(c/2)) = 0$, we conclude that $\gamma_2 = 0$ provided $c \in (0, c_0) \setminus \{c_{\max}\}$. If $c = c_{\max}$, then Φ_2 is a nontrivial solution of (2.11). For the

corresponding nontrivial solution φ of (2.7) we derive

$$\begin{aligned} \varphi(x) &= \gamma_2 \frac{c_{\max}}{2} \int_0^x \left(1 + G^{-1} \left(\frac{c_{\max}}{2} - c_{\max} \xi \right)^2 \right)^{5/4} (1 - 2\xi) d\xi \\ &= \frac{\gamma_2}{c_{\max}} \int_{G^{-1}(\frac{c_{\max}}{2} - c_{\max}x)}^{G^{-1}(\frac{c_{\max}}{2})} G(\eta) d\eta, \end{aligned}$$

where we made use of the boundary conditions $\varphi(0) = \varphi(1) = 0$. \square

Proof of Theorem 1.2. The proof is now immediate. By the preceding lemma we have that on $[0, c_{\max})$, 0 is not an eigenvalue of (2.7). Since $D^2\tilde{W}_0(u)(\varphi, \varphi)$ is positive definite in $H^2 \cap H_0^1(0, 1)$ with respect to the $L^2(0, 1)$ -norm, by continuity, the same holds true for $D^2\tilde{W}_\alpha(u)(\varphi, \varphi)$ for $c \in [0, c_{\max})$, which is the stated linearized stability of the corresponding small solutions of (1.1). \square

As an immediate consequence of Theorem 1.2 we obtain a global existence result for the geometric flow associated with (1.1), namely,

$$V = -\kappa_{ss} - \frac{1}{2}\kappa^3 \quad \text{on } \Gamma(t).$$

Here, V denotes the upward normal velocity of the evolving graphs

$$\Gamma(t) = \{(x, v(x, t)) \mid x \in [0, 1]\}.$$

The above evolution law then leads to the parabolic initial-boundary value problem (2.18) below. The principle of linearized stability as it was proved in great generality by Latushkin, Prüss, and Schnaubelt [11, Proposition 16] can be applied to our situation and allows us to obtain global existence and asymptotic stability for initial data close to a small solution to (1.1).

COROLLARY 2.5. *Assume that $c \in [0, c_{\max})$, and let $\alpha = h(c) = \frac{c}{\sqrt[4]{1+G^{-1}(\frac{c}{2})^2}}$ and*

$$(2.17) \quad u(x) = \frac{2}{c^4 \sqrt[4]{1 + G^{-1}(\frac{c}{2} - cx)^2}} - \frac{2}{c^4 \sqrt[4]{1 + G^{-1}(\frac{c}{2})^2}}$$

be the corresponding small solution of (1.1). We fix some $p > 5$. Then, there exist $\delta, \rho, C > 0$ such that for $v_0 \in W^{4,p}(0, 1)$ with $v_0(0) = v_0(1) = 0$, $\kappa_{v_0}(0) = \kappa_{v_0}(1) = -\alpha$, and

$$\|v_0 - u\|_{W^{4,p}(0,1)} \leq \delta,$$

there exists a global solution $v \in L^p(0, \infty, W^{4,p}(0, 1)) \cap W^{1,p}(0, \infty, L^p(0, 1))$ of the initial Navier boundary value problem

$$(2.18) \quad \begin{cases} \frac{v_t(t,x)}{\sqrt{1+v_x(t,x)^2}} + \frac{1}{\sqrt{1+v_x(t,x)^2}} \frac{d}{dx} \left(\frac{\kappa_{v,x}(t,x)}{\sqrt{1+v_x(t,x)^2}} \right) + \frac{1}{2}\kappa_v^3(t,x) = 0, & (t,x) \in [0, \infty) \times [0, 1], \\ v(t,0) = v(t,1) = 0, \quad \kappa_v(t,0) = \kappa_v(t,1) = -\alpha, & t \in [0, \infty), \\ v(0,x) = v_0(x), & x \in [0, 1]. \end{cases}$$

One has exponential convergence toward the steady state u :

$$(2.19) \quad \|v(t, \cdot) - u(\cdot)\|_{W^{4,p}(0,1)} \leq C \exp(-\rho t) \quad (t \geq 1).$$

With similar but simpler techniques and calculations one finds that the unique solution (cf. [4, Theorem 2]), being symmetric about $x = 1/2$ of the Dirichlet problem

$$(2.20) \quad \begin{cases} \frac{1}{\sqrt{1+u_x(x)^2}} \frac{d}{dx} \left(\frac{\kappa_x(x)}{\sqrt{1+u_x(x)^2}} \right) + \frac{1}{2} \kappa^3(t, x) = 0, & x \in [0, 1], \\ u(0) = u(1) = 0, \quad u_x(0) = -u_x(1) = \beta, \end{cases}$$

$\beta \in \mathbb{R}$, is (linearized) stable. Analogously, a global existence result follows for the initial Dirichlet boundary value problem

$$(2.21) \quad \begin{cases} \frac{v_t(t, x)}{\sqrt{1+v_x(t, x)^2}} + \frac{1}{\sqrt{1+v_x(t, x)^2}} \frac{d}{dx} \left(\frac{\kappa_{v, x}(t, x)}{\sqrt{1+v_x(t, x)^2}} \right) + \frac{1}{2} \kappa_v^3(t, x) = 0, & (t, x) \in [0, \infty) \times [0, 1], \\ v(t, 0) = v(t, 1) = 0, \quad v_x(t, 0) = -v_x(t, 1) = \beta, & t \in [0, \infty), \\ v(0, x) = v_0(x), & x \in [0, 1], \end{cases}$$

provided the initial datum v_0 obeys the same boundary data and is sufficiently close to the stationary solution u of (2.20) with respect to the $W^{4,p}$ -norm ($p > 5$).

3. Morse index of the large solution. For $c \in (0, c_0)$ we consider as in (2.5)

$$u_c(x) = \frac{2}{c^4 \sqrt{1 + G^{-1} \left(\frac{c}{2} - cx \right)^2}} - \frac{2}{c^4 \sqrt{1 + G^{-1} \left(\frac{c}{2} \right)^2}}.$$

In order to prove Theorem 1.3, we have to show that exactly one eigenvalue of the quadratic form

$$\varphi \mapsto D^2 \tilde{W}_\alpha(u_c)(\varphi, \varphi), \quad \alpha = h(c),$$

passes through 0 when c passes through c_{\max} and that for $c \in (c_{\max}, c_0)$, 0 is not an eigenvalue of $D^2 \tilde{W}_\alpha(u_c)$, i.e., of (2.7). The latter was already done in Lemma 2.4. Moreover, its proof yields that there is at most one eigenvalue, which crosses 0 in $c = c_{\max}$. It remains to show that for $c > c_{\max}$ and suitable $\varphi \in H^2 \cap H_0^1(0, 1)$, one indeed has $D^2 \tilde{W}_\alpha(u_c)(\varphi, \varphi) < 0$. Making use of the same transformations and notations of section 2 and restricting ourselves to symmetric φ , we find

$$\begin{aligned} D^2 \tilde{W}_\alpha(u_c)(\varphi, \varphi) &= 2 \int_0^1 \frac{\chi'(x)^2}{\left(1 + G^{-1} \left(\frac{c}{2} - cx\right)^2\right)^{5/2}} dx \\ &\quad - 5c^2 \int_0^1 \frac{G^{-1} \left(\frac{c}{2} - cx\right)^2 - 1}{\left(1 + G^{-1} \left(\frac{c}{2} - cx\right)^2\right)^2} \chi(x)^2 dx - 12h(c) \frac{G^{-1} \left(\frac{c}{2}\right)}{\left(1 + G^{-1} \left(\frac{c}{2}\right)^2\right)^2} \chi(1)^2 \\ &= 2c \int_{-G^{-1} \left(\frac{c}{2}\right)}^{G^{-1} \left(\frac{c}{2}\right)} \frac{\psi'(y)^2}{(1 + y^2)^{5/4}} dy - 5c \int_{-G^{-1} \left(\frac{c}{2}\right)}^{G^{-1} \left(\frac{c}{2}\right)} \frac{(y^2 - 1)}{(1 + y^2)^{13/4}} \psi(y)^2 dy \\ &\quad - 12 \frac{c G^{-1} \left(\frac{c}{2}\right)}{\left(1 + G^{-1} \left(\frac{c}{2}\right)^2\right)^{9/4}} \psi \left(G^{-1} \left(\frac{c}{2}\right)\right)^2. \end{aligned}$$

We choose

$$\psi_c(y) := (1 + y^2)^{3/2} \Phi_2(y) = (1 + y^2)^{5/4} G(y)$$

and obtain for the corresponding $\varphi_c \in H^2 \cap H_0^1(0, 1)$

$$\begin{aligned}
 & \frac{1}{4c} D^2 \tilde{W}_\alpha(u_c)(\varphi_c, \varphi_c) \\
 &= \int_0^{G^{-1}(\frac{c}{2})} \left(\left((1+y^2)^{-1/4} + \frac{5}{2}yG(y) \right)^2 - \frac{5}{2}(y^2-1)G(y)^2 \right) \frac{dy}{(1+y^2)^{3/4}} \\
 (3.1) \quad & - \frac{3}{4}c^2 G^{-1} \left(\frac{c}{2} \right) \left(1 + G^{-1} \left(\frac{c}{2} \right)^2 \right)^{1/4}.
 \end{aligned}$$

According to Theorem 1.2, we know that this expression is equal to 0 for $c = c_{\max}$. Writing $c = 2G(d)$, we see that the asymptotic behavior of the right-hand side is dominated by

$$\frac{c_0^2}{4} \left(\frac{25}{4} \cdot \frac{2}{3} - \frac{5}{2} \cdot \frac{2}{3} - 3 \right) d^{3/2} = -\frac{c_0^2}{8} d^{3/2} \rightarrow -\infty$$

for $d \rightarrow \infty$, i.e., $c \nearrow c_0$. This shows, together with Lemma 2.4, that $D^2 \tilde{W}_\alpha(u_c)(\varphi_c, \varphi_c) < 0$ for $c \in (c_{\max}, c_0)$ and concludes the proof of Theorem 1.3.

The right-hand side of (3.1) is plotted in Figure 3.1. Since $\varphi_c \rightarrow 0$ for $c \searrow 0$, the curve starts in $(0, 0)$, although there, $D^2 \tilde{W}_\alpha(u_0)$ is positive definite.

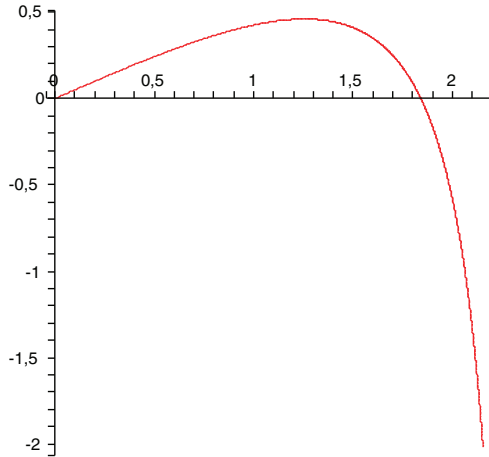


FIG. 3.1. $c \mapsto \frac{1}{4c} D^2 \tilde{W}_\alpha(u_c)(\varphi_c, \varphi_c)$.

4. Global minima and symmetry. The aim of this section is to examine whether the small solutions which were found to be local minima in section 2 are also global minima for the functional \tilde{W}_α . In what follows it will be convenient to write

$$\begin{aligned}
 \tilde{W}_\alpha(v) &= \int_0^1 (\kappa(x)^2 + 2\alpha\kappa(x)) \sqrt{1 + v'(x)^2} dx, \\
 &= \int_0^1 \kappa(x)^2 \sqrt{1 + v'(x)^2} dx + 2\alpha [\arctan(v'(x))]_0^1 =: W(v) + BC_\alpha(v).
 \end{aligned}$$

We remark that all quantities are geometric and so are invariant under rotation. Moreover, when stretching a curve by a factor k , W is multiplied by a factor $1/k$ while BC_α remains unchanged.

We shall see that the task of minimizing \tilde{W}_α can be reduced to a minimization problem for a function of two variables. As a byproduct of the analysis of this function we find that in order to determine $\inf_{v \in H^2 \cap H_0^1} \tilde{W}_\alpha(v)$ it is sufficient to minimize over all symmetric functions. The reduction to a two-dimensional problem is achieved in two steps. We begin by showing that it is enough to consider concave functions.

LEMMA 4.1. *Suppose that $u \in H^2 \cap H_0^1(0, 1)$ is not concave. Then there exists a concave function $v \in H^2 \cap H_0^1(0, 1)$ with $\tilde{W}_\alpha(v) < \tilde{W}_\alpha(u)$.*

Proof. It is natural to think of v as the concave envelope of u , so that we are led to consider the following obstacle problem: find $v \in K$ such that

$$(4.1) \quad \forall \eta \in K, \quad \int_0^1 v'(\eta' - v') \geq 0,$$

where $K = \{\eta \in H_0^1(0, 1) \mid \eta \geq u \text{ a.e. in } (0, 1)\}$. It is shown in Chapter IV of [7] that v can be obtained as the limit of a sequence $(v_\varepsilon)_{\varepsilon > 0}$, where $v_\varepsilon \in H^2 \cap H_0^1(0, 1)$ solves

$$(4.2) \quad -v_\varepsilon'' = (-u'')^+ \vartheta_\varepsilon(v_\varepsilon - u) \quad \text{in } (0, 1).$$

Here, $\vartheta_\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ satisfies

$$\vartheta_\varepsilon(t) = \begin{cases} 1, & t < 0, \\ 1 - \frac{t}{\varepsilon}, & 0 \leq t \leq \varepsilon, \\ 0, & t > \varepsilon. \end{cases}$$

It follows from the analysis in [7] that $v_\varepsilon \rightarrow v$ in $H^1(0, 1)$ and $v_\varepsilon'' \rightarrow v''$ in $L^2(0, 1)$ as $\varepsilon \rightarrow 0$, so that $v \in H^2 \cap H_0^1(0, 1)$ and $v'' \leq 0$ a.e. in $(0, 1)$; in particular, v is concave. Denoting by $I = \{x \in [0, 1] \mid v(x) = u(x)\}$ the coincidence set, we have that $v'' = 0$ a.e. in $[0, 1] \setminus I$. Furthermore, using (4.2)

$$\begin{aligned} W(v) &= \int_0^1 \frac{|v''|^2}{(1 + (v')^2)^{\frac{5}{2}}} = \int_I \frac{|v''|^2}{(1 + (u')^2)^{\frac{5}{2}}} \leq \liminf_{\varepsilon \rightarrow 0} \int_I \frac{|v_\varepsilon''|^2}{(1 + (u')^2)^{\frac{5}{2}}} \\ &\leq \int_I \frac{|(-u'')^+|^2}{(1 + (u')^2)^{\frac{5}{2}}} \leq \int_0^1 \frac{|(-u'')^+|^2}{(1 + (u')^2)^{\frac{5}{2}}} \leq \int_0^1 \frac{|u''|^2}{(1 + (u')^2)^{\frac{5}{2}}} = W(u). \end{aligned}$$

If we had $W(v) = W(u)$, then the above argument would imply that $(-u'')^- = 0$ a.e. in $(0, 1)$ and therefore $u'' \leq 0$ a.e. in $(0, 1)$, contradicting our assumption that u is not concave. Hence $W(v) < W(u)$; since $v \geq u$ we have that $u'(0) \leq v'(0)$ and $u'(1) \geq v'(1)$, and therefore $\tilde{W}_\alpha(v) < \tilde{W}_\alpha(u)$. \square

In what follows we shall make use of the prototype solution

$$(4.3) \quad U_0(x) = \frac{2}{c_0 \sqrt[4]{1 + G^{-1} \left(\frac{c_0}{2} - c_0 x\right)^2}}.$$

Formally, it is the large solution of the Navier boundary value problem (1.1) for $\alpha = 0$. However, one should observe that this solution is no longer smooth as a graph near $x = 0$ and $x = 1$, and for this reason, it was not included in Proposition 1.1.

Suppose that $0 \leq x_0 < x_1 \leq 1$ are two points with $x_1 - x_0 < 1$. Then $U_0|_{[x_0, x_1]}$ can be written as a graph over the segment connecting $(x_0, U_0(x_0))$ and $(x_1, U_0(x_1))$. We denote by $u_{x_0, x_1} : [0, 1] \rightarrow \mathbb{R}$ the strictly concave function which is obtained by translating, rotating, and rescaling the above graph to the unit interval $[0, 1]$. Note

that $u_{x_0, x_1} \in H^2 \cap H_0^1(0, 1)$. Our next lemma essentially reduces the minimization of \tilde{W}_α to a two-dimensional minimization problem.

LEMMA 4.2. *Suppose that $u \in H^2 \cap H_0^1(0, 1) \setminus \{0\}$ is concave. Then there exist $0 \leq x_0 < x_1 \leq 1$, $x_1 - x_0 < 1$ such that $v = u_{x_0, x_1}$ satisfies $BC_\alpha(u) = BC_\alpha(v)$ and either $W(v) \leq W(u)$, $u'(0) = v'(0)$, or $W(v) < W(u)$, $u'(0) \neq v'(0)$.*

Proof. Let us denote by β_ℓ and β_r the boundary angles of $\text{graph}(u)$ on the left and on the right, respectively. Since u is assumed to be concave and nontrivial we have $\beta_\ell, \beta_r \in (0, \frac{\pi}{2})$. Consider

$$\mathcal{K} := \text{graph}(U_0) \cup \{(0, y) : y \leq 0\} \cup \{(1, y) : y \leq 0\}.$$

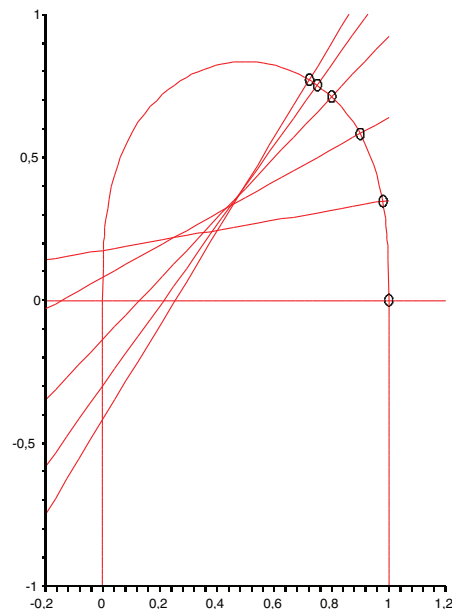


FIG. 4.1. Left angle β_ℓ ; right angle $\pi/2$.

This is neither a graph nor a solution of the Willmore equation. However, it is a regular H^2 -curve, locally an H^2 -graph over the x - or the y -axis, respectively, and it has minimal Willmore energy c_0^2 among all concave curves connecting any point from $\{(0, y) : y \leq 0\}$ with any point from $\{(1, y) : y \leq 0\}$ with tangential directions $(0, 1)$ and $(0, -1)$, respectively. This minimality follows similarly as in [4, end of section 5].

Claim. There exist two points $P = (x_P, y_P)$, $Q = (x_Q, y_Q) \in \mathcal{K}$, $P \neq Q$, such that the segment $[P, Q]$ intersects \mathcal{K} under the angles β_ℓ at P and β_r at Q .

To see this, we start with the point $(x_1, y_1) = (1, 0)$ and the orthogonal straight line through this point. This line intersects the left part of \mathcal{K} in (x_0, y_0) under a right angle. Now we move the point (x_1, y_1) and the corresponding orthogonal straight line counterclockwise. The corresponding (x_0, y_0) finally moves down, and the intersection angle (at least finally) decreases and becomes arbitrarily small. In particular, the left angle β_ℓ is attained. See Figure 4.1. Now we keep this angle fixed and move the point (x_0, y_0) clockwise. We consider (x_1, y_1) on the right part of \mathcal{K} as an intersection point with the straight line building the angle β_ℓ with \mathcal{K} in (x_0, y_0) . At the beginning this right angle is $\pi/2$, while it becomes arbitrarily small when (x_0, y_0) moves clockwise. In particular, β_r is attained as the angle on the right and the claim is proved.

In view of the above-mentioned minimality property of \mathcal{K} , \mathcal{K}' enjoys a similar minimality among those arcs with boundary angles β_ℓ, β_r . We infer that

$$(4.4) \quad W(\mathcal{K}') \leq \frac{1}{|P - Q|} W(u),$$

where \mathcal{K}' denotes the subarc of \mathcal{K} between P and Q . Observing that by construction y_P and y_Q cannot both be negative, we may distinguish two cases.

Case 1. $y_P \geq 0$ and $y_Q \geq 0$. Setting $x_0 = x_P, x_1 = x_Q$, we have $x_1 - x_0 < 1$ since $\beta_\ell, \beta_r \in (0, \frac{\pi}{2})$. The function $v = u_{x_0, x_1}$ then satisfies

$$W(v) = |P - Q|W(\mathcal{K}') \leq W(u)$$

as well as $v'(0) = u'(0)$ and $v'(1) = u'(1)$.

Case 2. Either $y_P < 0$ or $y_Q < 0$. If $y_P < 0$, then $y_Q > 0$ since $\beta_r < \frac{\pi}{2}$, and we let $x_0 = 0, x_1 = x_Q$, and $v = u_{x_0, x_1}$. Denoting by $L(x_0, x_1)$ the length of the segment connecting $(x_0, U_0(x_0))$ and $(x_1, U_0(x_1))$, we have

$$W(v) = L(x_0, x_1)W(\mathcal{K}') \leq \frac{L(x_0, x_1)}{|P - Q|} W(u) < W(u)$$

since $u \neq 0$ and by construction any point on $\text{graph}(U_0)$ is strictly closer to $(0, 0)$ than to any other point on $\{(0, y) \mid y < 0\}$. A similar argument applies if $y_Q < 0$. Finally note that while $BC_\alpha(u) = BC_\alpha(v)$, we have $u'(0) \neq v'(0)$ in this case. \square

We deduce from Lemmas 4.1 and 4.2 that when determining $\inf \tilde{W}_\alpha$ over $H^2 \cap H_0^1(0, 1)$ it is sufficient to calculate the Willmore energy for functions $v = u_{x_0, x_1}$ with $0 \leq x_0 < x_1 \leq 1$ and $x_1 - x_0 < 1$. The integrand for W on $[x_0, x_1]$ is c_0^2 , so the integral is $c_0^2 \cdot (x_1 - x_0)$. The length of the base line is $((x_1 - x_0)^2 + (U_0(x_1) - U_0(x_0))^2)^{1/2}$. As for BC_α , we have $2\alpha(\arctan(U_0'(x_1)) - \arctan(U_0'(x_0)))$. After rotation and rescaling we come up with

$$\begin{aligned} \tilde{W}_\alpha(u_{x_0, x_1}) &= c_0^2 \cdot (x_1 - x_0) \left((x_1 - x_0)^2 + (U_0(x_1) - U_0(x_0))^2 \right)^{1/2} \\ &\quad + 2\alpha(\arctan(U_0'(x_1)) - \arctan(U_0'(x_0))) \\ &= c_0^2 \cdot (x_1 - x_0) \\ &\quad \cdot \left((x_1 - x_0)^2 + \frac{4}{c_0^2} \left(\frac{1}{\sqrt[4]{1 + G^{-1}(c_0/2 - c_0x_1)^2}} - \frac{1}{\sqrt[4]{1 + G^{-1}(c_0/2 - c_0x_0)^2}} \right)^2 \right)^{1/2} \\ &\quad + 2\alpha(\arctan(G^{-1}(c_0/2 - c_0x_1)) - \arctan(G^{-1}(c_0/2 - c_0x_0))). \end{aligned}$$

We now introduce the new variables

$$(4.5) \quad d_0 := G^{-1}(c_0/2 - c_0x_0), \quad d_1 := -G^{-1}(c_0/2 - c_0x_1), \quad d_1 > -d_0,$$

so that

$$x_0 = \frac{1}{2} - \frac{1}{c_0}G(d_0), \quad x_1 = \frac{1}{2} + \frac{1}{c_0}G(d_1).$$

Defining

$$\hat{W}_\alpha(d_0, d_1) := \tilde{W}_\alpha(u_{x_0, x_1})$$

(see Figure 4.2), we end up with

$$\begin{aligned} & \hat{W}_\alpha(d_0, d_1) \\ &= (G(d_0) + G(d_1)) \left((G(d_0) + G(d_1))^2 + 4((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4})^2 \right)^{1/2} \\ & \quad - 2\alpha(\arctan(d_0) + \arctan(d_1)). \end{aligned}$$

The following result summarizes what we have achieved so far.

THEOREM 4.3. *Let $\alpha \geq 0$. Then*

$$\inf_{v \in H^2 \cap H_0^1(0,1)} \tilde{W}_\alpha(v) = \inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1).$$

Let us remark that ideas similar to those employed to obtain Theorem 4.3 were used in [2] in order to prove an existence result for axially symmetric Willmore surfaces satisfying Dirichlet boundary conditions. The corresponding Navier problem, however, is still open.

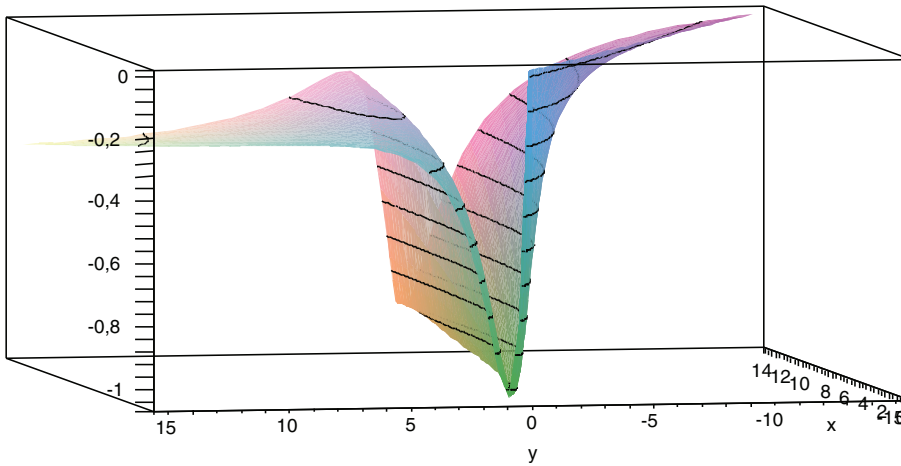


FIG. 4.2. Cross section of the graph of \hat{W}_1 along the axis $d_0 = d_1$.

It remains to discuss the two-dimensional function $\hat{W}_\alpha(d_0, d_1)$, ($d_1 \geq -d_0$). Here, the key step is proving positivity for the following expression.

LEMMA 4.4. *For $d_1 > -d_0$ we have that*

$$\begin{aligned} & (G(d_0) + G(d_1)) \cdot \left(\frac{d_0}{\sqrt[4]{1 + d_0^2}} + \frac{d_1}{\sqrt[4]{1 + d_1^2}} \right) - (G(d_0) + G(d_1))^2 \\ & \quad - 2 \left(\frac{1}{\sqrt[4]{1 + d_0^2}} - \frac{1}{\sqrt[4]{1 + d_1^2}} \right)^2 > 0. \end{aligned}$$

Proof. By the fundamental theorem of calculus and since G is odd, we have

$$\begin{aligned} G(d_0) + G(d_1) &= G(d_0) - G(-d_1) = \int_{-d_1}^{d_0} \frac{1}{(1 + \tau^2)^{5/4}} d\tau, \\ \frac{d_0}{\sqrt[4]{1 + d_0^2}} + \frac{d_1}{\sqrt[4]{1 + d_1^2}} &= \left[\frac{\tau}{(1 + \tau^2)^{1/4}} \right]_{-d_1}^{d_0} = \int_{-d_1}^{d_0} \frac{1 + \frac{1}{2}\tau^2}{(1 + \tau^2)^{5/4}} d\tau, \end{aligned}$$

$$\frac{1}{\sqrt[4]{1+d_0^2}} - \frac{1}{\sqrt[4]{1+d_1^2}} = \left[\frac{1}{(1+\tau^2)^{1/4}} \right]_{-d_1}^{d_0} = -\frac{1}{2} \int_{-d_1}^{d_0} \frac{\tau}{(1+\tau^2)^{5/4}} d\tau.$$

One may observe that $d_1 > -d_0$ is equivalent to $-d_1 < d_0$. The first two terms in the expression under consideration combine as follows:

$$\begin{aligned} & (G(d_0) + G(d_1)) \cdot \left(\frac{d_0}{\sqrt[4]{1+d_0^2}} + \frac{d_1}{\sqrt[4]{1+d_1^2}} \right) - (G(d_0) + G(d_1))^2 \\ &= \frac{1}{2} \left(\int_{-d_1}^{d_0} \frac{1}{(1+\tau^2)^{5/4}} d\tau \right) \cdot \left(\int_{-d_1}^{d_0} \frac{\tau^2}{(1+\tau^2)^{5/4}} d\tau \right). \end{aligned}$$

We now apply the Cauchy–Schwarz inequality and make use of $\tau \mapsto \frac{1}{(1+\tau^2)^{5/8}}$ and $\tau \mapsto \frac{\tau}{(1+\tau^2)^{5/8}}$ being linearly independent to obtain

$$\begin{aligned} & 2 \left(\frac{1}{\sqrt[4]{1+d_0^2}} - \frac{1}{\sqrt[4]{1+d_1^2}} \right)^2 = \frac{1}{2} \left(\int_{-d_1}^{d_0} \frac{\tau}{(1+\tau^2)^{5/4}} d\tau \right)^2 \\ & < \frac{1}{2} \left(\int_{-d_1}^{d_0} \frac{1}{(1+\tau^2)^{5/4}} d\tau \right) \cdot \left(\int_{-d_1}^{d_0} \frac{\tau^2}{(1+\tau^2)^{5/4}} d\tau \right) \\ &= (G(d_0) + G(d_1)) \cdot \left(\frac{d_0}{\sqrt[4]{1+d_0^2}} + \frac{d_1}{\sqrt[4]{1+d_1^2}} \right) - (G(d_0) + G(d_1))^2, \end{aligned}$$

thereby proving the claim. \square

Next we show that in the open interior of the domain of definition of the two-dimensional energy function \hat{W}_α , critical points may occur at most on the diagonal, i.e., on symmetric graphs in the original context.

LEMMA 4.5. *Let $\alpha \geq 0$ and assume that*

$$\begin{aligned} & (d_0, d_1) \mapsto \hat{W}_\alpha(d_0, d_1) \\ &= (G(d_0) + G(d_1)) \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1+d_0^2)^{-1/4} - (1+d_1^2)^{-1/4} \right)^2 \right)^{1/2} \\ & \quad - 2\alpha (\arctan(d_0) + \arctan(d_1)) \end{aligned}$$

has a critical point (d_0, d_1) with $d_1 > -d_0$. Then

$$d_0 = d_1.$$

Proof. In a critical point of \hat{W}_α , we have that

$$\begin{aligned} 0 &= \frac{\partial}{\partial d_0} \hat{W}_\alpha(d_0, d_1) \\ &= \frac{1}{2} (G(d_0) + G(d_1)) \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1+d_0^2)^{-1/4} - (1+d_1^2)^{-1/4} \right)^2 \right)^{-1/2} \\ & \quad \cdot \left(2(G(d_0) + G(d_1))(1+d_0^2)^{-5/4} - 4d_0(1+d_0^2)^{-5/4} \left((1+d_0^2)^{-1/4} - (1+d_1^2)^{-1/4} \right) \right) \\ & \quad + (1+d_0^2)^{-5/4} \left((G(d_0) + G(d_1))^2 + 4 \left((1+d_0^2)^{-1/4} - (1+d_1^2)^{-1/4} \right)^2 \right)^{1/2} - \frac{2\alpha}{1+d_0^2}; \end{aligned}$$

$$\begin{aligned}
 0 &= \frac{\partial}{\partial d_1} \hat{W}_\alpha(d_0, d_1) \\
 &= \frac{1}{2} (G(d_0) + G(d_1)) \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{-1/2} \\
 &\quad \cdot \left(2(G(d_0) + G(d_1))(1 + d_1^2)^{-5/4} + 4d_1(1 + d_1^2)^{-5/4} \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right) \right) \\
 &\quad + (1 + d_1^2)^{-5/4} \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{1/2} \\
 &\quad - 2\alpha \frac{1}{1 + d_1^2}.
 \end{aligned}
 \tag{4.6}$$

Equivalently,

$$\begin{aligned}
 0 &= (G(d_0) + G(d_1)) \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{-1/2} \\
 &\quad \cdot \left((G(d_0) + G(d_1))(1 + d_0^2)^{-1/4} - 2d_0(1 + d_0^2)^{-1/4} \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right) \right) \\
 &\quad + (1 + d_0^2)^{-1/4} \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{1/2} - 2\alpha; \\
 0 &= (G(d_0) + G(d_1)) \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{-1/2} \\
 &\quad \cdot \left((G(d_0) + G(d_1))(1 + d_1^2)^{-1/4} + 2d_1(1 + d_1^2)^{-1/4} \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right) \right) \\
 &\quad + (1 + d_1^2)^{-1/4} \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{1/2} - 2\alpha.
 \end{aligned}$$

Subtracting both equations yields

$$\begin{aligned}
 0 &= \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right) \\
 &\quad \cdot \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right)^{-1/2} \\
 &\quad \cdot \left\{ (G(d_0) + G(d_1))^2 - 2(G(d_0) + G(d_1)) \left(\frac{d_0}{\sqrt[4]{1 + d_0^2}} + \frac{d_1}{\sqrt[4]{1 + d_1^2}} \right) \right. \\
 &\quad \left. + \left((G(d_0) + G(d_1))^2 + 4 \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right)^2 \right) \right\}.
 \end{aligned}$$

By Lemma 4.4, the curly bracket is strictly negative, since we assume that $d_1 > -d_0$. We conclude that

$$0 = \left((1 + d_0^2)^{-1/4} - (1 + d_1^2)^{-1/4} \right),$$

which yields that $d_0 = d_1$. \square

We are now in position to solve the two-dimensional minimization problem.

PROPOSITION 4.6. *Let $0 < \alpha \leq \alpha_{\max}$. There exists $\alpha^* = 1.132372323\dots \in (0, \alpha_{\max})$ such that*

$$\inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1) = \begin{cases} \hat{W}_\alpha(G^{-1}(\frac{c}{2}), G^{-1}(\frac{c}{2})), & 0 < \alpha \leq \alpha^*, \\ c_0^2 - 2\alpha\pi, & \alpha^* < \alpha \leq \alpha_{\max}, \end{cases}$$

where $c \in (0, c_{\max})$ solves $h(c) = \alpha$. In the first case $d_0 = d_1 = G^{-1}(\frac{c}{2})$ is the only point for which the minimum is attained, while it is not attained for $\alpha^* < \alpha \leq \alpha_{\max}$.

Proof. In view of Lemma 4.5 and the symmetry of \hat{W}_α ,

$$\inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1)$$

is the minimum between

$$(4.7) \quad \inf_{d \in (0, \infty)} \hat{W}_\alpha(d, d),$$

$$(4.8) \quad \inf_{d \in \mathbb{R}} \hat{W}_\alpha(d, -d) = 0,$$

and

$$(4.9) \quad \inf_{d \in \mathbb{R}} \hat{W}_\alpha(d, \infty).$$

Since

$$\hat{W}_\alpha(d, d) = 4G(d)^2 - 4\alpha \arctan(d)$$

is certainly negative for $d > 0$ close to 0, we see that $\inf_{d \in (0, \infty)} \hat{W}_\alpha(d, d) < 0$, so we need not consider (4.8). As for (4.9) we have

$$\hat{W}_\alpha(d, \infty) = \left(G(d) + \frac{c_0}{2}\right) \cdot \left(\left(G(d) + \frac{c_0}{2}\right)^2 + 4(1 + d^2)^{-1/2}\right)^{1/2} - 2\alpha \left(\arctan(d) + \frac{\pi}{2}\right).$$

It is sufficient to discuss local minima, since $\hat{W}_\alpha(\infty, \infty)$ is already covered by (4.7) and $\hat{W}_\alpha(-\infty, \infty) = 0$ by (4.8). Passing to the $c = 2G(d)$ -variable, we see that \hat{W}_α attains its minimum on $\{(d_0, d_1) : d_0 \in [-\infty, \infty], d_1 \in [-d_0, \infty]\}$. For fixed $d_0 \in \mathbb{R}$, we infer from (4.6) that for d_1 large enough, $\frac{\partial \hat{W}_\alpha}{\partial d_1} > 0$. This follows since the slowest term $4d_1(1 + d_1^2)^{-5/4}(1 + d_0^2)^{-1/4}$ decays of order $-3/2$ and has a positive coefficient. Hence, the minimum is not attained on $\mathbb{R} \times \{\infty\}$ but either in (∞, ∞) or in the interior of our domain. This proves that

$$(4.10) \quad \inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1) = \inf_{d \in (0, \infty)} \hat{W}_\alpha(d, d).$$

It remains to evaluate the right-hand side of (4.10). Let

$$\phi(d) := \hat{W}_\alpha(d, d) = 4G(d)^2 - 4\alpha \arctan(d).$$

We have

$$\phi'(d) = \frac{8G(d)}{(1 + d^2)^{5/4}} - \frac{4\alpha}{1 + d^2} = \frac{4}{1 + d^2} (h(2G(d)) - \alpha),$$

with h defined in (2.6). Thus, $\phi'(d) = 0$ iff $d = G^{-1}(\frac{c}{2})$, where c is one of the solutions of $h(c) = \alpha$. Only the solution $c \in (0, c_{\max})$ is a local minimum so that

$$\inf_{d \in (0, \infty)} \phi(d) = \min \left(c^2 - 4\alpha \arctan \left(G^{-1} \left(\frac{c}{2} \right) \right), c_0^2 - 2\alpha\pi \right);$$

we take into account that $\phi(0) = 0$ and $\phi(d) < 0$ for small $d > 0$. In order to calculate the last minimum we introduce the following auxiliary function $f : [0, c_{\max}] \rightarrow \mathbb{R}$:

$$f(c) := c_0^2 - 2h(c)\pi - c^2 + 4h(c) \arctan G^{-1} \left(\frac{c}{2} \right).$$

We find that $f(0) = c_0^2 > 0$ and $f(c_{\max}) = -0.6674542140\dots < 0$, and a short calculation shows that

$$f'(c) = \left(4 \arctan G^{-1} \left(\frac{c}{2} \right) - 2\pi \right) h'(c) < 0, \quad c \in (0, c_{\max}),$$

so that f has a unique zero

$$c^* = 1.274998908\dots \in [0, c_{\max}] \text{ with } \alpha^* := h(c^*) = 1.132372323\dots$$

This proves the formula for $\inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1)$. The uniqueness of the minimum for $0 \leq \alpha \leq \alpha^*$ follows from Lemma 2.2. \square

We are now in position to prove Theorems 1.4 and 1.5. The second result is an immediate consequence of (4.10) and Theorem 4.3. As for Theorem 1.4, we focus on the case $0 < \alpha \leq \alpha^*$. Let $c \in (0, c_{\max})$ be the unique solution of $h(c) = \alpha$ with corresponding small solution u_c . Clearly,

$$\begin{aligned} \tilde{W}_\alpha(u_c) &= \hat{W}_\alpha \left(G^{-1} \left(\frac{c}{2} \right), G^{-1} \left(\frac{c}{2} \right) \right) = \inf_{(d_0, d_1) \in \mathbb{R}^2, d_1 \geq -d_0} \hat{W}_\alpha(d_0, d_1) \\ &= \inf_{v \in H^2 \cap H_0^1(0,1)} \tilde{W}_\alpha(v) \end{aligned}$$

by Proposition 4.6 and Lemma 4.2. It remains to show that u_c is the only function in $H^2 \cap H_0^1(0, 1)$ for which the minimum is attained. Suppose that $u \in H^2 \cap H_0^1(0, 1)$ satisfies $\tilde{W}_\alpha(u) = \inf_{v \in H^2 \cap H_0^1(0,1)} \tilde{W}_\alpha(v)$. In view of Lemma 4.1, u is necessarily concave. Let $v = u_{x_0, x_1} \in H^2 \cap H_0^1(0, 1)$ be the function appearing in Lemma 4.2 with d_0, d_1 given by (4.5). Using the minimality of u , Proposition 4.6, and Lemma 4.2, we obtain

$$\tilde{W}_\alpha(u) \leq \tilde{W}_\alpha(u_c) = \hat{W}_\alpha \left(G^{-1} \left(\frac{c}{2} \right), G^{-1} \left(\frac{c}{2} \right) \right) \leq \hat{W}_\alpha(d_0, d_1) = \tilde{W}_\alpha(v) \leq \tilde{W}_\alpha(u).$$

This implies that $\hat{W}_\alpha(G^{-1}(\frac{c}{2}), G^{-1}(\frac{c}{2})) = \hat{W}_\alpha(d_0, d_1)$ and hence by Proposition 4.6 that $d_0 = d_1 = G^{-1}(\frac{c}{2})$ so that $v = u_c$. In particular, we infer with the help of Lemma 4.2 that $u'(0) = v'(0) = u'_c(0)$ and $u'(1) = v'(1) = u'_c(1)$. As a consequence we have $BC_\alpha(u) = BC_\alpha(u_c)$, and therefore $W(u) = W(u_c)$. However, in view of Theorem 2 in [4], u_c is the *unique* minimum of W in the class $M_\beta = \{w \in H^2 \cap H_0^1(0, 1) \mid w'(0) = -w'(1) = \beta\}$ ($\beta = u'_c(0)$) so that we must have $u = u_c$. This completes the proof of Theorem 1.4. \square

For selected values of α , Figure 4.3 shows plots of the function $c \mapsto \tilde{W}_\alpha(u_c)$ on the interval $[0, c_0)$.

The uniqueness part of Theorem 1.4 guarantees that a global minimizer of \tilde{W}_α is symmetric. A more difficult question is whether any solution of (1.1) is necessarily symmetric. The following example (see also [12, p. 461]) shows that such a result certainly does not hold if one extends the class of admissible functions to include graphs with singularities in their first derivatives. Let

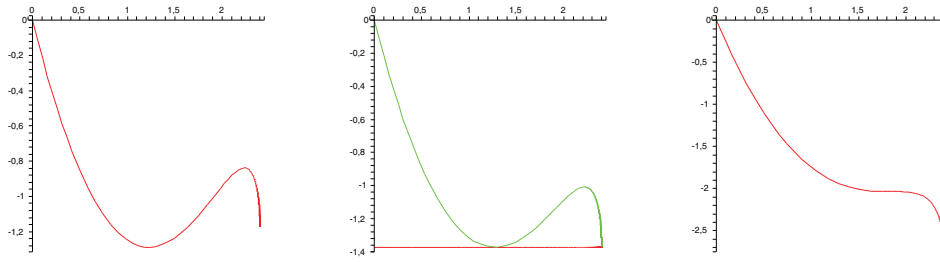


FIG. 4.3. Graphs of the function $c \mapsto \tilde{W}_\alpha(u_c)$ for $\alpha = 1.1$, $\alpha = \alpha^*$, and $\alpha = 1.34$ (left to right).

$$u(x) := \begin{cases} \frac{1}{2}U_0(2x), & 0 \leq x \leq \frac{1}{2}, \\ -\frac{1}{2}U_0(2-2x), & \frac{1}{2} \leq x \leq 1, \end{cases}$$

where the function U_0 was defined in (4.3). Clearly, as a curve in \mathbb{R}^2 , the graph of u is a nonsymmetric solution of (1.1) for $\alpha = 0$. Note, however, that $u \notin H^2(0, 1) \cap H_0^1(0, 1)$ so that u is not a smooth solution of (1.1).

Acknowledgment. The authors thank N. Masel (Minsk) for pointing out an error in an earlier version of Lemma 2.3.

REFERENCES

- [1] M. BAUER AND E. KUWERT, *Existence of minimizing Willmore surface of prescribed genus*, Int. Math. Res. Not., 2003, No. 10 (2003), pp. 553–576.
- [2] A. DALL’ACQUA, K. DECKELNICK, AND H.-CH. GRUNAU, *Classical solutions to the Dirichlet problem for Willmore surfaces of revolution*, Adv. Calc. Var., to appear.
- [3] K. DECKELNICK AND G. DZIUK, *Error analysis of a finite element method for the Willmore flow of graphs*, Interfaces Free Bound., 8 (2006), pp. 21–46.
- [4] K. DECKELNICK AND H.-CH. GRUNAU, *Boundary value problems for the one-dimensional Willmore equation*, Calc. Var. Partial Differential Equations, 30 (2007), pp. 293–314.
- [5] G. DZIUK, E. KUWERT, AND R. SCHÄTZLE, *Evolution of elastic curves in \mathbb{R}^n : Existence and computation*, SIAM J. Math. Anal., 33 (2002), pp. 1228–1245.
- [6] L. EULER, *Opera Omnia*, Ser. 1, 24, Orell Füssli, Zürich, 1952.
- [7] D. KINDERLEHRER AND G. STAMPACCHIA, *An Introduction to Variational Inequalities and Their Applications*, Classics in Appl. Math. 31, SIAM, Philadelphia, 2000.
- [8] E. KUWERT AND R. SCHÄTZLE, *The Willmore flow with small initial energy*, J. Differential Geom., 57 (2001), pp. 409–441.
- [9] E. KUWERT AND R. SCHÄTZLE, *Gradient flow for the Willmore functional*, Comm. Anal. Geom., 10 (2002), pp. 307–339.
- [10] E. KUWERT AND R. SCHÄTZLE, *Removability of point singularities of Willmore surfaces*, Ann. of Math. (2), 160 (2004), pp. 315–357.
- [11] YU. LATUSHKIN, J. PRÜSS, AND R. SCHNAUBELT, *Stable and unstable manifolds for quasilinear parabolic systems with fully nonlinear boundary conditions*, J. Evol. Equ., 6 (2006), pp. 537–576.
- [12] A. LINNÉR, *Explicit elastic curves*, Ann. Global Anal. Geom., 16 (1998), pp. 445–475.
- [13] U. F. MAYER AND G. SIMONETT, *A numerical scheme for axisymmetric solutions of curvature-driven free boundary problems, with applications to the Willmore flow*, Interfaces Free Bound., 4 (2002), pp. 89–109.
- [14] J. C. C. NITSCHKE, *Boundary value problems for variational integrals involving surface curvatures*, Quart. Appl. Math., 51 (1993), pp. 363–387.
- [15] A. POLDEN, *Curves and Surfaces of Least Total Curvature and Fourth-Order Flows*, Ph.D. thesis, University of Tübingen, Tübingen, Germany, 1996.
- [16] R. SCHÄTZLE, *The Willmore Boundary Value Problem*, preprint, University of Tübingen, Tübingen, Germany, 2006.

- [17] L. SIMON, *Existence of surfaces minimizing the Willmore functional*, *Comm. Anal. Geom.*, 1 (1993), pp. 281–326.
- [18] G. SIMONETT, *The Willmore flow near spheres*, *Differential Integral Equations*, 14 (2001), pp. 1005–1014.
- [19] T. J. WILLMORE, *Total Curvature in Riemannian Geometry*, *Ellis Horwood Ser. Math. Appl.*, Ellis Horwood Limited, Chichester, Halsted Press, New York, 1982.

ON THE ENERGY OF SUPERCONDUCTORS IN LARGE AND SMALL DOMAINS*

MATTHIAS KURZKE[†] AND DANIEL SPIRN[‡]

Abstract. We study the Ginzburg–Landau energy functional for superconductors in an applied magnetic field. We focus on asymptotically large or small domains and establish the asymptotic behavior of the energy as a function of the Ginzburg–Landau parameter, applied magnetic field, and domain size. For a large class of domain sizes, we calculate the critical field strength where vortex nucleation becomes energetically favorable, and describe the vorticity of minimizers. For supercritical magnetic field strengths, we recover the energy of a classical Abrikosov vortex lattice. Our findings generalize several known results of Sandier and Serfaty for domains of fixed size.

Key words. Ginzburg–Landau, vortices, asymptotics, critical field

AMS subject classifications. 82D55, 35B25, 35J60

DOI. 10.1137/070687992

1. Introduction. In the widely studied Ginzburg–Landau model, superconducting materials are modeled via the free energy functional

(1.1)

$$G_{phys}(\psi, \hat{A}) = G_0 + \int_{\Omega} \frac{|\operatorname{curl} \hat{A} - H_{ex}|^2}{8\pi} + \frac{1}{2m_*} \left| \left(\hbar \nabla - \frac{ie_* \hat{A}}{c} \right) \psi \right|^2 + \alpha |\psi|^2 + \frac{\beta}{2} |\psi|^4,$$

where Ω is the region occupied by the superconductor, ψ a complex-valued order parameter, \hat{A} the vector potential of the magnetic field, H_{ex} an applied magnetic field, e_* and m_* physical constants corresponding to the charge and mass of superconducting carriers, c the speed of light, and \hbar the Planck constant. G_0 is the energy of the normal state and independent of (ψ, \hat{A}) . The quantities α and β are temperature-dependent constants. We will assume subcritical temperatures, so $\alpha < 0$ and $\beta > 0$. An introduction to the physics of (1.1) explaining the meaning of these constants can be found in [26].

We will work with only two-dimensional domains (corresponding to cylindrical symmetry in the x_3 -direction) and will assume that the domain size is s . As we want to study regimes in which the domain size becomes small or large, we renormalize by assuming that $U = \frac{\Omega}{s}$ is a fixed domain of size $O(1)$.

Setting $u(x) = \sqrt{\frac{\beta}{|\alpha|}} \psi(sx)$, $A(x) = \frac{e_* s}{\hbar c} \hat{A}(sx)$, and $h_{ex} = \frac{e_*}{\hbar c} H_{ex}$, we rewrite the energy as

$$(1.2) \quad G_0 + K \left(\frac{1}{2} \int_U \frac{1}{\ell^2} |\operatorname{curl} A - h_{ex}|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2 + |(\nabla - iA)u|^2 \right),$$

where $K = \frac{\hbar^2 |\alpha|}{2m_* \beta}$, $\ell^2 = \frac{2\pi s^2 e_*^2 |\alpha|}{m_* \beta c^2}$, and $\varepsilon^2 = \frac{\hbar^2}{4m_* |\alpha| s^2}$.

*Received by the editors April 11, 2007; accepted for publication (in revised form) September 8, 2008; published electronically January 23, 2009.

<http://www.siam.org/journals/sima/40-5/68799.html>

[†]Institut für angewandte Mathematik, Universität Bonn, Germany D-53115 (kurzke@iam.uni-bonn.de).

[‡]School of Mathematics, University of Minnesota, Minneapolis, MN 55455 (spirn@math.umn.edu). This author was supported in part by NSF grant DMS-0510121.

Dropping K and G_0 by an affine transformation of the energy, we arrive at the functional

$$(1.3) \quad G_\varepsilon(u, A) = \frac{1}{2} \int_U |\nabla_A u|^2 + \frac{1}{\ell^2} |\operatorname{curl} A - h_{ex}|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2 dx.$$

Using the commonly used comparison length scales of *penetration depth*,

$$(1.4) \quad \lambda = \sqrt{\frac{m_* \beta c^2}{4\pi |\alpha| e_*^2}},$$

and *coherence length*,

$$(1.5) \quad \xi = \frac{\hbar}{\sqrt{2m_* |\alpha|}},$$

we find that our dimensionless parameters ε and ℓ can be expressed as

$$(1.6) \quad \varepsilon = \frac{\xi}{\sqrt{2}s}, \quad \ell = \frac{s}{\sqrt{2}\lambda}.$$

We will consider ε small (i.e., domains of size $s \gg \xi$) and various scalings of ℓ .

We will study *type II* superconductors, those where the Ginzburg–Landau parameter $\kappa = \frac{\lambda}{\xi}$ is large, so the coherence length will be the smallest of the three physical length scales. This corresponds to $\ell \ll \frac{1}{\varepsilon}$, which will be implied by other assumptions throughout this article. Under this limit minimizers will start to energetically favor the formation of vortices once the applied magnetic field grows large enough.

The asymptotics of (1.3) with $\ell \equiv 1$, i.e.,

$$(1.7) \quad \frac{1}{2} \int_U |\nabla_A u|^2 + |\operatorname{curl} A - h_{ex}|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2 dx$$

under the asymptotic limit $\varepsilon \rightarrow 0$, has been widely studied in the past decade and a half.

1.1. Background results. In the last 15 years, there has been considerable progress made in the mathematical understanding of the Ginzburg–Landau model. A major step has been the groundbreaking work of Bethuel, Brezis, and Hélein [5] on the related functional without gauge term, the so-called BBH functional

$$(1.8) \quad E_\varepsilon(u) = \frac{1}{2} \int_U |\nabla u|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2,$$

and much of the analysis for the full gauge-invariant functional G_ε is based on analysis of E_ε . We can only sketch some of the developments for the static Ginzburg–Landau model with magnetic field, and refer the reader to the recent monograph [18] by Sandier and Serfaty, which contains a thorough treatment of vortex solutions and critical fields for vortex nucleation for (1.7).

We mention some works that are of particular relevance to the topic of this article. The first rigorous treatment of (1.7) in the $\varepsilon \rightarrow 0$ limit can be found in Bethuel and Riviere [6], who discovered many important features of the standard Ginzburg–Landau functional. Serfaty [22, 23] built on this work and gave the first rigorous treatment of the critical field question by a study of local minimizers close to the critical field.

The technique used assumptions on the BBH energy (1.8) to obtain an a priori bound on the number of vortices. Using the “vortex ball construction” of Sandier [17] and Jerrard [10], a key ingredient in most of the later research, Sandier and Serfaty [21] were then able to show that the global minimizer below the critical field is indeed vortex-free and $|u_\varepsilon|$ is bounded away from zero.

The structure of global minimizers with an unbounded number of vortices and with external field of order $h_{ex} = O(|\log \varepsilon|)$ was analyzed by Sandier and Serfaty in [20]. This result, combined with a Jacobian compactness theorem [11], was rephrased by Jerrard and Sonar [12] in the framework of Γ -convergence. The limit problem is equivalent to a certain obstacle problem, and the limiting vorticity (after rescaling with $|\log \varepsilon|$) is constant in the set where the obstacle is active.

For asymptotically larger applied magnetic fields ($|\log \varepsilon| \ll h_{ex} \ll \frac{1}{\varepsilon^2}$), the vortices fill the whole domain as an Abrikosov-type lattice with uniform limiting density of vortices. This was established by Sandier and Serfaty in [19].

There are few results as of yet on the influence of the domain size on the behavior of the functional. The asymptotic expansion results of Chapman et al. [8] contain a scaling argument similar to ours above, and then study domains of size $s = O(\xi)$ with $\kappa \gg 1$, which corresponds to very small domains.

Aftalion and Dancer [1] studied critical points of the Ginzburg–Landau energy. For small domains ($\ell < C \min(1, \frac{C}{\kappa})$), they showed that any solution that is not the normal solution (where $u \equiv 0$) will be bounded away from zero, regardless of the external field. For the special case where the domain is a ball, $U_\ell = B_\ell(0)$, they showed that solutions in small domains are necessarily radially symmetric, and there exists a critical field of order $O(\frac{1}{\ell})$ such that, above this field, only the normal solution exists, while a unique superconducting solution exists below this threshold.

A numerical study was performed by Aftalion and Du [2], who studied the response of a superconductor to the raising and lowering of the external field depending on Ginzburg–Landau parameter κ and domain size. They found bifurcation diagrams in several distinct regimes, including a critical line separating type I and type II behavior.

Recently, there have been results by Aydi [3] and Aydi and Sandier [4] on minimizing sequences of the Ginzburg–Landau energy in periodic boxes. In these papers the authors establish a detailed description of the critical field strength for vortex nucleation; vortices are located at minimizers of an explicit renormalized energy. Furthermore, the limiting induced magnetic field satisfies a periodic London equation with point sources at the vortex locations. These results can be understood as a study of the interior behavior of an Abrikosov lattice in the large domain limit.

There have also been a few results that study (1.7) with applied magnetic fields and domain-dependence between h_{c2} and h_{c3} , the regime associated with surface superconductivity. However, we restrict ourselves to field strengths asymptotically below h_{c2} ; hence we do not attempt to review results within this class of field strengths.

There are similarities between the Ginzburg–Landau energy (1.3) and the Chern–Simons–Higgs energy

$$(1.9) \quad G_{csh}(u, A) = \frac{1}{2} \int_U |\nabla_A u|^2 + \frac{\mu^2}{4} \frac{|\text{curl } A - h_{ex}|^2}{|u|^2} + \frac{1}{\varepsilon^2} |u|^2 (1 - |u|^2)^2 dx$$

for an applied magnetic field, h_{ex} , and a bounded simply connected domain, $U \subset \mathbb{R}^2$. The Chern–Simons–Higgs model is an anyon theory that is of interest in connection with high-temperature superconductors and the quantum Hall effect. For an overview

of the study in the self-dual case $\mu = \varepsilon$, see Yang [27].

In [14, 15] the authors proved several results of a similar nature to those found here: For $h_{ex} = H|\log \varepsilon|$ and $G_{csh}(u_\varepsilon, A_\varepsilon) = O(|\log \varepsilon|^2)$, we were able to show Γ -convergence results for the cases $\mu = \mu_\varepsilon \rightarrow \mu_0 \in (0, \infty]$. These enabled us to calculate the critical field for vortex nucleation. The main ingredient in these results is a compactness proof that relates the Jacobian of u , $J(u) = \det \nabla u$, to the energy

$$E_{csh}(u) = \frac{1}{2} \int_U |\nabla u|^2 + \frac{1}{\varepsilon^2} |u|^2 (1 - |u|^2).$$

Using this compactness result from [14] and an energy decomposition, we showed Γ -convergence for finite μ in [14] and for $\mu \rightarrow \infty$ in [15]. For $\mu \rightarrow 0$, we gave an explicit counterexample that illustrates why this method, using a decomposition and bounds for E_{csh} , fails in this case. However, we were later able to show that for h_{ex} much larger than the critical field, and under certain restrictions on μ , the energy of minimizers scales in the same fashion as the energy of an Abrikosov-type lattice just as for the Ginzburg–Landau energy (1.7); see [16].

All of our results here carry over from the Chern–Simons–Higgs energy (1.9) under the assumptions above. In particular, we can extend the results of [14, 15, 16, 25] and understand vortices in non-self-dual Chern–Simons–Higgs for a wider range of parameters and in more detail. Results for (1.3) in the next subsection can be related to results for (1.9) by simply setting $\ell = \frac{2}{\mu}$.

1.2. Main results. In this subsection, we list our main theorems on the behavior of minimizers for various parameter regimes. These results, most of which are generalizations of known results from the last section, provide a partial solution to Open Problem 1 of [18].

Our first result is the calculation of the first critical field where minimizers of the Ginzburg–Landau energy start to have vortices. This field is $O(|\log \varepsilon|)$ if the domains stay bounded and $O(\ell^2 |\log \varepsilon|)$ if the domains are unbounded and ℓ is bounded by a power of $|\log \varepsilon|$.

THEOREM 1.1. *There exists a sequence of critical fields $h_{c_1}(\varepsilon)$ such that any minimizer of the Ginzburg–Landau functional with $h_{ex} < h_{c_1}(\varepsilon) - o(|\log \varepsilon|)$ is vortex-free, while any minimizer with $h_{ex} > h_{c_1}(\varepsilon) + o(|\log \varepsilon|)$ has vortices.*

As $\varepsilon \rightarrow 0$, the critical field $h_{c_1}(\varepsilon)$ satisfies the following expansion: If $\ell_\varepsilon \rightarrow \ell_0$ with $0 \leq \ell_0 < \infty$, then

$$(1.10) \quad \frac{h_{c_1}(\varepsilon)}{|\log \varepsilon|} \rightarrow H_1(\ell_0),$$

where

$$(1.11) \quad H_1(\ell_0) = \frac{1}{2 \max_{\overline{U}} |y_{\ell_0}|},$$

where y_{ℓ_0} is the solution of

$$-\Delta y_{\ell_0} + \ell_0^2 y_{\ell_0} + 1 = 0$$

with Dirichlet boundary conditions $y_{\ell_0} = 0$ on ∂U .

Finally, if $\ell_\varepsilon \rightarrow \infty$ and $\ell_\varepsilon \leq |\log \varepsilon|^\gamma$ for any fixed $\gamma > 0$, then

$$(1.12) \quad \frac{h_{c_1}(\varepsilon)}{\ell_\varepsilon^2 |\log \varepsilon|} \rightarrow \frac{1}{2}$$

as $\varepsilon \rightarrow 0$. Therefore, the critical field scales as $h_{c_1} = \frac{\ell_\varepsilon^2 |\log \varepsilon|}{2}$ in this regime of domain sizes.

For small or bounded domains Theorem 1.1 follows from adapting the proof of Sandier and Serfaty [21], where $\ell \equiv 1$. Formally examining the resulting critical field (1.11), one finds $H_1(\ell_0) \rightarrow \frac{\ell_0^2}{2}$ as $\ell_0 \rightarrow \infty$, so we expect that $h_{c_1} = \frac{\ell_\varepsilon^2 |\log \varepsilon|}{2}$ for any $\ell_\varepsilon \rightarrow +\infty$ and $\ell_\varepsilon \ll \frac{1}{\sqrt{\varepsilon |\log \varepsilon|}}$; see the discussion before Lemma 2.7. We give a proof for the case of large (but not too large) domains in section 2; see Proposition 2.1 for details.

The following results can be used to characterize the minimizers of $(u_\varepsilon, A_\varepsilon)$ for external fields of order $O(|\log \varepsilon|)$ and small or bounded domains. The first step is a Γ -convergence result that relates $G_\varepsilon(u_\varepsilon, A_\varepsilon)$ to a simpler functional that no longer involves ε . We skip some of the detailed convergence statements for ease of presentation. The full statement is given in Theorem 3.1.

THEOREM 1.2. *As $\varepsilon \rightarrow 0$, the functional $\frac{1}{|\log \varepsilon|^2} G_\varepsilon$ is Γ -convergent to $G(v, a)$, where the limit functional G is given by*

$$(1.13) \quad G(v, a) := \begin{cases} \frac{1}{2} \int_U |v - a|^2 + \frac{1}{\ell_0^2} |\operatorname{curl} a - H|^2 + \frac{1}{2} \|\operatorname{curl} v\|_{\mathcal{M}} & \text{if } \ell_0 > 0, \\ \frac{1}{2} \int_U |v - a|^2 + \frac{1}{2} \|\operatorname{curl} v\|_{\mathcal{M}} & \text{if } \ell_0 = 0 \text{ and } \operatorname{curl} a = H, \\ +\infty & \text{otherwise,} \end{cases}$$

under a convergence that includes $(\frac{1}{|\log \varepsilon|}(iu_\varepsilon, \nabla u_\varepsilon) - \frac{1}{|\log \varepsilon|}|u|^2 A_\varepsilon) \rightharpoonup (v - a)$ and $\frac{1}{|\log \varepsilon|} A_\varepsilon \rightharpoonup \operatorname{curl} a$ in L^2 .

Since Γ -convergence and the compactness we have here imply that minimizers of G_ε and of G approximate each other, we study minimizers of G to gain insight into the structure of minimizers of G_ε .

THEOREM 1.3. *If (v_0, a_0) is a minimizer of (1.13) and $\ell_0 > 0$, then $z_0 = \ell_0^{-2}(\operatorname{curl} a_0 - H)$ is the unique minimizer of the following obstacle problem: Minimize*

$$(1.14) \quad F_{\ell_0, H}(z) = \frac{1}{2} \int_U |\nabla z|^2 + \ell_0^2 z^2 + 2zH$$

in the admissible class

$$(1.15) \quad \mathcal{K} = \left\{ z \in H_0^1(U) : z \geq -\frac{1}{2} \text{ a.e. in } U \right\}.$$

The limit (v_0, a_0) satisfies the following additional properties:

$$\begin{aligned} \ell_0^{-2} \operatorname{curl}(\operatorname{curl} a_0 - H) + a_0 &= v_0 \quad \text{in } U, \\ \operatorname{curl} a_0 - H &= 0 \quad \text{on } \partial U, \\ -\frac{1}{2} &\leq z_0 \leq 0, \\ \operatorname{curl} v_0 &\geq 0, \\ \operatorname{spt}(\operatorname{curl} v_0) &\subset \left\{ z_0 = -\frac{1}{2} \right\}. \end{aligned}$$

In the case where $\ell_0 = 0$, we have $\operatorname{curl} a_0 = H$ and obtain a slightly different obstacle problem: Let y_0 be the solution of $-\Delta y_0 = \operatorname{curl} v_0 - H$ with zero boundary conditions. Then y_0 is the unique minimizer of

$$(1.16) \quad F_{0, H}(y) = \frac{1}{2} \int_U |\nabla y|^2 + 2yH$$

in the admissible class

$$(1.17) \quad \mathcal{K} = \left\{ y \in H_0^1(U) : y \geq -\frac{1}{2} \text{ a.e. in } U \right\}.$$

Moreover, $\text{curl } v_0 \geq 0$ and $\text{spt}(\text{curl } v_0) \subset \{y_0 = -\frac{1}{2}\}$.

This theorem, proved later in section 3, implies again the results on the first critical field: When the obstacle is not active, the minimizer satisfies $\text{curl } v = 0$. This happens if and only if $H < H_1(\ell_0)$ with the same fields as above. However, since we rescaled the vorticity to obtain convergence, this only shows that an approximating sequence $(u_\varepsilon, A_\varepsilon)$ has at most $o(|\log \varepsilon|)$ vortices, a result that is weaker than the “no vortices below the critical field” obtained in Theorem 1.1.

Finally, we study minimizers of the Ginzburg–Landau functional with a very large (supercritical) applied external magnetic field and obtain energy asymptotics of a uniform vortex lattice, as follows.

THEOREM 1.4. *Assume that*

$$(1.18) \quad \max\{1, \ell^2\} |\log \varepsilon| \ll h_{ex} \ll \frac{1}{\varepsilon^2};$$

then minimizing sequence $\{u_\varepsilon, A_\varepsilon\}$ satisfies

$$(1.19) \quad G_\varepsilon(u_\varepsilon, A_\varepsilon) = \frac{1}{2} |U| h_{ex} \log \frac{1}{\varepsilon \sqrt{h_{ex}}} (1 - o_\varepsilon(1)).$$

Furthermore, the vortex density is uniform in the limit; see Proposition 5.2 for a precise statement.

The proof of Theorem 1.4 is given in Proposition 5.2 and relies on the upper bound Proposition 4.1. For the lower bound, we have followed the simple proof of Sandier and Serfaty [18] and reduced the theorem by a well-chosen rescaling to the lower bound part of the Γ -convergence theorem in the $|\log \varepsilon|^2$ scaling.

For the upper bound, we use a simple approach with Fourier series (Proposition 4.1). A refined version of the upper bound, (4.2), motivates our conjecture on the behavior close to the critical field for large domains.

Remark 1.5. For $\ell_\varepsilon \rightarrow \infty$ and $h_{ex} = O(\ell^2 |\log \varepsilon|)$, we do not yet have a rigorous result on the structure of minimizers. However, we expect from formal calculations that a uniform lattice, such as those constructed in section 4, should be minimizing. This is further supported by the analysis of Aydi [3], who showed such a result in the periodic setting. Aydi’s upper bound and ours are essentially equivalent.

1.3. Discussion. We conclude the introduction with several unresolved questions regarding asymptotics of (1.3). There is still work to do to complete the answer to Problem 1 in [18]. In particular, a complete phase diagram for the minimizing behavior depending on κ , ℓ , and h_{ex} should be given, including the cross-over between type I and type II behavior that happens for $\kappa \ell = O(1)$, and the results of Aftalion and Du [2] should be made fully rigorous. For such a study, it would also be necessary to understand local minimizers and hysteresis phenomena for slowly changing fields.

It is an interesting problem to further study beginning vortex nucleation close to the critical field in the large domain limit, $\ell \rightarrow +\infty$. Based on the construction of Proposition 4.1 and the structure of the Meissner state, we expect that minimizers exhibit a uniform vortex lattice that fills the whole domain. However, vortices will be far apart and interact only weakly, making this a subtle problem. Finally, it would

be interesting to study (1.3) with applied fields in the “intermediate range,” recently undertaken for (1.7) in [18]. For states with few vortices in large domains, we similarly expect very slow motion for the gradient flow, as vortices will move in an almost flat potential.

2. Critical field calculation. In this section we establish Theorem 1.1. The proof of Theorem 1.1 for $\ell \rightarrow [0, +\infty)$ follows from a direct insertion of ℓ^{-2} into the magnetic field term of (1.7) and following the proof found in [22, 21]. However, when $\ell \rightarrow +\infty$, a simple scaling argument fails, and we need to be more careful. In the following we show that for a substantial class of large-domain asymptotics, the critical field strength is indeed

$$h_{c_1} = \frac{\ell^2}{2} |\log \varepsilon| + o(|\log \varepsilon|),$$

as suggested by the formal analysis of the scaled renormalized energy.

PROPOSITION 2.1. *Let $\ell \rightarrow +\infty$ with $\ell \leq C|\log \varepsilon|^\gamma$ for any fixed $\gamma \in \mathbb{R}^+$, and suppose that (u, A) is a minimizer of the Ginzburg–Landau energy (1.3). Then the first critical field for vortex nucleation is $h_{c_1} = \frac{\ell^2}{2} |\log \varepsilon| = \frac{1}{2} |\log \varepsilon|^{2\gamma+1}$. In particular for $h_{ex} < \frac{\ell^2}{2} |\log \varepsilon|$, any minimizer will satisfy $|u| \geq \frac{3}{4}$ for all ε sufficiently small, and for $h_{ex} > \frac{\ell^2}{2} |\log \varepsilon|$ any minimizer must have a vortex.*

Remark 2.2. Although we establish the conjectured critical field for $\ell = |\log \varepsilon|^\gamma$, we believe the critical field should be true over length scales up to $\ell \leq \frac{C}{\sqrt{\varepsilon} |\log \varepsilon|}$. In particular, the more refined vortex ball estimates found in [13, 18] should be powerful enough to handle larger domains, but in the interest of brevity we consider only ℓ 's that satisfy $1 \ll \ell \leq C|\log \varepsilon|^\gamma$. We establish Proposition 2.1 by using the explicit vortex structure that exists for these intermediate-sized domains.

The Euler–Lagrange equations of (1.3),

$$(2.1) \quad \begin{aligned} 0 &= \nabla_A^2 u + \frac{1}{\varepsilon^2} u (1 - |u|^2), \\ 0 &= \operatorname{curl} \operatorname{curl} A + \ell^2 j_A(u), \end{aligned}$$

in U and $n \cdot \nabla_A u = 0$ and $\operatorname{curl} A = h_{ex}$ on ∂U . Setting the Coulomb gauge, we see that

$$\operatorname{div} A = 0 \text{ in } U, \quad n \cdot A = 0 \text{ on } \partial U.$$

Solutions to (2.1) satisfy the maximum principle

$$(2.2) \quad \|u\|_{L^\infty(U)} \leq 1,$$

the proof of which can be found in [6, 19].

The key to establishing Proposition 2.1 is a good energy decomposition. In order to establish this decomposition we use the following result of Sandier and Serfaty that supplies the vortex structure for our range of ℓ 's. Their result, based on the method of Jerrard [10] and Jerrard and Soner [11] is the following.

PROPOSITION 2.3 (see Sandier and Serfaty[21]). *Let $u : U \rightarrow \mathbb{C}$ be such that $|\nabla u| \leq \frac{C}{\varepsilon}$ and that $E_\varepsilon(u) \leq C|\log \varepsilon|^M$ for $M \geq 2$ a fixed number. Then, for any $\alpha > 0$ there exist disjoint balls $\{B_{r_i}\}_{i \in I}$ of radii r_i such that for sufficiently small ε ,*

1. $\{|u| < \frac{3}{4}\} \subset \cup_{i \in I} B_{r_i}$,

- 2. $\text{card } I \leq C|\log \varepsilon|^M,$
- 3. $r_i \leq \frac{C}{|\log \varepsilon|^\alpha},$
- 4. if $\overline{B_{r_i}} \subset U$ and $d_i = \text{deg}(u, \partial B_{r_i}),$ then

$$(2.3) \quad \int_{B_{r_i}} e_\varepsilon(u) dx \geq \pi |d_i| (|\log \varepsilon| - O(\log |\log \varepsilon|)).$$

Remark 2.4. The result in [21] is restricted to energies of the size $E_\varepsilon(u) \leq K|\log \varepsilon|^2;$ however, the same proof holds for the higher energies in the assumptions found in Proposition 2.1.

We now state our energy decomposition, in the spirit of Bethuel and Riviere [6], Serfaty [22], and Sandier and Serfaty [19].

PROPOSITION 2.5. *Let (u, A) be a minimizer, where A satisfies the Coulomb gauge and $1 \ll \ell \leq C|\log \varepsilon|^\gamma.$ Let $A = h_{ex} \nabla^\perp \xi_\ell + \nabla^\perp \zeta,$ where ξ_ℓ satisfies*

$$(2.4) \quad \begin{aligned} -\frac{1}{\ell^2} \Delta^2 \xi_\ell + \Delta \xi_\ell &= 0 \text{ in } U, \\ \Delta \xi_\ell &= 1 \text{ on } \partial U, \\ \xi_\ell &= 0 \text{ on } \partial U. \end{aligned}$$

Then

$$(2.5) \quad \begin{aligned} G_\varepsilon(u, A) &\geq \sum_{i \in I} \int_{B_{r_i}} e_\varepsilon(u) dx + \frac{1}{\ell^2} \int_U |\Delta \zeta|^2 + G^0 \\ &\quad + 2\pi h_{ex} \sum_{i \in I} d_i \xi_\ell(a_i) - \frac{C}{|\log \varepsilon|}, \end{aligned}$$

where $G^0 = G_\varepsilon(1, h_{ex} \nabla^\perp \xi_\ell).$ Here the vortex balls B_{r_i} and degrees d_i are defined via Proposition 2.3

We prove several intermediate lemmas before attempting the proof of Proposition 2.5. The first facts we establish are on the scaled London equation. This limiting equation for the stream function of the magnetic field potential is the expected Meissner solution.

LEMMA 2.6. *Let ξ_ℓ be the solution of (2.4) with $\ell \gg 1;$ then*

$$(2.6) \quad -\frac{1}{\ell^2} \leq \xi_\ell \leq 0$$

and

$$(2.7) \quad \sup_{x \in \overline{U}} |\xi_\ell| = \frac{1}{\ell^2} (1 - o_\ell(1)).$$

Further,

$$(2.8) \quad \|\xi_\ell\|_{H^2} \leq C \quad \text{and} \quad \|\nabla \xi_\ell\|_{L^\infty} \leq C,$$

where C depends only on $U.$

Proof. These results are similar to results found in [6, 22, 24] for $\ell \equiv 1.$ If $\Delta \xi_\ell = h_\ell$ in U and $\xi_\ell = 0$ on $\partial U,$ then h_ℓ satisfies

$$(2.9) \quad \begin{aligned} -\frac{1}{\ell^2} \Delta h_\ell + h_\ell &= 0 \text{ in } U, \\ h_\ell &= 1 \text{ on } \partial U. \end{aligned}$$

If we let $\chi = \xi_\ell - \frac{1}{\ell^2}(h_\ell - 1)$, then $\Delta\chi = \Delta\xi_\ell - \frac{1}{\ell^2}\Delta(h_\ell - 1) = h_\ell - \frac{1}{\ell^2}\Delta h_\ell = 0$ in U and $\chi = 0$ on ∂U . Therefore, $\chi \equiv 0$, and thus

$$(2.10) \quad \xi_\ell = \frac{1}{\ell^2}(h_\ell - 1).$$

Applying the maximum principle to (2.9) yields $0 < h_\ell < 1$. In particular, if a minimum occurs at a point x_m in the interior of U , then $0 < \frac{1}{\ell^2}\Delta h_\ell(x_m) = h_\ell(x_m)$, and by the boundary condition we see $h_\ell \geq 0$. On the other hand, if the maximum occurs at a point x_M in the interior of U , then $0 > \frac{1}{\ell^2}\Delta h_\ell(x_M) = h_\ell(x_M)$, and by the boundary condition $h_\ell \leq 1$. Applying this to (2.10) yields (2.6).

Next, using the boundary conditions on ξ_ℓ ,

$$\begin{aligned} 0 &= \int_U \xi_\ell \left[-\frac{1}{\ell^2}\Delta^2 \xi_\ell + \Delta \xi_\ell \right] = \int_U \frac{1}{\ell^2} \nabla \xi_\ell \cdot \nabla \Delta \xi_\ell - |\nabla \xi_\ell|^2 \\ &= \int_{\partial U} \frac{1}{\ell^2} \partial_n \xi_\ell - \int_U \left[\frac{1}{\ell^2} |\Delta \xi_\ell|^2 + |\nabla \xi_\ell|^2 \right]. \end{aligned}$$

Thus by (2.9), (2.10), and the bound on h_ℓ ,

$$\int_U \frac{1}{\ell^2} |\Delta \xi_\ell|^2 + |\nabla \xi_\ell|^2 = \frac{1}{\ell^2} \int_U \Delta \xi_\ell = \frac{1}{\ell^2} \int_U h_\ell \leq \frac{1}{\ell^2} |U|.$$

This implies $\|\xi_\ell\|_{H^2(U)} \leq 2\sqrt{|U|}$. Since $0 < \Delta \xi_\ell < 1$, then $\|\xi_\ell\|_{W^{2,p}} \leq C_p$ for any chosen $p > 2$. Therefore, by Sobolev embedding, we have (2.8).

Set $z_\ell = \partial_\ell h_\ell$; then z_ℓ satisfies

$$\Delta z_\ell - \ell^2 z_\ell = 2\ell h_\ell \geq 0$$

in U and $z_\ell = 0$ on ∂U . By the maximum principle $z_\ell \leq 0$; hence, h_ℓ is monotonically decreasing in ℓ for all $x \in U$. Since h_ℓ is bounded below by -1 , then $h_\ell(x) = -1 + o_\ell(1)$ for all $x \in U' \Subset U$. Thus $\max |\xi_\ell| = \frac{1}{\ell^2}(1 - o_\ell(1))$ by (2.10). \square

In order to use Proposition 2.3 we need to establish a bound on $E_\varepsilon(u)$; see (1.8). As we see below, the BBH energy can be much larger than the Ginzburg–Landau energy $G_\varepsilon(u, A)$, since the magnetic field term in the energy can absorb large induced fields generated by a large number of vortices. We have the following result.

LEMMA 2.7. *Let (u, A) be a minimizer of the Ginzburg–Landau energy. Suppose $h_{ex} \leq C\ell^2|\log \varepsilon| \leq \frac{C}{\varepsilon}$ and $G_\varepsilon(u, A) \leq K\ell^2|\log \varepsilon|^2$; then*

$$(2.11) \quad E_\varepsilon(u) \leq C\ell^4|\log \varepsilon|^2$$

and

$$(2.12) \quad \|A\|_{H^1(U)} \leq C\ell^2|\log \varepsilon| \quad \text{and} \quad \|A\|_{H^2(U)} \leq C\ell^3|\log \varepsilon|.$$

Proof. We first establish a uniform H^1 estimate on A . From the assumption on the energy,

$$\int |h - h_{ex}|^2 \leq K\ell^4|\log \varepsilon|^2,$$

and hence from the bound on h_{ex} we see that

$$\|h\|_{L^2(U)} \leq C\ell^2|\log \varepsilon|.$$

Since $\operatorname{div} A = 0$ and $n \cdot A = 0$ on ∂U , there exists ξ such that $\nabla^\perp \xi = A$ and $\xi = 0$ on ∂U . From standard elliptic estimates we get $\|\xi\|_{H^2(U)} \leq C\ell^2|\log \varepsilon|$. Thus

$$\|A\|_{H^1(U)} \leq C\ell^2|\log \varepsilon|.$$

Decomposing $\frac{1}{2}|\nabla_A u|^2 = \frac{1}{2}|\nabla u|^2 - A \cdot j(u) + \frac{1}{2}A^2|u|^2$, we control the cross term via

$$\begin{aligned} A \cdot j(u) &\leq \frac{1}{4} \left| \frac{j(u)}{|u|} \right|^2 + A^2|u|^2 \\ &\leq \frac{1}{4} \left| \frac{j(u)}{|u|} \right|^2 + A^2 + 2\varepsilon^2 A^4 + \frac{1}{8\varepsilon^2} (1 - |u|^2)^2. \end{aligned}$$

Therefore, from the algebraic bounds, the estimate on A , and Sobolev embedding,

$$\begin{aligned} G_\varepsilon(u, A) &\geq E_\varepsilon + \frac{1}{2\ell^2} \|h - h_{ex}\|^2 - \left[\frac{1}{2} E_\varepsilon(u) + \int_U A^2 + 2\varepsilon^2 A^4 \right] \\ &\geq \frac{1}{2} E_\varepsilon(u) - C \left[\|A\|_{H^1(U)}^2 + \varepsilon^2 \|A\|_{H^1(U)}^4 \right] \\ &\geq \frac{1}{2} E_\varepsilon(u) - C\ell^4 |\log \varepsilon|^2 - C\varepsilon^2 \ell^8 |\log \varepsilon|^4. \end{aligned}$$

The upper bound on $E_\varepsilon(u)$ follows.

In order to establish higher bounds on A we use the Euler–Lagrange equation $-\nabla^\perp h = \ell^2 j_A(u)$. Therefore,

$$\|\nabla h\|_{L^2(U)} \leq \ell^2 \|\nabla_A u\|_{L^2(U)} \|u\|_{L^\infty(U)} \leq C\ell^3 |\log \varepsilon|,$$

and hence (2.8). \square

The fact that $E_\varepsilon(u)$ can have a much larger energy than $G_\varepsilon(u, A)$ is an essential difference in the large- ℓ asymptotics. It implies a more complicated global vortex structure. Given the energy bound on $E_\varepsilon(u)$, we can split apart the full Ginzburg–Landau energy into its chief components. We start with an initial energy splitting.

LEMMA 2.8. *We can decompose $G_\varepsilon(u, A) = G_\varepsilon(u, \nabla^\perp \xi) = G_\varepsilon(u, h_{ex} \nabla^\perp \xi_\ell + \nabla^\perp \zeta)$ as*

$$\begin{aligned} (2.13) \quad G_\varepsilon(u, A) &\geq \frac{1}{2} \int_U |\nabla u - iu \nabla^\perp \xi|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2 + \frac{1}{2\ell^2} \int_U |\Delta \xi|^2 \\ &\quad + \int_U (\nabla u, -ih_{ex} \nabla^\perp \xi_\ell u) + G^0 - \frac{C}{|\log \varepsilon|}. \end{aligned}$$

Proof. We decompose the energy in a series of steps.

Our first step is to compute the approximate energy of the Meissner state via the method of Serfaty [22]. Since $\operatorname{div} A = 0$ and $n \cdot A = 0$ in ∂U , we can write $A = \nabla^\perp \xi$ with $\xi = 0$ on ∂U and so $\Delta \xi = h$. We further decompose $\nabla^\perp \xi$ as $\nabla^\perp \xi = h_{ex} \nabla^\perp \xi_\ell + \nabla^\perp \zeta$, where $\zeta = \Delta \zeta = 0$ on ∂U and where ξ_ℓ satisfies (2.4). Consider now the Meissner energy associated with $G^0 = G_\varepsilon(1, h_{ex} \nabla^\perp \xi_\ell)$. We compute the form of the Meissner energy, setting $(u_0, A_0) = (1, h_{ex} \nabla^\perp \xi_\ell)$. Then

$$\begin{aligned} G_\varepsilon(u_0, A_0) &= \frac{1}{2} \int_U |\nabla_{A_0} u_0|^2 + \frac{1}{\ell^2} |\operatorname{curl} A_0 - h_{ex}|^2 + \frac{1}{2\varepsilon^2} (1 - |u_0|^2)^2 \\ &= \frac{h_{ex}^2}{2} \int_U |\nabla \xi_\ell|^2 + \frac{1}{\ell^2} |\Delta \xi_\ell - 1|^2. \end{aligned}$$

Multiplying ξ_ℓ against $-\frac{1}{\ell^2}\Delta^2\xi_\ell + \Delta\xi_\ell$ and integrating over U yields

$$\int_U |\nabla\xi_\ell|^2 + \frac{1}{\ell^2} |\Delta\xi_\ell|^2 dx = \frac{1}{\ell^2} \int_U \Delta\xi_\ell.$$

We use the above identity to rewrite the Meissner energy as

$$\begin{aligned} G^0 &= \frac{1}{2} \int_U h_{ex}^2 |\nabla\xi_\ell|^2 + \frac{h_{ex}^2}{\ell^2} |\Delta\xi_\ell|^2 + \frac{h_{ex}^2}{\ell^2} - 2\frac{h_{ex}^2}{\ell^2} \Delta\xi_\ell \\ (2.14) \quad &= -\frac{h_{ex}^2}{2} \int_U \left[|\nabla\xi_\ell|^2 + \frac{1}{\ell^2} |\Delta\xi_\ell|^2 \right] + \frac{h_{ex}^2}{2\ell^2} |U|. \end{aligned}$$

Therefore the Meissner energy is of order $O(\frac{h_{ex}^2}{\ell^2})$.

Next we write

$$\begin{aligned} \frac{1}{2} \int_U |\nabla_A u|^2 &= \frac{1}{2} \int_U |\nabla u - ih_{ex} \nabla^\perp \xi_\ell u - i\nabla^\perp \zeta u|^2 \\ &= \frac{1}{2} \int_U |\nabla u - i\nabla^\perp \zeta u|^2 + \frac{h_{ex}^2}{2} \int_U |u|^2 |\nabla^\perp \xi_\ell|^2 \\ &\quad + \int_U (\nabla u - i\nabla^\perp \zeta u, -ih_{ex} \nabla^\perp \xi_\ell u) \end{aligned}$$

and

$$\begin{aligned} &\int_U (\nabla u - i\nabla^\perp \zeta u, -ih_{ex} \nabla^\perp \xi_\ell u) \\ &= \int_U (\nabla u, -ih_{ex} \nabla^\perp \xi_\ell u) + h_{ex} \int_U |u|^2 \nabla\xi_\ell \cdot \nabla\zeta. \end{aligned}$$

Therefore, we can write the Ginzburg–Landau energy as

$$\begin{aligned} G_\varepsilon(u, A) &= \frac{1}{2} \int_U |\nabla u - i\nabla^\perp \zeta u|^2 + \frac{1}{2\varepsilon^2} (1 - |u|^2)^2 \\ (2.15) \quad &+ \int_U (\nabla u, -ih_{ex} \nabla^\perp \xi_\ell u) \\ &+ \frac{h_{ex}^2}{2} \int_U (|u|^2 - 1) |\nabla^\perp \xi_\ell|^2 + h_{ex} \int_U (|u|^2 - 1) \nabla\xi_\ell \cdot \nabla\zeta \\ &+ \frac{1}{2\ell^2} \int_U |h - h_{ex}|^2 + \frac{h_{ex}^2}{2} \int_U |\nabla^\perp \xi_\ell|^2 + h_{ex} \int_U \nabla\xi_\ell \cdot \nabla\zeta. \end{aligned}$$

The terms in the third line of (2.15) are small since

$$\begin{aligned} h_{ex}^2 \int_U (|u|^2 - 1) |\nabla\xi_\ell|^2 &\leq Ch_{ex}^2 \|\nabla\xi_\ell\|_{L^\infty}^2 \|1 - |u|^2\|_{L^2} \leq C\varepsilon h_{ex}^2 E_\varepsilon^{\frac{1}{2}}(u) \\ &\leq C\varepsilon \ell^6 |\log \varepsilon|^3 \leq C\varepsilon |\log \varepsilon|^{6\gamma+3} \leq \frac{C}{|\log \varepsilon|} \end{aligned}$$

and

$$\begin{aligned} h_{ex} \int_U (|u|^2 - 1) \nabla\xi_\ell \cdot \nabla\zeta &\leq Ch_{ex} \|\nabla\xi_\ell\|_{L^\infty} \|\nabla\zeta\|_{L^2} \|1 - |u|^2\|_{L^2} \\ &\leq C\varepsilon h_{ex} \|A - h_{ex} \nabla^\perp \xi_\ell\|_{L^2} E_\varepsilon^{\frac{1}{2}}(u) \\ &\leq C\varepsilon \ell^2 |\log \varepsilon| (\ell^2 |\log \varepsilon| + \ell^2 |\log \varepsilon|) \ell^2 |\log \varepsilon| \\ &\leq C\varepsilon |\log \varepsilon|^{6\gamma+3} \leq \frac{C}{|\log \varepsilon|}. \end{aligned}$$

For the fourth line of (2.15) we have

$$\begin{aligned} & \frac{1}{2\ell^2} \int_U |h - h_{ex}|^2 + \frac{h_{ex}^2}{2} \int_U |\nabla^\perp \xi_\ell|^2 + h_{ex} \int_U \nabla \xi_\ell \cdot \nabla \zeta \\ &= \frac{1}{2\ell^2} \int_U |h_{ex} \Delta \xi_\ell - h_{ex} + \Delta \zeta|^2 + \frac{h_{ex}^2}{2} \int_U |\nabla \xi_\ell|^2 + h_{ex} \int_U \nabla \xi_\ell \cdot \nabla \zeta \\ &= \frac{h_{ex}^2}{2\ell^2} \int_U |\Delta \xi_\ell - 1|^2 + \frac{h_{ex}^2}{2} \int_U |\nabla \xi_\ell|^2 + \frac{1}{2\ell^2} \int_U |\Delta \zeta|^2 \\ & \quad + \frac{h_{ex}}{\ell^2} \int_U (\Delta \xi_\ell - 1) \Delta \zeta + h_{ex} \int_U \nabla \xi_\ell \cdot \nabla \zeta. \end{aligned}$$

Multiplying ζ against (2.4) and integrating over U , we have

$$0 = -\frac{1}{\ell^2} \int_U \Delta \zeta (\Delta \xi_\ell - 1) - \int_U \nabla \xi_\ell \cdot \nabla \zeta,$$

and then

$$\begin{aligned} (2.16) \quad & \frac{1}{2\ell^2} \int_U |h - h_{ex}|^2 + \frac{h_{ex}^2}{2} \int_U |\nabla^\perp \xi_\ell|^2 + h_{ex} \int_U \nabla \xi_\ell \cdot \nabla \zeta \\ &= G^0 + \frac{1}{2\ell^2} \int_U |\Delta \zeta|^2. \end{aligned}$$

Combining (2.15), the bounds on the third line, and (2.16) yields (2.13). \square

We can now prove Proposition 2.5 by carefully extracting the concentration of the Ginzburg–Landau energy against the magnetic field potential ξ_ℓ . Note that there are potentially an unbounded number of vortices, so we need to extract good decay on each vortex ball.

Proof of Proposition 2.5. We follow the approach in [21] for $\ell \equiv 1$. The first step is to establish the concentration in the cross term $\int \nabla^\perp \xi_\ell \cdot j(u)$. In particular, we claim

$$(2.17) \quad \left| \int_U (\nabla u, -ih_{ex} \nabla^\perp \xi_\ell u) - 2\pi h_{ex} \sum_{i \in I} d_i \xi_\ell(a_i) \right| \leq \frac{C}{|\log \varepsilon|},$$

where a_j is the center of the vortex ball B_{r_i} and I is the vortex ball collection.

Step 1. Since $E_\varepsilon(u) \leq C\ell^4 |\log \varepsilon|^2 \leq C|\log \varepsilon|^{4\gamma+2}$ and $h_{ex} = C\ell^2 |\log \varepsilon| = C|\log \varepsilon|^{2\gamma+1}$, then by Proposition 2.3 we have balls $\{B_{r_i}\}_{i \in I}$ such that

$$\left\{ |u| < \frac{3}{4} \right\} \subset \cup_{i \in I} B_{r_i}, \quad \text{card } I \leq C|\log \varepsilon|^{4\gamma+2}, \quad r_i \leq \frac{C}{|\log \varepsilon|^{10\gamma+6}},$$

if $\overline{B_{r_i}} \subset U$ and $d_i = \text{deg}(u, \partial B_{r_i})$, then $\int_{B_{r_i}} e_\varepsilon(u) dx \geq \pi |d_i| (|\log \varepsilon| - O(\log |\log \varepsilon|))$,

where we chose $\alpha = 10\gamma + 6$ in Proposition 2.3. Therefore,

$$\begin{aligned} \left| \int_{\cup_i B_{r_i}} (\nabla u, -ih_{ex} \nabla^\perp \xi_\ell u) \right| &\leq (\text{card } I) h_{ex} \|\nabla u\|_{L^2} \max_{i \in I} r_i \\ &\leq C|\log \varepsilon|^{8\gamma+4-\alpha} \leq \frac{C}{|\log \varepsilon|}. \end{aligned}$$

Setting $\tilde{U} = U \setminus \cup_{i \in I} B_{r_i}$, for $v = \frac{u}{|u|}$ we have

$$\int_{\tilde{U}} (\nabla u, -i h_{ex} \nabla^\perp \xi_\ell u) = h_{ex} \int_{\tilde{U}} \nabla \xi_\ell \times j(v) + h_{ex} \int_{\tilde{U}} \nabla \xi_\ell \times (j(u) - j(v)).$$

The second term is small, using (2.8) and (2.11), since

$$\begin{aligned} h_{ex} \int_{\tilde{U}} \nabla \xi_\ell \times (j(u) - j(v)) &\leq h_{ex} \|\nabla^\perp \xi_\ell\|_{L^\infty} \int_{\tilde{U}} \left| j(u) - \frac{j(u)}{|u|^2} \right| \\ &\leq C h_{ex} \int_{\tilde{U}} \frac{|j(u)|}{|u|} \frac{||u|^2 - 1|}{|u|} \leq C \varepsilon h_{ex} E_\varepsilon(u) \\ &\leq C \varepsilon |\log \varepsilon|^{6\gamma+3} = \frac{C}{|\log \varepsilon|}, \end{aligned}$$

where we used $|u| \geq \frac{1}{2}$ in \tilde{U} in the second line. Therefore,

$$(2.18) \quad \left| \int_{\tilde{U}} (\nabla u, -i h_{ex} \nabla^\perp \xi_\ell u) - h_{ex} \int_{\tilde{U}} \nabla \xi_\ell \times j(v) \right| \leq \frac{C}{|\log \varepsilon|}.$$

Next for $J = \{i \text{ such that } \overline{B_{r_i}} \subset U\}$ we claim we can extract the following bound:

$$(2.19) \quad \left| h_{ex} \int_{B_{r_i}} \nabla \xi_\ell \times j(v) - 2\pi h_{ex} d_i \xi_\ell(a_i) \right| \leq \frac{C}{|\log \varepsilon|^{4\gamma+3}}.$$

For $\Omega_i = B_{r_i} \cap \{x \in U \text{ such that } |u| \leq \frac{1}{2}\}$, $\Omega_i \cap \partial B_{r_i} = \emptyset$. Since $|u| \geq \frac{1}{2}$ in Ω_i , by Stokes' theorem,

$$\begin{aligned} &h_{ex} \left| \int_{\partial B_{r_i}} (\xi_\ell - \xi_\ell(a_i)) j(v) \cdot \tau - \int_{\partial \Omega_i} (\xi_\ell - \xi_\ell(a_i)) j(v) \cdot \tau \right| \\ &= h_{ex} \left| \int_{B_i \setminus \Omega_i} \nabla \xi_\ell \times j(u) \right| \leq C h_{ex} \|\nabla \xi_\ell\|_{L^\infty} \|\nabla u\|_{L^2} r_i \\ &\leq C |\log \varepsilon|^{4\gamma+2-\alpha} \leq \frac{C}{|\log \varepsilon|^{4\gamma+3}}; \end{aligned}$$

thus

$$(2.20) \quad \left| h_{ex} \int_{\partial B_{r_i}} (\xi_\ell - \xi_\ell(a_i)) j(v) \cdot \tau - \int_{\partial \Omega_i} (\xi_\ell - \xi_\ell(a_i)) j(v) \cdot \tau \right| \leq \frac{C}{|\log \varepsilon|^{4\gamma+3}}.$$

On the other hand, since $|u| = \frac{1}{2}$ on $\partial\Omega_i$ we find

$$\begin{aligned} & h_{ex} \left| \int_{\partial\Omega_i} (\xi_\ell - \xi_\ell(a_i)) j(v) \cdot \tau \right| \\ &= h_{ex} \left| \int_{\partial\Omega_i} (\xi_\ell - \xi_\ell(a_i)) \frac{j(u)}{|u|^2} \cdot \tau \right| = 4h_{ex} \left| \int_{\partial\Omega_i} (\xi_\ell - \xi_\ell(a_i)) j(u) \cdot \tau \right| \\ &= 4h_{ex} \left| \int_{\Omega_i} \operatorname{curl} [(\xi_\ell - \xi_\ell(a_i)) j(u)] \right| \\ &\leq 4h_{ex} \left| \int_{\Omega_i} \nabla \xi_\ell \times j(u) \right| + 8h_{ex} \left| \int_{\Omega_i} (\xi_\ell - \xi_\ell(a_i)) J(u) \right| \\ &\leq Ch_{ex} \|\nabla \xi_\ell\|_{L^\infty} \|\nabla u\|_{L^2} r_i + Ch_{ex} \|\nabla \xi_\ell\|_{L^\infty} \|\nabla u\|_{L^2}^2 r_i \\ &\leq C|\log \varepsilon|^{4\gamma+2-\alpha} + C|\log \varepsilon|^{6\gamma+3-\alpha} \leq \frac{C}{|\log \varepsilon|^{4\gamma+3}}, \end{aligned}$$

and consequently, since $\operatorname{card} I \leq C|\log \varepsilon|^{4\gamma+2}$ and $\int_{\partial B_{r_i}} \xi_\ell(a_i) j(v) \cdot \tau = 2\pi d_i$,

$$(2.21) \quad \sum_{i \in J} \left| h_{ex} \int_{\partial B_{r_i}} \xi_\ell j(v) \cdot \tau - h_{ex} \int_{\partial B_{r_i}} \xi_\ell(a_i) j(v) \cdot \tau \right| \leq \frac{C}{|\log \varepsilon|}.$$

Finally, for the balls that intersect ∂U , $I \setminus J$. Since $\xi_\ell = 0$ on ∂U , then for $\Omega_i = B_i \cap \{x \in U \text{ such that } |u| \leq \frac{1}{2}\}$ we follow the above argument and see

$$\begin{aligned} \left| h_{ex} \int_{\partial B_{r_i} \cap U} \xi_\ell j(v) \cdot \tau \right| &\leq \left| h_{ex} \int_{\partial(\Omega_i \cap U)} \xi_\ell j(v) \cdot \tau \right| + \frac{C}{|\log \varepsilon|^{4\gamma+3}} \\ &\leq 4h_{ex} \left| \int_{\Omega_i \cap U} \nabla \xi_\ell \times j(v) + 2\xi_\ell J(v) \right| + \frac{C}{|\log \varepsilon|^{4\gamma+3}} \\ &\leq \frac{C}{|\log \varepsilon|^{4\gamma+3}}. \end{aligned}$$

Combining this estimate along with $\operatorname{card} I \leq C|\log \varepsilon|^{4\gamma+2}$ and (2.21) yields estimate (2.17).

Step 2. We bound $\int_U |\nabla u - i\nabla^\perp \zeta u|^2 \geq \int_{\cup_{i \in I} B_{r_i}} |\nabla u|^2 - \frac{C}{|\log \varepsilon|}$. In particular,

$$\begin{aligned} \int_U |\nabla u - i\nabla^\perp \zeta u|^2 &\geq \int_{\cup_{i \in I} B_{r_i}} |\nabla u - i\nabla^\perp \zeta u|^2 \\ &= \int_{\cup_{i \in I} B_{r_i}} |\nabla u|^2 - 2\nabla^\perp \zeta \cdot j(u) + |\nabla \zeta|^2 |u|^2. \end{aligned}$$

From (2.12) we see $\|A\|_{L^\infty(U)} \leq C\|A\|_{H^2(U)} \leq C|\log \varepsilon|^{3\gamma+2}$, and thus

$$\begin{aligned} \left| \int_{\cup_{i \in I} B_{r_i}} \nabla^\perp \zeta \cdot j(u) \right| &\leq (\operatorname{card} I) \|A - h_{ex} \nabla^\perp \xi_\ell\|_{L^\infty} \|\nabla u\|_{L^2} \max_{i \in I} r_i \\ &\leq C|\log \varepsilon|^{4\gamma+2} \left(|\log \varepsilon|^{3\gamma+2} + |\log \varepsilon|^{2\gamma+1} \right) |\log \varepsilon|^{2\gamma+1} |\log \varepsilon|^{-\alpha} \\ &\leq C|\log \varepsilon|^{9\gamma+5-\alpha} \leq \frac{C}{|\log \varepsilon|}, \end{aligned}$$

and so

$$(2.22) \quad \int_U |\nabla u - i\nabla^\perp \zeta u|^2 \geq \int_{\cup_{i \in I} B_{r_i}} |\nabla u|^2.$$

Combining (2.13) with (2.17) and (2.22) yields (2.5). \square

We are finally in the position to establish the next claim.

Proof of Proposition 2.1. The first part of the proof establishes that a minimizing sequence must be in the Meissner state when $h_{ex} < \frac{\ell^2 |\log \varepsilon|}{2}$.

Step 1. From Proposition 2.5 and the minimality of (u, A)

$$\begin{aligned} G^0 &\geq G_\varepsilon(u, A) \\ &\geq \sum_{i \in I} \int_{B_{r_i}} e_\varepsilon(u) dx + \frac{1}{2\ell^2} \int_U |\Delta \zeta|^2 + G^0 + 2\pi h_{ex} \sum_{i \in I} d_i \xi_\ell(a_i) - \frac{C}{|\log \varepsilon|}. \end{aligned}$$

Therefore, since $\xi_\ell \leq 0$, we use (2.7) and lower bound in Proposition 2.3.4 to get

$$\begin{aligned} \pi \sum_{i \in I} |d_i| (|\log \varepsilon| + O(\log |\log \varepsilon|)) &\leq 2\pi h_{ex} \sum_{i \in I} d_i |\xi_\ell(a_i)| \\ &\leq 2\pi h_{ex} \left(\sum_{i \in I} |d_i| \right) \max |\xi_\ell| \\ &\leq \left(\sum_{i \in I} |d_i| \right) 2\pi \frac{h_{ex}}{\ell^2} (1 - o_\ell). \end{aligned}$$

So if $\sum_{i \in I} |d_i| \neq 0$, then

$$h_{ex} \geq \frac{\ell^2}{2} (|\log \varepsilon| + O(\log |\log \varepsilon|)).$$

Hence, for $h_{ex} < \frac{\ell^2}{2} |\log \varepsilon|$, either $\deg(u, \partial B_{r_i}) = 0$ or $B_{r_i} \cap U \neq \emptyset$. It is straightforward to show from this point that $|u| \geq \frac{3}{4}$ in U ; see [5, 6].

Step 2. We now complete the proof of the critical field strength. In particular, we show that if $h_{ex} > \frac{\ell^2 |\log \varepsilon|}{2}$, then there must be a vortex. We prove this by contradiction. Let $(u_\varepsilon, A_\varepsilon)$ be a minimizing sequence with $\sum_{j \in J} |d_j| = 0$; then we claim $G_\varepsilon(u, A) \geq G_\varepsilon(1, h_{ex} \nabla^\perp \xi_\ell) - \frac{C}{|\log \varepsilon|}$.

In order to get better bounds on $\nabla^\perp \zeta = A - h_{ex} \nabla^\perp \xi_\ell$, we replace lower bound (2.22) with

$$\frac{1}{2} \int_U |\nabla u - i\nabla^\perp \zeta u|^2 \geq \frac{1}{2} \int_U |\nabla u|^2 - 2j(u) \cdot \nabla^\perp \zeta + |\nabla \zeta|^2 - \frac{C}{|\log \varepsilon|},$$

where we used the argument for the estimate of the third line of (2.15) in the proof of Lemma 2.8. By (2.8) and (2.12) we see that ζ is continuous. Since there are no nontrivial-degree vortex balls, then by an argument identical to the proof of (2.17) we have the lower bound

$$G_\varepsilon(u, A) \geq E_\varepsilon(u) + G^0 + \frac{1}{2} \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 - \frac{C}{|\log \varepsilon|}.$$

Since (u, A) is an energy minimizer, $G^0 \geq G_\varepsilon(u, A)$, and so $E_\varepsilon(u) + \frac{1}{2} \int_U |\nabla \zeta|^2 \leq \frac{C}{|\log \varepsilon|}$. Even more so, the boundary condition $\zeta = 0$ implies $\zeta \rightarrow 0$ and $E_\varepsilon(u) \rightarrow 0$ as $\varepsilon \rightarrow 0$. We see that

$$(2.23) \quad G_\varepsilon(u, A) \geq G^0 - \frac{C}{|\log \varepsilon|}$$

when $\sum_{j \in J} |d_j| = 0$.

To prove that $G_\varepsilon(u, A)$ is no longer the Meissner state, we construct a sequence of functions $(u_\varepsilon, A_\varepsilon)$ which have lower energy than the Meissner energy when $h_{ex} > \frac{\ell^2 |\log \varepsilon|}{2}$. Set $A_\varepsilon = h_{ex} \nabla^\perp \xi_\ell + \nabla^\perp \zeta$, where ξ_ℓ is defined in (2.4) and

$$(2.24) \quad -\frac{1}{\ell^2} \Delta^2 \zeta + \Delta \zeta = 2\pi \delta_a \text{ in } U, \quad \Delta \zeta = \xi_\ell = 0 \text{ on } \partial U.$$

To define $u_\varepsilon = \rho_\varepsilon e^{i\varphi_\varepsilon}$ we set $\nabla \varphi_\varepsilon = A_\varepsilon + \frac{1}{\ell^2} \nabla^\perp \text{curl } A_\varepsilon$ and

$$\rho_\varepsilon = \begin{cases} 0 & |x - a| \leq \frac{\varepsilon}{2}, \\ 1 & |x - a| \geq \varepsilon. \end{cases}$$

Then for any $B_R \supset \{a\}$, $\int_{\partial B_R} \partial_\tau \varphi_\varepsilon = \int_{B_R} h_\varepsilon - \frac{1}{\ell^2} \Delta h_\varepsilon = 2\pi$, which correctly quantizes the phase. A straightforward calculation shows that $E_\varepsilon(u_\varepsilon) \leq \pi \log \frac{\text{diam } U}{\varepsilon} + C \leq \pi |\log \varepsilon| + C$, where C is a fixed constant. The arguments in section 4 contain more refined upper bound calculations; however, they are similar in spirit.

Following Step 1 of the proof of Proposition 2.5 yields

$$\begin{aligned} & \frac{1}{2} \int_U A_\varepsilon^2 |u_\varepsilon|^2 + \frac{1}{\ell^2} |\text{curl } A_\varepsilon - h_{ex}|^2 \\ & \leq \frac{h_{ex}^2}{2} \int_U |\nabla \xi_\ell|^2 + \frac{1}{\ell^2} |\Delta \xi_\ell - h_{ex}|^2 + \frac{1}{2} \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 + \frac{C}{|\log \varepsilon|}. \end{aligned}$$

Again we decompose $\|\nabla u - i\nabla^\perp \zeta u\|_{L^2(U)}^2 = \|\nabla u\|_{L^2(U)}^2 - 2 \int j(u) \cdot \nabla^\perp \zeta + \|\nabla^\perp \zeta u\|_{L^2(U)}^2$, and a similar calculation as in Step 1 shows

$$\begin{aligned} G_\varepsilon(u_\varepsilon, A_\varepsilon) & \leq E_\varepsilon(u_\varepsilon) - \int_U j(u_\varepsilon) \cdot \nabla^\perp (h_{ex} \xi_\ell + \zeta) + G^0 \\ & \quad + \frac{1}{2} \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 + \frac{C}{|\log \varepsilon|} \\ & \leq G^0 + \pi |\log \varepsilon| - \frac{2\pi h_{ex}}{\ell^2} (1 - o_\ell(1)) + C \\ & \quad + \frac{1}{2} \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 + 2\pi \zeta(a), \end{aligned}$$

where we used (2.7) in the last inequality. Multiplying (2.24) by ζ and integrating over U shows $2\pi \zeta(a) = - \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 < 0$; hence $\frac{1}{2} \int_U |\nabla \zeta|^2 + \frac{1}{\ell^2} |\Delta \zeta|^2 + 2\pi \zeta(a) < 0$. Therefore,

$$G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq G^0 + \pi |\log \varepsilon| - \frac{2\pi h_{ex}}{\ell^2} (1 - o_\ell(1)) + C.$$

Since $h_{ex} > \frac{\ell^2 |\log \varepsilon|}{2}$, there exists $\delta > 0$, bounded away from zero, such that $\pi |\log \varepsilon| - \frac{2\pi h_{ex}}{\ell^2} (1 - o_\ell(1)) + C < C - \delta |\log \varepsilon| < -\frac{|C|}{2}$ for ε small enough; thus

$$G_\varepsilon(u_\varepsilon, A_\varepsilon) < G^0 - \frac{|C|}{2} < G_\varepsilon(u, A).$$

Therefore, a vortex-less configuration *cannot* be minimizing in the $h_{ex} > \frac{\ell^2 |\log \varepsilon|}{2}$ regime. \square

Remark 2.9. For values of h_{ex} well above the critical field, we expect the minimizers to be similar to the functions constructed in the proof of (4.2) in section 4.

Remark 2.10. The proof of the critical field for $\ell_0 \in [0, +\infty)$ proceeds in the same way as for the proof of Proposition 2.1 and can be done by a suitable modification of the method in [21]. Since we handled the more difficult case $\ell \rightarrow +\infty$ such that $\ell_\varepsilon = |\log \varepsilon|^\gamma$ for some $\gamma \in \mathbb{R}^+$, we leave out the proof for the case $\ell_\varepsilon \rightarrow \ell_0 \in [0, +\infty)$.

3. Obstacle problem for small and bounded domains. In this section, we study the functional (1.3) where $\ell_\varepsilon \rightarrow \ell_0 \in [0, \infty)$, i.e., for domain sizes s that are smaller than or comparable with the penetration depth ($s \ll \lambda$ or $s = O(\lambda)$), in the critical scaling of energy and magnetic field.

The following result is a generalization of Theorem 1.3 in [12] (where it is proved for $\ell = 1$). Closely related results in the context of the Chern–Simons–Higgs energy were shown by the authors in [14, Theorem 1.3] and [15, Theorem 3]. We state the theorem in its gauge-invariant form.

THEOREM 3.1. *Let $(u_\varepsilon, A_\varepsilon)$ be a sequence with $G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq K|\log \varepsilon|^2$ and assume that h_{ex} satisfies $\frac{h_{ex}}{|\log \varepsilon|} \rightarrow H$ for some $H \geq 0$ and $\ell_\varepsilon \rightarrow \ell_0 \in [0, \infty)$. Define the following rescaled quantities:*

$$\begin{aligned} a_\varepsilon &:= \frac{1}{|\log \varepsilon|} A_\varepsilon, \\ v_\varepsilon &:= \frac{1}{|\log \varepsilon|} (iu_\varepsilon, \nabla u_\varepsilon), \\ w_\varepsilon &:= v_\varepsilon - |u_\varepsilon|^2 a_\varepsilon. \end{aligned}$$

Then $\text{curl } a_\varepsilon$ is weakly compact in $L^2(U)$, and w_ε is weakly compact in L^p for $1 \leq p < 2$. Furthermore, $\frac{w_\varepsilon}{|u_\varepsilon|}$ converges weakly in L^2 if and only if w_ε converges weakly, and the weak limits are equal.

Any weak limit of $(w_\varepsilon, \text{curl } a_\varepsilon)$ can be expressed in the form $(v - a, \text{curl } a)$ for some $(v, a) \in L^2(U; \mathbb{R}^2) \times H^1(U; \mathbb{R}^2)$ such that $\text{curl } v$ is a Radon measure. In addition, we have the following $\Gamma - \liminf$ inequality:

$$\liminf_{\varepsilon \rightarrow 0} \frac{1}{|\log \varepsilon|^2} G_\varepsilon(u_\varepsilon, A_\varepsilon) \geq G(v, a),$$

where the limit functional G is given by

$$(3.1) \quad G(v, a) := \begin{cases} \frac{1}{2} \int_U |v - a|^2 + \frac{1}{\ell_0^2} |\text{curl } a - H|^2 + \frac{1}{2} \|\text{curl } v\|_{\mathcal{M}} & \text{if } \ell_0 > 0, \\ \frac{1}{2} \int_U |v - a|^2 + \frac{1}{2} \|\text{curl } v\|_{\mathcal{M}} & \text{if } \ell_0 = 0 \text{ and } \text{curl } a = H, \\ +\infty & \text{otherwise.} \end{cases}$$

Conversely, for every $(v, a) \in L^2(U; \mathbb{R}^2) \times H^1(U; \mathbb{R}^2)$ such that $\text{curl } v$ is a Radon measure there exist approximating sequences $(\tilde{u}_\varepsilon, \tilde{A}_\varepsilon)$ such that the convergences above hold and such that

$$(3.2) \quad \lim_{\varepsilon \rightarrow 0} \frac{1}{|\log \varepsilon|^2} G_\varepsilon(\tilde{u}_\varepsilon, \tilde{A}_\varepsilon) = G(v, a).$$

Proof. It suffices to check the theorem for sequences $(u_\varepsilon, A_\varepsilon)$ that satisfy the Coulomb gauge condition $\text{div } A_\varepsilon = 0$ in U , $A \cdot \nu = 0$ on ∂U , since $G(u_\varepsilon, A_\varepsilon) =$

$G(u_\varepsilon e^{i\chi}, A_\varepsilon + \nabla\chi)$ and the quantities w_ε and $\text{curl} a_\varepsilon$ are invariant under this gauge transformation. The limit functional $G(v, a)$ also has the gauge invariance $G(v + \nabla\chi, a + \nabla\chi) = G(v, a)$.

From the energy bound $G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq K|\log \varepsilon|^2$ we infer that

$$\int_U |\text{curl} a_\varepsilon - H|^2 \leq 2K\ell_\varepsilon^2 \leq C,$$

since ℓ_ε is bounded, and together with $\text{div} a_\varepsilon = 0$ this implies $\|a_\varepsilon\|_{H^1(U)} \leq C$, and via Sobolev embedding $\|A_\varepsilon\|_{L^p(U)} \leq C_p|\log \varepsilon|$ for $p \geq 1$.

We can now establish that the BBH energy $E_\varepsilon(u_\varepsilon)$ is bounded, using the following decomposition:

$$|\nabla_A u|^2 = |\nabla u|^2 - 2j(u) \cdot A + |u|^2|A|^2,$$

which implies that

$$E_\varepsilon(u_\varepsilon) \leq G_\varepsilon(u_\varepsilon, A_\varepsilon) + \int_U |j(u_\varepsilon) \cdot A_\varepsilon|.$$

As in [12], we can estimate the cross term via

$$\begin{aligned} |j(u_\varepsilon) \cdot A_\varepsilon| &\leq \frac{1}{4}|\nabla u_\varepsilon|^2 + |u_\varepsilon|^2|A_\varepsilon|^2 \\ &\leq \frac{1}{4}|\nabla u_\varepsilon|^2 + (|u_\varepsilon|^2 - 1)|A_\varepsilon|^2 + |A_\varepsilon|^2 \\ &\leq \frac{1}{4}|\nabla u_\varepsilon|^2 + \frac{1}{8\varepsilon^2}(1 - |u_\varepsilon|^2)^2 + 2\varepsilon^2|A_\varepsilon|^4 + |A_\varepsilon|^2 \\ &\leq \frac{1}{2}E_\varepsilon(u_\varepsilon) + C\varepsilon^2|\log \varepsilon|^4 + C|\log \varepsilon|^2, \end{aligned}$$

and it follows that $E_\varepsilon(u_\varepsilon) \leq C|\log \varepsilon|^2$. We are therefore able to use the compactness results of [12] that show compactness for v_ε and the estimate

$$\liminf_{\varepsilon \rightarrow 0} E_\varepsilon(u_\varepsilon) \geq \frac{1}{2} \int_U |v|^2 + \frac{1}{2} \|\text{curl} v\|_{\mathcal{M}}.$$

It is then not difficult to show the lower bound for the full energy using the same decomposition as above and the weak convergence of a_ε implied by the bounds.

The Γ -limsup property (3.2) can be shown as follows: Given a limit (v, a) with $\text{div} a = 0$, we set $\tilde{A}_\varepsilon = a|\log \varepsilon|$ and construct \tilde{u}_ε as in [12, section 7]. It is then easy to see that the claimed convergence holds, using the Γ -convergence result for E_ε from [12] and the same decomposition as above. \square

Remark 3.2. Note that compactness for v_ε only holds due to our choice of gauge. The representative $\tilde{u}_\varepsilon = u_\varepsilon e^{i\chi_\varepsilon}$ corresponds to $\tilde{v}_\varepsilon = v_\varepsilon + \frac{1}{|\log \varepsilon|} \nabla\chi_\varepsilon$, and so v_ε and \tilde{v}_ε need not have the same compactness properties. The limit functional $G(v, a)$ also has the gauge invariance $G(v + \nabla\chi, a + \nabla\chi) = G(v, a)$. If $\ell_\varepsilon \rightarrow \infty$, the compactness argument for a_ε fails, since we only know that $\int_U |\text{curl} a_\varepsilon - H|^2 \leq K\ell_\varepsilon^2$, so this sequence need not be bounded. The example given (for $H = 0$) in [15, Theorem 5], which can be used for (1.3), shows that v_ε also need not be compact in this case, even if $\text{div} A_\varepsilon = 0$. In fact we construct a sequence of $(v_\varepsilon, a_\varepsilon)$ with bounded energy but $\|v_\varepsilon\|_{L^2(U)} \gtrsim \log(\ell_\varepsilon \wedge \frac{1}{\varepsilon|\log \varepsilon|^{1/2}}) \rightarrow +\infty$ by constructing a set of vortices that

concentrate about a single point. Therefore, the energy splitting approach of [12] is insufficient to treat the case of large domains.

As in [20], we can characterize the minimizers of the limit functional. We obtain, following the presentation of [12], Theorem 1.3.

Proof of Theorem 1.3. We prove only the part for $\ell_0 = 0$; the first half can be shown by a completely straightforward insertion of ℓ^{-2} into the argument of [12]. Our proof of the second half also follows the structure of their argument.

For a Radon measure $\mu \in H^{-1}$, define the vector field $v^\mu \in L^2(U; \mathbb{R}^2)$ by $\text{curl } v^\mu = \mu$ and $\text{div } v^\mu = 0$. We decompose μ as $\mu^{ac} + \mu^{sing}$ into an absolutely continuous and a singular part. Setting $g(t, \mu) = G(v_0 + tv_\mu, a)$, we calculate

$$0 \leq \lim_{t \rightarrow 0^+} \frac{g(t, \mu) - g(0, \mu)}{t} = \int_U (v_0 - a_0, v^\mu) + \frac{1}{2} \int_U \text{sgn}(\mu_0) d\mu^{ac} + \frac{1}{2} \|\mu^{sing}\| (U).$$

Integrating by parts and using the definition of y_0 , we see that

$$\int_U (v_0 - a_0, v^\mu) = \int_U y_0 \mu,$$

so we obtain

$$(3.3) \quad 0 \leq \int_U \left(y_0 + \frac{1}{2} \text{sgn}(\mu_0) \right) d\mu^{ac} + \int_U \left(y_0 + \frac{1}{2} \text{sgn}(\mu) \right) d\mu^{sing}$$

and similarly by one-sided differentiation in the opposite direction,

$$(3.4) \quad 0 \geq \int_U \left(y_0 + \frac{1}{2} \text{sgn}(\mu_0) \right) d\mu^{ac} + \int_U \left(y_0 - \frac{1}{2} \text{sgn}(\mu) \right) d\mu^{sing}.$$

Together, (3.3) and (3.4) imply, due to the arbitrariness of μ^{ac} and μ^{sing} , that $|y_0| \leq \frac{1}{2}$ everywhere and $y_0 = -\frac{1}{2} \text{sgn}(\mu_0)$ in $\text{spt } \mu_0$. It follows that for any smooth function φ with $\varphi(z) = 0$ for $z \leq 0$ and $\varphi'(z) \geq 0$ there holds

$$\int_U \varphi(y_0) d\mu_0 = -\varphi\left(\frac{1}{2}\right) \mu_0^-(U),$$

where μ_0^- denotes the negative part in the Hahn decomposition of μ_0 . Since $\mu_0 = -\Delta y_0 + H$, we can integrate by parts and obtain

$$\int_U \varphi(y_0) d\mu_0 = \int_U \varphi'(y_0) |\nabla y_0|^2 + \varphi(y_0) H \geq 0,$$

and we conclude that $\mu_0^-(U) = 0$ and so $\mu_0 \geq 0$.

To see that y_0 is a solution of the obstacle problem, we take any $y \in \mathcal{X}$ and compare using $|v|^2 - |w|^2 \geq 2(v - w) \cdot w$ and integration by parts:

$$F_{0,H}(y) - F_{0,H}(y_0) \geq \int_U \nabla(y - y_0) \cdot \nabla y_0 + (y - y_0)H = \int_U (y - y_0) d\mu_0.$$

Now $y_0 = -\frac{1}{2}$ on $\text{spt}(\mu_0)$, so $(y - y_0) \geq 0$ on $\text{spt}(\mu_0)$ for all $y \in \mathcal{X}$. It follows that y_0 is a minimizer of the obstacle problem. Standard theory [9] can now be used to show uniqueness of y_0 . \square

Remark 3.3. Another proof of Theorem 1.3 can be given in the framework of Brezis and Serfaty [7], who examine the obstacle problem arising from (1.7).

COROLLARY 3.4. *Let (v_0, a_0) be a minimizer of $G(v, a)$. Then $\operatorname{curl} v_0 = 0$ for $H < H_1(\ell_0)$ and $\operatorname{curl} v_0 \neq 0$ for $H > H_1(\ell_0)$, where $H_1(\ell_0)$ is given by*

$$(3.5) \quad H_1(\ell_0) = \frac{1}{2 \max_{\overline{U}} |y_{\ell_0}|},$$

where y_{ℓ_0} is the solution of

$$-\Delta y_{\ell_0} + \ell_0^2 y_{\ell_0} + 1 = 0$$

with Dirichlet boundary conditions $y_{\ell_0} = 0$ on ∂U .

Remark 3.5. We reiterate that the function H_1 in (3.5) satisfies $\frac{2H_1(\ell_0)}{\ell_0^2} \rightarrow 1$ as $\ell_0 \rightarrow \infty$.

Remark 3.6. In the case where $U = B_1(0)$ is a ball, the function H_1 can be written down explicitly since the solutions of $-\Delta y + \alpha y + 1 = 0$ with Dirichlet boundary conditions are given by known special functions. Denoting by I_0 the modified Bessel function of zeroth order, we have that

$$H_1(\ell_0) = \frac{\ell_0^2 I_0(\ell_0)}{2(I_0(\ell_0) - 1)}.$$

Since $I_0(x) \sim \frac{e^x}{\sqrt{2\pi x}}$ as $x \rightarrow \infty$ and $I_0(x) = 1 + \frac{x^2}{4} + O(x^4)$ as $x \rightarrow 0$, it is easy to see that this matches the claimed behavior at zero and infinity.

4. Upper bound for vortex lattices. In this section, we construct good comparison sequences that correspond to vortex lattices and calculate their energy.

PROPOSITION 4.1. *Assume $\varepsilon < \frac{C}{\sqrt{h_{ex}}}$ and $\varepsilon \rightarrow 0$. There exists a sequence of functions $(u_\varepsilon, A_\varepsilon)$ such that the Ginzburg–Landau energy satisfies*

$$(4.1) \quad G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq h_{ex} \frac{|U|}{2} \left(\log \frac{1}{\sqrt{h_{ex}} \varepsilon} + C \right).$$

If $h_{ex} - \frac{\ell^2 |\log \varepsilon|}{2} = S \gg 1$ and $h_{ex} \leq \frac{1}{\varepsilon^2}$, then there exists a sequence of functions with Ginzburg–Landau energy

$$(4.2) \quad G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq \frac{|U|}{2} \left(S \left(\left| \log \left(\varepsilon \max(\sqrt{S}, \ell) \right) \right| + C \right) + \frac{\ell^2}{4} |\log \varepsilon|^2 \right).$$

Remark 4.2. The bound given in (4.2) is better than that of (4.1). It is essentially equivalent to the one given for $\ell = 1$ in Proposition 5.8 of Aydi’s thesis [3]. Aydi’s upper bound, adapted to our setting, reads as

$$G_\varepsilon(u_\varepsilon, A_\varepsilon) \leq \frac{|U|}{2} \left(\left(\frac{1}{2} |\log \varepsilon| + S \right)^2 - \frac{1}{2} S^2 \right),$$

and, expanding the square, this is essentially (4.2).

Remark 4.3. If $\ell^2 \geq Kh_{ex}$, then

$$G_\varepsilon(1, 0) \leq h_{ex} \frac{|U|}{2K}.$$

In particular, there is a constant $K > 0$ such that our vortex lattice construction is not minimizing for $\ell^2 \geq Kh_{ex}$.

Remark 4.4. Under the assumptions for the upper bound (4.2), the trivial Meissner-like state $(u, A) = (1, 0)$ has the energy

$$\frac{|U|}{2} \left(\frac{S^2}{\ell^2} + S|\log \varepsilon| + \frac{\ell^2}{4} |\log \varepsilon|^2 \right).$$

Since $|\log(\varepsilon \max(\ell, \sqrt{S}))| \leq |\log \varepsilon| - C$, the vortex lattice state with $O(S)$ vortices is energetically favorable compared to $(u, A) = (1, 0)$. Consult Proposition 2.1 for a more detailed statement regarding the first critical field for vortex nucleation.

We now turn to the proof of Proposition 4.1. We present a new elementary approach using Fourier series to estimate the energy of a vortex lattice. This approach shows the cross-over that happens at $L = \ell^{-1}$; this is related to the decay properties of the Bessel function-type solutions. On the unit cell of our lattice, we investigate solutions of

$$(4.3) \quad -\alpha \Delta h + h = 2\pi \delta_\varepsilon \quad \text{in } K_L$$

with homogeneous Neumann boundary conditions. Here $K_L = \left(-\frac{L}{2}, \frac{L}{2}\right)^2$ for some $L > 0$. This is equivalent to looking at $L\mathbb{Z}^2$ -periodic solutions in \mathbb{R}^2 . For δ_ε we use the Dirac sequence

$$\delta_\varepsilon(x) = \frac{1}{4\varepsilon^2} \chi_{(-\varepsilon, \varepsilon)}(x_1) \chi_{(-\varepsilon, \varepsilon)}(x_2),$$

where χ_A is the characteristic function of a set $A \subset \mathbb{R}$. We assume $2\varepsilon < L$. We obtain the following results on the lattice.

PROPOSITION 4.5. *There exists a $C > 0$ such that for any L, ε with $\varepsilon < \frac{L}{2}$ there exists a periodic function h such that $-\ell^{-2} \Delta h + h = 2\pi \delta_\varepsilon$ and*

$$(4.4) \quad \int_{K_L} \ell^{-4} |\nabla h|^2 + \ell^{-2} |h - h_{ex}|^2 \leq 2\pi \log \frac{1}{\max(\ell, L^{-1})\varepsilon} + C + L^2 \left(h_{ex} - \frac{2\pi}{L^2} \right)^2.$$

Proof. We calculate the energy

$$\int_{K_L} \alpha^2 |\nabla h|^2 + \ell^{-2} |h - h_{ex}|^2.$$

It will become apparent later that we should use $\alpha = \ell^{-2}$.

We use double Fourier series as follows. For $f \in L^2(K_L)$ and $k \in \mathbb{Z}^2$, set

$$a_k = \frac{1}{L^2} \int_{K_L} f(x) e^{-i\gamma k \cdot x},$$

where $\gamma = \frac{2\pi}{L}$. Then f can be reconstructed as

$$f(x) = \sum_{k \in \mathbb{Z}^2} a_k e^{i\gamma k \cdot x}.$$

By Plancherel’s theorem we have

$$\int_{K_L} |f|^2 = L^2 \sum_{k \in \mathbb{Z}^2} |a_k|^2.$$

It is standard that ∇f corresponds to the series $(i\gamma k a_k)$, and Δf to the series $(-\gamma^2 |k|^2 a_k)$. Solving (4.3) therefore corresponds to

$$(\alpha\gamma^2 |k|^2 + 1)a_k = b_k,$$

where b_k are the Fourier coefficients for δ_ε .

We calculate these coefficients. Set $k = (k_1, k_2)$. If $k_1 k_2 \neq 0$, then

$$b_k = \frac{2\pi}{4\varepsilon^2 L^2} \frac{4 \sin(\gamma k_1 \varepsilon) \sin(\gamma k_2 \varepsilon)}{\gamma^2 k_1 k_2}.$$

In other cases we have

$$b_{(m,0)} = b_{(0,m)} = \frac{2\pi}{4\varepsilon^2 L^2} \frac{4\varepsilon \sin(\gamma m \varepsilon)}{\gamma m}$$

and finally $b_0 = \frac{2\pi}{L^2}$ corresponding to $\int_{K_L} h = 2\pi$.

To simplify notation we write this using $\text{sinc}(x)$, the continuous continuation of $\frac{\sin x}{x}$, which yields

$$b_k = \frac{2\pi}{L^2} \text{sinc}(\gamma k_1 \varepsilon) \text{sinc}(\gamma k_2 \varepsilon).$$

Since $a_0 = b_0 = \frac{2\pi}{L^2}$, we have that

$$\int |h - h_{ex}|^2 = L^2 \sum_{k \neq 0} |a_k|^2 + L^2 \left(\frac{2\pi}{L^2} - h_{ex} \right)^2.$$

We want to calculate

$$E = L^2 \sum_{k \in \mathbb{Z}^2 \setminus 0} (\alpha^2 \gamma^2 |k|^2 + \alpha) |a_k|^2 + L^2 \left(\frac{2\pi}{L^2} - h_{ex} \right)^2.$$

Using the expressions obtained for b_k above, it follows that we have

$$E = \frac{4\pi^2}{L^2} \sum_{k \in \mathbb{Z}^2 \setminus 0} \frac{\alpha}{\alpha |\gamma|^2 |k|^2 + 1} \text{sinc}^2(\gamma \varepsilon k_1) \text{sinc}^2(\gamma \varepsilon k_2) + L^2 \left(\frac{2\pi}{L^2} - h_{ex} \right)^2.$$

We split up the double sum as follows. First, consider $k = (k_1, k_2)$ with $1 \leq |k| \leq \frac{1}{\gamma \varepsilon}$. For these terms we estimate $|\text{sinc}| \leq 1$. We label this part of the energy E_1 , and so

$$E_1 \leq \frac{4\pi^2}{L^2} \sum_{K_1} \frac{\alpha}{\alpha \gamma^2 |k|^2 + 1}.$$

Now we compare the sum with an integral. For any decreasing function f , we have

$$\sum_{1 \leq |k| \leq A} f(|k|) \leq \int_{1 - \frac{1}{\sqrt{2}}}^{A + \frac{1}{\sqrt{2}}} f(r) 2\pi r dr$$

and so

$$E_1 \leq \frac{4\pi^2}{L^2} \int_{1-c}^{\frac{1}{\gamma \varepsilon} + c} \frac{2\pi \alpha r}{1 + \alpha \gamma^2 r^2} dr = \frac{4\pi^2}{L^2 \gamma^2} \left(\frac{2\pi}{2} \log(\alpha \gamma^2 x^2 + 1) \right) \Big|_{x=1-c}^{x=\frac{1}{\gamma \varepsilon} + c},$$

where $c = \frac{1}{\sqrt{2}}$. As $\frac{4\pi^2}{L^2\gamma^2} = 1$, we obtain

$$E_1 \leq \frac{2\pi}{2} \log \frac{\frac{\alpha}{\varepsilon^2} + \frac{2c\alpha\gamma}{\varepsilon} + c^2 + 1}{\alpha\gamma^2(1-c)^2 + 1}.$$

We distinguish two cases. If $\alpha\gamma^2 \leq 1$, we estimate the denominator as ≥ 1 and obtain

$$E_1 \leq \frac{2\pi}{2} \log \frac{\alpha}{\varepsilon^2} + C \leq 2\pi \log \frac{\sqrt{\alpha}}{\varepsilon} + C.$$

In the case where $\alpha\gamma^2 \geq 1$, we estimate the denominator as $\geq C\alpha\gamma^2$ and obtain

$$E_1 \leq \frac{2\pi}{2} \log \frac{1}{\gamma^2\varepsilon^2} + C \leq 2\pi \log \frac{1}{\gamma\varepsilon} + C$$

if $\alpha\gamma > \varepsilon$.

We still need to deal with the frequencies k with $|k| \geq \frac{1}{\gamma\varepsilon}$. For this we use that $\text{sinc}^2(x)\text{sinc}^2(y) \leq \frac{2}{r^2}$, which can be seen as follows. Assume without loss of generality that $|x| \leq |y|$. Then $r^2 = x^2 + y^2 \leq 2y^2$. Estimating $|\text{sinc}(x)| \leq 1$ and $|\text{sinc}(y)| \leq \frac{1}{y}$, we see that $\text{sinc}^2(x)\text{sinc}^2(y) \leq \frac{1}{y^2} \leq \frac{2}{r^2}$, as claimed.

To calculate the energy contribution E_2 of those k with $|k| \geq \frac{1}{\gamma\varepsilon}$, we again replace the sum by an integral. Using the sinc bound, we see that

$$E_2 \leq \frac{4\pi^2}{L^2} \int_{\frac{1}{\gamma\varepsilon}-c}^{\infty} \frac{4\pi\alpha r}{(\alpha r^2\gamma^2 + 1)\gamma^2 r^2\varepsilon^2}.$$

We estimate this using $4\pi^2 L^{-2}\gamma^{-2} = 1$ as

$$E_2 \leq C \int_{\frac{1}{\gamma\varepsilon}-c}^{\infty} \frac{1}{\gamma^2\varepsilon^2 r^3} dr \leq \frac{C}{\gamma^2\varepsilon^2} \cdot \frac{1}{(\frac{1}{\gamma\varepsilon} - c)^2},$$

and for $\frac{1}{\gamma\varepsilon} > 2c$ we obtain that $E_2 \leq C$. The claim then follows using the definitions. \square

Choosing $L = \frac{2\pi}{h_{ex}}$ in (4.4) implies the following upper bound.

COROLLARY 4.6. *If $\varepsilon < \frac{1}{2\sqrt{h_{ex}}}$, then there exists a periodic function h with period $L = \sqrt{\frac{2\pi}{h_{ex}}}$ such that $-\ell^{-2}\Delta h + h = \delta_\varepsilon$ and*

$$\int_{K_L} \ell^{-4}|\nabla h|^2 + \ell^{-2}|h - h_{ex}|^2 \leq 2\pi \log \frac{1}{\max(\ell, \sqrt{h_{ex}})\varepsilon} + C.$$

Furthermore, for any $L > 2\varepsilon$ there exists h with $-\ell^{-2}\Delta h + h = \delta_\varepsilon$ and

$$\int_{K_L} \ell^{-4}|\nabla h|^2 + \ell^{-2}|h - h_{ex}|^2 \leq 2\pi \log \frac{1}{\max(\ell, \sqrt{h_{ex}})\varepsilon} + C + L^2\ell^{-2} \left(\frac{2\pi}{L^2} - h_{ex}\right)^2.$$

Remark 4.7. This can be easily extended to $\varepsilon \leq \frac{C}{\sqrt{h_{ex}}}$ for any C that is bounded independently of ε, ℓ , and h_{ex} by choosing $\tilde{\varepsilon} = \frac{2\varepsilon}{C}$ and constructing with $\tilde{\varepsilon}$ instead of ε .

To construct a pair (u, A) from h , we do the following. To define the modulus ρ , we set

$$\rho(r) = \begin{cases} 0, & r < \varepsilon\sqrt{2}, \\ \frac{r-\varepsilon\sqrt{2}}{\varepsilon}, & \varepsilon\sqrt{2} < r < \varepsilon(1 + \sqrt{2}), \\ 1, & r > \varepsilon(1 + \sqrt{2}). \end{cases}$$

We take any A with $\text{curl } A = h$. Outside $B_{\varepsilon\sqrt{2}}$, we define u as $\rho e^{i\varphi}$, where $\nabla\varphi - A = \alpha\nabla^\perp h$. This is possible since for any simple closed curve $\Gamma \subset K_L \setminus B_{\varepsilon\sqrt{2}}$ with $\Gamma = \partial G$ we have

$$(4.5) \quad \int_\Gamma \frac{\partial\varphi}{\partial\tau} = \int_\Gamma A \cdot \tau - \alpha \frac{\partial h}{\partial\nu} = \int_G (-\Delta h + h) = \begin{cases} 2\pi & \text{if } G \supset B_{\varepsilon\sqrt{2}}, \\ 0 & \text{otherwise.} \end{cases}$$

On K_L , we can therefore estimate

$$(4.6) \quad \begin{aligned} \frac{1}{2} \int_{K_L} |(\nabla - iA)u|^2 + \ell^{-2}|h - h_{ex}|^2 + \frac{1}{2\varepsilon^2}(1 - \rho^2)^2 \\ \leq \frac{1}{2} \int_{K_L} \rho^2 \ell^{-4} |\nabla h|^2 + \ell^{-2}|h - h_{ex}|^2 + |\nabla\rho|^2 + \frac{1}{2\varepsilon^2}(1 - \rho^2)^2 \\ \leq C + \pi \log \frac{1}{\max(\ell, \sqrt{h_{ex}})\varepsilon}. \end{aligned}$$

We are now in the position to establish the following argument.

Proof of Proposition 4.1. This will be done in two steps.

Step 1. We use the above construction to build an h in \mathbb{R}^2 and to define a periodic ρ_ε corresponding to the lattice. As the equivalent of (4.5) holds in all of \mathbb{R}^2 , we can define (u, A) in all of \mathbb{R}^2 such that (4.6) holds on every cell of the lattice. All we need to do is choose a proper origin for our lattice: For any $a \in K_L$ we can set $(u^a(z), A^a(z)) = (u(z - a), A(z - a))$, which has energy density $gl^a(z) = gl(z - a)$, where $gl(z) = \frac{1}{2}|(\nabla - iA)u|^2(z) + \ell^{-2}|\text{curl } A(z) - h_{ex}|^2 + \frac{1}{2\varepsilon^2}(1 - |u(z)|^2)^2$. Integrating over the unit cell, we see that

$$\int_{K_L} G_\varepsilon(u_\varepsilon^a, A_\varepsilon^a; U) da = \int_{K_L} \int_U gl^a(z) dz da = |U|G_\varepsilon(u, A; K_L).$$

The mean value theorem shows that there exists some a such that $G(u^a, A^a; U) \leq \frac{|U|}{|K_L|}G(u, A; K_L)$, and since $|K_L| = \frac{h_{ex}}{2\pi}$, this finishes the proof of (4.1).

Step 2. We follow the argument in Step 1; however, we choose a lattice of size $L = \sqrt{\frac{2\pi}{S}}$, which is optimal up to logarithmic terms. Since $h_{ex} \leq \varepsilon^{-2}$ and $S \gg 1$, we have $2\varepsilon \leq L \ll 1$, and we can follow the same construction as above and obtain for the energy after choosing a suitable origin

$$\begin{aligned} G(u_\varepsilon, A_\varepsilon; U) &\leq \frac{|U|}{2|K_L|} \left(2\pi \log \frac{1}{\max(\sqrt{S}, \ell)\varepsilon} + C + L^2 \ell^{-2} \left| h_{ex} - \frac{2\pi}{L^2} \right| \right) \\ &= \frac{|U|}{2} \left(\frac{S}{2\pi} \left(2\pi \log \frac{1}{\max(\sqrt{S}, \ell)\varepsilon} + C \right) + \ell^{-2} \left(\frac{1}{2} \ell^2 |\log \varepsilon| \right)^2 \right) \\ &= \frac{|U|}{2} \left(S \log \frac{1}{\max(\sqrt{S}, \ell)\varepsilon} + C + \frac{1}{4} \ell^2 |\log \varepsilon|^2 \right) \end{aligned}$$

since $h_{ex} - \frac{2\pi}{L^2} = h_{ex} - S = \frac{1}{2} \ell^2 |\log \varepsilon|$. This completes the proof of (4.2). \square

5. Lower bound for vortex lattices. The lower bound counterpart to Proposition 4.1 can be obtained similarly to the calculation in [18] (see also [19]). We can actually get a uniform vortex density estimate.

The heart of the lower bound argument is the following blow-up estimate.

PROPOSITION 5.1. *Assume that $\ell^2|\log \varepsilon| \ll h_{ex} \ll \frac{1}{\varepsilon^2}$, and let $(u_\varepsilon, A_\varepsilon)$ be a minimizer of the Ginzburg–Landau energy (1.3). Then for every $H > 0$ we can find λ with $\ell \ll \lambda \ll \frac{1}{\varepsilon}$ such that for every x such that the ball $B(x, \frac{1}{\lambda}) \subset U$ there holds*

$$(5.1) \quad G_\varepsilon\left(u_\varepsilon, A_\varepsilon; B\left(x, \frac{1}{\lambda}\right)\right) \geq \frac{\alpha_H|B(x, \frac{1}{\lambda})|}{2} h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}(1 - o_\varepsilon(1)),$$

where α_H satisfies $\alpha_H \rightarrow 1$ for $H \rightarrow \infty$.

Proof. We follow the proof of Proposition 8.2 in [18] and rescale by λ . Setting $u(x) = u'(\lambda x)$, $A'(x) = \lambda A(\lambda x)$, $\varepsilon' = \lambda\varepsilon$, and $h'_{ex} = \frac{h_{ex}}{\lambda^2}$, the claim is equivalent to proving for a ball B_1 of radius 1 that

$$(5.2) \quad \begin{aligned} \frac{1}{2} \int_{B_1} |(\nabla - iA')u'|^2 + \frac{\lambda^2}{\ell^2} |\operatorname{curl} A' - h'_{ex}|^2 + \frac{1}{2\varepsilon'^2} (1 - |u'|^2)^2 \\ \geq \frac{\alpha_H|B_1|}{2} h'_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}(1 - o_\varepsilon(1)). \end{aligned}$$

Since $\varepsilon^2\ell^2|\log \varepsilon| \ll \varepsilon^2 h_{ex} \ll 1$ and using the asymptotic behavior of the function $x \mapsto Hx^2 \log \frac{1}{x}$ near 0, we can find λ such that $\varepsilon^2 h_{ex} = H(\varepsilon\lambda)^2 \log \frac{1}{\varepsilon\lambda}$. It also follows that λ satisfies $\ell \ll \lambda \ll \frac{1}{\varepsilon}$.

With this choice of λ , we have $\varepsilon' \rightarrow 0$ and $h'_{ex} = H|\log \varepsilon'|$. We also have $\log \frac{1}{\varepsilon\sqrt{h_{ex}}} = |\log \varepsilon'|(1 - o_\varepsilon(1))$, and so it suffices to prove that (dropping some ε 's)

$$(5.3) \quad \begin{aligned} G'_{\varepsilon'}(u', A'; x) = \frac{1}{2} \int_{B_1(x)} |(\nabla - iA')u'|^2 + \frac{\lambda^2}{\ell^2} |\operatorname{curl} A' - h'_{ex}|^2 + \frac{1}{2\varepsilon'^2} (1 - |u'|^2)^2 \\ \geq \frac{H\alpha_H|B_1|}{2} |\log \varepsilon'|^2 (1 - o_\varepsilon(1)) \end{aligned}$$

for some α_H with $\alpha_H \rightarrow 1$ as $H \rightarrow \infty$.

Depending on the blow-up origin x , we distinguish two cases. Either $G'_{\varepsilon'}(u', A'; x) \gg |\log \varepsilon'|^2$ (in which case (5.3) is true trivially) or $G'_{\varepsilon'}(u', A'; x) \leq K|\log \varepsilon'|^2$. In the latter case we can use the Gamma-convergence result Theorem 3.1 to see that $\frac{1}{|\log \varepsilon'|}(iu'_\varepsilon, \nabla u'_\varepsilon) \rightharpoonup v$ and $\frac{1}{|\log \varepsilon'|}A'_\varepsilon \rightharpoonup a$. As $\lambda \gg \ell$, we are in the case $\ell_0 = 0$.

Letting (v_0, a_0) denote the minimizer of the functional G defined in (3.1), we now use the characterization of the limit in Theorem 1.3. It follows that $\operatorname{curl} a_0 = H$. Using Lemma 5.3 below, $\operatorname{curl} v_0 = H$ on $B_{r_0(H)}$ with $r_0(H) \rightarrow 1$ for $H \rightarrow \infty$. As $\operatorname{curl} v_0 \geq 0$ by Theorem 1.3, it follows that

$$(5.4) \quad \begin{aligned} \liminf_{\varepsilon \rightarrow 0} \frac{1}{|\log \varepsilon'|^2} G'_{\varepsilon'}(u', A'; x) \geq G(v, a) \geq G(v_0, a_0) \geq \frac{1}{2} \|\operatorname{curl} v_0\|_{\mathcal{M}} \\ \geq H \frac{|B_{r_0(H)}|}{2} = H\alpha_H \frac{|B_1|}{2}, \end{aligned}$$

and $\alpha_H \rightarrow 1$ by Lemma 5.3. \square

Combining this with the result of the previous section, we obtain the following claim.

PROPOSITION 5.2. Assume that $\ell^2|\log \varepsilon| \ll h_{ex} \ll \frac{1}{\varepsilon^2}$, and let $(u_\varepsilon, A_\varepsilon)$ be a minimizer of the Ginzburg–Landau energy (1.3). Then the energy density

$$g_\varepsilon = \frac{1}{2} \left(|\nabla_{A_\varepsilon} u_\varepsilon|^2 + \frac{1}{\ell^2} |\operatorname{curl} A_\varepsilon - h_{ex}|^2 + \frac{1}{2\varepsilon^2} (1 - |u_\varepsilon|^2)^2 \right)$$

satisfies

$$\frac{g_\varepsilon}{h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}} \rightharpoonup \mathcal{L}^2$$

in the sense of measures, and the energy satisfies

$$(5.5) \quad G_\varepsilon(u_\varepsilon, A_\varepsilon) = \frac{1}{2} |U| h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}} (1 - o_\varepsilon(1)).$$

Proof. As in [18], this follows by integrating (5.1) over the domain. If $W \subset\subset U$ is a compactly contained subdomain, then we can use Fubini’s theorem and estimate

$$(5.6) \quad \begin{aligned} \int_W G_\varepsilon \left(u_\varepsilon, A_\varepsilon; B \left(x, \frac{1}{\lambda} \right) \cap W \right) dx &= \int_{x \in W} \int_{y \in B(x, \frac{1}{\lambda}) \cap W} g_\varepsilon dy dx \\ &= \int_{y \in B(x, \frac{1}{\lambda}) \cap W} \int_{x \in W} g_\varepsilon dx dy = \int_{y \in W} \left| B \left(y, \frac{1}{\lambda} \right) \cap W \right| g_\varepsilon dy \leq \frac{\pi}{\lambda^2} G_\varepsilon(u_\varepsilon, A_\varepsilon; W). \end{aligned}$$

Using now (5.1) and Fatou’s lemma, we continue to estimate

$$(5.7) \quad \begin{aligned} \liminf_{\varepsilon \rightarrow 0} \frac{G_\varepsilon(u_\varepsilon, A_\varepsilon; W)}{h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}} &\geq \liminf_{\varepsilon \rightarrow 0} \int_{x \in W} \frac{\lambda^2 G_\varepsilon(u_\varepsilon, A_\varepsilon; B(x, \frac{1}{\lambda} \cap W))}{\pi h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}} \\ &\geq \int_{x \in W} \liminf_{\varepsilon \rightarrow 0} \left(\chi_{\{x: B(x, \frac{1}{\lambda}) \subset W\}} \frac{G_\varepsilon(u_\varepsilon, A_\varepsilon; B(x, \frac{1}{\lambda}))}{|B(x, \frac{1}{\lambda})| h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}} \right) \geq \frac{1}{2} |W|. \end{aligned}$$

Using the upper bound of the previous section, we can deduce the existence of a weak limit g of $\frac{g_\varepsilon}{h_{ex} \log \frac{1}{\varepsilon\sqrt{h_{ex}}}}$. Using an approximation result, we see from (5.7) that

$g \geq \frac{1}{2} \mathcal{L}^2$, and by the upper bound, equality follows. \square

LEMMA 5.3. For the domain $U = B_1(0)$, the unique minimizer y_0 of

$$F_{0,H}(y) = \frac{1}{2} \int_{B_1(0)} |\nabla y|^2 + 2yH$$

in the admissible class

$$\mathcal{X} = \left\{ y \in H_0^1(B_1(0)) : y \geq -\frac{1}{2} \text{ a.e. in } B_1(0) \right\}$$

satisfies $y_0 \equiv -\frac{1}{2}$ in $B_{r_0(H)}$ with $r_0(H) \rightarrow 1$ as $H \rightarrow \infty$. Furthermore, the minimizer (v_0, a_0) of the functional G given in (3.1) satisfies $\operatorname{curl} v_0 = H$ on $B_{r_0}(0)$, and $\alpha_H = \frac{\|\operatorname{curl} v_0\|_{\mathcal{M}}}{H|B_1(0)|} \rightarrow 1$ as $H \rightarrow \infty$.

Proof. Using the uniqueness and regularity properties [9], we can calculate the solution of the obstacle problem explicitly: It is given by

$$y_0(r) = \max \left\{ H \left(\frac{r^2}{4} - \frac{1}{4} + \frac{r_0^2}{2} \log r \right), -\frac{1}{2} \right\},$$

where r_0 is chosen such that

$$H \left(\frac{r_0^2}{4} - \frac{1}{4} + \frac{r_0^2}{2} \log r_0 \right) = -\frac{1}{2},$$

which we rewrite as

$$(5.8) \quad 1 - \frac{2}{H} = r_0^2(1 + 2 \log r_0).$$

From (5.8), it is easy to see that $r_0 \rightarrow 1$ as $H \rightarrow \infty$.

In B_{r_0} , $y_0 \equiv -\frac{1}{2}$, so $\Delta y_0 = 0$, and from Theorem 1.3 we infer $\text{curl } v_0 = H$ in B_{r_0} , and $\text{curl } v_0 \geq 0$ elsewhere, which shows the claim. \square

Acknowledgments. The authors wish to thank the referees for many helpful suggestions, including the approach to Proposition 5.2, that significantly improved the paper.

REFERENCES

- [1] A. AFTALION AND E.N. DANCER, *On the symmetry and uniqueness of solutions of the Ginzburg-Landau equations for small domains*, Commun. Contemp. Math., 3 (2001), pp. 1–14.
- [2] A. AFTALION AND Q. DU, *The bifurcation diagrams for the Ginzburg-Landau system of superconductivity*, Phys. D, 163 (2002), pp. 94–105.
- [3] H. AYDI, *Vorticité dans le Modèle de Ginzburg-Landau de la Supraconductivité*, Ph.D. thesis, Department of Mathematics, Paris XII, Paris, France, 2004.
- [4] H. AYDI AND E. SANDIER, *Vortex analysis of the periodic Ginzburg-Landau model*, Ann. Inst. H. Poincaré Anal. Non Linéaire, to appear.
- [5] F. BETHUEL, H. BREZIS, AND F. HÉLEIN, *Ginzburg-Landau Vortices*, Progr. Nonlinear Differential Equations Appl. 13, Birkhäuser Boston, Cambridge, MA, 1994.
- [6] F. BETHUEL AND T. RIVIERE, *Vortices for a variational problem related to superconductivity*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 12 (1995), pp. 243–303.
- [7] H. BREZIS AND S. SERFATY, *A variational formulation for the two-sided obstacle problem with measure data*, Comm. Contemp. Math., 4 (2002), pp. 357–374.
- [8] G. CHAPMAN, Q. DU, M. GUNZBURGER, AND J. PETERSON, *Simplified Ginzburg-Landau models for superconductivity valid for high kappa and high fields*, Adv. Math. Sci. Appl., 5 (1995), pp. 193–218.
- [9] A. FRIEDMAN, *Variational Principles and Free-boundary Problems*, John Wiley & Sons, New York, 1982.
- [10] R.L. JERRARD, *Lower bounds for generalized Ginzburg-Landau functionals*, SIAM J. Math. Anal., 30 (1999), pp. 721–746.
- [11] R.L. JERRARD AND H.M. SONER, *The Jacobian and the Ginzburg-Landau energy*, Calc. Var. Partial Differential Equations, 14 (2002), pp. 151–191.
- [12] R.L. JERRARD AND H.M. SONER, *Limiting behavior of the Ginzburg-Landau functional*, J. Funct. Anal., 192 (2002), pp. 524–561.
- [13] R.L. JERRARD AND D. SPIRN, *Refined Jacobian estimates and the Ginzburg-Landau energy*, Indiana Univ. Math. J., 56 (2007), pp. 135–186.
- [14] M. KURZKE AND D. SPIRN, *Gamma limit of the nonself-dual Chern-Simons-Higgs energy*, J. Funct. Anal., 244 (2008), pp. 535–588.
- [15] M. KURZKE AND D. SPIRN, *Scaling limits of the Chern-Simons-Higgs energy*, Commun. Contemp. Math., 10 (2008), pp. 1–16.
- [16] M. KURZKE AND D. SPIRN, *On the energy of a Chern-Simons-Higgs vortex lattice*, CRM Proc. Lecture Notes, 44 (2008), pp. 127–152.
- [17] E. SANDIER, *Lower bounds for the energy of unit vector fields and applications*, J. Funct. Anal., 152 (1998), pp. 379–403.
- [18] E. SANDIER AND S. SERFATY, *Vortices in the Magnetic Ginzburg-Landau Model*, Progr. Nonlinear Differential Equations Appl. 70, Birkhäuser Boston, Cambridge, MA, 2007.
- [19] E. SANDIER AND S. SERFATY, *On the energy of type-II superconductors in the mixed phase*, Rev. Math. Phys., 12 (2000), pp. 1219–1257.
- [20] E. SANDIER AND S. SERFATY, *A rigorous derivation of a free boundary problem arising in superconductivity*, Ann. Sci. École Norm. Sup. (4), 33 (2000), pp. 561–592.

- [21] E. SANDIER AND S. SERFATY, *Global Minimizers for the Ginzburg-Landau functional below the first critical field*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 17 (2000), pp. 119–145.
- [22] S. SERFATY, *Local minimizers for the Ginzburg-Landau energy near critical magnetic field. I.*, Commun. Contemp. Math., 1 (1999), pp. 213–254.
- [23] S. SERFATY, *Local minimizers for the Ginzburg-Landau energy near critical magnetic field. II.*, Commun. Contemp. Math., 1 (1999), pp. 295–333.
- [24] D. SPIRN, *Vortex motion law for the Schrödinger–Ginzburg–Landau equations*, SIAM J. Math. Anal., 34 (2003), pp. 1435–1476.
- [25] D. SPIRN AND X. YAN, *Minimizers near the first critical field for the non self-dual Chern-Simons-Higgs energy*, Calc. Var. Partial Differential Equations, to appear.
- [26] M. TINKHAM, *Introduction to Superconductivity*, Dover Publications, New York, 2004.
- [27] Y. YANG, *Solitons in Field Theory and Nonlinear Analysis*, Springer Monogr. Math., Springer-Verlag, New York, 2001.

INVERSE SPECTRAL AND SCATTERING THEORY FOR THE HALF-LINE LEFT-DEFINITE STURM–LIOUVILLE PROBLEM*

C. BENNEWITZ[†], B. M. BROWN[‡], AND R. WEIKARD[§]

Abstract. The problem of integrating the Camassa–Holm equation leads to the scattering and inverse scattering problem for the Sturm–Liouville equation $-u'' + \frac{1}{4}u = \lambda wu$, where w is a weight function which may change sign but where the left-hand side gives rise to a positive quadratic form so that one is led to a left-definite spectral problem. In this paper the spectral theory and a generalized Fourier transform associated with the equation $-u'' + \frac{1}{4}u = \lambda wu$ posed on a half-line are investigated. An inverse spectral theorem and an inverse scattering theorem are established. A crucial ingredient of the proofs of these results is a theorem of Paley–Wiener type which is shown to hold true. Additionally, the accumulation properties of eigenvalues are investigated.

Key words. inverse scattering problems, inverse spectral problems, left-definite problems, Sturm–Liouville, Camassa–Holm equation

AMS subject classifications. 37K15, 34A55, 34B24, 34L25, 35Q53

DOI. 10.1137/080724575

1. Introduction. Standard Sturm–Liouville theory deals with the eigenvalue problem

$$(1.1) \quad -(pu')' + qu = \lambda wu,$$

together with appropriate boundary conditions, in the space L_w^2 of functions square integrable with respect to the weight w , *i.e.*, the norm-square of the space is $\|u\|^2 = \int |u|^2 w$. A basic assumption for this to be possible is that $w \geq 0$. In some situations of interest this is not the case, but instead one has $p > 0$, $q \geq 0$. One may then use as a norm-square the integral $\int (p|u'|^2 + q|u|^2)$, and a problem of this type is usually called *left-definite*. A left-definite problem of current interest is the spectral problem associated with the Camassa–Holm equation, which is of the form

$$(1.2) \quad -u'' + \frac{1}{4}u = \lambda wu.$$

The Camassa–Holm equation is an integrable system in a similar sense as the Korteweg–de Vries (KdV) equation. It was first derived as an abstract bi-Hamiltonian system by Fuchssteiner and Fokas [22]. Subsequently, it was shown by Camassa and Holm [11] that it may serve as an integrable model for shallow water waves. In that paper Camassa and Holm also showed that the solitons are peaked and called them peakons (see also Fokas and Liu [21] and Johnson [23]). In contrast to the KdV equation the Camassa–Holm equation may model breaking waves, *i.e.*, smooth

*Received by the editors May 19, 2008; accepted for publication (in revised form) August 19, 2008; published electronically January 23, 2009. This paper was written with partial support from the Mittag-Leffler Institute in Stockholm, Sweden, the Newton Institute in Cambridge, UK, and the National Science Foundation under grant DMS-0304280.

<http://www.siam.org/journals/sima/40-5/72457.html>

[†]Department of Mathematics, Lund University, Box 118, SE-221 00 Lund, Sweden (christer.bennewitz@math.lu.se).

[‡]School of Computer Science, Cardiff University, Cardiff, P.O. Box 916, Cardiff CF2 3XF, UK (Malcolm.Brown@cs.cardiff.ac.uk).

[§]Department of Mathematics, University of Alabama at Birmingham, Birmingham, AL 35226-1170 (rudi@math.uab.edu).

initial data may develop singularities in finite time; cf. Constantin and Escher [15] and Constantin [13] (see also Bressan and Constantin [10] for a way to resolve the singularities due to wave breaking). This, however, happens only when w changes sign and it is this fact which motivates us to consider (1.2) without the assumption that w is positive. The well developed theory of scattering and inverse scattering for the Schrödinger equation is of crucial importance to the theory of the KdV equation. In the same way scattering/inverse scattering theory for (1.2) is important for dealing with the Camassa–Holm equation. Unfortunately, no such theory is available unless $w \geq 0$, and even then current theory requires more smoothness of w than is convenient to assume, in view of the lack of smoothness for the corresponding peakons.

The problem of inverse scattering for (1.2) is considerably more difficult than for the Schrödinger equation, which may be viewed as a rather mild perturbation of the equation $-u'' = \lambda u$. In case of (1.2) the perturbation is of the equation $-u'' + \frac{1}{4}u = \lambda u$, and thus changes the coefficient containing the eigenvalue parameter λ . It appears that the methods used so far for dealing with the Schrödinger equation are no longer applicable.

In this paper we will prove some uniqueness results for inverse spectral theory and inverse scattering for the left-definite case which apply to (1.2) posed on a half-line. One would also like to have results for the full-line, but this appears to be more difficult. One exception is the case of odd initial data for the Camassa–Holm equation on the full-line because the problem can be reduced to one on a half-line. We mention here that the half-line case was also investigated by Boutet de Monvel and Shepelsky [8], [9], who employ Riemann–Hilbert techniques but assume that w is positive. Our approach is via the inverse spectral theory for the left-definite problem, which also is not very well developed. Even the spectral theory for left-definite problems is not widely known (but see for example [1]), in the level of detail necessary for dealing with the inverse problem. We will therefore start by presenting a reasonably comprehensive spectral theory, then prove some uniqueness theorems for the inverse spectral problem, and finally a uniqueness theorem for inverse half-line scattering.

Spectral theory for left-definite Sturm–Liouville problems seems to have been initiated by Weyl [28], who called such problems *polar*. Later many authors have dealt with more or less general left-definite problems. In particular we mention a series of papers by Niessen, Schneider, and their collaborators on singular left-definite so-called S-hermitian systems; see, *e.g.*, [26]. See also [1] and the references cited there. For a more recent contribution, see Kong, Wu, and Zettl [24]. However, papers in inverse spectral theory for left-definite problems are much more scarce; one example is Binding, Browne, and Watson [7].

Because of the connection with the Camassa–Holm equation the inverse scattering problem for (1.2) has attracted some attention. From the physical point of view the full-line case where w decays at infinity and the periodic case are most interesting. The former was treated by Fokas [20] and Constantin and various co-authors, for example in [14], [16], and [17]. The latter was addressed by Constantin and McKean [18], Constantin [12], and Vaninsky [27]. The full-line case with odd initial data reduces to a half-line case, but the half-line case is also of interest independently.

It will be convenient to deal only with the equation

$$(1.3) \quad -u'' + qu = \lambda wu.$$

There is no loss of generality in doing this, since the change of variable $t = \int_0^x 1/p$ will, as is readily seen, turn (1.1) into an equation of this form.

The plan of the paper is as follows. In section 2 we give a general spectral theory for left-definite problems on intervals with at least one regular endpoint, modelled on standard Titchmarsh–Weyl theory. One may extend this to intervals with two singular endpoints, in the same way as one can extend the right-definite theory, but since we will have no use of it here we have abstained from this.

In section 3 we deal with the generalized Fourier transform associated with a left-definite problem. To simplify the discussion we have restricted ourselves to one case, when so-called finite functions are dense in the Hilbert space associated with the equation. There are no fundamental difficulties involved in dealing with the general situation, but again we have no need of it in the applications we are thinking of.

Section 4 discusses uniqueness of the inverse spectral problem. Unfortunately we have neither a characterization nor a reconstruction algorithm, but the fundamental uniqueness theorem is quite general.

In section 5 we prove a theorem of Paley–Wiener type which is crucial for our approach to the inverse spectral theory, and section 6 deals with the uniqueness theorem for the half-line inverse scattering of a left-definite problem. Section 7 is devoted to some results about the number of eigenvalues for a left-definite problem under scattering conditions. Some elementary, but rather lengthy, calculations needed in section 4 have been relegated to the appendix.

2. Spectral theory. We shall consider (1.3) on an interval $[0, b)$ and assume that q and w are real-valued and integrable on compact subsets of $[0, b)$, that $q \geq 0$, and that neither q nor w vanish a.e. Let \mathcal{H}_1 be the set of locally absolutely continuous functions u defined in $[0, b)$ such that $u' \in L^2(0, b)$ and $q|u|^2 \in L^1(0, b)$. As we shall see presently \mathcal{H}_1 is a Hilbert space with scalar product

$$\langle u, v \rangle = \int_0^b (u' \overline{v'} + qu \overline{v})$$

and norm $\|u\| = \sqrt{\langle u, u \rangle}$. In order to show completeness of \mathcal{H}_1 and discuss how to find self-adjoint realizations corresponding to (1.3) we first note the following simple result.

LEMMA 2.1. *For any $a \in [0, b)$ there exists a constant C_a such that*

$$(2.1) \quad |u(x)| \leq C_a \|u\|$$

for any $x \in [0, a]$ and any $u \in \mathcal{H}_1$.

Proof. By the fundamental theorem of calculus and the Cauchy–Schwarz inequality $|u(x)| \leq |u(y)| + |y - x|^{1/2} (\int_0^b |u'|^2)^{1/2}$. If $c \in [a, b)$ is such that $\int_0^c q > 0$, multiplication by $q(y)$ and integrating with respect to y gives

$$|u(x)| \int_0^c q \leq \int_0^c q|u| + c^{1/2} \int_0^c q \left(\int_0^b |u'|^2 \right)^{1/2}.$$

Using Cauchy–Schwarz again we obtain (2.1) with $C_a = (c + 1/\int_0^c q)^{1/2}$. □

PROPOSITION 2.2. *The space \mathcal{H}_1 is complete.*

Proof. By (2.1) a Cauchy sequence u_1, u_2, \dots in \mathcal{H}_1 converges locally uniformly to a continuous function u . Furthermore, $\sqrt{q}u_j$ and u'_j converge in $L^2[0, b)$ to $\sqrt{q}u$ and, say, v , respectively. Now

$$u_j(x) - u_j(0) = \int_0^x u'_j.$$

Letting $j \rightarrow \infty$ we obtain $u(x) = u(0) + \int_0^x v$. Thus u is absolutely continuous with derivative v and u_j converges to u in \mathcal{H}_1 . \square

Denote the set of integrable functions with compact support in $(0, b)$ by L_0 . Then, if $u \in \mathcal{H}_1$ and $v \in L_0$, it follows that $|\int u\bar{v}| \leq C_a \int |v| \|u\|$ if $\text{supp } v \subset [0, a]$, so that the linear form $\mathcal{H}_1 \ni u \mapsto \int u\bar{v}$ is bounded. By Riesz's representation theorem we may therefore find a unique $v^* \in \mathcal{H}_1$ so that $\int u\bar{v} = \langle u, v^* \rangle$. Clearly v^* depends linearly on v , so we obtain a (bounded) operator $G_0 : L_0 \rightarrow \mathcal{H}_1$ such that

$$\langle u, G_0v \rangle = \int_0^b u\bar{v} \text{ for } u \in \mathcal{H}_1, v \in L_0.$$

The operator G_0 is central for the left-definite spectral theory of (1.3).

PROPOSITION 2.3. *The operator G_0 is an integral operator $G_0u(x) = \int u g_0(x, \cdot)$, it is injective, and its restriction to $L_0 \cap \mathcal{H}_1$ is symmetric with range dense in \mathcal{H}_1 .*

Proof. By (2.1) the map $\mathcal{H}_1 \ni u \mapsto u(x)$ is for each fixed $x \in [0, b)$ a bounded linear form, so there exists an element $\overline{g_0(x, \cdot)} \in \mathcal{H}_1$ so that $u(x) = \langle u, \overline{g_0(x, \cdot)} \rangle$ for $u \in \mathcal{H}_1$, and therefore $G_0v(x) = \langle G_0v, \overline{g_0(x, \cdot)} \rangle = \int_0^b v g_0(x, \cdot)$ for any $v \in L_0$. Thus G_0 is an integral operator with kernel $g_0(x, y)$ (actually, as we shall see in Proposition 2.7, g_0 is real-valued). If u and $v \in L_0 \cap \mathcal{H}_1$, then

$$\langle G_0u, v \rangle = \overline{\langle v, G_0u \rangle} = \int_0^b u\bar{v} = \langle u, G_0v \rangle,$$

so the restriction of G_0 to $L_0 \cap \mathcal{H}_1$ is symmetric.

Let $[c, d] \subset (0, b)$ and $u_j(x) = \min(1, j(x-c), j(d-x))$ for $x \in [c, d]$ and $u_j(x) = 0$ otherwise. Then $u_j \in L_0 \cap \mathcal{H}_1$ and tends boundedly to the characteristic function of $[c, d]$ as $j \rightarrow \infty$, so if $G_0v = 0$, it follows from $0 = \langle G_0v, u_j \rangle = \int v\bar{u}_j$ that $\int_c^d v = 0$ for all $[c, d] \subset (0, b)$. Thus $v = 0$ a.e. so that G_0 is injective. On the other hand, if $u \in \mathcal{H}_1$ is orthogonal to G_0v for all $v \in L_0 \cap \mathcal{H}_1$, we may put $v = u_j$, so that $0 = \langle u, G_0u_j \rangle \rightarrow \int_c^d u$. It follows that $u = 0$ so the range of G_0 restricted to $L_0 \cap \mathcal{H}_1$ is dense and the proof is complete. \square

We shall have to briefly use the theory of symmetric relations as presented in [1, section 1], and define maximal and minimal relations corresponding to (1.3). We start by setting

$$T_c = \{(G_0(wv), v) \mid v \in L_0 \cap \mathcal{H}_1\}.$$

Then, since w is real-valued, T_c is a symmetric relation in \mathcal{H}_1 for

$$\langle G_0(wu), v \rangle = \overline{\langle v, G_0(wu) \rangle} = \int_0^b wu\bar{v} = \langle u, G_0(wv) \rangle.$$

Proposition 2.3 implies that T_c is the graph of a densely defined symmetric operator in \mathcal{H}_1 if $\text{supp } w = [0, b)$, but at this point we do not want to exclude the possibility of w vanishing on an open set. We define the *minimal relation* T_0 as the closure (in $\mathcal{H}_1 \oplus \mathcal{H}_1$) of T_c , and the *maximal relation* T_1 as the adjoint of this, i.e.,

$$T_1 = \{(u, f) \in \mathcal{H}_1 \oplus \mathcal{H}_1 \mid \langle u, v \rangle = \langle f, G_0(wv) \rangle \text{ for all } v \in L_0 \cap \mathcal{H}_1\}.$$

We must show that T_1 is a differential relation.

PROPOSITION 2.4. *We have $(u, f) \in T_1$ if and only if u and $f \in \mathcal{H}_1$, u' is locally absolutely continuous, and $-u'' + qu = wf$.*

Proof. First note that if u and $f \in \mathcal{H}_1$, then the definition of G_0 shows that

$$(2.2) \quad \langle u, v \rangle - \langle f, G_0(wv) \rangle = \int_0^b (u' \overline{v'} + qu \overline{v} - wf \overline{v})$$

for any $v \in L_0 \cap \mathcal{H}_1$. If in addition u' is locally absolutely continuous and satisfies $-u'' + qu = wf$, integrating by parts gives

$$\langle u, v \rangle - \langle f, G_0(wv) \rangle = \int_0^b (-u'' + qu - wf) \overline{v} = 0.$$

This proves one direction of the proposition.

In proving the other direction the assumption is that the quantity (2.2) is zero. But since $C_0^\infty(0, b) \subset \mathcal{H}_1$ this means that the distributional derivative of u' is $qu - wf$ so that u' is locally absolutely continuous and u satisfies the differential equation.

To give a proof without the use of distribution theory we prove a variant of the classical du Bois-Reymond lemma. If $v \in L_0 \cap \mathcal{H}_1$, integration by parts in (2.2) gives

$$(2.3) \quad \int_0^b \left\{ u' - \int_0^x (qu - wf) - C \right\} \overline{v'} = 0$$

for any constant C . Now let $[c, d] \subset (0, b)$ and choose $C = \frac{1}{d-c} \int_c^d \{u' - \int_0^x (qu - wf)\}$. Put $v(y) = 0$ for $y \notin [c, d]$ and

$$v(y) = \int_c^y \left\{ u'(x) - \int_0^x (qu - wf) - C \right\} dx$$

for $y \in [c, d]$. Then $v \in L_0 \cap \mathcal{H}_1$ and (2.3) gives

$$\int_c^d \left| u' - \int_0^x (qu - wf) - C \right|^2 = 0$$

so that $u' - \int_0^x (qu - wf)$ is constant in $[c, d]$. Thus u' is locally absolutely continuous, and differentiation gives $-u'' + qu = wf$. \square

Let $\mathcal{D}_\lambda = \{(u, \lambda u) \in T_1\}$ and let D_λ be the projection of \mathcal{D}_λ onto its first components, i.e., $u \in D_\lambda$ means that $u \in \mathcal{H}_1$ and u satisfies $-u'' + qu = \lambda wu$. We then have

$$T_1 = T_0 \dot{+} \mathcal{D}_\lambda \dot{+} \mathcal{D}_{\overline{\lambda}}$$

as a direct sum for any nonreal λ . Here $\dim \mathcal{D}_\lambda = \dim D_\lambda$ is constant in each of the upper and lower half-planes, and these dimensions will be called the *deficiency indices* of T_1 . See [1, Theorem 1.4] for this simple generalization of the von Neumann formula for symmetric operators and its consequences. It is clear that $\dim D_\lambda \leq 2$, and that $\dim D_{\overline{\lambda}} = \dim D_\lambda$, since $\overline{u} \in D_{\overline{\lambda}}$ if and only if $u \in D_\lambda$. Thus deficiency indices are always equal, and there are always self-adjoint extensions of T_0 , which will at the same time be restrictions of T_1 , and therefore realizations of (1.3). It is of course of interest to have criteria in terms of the coefficients q and w for different values of the deficiency indices $\dim D_\lambda$. In surprising contrast to the right-definite case, we have the following simple and explicit criteria.

THEOREM 2.5. *Suppose $\text{Im } \lambda \neq 0$ and let W be an antiderivative of w . Then $\dim D_\lambda = 2$ if $b < \infty$ and $q + W^2 \in L^1[0, b)$. Otherwise $\dim D_\lambda = 1$ for $\text{Im } \lambda \neq 0$.*

The theorem is a special case of [2, Theorem 2.3]. See also [5]. In the right-definite case a simple variation of constants argument shows that if $\dim D_\lambda = 2$ for one real or nonreal value of λ , then this holds for all $\lambda \in \mathbb{C}$. A similar argument shows that this remains true in the left-definite case, with the exception that it is possible that $\dim D_0 = 2$ even if $\dim D_\lambda < 2$ for all $\lambda \neq 0$. This is to be expected, since D_0 does not depend on the choice of w . We characterize $\dim D_0$ completely in the following theorem, which also brings out the significance of the space D_0 . We use the expression *finite function* in \mathcal{H}_1 to denote a function which vanishes near b .

THEOREM 2.6.

- (1) *The set D_0 is the orthogonal complement in \mathcal{H}_1 of $L_0 \cap \mathcal{H}_1$ and has dimension 1 or 2.*
- (2) *$\dim D_0 = 2$ if and only if $b < \infty$ and $q \in L^1[0, b)$.*
- (3) *If $b < \infty$ and $q \in L^1[0, b)$, then v and v' have finite limits at b for all $v \in D_0$, and these limits uniquely determine v .*
- (4) *If $b < \infty$ and $q \in L^1[0, b)$, then every $u \in \mathcal{H}_1$ has a limit at b which is a bounded linear form on \mathcal{H}_1 .*
- (5) *If $\dim D_0 = 1$ and $D_0 \ni v \neq 0$, then $v(0)\overline{v'(0)} < 0$ and $u(x)\overline{v'(x)} \rightarrow 0$ as $x \rightarrow b$ for any $u \in \mathcal{H}_1$.*
- (6) *Finite functions are dense in \mathcal{H}_1 if and only if $\dim D_0 = 1$.*

Most of this is also a special case of the results of [2] and [5], but we give a simple proof, an elaboration of which can also prove Theorem 2.5.

Proof. We have $u \in D_0$ precisely if $\langle u, 0 \rangle \in T_1$, which holds precisely if $\langle u, v \rangle = \langle u, v \rangle - \langle 0, G_0(wv) \rangle = 0$ for all $v \in L_0 \cap \mathcal{H}_1$, proving the first claim. Since there are elements $v \in \mathcal{H}_1$ with $v(0) \neq 0$, and since $u(0) = 0$ for every $u \in L_0 \cap \mathcal{H}_1$, it follows from (2.1) for $x = 0$ that $\dim D_0 \geq 1$ and we have proved (1).

If b is finite and q integrable, standard existence and uniqueness theorems show that all solutions of $-v'' + qv = 0$ are continuously differentiable with absolutely continuous derivative in $[0, b]$, and thus in \mathcal{H}_1 , and that they are uniquely determined by the values of v and v' at b . In this case the proof of Lemma 2.1 clearly also works for $a = b$, so we have proved (3), (4), and one direction of (2).

Now let $u \in \mathcal{H}_1$ and $v \in D_0$. Integration by parts gives

$$(2.4) \quad \int_0^x (u'\overline{v'} + qu\overline{v}) + u(0)\overline{v'(0)} = u(x)\overline{v'(x)}.$$

Thus $u(x)\overline{v'(x)}$ has a limit at b . If this is not 0, then $(u(x)\overline{v'(x)})^{-1}$ is bounded close to b . Therefore $u'/u = u'\overline{v'}/(u\overline{v'})$ is integrable near b , so that u has a nonzero limit at b . Since $q|u|^2$ is integrable it follows that $q \in L^1(0, b)$. Similarly, $v''/v' = qv/v' = qv\overline{u}/(v'\overline{u})$ is integrable near b , so v' has a nonzero limit at b . Since $|v'|^2$ is integrable it follows that b is finite.

Now, setting $u = v \neq 0$ in (2.4) the integral is increasing, ≥ 0 , and not constant, so if $v(0)\overline{v'(0)} \geq 0$, then $v(x)\overline{v'(x)}$ cannot tend to 0 at b . However, if $\dim D_0 = 2$, we may choose $v \in D_0$ with $v'(0) = 0$, so it follows that $q \in L^1(0, b)$ and b finite, completing the proof of (2).

On the other hand, if $\dim D_0 = 1$, then $u(x)\overline{v'(x)}$ must tend to zero for any $u \in \mathcal{H}_1$. In particular, for $u = v$ one therefore has $v(0)\overline{v'(0)} < 0$ for any nonzero $v \in D_0$ which proves (5).

Finally, if $u \in \mathcal{H}_1$ is finite and $v \in D_0$, integration by parts shows that $\langle u, v \rangle = -u(0)\overline{v'(0)}$, so the orthogonal complement of the finite functions consists of those

$v \in D_0$ for which $v'(0) = 0$. According to (5) this implies $v = 0$ if $\dim D_0 = 1$ and the proof is complete. \square

It is now possible to give a detailed description of the kernel g_0 .

PROPOSITION 2.7. *The kernel $g_0(x, y)$ is real-valued and symmetric in x, y . As a function of y it satisfies (1.3) with $\lambda = 0$ for $y \neq x$, and there are real-valued functions ψ_0 and φ_0 which solve (1.3) with $\lambda = 0$, such that if $u \in \mathcal{H}_1$, then*

- (1) $\psi_0 \in \mathcal{H}_1$, $\psi'_0(0) = 1$ and $\psi'_0(x)u(x) \rightarrow 0$ as $x \rightarrow b$,
- (2) $\varphi_0(0) = -1$, $\varphi'_0(0) = 0$,
- (3) $g_0(x, y) = \varphi_0(\min(x, y))\psi_0(\max(x, y))$.

Proof. The existence of the solution φ_0 is not in question, and if a solution with the properties of ψ_0 exists, it is easy to verify that the kernel $\varphi_0(\min(x, y))\psi_0(\max(x, y))$ has the properties required of $g_0(x, y)$.

The existence of ψ_0 follows from Theorem 2.6. Indeed, if $\dim D_0 = 2$, the element $v \in D_0$ with $v(b) = 1$, $v'(b) = 0$ is real-valued and must have $v(0)v'(0) < 0$ by (2.4), so $v'(0) \neq 0$, and an appropriate multiple will have the properties required of ψ_0 .

On the other hand, if $\dim D_0 = 1$, any nonzero $v \in D_0$ satisfies $v(0)v'(0) < 0$ so $v'(0) \neq 0$, and an appropriate multiple will satisfy the requirements for ψ_0 . Note that this solution is real-valued, since its real and imaginary parts also are in D_0 , and are thus proportional, and the initial condition guarantees that the imaginary part vanishes. \square

Now let T be a self-adjoint restriction of T_1 and assume that (u, f) and $(v, g) \in T$. Integrating by parts we then obtain

$$(2.5) \quad \int_0^x (u'\bar{g} + qu\bar{g}) - \int_0^x (f'\bar{v}' + qf\bar{v}) = (u'\bar{g} - f\bar{v}')|_0^x.$$

As $x \rightarrow b$ this vanishes, since the left-hand side tends to $\langle u, g \rangle - \langle f, v \rangle$. Thus the condition for symmetry is that

$$(u'\bar{g} - f\bar{v}')|_0^b = 0.$$

Comparing this with $(u'\bar{v} - u\bar{v}')|_0^b = 0$, which is the similar condition in the right-definite case, we see that only exceptionally would self-adjoint boundary conditions in the left-definite case also be self-adjoint boundary conditions in the right-definite case.

Separated boundary conditions are those that make $u'\bar{g} - f\bar{v}'$ vanish at each endpoint separately, and are thus at 0 of the form

$$(2.6) \quad f(0) \cos \alpha + u'(0) \sin \alpha = 0,$$

for some $\alpha \in [0, \pi)$. Again comparing with the right-definite case, where the condition is $u(0) \cos \alpha + u'(0) \sin \alpha = 0$, the conditions coincide only in the case $\alpha = \pi/2$, the Neumann boundary condition. However, for eigenfunctions, where $f = \lambda u$, it is clear that also $\alpha = 0$, the Dirichlet boundary condition, gives the same spectra outside of $\lambda = 0$.

We shall not need a detailed description of self-adjoint boundary conditions at a singular endpoint. However, one may always impose the condition (2.6) at 0. It is easy to see that the corresponding restriction of T_1 has a symmetric adjoint, which is a strict extension of T_0 . If the deficiency indices of T_0 equal 1, this is sufficient to obtain a self-adjoint restriction T of T_1 , and all self-adjoint realizations are of this form. Otherwise, a condition needs to be imposed also at b . From (2.5) it follows

immediately that every $(u, f) \in T_1$ satisfying such a condition at b must satisfy $\text{Im}(u'(x)\overline{f(x)}) \rightarrow 0$ as $x \rightarrow b$.

Assuming now that we have a self-adjoint relation T , the spectral theorem looks as follows (see [1, Theorem 1.15]). Consider the set $\mathcal{H}_\infty = \{u \in \mathcal{H}_1 \mid (0, u) \in T\}$. Then \mathcal{H}_∞ is a subspace of \mathcal{H}_1 , and setting $\mathcal{H} = \mathcal{H}_1 \ominus \mathcal{H}_\infty$ the domain D_T of T (i.e., the set of first components of T) is a dense subset of \mathcal{H} , and $T \cap \mathcal{H} \oplus \mathcal{H}$ is the graph of a self-adjoint operator in \mathcal{H} . We will denote this operator by T as well, and may now apply the usual spectral theorem to T . If the resolution of the identity for the operator T is $\{E_t\}_{t \in \mathbb{R}}$, we extend the domain of the projection E_t to all of \mathcal{H}_1 by setting $E_t \mathcal{H}_\infty = 0$. Clearly one may view \mathcal{H}_∞ as an eigenspace for the relation T belonging to the eigenvalue ∞ , so adjoining the orthogonal projection onto \mathcal{H}_∞ to $\{E_t\}_{t \in \mathbb{R}}$ gives a resolution of the identity in \mathcal{H}_1 for the relation T . In the present case one may give a rather complete description of \mathcal{H}_∞ .

PROPOSITION 2.8. *The space \mathcal{H}_∞ consists of those elements $g \in \mathcal{H}_1$ for which $wg = 0$ a.e., and for which $(0, g)$ satisfies the boundary conditions that define T . In particular, if $wg = 0$ a.e. and $g \in L_0 \cap \mathcal{H}_1$, then $g \in \mathcal{H}_\infty$.*

Proof. Now $g \in \mathcal{H}_\infty$ means that $(0, g) \in T$, which therefore satisfies the boundary conditions defining T . In particular, $0 = \langle g, G_0(wf) \rangle - \langle 0, f \rangle = \langle g, G_0(wf) \rangle = \int g \overline{f} w$ for any $f \in L_0 \cap \mathcal{H}_1$. It follows, as in the proof of Proposition 2.3, that $wg = 0$ a.e.

Conversely, if $(0, g)$ satisfies the boundary conditions and $gw = 0$ a.e., then if $(u, f) \in T$, an integration by parts gives

$$\langle u, g \rangle - \langle f, 0 \rangle = \lim_{x \rightarrow b} (u' \overline{g} - f \cdot 0)|_0^x = 0,$$

i.e., $(0, g) \in T$, so the proof is complete. □

We remark that if an endpoint is regular, then the boundary condition implied by $u \in \mathcal{H}_\infty$ is in most cases the vanishing of u in that endpoint. For separated boundary conditions an exception occurs when the boundary condition is of Neumann type (i.e., when $\alpha = \pi/2$ in (2.6)). If we have Neumann conditions at both ends, or at one end when deficiency indices equal 1, there are *no* boundary conditions for elements of \mathcal{H}_∞ .

We will base our derivation of the expansion theorem for the operator T on a detailed description of the resolvent $R_\lambda = (T - \lambda)^{-1}$. Thus R_λ is defined on \mathcal{H} , but we extend its domain to \mathcal{H}_1 by setting $R_\lambda \mathcal{H}_\infty = 0$. The range of R_λ is of course D_T , which is a dense set in \mathcal{H} . Using the kernel g_0 for the evaluation operator on \mathcal{H}_1 introduced in the proof of Proposition 2.3, we have $R_\lambda u(x) = \langle R_\lambda u, \overline{g_0(x, \cdot)} \rangle = \langle u, \overline{R_\lambda^{-1} g_0(x, \cdot)} \rangle$, since the adjoint of R_λ is R_λ^{-1} . Thus we may view $G(x, \cdot, \lambda) = \overline{R_\lambda^{-1} g_0(x, \cdot)}$ as Green's function for our operator; note, however, that G is not the kernel of a standard integral operator. It will turn out to be convenient to introduce the kernel $g(x, y, \lambda) = G(x, y, \lambda) + g_0(x, y)/\lambda$, so that we obtain

$$(2.7) \quad R_\lambda u(x) = \langle u, \overline{g(x, \cdot, \lambda)} \rangle - u(x)/\lambda.$$

Note that $G(x, \cdot, \lambda) \in \mathcal{H}$ but this is not true of $g(x, \cdot, \lambda)$ unless $\mathcal{H}_\infty = \{0\}$. We shall need a precise description of $g(x, y, \lambda)$. To do this we must introduce solutions of (1.3) satisfying initial conditions at 0, so let $\varphi(x, \lambda), \theta(x, \lambda)$ be solutions of (1.3) for $\lambda \neq 0$ satisfying

$$(2.8) \quad \begin{cases} \lambda \varphi(0, \lambda) = -\sin \alpha \\ \varphi'(0, \lambda) = \cos \alpha \end{cases}, \quad \begin{cases} \lambda \theta(0, \lambda) = \cos \alpha \\ \theta'(0, \lambda) = \sin \alpha \end{cases}.$$

This means that φ satisfies the boundary condition (2.6) and θ another similar boundary condition at 0. We have the following theorem.

THEOREM 2.9. *Suppose T is a self-adjoint realization of (1.3) given by (2.6) and, if needed, an appropriate condition at b . Then there exists a function $m(\lambda)$ defined for $\text{Im } \lambda \neq 0$, the Titchmarsh–Weyl m -function for T , depending only on λ and such that $\psi(x, \lambda) = \theta(x, \lambda) + m(\lambda)\varphi(x, \lambda)$, called the Weyl solution for T , is in \mathcal{H}_1 and satisfies the boundary condition at b , if any. Furthermore*

$$g(x, y, \lambda) = \varphi(\min(x, y), \lambda)\psi(\max(x, y), \lambda).$$

Proof. For nonreal λ neither φ nor θ can be in \mathcal{H}_1 and satisfy the boundary condition at b , since that would make λ a nonreal eigenvalue for a self-adjoint problem. Thus there is a solution $\psi(x, \lambda) = \theta(x, \lambda) + m(\lambda)\varphi(x, \lambda)$ in \mathcal{H}_1 which also satisfies the boundary condition at b , since if $\dim D(\lambda) = 2$, one linear, homogeneous condition still leaves a one-dimensional space, whereas if $\dim D(\lambda) = 1$, no boundary condition is imposed at b .

Define, for fixed x and $\lambda \notin \mathbb{R}$, the function

$$F(y) = \varphi(\min(x, y), \lambda)\psi(\max(x, y), \lambda) - \lambda^{-1}g_0(x, y).$$

Since $\psi(\cdot, \lambda)$ and ψ_0 are in \mathcal{H}_1 so is F . We claim that $F \in D_T$. In fact, one easily checks that F' is locally absolutely continuous and that F satisfies $-F'' + qF = \lambda wF + wg_0(x, \cdot)$. It is also easy to check that F satisfies the boundary condition (2.6).

Finally, for $y > x$ the function F is a linear combination of $\psi(\cdot, \lambda)$ and ψ_0 . The former satisfies the boundary condition at b by construction, and ψ_0 satisfies the boundary condition at b by Theorem 2.6(5), since if $(u, f) \in T$, then $\psi'_0 \bar{f} - 0\bar{u}' = \psi'_0 \bar{f} \rightarrow 0$ at b . All this means that $F = R_\lambda g_0(x, \cdot) = \overline{R_{\bar{\lambda}} g_0(x, \cdot)} = G(x, \cdot, \lambda)$ so that $g(x, y, \lambda)$ is as claimed. \square

THEOREM 2.10. *The function m is analytic outside \mathbb{R} , it maps the upper half plane into itself, and it satisfies $\overline{m(\lambda)} = m(\bar{\lambda})$.*

Proof. Since R_λ is analytic outside \mathbb{R} in the strong operator topology $R_\lambda u(x)$ is, by (2.1), pointwise analytic. It follows that $g(x, \cdot, \lambda)$ is weakly analytic for each x , and thus, again by (2.1), $g(x, y, \lambda)$ is analytic outside \mathbb{R} for each x and y . Since $\varphi(x, \lambda)$ and $\theta(x, \lambda)$ also are analytic and since an integration by parts shows that they are nonzero for $x > 0$ and $\lambda \notin \mathbb{R}$, it follows that $m(\lambda)$ is analytic in $\mathbb{C} \setminus \mathbb{R}$.

If (v, g) defines a boundary condition at b , then so does either its real part or its imaginary part, which is easily seen. Therefore, since $\psi(x, \lambda)$ satisfies (1.3) and the boundary condition at b , so does $\overline{\psi(x, \bar{\lambda})}$, and is thus a multiple of $\psi(x, \lambda)$. Now $\overline{\varphi(x, \bar{\lambda})} = \varphi(x, \lambda)$, $\overline{\theta(x, \bar{\lambda})} = \theta(x, \lambda)$ and $\psi(x, \lambda) = \theta(x, \lambda) + m(\lambda)\varphi(x, \lambda)$ so it follows that $\overline{m(\bar{\lambda})} = m(\lambda)$.

Integrating by parts we have

$$\text{Im } \lambda \int_0^x (|\psi'(\cdot, \lambda)|^2 + q|\psi(\cdot, \lambda)|^2) = \text{Im}(\overline{\psi'(\cdot, \lambda)}\lambda\psi(\cdot, \lambda))\Big|_0^x.$$

Since ψ satisfies a boundary condition at b , the integrated term vanishes as $x \rightarrow b$. At 0 the integrated term evaluates to $-\text{Im } m(\lambda)$, so we obtain

$$(2.9) \quad \|\psi(\cdot, \lambda)\|^2 = \text{Im } m(\lambda) / \text{Im } \lambda.$$

Thus m maps the upper and lower half-planes into themselves. \square

A function with the properties of m is a so-called Nevanlinna or Herglotz function, and has a unique representation

$$(2.10) \quad m(\lambda) = A + B\lambda + \int_{\mathbb{R}} \left(\frac{1}{t - \lambda} - \frac{t}{t^2 + 1} \right) d\rho,$$

where $A \in \mathbb{R}$, $B \geq 0$, and $d\rho$ is a positive measure with $\int_{\mathbb{R}} \frac{d\rho(t)}{1+t^2} < \infty$. We will call the measure $d\rho$ the *spectral measure* for T , for reasons that will become clear presently.

We finally note the following proposition.

PROPOSITION 2.11. *Unless $\alpha = \pi/2$ and $0 \notin \text{supp } w$ the functions ψ_0 and $\psi(\cdot, \lambda)$ are in \mathcal{H} .*

Proof. Suppose $g \in \mathcal{H}_\infty$. An integration by parts then gives

$$\langle g, \psi \rangle = -g(0)\overline{\psi'(0)},$$

where $\psi = \psi_0$ or $\psi(\cdot, \lambda)$. The boundary condition at 0 requires $g(0) = 0$ unless $\alpha = \pi/2$, and even then $g(0) = 0$ unless $w = 0$ in a neighbourhood of 0. \square

3. The Fourier transform. We shall call functions that vanish near b *finite* and from now on make the following simplifying assumption.

Assumption 3.1. Assume that finite functions are dense in \mathcal{H}_1 .

According to Theorem 2.6 this means exactly that either $q \notin L^1(0, b)$ or else $b = \infty$. Note that, according to Theorem 2.5, the assumption implies that the deficiency indices of T_1 equal 1.

The spectral measure introduced in the previous section gives rise to a Hilbert space L^2_ρ with scalar product $\langle \hat{u}, \hat{v} \rangle_\rho = \int_{-\infty}^\infty \hat{u}\bar{\hat{v}} d\rho$. We shall define a generalized Fourier transform $\mathcal{F} : \mathcal{H}_1 \rightarrow L^2_\rho$ with the following properties.

THEOREM 3.2.

- (1) *The map $u \mapsto \int_0^b (u'\varphi'(\cdot, t) + qu\varphi(\cdot, t))$, defined for finite $u \in \mathcal{H}_1$, extends by continuity to a map $\mathcal{F} : \mathcal{H}_1 \rightarrow L^2_\rho$ called the generalized Fourier transform. The image of $u \in \mathcal{H}_1$ is denoted by $\mathcal{F}(u)$ or \hat{u} . We write this as $\hat{u}(t) = \langle u, \varphi(\cdot, t) \rangle$ although the integral in general does not converge pointwise.*
- (2) *The mapping $\mathcal{F} : \mathcal{H}_1 \rightarrow L^2_\rho$ has kernel \mathcal{H}_∞ and is unitary between \mathcal{H} and L^2_ρ so that Parseval's formula $\langle u, v \rangle = \langle \hat{u}, \hat{v} \rangle_\rho$ holds if at least one of u and v is in \mathcal{H} .*
- (3) *If $u \in D_T$, then $\mathcal{F}(Tu)(t) = t\hat{u}(t)$. Conversely, if \hat{u} and $t\hat{u}(t)$ are in L^2_ρ , then $\mathcal{F}^{-1}(\hat{u}) \in D_T$.*
- (4) *Suppose $\alpha \neq 0$ in (2.6). Then $\varphi(x, \cdot) \in L^2_\rho$ for each x and $\int_{-\infty}^\infty \hat{u}\varphi(x, \cdot) d\rho = \langle \hat{u}, \varphi(x, \cdot) \rangle_\rho$ converges in \mathcal{H} , and hence locally uniformly in x , for $\hat{u} \in L^2_\rho$. This is the adjoint of $\mathcal{F} : \mathcal{H}_1 \rightarrow L^2_\rho$ and thus the inverse of \mathcal{F} restricted to \mathcal{H} . If M is a Borel set in \mathbb{R} , then*

$$(3.1) \quad E_M u(x) = \int_M \hat{u}\varphi(x, \cdot) d\rho.$$

If $\alpha = 0$, the same is true, except that we must replace $\varphi(\cdot, t)$ for $t = 0$ by the function ψ_0 of Proposition 2.7. Note that ψ_0 is the eigenfunction for the eigenvalue 0 in this case.

We first consider the Fourier transform for finite functions $u \in \mathcal{H}_1$, for every $\lambda \in \mathbb{C}$ setting

$$\hat{u}(\lambda) = \langle u, \varphi(\cdot, \bar{\lambda}) \rangle.$$

It is clear that \hat{u} is an entire function, since integration by parts shows that

$$\hat{u}(\lambda) = \langle u, \varphi(\cdot, \bar{\lambda}) \rangle = \int_0^b u \lambda \varphi(\cdot, \lambda) w - u(0) \cos \alpha,$$

and by (2.8) $\lambda \varphi(x, \lambda)$ is an entire function of λ , locally uniformly in x .

LEMMA 3.3. *For finite u and $v \in \mathcal{H}_1$ we have \hat{u} and $\hat{v} \in L^2_\rho$. If E_Δ is the spectral projection for T associated with an interval Δ , then $\langle E_\Delta u, v \rangle = \int_\Delta \hat{u} \bar{\hat{v}} d\rho$.*

Proof. We have $\langle R_\lambda u, v \rangle = \hat{u}(\lambda) \hat{v}(\bar{\lambda}) m(\lambda) + g(\lambda)$, where g is entire, as is easily verified by direct calculation. Integrating around a rectangle γ with corners at $c \pm i$ and $d \pm i$ we therefore have $\int_\gamma \langle R_\lambda u, v \rangle d\lambda = \int_\gamma \hat{u}(\lambda) \hat{v}(\bar{\lambda}) m(\lambda) d\lambda$ whenever one of the integrals exists. By the spectral theorem the first integral equals $\int_\gamma \int_\mathbb{R} \frac{d\langle E_t u, v \rangle}{t - \lambda} d\lambda$, so if the integral is absolutely convergent, changing the order of integration gives $-2\pi i \langle E_{(c,d)} u, v \rangle$ if c and d are points of continuity for $\langle E_t u, v \rangle$.

Similarly, using the Nevanlinna representation (2.10), the other integral equals $-2\pi i \int_c^d \hat{u}(t) \bar{\hat{v}}(t) d\rho(t)$ if it is absolutely convergent and c, d are points of continuity for ρ .

The absolute convergence of the double integrals is ensured if $\langle E_t u, v \rangle$ and ρ are differentiable at c and d as is easily seen. For more details of the identical calculation carried out for the right-definite case, see [6, Lemmas 14.3, 14.4].

As functions of bounded variation $\langle E_t u, v \rangle$ and ρ are both differentiable a.e., so the second claim of the lemma is true if the endpoints of Δ belong to this dense set of points, and so in general by continuity. In particular, letting $c \rightarrow -\infty, d \rightarrow \infty$ through such points it follows that $\langle E_\mathbb{R} u, u \rangle = \langle \hat{u}, \hat{u} \rangle_\rho$, so that $\hat{u}, \hat{v} \in L^2_\rho$. \square

Since finite functions are dense in \mathcal{H}_1 , and since $E_\mathbb{R}$ has kernel \mathcal{H}_∞ , we now obtain Theorem 3.2(1) by continuity and also (2) except for the surjectivity of \mathcal{F} . To prove this we need the following lemmas.

LEMMA 3.4. *The transform of $R_\lambda u$ is $\hat{u}(t)/(t - \lambda)$.*

Proof. According to the spectral theorem we have $\langle R_\lambda u, v \rangle = \int_\mathbb{R} \frac{d\langle E_t u, v \rangle}{t - \lambda}$ and by Lemma 3.3 we have $\langle E_t u, v \rangle = \int_{-\infty}^t \hat{u} \bar{\hat{v}} d\rho$ so that

$$\langle R_\lambda u, v \rangle = \int_\mathbb{R} \frac{\hat{u}(t)}{t - \lambda} \bar{\hat{v}}(t) d\rho(t).$$

We also have $R_\lambda - R_{\bar{\lambda}} = (\lambda - \bar{\lambda}) R_{\bar{\lambda}} R_\lambda$ and $\langle R_\lambda u, R_\lambda u \rangle = \langle R_{\bar{\lambda}} R_\lambda u, u \rangle$ so

$$\langle R_\lambda u, R_\lambda u \rangle = \frac{1}{\lambda - \bar{\lambda}} (\langle R_\lambda u, u \rangle - \langle R_{\bar{\lambda}} u, u \rangle) = \left\| \frac{\hat{u}(t)}{t - \lambda} \right\|_\rho^2.$$

Expanding $\left\| \frac{\hat{u}(t)}{t - \lambda} - \mathcal{F}(R_\lambda u) \right\|_\rho^2$ and using Parseval's formula and the above yields 0, thus proving the lemma. \square

LEMMA 3.5. *The operator T has eigenvalue 0 if and only if $\alpha = 0$, in which case the eigenfunction is ψ_0 , and for any $u \in \mathcal{H}_1$ we then have $\hat{u}(0) = -u(0)$.*

Furthermore, the measure $d\rho$ has mass at 0 ($\{0\}$ is not a nullset with respect to $d\rho$) precisely if $\alpha = 0$. In this case $\hat{\psi}_0 = \chi_{\{0\}}/\rho\{0\}$, where $\chi_{\{0\}}$ is the characteristic function of the singleton $\{0\}$ and $\rho\{0\}$ the spectral measure of this set.

Proof. According to Theorem 2.6 the only nontrivial solutions of (1.3) for $\lambda = 0$ in \mathcal{H}_1 are multiples of a solution u for which $u'(0)u(0) < 0$, so that $u'(0) \neq 0$. These solutions satisfy the boundary condition (2.6) precisely if $\alpha = 0$, which proves the first

claim. If u is any finite function, integrating by parts gives $\hat{u}(0) = \langle u, \varphi(\cdot, 0) \rangle = -u(0)$. This holds in general by continuity, $u(0)$ being a bounded linear form on \mathcal{H}_1 by (2.1), and $\hat{u}(0)$ on L^2_ρ since $d\rho$ has mass at 0, as we shall see presently.

Now $u \in D_T$ and $Tu = 0$ precisely if $u + \lambda R_\lambda u = 0$, and the Fourier transform of $u + \lambda R_\lambda u$ is $(1 + \frac{\lambda}{t-\lambda})\hat{u}(t) = \frac{t\hat{u}(t)}{t-\lambda}$. If this is 0, then $\hat{u} = 0$ a.e. with respect to $d\rho$ except possibly at $t = 0$. Thus, if $\alpha = 0$, then $\{0\}$ cannot be a nullset with respect to $d\rho$. It also follows that $\hat{\psi}_0$ is a multiple of the characteristic function of the set $\{0\}$. On the other hand, since $\dim D_\lambda = 1$, Weyl solutions for different α are proportional so it immediately follows that

$$(3.2) \quad m_0(\lambda) = \frac{\psi'(0, \lambda)}{\lambda\psi(0, \lambda)} = \frac{\sin \alpha + m_\alpha(\lambda) \cos \alpha}{\cos \alpha - m_\alpha(\lambda) \sin \alpha},$$

where m_α denotes the m -function associated with the boundary condition parameter α . Now $m_0(i\nu) \rightarrow \infty$ as $\nu \downarrow 0$, as a consequence of the mass at 0, so that $m_\alpha(i\nu) \rightarrow \cot \alpha$ for $\alpha \neq 0$. For $\alpha \neq 0$ the spectral measure therefore has no mass at 0.

It only remains to prove the formula for $\hat{\psi}_0$. By Parseval's formula (note that $\psi_0 \in \mathcal{H}$ by Proposition 2.11) we have $\hat{\psi}_0(0) = -\psi_0(0) = \|\psi_0\|^2 = \|\hat{\psi}_0\|^2_\rho = |\hat{\psi}_0(0)|^2 \rho\{0\}$. Hence $-\psi_0(0) = \hat{\psi}_0(0) = 1/\rho\{0\}$. \square

It is now easy to prove that \mathcal{F} is surjective.

LEMMA 3.6. *The Fourier transform $\mathcal{H} \rightarrow L^2_\rho$ is surjective.*

Proof. Suppose that $\hat{u} \in L^2_\rho$ is orthogonal to all Fourier transforms \hat{v} . Since $\hat{v}(t)/(t - \lambda)$ is also a transform, for any nonreal λ , we have $\int \frac{1}{t-\lambda} \hat{u}(t) \overline{\hat{v}(t)} d\rho(t) = 0$ for all nonreal λ . Thus the Stieltjes transform of the measure $\hat{u} \overline{\hat{v}} d\rho$ is 0, so by the uniqueness of the Stieltjes transform it follows that this measure is the zero measure.

Now, if \hat{v} is the transform of a finite function in \mathcal{H}_1 , then it is an entire function, so to prove that t is outside the support of $\hat{u} d\rho$ it is enough to show that there is such a \hat{v} for which $\hat{v}(t) \neq 0$. If $t \neq 0$ and $\hat{v}(t) = 0$ for all compactly supported $v \in \mathcal{H}_1$, then as in the proof of Proposition 2.4 it follows that $\varphi(\cdot, t)$ satisfies (1.3) both for $\lambda = 0$ and $\lambda = t$, so that $\varphi(\cdot, t)w = 0$ a.e., which is not possible since it implies that $\varphi(\cdot, t) = 0$ in a set of positive Lebesgue measure. It therefore follows that $\hat{u} d\rho$ vanishes outside 0. But according to Lemma 3.5 this proves that the measure is zero, unless $\alpha = 0$. However, also in this case $\hat{u} = 0$ since otherwise \hat{u} would be the transform of an eigenfunction. \square

We next turn to Theorem 3.2(3).

LEMMA 3.7. *If $u \in D_T$, then $\mathcal{F}(Tu)(t) = t\hat{u}(t)$. Conversely, if \hat{u} and $t\hat{u}(t)$ are in L^2_ρ , then $\mathcal{F}^{-1}(\hat{u}) \in D_T$.*

Proof. We have $u \in D_T$ if and only if for some $v \in \mathcal{H}_1$ we have $u = R_\lambda(v - \lambda u)$, i.e., if and only if $\hat{u}(t) = (\hat{v}(t) - \lambda\hat{u}(t))/(t - \lambda)$ or $t\hat{u}(t) = \hat{v}(t)$ for some $\hat{v} \in L^2_\rho$. \square

We obtain the following corollary which will be useful later on.

COROLLARY 3.8. *If $u \in D_T$, then \hat{u} is integrable with respect to $d\rho$.*

Proof. The functions $t\hat{u}(t)$, \hat{u} , and $1/(t - i)$ are all in L^2_ρ , so that $\hat{u}(t) = (t\hat{u}(t) - i\hat{u}(t))/(t - i)$ is integrable with respect to $d\rho$. \square

To finish the proof of Theorem 3.2 it only remains to consider the inverse transform.

LEMMA 3.9. *If $\alpha \neq 0$, the integral $\langle \hat{u}, \varphi(x, \cdot) \rangle_\rho$ converges in \mathcal{H} and locally uniformly for every $\hat{u} \in L^2_\rho$. If $\hat{u} = \mathcal{F}(u)$ for some $u \in \mathcal{H}_1$, then the integral is the orthogonal projection of u onto \mathcal{H} .*

If $\alpha = 0$, the same statement is true if one replaces $\varphi(\cdot, 0)$ by ψ_0 in the integral.

Remark 3.10. A simple integration by parts shows that every finite function is orthogonal to φ_0 . Now suppose $\alpha = 0$ and let $\theta_0 = \varphi(\cdot, 0)$ so that θ_0 solves (1.3) for $\lambda = 0$ with initial data $\theta_0(0) = 0, \theta_0'(0) = 1$. Then, when calculating the Fourier transform at 0 we may replace $\varphi(\cdot, 0)$ by any function $\theta_0 + A\varphi_0$ for A constant, with no change to the Fourier transform.

In particular we may choose $A = -\psi_0(0) = 1/\rho\{0\}$, according to Lemma 3.5, so that $\theta_0 + A\varphi_0 = \psi_0$. This might seem a more natural choice of kernel for the Fourier transform, in view of the fact that it must be used for the inverse transform, and that ψ_0 is an eigenfunction to the eigenvalue 0, but would thus not actually change the Fourier transform.

Proof of Lemma 3.9. We have $u(x) = \langle u, g_0(x, \cdot) \rangle = \langle \hat{u}, e(x, \cdot) \rangle_\rho$ for $u \in \mathcal{H}$, where $e(x, t) = \mathcal{F}(g_0(x, \cdot))(t)$. If $u \in \mathcal{H}_1$, we instead get the projection of u onto \mathcal{H} , so that the integral operator $\hat{u} \mapsto \langle \hat{u}, e(x, \cdot) \rangle_\rho$ is the adjoint of \mathcal{F} . We must prove that $e(x, t) = \varphi(x, t)$, so suppose \hat{u} has compact support and consider $\tilde{u}(x) = \langle \hat{u}, \varphi(x, \cdot) \rangle_\rho$ which satisfies the equation $-\tilde{u}'' + q\tilde{u} = w(x)\langle \hat{u}, t\varphi(x, \cdot) \rangle_\rho$, differentiating under the integral sign. Since \hat{u} has compact support $u \in D_T$, so that $-u'' + qu = w(x)\langle t\hat{u}(t), e(x, t) \rangle_\rho$. Thus $u_1 = u - \tilde{u}$ satisfies $-u_1'' + qu_1 = w(x)\langle t\hat{u}(t), e(x, t) - \varphi(x, t) \rangle_\rho$.

Now, if v is finite, then

$$\langle \tilde{u}, v \rangle = \iint \hat{u}(t)(\varphi'(\cdot, t)\bar{v}' + q\varphi(\cdot, t)\bar{v}) d\rho(t) = \langle \hat{u}, \hat{v} \rangle_\rho = \langle u, v \rangle,$$

since the double integral is absolutely convergent. Hence u_1 is orthogonal to all finite v so it satisfies $-u_1'' + qu_1 = 0$. It follows that $w(x)\langle t\hat{u}(t), e(x, t) - \varphi(x, t) \rangle_\rho = 0$ a.e., so that $\langle t\hat{u}(t), e(x, t) - \varphi(x, t) \rangle_\rho = 0$ on a set of positive measure. But this function also satisfies (1.3) for $\lambda = 0$, as is seen by replacing \hat{u} by $t\hat{u}(t)$ in the previous calculations. It follows that $t(e(x, t) - \varphi(x, t)) = 0$ for a.a. t with respect to $d\rho$, so that $e(x, t) = \varphi(x, t)$ except possibly if $t = 0$ and $\alpha = 0$.

However, 0 is an eigenvalue for $\alpha = 0$ and the eigenfunction ψ_0 has transform $\chi_{\{0\}}/\rho\{0\}$ according to Lemma 3.5, so we must choose $e(x, 0) = \psi_0(x)$. \square

The proof of Theorem 3.2 is now complete if we note that from $\langle E_t u, v \rangle = \int_{-\infty}^t \hat{u}\bar{v}$ follows that the transform of $E_t u$ is \hat{u} multiplied by the characteristic function of $(-\infty, t]$. The formula $E_M u(x) = \int_M \hat{u}\varphi(x, \cdot) d\rho$ therefore follows from the inversion formula.

In Lemma 3.5 we calculated the Fourier transform of ψ_0 in the case $\alpha = 0$. We shall need to find a few more Fourier transforms.

LEMMA 3.11. *If $\lambda \notin \mathbb{R}$, the Fourier transform of $\psi(\cdot, \lambda)$ is $\hat{\psi}(t, \lambda) = 1/(t - \lambda)$. Furthermore, the Fourier transform of ψ_0 equals $\hat{\psi}_0(t) = \sin \alpha/t$ for $\alpha \neq 0$ and $1/\rho\{0\}$ times the characteristic function of the set $\{0\}$ for $\alpha = 0$.*

Proof. We have already calculated $\hat{\psi}_0$ for $\alpha = 0$ in Lemma 3.5. If $\alpha \neq 0$, we note that $\psi_0(x) = -g_0(0, x)$ so its Fourier transform is $-e(0, t) = -\varphi(0, t) = \sin \alpha/t$.

According to (2.7), Theorem 2.9, and Lemma 3.9, for $u \in \mathcal{H}$ we have

$$\begin{aligned} -\sin \alpha \langle \hat{u}, \hat{\psi}(\cdot, \lambda) \rangle_\rho &= \bar{\lambda}\varphi(0, \bar{\lambda})\langle u, \psi(\cdot, \lambda) \rangle \\ &= \bar{\lambda}R_{\bar{\lambda}}u(0) + u(0) = \left\langle \left(\frac{\bar{\lambda}}{t - \bar{\lambda}} + 1 \right) \hat{u}(t), e(0, t) \right\rangle_\rho \\ &= \left\langle \hat{u}(t), \frac{te(0, t)}{t - \lambda} \right\rangle_\rho = -\sin \alpha \left\langle \hat{u}(t), \frac{1}{t - \lambda} \right\rangle_\rho \end{aligned}$$

so that we have $\hat{\psi}(t, \lambda) = 1/(t-\lambda)$ if $\alpha \neq 0$. If $\alpha = 0$, we assume \hat{u} has compact support so that we may differentiate $u(x) = \langle \hat{u}, e(x, \cdot) \rangle_\rho$ under the integral sign to obtain

$$\begin{aligned} \langle \hat{u}, \hat{\psi}(\cdot, \lambda) \rangle_\rho &= \varphi'(0, \bar{\lambda}) \langle u, \psi(\cdot, \lambda) \rangle = (R_{\bar{\lambda}}u)'(0) \\ &= \left\langle \frac{\hat{u}(t)}{t-\bar{\lambda}}, e'_x(0, t) \right\rangle_\rho = \left\langle \hat{u}(t), \frac{e'_x(0, t)}{t-\bar{\lambda}} \right\rangle_\rho = \left\langle \hat{u}(t), \frac{1}{t-\bar{\lambda}} \right\rangle_\rho. \end{aligned}$$

Thus, also in this case we obtain $\hat{\psi}(t, \lambda) = 1/(t-\lambda)$. □

COROLLARY 3.12. *Suppose $u \in \mathcal{H}$. Then $\langle u, \psi(\cdot, t\lambda) \rangle \rightarrow 0$ as $t \rightarrow \infty$, locally uniformly for $\lambda \notin \mathbb{R}$. By (2.1) this means that $\psi(x, t\lambda) \rightarrow 0$ as $t \rightarrow \infty$, locally uniformly in x and $\lambda \notin \mathbb{R}$.*

In fact, unless $0 \notin \text{supp } w$ and $\alpha = \pi/2$ we have $\psi(\cdot, t\lambda) \rightarrow 0$ in \mathcal{H} , locally uniformly in $\lambda \notin \mathbb{R}$ as $t \rightarrow \infty$.

Proof. We have $\langle u, \psi(\cdot, \lambda) \rangle = \langle \hat{u}, \hat{\psi}(\cdot, \lambda) \rangle_\rho$. With the extra assumptions Proposition 2.11 shows that $\psi(\cdot, \lambda) \in \mathcal{H}$ so that $\|\psi(\cdot, \lambda)\| = \|\hat{\psi}(\cdot, \lambda)\|_\rho$.

It follows immediately by dominated convergence from Lemma 3.11 that the claims are true. □

Remark 3.13. All of the theory of sections 2 and 3 extends with no essential change to the case when w is just a measure, or even an element of $H_{\text{loc}}^{-1}(0, b)$.

4. Uniqueness of the inverse problem. We shall here deal with the following question: *To what extent is the operator T , i.e., the interval $[0, b)$, the coefficients q and w , and the boundary condition parameter α , determined by the spectral measure $d\rho$?* To answer this question we introduce the concept of a *Liouville transform* as a map $v \mapsto u$ given by $u(x) = f(x)v(g(x))$, where f and g are fixed functions. We suppose that g is strictly increasing and continuous, and that f is never 0. It is then easy to see that the inverse of a Liouville transform is also a Liouville transform, as is the composition of two Liouville transforms.

Now consider another relation \check{T} of the same type as T , with Hilbert space $\check{\mathcal{H}}_1$, interval $[0, \check{b})$, boundary condition parameter $\check{\alpha}$, and coefficients \check{q} and \check{w} . We will assume, as we do for \mathcal{H}_1 , that finite functions are dense in $\check{\mathcal{H}}_1$.

THEOREM 4.1. *Suppose that $\alpha = \check{\alpha}$, or $0 < \alpha = \pi/2 - \check{\alpha} < \pi/2$, or $\pi/2 < \alpha = 3\pi/2 - \check{\alpha} < \pi$ and that there is a continuously differentiable bijection g from $[0, b)$ to $[0, \check{b})$ with the following properties: $g, g',$ and g'' are locally absolutely continuous, $g' > 0, g(0) = g''(0) = 0, g'(0) = (\sin \check{\alpha} / \sin \alpha)^2$ if $\alpha \neq 0 \neq \check{\alpha}, g'(0) = 1$ if $\alpha = \check{\alpha} = 0$, and the coefficients of T and \check{T} satisfy $\check{q}(g(x)) = (-f(x)f''(x) + q(x)f(x)^2)/g'(x)$ and $\check{w}(g(x)) = w(x)/g'(x)^2$, where $f(x) = g'(x)^{-1/2}$.*

Then the spectral measures associated with T and \check{T} are identical.

Proof. The functions g and f give rise to Liouville transform \mathcal{L} from functions defined on $[0, \check{b})$ to functions defined on $[0, b)$, in particular to a transform from $\check{\mathcal{H}}_1$ to \mathcal{H}_1 . We will first show that this latter transform is unitary. To that end assume that \check{u} and \check{v} are in $\check{\mathcal{H}}_1$ and that at least one of them is a finite function. Obviously $\mathcal{L}\check{u}$ and $\mathcal{L}\check{v}$ are locally absolutely continuous. Furthermore we obtain after a partial integration

$$\begin{aligned} \langle \mathcal{L}\check{u}, \mathcal{L}\check{v} \rangle_{\mathcal{H}_1} &= \int_0^b (g'(\check{u}'\check{v}') \circ g + (-ff'' + qf^2)(\check{u}\check{v}) \circ g) \\ &= \int_0^{\check{b}} (\check{u}'\check{v}' + \check{q}\check{u}\check{v}) = \langle \check{u}, \check{v} \rangle_{\check{\mathcal{H}}_1}. \end{aligned}$$

This proves first that $\mathcal{L}\check{u} \in \mathcal{H}_1$ whenever \check{u} is a finite function in $\check{\mathcal{H}}_1$ and second that \mathcal{L} is an isometry from the finite functions in $\check{\mathcal{H}}_1$ onto the finite functions in \mathcal{H}_1 . As an isometry \mathcal{L} can be extended to a unitary operator from $\check{\mathcal{H}}_1$ to \mathcal{H}_1 .

Next, a straightforward computation, using that $2f'g' + fg'' = 0$, shows that $-u'' + qu = wr$ if $u = \mathcal{L}\check{u}$, $r = \mathcal{L}\check{r}$, and $-\check{u}'' + \check{q}\check{u} = \check{w}\check{r}$. In particular, $(\check{u}, \check{r}) \in \check{T}$ implies that $(\mathcal{L}\check{u}, \mathcal{L}\check{r}) \in T$ and $\mathcal{L}\check{\psi}(\cdot, \lambda)$ must be a multiple of $\psi(\cdot, \lambda)$.

Also, since $\check{\varphi}(\cdot, \lambda)$ satisfies the differential equation $-\check{u}'' + \check{q}\check{u} = \lambda\check{u}$ the function $\mathcal{L}\check{\varphi}(\cdot, \lambda)$ satisfies $-u'' + qu = \lambda wu$. Our assumptions on α , $\check{\alpha}$, $g'(0)$, and $g''(0)$ imply that $f(0) = \sin \alpha / \sin \check{\alpha} = \cos \check{\alpha} / \cos \alpha$ and that $f'(0) = 0$. Therefore we find $\lambda(\mathcal{L}\check{\varphi}(\cdot, \lambda))(0) = \lambda f(0)\check{\varphi}(0, \lambda) = -\sin \alpha$ and $(\mathcal{L}\check{\varphi}(\cdot, \lambda))'(0) = \check{\varphi}'(0, \lambda)/f(0) = \cos \alpha$ which shows that $\varphi(\cdot, \lambda) = \mathcal{L}\check{\varphi}(\cdot, \lambda)$. The situation is a little more complicated for the relationship between θ and $\check{\theta}$ where one finds that

$$\mathcal{L}\check{\theta}(\cdot, \lambda) = \theta(\cdot, \lambda) + (\tan \check{\alpha} - \tan \alpha)\varphi(\cdot, \lambda).$$

By the linearity of \mathcal{L} we have

$$\mathcal{L}\check{\psi}(\cdot, \lambda) = \theta(\cdot, \lambda) + (\tan \check{\alpha} - \tan \alpha + \check{m})\varphi(\cdot, \lambda) = \psi(\cdot, \lambda).$$

This proves that $\check{m} + \tan \check{\alpha} = m + \tan \alpha$ and hence that $\check{\rho} = \rho$. □

In the rest of this section we will make the following additional assumption about (1.3).

Assumption 4.2. The coefficients w and \check{w} satisfy $\text{supp } w = [0, b)$, $\text{supp } \check{w} = [0, \check{b})$.

Note that this does not mean that $w \neq 0$ a.e.; w could vanish on a nowhere dense set of strictly positive measure. However, it does mean that $\mathcal{H}_\infty = \{0\}$, $\mathcal{H} = \mathcal{H}_1$.

Remark 4.3. One may also allow w to be an arbitrary measure. However, then in the definition of the function h below, and in the statement of Lemma 5.1, w should be replaced by the density of the absolutely continuous part of the measure w , and Assumption 4.2 will have to be made on this density. If this is done, the results in the rest of the paper are still true, *mutatis mutandis*, with essentially the same proofs.

Now define the functions $h(x) = \int_0^x \sqrt{|w|}$ on $[0, b)$ and $\check{h}(x) = \int_0^x \sqrt{|\check{w}|}$ on $[0, \check{b})$, respectively. By Assumption 4.2 these are strictly increasing, locally absolutely continuous functions.

Our main theorem is the following.

THEOREM 4.4. *Suppose that T and \check{T} have the same spectral measure dp . Then there is a unitary Liouville transform \mathcal{U} taking \check{T} into T , in the sense that $\mathcal{H} \ni u \mapsto \mathcal{U}u \in \check{\mathcal{H}}$ through $u(x) = f(x)\mathcal{U}u(g(x))$ and $\mathcal{U}T = \check{T}\mathcal{U}$. Here $g(x) = \check{h}^{-1} \circ h(x)$ and $f(x) = (g'(x))^{-1/2}$.*

The functions f and g are continuously differentiable, f is strictly positive, and f' is locally absolutely continuous with $f'(0) = 0$. Also $\alpha = \check{\alpha}$, in which case $f(0) = 1$, or else $0 < \alpha = \pi/2 - \check{\alpha} < \pi/2$ or $\pi/2 < \alpha = 3\pi/2 - \check{\alpha} < \pi$, in which case $f(0) = |\tan \alpha|$.

The relations between the coefficients are $\check{w}(g(x)) = w(x)/(g'(x))^2$ and $\check{q}(g(x)) = (-f''(x) + q(x)f(x))/(f(x)(g'(x))^2)$.

It is clear from Theorem 4.1 that Theorem 4.4 is optimal in the sense that it is not possible to deduce more about the relation between T and \check{T} from the equality of their spectral measures than is done in Theorem 4.4. Sufficient additional information, however, will imply that T and \check{T} are identical. We give two corollaries of this type.

COROLLARY 4.5. *Suppose T and \check{T} have the same spectral measure and that $|w| = |\check{w}|$ in $[0, \min(b, \check{b})]$. Then $T = \check{T}$, i.e., $b = \check{b}$, $\alpha = \check{\alpha}$, $q = \check{q}$, and $w = \check{w}$.*

Proof. The assumptions together with Theorem 4.4 show that $g(x) = x$ so that $b = \check{b}$, and that $f(x) = 1$, so that T and \check{T} are identical. □

Note that only the absolute value of w need be known, so that all information about sign changes in w is encoded in the spectral measure. Also note that if $|w| = |\check{w}|$ only in $[0, a)$, where $0 < a < \min(b, \check{b})$, we still have $\alpha = \check{\alpha}$ and $q = \check{q}$, $w = \check{w}$ in $[0, a)$.

COROLLARY 4.6. *Suppose T and \check{T} have the same spectral measure, that $q = \check{q}$ on $[0, \min(b, \check{b}))$, and that either $b = \check{b}$ or $\alpha = \check{\alpha}$. Then $T = \check{T}$, i.e., $b = \check{b}$, $\alpha = \check{\alpha}$, $q = \check{q}$, and $w = \check{w}$.*

We will postpone the proof and first prove Theorem 4.4. To do this we will use a theorem of Paley–Wiener type. For its statement it will be convenient to introduce a special class of entire functions.

DEFINITION 4.7. *Let \mathcal{A} be the set of entire functions \hat{u} of order $\leq 1/2$ which satisfy*

$$(4.1) \quad \limsup_{t \rightarrow \infty} t^{-1} \ln |\hat{u}(t^2 \lambda)| \leq \int_0^a \operatorname{Re} \sqrt{-\lambda w}$$

for some $a \in (0, b)$ and all $\lambda \in \mathbb{C} \setminus \mathbb{R}$. Here the branch of the square root is that with a positive real part.

THEOREM 4.8. *Let \hat{u} be the generalized Fourier transform of $u \in \mathcal{H}$. Then \hat{u} has at most one entire continuation in \mathcal{A} , and if $\operatorname{supp} u = a < b$, such a continuation is given by*

$$\hat{u}(\lambda) = \int_0^a (u' \varphi'(\cdot, \lambda) + qu \varphi(\cdot, \lambda))$$

in which case (4.1) holds with equality for all $\lambda \in \mathbb{C}$.

Conversely, if \hat{u} has an entire continuation of order $\leq 1/2$ satisfying (4.1) for λ on at least two different rays from the origin, then $\operatorname{supp} u \subset [0, a]$.

We will postpone the proof of Theorem 4.8 to the next section and instead turn to the proof of Theorem 4.4.

LEMMA 4.9. *Let $g : [0, b) \rightarrow [0, \check{b})$ be increasing and $g(0) = 0$. Suppose $\mathcal{U} : \mathcal{H}_1 \rightarrow \check{\mathcal{H}}_1$ is linear with the properties that $(\mathcal{U}u)(0) = 0$ if $u(0) = 0$, that $\operatorname{supp} \mathcal{U}u \subset [0, g(x)]$ if $\operatorname{supp} u \subset [0, x]$, and that $\operatorname{supp} \mathcal{U}u \subset [g(x), \check{b})$ if $\operatorname{supp} u \subset [x, b)$. Then there exists a function f such that $(\mathcal{U}u)(g(x)) = f(x)u(x)$ for all $u \in \mathcal{H}_1$.*

Proof. Fix $x \in [0, b)$. Suppose $u, v \in \mathcal{H}_1$ and that $u(x) = v(x)$. We will first show that $(\mathcal{U}(u - v))(g(x)) = 0$. If $x = 0$, this is by assumption.

For $x > 0$ we define¹ $u_- = \chi_{[0, x]}(u - v)$ and $u_+ = \chi_{[x, b)}(u - v)$. These are elements of \mathcal{H} . Thus $\operatorname{supp} \mathcal{U}u_- \subset [0, g(x)]$ and $\operatorname{supp} \mathcal{U}u_+ \subset [g(x), b)$ so that the functions $\mathcal{U}u_{\pm}$ vanish in $g(x)$. Adding them gives $\mathcal{U}(u - v)(g(x)) = 0$ as desired.

It follows that the value of $\mathcal{U}u$ at $g(x)$ only depends on the value of u at x . Thus, for each fixed $x \in [0, b)$, the map $u(x) \mapsto \mathcal{U}u(g(x))$ is well-defined and linear on \mathbb{C} , so we may find $f(x)$ so that $\mathcal{U}u(g(x)) = f(x)u(x)$. \square

We will also need the following lemma.

LEMMA 4.10. *Put $m(x, \lambda) = \psi'(x, \lambda)/(\lambda \psi(x, \lambda))$. Then $m(x, \lambda) \rightarrow 0$ and $\lambda m(x, \lambda) \rightarrow \infty$ for every $x \in [0, b)$ as $\lambda \rightarrow \infty$ along any nonreal ray starting from the origin.*

Proof. First note that $m(x, \lambda)$ is the m -function for (1.3) on the interval $[x, b)$, with the Dirichlet boundary condition ($\alpha = 0$) at x . The first claim is then an immediate consequence of [3, Theorem 3.6].

¹ χ_I denotes the characteristic function of an interval I .

To prove the second claim, first assume that q does not have compact support, so that it does not vanish identically on $[x, b)$. Now note that, according to (3.2), $\tilde{m}(\lambda) = -1/m(x, \lambda)$ is the m -function for the Neumann boundary condition ($\alpha = \pi/2$) at x , so we need to show this to be $o(|\lambda|)$. Now, in the Nevanlinna representation (2.10) it is easy to see that the integral is always $o(|\lambda|)$, so we simply need to prove that $B = 0$ in the representation of \tilde{m} . Denote the corresponding Weyl solution by $\tilde{\psi}$ and the spectral measure by $d\tilde{\rho}$. Using (2.9) and Lemma 3.11 we obtain

$$\|\tilde{\psi}(\cdot, \lambda)\|_{[x,b)}^2 = \frac{\text{Im } \tilde{m}(\lambda)}{\text{Im } \lambda} = B + \int_{-\infty}^{\infty} \frac{d\tilde{\rho}(t)}{|t - \lambda|^2} = B + \|\hat{\psi}(\cdot, \lambda)\|_{\tilde{\rho}}^2.$$

However, by Proposition 2.11, Parseval’s formula is correct for $\tilde{\psi}$, so that $B = 0$ and we are done in the case when q does not have compact support.

Now suppose q vanishes identically in $[x, b)$. Consider an auxiliary equation for which q does not have compact support, but which has the same coefficients as (1.3) up to some point c , $x < c < b$. For this equation the above proof of the lemma is valid. Moreover, let $\tilde{\theta}$ and $\tilde{\varphi}$ denote functions analogous to θ and φ for $\alpha = 0$, but with initial data given in the point x . In view of (2.9) both the original $m(x, \lambda)$ and the corresponding function for the auxiliary equation are in the “Weyl disk” defined by

$$\int_x^c |\tilde{\theta}' + m\tilde{\varphi}'|^2 \leq \frac{\text{Im } m}{\text{Im } \lambda},$$

so their distance is bounded by the diameter of the disk, which is exponentially small as λ becomes large (see [3, Theorem 6.3] for this result). Since $m(x, \lambda)$ is a nontrivial Nevanlinna function it cannot tend to 0 faster than a multiple of $1/|\lambda|$ for large $|\lambda|$, so that asymptotically $m(x, \lambda)$ is the same as the corresponding function for the auxiliary equation. Thus the lemma is actually valid in all cases. \square

Proof of Theorem 4.4. Note first that by Lemma 3.5 we must have either $\alpha = \check{\alpha} = 0$ or else $\alpha \neq 0 \neq \check{\alpha}$.

Let \mathcal{H} and $\check{\mathcal{H}}$ denote the Hilbert spaces and \mathcal{F} and $\check{\mathcal{F}}$ the generalized Fourier transforms associated with the two equations, and put $\mathcal{U} = \check{\mathcal{F}}^{-1} \circ \mathcal{F} : \mathcal{H} \rightarrow \check{\mathcal{H}}$, which is unitary since the target space is L^2_ρ for both \mathcal{F} and $\check{\mathcal{F}}$. By Lemma 3.11 we have $\mathcal{U}\psi_0 = \check{\psi}_0$ if $\alpha = \check{\alpha}$, and if $\alpha \neq 0 \neq \check{\alpha}$, we have $\mathcal{U}\psi_0 = \frac{\sin \alpha}{\sin \check{\alpha}} \check{\psi}_0$. Since $\langle u, \psi_0 \rangle = -u(0)$ it follows that

$$(4.2) \quad u(0) = -\langle u, \psi_0 \rangle = -\langle \mathcal{U}u, \mathcal{U}\psi_0 \rangle = \frac{\sin \alpha}{\sin \check{\alpha}} \mathcal{U}u(0),$$

where the quotient of the sines is to be read as 1 for $\alpha = \check{\alpha} = 0$. In particular, $\mathcal{U}u(0) = 0$ if and only if $u(0) = 0$.

Now, applying Theorem 4.8 for the rays generated by $\pm i$, it is clear that if $\check{a} \in (0, \check{b})$ and $u \in \mathcal{H}$, then $\text{supp } u = a$ if $\text{supp } \mathcal{U}u = \check{a}$, where $h(a) = \check{h}(\check{a})$, provided there is such an $a \in (0, b)^2$ (see [4, p. 29] for more details). This will certainly be the case if \check{a} is sufficiently close to 0. Suppose for some $\check{a} \in (0, \check{b})$ we have $h(b) \leq \check{h}(\check{a})$. Then, since compactly supported functions are dense in \mathcal{H} , the range of \mathcal{U} would be orthogonal to all elements of $\check{\mathcal{H}}$ with supports in (\check{a}, \check{b}) , contradicting the fact that \mathcal{U} is unitary.

A similar reasoning applied to \mathcal{U}^{-1} shows that the mapping

$$g : [0, b) \ni a \mapsto \check{a} \in [0, \check{b})$$

²Note that $\text{Re } \sqrt{\pm iw} = \sqrt{|w|/2}$.

is bijective, and that $\sup \text{supp } \mathcal{U}u = \check{a}$ if $\sup \text{supp } u = a$. It follows that $\sup \text{supp } u = a$ if and only if $\sup \text{supp } \mathcal{U}u = g(a)$.

We also have $\inf \text{supp } u = a$ if and only if $\inf \text{supp } \mathcal{U}u = g(a)$. To see this, note that what we have already proved implies that if $\inf \text{supp } u = a > 0$, then $\mathcal{U}u$ is orthogonal to all elements of $\check{\mathcal{H}}$ with support in $[0, g(a)]$. This means that in this interval $\mathcal{U}u$ is a multiple of $\check{\varphi}_0$. However, since $u(0) = 0$ we also have $\mathcal{U}u(0) = 0$, so that the multiple is 0, and thus $\inf \text{supp } \mathcal{U}u \geq g(a)$. A similar reasoning applied to \mathcal{U}^{-1} proves the other direction.

We have now verified that \mathcal{U} and \mathcal{U}^{-1} both have the properties required in Lemma 4.9. This implies that there is a nonvanishing function f so that

$$(4.3) \quad u(x) = f(x)\mathcal{U}u(g(x)).$$

We must have f real-valued since \mathcal{F} and $\check{\mathcal{F}}^{-1}$, and thus \mathcal{U} , map real-valued functions to real-valued functions. We note that (4.2) implies that $f(0) = 1$ if $\alpha = \check{\alpha} = 0$ and $f(0) = \frac{\sin \alpha}{\sin \check{\alpha}} > 0$ if $\alpha \neq 0 \neq \check{\alpha}$. Now choose $\mathcal{U}u = 1$ in a neighborhood of $g(x)$. We then have $u = f$ in a neighborhood of x . Since $u \in \mathcal{H}$ is locally absolutely continuous, so is f . This also implies that f is strictly positive, since it cannot change sign and $f(0) > 0$. Similarly, choosing $\mathcal{U}u$ linear in a neighborhood of $g(x)$ it follows that also g is locally absolutely continuous.

According to Lemma 3.11 $\mathcal{U}\psi(\cdot, \lambda) = \check{\psi}(\cdot, \lambda)$, so we have $\psi(x, \lambda) = f(x)\check{\psi}(g(x), \lambda)$. Taking the logarithmic derivative we obtain

$$\frac{\psi'(x, \lambda)}{\psi(x, \lambda)} = \frac{f'(x)}{f(x)} + g'(x) \frac{\check{\psi}'(g(x), \lambda)}{\check{\psi}(g(x), \lambda)}.$$

Here the left member and the coefficient for $g'(x)$ are locally absolutely continuous, and the coefficient for $g'(x)$ is not independent of λ by Lemma 4.10. It follows that g' and f' are locally absolutely continuous, and differentiating, using the differential equations, we obtain

$$-\frac{f''}{f} + q - (g')^2 \check{q} \circ g - \lambda(w - (g')^2 \check{w} \circ g) = \frac{(f^2 g')'}{f^2} \frac{\check{\psi}'(g(\cdot), \lambda)}{\check{\psi}(g(\cdot), \lambda)}.$$

Here the right member is $o(|\lambda|)$ according to Lemma 4.10 so the coefficient of λ to the left vanishes. On the other hand, the right member is not independent of λ unless $(f^2 g')' = 0$, so that we obtain

$$\begin{aligned} \check{q} \circ g &= \frac{1}{f(g')^2} (-f'' + qf), \\ \check{w} \circ g &= (g')^{-2} w, \\ f^2 g' &= C \end{aligned}$$

for some constant C . Evaluating (4.3) and its derivative at 0 for $u = \psi(\cdot, \lambda)$ elementary calculations now shows³ that $C = 1$ and $f'(0) = 0$. One also deduces that either $\alpha = \check{\alpha}$ or else $0 < \alpha = \pi/2 - \check{\alpha} < \pi/2$ or $\pi/2 < \alpha = 3\pi/2 - \check{\alpha} < \pi$. In these calculations one uses that \check{m} is not a Möbius transform, which is clear since this would give

³See the appendix.

a transform space of dimension 1. This can only happen if w , and $d\rho$, is a point mass. \square

Finally we have to prove Corollary 4.6.

Proof of Corollary 4.6. The function $\tilde{f} = -\varphi_0$ solves $-\tilde{f}'' + q\tilde{f} = 0$ with initial data $\tilde{f}(0) = 1, \tilde{f}'(0) = 0$. Since $q \geq 0$ this solution is strictly positive on $[0, b)$, so we may put $\tilde{g}(x) = \int_0^x 1/\tilde{f}^2$. The pair of functions \tilde{f}, \tilde{g} gives us a Liouville transform F_0 mapping $[0, b)$ onto some interval $[0, c)$ and $[0, \check{b})$ onto $[0, \check{c})$, and transforming the equations into $-u_0'' = \lambda w_0 u_0$ and $-\check{u}_0'' = \lambda \check{w}_0 \check{u}_0$, respectively. Thus $F_0 F F_0^{-1}$, where F is the Liouville transform of Theorem 4.4, transforms one of these equations into the other.

Being a composition of Liouville transforms this is itself a Liouville transform given, say, by $u_0(x) = f_1(x)\check{u}_0(g_1(x))$. By construction we obtain $f_1(0) = f(0), f_1'(0) = 0$, and $f_1^2 g_1' \equiv 1$. Since both potentials are identically 0 it follows that $f_1'' = 0$. This means that $f_1 \equiv f(0)$ and $g_1(x) = x/(f(0))^2$.

If $\alpha = \check{\alpha}$, then by Theorem 4.4 $f(0) = 1$ so that $F_0 F F_0^{-1}$ is the identity, implying that also F is the identity. Similarly, if $b = \check{b}$, then $c = \check{c}$ so that $f(0) = 1$, unless $c = \check{c} = \infty$. We will show that c is always finite, and then it again follows that F is the identity.

Now $c = \int_0^b 1/\tilde{f}^2$, so we need to show that this integral is finite. Put $H = \tilde{f}'\tilde{f}$ which will be strictly positive sufficiently close to b by (2.4).

Differentiating $H' = (\tilde{f}')^2 + \tilde{f}''\tilde{f} = (\tilde{f}')^2 + q\tilde{f}^2 \geq (\tilde{f}')^2$. Thus $1/\tilde{f}^2 = (\tilde{f}')^2/H^2 \leq H'/H^2$ so that $\int_d^b 1/\tilde{f}^2 \leq 1/H(d) < \infty$ if d is sufficiently close to b . This completes the proof. \square

5. The Paley–Wiener theorem. The proof of Theorem 4.8 relies on the following lemma, which is taken from [3, Theorem 6.1, Corollary 6.2].

LEMMA 5.1. *The following asymptotic formulas hold, locally uniformly for $\lambda \in \mathbb{C} \setminus \mathbb{R}$ and $x > 0$. The square root refers to the branch with positive real part:*

$$\lim_{t \rightarrow \infty} t^{-1} \ln \varphi(x, t^2 \lambda) = \int_0^x \sqrt{-\lambda w},$$

$$\lim_{t \rightarrow \infty} t^{-1} \ln \psi(x, t^2 \lambda) = - \int_0^x \sqrt{-\lambda \check{w}}.$$

The next lemma implies the simple direction of Theorem 4.8.

LEMMA 5.2. *Suppose $u \in \mathcal{H}$ and $\text{supp } u \subset [0, a]$. Then $\hat{u}(\lambda)$ is entire of order $\leq 1/2$ and $\hat{u}(\lambda) = o(|\lambda \varphi(a + \varepsilon, \lambda)|)$ for every $\varepsilon > 0$ as $\lambda \rightarrow \infty$ along any nonreal ray originating at the origin.*

Proof. For finite u we have $\langle u, \varphi(\cdot, \bar{\lambda}) \rangle = -u(0) \cos \alpha + \int_0^b u \lambda \varphi(\cdot, \lambda) w$. Now write

$$\hat{u}(\lambda) = -u(0) \cos \alpha + \lambda \varphi(a + \varepsilon, \lambda) \int_0^a u \varphi(\cdot, \lambda) w / \varphi(a + \varepsilon, \lambda).$$

The function $\varphi(x, \lambda)/\varphi(a + \varepsilon, \lambda)$ tends to zero uniformly for $x \in [0, a]$ and $\lambda \varphi(a + \varepsilon, \lambda) \rightarrow \infty$ according to Lemma 5.1 as $\lambda \rightarrow \infty$ along a nonreal ray. The lemma follows. \square

The hard direction of Theorem 4.8 follows from the next lemma.

LEMMA 5.3. *Suppose $u \in \mathcal{H}$, that \hat{u} has an entire continuation of order $\leq 1/2$, and that $\hat{u}(\lambda) = \mathcal{O}(1/|\psi(a, \lambda)|)$ as $\lambda \rightarrow \infty$ along two different nonreal rays originating at the origin. Then $\text{supp } u \subset [0, a]$ and $\hat{u}(\lambda) = \langle u, \varphi(\cdot, \lambda) \rangle$.*

Proof. Let $\varepsilon > 0$ and consider $F(\lambda) = \langle R_\lambda u, v \rangle - \hat{u}(\lambda)\langle \psi(\cdot, \lambda), v \rangle$, where $v = G_0(wf)$ and $f \in \mathcal{H}$ has compact support in $(a + \varepsilon, b)$. In particular $v \in D_T$. We shall show that F has an entire continuation of order $\leq 1/2$ which tends to 0 along the given rays. By the Phragmén–Lindelöf principle it follows that F is bounded everywhere and is therefore constant by Liouville’s theorem, thus actually identically 0.

Now $F(\lambda) = \int_0^b (R_\lambda u - \hat{u}(\lambda)\psi(\cdot, \lambda))\overline{f}w$ so, arguing like in the proof of Proposition 2.3, it follows that $R_\lambda u - \hat{u}(\lambda)\psi(\cdot, \lambda)$ has support in $[0, a + \varepsilon]$. Applying the differential equation it follows that also u has support in $[0, a + \varepsilon]$. Since $\varepsilon > 0$ is arbitrary, in fact u has support in $[0, a]$. For $x > a$ the formula (2.7) gives $R_\lambda u(x) = \psi(x, \lambda)\langle u, \overline{\varphi(\cdot, \lambda)} \rangle$ so that $\psi(x, \lambda)(\hat{u}(\lambda) - \langle u, \overline{\varphi(\cdot, \lambda)} \rangle) = 0$. The lemma follows from this.

To prove that F is entire, Parseval’s formula and Lemma 3.11 show that

$$F(\lambda) = \int_{-\infty}^\infty \frac{\hat{u}(t) - \hat{u}(\lambda)}{t - \lambda} \overline{\hat{v}(t)} d\rho(t).$$

It is obvious that this is an entire function, at least if we can bound the integrand properly. To do this and see that the order is at most $1/2$, note that for $|t - \lambda| \leq 1$ we may estimate the integrand by $\sup_{|z| \leq 1} |\hat{u}'(\lambda + z)| |\hat{v}(t)|$. For $|t - \lambda| > 1$ we may estimate the integrand by $|\hat{u}(t)\hat{v}(t)| + |\hat{u}(\lambda)| |\hat{v}(t)|$. Hence we have locally uniformly dominated convergence of the integral and

$$|F(\lambda)| \leq \|u\| \|v\| + \left(\sup_{|z| \leq 1} |\hat{u}'(\lambda + z)| + |\hat{u}(\lambda)| \right) \int_{-\infty}^\infty |\hat{v}| d\rho,$$

which is the required estimate, the integral being finite by Corollary 3.8 and \hat{u} and therefore \hat{u}' being of order $\leq 1/2$.

Finally, to show that F tends to 0 along the rays, we first note that $\psi(x, \lambda)/\psi(a, \lambda)$ converges to 0 uniformly for $x \in [a + \varepsilon, b)$, according to Lemma 5.1. Assuming f has compact support in $[a + \varepsilon, b)$ we obtain $\int_0^b \psi(\cdot, \lambda)\overline{f}w = o(|\psi(a, \lambda)|)$. Since $R_\lambda \rightarrow 0$ strongly as $\text{Im } \lambda \rightarrow \infty$, it follows that F tends to 0 along the given rays. This finishes the proof. \square

Theorem 4.8 is a simple consequence of these lemmas.

Proof of Theorem 4.8. If $\text{supp } u \subset [0, a]$, it follows from Lemmas 5.2 and 5.1 that $\hat{u}(\lambda) = \langle u, \overline{\varphi(\cdot, \lambda)} \rangle$ is an entire continuation of \hat{u} of order $\leq 1/2$ such that

$$\limsup_{t \rightarrow \infty} t^{-1} \ln |\hat{u}(t^2\lambda)| \leq \lim_{t \rightarrow \infty} t^{-1} \ln |\varphi(a + \varepsilon, t^2\lambda)| = \int_0^{a+\varepsilon} \text{Re } \sqrt{-\lambda w}$$

for nonreal λ and all $\varepsilon > 0$.

On the other hand, suppose there is an entire continuation of \hat{u} of order $\leq 1/2$ and such that

$$\limsup_{t \rightarrow \infty} t^{-1} \ln |\hat{u}(t^2\lambda)| \leq \int_0^a \text{Re } \sqrt{-\lambda w}$$

for λ on two different rays from the origin. If one or both of these are real, an immediate application of the Phragmén–Lindelöf principle shows this to be true for all other rays as well, so we may assume them nonreal. By Lemma 5.1 this implies that $\hat{u}(\lambda) = \mathcal{O}(|\psi(a + \varepsilon, \lambda)|^{-1})$ for large λ on these rays if $0 < \varepsilon < b - a$. Lemma 5.3 now shows that $\text{supp } u \subset [0, a + \varepsilon]$ for small $\varepsilon > 0$ and thus for $\varepsilon = 0$. The uniqueness

of the continuation also follows from Lemma 5.3. If we have strict inequality on one ray, a simple argument using the Phragmén–Lindelöf principle (see [4, Lemma 3.6]) shows this to hold on all nearby rays as well, so that in fact $\sup \text{supp } u < a$. The proof is now complete. \square

6. Inverse scattering on the half-line. In this section we will show that scattering data for the half-line problem determines the coefficient w if q is known. We will of course have to assume that our equation is sufficiently close to a model equation, which, as usual, has constant coefficients.

Thus we consider (1.3) on $[0, \infty)$ with the following additional assumption, which will be in force throughout this section.

Assumption 6.1. There is a constant $q_0 \geq 0$ such that $q(x) - q_0$ and $w(x) - 1$ are both in $L^1(0, \infty)$.

Note that according to Theorems 2.5 and 2.6 finite functions are dense in \mathcal{H}_1 and, given the boundary condition (2.6), there is a unique self-adjoint realization T of (1.3) in \mathcal{H}_1 .

We will need the following standard result.

PROPOSITION 6.2. *For $\text{Im } k \geq 0, k \neq 0$, there exists a solution $f(\cdot, k)$ of (1.3) with $\lambda = k^2 + q_0$ having the following properties: (1) $f(x, \cdot)$ and $f'(x, \cdot)$ are analytic for $\text{Im } k > 0$ and continuous for $\text{Im } k \geq 0, k \neq 0$; (2) $f(x, k) \sim e^{ikx}$ and $f'(x, k) \sim ik e^{ikx}$ as $x \rightarrow \infty$.*

This is standard. It is easily proved by first writing the equation for $g(x, k) = f(x, k)e^{-ikx}$ as $g'' + 2ikg' = (q - q_0 - (k^2 + q_0)(w - 1))g$ and then solving this equation by successive approximations from its desired initial values $g(\infty) = 1, g'(\infty) = 0$ at ∞ using the estimate $|e^{2ik(t-x)} - 1| \leq 2$. See, for instance, Deift and Trubowitz [19].

If $\text{Im } k > 0$, then $f(\cdot, k) \in \mathcal{H}_1$. Thus, if $\lambda \notin \mathbb{R}$ (*i.e.*, also $\text{Re } k \neq 0$), then

$$f(x, k) = F(k)\psi(x, \lambda)$$

for some function F defined in $\text{Im } k > 0, \text{Re } k \neq 0$.

Let $[u, v] = u'v - uv'$ denote the Wronskian of the functions u and v and recall that Wronskians of solutions to (1.3) are independent of x . Since

$$(6.1) \quad [\lambda\varphi(\cdot, \lambda), f(\cdot, k)] = F(k)[\lambda\varphi(\cdot, \lambda), \psi(\cdot, \lambda)] = F(k)$$

is analytic for $\text{Im}(k) > 0$ we find that F is analytic and can be extended analytically to the positive imaginary axis. Moreover, since $[\lambda\varphi(\cdot, \lambda), f(\cdot, k)]$ is continuous in $\text{Im}(k) \geq 0, k \neq 0$, the function F extends continuously to the positive and negative real line. The zeros of F are located exactly where φ and f are linearly dependent, *i.e.*, when $\lambda = q_0 + k^2$ is an eigenvalue.

Equation (6.1) gives also that $F(-k) = \overline{F(k)}$ for real $k \neq 0$ and that F has no zeros on either the positive or the negative real line since $\varphi(\cdot, \lambda)$ is real for real λ and the real and imaginary parts of $f(x, k) \sim e^{ikx}$ are linearly independent.

For $k > 0$ and thus $\lambda = k^2 + q_0 > q_0$ define

$$\psi_{\pm}(\cdot, \lambda) = \lim_{\epsilon \rightarrow 0} \psi(\cdot, (\pm k + i\epsilon)^2 + q_0)$$

and

$$m_{\pm}(\lambda) = \lim_{\epsilon \rightarrow 0} m((\pm k + i\epsilon)^2 + q_0).$$

Since $\overline{m(\lambda)} = m(\overline{\lambda})$ when λ is not real we find that $m_+(\lambda) = \overline{m_-(\lambda)}$ when λ is real. Therefore

$$\frac{2ik\lambda}{|F(k)|^2} = \lambda[\psi_+(\cdot, \lambda), \psi_-(\cdot, \lambda)] = m_+(\lambda) - m_-(\lambda) = 2i \operatorname{Im} m_+(\lambda)$$

when $k > 0$ so that $\lambda > q_0$. This in turn implies

$$\pi\rho'(\lambda) = \operatorname{Im} m(\lambda + i0) = \frac{k\lambda}{|F(k)|^2}$$

for $\lambda > q_0$. Thus the restriction of F to the positive real line determines the spectral measure on the interval (q_0, ∞) . It follows from this that the spectrum of T is absolutely continuous⁴ in (q_0, ∞) .

In the interval $(-\infty, q_0)$, where λ corresponds to the positive half of the imaginary axis for k , the spectrum is discrete since F is analytic there. There might also be an eigenvalue for $k = 0, \lambda = q_0$. Suppose $\lambda \neq 0$ is an eigenvalue. Then $\varphi(\cdot, \lambda)$ is a corresponding eigenfunction, and its Fourier transform $\hat{\varphi}(\lambda)$ is a multiple of the characteristic function of the set $\{\lambda\}$. The inversion formula (3.1) gives $\varphi(x, \lambda) = \hat{\varphi}(\lambda)\varphi(x, \lambda)\rho\{\lambda\}$, where $\rho\{\lambda\}$ is the spectral measure of the set $\{\lambda\}$. Thus $\hat{\varphi}(\lambda) = 1/\rho\{\lambda\}$. Parseval's formula gives $\|\varphi(\cdot, \lambda)\|^2 = |\hat{\varphi}(\lambda)|^2\rho\{\lambda\} = 1/\rho\{\lambda\}$. On the interval $(-\infty, q_0]$ we therefore know the spectral measure if we know all eigenvalues λ and the corresponding *normalization constants* $\|\varphi(\cdot, \lambda)\|^2$. Similarly, if $\alpha = 0$, then by Lemma 3.5 also $\lambda = 0$ is an eigenvalue, and $1/\rho\{0\}$ is the normalization constant for the eigenfunction ψ_0 . We obtain the following theorem.

THEOREM 6.3. *Given the absolute value of the coefficient $F(k)$ for positive k , all eigenvalues, the corresponding normalization constants, and either q or $|w|$, the coefficients q and w and the boundary value parameter α are uniquely determined.*

Proof. We have already seen that the given data determine the spectral measure, and may now apply Corollaries 4.5 and 4.6 to draw the desired conclusion. \square

7. Eigenvalues. This section is devoted to the proof of the following theorem. Part of the proof is an adaptation of Marchenko [25].

THEOREM 7.1. *Assume that q and w satisfy Assumption 6.1. Then we have the following:*

- (1) *The eigenvalues of T are isolated and can accumulate only at q_0 or negative infinity.*
- (2) *There will be infinitely many negative eigenvalues if and only if w is negative on a set of positive measure.*

If in addition we have $\int_0^\infty t|q(t) - q_0w(t)|dt < \infty$, we also have the following:

- (3) *Eigenvalues will not accumulate at q_0 .*
- (4) *q_0 is not an eigenvalue unless $q_0 = 0$ and $\alpha = 0$.*

To prove this we need the following strengthening of Proposition 6.2.

PROPOSITION 7.2. *Suppose q and w satisfy Assumption 6.1 and the integral $\int_0^\infty t|q(t) - q_0w(t)|dt$ is finite. Then, for every $x \in [0, \infty)$, the function $f(x, \cdot)$ and its x -derivative, which were previously defined for $\operatorname{Im}(k) \geq 0, k \neq 0$, extend continuously to $k = 0$.*

The additional assumption and the improved estimate

$$|e^{2ik(t-x)} - 1| \leq \min(2|k|t, 2)$$

⁴For $q_0 < s < t$ we have $\int_s^t \operatorname{Im} m(\mu + i\varepsilon) d\mu \rightarrow \pi(\rho(t) - \rho(s))$ as $\varepsilon \downarrow 0$. But the left-hand side converges to $\int_s^t \operatorname{Im} m(\mu + i0) d\mu$ so ρ is absolutely continuous.

allow us to perform the successive approximations also near $k = 0$. The proposition follows from this.

Proof of Theorem 7.1. If $\mu = k^2 + q_0 < q_0$ is an eigenvalue of T , then, since F is analytic in the upper half-plane, eigenvalues are isolated and hence cannot accumulate at any point in $(-\infty, q_0)$. This proves (1).

To prove the second statement we make first the assumption that $q_0 > 0$ and $\alpha \neq 0$. By Lemma 3.5 zero is then not in the spectrum of T so that the range of T is \mathcal{H} and we may define a bilinear form Q on \mathcal{H} by setting

$$Q(u, v) = \int_{\mathbb{R}} \frac{1}{t} \hat{u}(t) \overline{\hat{v}(t)} d\rho(t).$$

Note that $Q(u, v) = 0$ if the supports of \hat{u} and \hat{v} do not intersect, which happens, for instance, if u and v are eigenvectors for different eigenvalues. Furthermore, by Lemma 3.7 $Q(u, Tv) = \int_{\mathbb{R}} \hat{u}(t) \hat{v}(t) d\rho(t) = \langle u, v \rangle$. An integration by parts gives

$$\int_0^x (u' \overline{v'} + qu \overline{v}) = u(x) \overline{v'(x)} - u(0) \overline{v'(0)} + \int_0^x wu \overline{Tv}$$

for $u \in \mathcal{H}$ and $v \in D_T$. Hence if v is in the range of T and u is finite, or if u and v are exponentially decaying eigenfunctions, then we obtain

$$(7.1) \quad Q(u, v) = \int_0^\infty wu \overline{v} + \cot(\alpha)u(0) \overline{v(0)}$$

taking into account the boundary condition satisfied by $(T^{-1}v, v)$.

Now assume that $w \geq 0$. If $\cot(\alpha) \geq 0$, there can be no negative eigenvalue since $Tv = \lambda v$, $\lambda < 0$, $\|v\| \neq 0$ would imply that

$$0 \leq \int_0^\infty w|v|^2 + \cot \alpha |v(0)|^2 = \frac{1}{\lambda} Q(v, Tv) = \frac{1}{\lambda} \|v\|^2 < 0,$$

giving a contradiction. If $\cot \alpha < 0$, there can be at most one negative eigenvalue as we shall show now. If there were two distinct negative eigenvalues λ_1 and λ_2 with associated eigenvectors v_1 and v_2 , we could assume that $v_1(0) = v_2(0)$. This would entail that

$$0 \leq \int_0^\infty w|v_1 - v_2|^2 = Q(v_1 - v_2, v_1 - v_2) = Q(v_1, v_1) + Q(v_2, v_2) < 0$$

since eigenfunctions decay exponentially so that we are allowed to employ (7.1).

Next assume $w < 0$ on a set of positive Lebesgue measure. We shall show that there are infinitely many negative eigenvalues. For any integer n one can choose elements u_1, \dots, u_n in \mathcal{H} , compactly supported in $(0, \infty)$, such that $Q(u_j, u_j) < 0$ and $Q(u_j, u_k) = 0$ if $j \neq k$. To achieve this one may for instance choose first bounded sets A_1, \dots, A_n of positive measure and positive distances from zero and each other on which w is negative. Then one lets u_j be a suitable mollification of the characteristic function of A_j . Equation (7.1) now guarantees that they have the desired properties.

Thus $Q(u, u) < 0$ whenever u is in the linear span B of u_1, \dots, u_n . Let P be the orthogonal projection of B into the negative spectral subspace of \mathcal{H} , i.e., $Pu = \mathcal{F}^{-1}(u\chi)$, where χ is the characteristic function of $(-\infty, 0)$. Suppose now that n is

larger than the number of negative eigenvalues. Then the kernel of P cannot be trivial so that there is a nontrivial $u \in B$ such that \hat{u} is supported in $[0, \infty)$. Hence

$$0 > Q(u, u) = \int_{\mathbb{R}} \frac{1}{t} |\hat{u}(t)|^2 d\rho(t) \geq 0.$$

Since this is impossible the number of negative eigenvalues must be infinite.

If we only have $q_0 \geq 0$, but still $\alpha \neq 0$, then Q remains defined for functions u, v with Fourier transforms bounded near 0, since in this case $1/t \in L^2_\rho$ by Lemma 3.11. But the Fourier transforms of eigenfunctions to nonzero eigenvalues are supported away from 0, and the Fourier transform of a finite function is entire and thus locally bounded. Also, u_j is in the range of T . To see this, solve $-y'' + qy = wu_j$ with 0 initial data at a point to the right of $\text{supp } u_j$ which yields a finite function y . Adding an appropriate multiple of ψ_0 (Proposition 2.7) gives a function in D_T . Thus the proof applies also in this case.

Allowing also $\alpha = 0$ the form Q is still defined if $\hat{u}(t)\overline{\hat{v}(t)}/t$ is continuous at 0. This is the case if u and v are eigenfunctions to negative eigenvalues. Also, if u is a finite function orthogonal to the eigenfunction ψ_0 , then $\hat{u}(0) = 0$; so Q is defined for such functions. This last condition is just one linear condition on the space B , so the remainder can still have arbitrarily large dimension. All of the u_j are in the range of T , since the boundary condition now reads $u_j(0) = 0$. Thus the proof applies also in this case, and the proof of (2) is finished.

Now assume that $\int_0^\infty t|q(t) - q_0w(t)|dt$ is finite, and that, contrary to our claim, there is a sequence $\mu_n = k_n^2 + q_0 < q_0$ of eigenvalues converging to q_0 . Since eigenfunctions are orthogonal and satisfy the boundary condition an integration by parts shows

$$(7.2) \quad \int_0^\infty wf(\cdot, k_n)\overline{f(\cdot, k_m)} = -f(0, k_n)\overline{f(0, k_m)} \cot \alpha$$

if $n \neq m$. If $\alpha = 0$, the right-hand side has to be replaced by zero.

Since $\int_0^\infty t|q(t) - q_0w(t)|dt < \infty$, our construction of f shows that $f(x, k) \sim e^{ikx}$ as $x \rightarrow \infty$, uniformly for $k \in i[0, 1]$. This shows first that (7.2) is bounded as n and m tend to infinity, and second that we may find a positive c such that $|f(x, k) - e^{ikx}| \leq e^{-|k|x}/4$ if $x \geq c$, $k \in i[0, 1]$. Simple estimates then show that

$$\frac{7}{16}e^{-(|k_n|+|k_m|x)} \leq \text{Re}(f(x, k_n)\overline{f(x, k_m)}) \leq \frac{25}{16}e^{-(|k_n|+|k_m|x)}$$

if n and m are large. Since $w - 1$ is integrable this shows that the integral

$$\int_c^\infty \text{Re}(f(x, k_n)\overline{f(x, k_m)})w \rightarrow +\infty$$

as n, m tend to infinity. Now, since $f(x, k)$ is uniformly continuous on $[0, c] \times i[0, 1]$ it follows that the integral over $[0, c]$ is bounded, so the integral over $[0, \infty)$ tends to infinity, contradicting the previously established boundedness and proving (3).

Finally, if $q_0 = 0$, we already know q_0 is an eigenvalue if and only if $\alpha = 0$. On the other hand, if $q_0 > 0$, then $f(\cdot, q_0)$ is asymptotic to 1, and any other solution to (1.3) is asymptotically linear, as is easily seen from the well-known reduction of order method. Thus no such solution is in \mathcal{H} and there is no eigenfunction with eigenvalue q_0 . This proves (4). \square

Remark 7.3. If we allow w to be a general measure, then the negative part of w could be a finite sum of Dirac measures. In this case one may in the same way show that the number of negative eigenvalues is equal to the number of these Dirac measures if $\alpha \neq 0$, $\cot \alpha \geq 0$, and $q_0 > 0$, with suitable modifications in the other cases.

8. Appendix. Here we present some calculations which were omitted from the proof of Theorem 4.4.

For $x = 0$ the relation $\psi(x, \lambda) = f(x)\check{\psi}(g(x), \lambda)$ gives

$$(8.1) \quad \cos \alpha - m(\lambda) \sin \alpha = f(0)\{\cos \check{\alpha} - \check{m}(\lambda) \sin \check{\alpha}\},$$

while $\psi'(x, \lambda) = f'(x)\check{\psi}(g(x), \lambda) + f(x)g'(x)\check{\psi}'(g(x), \lambda)$ for $x = 0$ gives

$$(8.2) \quad \sin \alpha + m(\lambda) \cos \alpha = \frac{f'(0)}{\lambda}\{\cos \check{\alpha} - \check{m}(\lambda) \sin \check{\alpha}\} + \frac{C}{f(0)}\{\sin \check{\alpha} + \check{m}(\lambda) \cos \check{\alpha}\}.$$

From (8.1), (8.2) we obtain

$$1 = \left\{ f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha \right\} \{\cos \check{\alpha} - \check{m}(\lambda) \sin \check{\alpha}\} + \frac{C \sin \alpha}{f(0)} \{\sin \check{\alpha} + \check{m}(\lambda) \cos \check{\alpha}\}$$

and

$$m(\lambda) = \left\{ -f(0) \sin \alpha + \frac{f'(0)}{\lambda} \cos \alpha \right\} \{\cos \check{\alpha} - \check{m}(\lambda) \sin \check{\alpha}\} + \frac{C \cos \alpha}{f(0)} \{\sin \check{\alpha} + \check{m}(\lambda) \cos \check{\alpha}\},$$

which after rearranging gives

$$(8.3) \quad 1 - \left(f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha \right) \cos \check{\alpha} - \frac{C \sin \alpha \sin \check{\alpha}}{f(0)} = \check{m}(\lambda) \left\{ - \left(f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha \right) \sin \check{\alpha} + \frac{C \sin \alpha \cos \check{\alpha}}{f(0)} \right\}$$

and

$$(8.4) \quad \left(f(0) \sin \alpha - \frac{f'(0)}{\lambda} \cos \alpha \right) \cos \check{\alpha} - \frac{C \cos \alpha \sin \check{\alpha}}{f(0)} = \check{m}(\lambda) \left\{ \left(f(0) \sin \alpha - \frac{f'(0)}{\lambda} \cos \alpha \right) \sin \check{\alpha} + \frac{C \cos \alpha \cos \check{\alpha}}{f(0)} \right\} - m(\lambda).$$

In (8.3) the left member and the coefficient of \check{m} are linear in $1/\lambda$, while $\check{m}(\lambda)$ is not constant or a Möbius transform (this would give a one-dimensional transform space). From (8.3) we therefore obtain

$$\begin{aligned} \left(f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha \right) \cos \check{\alpha} &= 1 - \frac{C \sin \alpha \sin \check{\alpha}}{f(0)}, \\ \left(f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha \right) \sin \check{\alpha} &= \frac{C \sin \alpha \cos \check{\alpha}}{f(0)}, \end{aligned}$$

which gives

$$f(0) \cos \alpha + \frac{f'(0)}{\lambda} \sin \alpha = \cos \check{\alpha},$$

$$\frac{C \sin \alpha}{f(0)} = \sin \check{\alpha}.$$

From this it is (again) clear that $\sin \alpha = 0$ if and only if $\sin \check{\alpha} = 0$, so that we have two cases.

- $\alpha = \check{\alpha} = 0$. We obtain $f(0) = 1$, and insertion in (8.4) shows that $\frac{f'(0)}{\lambda} = m(\lambda) - C\check{m}(\lambda)$. The right member is $(1-C)m(\lambda)$ since $m(i\nu)$ and $\check{m}(i\nu) \rightarrow 0$ as $\nu \rightarrow +\infty$ by Lemma 4.10, and m, \check{m} have the same spectral measure. Again by Lemma 4.10 it follows that $C = 1$, and thus $f'(0) = 0$.
- $\alpha \neq 0 \neq \check{\alpha}$. We obtain $f'(0) = 0$, $f(0) = C \sin \alpha / \sin \check{\alpha}$, and $C \sin(2\alpha) = \sin(2\check{\alpha})$. But we know that $f(0) = \sin \alpha / \sin \check{\alpha}$ so that $C = 1$. Insertion in (8.4) gives $m(\lambda) - \check{m}(\lambda) = \cot \alpha - \cot \check{\alpha}$. Since $\sin(2\alpha) = \sin(2\check{\alpha})$ we have either $\alpha = \check{\alpha}$ or $0 < \alpha = \pi/2 - \check{\alpha} < \pi/2$ or $\pi/2 < \alpha = 3\pi/2 - \check{\alpha} < \pi$. If $\alpha = \check{\alpha}$, we obtain $f(0) = 1$ and $m(\lambda) = \check{m}(\lambda)$. In the other cases we obtain $f(0) = |\tan \alpha|$ and $m(\lambda) - \check{m}(\lambda) = 2 \cot(2\alpha)$.

REFERENCES

- [1] C. BENNEWITZ, *Spectral theory for pairs of differential operators*, Ark. Mat., 15 (1977), pp. 33–61.
- [2] C. BENNEWITZ, *A generalisation of Niessen's limit-circle criterion*, Proc. Roy. Soc. Edinburgh Sect. A, 78 (1977/78), pp. 81–90.
- [3] C. BENNEWITZ, *Spectral asymptotics for Sturm-Liouville equations*, Proc. London Math. Soc. (3), 59 (1989), pp. 294–338.
- [4] C. BENNEWITZ, *A Paley-Wiener theorem with applications to inverse spectral theory*, in Advances in Differential Equations and Mathematical Physics (Birmingham, AL, 2002), Contemp. Math. 327, AMS, Providence, RI, 2003, pp. 21–31.
- [5] C. BENNEWITZ AND B. M. BROWN, *A limit point criterion with applications to nonselfadjoint equations*, J. Comput. Appl. Math., 148 (2002), pp. 257–265.
- [6] C. BENNEWITZ AND W. N. EVERITT, *The Titchmarsh-Weyl eigenfunction expansion theorem for Sturm-Liouville differential equations*, in Sturm-Liouville Theory, Birkhäuser, Basel, 2005, pp. 137–171.
- [7] P. A. BINDING, P. J. BROWNE, AND B. A. WATSON, *Inverse spectral problems for left-definite Sturm-Liouville equations with indefinite weight*, J. Math. Anal. Appl., 271 (2002), pp. 383–408.
- [8] A. BOUTET DE MONVEL AND D. SHEPELSKY, *The Camassa-Holm equation on the half-line*, C. R. Math. Acad. Sci. Paris, 341 (2005), pp. 611–616.
- [9] A. BOUTET DE MONVEL AND D. SHEPELSKY, *The Camassa-Holm equation on the half-line: A Riemann-Hilbert approach*, J. Geom. Anal., 18 (2008), pp. 285–323.
- [10] A. BRESSAN AND A. CONSTANTIN, *Global conservative solutions of the Camassa-Holm equation*, Arch. Ration. Mech. Anal., 183 (2007), pp. 215–239.
- [11] R. CAMASSA AND D. D. HOLM, *An integrable shallow water equation with peaked solitons*, Phys. Rev. Lett., 71 (1993), pp. 1661–1664.
- [12] A. CONSTANTIN, *On the inverse spectral problem for the Camassa-Holm equation*, J. Funct. Anal., 155 (1998), pp. 352–363.
- [13] A. CONSTANTIN, *Existence of permanent and breaking waves for a shallow water equation: A geometric approach*, Ann. Inst. Fourier (Grenoble), 50 (2000), pp. 321–362.
- [14] A. CONSTANTIN, *On the scattering problem for the Camassa-Holm equation*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 457 (2001), pp. 953–970.
- [15] A. CONSTANTIN AND J. ESCHER, *Wave breaking for nonlinear nonlocal shallow water equations*, Acta Math., 181 (1998), pp. 229–243.
- [16] A. CONSTANTIN, V. S. GERDJKOV, AND R. I. IVANOV, *Inverse scattering transform for the Camassa-Holm equation*, Inverse Problems, 22 (2006), pp. 2197–2207.

- [17] A. CONSTANTIN AND J. LENELLS, *On the inverse scattering approach to the Camassa-Holm equation*, J. Nonlinear Math. Phys., 10 (2003), pp. 252–255.
- [18] A. CONSTANTIN AND H. P. MCKEAN, *A shallow water equation on the circle*, Comm. Pure Appl. Math., 52 (1999), pp. 949–982.
- [19] P. DEIFT AND E. TRUBOWITZ, *Inverse scattering on the line*, Comm. Pure Appl. Math., 32 (1979), pp. 121–251.
- [20] A. S. FOKAS, *On a class of physically important integrable equations*, Phys. D, 87 (1995), pp. 145–150.
- [21] A. S. FOKAS AND Q. M. LIU, *Asymptotic integrability of water waves*, Phys. Rev. Lett., 77 (1996), pp. 2347–2351.
- [22] B. FUCHSSTEINER AND A. S. FOKAS, *Symplectic structures, their Bäcklund transformations and hereditary symmetries*, Phys. D, 4 (1981/82), pp. 47–66.
- [23] R. S. JOHNSON, *Camassa-Holm, Korteweg-de Vries and related models for water waves*, J. Fluid Mech., 455 (2002), pp. 63–82.
- [24] Q. KONG, H. WU, AND A. ZETTL, *Singular left-definite Sturm-Liouville problems*, J. Differential Equations, 206 (2004), pp. 1–29.
- [25] V. A. MARCHENKO, *Sturm-Liouville Operators and Applications*, Oper. Theory Adv. Appl. 22, Birkhäuser Verlag, Basel, 1986.
- [26] H. D. NIESSEN AND A. SCHNEIDER, *Spectral theory for left-definite singular systems of differential equations*, in Spectral Theory and Asymptotics of Differential Equations (Proc. Conf., Scheveningen, 1973), North-Holland Math. Stud. 13, North-Holland, Amsterdam, 1974, pp. 29–43.
- [27] K. L. VANINSKY, *Equations of Camassa-Holm type and Jacobi ellipsoidal coordinates*, Comm. Pure Appl. Math., 58 (2005), pp. 1149–1187.
- [28] H. WEYL, *Über gewöhnliche lineare Differentialgleichungen mit singulären Stellen und ihre Eigenfunktionen (2. Note)*, Gött. Nachr., (1910), pp. 442–467.

ABSOLUTELY CONTINUOUS LAWS OF JUMP-DIFFUSIONS IN FINITE AND INFINITE DIMENSIONS WITH APPLICATIONS TO MATHEMATICAL FINANCE*

BARBARA FORSTER[†], EVA LÜTKEBOHMERT[‡], AND JOSEF TEICHMANN[†]

Abstract. In mathematical Finance calculating the Greeks by Malliavin weights has proved to be a numerically satisfactory procedure for finite-dimensional Itô-diffusions. The existence of Malliavin weights relies on absolute continuity of laws of the projected diffusion process and a sufficiently regular density. In this article we first prove results on absolute continuity for laws of projected jump-diffusion processes in finite and infinite dimensions and a general result on the existence of Malliavin weights in finite dimension. In both cases we assume Hörmander conditions and hypotheses on the invertibility of the so-called linkage operators. The purpose of this article is to show that for the construction of numerical procedures for the calculation of the Greeks in fairly general jump-diffusion cases one can proceed as in a pure diffusion case. We also show how the given results apply to infinite-dimensional questions in mathematical Finance. There we start from the Vasiček model, and add—by pertaining no arbitrage—a jump-diffusion component. We prove that we can obtain in this case an interest rate model, where the law of any projection is absolutely continuous with respect to Lebesgue measure on \mathbb{R}^M .

Key words. Malliavin calculus, compound Poisson process, Hörmander condition, Greeks, Malliavin weight, stochastic partial differential equation, jump-diffusion, interest rate theory

AMS subject classifications. 60H07, 60H15, 62P05

DOI. 10.1137/070708822

1. Introduction. We shall consider in this article the question of whether the law of $l(X_t^x)$, for a finite-dimensional projection $l : H \rightarrow \mathbb{R}^M$, is absolutely continuous with respect to Lebesgue measure on \mathbb{R}^M , where X_t^x is the solution of the stochastic (partial) differential equation (SPDE)

$$(1.1) \quad dX_t^x = (AX_{t-}^x + \alpha(X_{t-}^x))dt + \sum_{i=1}^d V_i(X_{t-}^x)dB_t^i + \sum_{j=1}^m \delta_j(X_{t-}^x)dL_t^j,$$

$$(1.2) \quad X_0^x = x \in H$$

and H is a possibly infinite-dimensional separable Hilbert space. We refer to the previous equation loosely speaking as a jump-diffusion on the Hilbert space H , pointing out that the involved Lévy processes are of finite type. For sake of simplicity we shall always work with the cadlag integrand $t \mapsto X_{t-}^x$, even though for the dt and dB_t integrals this is superfluous. In the infinite-dimensional setting we are not aware of results on absolute continuity of the projected process in the jump-diffusion case.

*Received by the editors November 20, 2007; accepted for publication (in revised form) September 23, 2008; published electronically January 23, 2009.

<http://www.siam.org/journals/sima/40-5/70882.html>

[†]Department of Mathematical Methods in Economics, Vienna University of Technology, Research Group e105 Financial and Actuarial Mathematics, Wiedner Hauptstrasse 8-10, A-1040 Wien, Austria (bforster@fam.tuwien.ac.at, jteichma@fam.tuwien.ac.at). The first author acknowledges the support from the FWF-project P15889 “Utility Maximization in Incomplete Financial Markets” and from the FWF-Wissenschaftskolleg W 8 “Differential Equation Models in Science and Engineering.” The third author acknowledges the support from the RTN network HPRN-CT-2002-00281 and from the FWF-grant Y328.

[‡]Institute of Social Sciences and Economics, University of Bonn, Adenauerallee 24-42, D-53113 Bonn, Germany (eva.luetkebohmert@uni-bonn.de).

Related work has been done in [5] for the construction of first variation processes. In the diffusion case we refer the reader to the work [3] and the references therein and in particular to the recently published inspiring results of Jonathan Mattingly; see, for instance [1]. We point out that we deal here with SPDEs and stochastic differential equations (SDEs) at the same time, where the latter case appears in this setting when the state (Hilbert) space H is finite-dimensional.

Very satisfying results in the finite-dimensional setting with Lévy processes of infinite type have been obtained in [11] through a generalization of the Norris lemma to Lévy processes. These results have been built upon our results presented in this work for the finite activity case (see section 7). Substantial work with respect to absolute continuity and smoothness of the density has already been published in the 1980s, the most prominent being [7] and [6]. Therein several questions of extension of hypoellipticity results (and Malliavin Calculus) to jump-processes are discussed and completely solved; however, the problem of a hypoelliptic diffusion together with a finite-activity jump-structure remained open. It has to be pointed out here that—in contrast to [6]—we do not need any extension of Malliavin Calculus to jump-processes for our results (see also the discussion in Remark 5). This gap was filled by the announced results of [27], but several proofs therein are extremely short. Recently, motivated by questions from financial mathematics (see section 8 for an outline of the problem), there has been increasing interest in those results; see the works [2], [13], [14], and the references therein. This article aims to work out the most general finite activity case under Hörmander conditions on the diffusion part.

From the point of view of existence and uniqueness for jump-diffusions in infinite dimensions, our main reference is [15] and the references therein. Since we consider jump-diffusions as concatenated diffusions on Hilbert space, we mention [12] as the main reference for existence and uniqueness results but also [4] and [10] for many interesting constructions and ideas.

There are two applications added to this work. The first is the Heath–Jarrow–Morton (HJM) equation (as presented in [15]), where we show that the innocent Vasiček model (see, for instance, the seminal work [9]) with a certain jump-structure triggered by a one-dimensional Poisson process yields, under no-arbitrage assumptions, a model where not only do no finite-dimensional realizations exist but also where every projection into a finite-dimensional subspace admits a density (compare with the notion of generic interest rate evolutions from [3]). The second application is concerned with concrete formulas for the calculation of Malliavin weights. There our message is that one can think Poisson-trajectory-wise; i.e., the results from [18] or [17] can be literally applied by replacing the diffusion process by the respective jump-diffusion process.

When we analyze jump-diffusions with values in Hilbert spaces, loosely speaking the following facts hold true:

- Between two consecutive jumps of the jump-diffusion we are given an ordinary diffusion.
- At a jump we add to the left limit the jump size (which usually depends on the left limit, too). In [25, Chapter V.10, Hypothesis (H3)], this operation is formalized by the so-called *linkage operators* $x \mapsto x + \mu\delta^j(x)$, which encode what happens at a jump of size μ at x . We shall apply this notion here, too.

Hence the following picture arises:

- In order to obtain absolute continuity of the projected diffusion process, we need the Hörmander condition to be in force. Otherwise we cannot expect—

conditioned on the event of positive probability that no jump occurs—that the law of $1(X_t^x)$ is absolutely continuous with respect to Lebesgue measure.

- In order to preserve absolute continuity we need the linkage operators to be invertible in a proper sense.

Remark 1. Both conditions are “sine qua non,” since it is easy to imagine counterexamples.

In section 2 we fix the general setting of this article. We shall deal with Lévy processes of finite type as drivers of the SDEs, even though we believe that one should be able to prove similar results in the case of many small jumps, too. We also state the main assumptions of this work in section 2 for later use.

In section 3 we state a “folklore” decomposition theorem, which tells that solving a jump-diffusion S(P)DE is the same as solving associated diffusion S(P)DEs and concatenating the solutions by linkage operators at jumps. In section 4 we show that we can also prove results on first variation processes in the spirit of the decomposition theorem. We prove that under our analytic requirements there is in fact a sufficiently regular first variation process.

In section 5 we show by means of Malliavin calculus for a d -dimensional Brownian motion that the law of a projected jump-diffusion is absolutely continuous with respect to Lebesgue measure. In section 6 we introduce a class of examples from mathematical Finance, where we see very directly the phenomenon of absolute continuity arising from the introduction of jumps and the no-arbitrage condition. This shows again that finite-dimensional realizations, as constructed, for instance, in [8], are a rare case in infinite dimensions. In section 7 we restrict our attention to the finite-dimensional setting to show that the density of the absolutely continuous law is in fact smooth by proving that the inverse of the covariance matrix has p th moments for all $p \geq 1$. In section 8 we apply the invertibility of the covariance matrix to the calculation of Greeks. The appendix shows an important estimate implicitly present in the Norris lemma as presented in Nualart’s book [24]. A similar result (which could be directly used in section 7 for the proof of the main theorem) can be found in [21, Corollary 3.25]. The article [21] is most likely the source of the first appearance of the precise polynomial time-dependence in the estimate of the L^p -norm of the inverse of the covariance matrix. We have been choosing here the path via the Norris lemma; we explain the estimate by redoing its proof in the appendix.

2. Setting and assumptions. Let $(\Omega, \mathcal{F}, P, (\mathcal{F}_t)_{t \geq 0})$ be a filtered probability space where the filtration $(\mathcal{F}_t)_{t \geq 0}$ satisfies the usual conditions. Let $(B_t)_{t \geq 0}$ be a d -dimensional Brownian motion and $(L_t^j)_{t \geq 0}$, $j = 1, \dots, m$, be m independent compound Poisson processes given by

$$L_t^j := \sum_{k=1}^{N_t^j} Z_k^j,$$

where N_t^j denotes a Poisson process with jump intensity $\widetilde{\lambda}_j > 0$ and $Z^j = (Z_k^j)_{k \geq 1}$ is an independently and identically distributed sequence of random variables with distribution μ_j for $j = 1, \dots, m$ such that each μ_j admits all moments. The compensated compound Poisson process reads as

$$L_t^j - E(L_t^j) = L_t^j - \lambda_j t,$$

where $\lambda_j = E(Z_1^j) \widetilde{\lambda}_j$ is the average jump size times the jump rate.

We could equally take an \mathbb{R}^m -valued Lévy process of finite type, i.e., introduce a dependence structure between the jumps of the components, and all theorems would equally hold true with slightly modified proofs, but we believe that this generalization does not bring further insight.

We assume that all sources of randomness are mutually independent and that the filtration $(\mathcal{F}_t)_{t \geq 0}$ is the natural filtration with respect to $(B_t, L_t^1, \dots, L_t^m)_{t \geq 0}$. Let H be a separable Hilbert space. We fix furthermore a strongly continuous semigroup S on H with generator A . Let α, V_1, \dots, V_d , the diffusion vector fields, and $\delta_1, \dots, \delta_m$, the jump vector fields, be C^∞ -bounded on H ; that is, the vector fields are infinitely often differentiable with bounded partial derivatives of all proper orders $n \geq 1$. We consider the mild cadlag solution $(X_t^x)_{t \geq 0}$ of an SDE

$$(2.1) \quad dX_t^x = (AX_{t-}^x + \alpha(X_{t-}^x))dt + \sum_{i=1}^d V_i(X_{t-}^x)dB_t^i + \sum_{j=1}^m \delta_j(X_{t-}^x)dL_t^j,$$

$$(2.2) \quad X_0^x = x \in H.$$

See [15] for all necessary details on existence and uniqueness of the previous equation.

The previous conditions are slightly more than standard for existence and uniqueness of mild solutions; i.e., in [15] the authors need Lipschitz conditions on the vector fields, whereas we assume them to be C^∞ -bounded. In order to speak about absolute continuity of projections to \mathbb{R}^M we shall need more assumptions; in particular, for conclusions drawn from the geometry of the given vector fields $\alpha, V_1, \dots, V_d, \delta_1, \dots, \delta_m$ several quite strong analytic requirements are necessary. We group the assumptions into three groups and indicate in each section which assumptions we shall need.

Let $\mathbf{l} : H \rightarrow \mathbb{R}^M$ be a projection; then we want to know whether the law of $\mathbf{l}(X_t^x)$ is absolutely continuous with respect to Lebesgue measure and if the density is smooth. Following the short discussion in the introduction, we need the Hörmander conditions to be in force, and we need to suppose invertibility on linkage operators.

We apply the following notation for Hilbert spaces $\text{dom}(A^k)$:

$$\begin{aligned} \text{dom}(A^k) &:= \{h \in H \mid h \in \text{dom}(A^{k-1}) \text{ and } A^{k-1}h \in \text{dom}(A)\}, \\ \|h\|_{\text{dom}(A^k)}^2 &:= \sum_{i=0}^k \|A^i h\|^2, \\ \text{dom}(A^\infty) &= \bigcap_{k \geq 0} \text{dom}(A^k), \end{aligned}$$

which we need in order to specify the analytic conditions.

ASSUMPTION 1. *We assume that the generator A of S generates a strongly continuous group. We assume furthermore that α, V_1, \dots, V_d , the diffusion vector fields, and $\delta_1, \dots, \delta_m$, the jump vector fields, are C^∞ -bounded on the Hilbert spaces $\text{dom}(A^k)$ for $k \geq 0$; that is, the vector fields are infinitely often differentiable with bounded partial derivatives of all proper orders $n \geq 1$ on the Hilbert space $\text{dom}(A^k)$ for $k \geq 0$.*

ASSUMPTION 2. *We take Assumption 1 for granted; i.e., we can consider all vector fields on the space $\text{dom}(A^k)$ for $k = 0, \dots, \infty$. For a proper statement of the Hörmander condition we apply the “geometrically relevant” drift*

$$V_0(x) = Ax + \alpha(x) - \frac{1}{2} \sum_{i=1}^d \text{T}V_i(x) \cdot V_i(x)$$

for $x \in \text{dom}(A)$ and call V_0 the Stratonovich drift of the diffusion. Recall the (tangent) directional derivative operator \mathbb{T} defined through

$$\mathbb{T}V(x) \cdot v = \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} V(x + \epsilon v).$$

Lie brackets can only be calculated on the Fréchet space $\text{dom}(A^\infty)$ and there we formulate the Hörmander condition. We assume that the distribution $\mathcal{D}(x)$ generated by the vector fields

$$(2.3) \quad \begin{aligned} &V_1(x), \dots, V_d(x), \quad [V_i(x), V_j(x)] \quad (i, j = 0, 1, \dots, d), \\ &[V_i(x), [V_j(x), V_k(x)]] \quad (i, j, k = 0, 1, \dots, d), \dots \end{aligned}$$

is dense in H for one $x \in \text{dom}(A^\infty)$.

ASSUMPTION 3. We assume that the inverse of $x \mapsto x + z\delta_j(x)$ exists and is C^∞ -bounded on each $\text{dom}(A^k)$ for $z \in \text{supp}(\mu_j)$, $j = 1, \dots, m$, and $k \geq 0$ (recall that μ_j was the distribution of the random variable Z_j).

Remark 2. As far as Assumption 1 is concerned, we do believe that the assertions of this paper also hold true for (most) strongly continuous semigroups. A proof based on an application of the Szökefalvi–Nagy theorem can be found in the recent preprint [26]; therefore, we could replace the assumption that A generates a strongly continuous group by the assumption that A generates a pseudocontractive, strongly continuous semigroup. This includes most of the second order partial differential operators. However, for this paper we do always assume the group property for the sake of simplicity.

Remark 3. The Hörmander condition could not be formulated without the analytic part of Assumption 1.

Example 1. In order to show examples of vector fields, which are C^∞ -bounded on $\text{dom}(A^k)$, consider the following structure. Let H be a separable Hilbert space and A the generator of a strongly continuous semigroup. We know that $\text{dom}(A^\infty)$ is a Fréchet space and an injective limit of the Hilbert spaces $\text{dom}(A^k)$ for $k \geq 0$. Following the analysis as developed in [16] (see also [20] and [19], where the analytic concepts have been originally developed), we can consider vector fields $V : U \subset H \rightarrow \text{dom}(A^\infty)$. If V is smooth in the sense explained in [16] and has the property that its derivatives of proper order $n \geq 1$ are bounded on $U \subset H$, then V is obviously a C^∞ -bounded vector field, and additionally $V|_{\text{dom}(A^\infty)}$ is a Banach map-vector field in the sense of [16]. Such vector fields constitute a class, where Assumptions 1–3 can be readily checked.

3. Decomposition theorem for jump-diffusions on Hilbert spaces. In order to properly understand how to apply the Malliavin calculus, we state the following rather obvious structure theorem on jump-diffusions, which simply takes into account that stochastic integration with respect to the Poisson process follows the rules of Lebesgue–Stieltjes integration (see, for instance, [25] for a general exposition). Here we need only that the vector fields are C^∞ -bounded on H in order to guarantee existence and uniqueness of the respective equations.

THEOREM 1. Let $(\Omega, \mathcal{F}, P, (\mathcal{F}_t)_{t \geq 0})$ be a filtered probability space, $(B_t)_{t \geq 0}$ be a d -dimensional Brownian motion, and $(L_t^j)_{t \geq 0}$ be m independent compound Poisson processes for $j = 1, \dots, m$, such that the filtration is the natural filtration with respect to $(B_t, L_t^1, \dots, L_t^m)_{t \geq 0}$. Let S be a strongly continuous semigroup with generator A on H . Let α, V_1, \dots, V_d , the diffusion vector fields, and $\delta_1, \dots, \delta_m$, the jump vector fields,

be C^∞ -bounded on H , and consider the cadlag solution $(X_t^x)_{0 \leq t \leq T}$ of an SPDE

$$(3.1) \quad dX_t^x = (AX_{t^-}^x + \alpha(X_{t^-}^x))dt + \sum_{i=1}^d V_i(X_{t^-}^x)dB_t^i + \sum_{j=1}^m \delta_j(X_{t^-}^x)dL_t^j,$$

$$(3.2) \quad X_0^x = x.$$

Let η denote a piecewise constant cadlag trajectory $\eta : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^m$ of the compound Poisson process L with finitely many jumps on compact intervals and starting at 0. We consider

$$(3.3) \quad dY_{s,t}^{x,\eta} = \left(AY_{s,t^-}^{x,\eta} + \alpha(Y_{s,t^-}^{x,\eta}) \right) dt + \sum_{i=1}^d V_i(Y_{s,t^-}^{x,\eta}) dB_t^i + \sum_{j=1}^m \delta_j(Y_{s,t^-}^{x,\eta}) d\eta^j(t),$$

$$(3.4) \quad Y_{s,s}^{x,\eta} = x.$$

Then $(Y_{s,t}^{x,\eta})_{s \geq 0}$ can be given explicitly in terms of the jump times τ_n of η for $n \geq 0$ and the diffusion process between two consecutive jumps:

$$\begin{aligned} Y_t^{x,\eta} &:= Y_{0,t}^{x,\eta} \text{ for } 0 \leq t < \tau_1, \\ Y_t^{x,\eta} &:= Y_{\tau_1,t}^{y,\eta} \Big|_{y=Y_{0,\tau_1}^{x,\eta} + \sum_{j=1}^m \delta_j(Y_{0,\tau_1}^{x,\eta}) \Delta \eta^j(\tau_1)} \text{ for } \tau_1 \leq t < \tau_2 \\ &\vdots \\ Y_t^{x,\eta} &:= Y_{\tau_{n-1},t}^{y,\eta} \Big|_{y=Y_{0,\tau_{n-1}}^{x,\eta} + \sum_{j=1}^m \delta_j(Y_{0,\tau_{n-1}}^{x,\eta}) \Delta \eta^j(\tau_{n-1})} \text{ for } \tau_{n-1} \leq t < \tau_n. \end{aligned}$$

Here we write $\Delta \eta(t) := \eta(t) - \eta(t^-)$ for $t \geq 0$. We define the process $(Y_t^{x,L})_{t \geq 0}$ by inserting the compound Poisson process L for η into $(Y_t^{x,\eta})_{t \geq 0}$. The resulting process $(Y_t^{x,L})_{t \geq 0}$ is then indistinguishable from $(X_t^x)_{t \geq 0}$.

Proof. For the proof we refer the reader to [25, Chapter V.10, Theorem 57], particularly with respect to the conditioning on the jump part. The proof remains unchanged in the infinite-dimensional setting; see [15] for the existence and uniqueness proof on separable Hilbert spaces. \square

Remark 4. For future use we shall always assume that the first jumping time of η is strictly positive, $\tau_1 > 0$, and that each time corresponds to the jump of exactly one coordinate process L^j , which is true for almost all trajectories of the compound Poisson process L . Notice that the dependence of $(Y_t^{x,\eta})_{t \geq 0}$ on the jump times of η is continuous but certainly not smooth since the jump times are inserted instead of the time of a hypoelliptic diffusion process.

Remark 5. Notice that one can also interpret the result in the following way: consider the solution $(X_t^x)_{t \geq 0}$ of (3.1) as an element of $L^2(\Omega_1 \times \Omega_2; H)$, where Ω_1 carries the Brownian motion part (with natural filtration), Ω_2 carries the Poisson part (with natural filtration), and $\Omega_1 \times \Omega_2$ is equipped with the respective product σ -algebra. Then we know by Fubini's theorem that

$$L^2(\Omega_1 \times \Omega_2; H) = L^2(\Omega_2; L^2(\Omega_1; H)).$$

The previous theorem only clarifies the jump-diffusion structure of the dependence on Ω_2 . In other words, between jumps we have ordinary diffusions, and at a jump we link by linkage operators.

4. First variation processes. In order to calculate Malliavin derivatives, which is crucial for arguments on absolute continuity, we need precise statements on first variation processes of jump-diffusions. For later purposes, but also in order to see results on the inverse of the first variation process easily, we write our equations in the Stratonovich notation. This is not innocent in infinite dimensions, since mild solutions are in general *not* semimartingales, and therefore the Stratonovich notation fails to be applicable in general. However, by Assumption 1 we are able to determine whether we are given a semimartingale, or not, by analyzing the initial value of the process. Indeed, for fixed $k \geq 0$, if $x \in \text{dom}(A^{k+1})$, then there is a mild solution taking values in $\text{dom}(A^{k+1})$ of the Itô SDE. However, this solution process has to coincide, by uniqueness, with the solution process obtained by considering the same equation on $\text{dom}(A^k)$ with an initial value in $\text{dom}(A^{k+1})$. Therefore, the mild solution in $\text{dom}(A^{k+1})$ is a strong solution in $\text{dom}(A^k)$. Therefore, we assume Assumptions 1 and 3 to be in force in this section.

By the previous arguments for the given SDE (3.1), we can switch to Stratonovich notation for $x \in \text{dom}(A)$ and obtain

$$dX_t^x = V_0(X_{t-}^x)dt + \sum_{i=1}^d V_i(X_{t-}^x) \circ dB_t^i + \sum_{j=1}^m \delta_j(X_{t-}^x) dL_t^j$$

with the Stratonovich drift given by

$$V_0(x) := Ax + \alpha(x) - \frac{1}{2} \sum_{i=1}^d \mathbb{T} V_i(x) \cdot V_i(x)$$

for $x \in \text{dom}(A)$. Recall the tangent (derivative) operator \mathbb{T} :

$$\mathbb{T} V(x) \cdot v = \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} V(x + \epsilon v).$$

We do also consider the SDE (3.3) with respect to one trajectory η and switch to Stratonovich notation there, too. The following theorem states the result on the first variation process along one trajectory η , which yields in what follows the same result by inserting the compound Poisson process L for η . Notice that the trajectory η is such that the first jumping time is strictly positive and that at each jumping time τ_n for $n \geq 1$ only one coordinate jumps.

THEOREM 2. *Assume Assumptions 1 and 3 hold. We fix $k \geq 0$. The first variation process $(J_{s \rightarrow t}(x, \eta))_{t \geq s}$ associated with $(Y_t^{x, \eta})_{t \geq 0}$ on $\text{dom}(A^k)$ is well defined and satisfies the SDE*

$$\begin{aligned} dJ_{s \rightarrow t}(x, \eta) \cdot h &= \left(AJ_{s \rightarrow t-}(x, \eta) \cdot h + \mathbb{T} \alpha(Y_{s, t-}^{x, \eta}) \cdot J_{s \rightarrow t-}(x, \eta) \cdot h \right) dt \\ &+ \sum_{i=1}^d \left(\mathbb{T} V_i(Y_{s, t-}^{x, \eta}) \cdot J_{s \rightarrow t-}(x, \eta) \cdot h \right) dB_t^i \\ &+ \sum_{j=1}^m \left(\mathbb{T} \delta_j(Y_{s, t-}^{x, \eta}) \cdot J_{s \rightarrow t-}(x, \eta) \cdot h \right) d\eta^j(t), \end{aligned}$$

$$(4.1) \quad J_{s \rightarrow s}(x, \eta) \cdot h = h$$

for $h, x \in \text{dom}(A^k)$ and $t \geq s$. The Itô equation has a unique global mild solution for $h, x \in \text{dom}(A^k)$, and $J_{s \rightarrow t}(x, \eta)$ defines a continuous linear operator on $\text{dom}(A^k)$, which is invertible if $x \in \text{dom}(A^{k+1})$.

The Stratonovich equation on $\text{dom}(A^k)$ in turn is well defined only for $h, x \in \text{dom}(A^{k+1})$. We apply the (formal) notation here:

$$\mathbb{T} V_0(x)v = Av + \mathbb{T} \alpha(x)v - \frac{1}{2} \mathbb{T} \left(x \mapsto \sum_{i=1}^d \mathbb{T} V_i(x) \cdot V_i(x) \right) v$$

for $x \in \text{dom}(A)$ and $v \in \text{dom}(A)$.

$$\begin{aligned} dJ_{s \rightarrow t}(x, \eta) \cdot h &= \left(\mathbb{T} V_0(Y_{s,t^-}^{x,\eta}) \cdot J_{s \rightarrow t^-}(x, \eta) \cdot h \right) dt \\ (4.2) \quad &+ \sum_{i=1}^d \left(\mathbb{T} V_i(Y_{s,t^-}^{x,\eta}) \cdot J_{s \rightarrow t^-}(x, \eta) \cdot h \right) \circ dB_t^i \\ &+ \sum_{j=1}^m \left(\mathbb{T} \delta_j(Y_{s,t^-}^{x,\eta}) \cdot J_{s \rightarrow t^-}(x, \eta) \cdot h \right) d\eta^j(t), \\ J_{s \rightarrow s}(x, \eta) \cdot h &= h \end{aligned}$$

for $h, x \in \text{dom}(A^{k+1})$ and $t \geq s$. The adjoint of the inverse

$$Z_t^{x,h} := (J_{s \rightarrow t}(x, \eta)^{-1})^* \cdot h,$$

if it exists, should satisfy the following Stratonovich equation at the point x in direction h :

$$\begin{aligned} dZ_t^{x,h} &= - \left(\mathbb{T} V_0(Y_{s,t^-}^{x,\eta})^* \cdot Z_{t^-}^{x,h} \right) dt - \sum_{i=1}^d \left(\mathbb{T} V_i(Y_{s,t^-}^{x,\eta})^* \cdot Z_{t^-}^{x,h} \right) \circ dB_t^i \\ (4.3) \quad &- \sum_{j=1}^m \left(\mathbb{T} \delta_j(Y_{s,t^-}^{x,\eta})^* \cdot Z_{t^-}^{x,h} \right) d\eta^j(t) \\ &+ \sum_{j=1}^m \left(\left(\mathbb{T} \delta_j(Y_{s,t^-}^{x,\eta}) \right)^2 \right)^* \cdot \left(\left(id_H + \Delta \eta^j(t) \mathbb{T} \delta_j(Y_{s,t^-}^{x,\eta}) \right)^{-1} \right)^* \\ &\cdot Z_{t^-}^{x,h} (\Delta \eta^j(t))^2 \end{aligned}$$

for $h, x \in \text{dom}(A^{k+1})$ and $t \geq s \geq 0$ (here we applied the notions of [25]).

Remark 6. The completely analogous theorem holds when we replace η by a compound Poisson process L . We do not state this theorem again, but we point out that we even have moment estimates for the respective processes, which is the only additional relevant information. To be precise, the first variation process $J_{s \rightarrow t}(x) \cdot h$, which equals $J_{s \rightarrow t}(x, L)$ by construction, has bounded second moments by [15].

Proof. Under our Assumption 1, the regularity in the initial values is clear by well-known results from [12] and the chain rule on Hilbert spaces (recall that the linkage operators are smooth). We are allowed to pass to the Stratonovich decomposition since we integrate semimartingales by Itô's formula on Hilbert spaces for $x, h \in \text{dom}(A^{k+1})$ due to the arguments of [3]: the core assertion here is that we can replace H by each $\text{dom}(A^k)$ for some $k \geq 0$, which means in turn if we start in $\text{dom}(A^{k+1})$ and obtain a mild solution there, it is indeed a strong solution considered on $\text{dom}(A^k)$ for $k \geq 0$. It remains to show the invertibility results on the respective first variation processes.

Left invertibility of the first variation $J_{s \rightarrow t}(y, \cdot)$ follows by Itô's formula since we have cadlag trajectories with finitely many jumps. Calculating the semimartingale decomposition of $(Z_t^{y,\cdot})^* \cdot J_{0 \rightarrow t}(x, \eta)$ given by (4.2) and (4.3) yields the result

$$(Z_t^{x,\cdot})^* J_{0 \rightarrow t}(x, \eta) = id_{\text{dom}(A^k)}.$$

Thus, the solution of (4.3) is the left inverse of $J_{s \rightarrow t}$.

We prove that the left inverse is also the right inverse by the same reasoning as in the proof of Proposition 2 in [3]. Therefore, we choose an orthonormal basis $(g_i)_{i \geq 1}$ of $\text{dom}(A^k)$ which lies in $\text{dom}(A^{k+1})$. Then we can compute the semimartingale decomposition of

$$\begin{aligned} & \sum_{i=1}^N \langle (Z_t^{x,h_1})^*, g_i \rangle_{\text{dom}(A^k)} \langle g_i, J_{s \rightarrow t}(x, \eta)^* \cdot h_2 \rangle_{\text{dom}(A^k)} \\ &= \sum_{i=1}^N \langle h_1, Z_t^{x,g_i} \rangle_{\text{dom}(A^k)} \langle J_{s \rightarrow t}(x, \eta) \cdot g_i, h_2 \rangle_{\text{dom}(A^k)} \end{aligned}$$

for $h_1, h_2 \in \text{dom}(A^{k+1})$ and $N \geq 1$. Applying the Stratonovich decomposition and by adjoining, we can free the g_i 's and pass to the limit, which yields the vanishing finite variation and martingale part. Hence

$$\begin{aligned} & \langle J_{s \rightarrow t}(x, \eta) (Z_t^{x,h_1})^*, h_2 \rangle_{\text{dom}(A^k)} \\ &= \lim_{N \rightarrow \infty} \sum_{i=1}^N \langle (Z_t^{x,h_1})^*, g_i \rangle_{\text{dom}(A^k)} \langle g_i, J_{s \rightarrow t}(x, \eta)^* \cdot h_2 \rangle_{\text{dom}(A^k)} \\ &= \langle h_1, h_2 \rangle_{\text{dom}(A^k)}, \end{aligned}$$

which is what a right inverse should satisfy. \square

5. Absolutely continuous laws in finite and infinite dimensions. In this section we assume Assumptions 1, 2, and 3. We want to determine by means of Malliavin calculus whether the law of $\mathbf{I}(Y_t^{x,\eta})$ is absolutely continuous with respect to Lebesgue measure for $t > 0$.

For details on Malliavin calculus see [22] and [24], where in particular the derivative operator and the Skorohod integral for Malliavin calculus with respect to a d -dimensional Brownian motion are defined. Notice that we do not need a Malliavin calculus with respect to the Poissonian trajectories, since we calculate Poisson-trajectory-wise.

Our first task is the calculation of the Malliavin derivative for a fixed cadlag path η . In a second step, we consider the composed problem, where we replace η by a compound Poisson process L as outlined before. Therefore, we first fix a piecewise constant cadlag trajectory $\eta : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^m$ of the process $(L_t^1, \dots, L_t^m)_{t \geq 0}$ with finitely many jumps on compact intervals starting at 0.

THEOREM 3. *We take Assumptions 1, 2, and 3 for granted, where $x \in \text{dom}(A^\infty)$ denotes the point where the Hörmander condition (2.3) holds true. Let $(Y_t^x)_{t \geq 0}$ denote the unique cadlag solution of (3.3). Then for projections $\mathbf{1} : H \rightarrow \mathbb{R}^M$ the law of $\mathbf{I}(Y_t^x)$ is absolutely continuous with respect to Lebesgue measure on \mathbb{R}^M for $t > 0$.*

Proof. Fix $t > 0$. We are able to write the Malliavin derivative of Y_t^x for each Poissonian trajectory η :

$$D_s^i(\mathbf{1} \circ Y_t^{x,\eta}) = \mathbf{1} \circ J_{0 \rightarrow t}(x) J_{0 \rightarrow s}(x)^{-1} V_i(Y_{s-}^{x,\eta}) \mathbf{1}_{[0,t]}(s).$$

We can calculate the *Malliavin covariance matrix* γ as

$$\langle \gamma(\mathbf{1} \circ Y_t^{x,\eta})\xi, \xi \rangle := \sum_{i=1}^d \int_0^t \langle \mathbf{1} \circ J_{0 \rightarrow t}(x) J_{0 \rightarrow s}(x)^{-1} V_i(Y_{s^-}^{x,\eta}), \xi \rangle^2 ds.$$

Consequently, the covariance matrix $\gamma(\mathbf{1} \circ Y_t^{x,\eta})$ can be calculated in the usual way via the reduced covariance matrix

$$\langle C_t \xi, \xi \rangle := \sum_{i=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V_i(Y_{s^-}^{x,\eta}), \xi \rangle^2 ds$$

through the relation

$$\gamma(\mathbf{1} \circ Y_t^{x,\eta}) = (\mathbf{1} \circ J_{0 \rightarrow t}(x)) C_t (\mathbf{1} \circ J_{0 \rightarrow t}(x))^*,$$

where $*$ denotes the adjoint operator with respect to the Hilbert space structures on H and \mathbb{R}^M . We assume n jumps of η on $[0, t]$, and we denote by $0 = \tau_0 < \tau_1 < \dots < \tau_n \leq t$ the sequence of jump times of η . For convenience, we denote the last point in time t by τ_{n+1} , even if $\tau_{n+1} = \tau_n$, which can in principle happen. Hence we can decompose:

$$\langle C_t \xi, \xi \rangle := \sum_{k=0}^n \sum_{i=1}^d \int_{\tau_k}^{\tau_{k+1}} \langle J_{0 \rightarrow s}(x)^{-1} V_i(Y_{s^-}^{x,\eta}), \xi \rangle^2 ds = \sum_{k=0}^n \langle C_t^k \xi, \xi \rangle.$$

Each of the summands determines a symmetric matrix C_t^k and can be interpreted as a reduced covariance matrix coming from a diffusion between τ_k and τ_{k+1} with initial value $Y_{\tau_k}^x$ for $k = 0, \dots, n$. We do not know whether the Hörmander condition is true everywhere. Therefore, we do not know whether C_t^k is a positive definite operator for each $k \geq 0$. From [3], Theorem 1, we do know, however, that C_t^0 is a positive definite operator and there exist null sets N_0 such that on N_0^c the matrix C_t^0 is invertible. Hence the law of $(\mathbf{1} \circ Y_t^{x,\eta})$ is absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^M , since $J_{0 \rightarrow t}(x)$ is invertible and therefore $\gamma(\mathbf{1} \circ Y_t^{x,\eta})$ has empty kernel (Theorem 2.1.2 in [24, p. 86]). \square

Remark 7. The same conclusions hold for $Y_t^{x,\eta}$: notice that $Y_t^{x,\eta} = Y_{t^-}^{x,\eta}$ if there is no jump at t . Otherwise, $Y_t^{x,\eta} = Y_{t^-}^{x,\eta} + \sum_{j=1}^{m_t} \delta_j(Y_{0,t^-}^{x,\eta}) \Delta \eta^j(t)$, but invertible diffeomorphisms transform absolutely continuous laws into absolutely continuous ones.

Now we extend this theorem to the jump-diffusion process $(X_t^x)_{t \geq 0}$, which is easy since, conditioned on one trajectory η , we do have an absolutely continuous law and this property is not perturbed by integration due to Fubini's theorem.

THEOREM 4. *We take Assumptions 1, 2, and 3 for granted, where $x \in \text{dom}(A^\infty)$ denotes the point where the Hörmander condition (2.3) holds true. Let $(X_t^x)_{t \geq 0}$ denote the unique cadlag solution of (3.1). Then for projections $\mathbf{1} = (l_1, \dots, l_k) : H \rightarrow \mathbb{R}^M$ the law of $\mathbf{1}(X_t^x)$ is absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^M for $t > 0$. Notice that $\mathbf{1}(X_t^x)$ and $\mathbf{1}(X_{t^-}^x)$ have the same distribution.*

Proof. The proof applies the following simple corollary of Fubini's theorem on \mathbb{R}^M with Lebesgue measure λ and a probability space (Ω, \mathcal{F}, P) : let ν be a probability measure on $\mathbb{R}^M \times \Omega$ such that there is random density $p : \mathbb{R}^M \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ with

$$\int_{\mathbb{R}^M \times \Omega} f(x, \eta) p(x, \eta) (\lambda \otimes P)(dx, d\eta) = \int_{\mathbb{R}^M \times \Omega} f(x, \eta) \nu(dx, d\eta);$$

then the marginal of ν on \mathbb{R}^M is absolutely continuous with respect to Lebesgue measure λ with density

$$p(x) = \int_{\Omega} p(x, \eta) P(d\eta)$$

for almost all $x \in \mathbb{R}^M$. In our case we know that the law of $(\mathbf{I}(X_t^x))$ is absolutely continuous for almost all trajectories η of the compound Poisson processes L (precisely those where $\tau_1 > 0$, where only one coordinate jumps at each jumping time, and finitely many jumps occur on compact intervals); the probability measure ν corresponds to the distribution of $(\mathbf{I}(X_t^x), L)$, where we choose Ω as the space of cadlag trajectories on $\mathbb{R}_{\geq 0}$ with values in \mathbb{R}^m . Finally, we have that the law of $(\mathbf{I}(X_t^x))$ is $p(x)\lambda(dx)$. \square

6. Applications of the infinite-dimensional result to interest rate theory. In mathematical Finance, the theory of interest rates deals with the market of interest rate related products like swaps, bonds, bills, etc. If one considers default-free products one can crystallize from the data of real markets the prices of default-free zero-coupon bonds for any maturity. A zero-coupon bond contract with maturity T (a calendar date) can be entered at calendar time $t \leq T$ and (certainly) pays 1 unit of currency at maturity time T . Therefore, bonds reflect the level of interest rate between time t and maturity time T . No coupons are paid between t and T , which explains the notion of “zero-coupon bond.” We denote the price of a default-free zero-coupon bond with maturity T at time $t \leq T$ by $P(t, T)$. Commonly one assumes that bond prices are at least C^1 with respect to T , i.e.,

$$P(t, T) = \exp\left(-\int_t^T f(t, r) dr\right)$$

for $0 \leq t \leq T$, where $f(t, T)$ denotes the forward rate. This leads to the concept of the short rate

$$R_t = f(t, t)$$

for $t \geq 0$, which corresponds to the level of interest rate for instantaneous transactions from t to $t + dt$. As usual in mathematical Finance, discounted default-free zero-coupon bonds are modeled by semimartingales, and one assumes the existence of an equivalent martingale measure for discounted price processes. This leads to the following fundamental formula with respect to the martingale measure

$$E\left(\exp\left(-\int_t^T R_s ds\right) \middle| \mathcal{F}_t\right) = P(t, T) = \exp\left(-\int_t^T f(t, r) dr\right).$$

The formula simply expresses the fact that the expected value of discounted value of the payoff $P(T, T) = 1$ conditional on today’s information equals today’s price with respect to the martingale measure. Assuming a jump-diffusion model for $(f(t, T))_{0 \leq t \leq T}$ for $T \geq 0$ for the forward rates together with Musiela’s parametrization $r(t, T - t) = f(t, T)$ for $0 \leq t \leq T$ leads to the famous Heath–Jarrow–Morton (HJM) equation of interest rate theory, which is an SDE taking values in a Hilbert space of forward rate curves H (and therefore an SPDE). We quote here as leading reference [15], where the no-arbitrage conditions for the HJM equation are discussed in all necessary detail. A

very readable introduction can also be found in [8], particularly for HJM equations with jumps.

The HJM equation has been analyzed from different points of view:

- The question of which HJM equations driven by finitely many Brownian motions admit finite-dimensional realizations has been treated in detail in [16]. This research was inspired by [9], where the geometric approach was introduced. The satisfying answer is that, under quite natural restrictions, finite-dimensional realizations do exist if and only if the corresponding factor processes are affine processes. This is the case, for instance, for Vasiček’s model or for the Cox–Ingersoll–Ross model of interest rate theory. In both cases the finite-dimensional realizations are, in fact, two-dimensional. In [8] finite-dimensional realizations are treated for HJM equations with jumps but under the strong restriction that the vector fields do not depend on the forward rate. In this case one can solve the HJM equation explicitly by variation of constants and read off the respective geometric properties of the solution process.
- The question of whether the solution process of an HJM equation always admits a density with respect to Lebesgue measure when projected to a finite-dimensional subspace has been treated in [3] and could be answered affirmatively under Hörmander-type conditions.

We ask here the question—having the theory of the previous sections in mind—of whether a structure of finite-dimensional realizations, such as for Vasiček’s model, can be perturbed so strongly through the introduction of jumps that the resulting HJM evolution is “hypoelliptic,” i.e., the assumptions of Theorem 3 are fulfilled. We can answer this question affirmatively in the case of a Vasiček model. In contrast to [8], we allow the vector fields to be state-dependent (therefore, we cannot hope for explicit solutions of the HJM equation, and we have to apply local methods from differential geometry to conclude).

We consider the HJM model with jumps

$$\begin{cases} dr_t &= \left(\frac{d}{dx} r_{t-} + \alpha_{HJM}(r_{t-}) + \beta_{HJM}(r_{t-}) \right) dt + \sigma(r_{t-}) dB_t + \delta(r_{t-}) dN_t, \\ r_0 &= r^* \in H, \end{cases}$$

on some Hilbert space H of forward rate curves as constructed in [3] and [15]. Here $(B_t)_{t \geq 0}$ is a standard Brownian motion, and $(N_t)_{t \geq 0}$ is a standard Poisson process with intensity $\tilde{\lambda} > 0$ and jump measure $\mu = \delta_1$ (hence $\tilde{\lambda} = \lambda$). Define

$$\begin{aligned} \Psi_1(z) &\equiv \ln \mathbb{E} [e^{zW_1}] = \ln \left(e^{\frac{z^2}{2}} \right) = \frac{z^2}{2} \quad \text{and} \\ \Psi_2(z) &\equiv \ln \mathbb{E} [e^{zN_1}] = \ln (\exp (\lambda (e^z - 1))) = \lambda (e^z - 1). \end{aligned}$$

Then we know from [15, equation (2.4)], that

$$\begin{aligned} \alpha_{HJM}(r)(x) &= -\sigma(r)(x) \Psi_1' \left(- \int_0^x \sigma(r)(y) dy \right) \\ &= \sigma(r)(x) \int_0^x \sigma(r)(y) dy \end{aligned}$$

and

$$\begin{aligned} \beta_{HJM}(r)(x) &= -\delta(r)(x)\Psi'_2\left(-\int_0^x \delta(r)(y)dy\right) \\ &= -\lambda\delta(r)(x)\exp\left(-\int_0^x \delta(r)(y)dy\right). \end{aligned}$$

For an explicit example we choose

$$\sigma(r)(x) = \sigma(x) > 0 \quad \text{and} \quad \delta(r)(x) = -\frac{d}{dx} \ln(B(r)(x)),$$

where the vector field B will be determined later. Then we have

$$\alpha_{HJM}(r)(x) = \sigma(x) \int_0^x \sigma(y)dy$$

and

$$\beta_{HJM}(r)(x) = \lambda \frac{d}{dx} \ln(B(r)(x)) \frac{(B(r)(x))}{(B(r)(0))} = \lambda \frac{\frac{d}{dx}(B(r)(x))}{(B(r)(0))}.$$

We choose B such that $(B(r)(x))$ is positive on H for $x \in \mathbb{R}$ and $B(r)(0) = 1$ for all $r \in H$, whence δ is well defined. Thus, for such $r \in U$ we have

$$\beta_{HJM}(r)(x) = \lambda \frac{d}{dx}(B(r)(x))$$

and

$$\delta(r)(x) = -\frac{d}{dx} \ln(B(r)(x)).$$

A particular choice in the spirit of Remark 1 is given through

$$B(r)(x) = \psi(x, l(r)),$$

where the maps $y \mapsto \frac{d}{dx} \psi(\cdot, y)$ and $y \mapsto \frac{1}{\psi(\cdot, y)}$ from \mathbb{R} to $\text{dom}(A^\infty) \subset H$ are supposed to be C^∞ -bounded with $\psi(0, y) = 1$ for all $y \in \mathbb{R}$. The map l denotes here a non-vanishing linear functional $l : H \rightarrow \mathbb{R}$. Hence δ and β are well defined C^∞ -bounded vector fields on the whole Hilbert space, and we have global existence of mild solutions.

The Vasiček model is defined by

$$\sigma(r) = \rho \exp(-ax)$$

for $\rho, a > 0$ without any jump component. By [9] and [16] we know that the Vasiček model admits finite-dimensional realizations, as for

$$V_0(r)(x) = \frac{d}{dx} r(x) + \alpha_{HJM}(r)(x)$$

we have

$$\dim(\{V_0, \sigma\}_{LA}(r)) \leq 2,$$

at any point $r \in \text{dom}((\frac{d}{dx})^\infty)$. Here the index LA stands for the Lie algebra generated by the vector fields V_0, σ on $\text{dom}(A^\infty)$. If we add a jump structure as described above and if we choose ψ generic, the two-dimensional structure (a regular finite-dimensional realization in the sense of [16]) is destroyed, since then the drift changes due to no-arbitrage. We obtain a dense Lie algebra if we choose the vector field B generically.

These results might be of interest for recent works in interest rate theory; see, for instance, [23], where under diffusion assumptions hypoellipticity is tested empirically. If one allows for jumps in an HJM model, the phenomenon of hypoellipticity seems to be more generic.

7. Smooth densities for the law X_T^x on \mathbb{R}^M . In the what follows we consider the cases $\dim H = M$ and $\mathbf{l} = \text{id}$, and we choose a coordinate representation $H = \mathbb{R}^M$. We then want to show that the p th power of the inverse of the Malliavin covariance matrix of X_t^x for $t > 0$ can be integrated even with respect to Poisson trajectory η . We therefore need an extension of the Hörmander condition which is called the uniform Hörmander condition.

Following [24], we define

$$\Sigma'_0 := \{V_1, \dots, V_d\},$$

$$\Sigma'_n := \left\{ [V_k, V], k = 1, \dots, d, V \in \Sigma'_{n-1}; [V_0, V] + \frac{1}{2} \sum_{i=1}^d [V_i, [V_i, V]], V \in \Sigma'_{n-1} \right\}$$

for $n \geq 1$. We assume that there exist j_0 and $c > 0$ such that

$$(7.1) \quad \inf_{\xi \in S^{M-1}} \sum_{j=0}^{j_0} \sum_{V \in \Sigma'_j} \langle V(x), \xi \rangle^2 \geq c$$

uniformly in $x \in \mathbb{R}^M$.

THEOREM 5. *Assume that $\dim H < \infty$. We take Assumptions 2 and 3 for granted but assume that the Hörmander condition (2.3) holds true uniformly on \mathbb{R}^M in the sense of (7.1). Let $(X_t^x)_{t \geq 0}$ denote the unique cadlag solution of (3.1), and fix $t > 0$. Then the random variable X_t^x admits a smooth density with respect to Lebesgue measure on \mathbb{R}^M . Furthermore, the covariance matrix of X_t^x is invertible with p -integrable inverse for all $p \geq 1$.*

Proof. We write the Malliavin derivative of X_t^x ,

$$D_s^i X_t^x = J_{0 \rightarrow t}(x) J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) \mathbf{1}_{[0,t]}(s),$$

and calculate the reduced covariance matrix

$$\langle C_t \xi, \xi \rangle = \sum_{i=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), \xi \rangle^2 ds.$$

We now apply the result from Theorem 1 and the condition on the trajectories of the compound Poisson process (2)

$$(7.2) \quad \sup_{\xi \in S^{M-1}} P(\langle C_t \xi, \xi \rangle < \epsilon)$$

$$= \sup_{\xi \in S^{M-1}} \sum_{n_1, \dots, n_m \geq 0} \left[\prod_{k=1}^m P(N_t^k = n_k) \right] P(\langle C_t \xi, \xi \rangle < \epsilon | N_t^j = n_j \text{ for } j = 1, \dots, m).$$

As in the proof of Theorem 3, we can decompose $\langle C_t \xi, \xi \rangle$ into

$$\langle C_t \xi, \xi \rangle = \sum_{k=0}^{\infty} \sum_{i=1}^d \int_{\tau_k \wedge t}^{\tau_{k+1} \wedge t} \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s-}^x), \xi \rangle^2 ds,$$

where $\tau_0 = 0 < \tau_1 < \dots < \tau_n \leq \dots$ denotes the sequence of jump times of $(N_t)_{0 \leq t \leq T}$. Hence, we obtain for $n = n_1 + \dots + n_m$

$$\begin{aligned} & \sup_{\xi \in S^{M-1}} P(\langle C_t \xi, \xi \rangle < \epsilon | N_t^j = n_j, j = 1, \dots, m) \\ & \leq \sup_{\xi \in S^{M-1}} P\left(\sum_{i=1}^d \int_{\tau_k \wedge t}^{\tau_{k+1} \wedge t} \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s-}^x), \xi \rangle^2 ds < \epsilon \mid N_t^j = n_j, j = 1, \dots, m\right) \end{aligned}$$

for all $0 \leq k \leq n$. Observing that $\max_{0 \leq k \leq n} (\tau_{k+1} - \tau_k)^{K(p)} \geq (\frac{t}{n})^{K(p)}$ (after all, we have only n jumps, so the maximal distance between two consecutive jumps is bigger than $\frac{t}{n}$), we finally obtain

$$P(\langle C_t \xi, \xi \rangle < \epsilon | N_t^j = n_j \text{ for } j = 1, \dots, m) \leq \epsilon^p$$

for $0 \leq \epsilon \leq (\frac{t}{n})^{K(p)} \epsilon_0(p)$ due to the calculations outlined in the appendix. Note that we can apply the calculations from the appendix, since $J_{0 \rightarrow s}(x)^{-1}$ is well defined and bounded due to boundedness of $(id + z d\delta_j)^{-1}$ for $z \in \text{supp}(\mu_j)$ and $j = 1, \dots, m$. Hence integration with respect to the measures μ_j is possible and yields finite bounds. Recall also that μ_j has moments of all orders; hence X_t^x is L^p , and so is $J_{0 \rightarrow s}(x)^{-1}$ (see [25] for all necessary details on SDEs).

Let $\Lambda = \inf_{\xi \in S^{M-1}} \langle C_t \xi, \xi \rangle$ be the smallest eigenvalue of the reduced covariance matrix C_t . Following the steps of [24, Lemma 2.3.1], we know that

$$P(\Lambda < \epsilon \mid N_t^j = n_j \text{ for } j = 1, \dots, m) \leq \text{const} \cdot \epsilon^p$$

for any $p \geq 2$ and $0 \leq \epsilon \leq (\frac{t}{n})^{K(p+2M)} \epsilon_0(p+2M) =: \epsilon_{max}$, where the constant depends on the p -norm of C_t . In what follows we shall denote any constant of this type by D . We denote by ρ the law of Λ conditioned on $N_t^j = n_j$ for $j = 1, \dots, m$. Consequently, for $j = 1, \dots, m$, we have by Fubini's theorem

$$\begin{aligned} E\left(\frac{1}{\Lambda^{p-1}} \mid N_t^j = n_j\right) &= E\left(\frac{1}{\Lambda^{p-1}} \cdot \mathbf{1}_{\{\Lambda > \epsilon_{max}\}} \mid N_t^j = n_j\right) \\ &+ E\left(\frac{1}{\Lambda^{p-1}} \cdot \mathbf{1}_{\{\Lambda \leq \epsilon_{max}\}} \mid N_t^j = n_j\right) \\ &\leq \frac{1}{\epsilon_{max}^{p-1}} + \int_0^{\epsilon_{max}} \frac{1}{z^{p-1}} \rho(dz) \\ &= \frac{1}{\epsilon_{max}^{p-1}} + \int_0^{\epsilon_{max}} (p-1) \int_z^\infty \frac{1}{t^p} dt \rho(dz) \\ &= \frac{1}{\epsilon_{max}^{p-1}} + (p-1) \int_0^{\epsilon_{max}} \frac{1}{z^p} \int_0^z \rho(dt) dz \end{aligned}$$

$$\begin{aligned}
 &+ (p - 1) \int_{\epsilon_{max}}^{\infty} \frac{1}{z^p} \int_0^{\epsilon_{max}} \rho(dt) dz \\
 &\leq \frac{D}{\epsilon_{max}^{p-1}} + D \underbrace{\int_0^{\epsilon_{max}} \frac{1}{z^p} z^p dz}_{=\epsilon_{max}} \\
 &\leq D \left(\frac{t}{n}\right)^{K(p+2M)} \epsilon_0(p + 2M) \\
 &+ \frac{D}{\left[\left(\frac{t}{n}\right)^{K(p+2M)} \epsilon_0(p + 2M)\right]^{p-1}}.
 \end{aligned}$$

Here we applied $\int_0^z \rho(dt) \leq D \times z^p$ as previously proved. Hence, through the decomposition (7.2),

$$\begin{aligned}
 E\left(\frac{1}{\Lambda^{p-1}}\right) &\leq \sum_{n_1, \dots, n_m > 0} \prod_{k=1}^m P(N_t^k = n_k) \\
 &\cdot D \cdot \left[\left(\frac{t}{n}\right)^{K(p+2M)} \epsilon_0(p + 2M) + \frac{1}{\left[\left(\frac{t}{n}\right)^{K(p+2M)} \epsilon_0(p + 2M)\right]^{p-1}} \right] < \infty;
 \end{aligned}$$

the result follows by $n = n_1 + \dots + n_m$ and by the following fact for any real number K :

$$\sum_{n_1, \dots, n_m > 0} \frac{\tilde{\lambda}_1^{n_1} \dots \tilde{\lambda}_m^{n_m}}{n_1! \dots n_m!} e^{-t\tilde{\lambda}_1 n_1 - \dots - t\tilde{\lambda}_m n_m} t^n (n_1 + \dots + n_m)^K < \infty. \quad \square$$

Remark 8. We could have also applied the beautiful results of [21, Corollary 3.25] to evaluate the L^p -norm of the inverse of the covariance matrix between two jumps. Both methods lead to the same result. We have been choosing our approach since we can root it as much as possible in the standard reference [24].

8. Calculating the Greeks in finite dimension. In what follows we consider the case $\dim H = M$ and $\mathbf{l} = \text{id}$ as in the previous section. Once we are given an invertible Malliavin covariance matrix with p -integrable inverse such as in Theorem 5, we can easily calculate derivatives with respect to initial values and obtain explicit formulas for so-called Malliavin weights (see [17] for successful applications of this method in mathematical Finance). We quickly sum up the main idea: in mathematical Finance the gradient of the function $x \mapsto E(f(X_t^x))$ has the meaning of hedging ratios, which control the hedging portfolios away from jumps. Hence for any hedging portfolio corresponding to prices $E(f(X_t^x))$ of a certain derivative at maturity $t > 0$ it is crucial to know $\nabla E(f(X_t^x))$ to perform hedging off jumps.

Very often pricing results in the applications of a weak approximation scheme for the process X , for instance, the Euler–Maruyama scheme. For the calculation of $\nabla E(f(X_t^x))$ in the direction of some vector $v \in H$, basically three methods can be applied:

- a finite difference method to approximate $\nabla E(f(X_t^x)) \cdot v$, resulting in the calculation of $\frac{E(f(X_t^{x+\epsilon v})) - E(f(X_t^x))}{\epsilon}$ ($v \in H$ denotes some vector) for small $\epsilon > 0$;
- a pathwise method applying the formula

$$\nabla E(f(X_t^x)) \cdot v = E(df(X_t^x) J_{0 \rightarrow t}(x) \cdot v),$$

- resulting in the weak numerical approximation of $(X_t^x, J_{0 \rightarrow t}(x))$; and
- the method of Malliavin weights applying the formula

$$\nabla E(f(X_t^x)) \cdot v = E(f(X_t^x)\pi^v),$$

resulting in the weak numerical approximation of (X_t^x, π^v) .

The first method is the most robust in the sense that it can be applied under very weak assumptions both on X_t^x and on the payoff f , but the rate of convergence might be very slow since the errors of Monte Carlo evaluations are amplified. The second method works for all reasonable jump-diffusion processes, but one needs Lipschitz conditions on the payoff f . The third method needs the assumptions of Theorem 5 on X_t^x , but no restrictions on the payoff f , which makes the third method attractive for several problems from mathematical Finance, where measurable, non-Lipschitz payoffs (e.g., digital options) are quite usual and hypoellipticity assumptions as in Theorem 5 are common, too.

The implementation of procedures for all three methods has been outlined in [18] in the pure diffusion case. We shall not work on this issue here, since our main message is that one can implement precisely the same methods in pure diffusion cases as in jump-diffusion cases. The important point is that the formulas have the same structure in both cases, a fact we shall point out in this section on several occasions.

We denote in what follows the Skorohod integral (resp., the divergence operator) by δ and its domain by $\text{dom}(\delta)$.

DEFINITION 1. Assume that $H = \mathbb{R}^M$ and fix $t > 0$ and a direction $v \in \mathbb{R}^M$. We define a set of Skorohod-integrable processes

$$\mathbb{A}_{t,x,v} = \left\{ a \in \text{dom}(\delta) \text{ such that } \sum_{i=1}^d \int_0^t J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) a_s^i ds = v \right\}$$

and call it the set of path-perturbations with target-value v .

Remark 9. In the previous definition, as in the whole section, assertions on Skorohod-integrability are meant Poissonian-trajectory-wise.

PROPOSITION 1. Assume that $H = \mathbb{R}^M$. We take Assumption 3 for granted. Fix $t > 0$ and a direction $v \in \mathbb{R}^M$. Assume furthermore uniform ellipticity; i.e., $M = d$ and there is $c > 0$ such that

$$\inf_{\xi \in S^{M-1}} \sum_{k=1}^M \langle V_k(x), \xi \rangle^2 \geq c.$$

Then $\mathbb{A}_{t,x,v} \neq \emptyset$ and there exists an integrable, real valued random variable π^v (which depends linearly on v) such that for all bounded random variables f we obtain

$$\frac{d}{d\epsilon} \Big|_{\epsilon=0} E(f(X_t^{x+\epsilon v})) = E(f(X_t^x)\pi^v).$$

Such a random variable π^v is called a Malliavin weight and can be obtained through an Itô integral.

Remark 10. The assertion of this theorem corresponds to Assumption (E) in [18] and to the assumptions of [17]. The assumptions are seen as too restrictive since not every problem in mathematical Finance has an elliptic volatility matrix. The formulas of [18] and [17] correspond precisely to the formulas obtained here, which

leads to the assertion that even in the presence of jumps one can apply the same (numerical) methods for the calculation of Greeks as in the pure diffusion cases.

Proof. Here the proof is particularly simple, since we can take a matrix $\sigma(x) := (V_1(x), \dots, V_M(x))$, which is uniformly invertible with bounded inverse. We define

$$a_s := \frac{1}{t} \sigma(X_{s^-}^x)^{-1} \cdot J_{0 \rightarrow s}(x) \cdot v$$

for $0 \leq s \leq t$ and obtain that $a \in \mathbb{A}_{t,x,v}$. Furthermore, as in [17] and [13], we obtain

$$\pi^v = \sum_{i=1}^M \int_0^t a_s^i dB_s^i,$$

since the Skorohod-integrable process a is in fact adapted, left-continuous, and hence Itô-integrable. \square

THEOREM 6. *Assume that $H = \mathbb{R}^M$. We take Assumptions 2 and 3 for granted but assume that the Hörmander condition (2.3) holds true uniformly on \mathbb{R}^M (see section 7). Fix $t > 0$ and a direction $v \in \mathbb{R}^M$. Then $\mathbb{A}_{t,x,v} \neq \emptyset$ and there exists an integrable, real valued random variable π^v (which depends linearly on v) such that for all bounded random variables f we obtain*

$$\frac{d}{d\epsilon} \Big|_{\epsilon=0} E(f(X_t^{x+\epsilon v})) = E(f(X_t^x) \pi^v).$$

We can choose π^v to be the Skorohod integral of any element $a \in \mathbb{A}_{t,x,v} \neq \emptyset$ and call it a Malliavin weight. Moreover, by the explicit construction of a in the proof, we can assert that π^v is the sum of an Itô integral and an integral with respect to Lebesgue measure; see, for instance, [18].

Remark 11. The assertion of this theorem corresponds to Assumption (E') in [18]. The assumptions (E) and (E') are fundamental for the third method in [18]. Again the formulas of [18] correspond to the formulas obtained here.

Proof. We take f bounded with bounded first derivative; then we obtain

$$\frac{d}{d\epsilon} \Big|_{\epsilon=0} E(f(X_t^{x+\epsilon v})) = E(df(X_t^x) J_{0 \rightarrow t}(x) \cdot v).$$

If there is $a \in \mathbb{A}_{t,x,v}$, we obtain

$$\begin{aligned} E(df(X_t^x) J_{0 \rightarrow t}(x) \cdot v) &= E \left(df(X_t^x) \sum_{i=1}^d \int_0^t J_{0 \rightarrow t}(x) J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) a_s^i ds \right) \\ &= E \left(\sum_{i=1}^d \int_0^t df(X_t^x) J_{0 \rightarrow t}(x) J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) a_s^i ds \right) \\ &= E \left(\sum_{i=1}^d \int_0^t D_s^i f(X_t^x) a_s^i ds \right) \\ &= E(f(X_t^x) \delta(a)). \end{aligned}$$

Here we cannot assert that the strategy is Itô-integrable, since it will be anticipative in general. In order to see that $\mathbb{A}_{t,x,v} \neq \emptyset$, we construct an element, namely,

$$a_s^i := \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), (C_t)^{-1} v \rangle,$$

where C_t denotes the reduced covariance matrix from Theorem 5. Indeed,

$$\begin{aligned} & \sum_{i=1}^d \left\langle \int_0^t J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) a_s^i ds, \xi \right\rangle \\ &= \sum_{i=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), \xi \rangle \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), (C_t)^{-1} v \rangle ds \\ &= \langle \xi, C_t (C_t)^{-1} v \rangle = \langle \xi, v \rangle \end{aligned}$$

for all $\xi \in \mathbb{R}^M$, since C_t is a symmetric random operator defined via

$$\langle \xi, C_t \xi \rangle = \sum_{i=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), \xi \rangle^2 ds$$

for $\xi \in \mathbb{R}^M$. \square

For any other derivative with respect to parameters ϵ , we consider a modified set, namely,

$$\mathbb{B}_{t,x,v} = \left\{ b \in \text{dom}(\delta) \left| \sum_{i=1}^d \int_0^t J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) b_s^i ds = J_{0 \rightarrow t}(x)^{-1} \frac{d}{d\epsilon} \Big|_{\epsilon=0} X_t^{x,\epsilon} \right. \right\}.$$

Here we are given a parameter-dependent process $X_t^{x,\epsilon}$, where all derivatives with respect to ϵ can be calculated nicely. Also in this case we can construct—if the reduced covariance matrix is invertible and regular enough—an element, namely,

$$b_s^i := \left\langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), (C_t)^{-1} J_{0 \rightarrow t}(x)^{-1} \frac{d}{d\epsilon} \Big|_{\epsilon=0} X_t^{x,\epsilon} \right\rangle.$$

This is a consequence of the reasoning

$$\begin{aligned} & \sum_{i=1}^d \left\langle \int_0^t J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x) b_s^i ds, \xi \right\rangle = \sum_{i=1}^d \int_0^t \left\langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), \xi \right\rangle \\ & \cdot \left\langle J_{0 \rightarrow s}(x)^{-1} V_i(X_{s^-}^x), (C_t)^{-1} J_{0 \rightarrow t}(x)^{-1} \frac{d}{d\epsilon} \Big|_{\epsilon=0} X_t^{x,\epsilon} \right\rangle ds \\ &= \left\langle \xi, C_t (C_t)^{-1} J_{0 \rightarrow t}(x)^{-1} \frac{d}{d\epsilon} \Big|_{\epsilon=0} X_t^{x,\epsilon} \right\rangle = \left\langle \xi, J_{0 \rightarrow t}(x)^{-1} \frac{d}{d\epsilon} \Big|_{\epsilon=0} X_t^{x,\epsilon} \right\rangle, \end{aligned}$$

due to the symmetry of C_t .

Appendix.

THEOREM 7. *Let $(\Omega, \mathcal{F}, P, (\mathcal{F}_t)_{t \geq 0})$ be a filtered probability space, and let $(B_t)_{t \geq 0}$ be a d -dimensional Brownian motion adapted to the filtration (which is not necessarily generated by the Brownian motion). Let V, V_1, \dots, V_d , the diffusion vector fields, be C^∞ -bounded on \mathbb{R}^M , and consider the continuous solution $(X_t^x)_{0 \leq t \leq T}$ of an SDE (in Stratonovich notation). V_0 denotes the Stratonovich corrected drift term,*

$$(A.1) \quad dX_t^x = V_0(X_t^x) dt + \sum_{i=1}^d V_i(X_t^x) \circ dB_t^i,$$

$$(A.2) \quad X_0^x = x.$$

Assume that the uniform Hörmander condition holds true (see the proof for the precise statement). Then for any $p \geq 1$ there exist numbers $\epsilon_0(p) > 0$ and an integer $K(p) \geq 1$ such that for each $0 < t < T$

$$\sup_{\xi \in S^{M-1}} P(\langle C_t \xi, \xi \rangle < \epsilon) \leq \epsilon^p$$

holds true for $0 \leq \epsilon \leq t^{K(p)} \epsilon_0(p)$. The result holds uniformly in x .

Remark 12. The time-dependence of the estimate $0 \leq \epsilon \leq t^{K(p)} \epsilon_0(p)$ is best explained by redoing the proof. It is heavily applied in section 7 and the main technical ingredient of the given proof. We could have also directly used the results from [21].

Proof. The proof of the theorem is a careful rereading of the Norris lemma and the classical proof of the Hörmander theorem in probability theory (see [22] or [24]). We shall sketch this path in what follows (see [24, pp. 120–123]).

1. Consider the random quadratic form

$$\langle C_t \xi, \xi \rangle = \sum_{i=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V_i(X_s^x), \xi \rangle^2 ds.$$

Following [24], we define

$$\begin{aligned} \Sigma'_0 &:= \{V_1, \dots, V_d\}, \\ \Sigma'_n &:= \left\{ [V_k, V], k = 1, \dots, d, V \in \Sigma'_{n-1}; [V_0, V] + \frac{1}{2} \sum_{i=1}^d [V_i, [V_i, V]], V \in \Sigma'_{n-1} \right\} \end{aligned}$$

for $n \geq 1$. We assume that there exist j_0 and $c > 0$ such that

$$\inf_{\xi \in S^{M-1}} \sum_{j=0}^{j_0} \sum_{V \in \Sigma'_j} \langle V(x), \xi \rangle^2 \geq c$$

uniformly in $x \in \mathbb{R}^M$.

2. We define $m(j) := 2^{-4j}$ for $0 \leq j \leq j_0$ and the sets

$$E_j := \left\{ \sum_{V \in \Sigma'_j} \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V(X_s^x), \xi \rangle^2 ds \leq \epsilon^{m(j)} \right\}.$$

We consider the decomposition

$$\begin{aligned} E_0 &= \{ \langle C_t \xi, \xi \rangle \leq \epsilon \} \subset (E_0 \cap E_1^c) \cup (E_1 \cap E_2^c) \cup \dots \cup (E_{j_0-1} \cap E_{j_0}^c) \cup F, \\ F &= E_0 \cap \dots \cap E_{j_0} \end{aligned}$$

and proceed with

$$P(F) \leq C \epsilon^{\frac{q\beta}{2}}$$

for $\epsilon \leq \epsilon_1$ and any $q \geq 2$ with a constant C depending on q and the norms of the derivatives of the vector fields V_0, \dots, V_d . Furthermore, $0 < \beta < m(j_0)$. The number ϵ_1 is determined by the following two (!) equations:

$$\begin{aligned} (j_0 + 1) \epsilon_1^{m(j_0)} &< \frac{c \epsilon_1^\beta}{4}, \\ \epsilon_1^\beta &< t. \end{aligned}$$

Hence ϵ_1 depends on j_0, c, t , and the choice of β via

$$\epsilon_1 < \min \left(t^{\frac{1}{\beta}}, \left(\frac{c}{4(j_0 + 1)} \right)^{\frac{1}{m(j_0) - \beta}} \right).$$

This little observation, in addition to the proof in [24], is key for our proof.

3. We obtain furthermore that

$$\begin{aligned} P(E_j \cap E_{j+1}^c) &= P \left(\sum_{V \in \Sigma'_j} \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V(X_s^x), \xi \rangle^2 ds \leq \epsilon^{m(j)}, \right. \\ &\quad \left. \sum_{V \in \Sigma'_{j+1}} \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V(X_s^x), \xi \rangle^2 ds > \epsilon^{m(j+1)} \right) \\ &\leq \sum_{V \in \Sigma'_j} P \left(\int_0^t \langle J_{0 \rightarrow s}(x)^{-1} V(X_s^x), \xi \rangle^2 ds \leq \epsilon^{m(j)}, \right. \\ &\quad \left. \sum_{k=1}^d \int_0^t \langle J_{0 \rightarrow s}(x)^{-1} [V_k, V](X_s^x), \xi \rangle^2 ds \right. \\ &\quad \left. + \int_0^t \left\langle J_{0 \rightarrow s}(x)^{-1} \left([V_0, V] + \frac{1}{2} \sum_{i=1}^d [V_i, [V_i, V]] \right) (X_s^x), \xi \right\rangle^2 ds > \frac{\epsilon^{m(j+1)}}{n(j)} \right), \end{aligned}$$

where $n(j) = \#\Sigma'_j$. Since we can find the bounded variation and the quadratic variation parts of the martingale $(\langle J_{0 \rightarrow s}(x)^{-1} V(X_s^x), \xi \rangle)_{0 \leq s \leq t}$ in the above expression, we are able to apply the Norris lemma (see [24, Lemma 2.3.2]). We observe that $8m(j + 1) < m(j)$; hence we can apply it with $q = \frac{m(j)}{m(j+1)}$.

4. We obtain for $p \geq 2$ —still by the Norris lemma—the estimate

$$P(E_j \cap E_{j+1}^c) \leq d_1 \left(\frac{\epsilon^{m(j+1)}}{n(j)} \right)^{rp} + d_2 \exp \left(- \left(\frac{\epsilon^{m(j+1)}}{n(j)} \right)^{-\nu} \right)$$

for $\epsilon \leq \epsilon_2$. Furthermore, $r, \nu > 0$ with $18r + 9\nu < q - 8$, and the numbers d_1, d_2 depend on the vector fields V_0, \dots, V_d , and on p, T . The number ϵ_2 can be chosen as $\epsilon_2 = \epsilon_3 t^{k_1}$, where ϵ_3 does not depend on t anymore.

5. Putting all this together, we take the minimum of ϵ_1 and ϵ_2 to obtain the desired dependence on t . \square

Acknowledgment. The authors are grateful to two unknown referees who helped to considerably improve the presentation and contents of this paper. In particular, we are grateful that they pointed out Corollary 3.25 of [21].

REFERENCES

[1] Y. BAKHTIN AND J. MATTINGLY, *Malliavin calculus for infinite-dimensional systems with additive noise*, J. Funct. Anal., 249 (2007), pp. 307–353.
 [2] V. BALLY, M.-P. BAVOUZET, AND M. MESSAOUD, *Integration by parts formula for locally smooth laws and applications to sensitivity computations*, Ann. Appl. Probab., 17 (2007), pp. 33–66.
 [3] F. BAUDOIN AND J. TEICHMANN, *Hypoellipticity in infinite dimensions and an application in interest rate theory*, Ann. Appl. Probab., 15 (2005), pp. 1765–1777.

- [4] Y. I. BELOPOL'SKAYA AND Y. L. DALECKY, *Stochastic Equations and Differential Geometry*, Mathematics and Its Applications (Soviet Series) 30, Kluwer Academic Publishers, Dordrecht, The Netherlands, 1990.
- [5] Y. I. BELOPOL'SKAYA AND Y. L. DALETSKY, *Smoothness of transition probabilities of Markov processes described by stochastic differential equations in Hilbert space*, *Dopov. / Dokl. Akad. Nauk Ukraïni*, 9 (1994), pp. 40–45 (in Russian).
- [6] K. BICHTLER, J.-B. GRAVEREAUX, AND J. JACOD, *Malliavin Calculus for Processes with Jumps*, Stochastics Monographs 2, Gordon and Breach Science Publishers, New York, 1987.
- [7] J.-M. BISMUT, *Calcul des variations stochastique et processus de sauts*, *Z. Wahrsch. Verw. Gebiete*, 63 (1983), pp. 147–235.
- [8] T. BJÖRK AND A. GOMBANI, *Minimal realizations of interest rate models*, *Finance Stoch.*, 3 (1999), pp. 413–432.
- [9] T. BJÖRK AND L. SVENSSON, *On the existence of finite-dimensional realizations for nonlinear forward rate models*, *Math. Finance*, 11 (2001), pp. 205–243.
- [10] R. A. CARMONA AND B. ROZOVSKII, EDs., *Stochastic Partial Differential Equations: Six Perspectives*, Math. Surveys Monogr. 64, AMS, Providence, RI, 1999.
- [11] T. CASS, *Smooth densities for stochastic differential equations with jumps*, *Stochastic Process Appl.*, to appear.
- [12] G. DA PRATO AND J. ZABCZYK, *Stochastic Equations in Infinite Dimensions*, Encyclopedia Math. Appl. 44, Cambridge University Press, Cambridge, UK, 1992.
- [13] M. H. A. DAVIS AND M. P. JOHANSSON, *Malliavin Monte Carlo Greeks for jump diffusions*, *Stochastic Process. Appl.*, 116 (2006), pp. 101–129.
- [14] Y. EL-KHATIB AND N. PRIVAULT, *Computations of Greeks in a market with jumps via the Malliavin calculus*, *Finance Stoch.*, 8 (2004), pp. 161–179.
- [15] D. FILIPOVIC AND S. TAPPE, *Existence of Lévy term structure models*, *Finance Stoch.*, 12 (2008), pp. 83–115.
- [16] D. FILIPOVIĆ AND J. TEICHMANN, *Existence of invariant manifolds for stochastic equations in infinite dimension*, *J. Funct. Anal.*, 197 (2003), pp. 398–432.
- [17] E. FOURNIÉ, J.-M. LASRY, J. LEBUCHOUX, P.-L. LIONS, AND N. TOUZI, *Applications of Malliavin calculus to Monte Carlo methods in finance*, *Finance Stoch.*, 3 (1999), pp. 391–412.
- [18] E. GOBET AND R. MUNOS, *Sensitivity analysis using Itô–Malliavin calculus and martingales, and application to stochastic optimal control*, *SIAM J. Control Optim.*, 43 (2005), pp. 1676–1713.
- [19] R. S. HAMILTON, *The inverse function theorem of Nash and Moser*, *Bull. Amer. Math. Soc. (N.S.)*, 7 (1982), pp. 65–222.
- [20] A. KRIEGL AND P. W. MICHOR, *The Convenient Setting of Global Analysis*, Math. Surveys Monogr. 53, AMS, Providence, RI, 1997.
- [21] S. KUSUOKA AND D. STROOCK, *Applications of the Malliavin calculus. II*, *J. Fac. Sci. Univ. Tokyo Sect. IA Math.*, 32 (1985), pp. 1–76.
- [22] P. MALLIAVIN, *Stochastic Analysis*, Grundlehren Math. Wiss. 313, Springer-Verlag, Berlin, 1997.
- [23] P. MALLIAVIN, M. E. MANCINO, AND M. C. RECCHIONI, *A non-parametric calibration of the HJM geometry: An application of Itô calculus to financial statistics*, *Jpn. J. Math.*, 2 (2007), pp. 55–77.
- [24] D. NUALART, *The Malliavin Calculus and Related Topics*, Probab. Appl. (New York), Springer-Verlag, New York, 1995.
- [25] P. PROTTER, *Stochastic Integration and Differential Equations. A New Approach*, Appl. Math. (New York), Springer-Verlag, Berlin, 1990.
- [26] J. TEICHMANN, *Another approach to some rough and stochastic partial differential equations*, Preprint, 2008.
- [27] X. Y. ZHOU, *On the existence of solutions with smooth density of stochastic differential equations in plane*, *Acta Math. Sinica (N.S.)*, 8 (1992), pp. 432–446.

BILLIARD SCATTERING ON ROUGH SETS: TWO-DIMENSIONAL CASE*

ALEXANDER PLAKHOV†

Abstract. The notion of a rough two-dimensional (convex) body is introduced, and to each rough body there is assigned a measure on \mathbb{T}^3 describing billiard scattering on the body. The main result is characterization of the set of measures generated by rough bodies. This result can be used to solve various problems of least aerodynamical resistance.

Key words. billiards, scattering on rough surfaces, Monge–Kantorovich optimal mass transportation, problems of minimal and maximal resistance, shape optimization

AMS subject classifications. 37N05, 49K30, 49Q10

DOI. 10.1137/070709700

1. Definition of a rough set and statement of main theorem.

1.1. Introductory remarks and review of literature. In this paper the notion of a rough two-dimensional (convex) body is given and some properties of rough bodies are established.

Let $B \subset \mathbb{R}^2$ be a convex bounded set with nonempty interior, that is, a bounded convex body. Consider the “set” obtained from B by moving off a set of “very small” area. Such a (heuristically defined) set is called a *rough body*: from the “macroscopic” point of view, it almost coincides with B , and, from the “microscopic” point of view, it contains some “flaws.” (One can imagine a detail of a mechanism that, after a period of exploitation, has some defects.) If the removed set adjoins the boundary ∂B , then one can expect that a flow of point particles incident on the rough body is reflected in another way as compared to reflection from B .

The notion of rough body arises naturally when studying Newton-like problems of the body of least resistance. The first problem of such kind was considered by Newton himself [1]. Recently there were several works made concerning the problem of least resistance in various classes of admissible bodies; see, e.g., [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [15]. The solution of a minimization problem for the case of *rotating bodies* can be naturally identified with a rough body [14]; see also the concluding remarks to this paper.

There are many papers on particle scattering by rough bodies (see, e.g., [16], [17], [18]); they describe bodies and flows of particles that occur in nature. On the contrary, we assume that a rough structure can be “manufactured,” and our aim is to describe all possible rough structures.

1.2. Definition of a rough body. It is supposed that the “microscopic structure” of the boundary of a rough body can be detected from observations of particle

*Received by the editors November 30, 2007; accepted for publication (in revised form) July 16, 2008; published electronically January 28, 2009. This work was supported by the *Centre for Research on Optimization and Control* (CEOC) from the “*Fundação para a Ciência e a Tecnologia*” (FCT), cofinanced by the European Community Fund FEDER/POCTI, and by FCT research project PTDC/MAT/72840/2006.

<http://www.siam.org/journals/sima/40-6/70970.html>

†Department of Mathematics, Aveiro University, Aveiro 3810, Portugal. Current address: Institute of Mathematical and Physical Sciences, Aberystwyth University, Aberystwyth SY23 3BZ, UK (axp@aber.ac.uk).

scattering on the body. From this point of view, two rough bodies are considered equal if they scatter flows of particles in an identical manner. Having these observations in mind, we give the definition of a rough body.

Let B be a bounded convex body. Denote by $n(\xi)$ the unit outer normal vector to ∂B at a regular point $\xi \in \partial B$, and denote by $(\partial B \times S^1)_+$ the set of pairs $(\xi, v) \in \partial B \times S^1$ such that $\langle n(\xi), v \rangle \geq 0$. Here and in what follows, $\langle \cdot, \cdot \rangle$ means the standard scalar product in \mathbb{R}^2 . The set $(\partial B \times S^1)_+$ is equipped with the measure μ which is defined by $d\mu(\xi, v) = \langle n(\xi), v \rangle d\xi dv$, where $d\xi$ and dv are the one-dimensional Lebesgue measures on ∂B and S^1 , respectively.

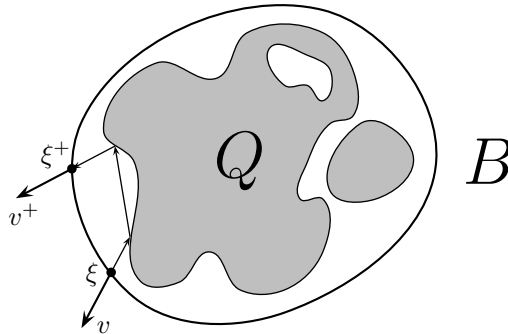


FIG. 1. A billiard trajectory in $\mathbb{R}^2 \setminus Q$.

Let Q be a set with piecewise smooth boundary contained in B ; consider the billiard in $\mathbb{R}^2 \setminus Q$. Note that Q is not necessarily connected. For $(\xi, v) \in (\partial B \times S^1)_+$, consider a billiard particle starting at the point ξ with the velocity $-v$. After several (maybe none) reflections from $\partial Q \setminus \partial B$, the particle will intersect ∂B again, at a point $\xi^+ = \xi_{Q,B}^+(\xi, v) \in \partial B$; denote by $v^+ = v_{Q,B}^+(\xi, v)$ the velocity at this point (see Figure 1). It may happen that the initial point ξ belongs to ∂Q ; in that case we have $\xi^+ = \xi$ and the vector v^+ is symmetric to v with respect to $n(\xi)$. It may also happen that at some moment the particle either gets into a singular point of ∂Q , or touches ∂Q at a regular point, or stays in $B \setminus Q$ forever and does not intersect ∂B again, or makes an infinite number of reflections in finite time. The set of corresponding points (ξ, v) has zero measure, and the corresponding values $\xi_{Q,B}^+(\xi, v)$ and $v_{Q,B}^+(\xi, v)$ are not defined.

Thus, there is defined the one-to-one mapping $T_{Q,B} : (\xi, v) \mapsto (\xi_{Q,B}^+(\xi, v), v_{Q,B}^+(\xi, v))$ of a full measure subset of $(\partial B \times S^1)_+$ onto itself. It has the following properties:

- T1. $T_{Q,B}$ preserves the measure μ .
- T2. $T_{Q,B}^{-1} = T_{Q,B}$.

The mapping $T_{Q,B}$ induces the measure $\nu_{Q,B}$ on $\mathbb{T}^3 = S^1 \times S^1 \times S^1$ in the following way. Let $A \subset \mathbb{T}^3$ be a Borel set; by definition,

$$\nu_{Q,B}(A) = \mu \left(\{(\xi, v) \in (\partial B \times S^1)_+ : (v, v_{Q,B}^+(\xi, v), n(\xi)) \in A\} \right).$$

In fact, the measure $\nu_{Q,B}$ contains information about particle scattering on Q . Imagine that an observer has no means to track the trajectory of particles inside B . Instead, for each incident particle there is registered the triple of vectors: the initial and final

velocities (measured at the points of first and second intersection with ∂B), and the normal vector to ∂B at the point of first intersection with ∂B . The normal vector at the second point of intersection is not registered; as will be seen later (Lemma 1), if the area of $B \setminus Q$ is small, then the difference between the normal vectors at these two points is also small. The measure $\nu_{Q,B}$ describes the distribution of triples.

DEFINITION 1. We say that a sequence of sets $\{Q_m, m = 1, 2, \dots\}$ represents a rough body if it has the following properties:

- M1. $Q_m \subset B$ and $\text{Area}(B \setminus Q_m) \rightarrow 0$ as $m \rightarrow \infty$.
- M2. The sequence of measures $\nu_{Q_m,B}$ weakly converges.

Two sequences of such sets are called equivalent if the corresponding limiting measures coincide. An equivalence class is called a body obtained by roughening B , or simply rough body, and denoted by \mathcal{B} , and the corresponding limiting measure is denoted by $\nu_{\mathcal{B}}$.

Note that the sets Q_m in this definition are not necessarily connected.

Remark. Since \mathbb{T}^3 is compact and the full measure of \mathbb{T}^3 satisfies $\nu_{Q,B}(\mathbb{T}^3) \leq 2\pi|\partial B|$, one concludes that the set of measures $\{\nu_{Q,B}\}$, with fixed B , is weakly precompact. That is, any sequence of measures $\{\nu_{Q_m,B}\}$ contains a weakly converging subsequence. In this sense one can say that a sequence, satisfying only condition M1, can represent more than one rough body.

We would also like to mention that, first, two rough bodies obtained one from another by translation are identified, according to our definition. Second, particle scattering on \mathcal{B} in a small neighborhood of $\xi \in \partial B$ can be detected if ξ is an extreme point of B , and cannot otherwise. Indeed, if ξ is an extreme point of B , then the scattering is described by the restriction of $\nu_{\mathcal{B}}$ on $\mathbb{T}^2 \times \mathcal{N}_{n(\xi)}$, with $\mathcal{N}_{n(\xi)}$ being a small neighborhood of $n(\xi)$ in S^1 . If, otherwise, ξ is not an extreme point of B , that is, belongs to an open linear segment contained in ∂B , the scattering can be determined only on the whole segment.

Actually, from the viewpoint of applications to the problems of optimal resistance in *homogeneous* and *rarefied* media (see section 4, containing concluding remarks and applications), these drawbacks are not so serious. Indeed, resistance of a body is invariant under translations (due to homogeneity). Besides, if the boundary of a body contains a linear segment, then one does not need to know scattering at each point of the segment; it suffices to know it on the whole segment (due to homogeneity and rarefaction).

The definition of a rough body could be made in a slightly different way, basing it on measures defined on $S^1 \times S^1 \times \partial B$. In that case the triple (v, v^+, ξ) should be registered, with ξ being the point of first intersection with ∂B . That definition would allow one to register particle scattering at each point of ∂B and to distinguish between bodies obtained by translation one from another. However, we prefer to adopt the former definition, since it seems to us mathematically more transparent and makes the arguments a bit easier.

1.3. Examples. Sometimes it is convenient to use another representation of the measure $\nu_{\mathcal{B}}$. Namely, consider the change of coordinates $(v, v^+, n) \mapsto (\varphi, \varphi^+, n)$, where $\varphi = \text{Arg } v - \text{Arg } n$, $\varphi^+ = \text{Arg } v^+ - \text{Arg } n$. Here $\text{Arg } v$ is the angle between a fixed vector and v measured, say, clockwise from this vector to v . If $(v, v^+, n) \in \text{spt } \nu_{\mathcal{B}}$, then φ and φ^+ belong to $[-\pi/2, \pi/2]$ modulo 2π . Introduce the shorthand notation $\square := [-\pi/2, \pi/2] \times [-\pi/2, \pi/2]$ and define the mapping $\varpi : \square \times S^1 \rightarrow \mathbb{T}^3$ by $\varpi(\varphi, \varphi^+, n) = (v, v^+, n)$. One has $\text{spt } \nu_{\mathcal{B}} \subset \varpi(\square \times S^1)$. Denote $\check{\nu}_{\mathcal{B}} := (\varpi^{-1})\# \nu_{\mathcal{B}}$. Sometimes this measure can be factorized: $\check{\nu}_{\mathcal{B}} = \eta_{\mathcal{B}} \otimes \tau_B$, where $\eta_{\mathcal{B}}$ is defined on \square

and τ_B is the surface measure on B ; so to say, the “roughness” is “homogeneous” along the body’s boundary. Consider several examples.

Example 1 (“smooth body”). The rough body represented by the sequence $Q_m = B$ is identified with B itself. The corresponding measure is $\check{\nu}_B = \eta_0 \otimes \tau_B$, where the measure η_0 has the density $\cos \varphi \cdot \delta(\varphi + \varphi^+)$; the support of η_0 is shown in Figure 2(b). In Figure 2(a), B is taken to be an ellipse.

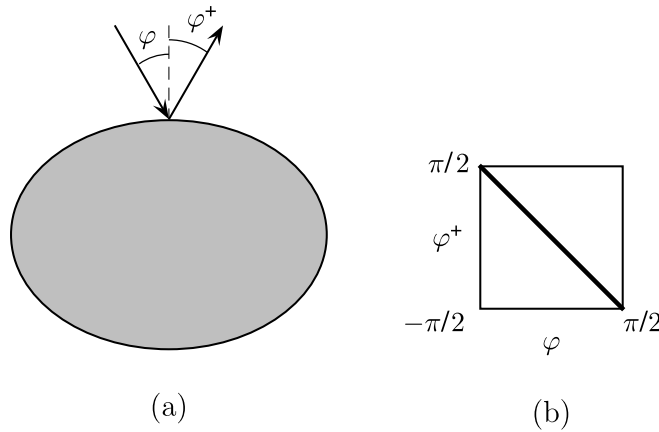


FIG. 2. A smooth body (a) and the support of the corresponding measure η_0 (b).

Example 2 (roughness formed by triangular hollows). Q_m is a $2m$ -polygon; the 270° angles alternate with the angles that are slightly smaller than 90° . All vertices corresponding to the angles smaller than 90° belong to ∂B . Any two sides that form a 270° angle are equal. The largest side length tends to zero as $m \rightarrow \infty$. Thus, the set Q_m is obtained by moving off m “hollows” from its convex hull, each of the hollows being an isosceles right triangle.

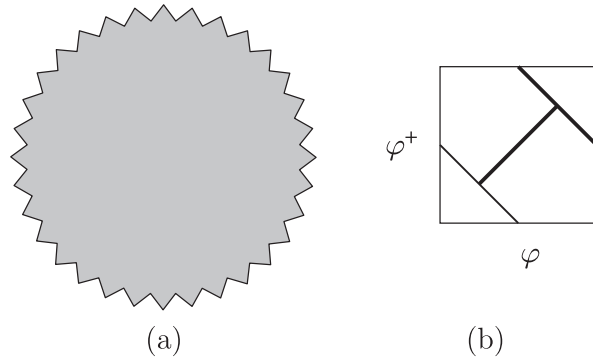


FIG. 3. A rough body with hollows being isosceles right triangles (a) and the support of the corresponding measure η_{∇} (b).

The corresponding measure is $\check{\nu}_B = \eta_{\nabla} \otimes \tau_B$, where the measure η_{∇} has the density $\cos \varphi \cdot [\chi_{[-\pi/2, -\pi/4]}(\varphi) \delta(\varphi + \varphi^+ + \frac{\pi}{2}) + \chi_{[-\pi/4, \pi/4]}(\varphi) \delta(\varphi - \varphi^+) + \chi_{[\pi/4, \pi/2]}(\varphi) \delta(\varphi + \varphi^+ - \frac{\pi}{2})] + |\sin \varphi| \cdot [\chi_{[-\pi/4, 0]}(\varphi) \delta(\varphi + \varphi^+ + \frac{\pi}{2}) - \chi_{[-\pi/4, \pi/4]}(\varphi) \delta(\varphi - \varphi^+) + \chi_{[0, \pi/4]}(\varphi) \delta(\varphi + \varphi^+ - \frac{\pi}{2})]$. Thus, the support of η_{∇} is the union of three segments; see Figure 3(b). The

middle segment $\varphi^+ = \varphi$ corresponds to double reflections, and the lateral segments, $\varphi^+ = -\varphi - \pi/2$ and $\varphi^+ = -\varphi + \pi/2$, correspond to single reflections, from the right or from the left side of a triangular hollow. In Figure 3(a), B is a circle.

Example 3 (roughness formed by rectangular hollows). The sets Q_m are obtained by removing a finite number of “rectangular hollows” from B . In other words, one has $Q_m = B \setminus (\cup_n \Omega_{m,n})$, where the removed sets $\Omega_{m,n}$ do not mutually intersect and each set $\partial\Omega_{m,n} \setminus \partial B$ is the union of three sides of a rectangle. The ratio (width)/(depth) of a hollow depends only on m and is denoted by h_m . Denote by $l_m = |\partial B \setminus \cup_n(\partial\Omega_{m,n})|/|\partial B|$ the relative length of the part of boundary ∂B not covered by hollows. We assume that $\lim_{m \rightarrow \infty} h_m = 0 = \lim_{m \rightarrow \infty} l_m$. In Figure 4(a), B is a square.

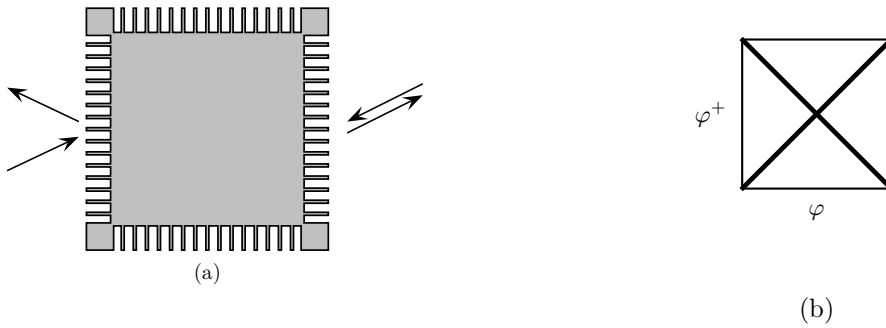


FIG. 4. A rough body where hollows are “thin rectangles” (a); the support of the corresponding measure η_{\square} (b).

The measure $\check{\nu}_B$ equals $\check{\nu}_B = \eta_{\square} \otimes \tau_B$. The density of the measure η_{\square} equals $\frac{1}{2} \cos \varphi \cdot (\delta(\varphi + \varphi^+) + \delta(\varphi - \varphi^+))$, and the support is the union of two diagonals, $\varphi^+ = \varphi$ and $\varphi^+ = -\varphi$; see Figure 4(b). The particles with an even (odd) number of reflections contribute to the first (second) diagonal.

1.4. Main theorem. According to Definition 1, each rough body is identified with a measure on \mathbb{T}^3 . The question is, What is the set of these measures? The following definition and theorem give the answer.

Let us first introduce some notation: $\pi_{v,n} : \mathbb{T}^3 \rightarrow \mathbb{T}^2$, $\pi_n : \mathbb{T}^3 \rightarrow S^1$, etc., are projections onto the corresponding subspaces: $\pi_{v,n}(v, v^+, n) = (v, n)$, $\pi_n(v, v^+, n) = n$, etc.; $\pi_d : \mathbb{T}^3 \rightarrow \mathbb{T}^3$ is the symmetry with respect to the plane $v = v^+$, that is, $\pi_d(v, v^+, n) = (v^+, v, n)$; $z_+ = \max\{0, z\}$ is the positive part of $z \in \mathbb{R}$; and u means Lebesgue measure on S^1 . Recall that τ_B is the surface measure on B and is defined on S^1 .

DEFINITION 2. We denote by \mathcal{M}_B the set of measures ν on \mathbb{T}^3 such that the following properties hold:

A1. The marginal measures $\pi_{v,n}^{\#} \nu$ and $\pi_{v^+,n}^{\#} \nu$ are

$$\pi_{v,n}^{\#} \nu = \langle v, n \rangle_+ \cdot u \otimes \tau_B, \quad \pi_{v^+,n}^{\#} \nu = \langle v^+, n \rangle_+ \cdot u \otimes \tau_B.$$

A2. $\pi_d^{\#} \nu = \nu$.

Denote also $\mathcal{M} = \cup_B \mathcal{M}_B$, the union being taken over all bounded convex bodies B .

Taking into account the Alexandrov theorem on characterization of surface measures, one concludes that \mathcal{M} is the set of measures ν on \mathbb{T}^3 such that

- (1) the marginal measure $\pi_n^\# \nu =: \tau$ satisfies the conditions
 - 1a. $\int_{S^1} n \, d\tau(n) = 0$;
 - 1b. for any $v \in S^1$ it holds that $\int_{S^1} \langle n, v \rangle^2 \, d\tau(n) \neq 0$;
- (2) the marginal measures $\pi_{v,n}^\# \nu$ and $\pi_{v^+,n}^\# \nu$ satisfy the conditions
 - 2a. $\pi_{v,n}^\# \nu = \langle v, n \rangle_+ \cdot u \otimes \tau$;
 - 2b. $\pi_{v^+,n}^\# \nu = \langle v^+, n \rangle_+ \cdot u \otimes \tau$.

Thus, these marginal measures coincide; the only difference is in the notation for the variables: v, n in case 2a and v^+, n in case 2b.

Now we can state the main theorem.

THEOREM. *The set of measures $\{\nu_B\}$, with B being all possible bodies obtained by roughening B , coincides with \mathcal{M}_B . Therefore, $\{\nu_B, B \text{ is a rough body}\} = \mathcal{M}$.*

In section 2, we formulate two auxiliary lemmas and using them prove the theorem. In section 3, the lemmas are proved. Section 4 contains concluding remarks and applications of the theorem to problems of optimal aerodynamic resistance. Appendices A and B contain proofs of some auxiliary technical results.

2. Statement of auxiliary lemmas and proof of theorem.

2.1. Statement of Lemma 1. Fix a bounded convex body B . Two points $\xi_1, \xi_2 \in \partial B$, $\xi_1 \neq \xi_2$, divide the curve ∂B into two arcs. Denote by $l(\xi_1, \xi_2)$ the length of the smallest arc and denote

$$c = c_B := \inf_{\substack{\xi_1, \xi_2 \in \partial B \\ \xi_1 \neq \xi_2}} \frac{|\xi_1 - \xi_2|}{l(\xi_1, \xi_2)}$$

one obviously has $0 < c < 1$.

Let $Q \subset B$; denote

$$\overline{|\xi - \xi^+|}_{Q,B} := \iint_{(\partial B \times S^1)_+} |\xi - \xi_{Q,B}^+(\xi, v)| \, d\mu(\xi, v)$$

and

$$\overline{|n - n^+|}_{Q,B} := \iint_{(\partial B \times S^1)_+} |n(\xi) - n(\xi_{Q,B}^+(\xi, v))| \, d\mu(\xi, v).$$

LEMMA 1. (a) *The following holds true:*

$$\overline{|\xi - \xi^+|}_{Q,B} \leq 2\pi \cdot \text{Area}(B \setminus Q).$$

(b) *For sufficiently small $\text{Area}(B \setminus Q)$,¹ one has*

$$\overline{|n - n^+|}_{Q,B} \leq \frac{2\pi\sqrt{8\pi}}{\sqrt{c}} \sqrt{\text{Area}(B \setminus Q)}.$$

2.2. Statement of Lemma 2. Let us first introduce the notion of a hollow.

DEFINITION 3. *Let $\Omega \subset \mathbb{R}^2$ be a closed bounded set with piecewise smooth boundary and $I \subset \partial\Omega$, where the following hold:*

- (i) *I is an interval contained in a straight line $\langle x, n \rangle = a$.*
- (ii) *$\Omega \setminus I$ is contained in the open half-plane $\langle x, n \rangle < a$. Here n is a fixed unit vector.*

Then the pair (Ω, I) is called a hollow oriented by n , or just an n -hollow.

¹That is, it is smaller than a positive value depending only on B .

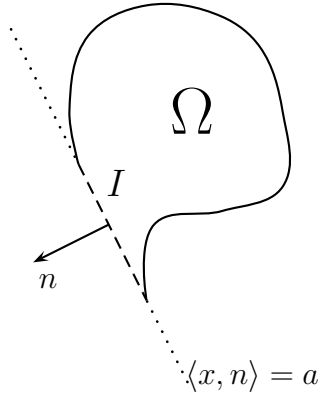


FIG. 5. A hollow.

In Figure 5 and in what follows, I is shown by a dashed line, and $\partial\Omega \setminus I$ is shown by a solid line.

Define the measure $\tilde{\mu}_I$ on $I \times S^1$ by $d\tilde{\mu}_I(\xi, v) = \frac{\langle n, v \rangle_+}{|I|} d\xi dv$, where $|I|$ means the length of I . Obviously, $\tilde{\mu}_I$ is supported on the set $(I \times S^1)_+ := \{(\xi, v) \in I \times S^1 : \langle n, v \rangle \geq 0\}$. Define the one-to-one mapping $(\xi, v) \mapsto (\Xi_{\Omega, I}^+(\xi, v), V_{\Omega, I}^+(\xi, v))$ of a full measure subset of $(I \times S^1)_+$ onto itself. Namely, consider the billiard in Ω . Let $(\xi, v) \in (I \times S^1)_+$; consider the billiard particle starting at the point ξ with the velocity $-v$. It makes several reflections from $\partial\Omega \setminus I$ and then reflects from I again, at a point $\Xi^+ = \Xi_{\Omega, I}^+(\xi, v)$. The velocity immediately before this reflection is denoted by $V^+ = V_{\Omega, I}^+(\xi, v)$. The mapping so defined preserves the measure $\tilde{\mu}_I$ and is an involution, that is, coincides with its inverse.

One can give an equivalent definition based on the mapping $\xi_{Q, B}^+(\xi, v), v_{Q, B}^+(\xi, v)$ just defined in section 1.2. Take a set Q such that Ω is a connected component of $\text{conv } Q \setminus Q$ and I is a connected component of $\partial(\text{conv } Q) \setminus \partial Q$. For $(\xi, v) \in (I \times S^1)_+$, let by definition $(\Xi_{\Omega, I}^+(\xi, v), V_{\Omega, I}^+(\xi, v)) := (\xi_{Q, \text{conv } Q}^+(\xi, v), v_{Q, \text{conv } Q}^+(\xi, v))$. This definition does not depend on the choice of Q .

DEFINITION 4. Let (Ω, I) be a hollow. The measure $\eta_{\Omega, I}$ on $\mathbb{T}^2 = S^1 \times S^1$ is defined as follows. For a Borel set $A \subset \mathbb{T}^2$, put

$$\eta_{\Omega, I}(A) := \tilde{\mu}_I(\{(\xi, v) \in (I \times S^1)_+ : (v, V_{\Omega, I}^+(\xi, v)) \in A\}).$$

We shall say that $\eta_{\Omega, I}$ is the measure generated by the hollow (Ω, I) .

Here we use the notation $\pi_v, \pi_{v^+} : \mathbb{T}^2 \rightarrow S^1$ for the projections onto the subspaces $\{v\}$ and $\{v^+\}$, respectively; $\pi_v(v, v^+) = v, \pi_{v^+}(v, v^+) = v^+$. We also denote by π_d the symmetry with respect to the diagonal $v = v^+$; $\pi_d(v, v^+) = (v^+, v)$.

DEFINITION 5. Denote by Λ_n the set of measures η on \mathbb{T}^2 such that

- (i) $d\pi_v^\# \eta(v) = \langle v, n \rangle_+ dv, d\pi_{v^+}^\# \eta(v^+) = \langle v^+, n \rangle_+ dv^+$;
- (ii) $\pi_d^\# \eta = \eta$.

Any measure $\eta_{\Omega, I}$ generated by an n -hollow belongs to Λ_n . Indeed, for any $A \subset S^1$ one has $\pi_v^\# \eta_{\Omega, I}(A) = \eta_{\Omega, I}(A \times S^1) = \tilde{\mu}_I(\{(\xi, v) \in (I \times S^1)_+ : v \in A\}) = \frac{1}{|I|} \iint_{I \times A} \langle n, v \rangle_+ d\xi dv = \int_A \langle n, v \rangle_+ dv$. This proves the first equality in (i).

Similarly, one has $\pi_{v^+}^\# \eta_{\Omega, I}(A) = \eta_{\Omega, I}(S^1 \times A) = \tilde{\mu}_I(\{(\xi, v) \in (I \times S^1)_+ : V_{\Omega, I}^+(\xi, v) \in A\})$. Since the mapping $(\xi, v) \mapsto (\Xi_{\Omega, I}^+(\xi, v), V_{\Omega, I}^+(\xi, v))$ preserves the measure, one gets the

value $\tilde{\mu}_I(\{(\xi, v) \in (I \times S^1)_+ : v \in A\})$, which in turns equals $\int_A \langle n, v \rangle_+ dv$. This proves the second equality in (i). Finally, the relation (ii) for $\eta_{\Omega, I}$ is a simple consequence of involutive and measure preserving properties of the mapping $(\xi, v) \mapsto (\Xi_{\Omega, I}^+, V_{\Omega, I}^+)$.

LEMMA 2. *The set of measures generated by n-hollows is weakly dense in Λ_n .*

2.3. Proof of the direct statement of theorem. Here we prove that for any body \mathcal{B} obtained by roughening B it holds that $\nu_{\mathcal{B}} \in \mathcal{M}_B$.

Let $Q \subset B$; define the measure $\nu'_{Q, B}$ on \mathbb{T}^3 by

$$\nu'_{Q, B}(A) := \mu \left(\{(\xi, v) \in (\partial B \times S^1)_+ : (v, v_{Q, B}^+(\xi, v), n(\xi_{Q, B}^+(\xi, v))) \in A\} \right),$$

where A is an arbitrary Borel subset of \mathbb{T}^3 . Thus, the definition of both $\nu_{Q, B}$ and $\nu'_{Q, B}$ is based on observations of vector triples (v, v^+, n) and (v, v^+, n^+) , respectively. Here n and n^+ are the outer normals to ∂B at the points where the particle gets in B and gets out of B . The measures $\nu_{Q, B}$ and $\nu'_{Q, B}$ have the following properties:

- (1) $\pi_{v, n}^\# \nu_{Q, B} = \langle v, n \rangle_+ \cdot u \otimes \tau_B,$
- (2) $\pi_{v^+, n^+}^\# \nu'_{Q, B} = \langle v^+, n^+ \rangle_+ \cdot u \otimes \tau_B,$
- (3) $\pi_d^\# \nu_{Q, B} = \nu'_{Q, B}.$

Consider a sequence $\{Q_m\}$ representing \mathcal{B} ; let us show that $\nu_{Q_m, B} - \nu'_{Q_m, B}$ weakly converges to zero as $m \rightarrow \infty$. It is enough to prove that for any continuous function f on \mathbb{T}^3 it holds that

$$(4) \quad \int_{\mathbb{T}^3} f(v, v^+, n) d\nu_{Q_m, B}(v, v^+, n) - \int_{\mathbb{T}^3} f(v, v^+, n^+) d\nu'_{Q_m, B}(v, v^+, n^+) \xrightarrow{m \rightarrow \infty} 0.$$

Taking into account the formulas for a change of variables

$$\int_{\mathbb{T}^3} f(v, v^+, n) d\nu_{Q, B}(v, v^+, n) = \int_{(\partial B \times S^1)_+} f(v, v_{Q, B}^+(\xi, v), n(\xi)) d\mu(\xi, v)$$

and

$$\int_{\mathbb{T}^3} f(v, v^+, n^+) d\nu'_{Q, B}(v, v^+, n^+) = \int_{(\partial B \times S^1)_+} f(v, v_{Q, B}^+(\xi, v), n(\xi_{Q, B}^+(\xi, v))) d\mu(\xi, v),$$

formula (4) takes the form

$$(5) \quad \lim_{m \rightarrow \infty} \int_{(\partial B \times S^1)_+} \left[f(v, v_{Q_m, B}^+(\xi, v), n(\xi_{Q_m, B}^+(\xi, v))) - f(v, v_{Q_m, B}^+(\xi, v), n(\xi)) \right] d\mu(\xi, v) = 0.$$

According to Lemma 1, the difference $n(\xi_{Q_m, B}^+(\xi, v)) - n(\xi)$ converges to zero in mean, and hence it converges to zero in measure; therefore the difference

$$f(v, v_{Q_m, B}^+(\xi, v), n(\xi_{Q_m, B}^+(\xi, v))) - f(v, v_{Q_m, B}^+(\xi, v), n(\xi))$$

also converges to zero in measure. It follows that formula (5) is true.

Thus, both $\nu_{Q_m, B}$ and $\nu'_{Q_m, B}$ weakly converge to ν_B . Substituting $Q = Q_m$ into formulas (1)–(3) and passing to limit as $m \rightarrow \infty$, one gets

$$\begin{aligned} \pi_{v, n}^\# \nu_B &= \langle v, n \rangle_+ \cdot u \otimes \tau_B, \\ \pi_{v^+, n}^\# \nu_B &= \langle v^+, n \rangle_+ \cdot u \otimes \tau_B, \\ \pi_d^\# \nu_B &= \nu_B, \end{aligned}$$

that is, $\nu_B \in \mathcal{M}_B$.

2.4. Proof of the inverse statement of theorem. Here it is proved that for any $\nu \in \mathcal{M}_B$ there exists a body \mathcal{B} obtained by roughening B such that $\nu_{\mathcal{B}} = \nu$. The proof is based on two statements.

STATEMENT 1. *Let B be a convex polygon. Then for any measure $\nu \in \mathcal{M}_B$ there exists a body \mathcal{B} obtained by roughening B such that $\nu_{\mathcal{B}} = \nu$.*

Proof. Let us enumerate the sides of the polygon B and denote by c_i the length of the i th side and by n_i the outer unit normal to this side. By δ_n , denote the probabilistic atomic measure on S^1 concentrated at $n \in S^1$, that is, $\delta_n(n) = 1$. The surface measure of B is $\tau_B = \sum c_i \delta_{n_i}$; this implies that any measure $\nu \in \mathcal{M}_B$ has the form $\nu = \sum c_i \eta_i \otimes \delta_{n_i}$, where $\eta_i \in \Lambda_{n_i}$.

According to Lemma 2, any measure η_i is the weak limit as $m \rightarrow \infty$ of measures $\eta_{\Omega_i^m, I_i^m}$ generated by a sequence of n_i -hollows (Ω_i^m, I_i^m) . Now take a sequence of sets Q_m such that $\text{conv } Q_m = B$ and each connected component of $B \setminus Q_m$ is the image of a set Ω_i^m under the composition of a homothety with positive ratio and a translation; additionally, the image of I_i^m under this transformation belongs to $(i$ th side of $B) \setminus \partial Q_m$. We also require that $\text{Area}(B \setminus Q_m) \rightarrow 0$ and $|(i$ th side of $B) \setminus \partial Q_m| =: c_i^m \rightarrow c_i$ as $m \rightarrow \infty$. In Appendix A it is shown how to construct such a sequence Q_m . The measure $\nu_{Q_m, B} = \tilde{\nu}_m + \sum_i \nu_i^m$ is the sum of the measure $\tilde{\nu}_m$ corresponding to reflections from $\partial B \cap \partial Q_m$ and the measures ν_i^m corresponding to particles getting into the ‘‘hollows on the i th side.’’ One has $\tilde{\nu}_m = \sum_i (c_i - c_i^m) \cdot \eta_0 \otimes \delta_{n_i}$ and $\nu_i^m = c_i^m \cdot \eta_{\Omega_i^m, I_i^m} \otimes \delta_{n_i}$. The norm of $\tilde{\nu}_m$ goes to zero and ν_i^m weakly converges to $c_i \eta_i \otimes \delta_{n_i}$ for any i ; it follows that $\nu_{Q_m, B}$ weakly converges to ν as $m \rightarrow \infty$. Therefore, the sequence Q_m represents a body \mathcal{B} obtained by roughening B , and $\nu_{\mathcal{B}} = \nu$. \square

STATEMENT 2. *For any measure $\nu \in \mathcal{M}_B$ there exist a sequence of convex polygons $B_k \subset B$ with $\text{Area}(B \setminus B_k) \rightarrow 0$ and a sequence of measures $\nu_k \in \mathcal{M}_{B_k}$ weakly converging to ν as $k \rightarrow \infty$.*

Proof. Consider a partition of the circumference S^1 into a finite number of arcs, $S^1 = \cup_i \mathcal{S}^i$. It induces the partition of ∂B into arcs $\partial B^i = \{\xi \in \partial B : n(\xi) \in \mathcal{S}^i\}$. Consider the polygon \check{B} inscribed into ∂B whose vertices are separation points of this partition. Denote by n_i the outer normal to the i th side of this polygon. Denote by s_{v_1, v_2} the operator of rotation on S^1 that takes v_1 to v_2 , and define the mapping $\Upsilon_i : \mathbb{T}^2 \times \mathcal{S}^i \rightarrow \mathbb{T}^2$ by $\Upsilon_i(v, v^+, n) = (s_{n, n_i} v, s_{n, n_i} v^+)$. Finally, consider the measure $\check{\nu} = \sum_i |b^i| \eta_i \otimes \delta_{n_i}$, where $|b^i|$ is the length of the i th side of the polygon, and the measure η_i on \mathbb{T}^2 is defined by $\eta_i(A) = \frac{1}{|\partial B^i|} \nu(\Upsilon_i^{-1}(A))$ for arbitrary Borel set $A \subset \mathbb{T}^2$. Here $|\partial B^i|$ is the length of the arc ∂B^i . One easily verifies that $\check{\nu}$ belongs to $\mathcal{M}_{\check{B}}$.

Now take a sequence of partitions of S^1 , $\{\mathcal{S}_k^i\}_i, k = 1, 2, \dots$, where the maximum arc length of a partition goes to zero as $k \rightarrow \infty$. Denote by $\{\partial B_k^i\}_i, k = 1, 2, \dots$, the sequence of induced partitions of ∂B , and take the sequence of polygons B_k generated by these partitions. One clearly has $\text{Area}(B \setminus B_k) \rightarrow 0$ and

$$(6) \quad \max_i \frac{|b_k^i|}{|\partial B_k^i|} \rightarrow 1 \quad \text{as } k \rightarrow \infty,$$

where $|b_k^i|$ is the length of the i th side of B_k . In the same way as above, one defines the mappings $\Upsilon_{ik} : \mathbb{T}^2 \times \mathcal{S}_k^i \rightarrow \mathbb{T}^2$ and the measures $\nu_k = \sum_i |b_k^i| \eta_{ik} \otimes \delta_{n_{ik}} \in \mathcal{M}_{B_k}$, where η_{ik} is given by $\eta_{ik}(A) := \frac{1}{|\partial B_k^i|} \nu(\Upsilon_{ik}^{-1}(A))$ and n_{ik} is the outer unit normal to the i th side of B_k .

It remains to show that ν_k weakly converges to ν . For any continuous function f on \mathbb{T}^3 one has

$$(7) \quad \begin{aligned} \iint_{\mathbb{T}^3} f(v, v^+, n) d\nu_k(v, v^+, n) &= \sum_i |b_k^i| \iint_{\mathbb{T}^2} f(v, v^+, n_{ik}) d\eta_{ik}(v, v^+) \\ &= \sum_i \frac{|b_k^i|}{|\partial B_k^i|} \iiint_{\mathbb{T}^2 \times \mathcal{S}_k^i} f(\Upsilon_{ik}(v, v^+, n), n_{ik}) d\nu(v, v^+, n). \end{aligned}$$

For each k define the mapping from \mathbb{T}^3 to \mathbb{T}^3 by the relations $(v, v^+, n) \mapsto (\Upsilon_{ik}(v, v^+, n), n_{ik})$ if $n \in \mathcal{S}_k^i$. It uniformly converges to the identity mapping as $k \rightarrow \infty$; hence the function \tilde{f}_k , defined by the relations $\tilde{f}_k(v, v^+, n) := f(\Upsilon_{ik}(v, v^+, n), n_{ik})$ if $n \in \mathcal{S}_k^i$, uniformly converges to f as $k \rightarrow \infty$. From here and from (6) it follows that the right-hand side in (7) converges to $\iiint_{\mathbb{T}^3} f(v, v^+, n) d\nu(v, v^+, n)$ as $k \rightarrow \infty$. Thus, the convergence $\int f d\nu_k \rightarrow \int f d\nu$ is proved. \square

The inverse statement of the theorem follows from Statements 1 and 2. Indeed, let $\nu \in \mathcal{M}_B$. Using Statement 2, find a sequence of convex polygons $B_k \subset B$ and a sequence $\nu_k \in \mathcal{M}_{B_k}$ weakly converging to ν . According to Statement 1, each measure ν_k is generated by a rough body. Consider the sequence of sets $Q_{kl} \subset B_k, l = 1, 2, \dots$, representing this body, and then from all of these sequences choose a diagonal sequence $\tilde{Q}_k = Q_{kl_k}$ such that the corresponding sequence of measures $\nu_{\tilde{Q}_k, B}$ weakly converges to ν and $\text{Area}(B \setminus \tilde{Q}_k)$ goes to zero as $k \rightarrow \infty$. The sequence \tilde{Q}_k represents a body \mathcal{B} obtained by roughening B and $\nu_{\mathcal{B}} = \nu$.

3. Proof of the lemmas.

3.1. Proof of Lemma 1. Consider the billiard in $\mathbb{R}^2 \setminus Q$. For $(\xi, v) \in (\partial B \times S^1)_+$, denote by $\tau(\xi, v)$ the time the billiard trajectory with the initial data $\xi, -v$ spends in $B \setminus Q$. In particular, if $\xi \in \partial B \cap \partial Q$, then one has $\tau(\xi, v) = 0$.

Denote by D the set of points $(x, w) \in (B \setminus Q) \times S^1$ that are accessible from $(\partial B \times S^1)_+$; that is, there exists $(\xi, v) \in (\partial B \times S^1)_+$ such that the billiard particle with the data $\xi, -v$ at the zero moment of time at some moment $0 \leq t \leq \tau(\xi, v)$ will pass through x with the velocity w . This description defines the change of coordinates in $D : (\xi, v, t) \mapsto (x, w); (\xi, v) \in (\partial B \times S^1)_+, t \in [0, \tau(\xi, v)]$, and the element of phase volume $d^2x dw$ in the new coordinates takes the form $d\mu(\xi, v) dt$. Hence, the phase volume of D equals $\iiint_D d^2x dw = \iint_{(\partial B \times S^1)_+} \tau(\xi, v) d\mu(\xi, v)$. Taking into account that $D \subset (B \setminus Q) \times S^1$ and the phase volume of $(B \setminus Q) \times S^1$ equals $2\pi \cdot \text{Area}(B \setminus Q)$, one gets

$$(8) \quad \iint_{(\partial B \times S^1)_+} \tau(\xi, v) d\mu(\xi, v) \leq 2\pi \cdot \text{Area}(B \setminus Q).$$

This is in fact a simple modification of the well-known *mean free path* formula (see, e.g., [19]).

One has $\tau(\xi, v) \geq |\xi - \xi_{\tilde{Q}, B}^+(\xi, v)|$: the time the particle spends in $B \setminus Q$ exceeds the distance between the initial and final points of the trajectory. This inequality and (8) imply (a).

The points ξ and $\xi_{Q,B}^+(\xi, v)$ divide the curve ∂B into two arcs; denote by $\gamma(\xi, v)$ the shortest one. One has $|\gamma(\xi, v)| = l(\xi, \xi_{Q,B}^+(\xi, v))$, and therefore $|\xi - \xi_{Q,B}^+(\xi, v)| \geq c|\gamma(\xi, v)|$. It follows that

$$(9) \quad c \iint_{(\partial B \times S^1)_+} |\gamma(\xi, v)| d\mu(\xi, v) \leq \iint_{(\partial B \times S^1)_+} |\xi - \xi_{Q,B}^+(\xi, v)| d\mu(\xi, v) \leq 2\pi \cdot \text{Area}(B \setminus Q).$$

Let $\varrho(y)$ be a natural parametrization of the curve ∂B , $\varrho : [0, |\partial B|] \rightarrow \partial B$. By $f(y)$ denote the measure of the values (ξ, v) such that the interval $\gamma(\xi, v)$ contains the point $\varrho(y)$; that is, $f(y) := \iint_{(\partial B \times S^1)_+} \mathbb{I}(\varrho(y) \in \gamma(\xi, v)) d\mu(\xi, v)$. Making a change of variables in the integral in the left-hand side of (9), one gets

$$\iint_{(\partial B \times S^1)_+} |\gamma(\xi, v)| d\mu(\xi, v) = \int_0^{|\partial B|} f(y) dy,$$

and therefore

$$(10) \quad \int_0^{|\partial B|} f(y) dy \leq \frac{2\pi}{c} \text{Area}(B \setminus Q).$$

One easily sees that $|f(y_1) - f(y_2)| \leq 4|y_1 - y_2|$ for any y_1 and y_2 and $f(y) \geq 0$. From here and from (10) it follows that for sufficiently small $\text{Area}(B \setminus Q)$ (namely, for $\text{Area}(B \setminus Q) \leq c|\partial B|^2/(2\pi)$) it holds that $f(y) \leq \sqrt{8\pi/c} \sqrt{\text{Area}(B \setminus Q)}$.

Recall that $\text{Arg}(v)$ is the angle the vector $v \neq 0$ forms with a fixed vector v_0 ; the angle is measured clockwise from v_0 to v and is defined modulo 2π . Introduce the shorthand notation $\xi^+ := \xi_{Q,B}^+(\xi, v)$, and denote by $\Delta \text{Arg}(\xi, v)$ the smallest in modulus of the values $\text{Arg}(n(\xi^+)) - \text{Arg}(n(\xi))$. In other words, $\Delta \text{Arg}(\xi, v)$ equals the smallest of the values

$$\int_{\gamma(\xi, v)} |d \text{Arg}(n_{\xi'})|, \quad \int_{\partial B \setminus \gamma(\xi, v)} |d \text{Arg}(n_{\xi'})|.$$

Taking into account that $|n(\xi^+) - n(\xi)| \leq |\Delta \text{Arg}(\xi, v)|$, one gets that

$$|n(\xi^+) - n(\xi)| \leq \int_{\gamma(\xi, v)} |d \text{Arg}(n_{\xi'})|,$$

and therefore

$$\overline{|n - n^+|}_{Q,B} \leq \iint_{(\partial B \times S^1)_+} \left(\int_{\gamma(\xi, v)} |d \text{Arg}(n_{\xi'})| \right) d\mu(\xi, v).$$

Making a change of variables in this integral, one obtains

$$\overline{|n - n^+|}_{Q,B} \leq \int_0^{|\partial B|} f(y) |d \text{Arg}(n_{\varrho(y)})| \leq 2\pi \sqrt{8\pi/c} \sqrt{\text{Area}(B \setminus Q)}.$$

Thus, (b) is also proved.

3.2. Proof of Lemma 2. Fix $n \in S^1$ and $m \in \mathbb{N}$. Let σ be an involutive permutation of $\{1, \dots, m\}$, that is, $\sigma^2 = \text{id}$. Divide the half-circumference $S_n^1 := \{v \in S^1 : \langle v, n \rangle \geq 0\}$ into m arcs $\mathcal{S}_{n,m}^1 = \mathcal{S}_n^1, \dots, \mathcal{S}_{n,m}^m = \mathcal{S}_n^m$ numbered clockwise such that, for any i , $\int_{\mathcal{S}_n^i} \langle v, n \rangle dv = 2/m$. For the sake of brevity we omit the subscript m when no confusion can arise.

DEFINITION 6. A measure η is called a (σ, n) -measure if $\eta \in \Lambda_n$ and $\text{spt } \eta \subset \cup_{i=1}^m (\mathcal{S}_n^i \times \mathcal{S}_n^{\sigma(i)})$, and therefore, for any i it holds that $\eta(\mathcal{S}_n^i \times \mathcal{S}_n^{\sigma(i)}) = 2/m$.

PROPOSITION 1. For any measure $\eta \in \Lambda_n$ there exists a sequence of involutive permutations σ_k on $\{1, \dots, m_k\}$, $k = 1, 2, \dots$, such that m_k tends to infinity and any sequence of (σ_k, n) -measures weakly converges to η as $k \rightarrow \infty$.

PROPOSITION 2. Let σ be an involutive permutation on $\{1, \dots, m\}$. Then the distance (in variation) between the set of measures generated by n -hollows and the set of (σ, n) -measures does not exceed $16/m$. In other words, whatever $\varepsilon > 0$, there exist a (σ, n) -measure η and an n -hollow (Ω, I) such that $\|\eta_{\Omega, I} - \eta\| < 16/m + \varepsilon$; here the norm means variation of measure.

This distance actually equals zero, but we need only the (weaker) claim of Proposition 2.

Lemma 2 follows from Propositions 1 and 2. Indeed, let $\eta \in \Lambda_n$. First, choose the sequence of permutations σ_k , according to Proposition 1, and then, using Proposition 2, for every k choose an n -hollow (Ω_k, I_k) such that the distance from η_{Ω_k, I_k} to the set of (σ_k, n) -measures does not exceed $17/m_k$. The sequence of chosen measures η_{Ω_k, I_k} weakly converges to η .

3.3. Proof of Proposition 1. Introduce on S_n^1 the angular coordinate $\varphi = \text{Arg } v - \text{Arg } n$; that is, φ changes between $-\pi/2$ and $\pi/2$ and increases clockwise. With this notation, to the arcs $\mathcal{S}_{n,m}^i$ correspond the segments $J_m^i = [\arcsin(-1 + 2(i-1)/m), \arcsin(-1 + 2i/m)]$. Define the measure λ on $[-\pi/2, \pi/2]$ by $d\lambda(\varphi) = \cos \varphi d\varphi$, and denote by Λ the set of measures η on $\square := [-\pi/2, \pi/2] \times [-\pi/2, \pi/2]$ such that (a) $\pi_\varphi^\# \eta = \lambda = \pi_{\varphi^+}^\# \eta$ and (b) $\pi_d^\# \eta = \eta$. Here $\pi_\varphi, \pi_{\varphi^+}$, and π_d are defined by $\pi_\varphi(\varphi, \varphi^+) = \varphi, \pi_{\varphi^+}(\varphi, \varphi^+) = \varphi^+,$ and $\pi_d(\varphi, \varphi^+) = (\varphi^+, \varphi)$. Reformulating Definition 6, we shall say that η is a σ -measure if $\eta \in \Lambda$ and $\text{spt } \eta \subset \cup_{i=1}^m (J_m^i \times J_m^{\sigma(i)})$. Notice that in the new notation the objects no longer depend on n : we write Λ instead of Λ_n , σ -measure instead of (σ, n) -measure, and hollow instead of n -hollow.

In this notation, Proposition 1 can be reformulated as follows: for any measure $\eta \in \Lambda$ there exists a sequence of involutive permutations σ_k on $\{1, \dots, m_k\}$, $k = 1, 2, \dots$, such that m_k tends to infinity and any sequence of σ_k -measures weakly converges to η as $k \rightarrow \infty$.

The idea of the proof is as follows. First, η is approximated by means of a rational matrix, and then this matrix is approximated by means of a larger matrix generated by a permutation.

Consider the partition of \square into smaller rectangles $\square_k^{ij} = J_k^i \times J_k^j, i, j = 1, \dots, k$. Choose rational nonnegative numbers c_k^{ij} such that $c_k^{ij} = c_k^{ji}, \sum_j c_k^{ij} = 2/k$ for any i , and $|\eta(\square_k^{ij}) - c_k^{ij}| \leq k^{-3}$ for any i and j . To do so, it suffices to take positive rational values \hat{c}_k^{ij} such that $\eta(\square_k^{ij}) - k^{-4} \leq \hat{c}_k^{ij} \leq \eta(\square_k^{ij})$ for $i > j$ and put $c_k^{ij} = \hat{c}_k^{ji}$ for $i < j$ and $c_k^{ii} = 2/k - \sum_{j \neq i} \hat{c}_k^{ij}$ for $i = j$. One has $\eta(J_k^i \times [-\pi/2, \pi/2]) = \sum_{j=1}^k \eta(\square_k^{ij}) = 2/k$; hence $c_k^{ii} - \eta(\square_k^{ii}) = \sum_{j \neq i} (\eta(\square_k^{ij}) - c_k^{ij}) \in [0, (k-1) \cdot k^{-4}] \subset [0, k^{-3}]$.

Any sequence of measures η_k satisfying the conditions $\eta_k(\square_k^{ij}) = c_k^{ij}, 1 \leq i, j \leq k$, weakly converges to η . Indeed, for any continuous function f on \square it holds that

$$\int_{\square} f d\eta_k - \int_{\square} f d\eta = \sum_{i,j=1}^k \int_{\square_k^{ij}} f (d\eta_k - d\eta) \leq k^{-1} \max f \rightarrow 0$$

as $k \rightarrow \infty$.

To complete the proof, it suffices to find an integer $m_k > k$ and an involutive permutation σ_k of $\{1, \dots, m_k\}$ such that any σ_k -measure, η_k , satisfies the equalities $\eta_k(\square_k^{ij}) = c_k^{ij}$, $i, j = 1, \dots, k$. Choose a positive integer N such that all of the values $a_{ij} := N \cdot c_k^{ij}$ are integer. The obtained matrix $A = (a_{ij})_{i,j=1}^k$ is symmetric, and for any i the value $\sum_{j=1}^k a_{ij} = 2N/k$ is a fixed positive integer. In Appendix B it is shown that there exist square matrices $B_{ij} = (b_{ij}^{\mu\nu})_{\mu,\nu}$ of size $2N/k$ such that $B_{ij}^T = B_{ji}$, the sum of elements in any matrix B_{ij} equals a_{ij} , and the block matrix $D = (B_{ij})$ composed of these matrices has exactly one unit in each row and each column, and other elements are zeros.

D is a symmetric square matrix of size $2N$; denote its elements by d_{ij} . Put $m_k = 2N$, and define the mapping σ_k on $\{1, \dots, 2N\}$ in such a way that $d_{i\sigma_k(i)} = 1$ for any i . The so defined mapping σ_k is a permutation; it is involutive since the matrix D is symmetric. Moreover, if η_k is a σ_k -measure, then for any i and j it holds that $\eta_k(\square_k^{ij}) = N^{-1} \sum_{\mu,\nu} b_{ij}^{\mu\nu} = c_k^{ij}$. The proposition is proved.

3.4. Proof of Proposition 2.

1. Whatever the n -hollow (Ω, I) , one introduces the reference system (x_1, x_2) in such a way that n coincides with $(0, -1)$, and the interval I belongs to the straight line $x_2 = 0$ and contains the origin $O = (0, 0)$. Like in the proof of Proposition 1, introduce the coordinate $\varphi = \text{Arg } v - \text{Arg } n$ on S_n^1 . One has $v = -(\sin \varphi, \cos \varphi)$, $\varphi \in [-\pi/2, \pi/2]$. The definitions of the segments $J_m^i = J^i$, the measure λ , the set of measures Λ , and the σ -measure are seen in the beginning of the previous subsection. The mapping $(\xi, v) \mapsto V_{\Omega, I}^+(\xi, v)$ in the new coordinates ξ, φ is written as $(\xi, \varphi) \mapsto \varphi_{\Omega, I}^+(\xi, \varphi)$. Finally, define the measure μ_I on $I \times [-\pi/2, \pi/2]$ by $d\mu_I(\xi, \varphi) = \frac{\cos \varphi}{|I|} d\xi d\varphi$.

Denote $\square' = (\cup_{i=2}^{m-1} J^i) \times (\cup_{i=2}^{m-1} J^i)$, $\square_1 = J^1 \times [-\pi/2, \pi/2]$, $\square_2 = J^m \times [-\pi/2, \pi/2]$, $\square_3 = (\cup_{i=2}^{m-1} J^i) \times J^1$, and $\square_4 = (\cup_{i=2}^{m-1} J^i) \times J^m$. Thus, one has $\square \setminus \square' = \square_1 \cup \square_2 \cup \square_3 \cup \square_4$; see Figure 6.

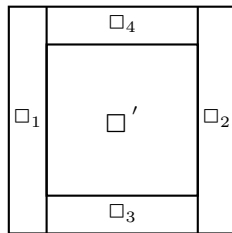


FIG. 6. Partition of the square into rectangles.

It suffices to construct a sequence of hollows $(\Omega_\varepsilon, I_\varepsilon)$, $\varepsilon > 0$, such that

- (P) for any $i \neq 1, m, \sigma(1), \sigma(m)$ the measure of the set of values $(\xi, \varphi) \in I_\varepsilon \times J^i$ such that $\varphi_{\Omega_\varepsilon, I_\varepsilon}^+(\xi, \varphi) \notin J^{\sigma(i)}$ goes to zero as $\varepsilon \rightarrow 0$.

Then, speaking of restrictions of measures on the subset \square' , one gets that the distance from the restrictions of measures $\eta_{\Omega_\varepsilon, I_\varepsilon}$ to the set of restrictions of σ -measures goes to zero as $\varepsilon \rightarrow 0$. On the other hand, for any measure $\eta \in \Lambda$ one has $\eta(\square_1) =$

$\eta(\square_2) = 2/m$, $\eta(\square_3) \leq 2/m$, $\eta(\square_4) \leq 2/m$, and hence $\eta(\square \setminus \square') \leq 8/m$; therefore the distance between the restrictions on $\square \setminus \square'$ of any two measures η_1 and η_2 from Λ does not exceed $16/m$: $\|\eta_1|_{\square \setminus \square'} - \eta_2|_{\square \setminus \square'}\| \leq 16/m$. It follows that the upper limit of distances from $\eta_{\Omega_\varepsilon, I_\varepsilon}$ to the set of σ -measures does not exceed $16/m$, and so Proposition 2 is proved.

2. The rest of this subsection is dedicated to the detailed description of the sequence of hollows $(\Omega_\varepsilon, I_\varepsilon)$ and to the proof of property (P) for them.

First, consider an auxiliary construction (see Figure 7). Take two different points F and F' above the line $l = \{x_2 = 0\}$, with $|OF| = 2 = |OF'|$. Denote by Φ and Φ' the angles the rays OF and OF' , respectively, formed with the vector $(0, 1)$. The angles are counted clockwise from $(0, 1)$. Thus, $F = 2(\sin \Phi, \cos \Phi)$ and $F' = 2(\sin \Phi', \cos \Phi')$. Assume, for further convenience, that F is situated on the left of F' ; thus, one has $-\pi/2 < \Phi < \Phi' < \pi/2$. (The case where F is situated on the right of F' is completely similar.) Select three positive numbers λ , λ' , and δ , and define two ellipses \mathcal{E} and \mathcal{E}' and two parabolas \mathcal{P} and \mathcal{P}' . The first ellipse has the foci O and F , the length of its large semiaxis is $\sqrt{1 + \lambda}$, of the small semiaxis, $\sqrt{\lambda}$, and the focal distance equals 2. The second ellipse has the foci O and F' , the lengths of its large and small semiaxes are $\sqrt{1 + \lambda'}$ and $\sqrt{\lambda'}$, respectively, and the focal distance is also 2. The parabolas \mathcal{P} and \mathcal{P}' have the foci F and F' , respectively, the common axis FF' , and the same focal distance δ . Thus, the parabolas are symmetric to each other with respect to the bisectrix of the triangle OFF' . The parameter δ is chosen sufficiently small, so that the point O lies in the exterior of both parabolas.

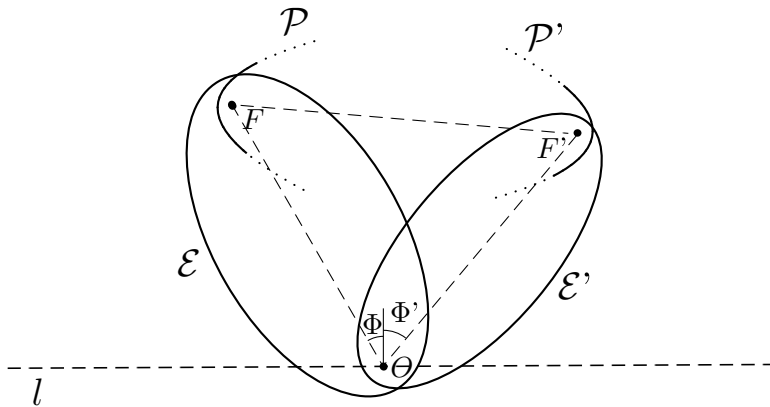


FIG. 7. Auxiliary construction.

In what follows, we shall distinguish between the billiard and *pseudobilliard* dynamics. The pseudobilliard dynamics is defined as follows. A particle starts at a point $(\xi, 0) \in l$ and moves with a velocity $(\sin \varphi, \cos \varphi)$ until it reflects from the interior side of \mathcal{E} . (Before the reflection it can intersect other curves \mathcal{E}' , \mathcal{P} , \mathcal{P}' , or even intersect \mathcal{E} from the outer side, without changing the velocity.) Then it moves again with constant velocity until it reflects from the interior side of \mathcal{P} . Then, in the same way, it reflects from the interior side of \mathcal{P}' , and then from the interior side of \mathcal{E}' , and, finally, it intersects l from above to below. Denote by $(\xi', 0)$ the point of intersection and by $-(\sin \varphi', \cos \varphi')$ the velocity at this point.

Consider the admissible set: the set of 7-tuples $(\varphi, \xi, \Phi, \Phi', \lambda, \lambda', \delta)$ such that all of the indicated reflections occur in the prescribed order. This set is open and nonempty.

Indeed, let $\delta(\Phi, \Phi')$ be the least of the values δ such that one of the parabolas (in fact, both of them simultaneously) passes through O . Put $\varphi = \Phi$, $\xi = 0$, and take arbitrary values $\lambda > 0$, $\lambda' > 0$, $-\pi/2 < \Phi < \Phi' < \pi/2$, $0 < \delta < \delta(\Phi, \Phi')$. The particle with initial data $\varphi = \Phi$, $\xi = 0$ first passes along the large semiaxis of \mathcal{E} , reflects from \mathcal{E} , returns along the same semiaxis, and reflects from \mathcal{P} . Then it moves with the velocity parallel to FF' , reflects from \mathcal{P}' , moves the large semiaxis of the ellipse \mathcal{E}' , reflects from it, and returns to O along the same semiaxis. Thus, the admissible set is nonempty. Under a small perturbation of the parameters $\varphi, \xi, \Phi, \Phi', \lambda, \lambda', \delta$, all of the reflections are maintained and the order of reflections remains the same. This implies that the admissible set is open.

This description determines the mapping $\varphi' = \varphi'(\varphi, \xi, \Phi, \Phi', \lambda, \lambda', \delta)$, $\xi' = \xi'(\varphi, \xi, \Phi, \Phi', \lambda, \lambda', \delta)^2$ from the admissible set to \mathbb{R}^2 . This mapping is infinitely differentiable. For $\varphi = \Phi$ and $\xi = 0$ one has

$$(11) \quad \varphi'(\Phi, 0, \Phi, \Phi', \lambda, \lambda', \delta) = \Phi'.$$

For $\xi = 0$ with arbitrary φ one has

$$(12) \quad \xi'(\varphi, 0, \Phi, \Phi', \lambda, \lambda', \delta) = 0,$$

and

$$\varphi'(\varphi, 0, \Phi, \Phi', \lambda, \lambda', \delta) \text{ does not depend on } \delta.$$

Indeed, a particle starting at O , after the reflection from \mathcal{E} passes through F , after reflecting from \mathcal{P} moves in parallel to FF' , after the reflection from \mathcal{P}' passes through F' , and, finally, after the reflection from \mathcal{E}' returns to O (see Figure 8). The initial and final velocities of the particle are, respectively, $(\sin \varphi, \cos \varphi)$ and $-(\sin \varphi', \cos \varphi')$. Denoting by α and α' the angles the second and fourth segments of the (5-segment) trajectory form, respectively, with OF and OF' , one has $\alpha = \alpha'$. The angle α is a function of φ , and φ' is a function of α' ; these functions depend only on the parameters of the ellipses \mathcal{E} and \mathcal{E}' , respectively, and do not depend on the parameter δ determining the shape of parabolas.

Using properties of ellipses, one derives the formulas connecting φ , α , and $\varphi' = \varphi'(\varphi, 0, \Phi, \Phi', \lambda, \lambda', \delta)$:

$$(13) \quad \sin(\varphi - \Phi) = \frac{\lambda \sin \alpha}{2 + \lambda - 2 \cos \alpha \sqrt{1 + \lambda}}, \quad \sin(\varphi' - \Phi') = -\frac{\lambda' \sin \alpha}{2 + \lambda' - 2 \cos \alpha \sqrt{1 + \lambda'}}.$$

It follows that

$$(14) \quad \left. \frac{\partial \varphi'}{\partial \varphi} \right|_{\substack{\varphi=\Phi \\ \xi=0}} = - \left(\frac{\sqrt{\lambda'}}{1 + \sqrt{\lambda'}} \frac{1 + \sqrt{\lambda}}{\sqrt{\lambda}} \right)^2.$$

With fixed $\Phi, \Phi', \lambda, \lambda'$, and δ the mapping $\varphi'(\varphi, \xi)$, $\xi'(\varphi, \xi)$ preserves the measure, $\cos \varphi d\varphi d\xi = \cos \varphi' d\varphi' d\xi'$, and hence

$$\cos \varphi = \pm \cos \varphi' \begin{vmatrix} \frac{\partial \varphi'}{\partial \varphi} & \frac{\partial \varphi'}{\partial \xi} \\ \frac{\partial \xi'}{\partial \varphi} & \frac{\partial \xi'}{\partial \xi} \end{vmatrix}.$$

²Note that throughout this paper the sign ' (prime) never means derivation.

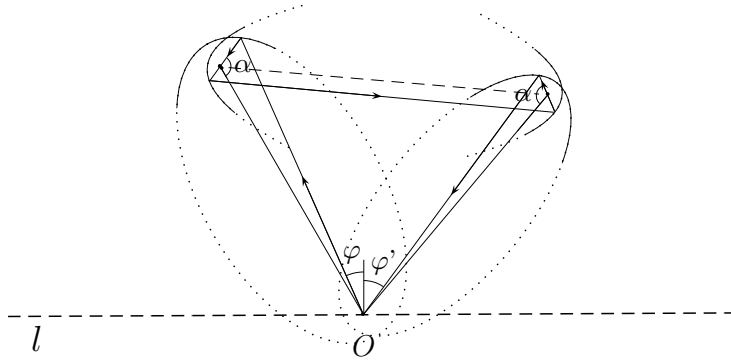


FIG. 8. Pseudobilliard dynamics.

Using (12), one gets that $\left. \frac{\partial \xi'}{\partial \varphi} \right|_{\xi=0} = 0$, and hence

$$\left| \begin{array}{cc} \frac{\partial \varphi'}{\partial \varphi} & \frac{\partial \varphi'}{\partial \xi} \\ \frac{\partial \xi'}{\partial \varphi} & \frac{\partial \xi'}{\partial \xi} \end{array} \right|_{\xi=0} = \frac{\partial \varphi'}{\partial \varphi} \frac{\partial \xi'}{\partial \xi} \Big|_{\xi=0};$$

therefore

$$(15) \quad \cos \varphi = \pm \cos \varphi' \frac{\partial \varphi'}{\partial \varphi} \frac{\partial \xi'}{\partial \xi} \Big|_{\xi=0}.$$

Putting $\varphi = \Phi$ and $\xi = 0$, and taking into account (11), (14), and (15), one gets

$$(16) \quad \cos \Phi = \pm \cos \Phi' \left(\frac{\sqrt{\lambda'}}{1 + \sqrt{\lambda'}} \frac{1 + \sqrt{\lambda}}{\sqrt{\lambda}} \right)^2 \frac{\partial \xi'}{\partial \xi} \Big|_{\substack{\varphi=\Phi \\ \xi=0}}.$$

Define the positive continuous functions $\lambda(\Phi')$ and $\lambda'(\Phi)$ by the relations

$$(17) \quad \left(\frac{\sqrt{\lambda}}{1 + \sqrt{\lambda}} \right)^2 = \frac{1}{2} \cos \Phi', \quad \left(\frac{\sqrt{\lambda'}}{1 + \sqrt{\lambda'}} \right)^2 = \frac{1}{2} \cos \Phi;$$

then one has

$$(18) \quad \left| \frac{\partial \xi'}{\partial \xi} \Big|_{\substack{\varphi=\Phi; \lambda=\lambda(\Phi') \\ \xi=0; \lambda'=\lambda'(\Phi)}} \right| = 1.$$

Additionally, taking into account (14) and (17), one gets

$$(19) \quad \frac{\cos \Phi'}{\cos \Phi} \frac{\partial \varphi'}{\partial \varphi} \Big|_{\substack{\varphi=\Phi; \lambda=\lambda(\Phi') \\ \xi=0; \lambda'=\lambda'(\Phi)}} = -1.$$

Recall that $\varphi' = \varphi'(\varphi, 0, \Phi, \Phi', \lambda, \lambda')$; that is, the restriction of the function φ' to the subspace $\xi = 0$ does not depend on δ . Hence the function $\left. \frac{\partial \varphi'}{\partial \varphi} \right|_{\xi=0}$ and, by formula (15), the function $\left. \frac{\partial \xi'}{\partial \xi} \right|_{\xi=0}$ also do not depend on δ . Put $\Phi_0 = \arcsin(1 - 2/m)$,

so that $J^1 = [-\pi/2, -\Phi_0]$ and $J^m = [\Phi_0, \pi/2]$, and put $\Delta\Phi = 2/m$. The set $\{(\Phi, 0, \Phi, \Phi', \lambda(\Phi'), \lambda'(\Phi)) : -\Phi_0 \leq \Phi, \Phi' \leq \Phi_0, \Phi' - \Phi \geq \Delta\Phi\}$ is compact and belongs to the (open) domain of the function φ' . Choose a sufficiently large integer value $k = k(\varepsilon)$, so that for

$$(20) \quad \begin{aligned} |\sin \varphi - \sin \Phi| < 2/(km), \quad \xi = 0, \quad -\Phi_0 \leq \Phi, \Phi' \leq \Phi_0, \\ \Phi' - \Phi \geq \Delta\Phi, \quad \lambda = \lambda(\Phi'), \quad \lambda' = \lambda'(\Phi) \end{aligned}$$

it holds true that

$$(21) \quad -\frac{\cos \Phi'}{\cos \Phi} \frac{\partial \varphi'}{\partial \varphi} \in [(1 + \varepsilon)^{-1}, 1 + \varepsilon].$$

Formulas (21) and (11) mean that, under conditions (20), φ' is also close to Φ' . Increasing k if necessary, ensure (under the same conditions) that

$$(22) \quad \frac{\cos \Phi'}{\cos \Phi} \frac{\cos \varphi}{\cos \varphi'} \in [(1 + \varepsilon)^{-1}, 1 + \varepsilon].$$

Taking into account (15), (21), and (22), one obtains that under conditions (20) it holds true that

$$(23) \quad \left| \frac{\partial \xi'}{\partial \xi} \right| \in [(1 + \varepsilon)^{-2}, (1 + \varepsilon)^2].$$

3. Now we proceed to the description of the hollow $(\Omega_\varepsilon, I_\varepsilon)$.

(a) If $2 \leq i \neq \sigma(i) \leq m - 1$, then divide the interval J^i into k subintervals $J^{i,j}$ of equal measure λ , going in increasing order: $J^i = \cup_{j=1}^k J^{i,j}$, $\lambda(J^{i,j}) = 2/(km)$ for any $j = 1, \dots, k$. Recall that $d\lambda(\varphi) = \cos \varphi d\varphi$ and the value $k = k(\varepsilon)$ is defined above. Without loss of generality assume that $k(\varepsilon) \rightarrow \infty$ as $\varepsilon \rightarrow 0$.

To each pair of intervals, $J^{i,j}$ and $J^{\sigma(i),j}$, we apply the construction described above; see Figure 9. Namely, draw arcs of ellipses $\mathcal{E}_{i,j} = AB$, $\mathcal{E}'_{i,j} = A'B'$ and arcs of parabolas $\mathcal{P}_{i,j}$, $\mathcal{P}'_{i,j}$. Without loss of generality suppose that $i < \sigma(i)$. The angles AOB and $A'OB'$ correspond to the angular intervals $J^{i,j}$ and $J^{\sigma(i),j}$, respectively. The foci $\bar{F} = F_{i,j}$ and $\bar{F}' = F'_{i,j}$ belong to the intervals OA and OA' , respectively. The endpoints of the arcs $\mathcal{P}_{i,j}$ and $\mathcal{P}'_{i,j}$ also belong to the intervals OA and OA' , respectively. The angle corresponding to the ray OA (and therefore to the left endpoint of the interval $J^{i,j}$) will be denoted by $\bar{\Phi} = \Phi_{i,j}$, and the angle corresponding to the ray OA' (and therefore to the right endpoint of the interval $J^{\sigma(i),j}$) will be denoted by $\bar{\Phi}' = \Phi'_{i,j}$. Denote $\bar{\lambda} = \lambda_{i,j} := \lambda(\bar{\Phi})$ and $\bar{\lambda}' = \lambda'_{i,j} := \lambda'(\bar{\Phi}')$, according to formula (17). Next, select a value $\delta = \delta_{i,j}$ and draw two curves (*lateral reflectors*) in such a way that (i) each of the curves contains an arc of parabola (the first curve contains $\mathcal{P}_{i,j}$ and the second one $\mathcal{P}'_{i,j}$), an arc of circumference centered at O , and three radial segments; (ii) these curves do not intersect the intervals whose endpoints belong to the set $\{F_{\alpha,\beta}, F'_{\gamma,\delta} : (\alpha, \beta) \neq (i, j), (\gamma, \delta) \neq (\sigma(i), j)\}$: this will guarantee free passage of particles from one parabola to another; and (iii) the λ -measure of the angular interval occupied by each lateral reflector does not exceed $\varepsilon/(km)$. In Figure 9, the angular reflectors are the curves joining the points A and C , and the points A' and C' .

Notice that $-\Phi_0 \leq \bar{\Phi}, \bar{\Phi}' \leq \Phi_0$, and $\bar{\Phi}' - \bar{\Phi} \geq \Delta\Phi$. Indeed, $\bar{\Phi}$ and $\bar{\Phi}'$ do not belong to the intervals $J^1 = [-\pi/2, -\Phi_0]$ and $J^m = [\Phi_0, \pi/2]$. On the other hand, one has $\bar{\Phi}' - \bar{\Phi} \geq \sin \bar{\Phi}' - \sin \bar{\Phi} = \lambda([\bar{\Phi}, \bar{\Phi}']) \geq 2/m = \Delta\Phi$.

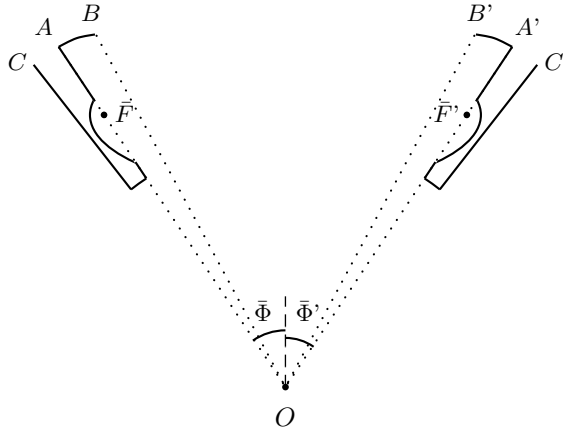


FIG. 9. Part of the hollow corresponding to the angular intervals $J^{i,j}$ and $J^{\sigma(i),j}$.

Introduce the shorthand notation $\varphi'(\varphi, \xi) = \varphi'(\varphi, \xi, \bar{\Phi}_{i,j}, \bar{\Phi}'_{i,j}, \lambda_{i,j}, \lambda'_{i,j}, \delta_{i,j})$. According to (21) and (23), for $\varphi \in J^{i,j}$ it holds true that

$$(24) \quad -\frac{\cos \bar{\Phi}'}{\cos \bar{\Phi}} \frac{\partial \varphi'}{\partial \varphi}(\varphi, 0) \in [(1 + \varepsilon)^{-1}, 1 + \varepsilon]$$

and

$$(25) \quad \left| \frac{\partial \xi'}{\partial \xi}(\varphi, 0) \right| \in [(1 + \varepsilon)^{-2}, (1 + \varepsilon)^2].$$

According to (11), one has $\varphi'(\bar{\Phi}, 0) = \bar{\Phi}'$; this equality and formula (24) imply that for $\varphi \in J^{i,j}$ and $\varphi' = \varphi'(\varphi, 0)$ one has

$$(26) \quad -\frac{\cos \bar{\Phi}'}{\cos \bar{\Phi}} \frac{\varphi' - \bar{\Phi}'}{\varphi - \bar{\Phi}} \in [(1 + \varepsilon)^{-1}, 1 + \varepsilon].$$

On the other hand, one has

$$(27) \quad \cos \bar{\Phi} |J^{i,j}| = \frac{2}{km} (1 + o(1)),$$

$$(28) \quad \cos \bar{\Phi}' |J^{\sigma(i),j}| = \frac{2}{km} (1 + o(1)),$$

with $o(1)$ being uniformly small over all i, j as $\varepsilon \rightarrow 0$, and $|J|$ being the Lebesgue measure of J . (Recall that the parameters $\bar{\Phi}, \bar{\Phi}', k$ and the intervals $J^{i,j}$ implicitly depend on ε .)

Choose closed intervals $\tilde{J}^{i,j} \subset J^{i,j}$ and $\tilde{J}^{\sigma(i),j} \subset J^{\sigma(i),j}$ satisfying the following conditions: (i) $\varphi'(\tilde{J}^{i,j} \times \{0\}) = \tilde{J}^{\sigma(i),j}$; (ii) some neighborhoods of $\tilde{J}^{i,j}$ and $\tilde{J}^{\sigma(i),j}$ belong to $J^{i,j}$ and $J^{\sigma(i),j}$, respectively; and (iii) the pseudobilliard trajectory with the initial data $(\varphi, 0)$, $\varphi \in \tilde{J}^{i,j}$ does not intersect the neighbor lateral reflectors (that is, the lateral reflectors corresponding to the intervals $\tilde{J}^{i,j+1}$ and $\tilde{J}^{\sigma(i),j-1}$ if $j \neq 1, k$; if $j = 1$, then $\tilde{J}^{\sigma(i),j-1}$ should be replaced with $\tilde{J}^{\sigma(i)-1,k}$, and if $j = k$, then $\tilde{J}^{i,j+1}$ should be replaced with $\tilde{J}^{i+1,1}$). Note in this regard that the neighbor lateral reflectors occupy a small part of the angular intervals $J^{i,j}$ and $J^{\sigma(i),j}$ (represented in the figure

by the arcs AB and $B'A'$). Other lateral reflectors will not be intersected by the choice of lateral reflectors.

By virtue of (26), (27), (28) and because of smallness of the angular intervals occupied by the lateral reflectors, $\tilde{J}^{i,j}$ and $\tilde{J}^{\sigma(i),j}$ may be chosen in such a way that the ratios $\lambda(\tilde{J}^{i,j})/\lambda(J^{i,j})$ and $\lambda(\tilde{J}^{\sigma(i),j})/\lambda(J^{\sigma(i),j})$ uniformly (with respect to i, j) tend to 1 as $\varepsilon \rightarrow 0$. Thus, a billiard particle going from O in a direction $\varphi \in \tilde{J}^{i,j}$ makes the same reflections and in the same order as under the pseudobilliard dynamics: first, reflection from $\mathcal{E}_{i,j}$, from $\mathcal{P}_{i,j}$, from $\mathcal{P}'_{i,j}$, and from $\mathcal{E}'_{i,j}$; finally, the particle goes back to O in the direction $\varphi'(\varphi, 0) \in \tilde{J}^{\sigma(i),j}$.

Choose $a_{i,j}$ in such a way that the following conditions are fulfilled: if $(\xi, \varphi) \in [-a_{i,j}, a_{i,j}] \times \tilde{J}^{i,j}$, then (i) the corresponding billiard trajectory does not intersect the lateral reflectors and the indicated order of reflections is preserved; (ii) $\varphi'(\varphi, \xi) \in J^{\sigma(i),j}$; (iii) $|\frac{\partial \xi'}{\partial \xi}(\varphi, \xi)| \in [(1 + \varepsilon)^{-3}, (1 + \varepsilon)^3]$. Analogously, choose $a_{\sigma(i),j}$ in such a way that the conditions are fulfilled: if $(\xi, \varphi) \in [-a_{\sigma(i),j}, a_{\sigma(i),j}] \times \tilde{J}^{\sigma(i),j}$, then (i) the billiard trajectory does not intersect the lateral reflectors and the order of its reflections is reversed; (ii) $\varphi'(\varphi, \xi) \in J^{i,j}$; (iii) $|\frac{\partial \xi'}{\partial \xi}(\varphi, \xi)| \in [(1 + \varepsilon)^{-3}, (1 + \varepsilon)^3]$. Note that the values $a_{i,j}$ and $a_{\sigma(i),j}$ implicitly depend on ε .

Select $a_\varepsilon \leq \min_{i,j} a_{ij}$ in such a way that $a_\varepsilon \rightarrow 0$ as $\varepsilon \rightarrow 0$, and denote $I_\varepsilon = (-a_\varepsilon, a_\varepsilon) \times \{0\}$, $\tilde{I}_\varepsilon = (-a_\varepsilon(1 + \varepsilon)^{-3}, a_\varepsilon(1 + \varepsilon)^{-3}) \times \{0\}$, and $\tilde{J}_\varepsilon^i = \tilde{J}^i := \cup_j \tilde{J}_\varepsilon^{i,j}$. The part of the boundary of Ω_ε related to the angular intervals $J^{i,j}$ and $J^{\sigma(i),j}$ under consideration is formed by the arcs of ellipses $\mathcal{E}_{i,j}$, $\mathcal{E}'_{i,j}$ and the corresponding lateral reflectors. Then a billiard particle with initial conditions $(\xi, \varphi) \in \tilde{I}_\varepsilon \times \tilde{J}^{i,j}$ after making four reflections will intersect l at a point $(\xi', 0) \in I_\varepsilon$, and the angle at the point of intersection will be $\varphi_{\Omega_\varepsilon, I_\varepsilon}^+(\xi, \varphi) = \varphi'(\varphi, \xi) \in J^{\sigma(i),j} \subset J^{\sigma(i)}$. Thus, the set of values $(\xi, \varphi) \in I_\varepsilon \times J^i$ such that $\varphi_{\Omega_\varepsilon, I_\varepsilon}^+(\xi, \varphi) \notin J^{\sigma(i)}$ is contained in the set $(I_\varepsilon \times J^i) \setminus (\tilde{I}_\varepsilon \times \tilde{J}_\varepsilon^i)$, whose measure is vanishing as $\varepsilon \rightarrow 0$.

(b) If $2 \leq i = \sigma(i) \leq m - 1$, then the corresponding part of the boundary is the arc of circumference of radius 2 with the center at O occupying the angular interval J^i , that is, the set $\{2(\sin \varphi, \cos \varphi), \varphi \in J^i\}$. Next, we will show that for all values $(\xi, 0) \in I_\varepsilon$, $\varphi \in J^i$, except for a portion of order $o(1)$, the corresponding billiard particle makes one reflection from the arc and then goes back to I_ε in the direction $\varphi' \in J^i$.

For all values $\varphi \in J^i$, except for the union of two intervals of vanishing length (each of the intervals is contained in J^i , has the length $2 \arctan(a_\varepsilon/4)$, and contains an endpoint of J^i), the particle starting at $(\xi, 0) \in I_\varepsilon$ in the direction φ will reflect from the indicated arc of circumference. Let $\psi \in J^i$ be the angular coordinate of the reflection point. By $(\xi', 0)$ denote the point at which the reflected particle intersects the straight line l . One easily verifies that

$$(29) \quad \frac{1}{\xi} + \frac{1}{\xi'} = \cos \psi.$$

One has

$$(30) \quad |\xi| < a_\varepsilon,$$

and hence

$$(31) \quad \frac{1}{|\xi'|} = \left| \cos \psi - \frac{1}{\xi} \right| > \frac{1}{a_\varepsilon} - 1.$$

From (29) it follows that $|\xi + \xi'|/|\xi\xi'| = |\cos \psi| \leq 1$, and, taking into account (30) and (31), one finds that $|\xi + \xi'| < a_\varepsilon^2/(1 - a_\varepsilon)$. This implies that for all values $(\xi, 0) \in I_\varepsilon$, except for a set of measure $O(a_\varepsilon^2)$, the second point of intersection of the billiard trajectory belongs to I_ε ; moreover, the velocity at this point, $\varphi_{\Omega_\varepsilon, I_\varepsilon}^+(\xi, \varphi)$, belongs to $\mathcal{N}_{2 \arctan(a_\varepsilon/4)}(J^i)$, the neighborhood of J^i of radius $2 \arctan(a_\varepsilon/4)$. This finally implies that for all $(\xi, \varphi) \in I_\varepsilon \times J^i$, except for a portion of order $O(a_\varepsilon)$, it holds that $\varphi_{\Omega_\varepsilon, I_\varepsilon}^+(\xi, \varphi) \in J^i$.

(c) The parts of the hollow's boundary, corresponding to J^1 and J^m , are formed by smooth curves joining the corresponding endpoints of I_ε and the points $2(\sin \Phi_0, -\cos \Phi_0)$ and $2(\sin \Phi_0, \cos \Phi_0)$, respectively. The unique condition on these curves is that they can be parametrized by the monotonically increasing angular coordinate. For those values $\sigma(1), \sigma(m)$ that coincide with neither 1 nor m take just the arcs of circumference of radius 2 corresponding to the angular intervals $J^{\sigma(1)}, J^{\sigma(m)}$.

Consider the union of all of the elliptic arcs $\mathcal{E}_{i,j}, \mathcal{E}_{i,j}^i$ introduced in item (a), all of the arcs of circumference defined in items (a) and (b), and the two curves introduced in this item (c). Let us call this union the *main element*. Each lateral reflector is a curve; select it in such a way that both its endpoints belong to the main element. Finally, the curve $\partial\Omega_\varepsilon \setminus I_\varepsilon$ is the union of all of the lateral reflectors and the part of the main element visible from O (that is, which is not shielded by the adjacent lateral reflectors). Thus, the definition of the hollow $(\Omega_\varepsilon, I_\varepsilon)$ is complete.

In Figure 10, there is shown a particular hollow $(\Omega_\varepsilon, I_\varepsilon)$ corresponding to the permutation $\sigma = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 5 & 4 & 3 & 2 & 1 \end{pmatrix}$. The angular intervals J^1, \dots, J^5 are separated by dotted lines. The family of hollows $(\Omega_\varepsilon, I_\varepsilon)$, with vanishingly small ε , has the following property: for almost all particles with the initial direction from J^2 (resp. J^3, J^4), the final direction will belong to J^4 (resp. J^3, J^2). In the figure, there is shown the trajectory of a particle with the initial direction $\varphi \in J^2$ and the final direction $\varphi^+ \in J^4$. The particle makes a reflection from an elliptic arc, then two reflections from (very small) parabolic arcs, and, finally, again from an elliptic arc. According to our notation, these arcs are $\mathcal{E}_{2,2}, \mathcal{P}_{2,2}, \mathcal{P}_{2,2}$, and $\mathcal{E}_{2,2}$.

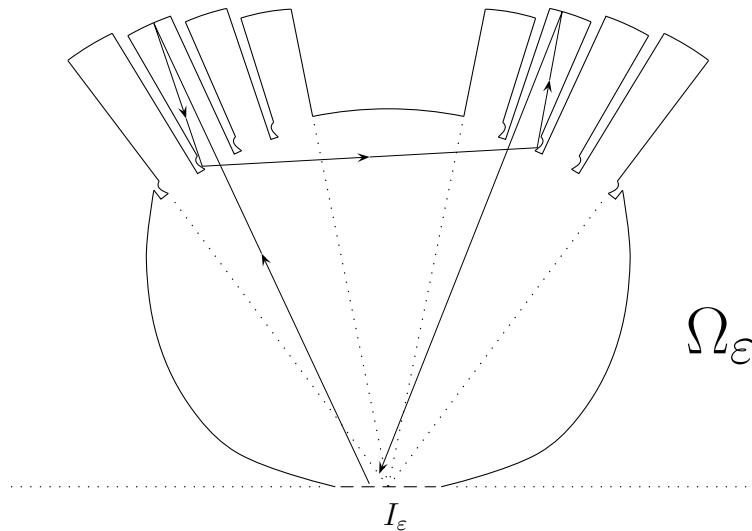


FIG. 10. A hollow $(\Omega_\varepsilon, I_\varepsilon)$ approximating a σ -measure.

4. Concluding remarks and applications. Physical bodies in the real world have atomic structure and therefore are disconnected. This is a reason for using (generally) disconnected sets Q_m in the definition of a rough body. In the future we intend to turn to propose and study the notion of a three-dimensional rough body, where the connectivity assumption is absolutely useless; this is another reason. By removing this assumption, the consideration in two dimensions (namely, the proof of Lemma 1) is made somewhat more difficult, but at the same time prerequisites for passing to the three-dimensional case are created.

In fact, the notions of “disconnected” (as everywhere in this paper) and “connected” rough bodies are equivalent. There is a natural one-to-one correspondence between the equivalence classes in the connected and disconnected cases,³ the former classes being subclasses of the latter ones under this correspondence.

Let us now consider applications of the main theorem to problems of the body of minimal or maximal aerodynamic resistance. A two-dimensional convex body B moves, at constant velocity, through a rarefied homogeneous medium in \mathbb{R}^2 , and at the same time it slowly rotates. The rotation is generally nonuniform; we assume that, during a sufficiently long observation period, in a reference system connected with the body the body’s velocity is distributed in S^1 according to a given density function ρ , with $\int_{S^1} \rho(v) dv = 1$. The medium particles do not mutually interact, and collisions of the particles with the body are absolutely elastic. The resistance of the medium to the motion of the body is a vector-valued function of time. After averaging it over a sufficiently long period of time, one gets a vector. We are interested in the projection of this vector onto the direction of motion; for the sake of brevity, it will be called mean resistance, or just resistance. The problem is as follows: given B , determine the roughness on it in such a way that the main resistance of the resulting rough body is minimal or maximal.

A prototype of such a mechanical system is an artificial satellite of the Earth on relatively low altitudes (100 ÷ 200 km), with restricted capacity of rotation angle control. The satellite’s motion is slowing down by the rest of the atmosphere; the problem is to minimize or maximize the effect of slowing down. The problems of resistance *maximization* may also arise when considering solar sail: a spacecraft driven by the pressure of solar photons.

The initial velocity of an incident particle (in the reference system connected with the body) is $-v$, and the final velocity is v^+ ; therefore, the momentum transmitted by the particle to the body is $v + v^+$. The projection of the transmitted momentum onto the direction of motion of the body equals $1 + \langle v, v^+ \rangle$. Averaging this value over all particles incident on the body within a sufficiently long time interval, one gets the mean resistance. The averaging amounts to integration over $\rho(v) d\nu_B(v, v^+, n)$; that is, the mean resistance of the rough body equals

$$R(\nu_B) = \iiint_{\mathbb{T}^3} (1 + \langle v, v^+ \rangle) \rho(v) d\nu_B(v, v^+, n).$$

Using the main theorem and Fubini’s theorem, one rewrites this formula in the form

$$(32) \quad R(\nu_B) = \int_{S^1} d\tau_B(n) \iint_{\mathbb{T}^2} (1 + \langle v, v^+ \rangle) \rho(v) d\eta_{B,n}(v, v^+),$$

³More precisely, we mean equivalence classes formed by sequences of connected/disconnected sets.

where $\eta_{B,n} \in \Lambda_n$. Thus, the minimization problem for $R(\nu_B)$ reduces to minimization, for any n , of the functional $\iint_{\mathbb{T}^2} (1 + \langle v, v^+ \rangle) \rho(v) d\eta(v, v^+)$ over all $\eta \in \Lambda_n$. Using the notation introduced in section 3.3, one comes to the problem:

$$(33) \quad \inf_{\eta \in \Lambda} \iint_{\square} (1 + \cos(\varphi - \varphi^+)) \varrho(\varphi) d\eta(\varphi, \varphi^+),$$

where $\varrho(\varphi) = \rho(v)$ for $\varphi = \text{Arg } v - \text{Arg } n$. This problem, in turn, by symmetrization of the cost function reduces to a particular Monge–Kantorovich problem:

$$(34) \quad \inf_{\eta \in \Lambda_{\lambda, \lambda}} \mathcal{F}(\eta), \quad \text{where} \quad \mathcal{F}(\eta) = \iint_{\square} c(\varphi, \varphi^+) d\eta(\varphi, \varphi^+),$$

where $c(\varphi, \varphi^+) = (1 + \cos(\varphi - \varphi^+)) \frac{\varrho(\varphi) + \varrho(\varphi^+)}{2}$ and $\Lambda_{\lambda, \lambda}$ is the set of measures η on \square having both marginal measures equal to λ : $\pi_{\varphi}^{\#} \eta = \lambda = \pi_{\varphi^+}^{\#} \eta$. Recall that λ is defined by $d\lambda(v) = \cos \varphi d\varphi$.

Problem (34) can be exactly solved in several particular cases. Consider the case of uniform motion, where the function ρ , and therefore ϱ , is constant, and thus one can take $c(\varphi, \varphi^+) = \frac{3}{8} (1 + \cos(\varphi - \varphi^+))$.⁴ Note that $\mathcal{F}(\eta_0) = 1$ and therefore resistance of the smooth body is equal to its perimeter: $R(\nu_B) = \int_{S^1} d\tau_B(n) \mathcal{F}(\eta_0) = |\partial B|$. (Recall that the measure η_0 belongs to Λ and is supported on the diagonal $\varphi^+ = -\varphi$.) The minimization problem (34) for constant ϱ was solved in [14]: one has $\inf_B R(\nu_B) = 0.9878 \dots |\partial B|$, the infimum being taken over all roughenings of B .

Note that the corresponding *maximization* problem for (34) has the trivial solution, which does not depend on the function ϱ : $\eta = \eta_{\star}$, the measure $\eta_{\star} \in \Lambda$ being supported on the diagonal $\varphi^+ = \varphi$. One has $\sup_B R(\nu_B) = \kappa |\partial B|$, where $\kappa = (\int_{-\pi/2}^{\pi/2} \varrho(\varphi) \cos \varphi d\varphi) / (\int_{-\pi/2}^{\pi/2} \varrho(\varphi) \cos^3 \varphi d\varphi) > 1$; in the case of uniform rotation one has $\kappa = 1.5$. The maximization problem was studied in more detail in [20].

Appendix A. The construction below is simple (see Figure 11), but its description is a bit cumbersome.

Take a point in the interior of B and connect it by segments with all vertices. The polygon is thus divided into several triangles; fix i and m and consider the triangle with the base \mathfrak{b}_i , the i th side of B . Denote by $d(\Omega_i^m)$ the diameter of the orthogonal projection of Ω_i^m onto the straight line containing I_i^m ; one obviously has $d(\Omega_i^m) \geq |I_i^m|$. Fix a positive number $\kappa < |I_i^m|/d(\Omega_i^m)$.

Take a rectangle Π^1 contained in the triangle and such that one side of Π^1 belongs to \mathfrak{b}_i . By δ_1 denote the total length of the part of \mathfrak{b}_i which is not occupied by this side.

For the sake of brevity, the image of a set under the composition of a homothety with positive ratio and a translation will be called a copy of this set. Take several copies of Ω_i^m (copies of first order) that do not mutually interact, that belong to Π^1 , whose corresponding copies of I_i^m belong to \mathfrak{b}_i , and whose portion of the side of Π^1 occupied by them is more than κ .

Next, take several rectangles that do not mutually intersect and do not intersect with the chosen copies of Ω_i^m , belong to Π^1 , and have one side contained in \mathfrak{b}_i . Denote by Π^2 the union of these rectangles and by δ_2 the total length of the part of the side of Π^1 which is not occupied by the rectangles from Π^2 and by the copies of I_i^m . Next,

⁴The normalization constant 3/8 is taken for further convenience.

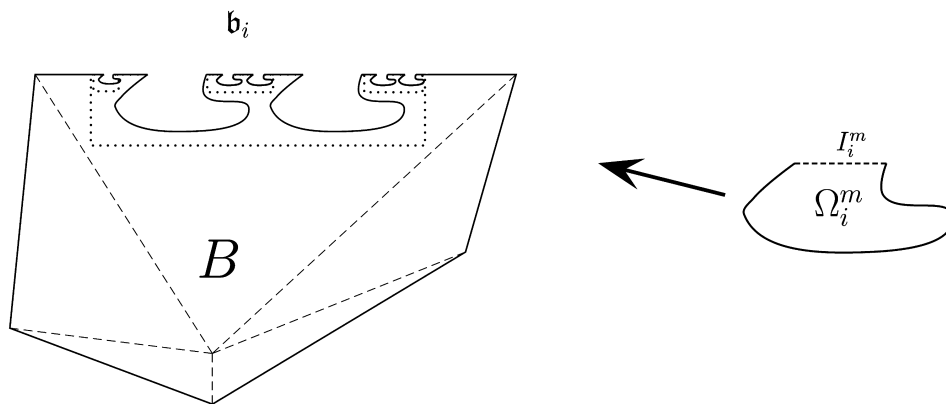


FIG. 11. Making hollows on the i th side of a polygon.

for each rectangle from Π^2 choose several copies of Ω_i^m (copies of second order) in the way completely similar to the described above (see Figure 11).

Continuing this process, one obtains a sequence Π^1, Π^2, \dots of unions of rectangles and collections of copies of Ω_i^m of first, second, \dots order. Choose the rectangles in such a way that $\delta_1 + \delta_2 + \dots < 1/m$ and $\text{Area}(\Pi^1) < 1/m$. Finally, choose k such that the total length of sides of rectangles from Π^{k+1} contained in b_i is less than $1/m$, and take the collection of copies of Ω_i^m of order $1, 2, \dots, k$ (we shall call it full collection). The total length of the part of b_i not occupied by the corresponding copies of I_i^m is less than $2/m$, and therefore it goes to zero as $m \rightarrow \infty$.

By definition, the desired set Q_m is B minus the union of full collections of copies of Ω_i^m over all i .

Appendix B. We prove here slightly more than needed.

STATEMENT 3. Let $A = (a_{ij})_{i,j=1}^k$ be a symmetric matrix, with a_{ij} being non-negative integers. Denote $n_i = \sum_{j=1}^k a_{ij}$. Then there exist matrices $B_{ij} = (b_{ij}^{\mu\nu})_{\mu,\nu}$ of size $n_i \times n_j$ such that $b_{ij}^{\mu\nu} \in \{0, 1\}$, $B_{ij}^T = B_{ji}$, the sum of elements in B_{ij} equals a_{ij} , and the block matrix $D = (B_{ij})$ contains exactly one unit in each row and each column.

Note that for some values $i = i_1, i_2, \dots$ it may happen that $n_i = 0$, that is, $a_{ij} = 0$ for all $j = 1, \dots, k$. Then the corresponding matrices B_{ij} have the size $0 \times n_j$, that is, are empty. In this case D coincides with the block matrix $D' = (B_{ij})$ having the rows i_1, i_2, \dots and columns i_1, i_2, \dots crossed out.

Proof. The proof is by induction on k . Let the statement be true for $k - 1$; prove it for k . Take the matrix $\tilde{A} = (a_{ij})_{i,j=2}^k$; there exists a block matrix $\tilde{B} = (\tilde{B}_{ij})_{i,j=2}^k$ satisfying the statement. Note that the order of \tilde{B}_{ij} is $\tilde{n}_i \times \tilde{n}_j$, where $\tilde{n}_i = \sum_{j=2}^k a_{ij} = n_i - a_{i1}$. Define the matrices B_{ij} as follows.

(a) Put $B_{11} = \text{diag}\{\underbrace{1, \dots, 1}_{a_{11}}, 0, \dots, 0\}$.

(b) Put $b_{12}^{a_{11}+1,1} = \dots = b_{12}^{a_{11}+a_{12},a_{12}} = 1$; $b_{13}^{a_{11}+a_{12}+1,1} = \dots = b_{13}^{a_{11}+a_{12}+a_{13},a_{13}} = 1$; \dots ; $b_{1k}^{a_{11}+\dots+a_{1,k-1}+1,1} = \dots = b_{1k}^{a_{11}+\dots+a_{1k},a_{1k}} = 1$; the other elements of the matrices B_{1j} , $j = 2, \dots, k$, are zeros. Thus, on the diagonal of B_{1j} starting from the element at the first column and the $(a_{11} + a_{12} + \dots + a_{1,j-1} + 1)$ th row, the first a_{1j} elements equal 1, and the remaining elements on this diagonal and all of the elements

off the diagonal are zeros. This defines the matrices B_{1j} , $j = 2, \dots, k$. The matrices B_{i1} , $i = 2, \dots, k$, are determined by the condition $B_{i1} = B_{1i}^T$.

(c) For $i \geq 2$, $j \geq 2$ define the matrix B_{ij} as follows. For $\mu \leq a_{1i}$ or $\nu \leq a_{1j}$, put $b_{ij}^{\mu\nu} = 0$, and for $\mu \geq a_{1i} + 1$, $\nu \geq a_{1j} + 1$, put $b_{ij}^{\mu\nu} = \tilde{b}_{ij}^{\mu-a_{1i}, \nu-a_{1j}}$. Thus, in the obtained matrix B_{ij} , the right lower corner coincides with the matrix \tilde{B}_{ij} , and all of the remaining elements are equal to zero. The number of rows of this matrix equals $a_{1i} + \tilde{n}_i = n_i$, and the number of columns equals $a_{1j} + \tilde{n}_j = n_j$. One obviously has $B_{ij}^T = B_{ji}$.

One easily verifies that $\sum_{\mu\nu} b_{ij}^{\mu\nu} = a_{ij}$ and that each row and each column of the obtained block matrix $D = (B_{ij})_{i,j=1}^k$ contains precisely one unit. \square

REFERENCES

- [1] I. NEWTON, *Philosophiae Naturalis Principia Mathematica*, William Dawson & Sons, London, 1686.
- [2] G. BUTTAZZO AND B. KAWOHL, *On Newton's problem of minimal resistance*, Math. Intelligencer, 15 (1993), pp. 7–12.
- [3] G. BUTTAZZO, V. FERONE, AND B. KAWOHL, *Minimum problems over sets of concave functions and related questions*, Math. Nachr., 173 (1995), pp. 71–89.
- [4] F. BROCK, V. FERONE, AND B. KAWOHL, *A symmetry problem in the calculus of variations*, Calc. Var. Partial Differential Equations, 4 (1996), pp. 593–599.
- [5] G. BUTTAZZO AND P. GUASONI, *Shape optimization problems over classes of convex domains*, J. Convex Anal., 4 (1997), pp. 343–351.
- [6] M. BELLONI AND B. KAWOHL, *A paper of Legendre revisited*, Forum Math., 9 (1997), pp. 655–668.
- [7] T. LACHAND-ROBERT AND M. A. PELETIER, *Newton's problem of the body of minimal resistance in the class of convex developable functions*, Math. Nachr., 226 (2001), pp. 153–176.
- [8] M. COMTE AND T. LACHAND-ROBERT, *Newton's problem of the body of minimal resistance under a single-impact assumption*, Calc. Var. Partial Differential Equations, 12 (2001), pp. 173–211.
- [9] M. COMTE AND T. LACHAND-ROBERT, *Existence of minimizers for Newton's problem of the body of minimal resistance under a single-impact assumption*, J. Anal. Math., 83 (2001), pp. 313–335.
- [10] D. HORSTMANN, B. KAWOHL, AND P. VILLAGGIO, *Newton's aerodynamic problem in the presence of friction*, NoDEA Nonlinear Differential Equations Appl., 9 (2002), pp. 295–307.
- [11] M. BELLONI AND A. WAGNER, *Newton's problem of minimal resistance in the class of bodies with prescribed volume*, J. Convex Anal., 10 (2003), pp. 491–500.
- [12] A. YU. PLAKHOV, *Newton's problem of a body of minimal aerodynamic resistance*, Dokl. Akad. Nauk, 390 (2003), pp. 314–317.
- [13] A. YU. PLAKHOV, *Newton's problem of the body of minimal resistance with a bounded number of collisions*, Russian Math. Surveys, 58 (2003), pp. 191–192.
- [14] A. YU. PLAKHOV, *Newton's problem of the body of minimum mean resistance*, Sb. Math., 195 (2004), pp. 1017–1037.
- [15] A. YU. PLAKHOV AND D. F. M. TORRES, *Newton's aerodynamic problem in media of chaotically moving particles*, Sb. Math., 196 (2005), pp. 885–933.
- [16] R. G. BARANTSEV, *Interaction of Rarefied Gases with Streamline Surfaces*, Mir, Moscow, 1975 (in Russian).
- [17] D. BLACKMORE AND J. G. ZHOU, *A general fractal distribution function for rough surface profiles*, SIAM J. Appl. Math., 56 (1996), pp. 1694–1719.
- [18] O. A. AKSENOVA AND I. A. KHALIDOV, *Fractal and statistical models of rough surface interacting with rarefied gas flow*, in AIP Conference Proceedings, RAREFIED GAS DYNAMICS: 24th International Symposium on Rarefied Gas Dynamics, Vol. 762, American Institute of Physics, College Park, MD, 2005, pp. 993–998.
- [19] N. CHERNOV, *Entropy, Lyapunov exponents, and mean free path for billiards*, J. Statist. Phys., 88 (1997), pp. 1–29.
- [20] A. PLAKHOV AND P. GOUVEIA, *Problems of maximal mean resistance on the plane*, Nonlinearity, 20 (2007), pp. 2271–2287.

PERMANENCE AND ASYMPTOTICALLY STABLE COMPLETE TRAJECTORIES FOR NONAUTONOMOUS LOTKA–VOLTERRA MODELS WITH DIFFUSION*

JOSÉ A. LANGA[†], JAMES C. ROBINSON[‡], ANÍBAL RODRÍGUEZ-BERNAL[§], AND
ANTONIO SUÁREZ[¶]

Abstract. Lotka–Volterra systems are the canonical ecological models used to analyze population dynamics of competition, symbiosis, or prey–predator behavior involving different interacting species in a fixed habitat. Much of the work on these models has been within the framework of infinite-dimensional dynamical systems, but this has frequently been extended to allow explicit time dependence, generally in a periodic, quasiperiodic, or almost periodic fashion. The presence of more general nonautonomous terms in the equations leads to nontrivial difficulties which have stalled the development of the theory in this direction. However, the theory of nonautonomous dynamical systems has received much attention in the last decade, and this has opened new possibilities in the analysis of classical models with general nonautonomous terms. In this paper we use the recent theory of attractors for nonautonomous PDEs to obtain new results on the permanence and the existence of forwards and pullback asymptotically stable global solutions associated to nonautonomous Lotka–Volterra systems describing competition, symbiosis, or prey–predator phenomena. We note in particular that our results are valid for prey–predator models, which are not order-preserving: even in the “simple” autonomous case the uniqueness and global attractivity of the positive equilibrium (which follows from the more general results here) is new.

Key words. Lotka–Volterra competition, symbiosis and prey–predator systems, nonautonomous dynamical systems, permanence, attracting complete trajectories

AMS subject classifications. 35B40, 35K55, 92D25, 37L05

DOI. 10.1137/080721790

1. Introduction. Partial differential equations have proved a very useful tool in the modelling of many ecological phenomena related to the dynamics between species interacting in a given habitat. Many authors have allowed explicit dependence on both space and time in the parameters of the equation, a natural way to take into account the spatial and temporal variations that influence real species interactions.

In this paper we consider a nonautonomous model for two species (u and v), evolving within a habitat Ω that is a bounded domain in \mathbf{R}^N , $N \geq 1$, with a smooth

*Received by the editors April 21, 2008; accepted for publication September 13, 2008; published electronically January 28, 2009.

<http://www.siam.org/journals/sima/40-6/72179.html>

[†]Dpto. Ecuaciones Diferenciales y Análisis Numérico. C/ Tarfia s/n, 41012. Sevilla. Spain (langa@us.es). Partly supported by Ministerio de Educación y Ciencia (Spain) under grants MTM2005-01412, Consejería de Innovación, Ciencia y Empresa (Junta de Andalucía, Spain) under the Proyecto de Excelencia FQM-02468.

[‡]Mathematics Institute, University of Warwick, Coventry CV4 7AL, UK (j.c.robinson@warwick.ac.uk). Partly supported by a Royal Society University Research Fellowship.

[§]Dpto. de Matemática Aplicada, Universidad Complutense de Madrid, Madrid 28040, Spain (arober@mat.ucm.es). Partly supported by Project MTM2006-08262, DGES, Spain and Grupo de Investigación UCM-CAM 920894, CADEDIF.

[¶]Dpto. Ecuaciones Diferenciales y Análisis Numérico. C/ Tarfia s/n, 41012, Sevilla, Spain (suarez@us.es). Partly supported by Ministerio de Educación y Ciencia (Spain) under grants MTM2006-07932, Consejería de Innovación, Ciencia y Empresa (Junta de Andalucía, Spain) under the Proyecto de Excelencia FQM-520.

boundary $\partial\Omega$, of the following type:

$$(1.1) \quad \begin{cases} u_t - d_1 \Delta u = uf(t, x, u, v) & x \in \Omega, t > s \\ v_t - d_2 \Delta v = vg(t, x, u, v) & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 & x \in \partial\Omega, t > s \\ u(s) = u_s, v(s) = v_s, \end{cases}$$

where f and g are regular functions, d_1, d_2 are positive constants, and \mathcal{B}_i denotes one of the boundary operators

$$(1.2) \quad \mathcal{B}u = u, \quad \text{or} \quad \mathcal{B}u = \frac{\partial u}{\partial \vec{n}}, \quad \text{or} \quad \mathcal{B}u = d \frac{\partial u}{\partial \vec{n}} + \sigma(x)u,$$

for the Dirichlet, Neumann, or Robin case, respectively, \vec{n} is the outward normal vector-field to $\partial\Omega$, and $\sigma(x)$ a C^1 function. Note that we take diffusion coefficient d_i and boundary potential $\sigma_i(x)$ for the case of Robin boundary condition \mathcal{B}_i . Also note that we allow all of the nine possible combinations of boundary conditions in (1.1).

A particularly interesting class of models of the form (1.1) are the nonautonomous Lotka–Volterra models:

$$(1.3) \quad \begin{cases} u_t - d_1 \Delta u = u(\lambda(t, x) - a(t, x)u - b(t, x)v) & x \in \Omega, t > s \\ v_t - d_2 \Delta v = v(\mu(t, x) - c(t, x)u - d(t, x)v) & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 & x \in \partial\Omega, t > s \\ u(s) = u_s, v(s) = v_s. \end{cases}$$

We refer, for example, to [6] for the biological meaning of the parameters $d_1, d_2, \lambda, \mu, a, b, c, d$ involved in (1.3).

In line with the ecological interpretation of these models, we will consider only positive solutions, and in the light of this we note here that $u_s, v_s \geq 0$ implies that the solution of (1.1) satisfies $u, v \geq 0$.

Note that our hypotheses on b and c allow different models of population dynamics: competition if $b, c > 0$, symbiosis if $b, c < 0$, and prey-predator if $b > 0$ and $c < 0$, although we do not allow sign-changing coefficients.

Of course, it is an important problem to determine the asymptotic behavior of solutions of the system (1.1). Since in general this is a very complicated task, one may try to solve simpler problems; e.g., one can try to determine whether or not the two species will survive in the long term or if, on the contrary, one of them will be driven to extinction. Survival of the species has been formalized in the notion of *permanence*; see Hale and Waltman [15] or Hutson and Schmitt [20]. Loosely speaking, the system (1.1) is said to be *permanent* if for any positive initial data u_s and v_s , within a finite time the values of the solution $(u(t, s, x; u_s, v_s), v(t, s, x; u_s, v_s))$, for $x \in \Omega$, enter and remain within a compact set in \mathbf{R}^2 that is strictly bounded away from zero in each component. Note, however, that this is an imprecise statement in the presence of Dirichlet boundary conditions.

Note that permanence is a form of *coexistence* of the species, since none is extinguished at any part of the habitat domain at any time.

A related situation, which implies that the system is permanent but gives more detail since it also indicates the expected final state of the system, is when there exists a solution, bounded away from zero, to which all other solutions tend asymptotically.

These two are the main topics with which we are concerned in this paper.

Before going further observe that both (1.1) and (1.3) always possess the *trivial* solution $(0, 0)$ and *semitrivial* solutions of the form $(u, 0)$ and $(0, v)$. In the latter

case the nontrivial component satisfies a scalar parabolic problem, of logistic type in the case of (1.3). The dynamics of these solutions have a deep impact on the global dynamics of general solutions. Indeed, if the system is permanent, this implies that semitrivial solutions must be unstable in some sense. On the other hand, if semitrivial solutions are stable, then it can be expected that some solutions of the system exhibit *extinction*, that is, one of the species (or both) approaches asymptotically the value zero.

Some results are already known along these lines. For example, in the autonomous case, assume that all the coefficients in (1.3) are constants and consider, for example, the problem with Dirichlet boundary conditions. In this case results about permanence for problem (1.3) depend on the values of λ and μ with respect to the first eigenvalue of certain associated linear elliptic problems, which we now describe. Given $d \in \mathbf{R}$, $d > 0$, and $f \in L^\infty(\Omega)$, we denote by $\Lambda(d, f)$ (we write $\Lambda_0 := \Lambda(d, 0)$) the first eigenvalue of the problem

$$\begin{cases} -d\Delta w = \sigma w + f(x)w & \text{in } \Omega, \\ w = 0 & \text{on } \partial\Omega, \end{cases}$$

and given $\gamma, \alpha \in \mathbf{R}$ with $\alpha > 0$, we denote by $\omega_{[d, \gamma, \alpha]}$ the unique positive solution of

$$\begin{cases} -d\Delta w = \gamma w - \alpha w^2 & \text{in } \Omega, \\ w = 0 & \text{on } \partial\Omega. \end{cases}$$

If λ and μ satisfy

$$(1.4) \quad \lambda > \Lambda(d_1, -b\omega_{[d_2, \mu, d]}) \quad \text{and} \quad \mu > \Lambda(d_2, -c\omega_{[d_1, \lambda, a]}),$$

then the autonomous version of the competition or prey-predator cases of (1.3), with Dirichlet boundary conditions in both components, are permanent and moreover there exists a positive equilibrium solution (Cantrell et al. [4], [6], [7], [8] and López-Gómez [27]).

Although the case of symbiosis, $b, c < 0$, is not treated in these papers, a similar result holds provided that

$$bc < ad,$$

a condition which is used to obtain *a priori bounds* for the solutions (see, for instance, Pao [30] or Theorem 9.8 in Delgado et al. [12], where moreover the coefficients a, b, c , and d depend on x).

Note that (1.4) is a condition that expresses the instability of semitrivial solutions.

However, in the competition case it is well known that if $\lambda \leq \Lambda_0$ or $\mu \leq \Lambda_0$, then one of the two species (or both of them) will be driven to extinction (see López-Gómez and Sabina [29] for an improvement of this result). Similar results can be obtained in the other cases; see [6] and [30]. Note that, in contrast with (1.4), the condition above expresses the stability of either one of the semitrivial solutions.

When nonautonomous terms are allowed in the equations, this is usually done under the assumption of periodicity, quasiperiodicity, or almost periodicity, and in this case similar results can be obtained to those for autonomous equations (see Hess [17], Hess and Lazer [18], Hetzer and Shen [19], and references therein). For the case of periodic coefficients, the use of the Poincaré map implies that the system resembles an autonomous one in many respects.

Cantrell and Cosner [5] assume general nonautonomous terms that are bounded by periodic functions, and using a comparison method give conditions on λ and μ that guarantee that (1.3) is permanent.

Note that most of the references cited in the papers above are concerned (besides periodicity or almost periodicity) with some particular choice of boundary conditions (typically Neumann, or even Dirichlet, in both components) and one of the competition, symbiosis or prey-predator cases. In the first two cases a common tool in the references is the use of order-preserving properties of the Lotka–Volterra system.

For example, in the case of almost periodic time dependence, Hetzer and Shen [19] proved similar results for the competition case, assuming that $d_1 = d_2$ and $\lambda = \mu$ are constant and both components of the system satisfy Dirichlet boundary conditions (no such restrictions are required in the case of Neumann BCs). In that paper, the limitation to almost periodic cases is due to the use of skew-product techniques which require, some way or another, some sort of time recurrence in the coefficients of the system.

Note that in [25] Langa et al. studied permanence for the competition case with Dirichlet boundary conditions when only the coefficient a is allowed to depend on time.

In this paper we allow general nonautonomous terms and do not restrict ourselves to (for example) almost periodic time dependence. As said before, we also consider all nine possible choices of boundary conditions and treat competition, symbiosis, and prey-predator models, since we do not rely on monotonicity properties of the system. Note that the only restriction that we impose on the coefficients is that $d_1 = d_2$ in the symbiotic case, a condition that we assume only in order to have explicit upper bounds on the solutions, but not for the permanence results. Also note that as we employ for the solutions of (1.1) or (1.3) the approach of *nonautonomous processes* rather than skew-product techniques, we have to pay attention to both the initial time, s , and the observation time for the solutions, $t > s$. This implies that concepts like permanence, stability, instability, and attractivity can be defined and analyzed in both pullback and forwards senses; see section 2 for further details and also [24]. Observe also that while pullback properties (e.g., permanence, attraction) are usually the most one can expect for general nonautonomous terms, in this case we can also show results on permanence and attractivity forwards in time; see Langa et al. [25, 23] for cases of pullback but not forwards permanence or attraction in nonautonomous reaction-diffusion equations.

In section 3, using results for the scalar nonautonomous logistic equations from, e.g., [25, 34], which we compile in section 3.1, we make use of the theory of attractors for nonautonomous PDEs as developed by Chepyzhov and Vishik [9] (see also Crauel et al. [11] or Kloeden and Schmalfuss [21]). Thus, we prove in section 3.2 that under the assumptions

$$\inf_{\mathbf{R} \times \Omega} a(t, x) > 0 \quad \text{and} \quad \inf_{\mathbf{R} \times \Omega} d(t, x) > 0$$

the system (1.3) has a nonautonomous attractor; see Theorem 3.5. The existence of a nonautonomous attractor in this case implies the presence of bounded complete trajectories, i.e., solutions defined for all time.

From here we derive in section 3.3 some sufficient conditions for the extinction of one (or both) of the species of the system. These conditions are far from optimal but qualitatively describe the stability of semitrivial solutions; see Proposition 3.6.

Then, in section 3.4 we give sufficient conditions reflecting the instability of semitrivial solutions that guarantee that (1.3) is permanent both in a pullback and in

a forwards sense. We want to stress here that these sufficient conditions involve only information about the behavior of the coefficients of the system as either $t \rightarrow -\infty$ or $t \rightarrow \infty$. Also, they are given in such a way that the result is robust with respect to perturbations of the coefficients.

The rest of the paper is then devoted to a more detailed analysis of the asymptotic behavior of the solutions of (1.3). After some preparatory material in sections 4 and 5, we will prove in section 6 that under appropriate conditions on the parameters all nonsemitrivial solutions of (1.3) have the same asymptotic behavior as $t \rightarrow \infty$. In particular all bounded complete trajectories in the nonautonomous attractor have the same asymptotic behavior as $t \rightarrow \infty$. For this we make use of the permanence results in section 3.4 and impose a smallness condition on the product of the coupling parameters:

$$\limsup_{t \rightarrow \infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow \infty} \|c\|_{L^\infty(\Omega)} < \rho_0$$

for some suitable constant $\rho_0 > 0$; see Theorem 6.1.

Moreover we show that, under a similar smallness condition on the coupling coefficients, now as $t \rightarrow -\infty$, if one of the bounded complete trajectories of (1.3) (which exists from the existence of the nonautonomous attractor) is bounded away from zero at $-\infty$, it is the unique such trajectory, and it also describes the unique pullback asymptotic behavior of all nonsemitrivial solutions of (1.3); see Theorem 6.2. When these two theorems can be applied together, there is a unique bounded complete trajectory $(u^*(t), v^*(t))$ that is both forwards and pullback attracting for (1.3); i.e., (u^*, v^*) is a bounded trajectory such that, for any $s \in \mathbf{R}$ and for any positive solution $(u(t, s), v(t, s))$ of (1.3) defined for $t > s$, one has

$$(1.5) \quad (u(t, s) - u^*(t), v(t, s) - v^*(t)) \rightarrow (0, 0) \quad \text{as } t \rightarrow \infty, \text{ or } s \rightarrow -\infty.$$

To obtain these results we need some nontrivial machinery for the linear scalar case, section 4.1, and some perturbation results about the exponential decay for solutions of linear parabolic nonautonomous systems, section 4.2. In particular, we find conditions guaranteeing that any bounded solution of

$$(1.6) \quad \begin{cases} u_t - d_1 \Delta u = p(t, x)u \\ v_t - d_2 \Delta v = q(t, x)v \end{cases}$$

gives rise to a solution that tends to zero as $t \rightarrow \infty$, when (1.6) is perturbed in a certain way; see Theorem 4.6. It is because we are able to study the linear part of the system in detail that we can obtain results for the nonlinear system.

Since we are able to treat the difference of two solutions of problem (1.3) within this framework, as a consequence of this argument we can apply our results to the Lotka–Volterra model in all three standard cases: competition, symbiosis, and prey-predator. It is noteworthy that these different situations are usually studied separately in the literature, but since we do not make any use of monotonicity arguments (which do not apply in the prey-predator case) we are able to give a unified treatment.

We close this paper in section 7 with a discussion of our results and some possibilities for further developments.

In the case in which all the coefficients are autonomous or periodic, our results in section 6 that we described above in (1.5) imply the uniqueness of the asymptotic behavior of all nonsemitrivial solutions.

Hence, in the autonomous case our results agree with all the classical results of uniqueness and stability of the nonsemitrivial steady states of (1.3) for the three cases of competition, symbiosis, and prey-predator (see, for instance, Theorem 4.4 in Furter and López-Gómez [13] and Corollary 4.3 in López-Gómez and Sabina de Lis [29] in the competition case, and Corollary 9.5 in Delgado et al. [12] in the symbiosis case).

Moreover, in the prey-predator case, with (1.5) we are able to conclude the uniqueness and *global stability* of a steady state, solving (for particular ranges of parameter values) one of the most interesting open problems in this field. We emphasize that this result is new even in the autonomous case, where until now only local stability has been proved; see Theorem 4.1 in Leung [26], see also Lakos [22], López-Gómez and Pardo [28], and Yamada [36].

2. Some notations and preliminaries. In this section we introduce some basic notations and terminology that will be used throughout the rest of the paper. In particular, we make precise the way systems (1.1) or (1.3) are said to be permanent.

2.1. Asymptotic behavior and complete trajectories for nonlinear systems. Note that if the solutions of (1.1) are global, then we can define a nonautonomous nonlinear *process* in some Banach space X appropriate for the solutions, i.e., a family of mappings $\{S(t, s)\}_{t \geq s} : X \rightarrow X$, $t, s \in \mathbf{R}$ satisfying:

- (a) $S(t, s)S(s, \tau)z = S(t, \tau)z$, for all $\tau \leq s \leq t$, $z \in X$,
- (b) $S(t, \tau)z$ is continuous in $t > \tau$ and z , and
- (c) $S(t, t)$ is the identity in X for all $t \in \mathbf{R}$.

$S(t, \tau)z$ arises as the value of the solution of our nonautonomous system at time t with initial condition z at initial time τ . For an autonomous system the solutions depend only on $t - \tau$, and we can write $S(t, \tau) = S(t - \tau, 0)$.

In order to describe the asymptotic behavior of nonautonomous systems like (1.1) and (1.3), we rely on the concept of a nonautonomous pullback attractor (Chepyzhov and Vishik [9], Kloeden and Schmalfuss [21]), which is the sensible generalization of an attractor for nonautonomous systems. For $A, B \subset X$ we denote the Hausdorff semidistance between A and B by $\text{dist}(A, B) = \sup_{a \in A} \inf_{b \in B} d(a, b)$.

DEFINITION 2.1. *We say that a family of compact sets $\{\mathcal{A}(t)\}_{t \in \mathbf{R}} \subset X$ is a pullback attractor associated to S if*

- (a) $S(t, \tau)\mathcal{A}(\tau) = \mathcal{A}(t)$, for all $t \geq \tau$ and
- (b) for all $t \in \mathbf{R}$ and $D \subset X$ bounded

$$\lim_{\tau \rightarrow -\infty} \text{dist}(S(t, \tau)D, \mathcal{A}(t)) = 0.$$

Observe that the attraction in (b) fixes the final time and moves the initial time backwards towards $-\infty$. We are not evolving one trajectory backwards in time, but rather we consider the current state of the system (at the fixed time t) which would result from the same initial condition starting at earlier and earlier times.

To guarantee the existence of such a pullback attractor, one is usually faced with the task of proving the existence of a pullback absorbing family, defined as follows.

DEFINITION 2.2. *Given $t_0 \in \mathbf{R}$, we say that $B(t_0) \subset X$ is pullback absorbing at time t_0 if for every bounded $D \subset X$ there exists a $T = T(t, D) \in \mathbf{R}$ such that*

$$S(t_0, \tau)D \subset B(t_0), \text{ for all } \tau \leq T.$$

A family $\{B(t)\}_{t \in \mathbf{R}}$ is pullback absorbing if $B(t_0)$ is pullback absorbing at time t_0 , for all $t_0 \in \mathbf{R}$.

The general result on the existence of nonautonomous pullback attractors is a generalization of the abstract theory for autonomous dynamical systems (Temam [37], Hale [14]).

THEOREM 2.3 (Crauel et al. [11], Schmalfuss [35]). *Assume that there exists a family of compact pullback absorbing sets. Then, there exists a pullback attractor $\{\mathcal{A}(t)\}_{t \in \mathbf{R}}$ that is minimal in the sense that if $\{C(t)\}_{t \in \mathbf{R}}$ is another family of closed pullback attracting sets, then $\mathcal{A}(t) \subset C(t)$ for all $t \in \mathbf{R}$.*

To have a more precise description of the dynamical objects within the pullback attractor, we make the following definition:

DEFINITION 2.4. *Let S be a process. We call the continuous map $w : \mathbf{R} \rightarrow X$ a complete trajectory if, for all $s \in \mathbf{R}$,*

$$S(t, s)w(s) = w(t) \quad \text{for all } t \geq s.$$

According to Chepyzhov and Vishik [9], when the family of absorbing sets is uniformly bounded, the pullback attractor can be characterized as

$$(2.1) \quad \mathcal{A}(t) = \{w(t) : w(\cdot) \text{ is a bounded complete trajectory for } S\}.$$

2.2. Pullback and forwards permanence for nonautonomous systems.

Consider the nonlinear system (1.1) and assume that f and g are regular functions. Hence, we can assume that for initial data $(u_s, v_s) \in C_{\mathcal{B}_1}(\overline{\Omega}) \times C_{\mathcal{B}_2}(\overline{\Omega})$ there exists a unique (local) smooth solution such that $(u, v) \in C_{\mathcal{B}_1}^1(\overline{\Omega}) \times C_{\mathcal{B}_2}^1(\overline{\Omega})$ for $t > s$, where, for $j = 0, 1$,

$$C_{\mathcal{B}}^j(\overline{\Omega}) = \begin{cases} C_0^j(\overline{\Omega}) & \text{for Dirichlet BCs,} \\ C^j(\overline{\Omega}) & \text{for Neumann or Robin BCs,} \end{cases}$$

with $C_0^j(\overline{\Omega})$ denoting functions in $C^j(\overline{\Omega})$ that are zero on $\partial\Omega$ and $C^0(\overline{\Omega}) = C(\overline{\Omega})$.

Note that in practice we will be interested only in nonnegative solutions and that if $u_s \geq 0$ and $v_s \geq 0$ in (1.1), then the local solution satisfies $u, v \geq 0$. In fact, the maximum principle implies that if both $u_s \geq 0$ and $v_s \geq 0$ are nontrivial, then u and v are strictly positive in Ω .

Although at this point we assume only local existence of solutions, it still makes sense to consider complete trajectories of (1.1), which roughly speaking are solutions defined for all times. These objects will play a central role in our analysis below, as can be seen from (2.1). More precisely, a restatement of Definition 2.4 gives the following.

DEFINITION 2.5. *A continuous function $U = \begin{pmatrix} u \\ v \end{pmatrix} : \mathbf{R} \rightarrow C_{\mathcal{B}_1}(\overline{\Omega}) \times C_{\mathcal{B}_2}(\overline{\Omega})$ is a complete trajectory of (1.1), if for all $s < t$ in \mathbf{R} , $(u(t), v(t))$ is the solution of (1.1) with initial data $u_s = u(s)$, $v_s = v(s)$.*

Now we define several concepts that will help us in making precise the concepts of pullback and forwards permanence for the solutions of (1.1) or (1.3). Note that the concepts below are related to the spaces $C_{\mathcal{B}_i}(\overline{\Omega})$ above. We start with the following.

DEFINITION 2.6. *A set of nonnegative functions $B \subset C(\overline{\Omega})$ is bounded away from zero if there exists a nonnegative nontrivial continuous function $\varphi_0(x) \geq 0$ in Ω (vanishing on $\partial\Omega$ in case of Dirichlet boundary conditions) such that*

$$u(x) \geq \varphi_0(x) \quad \text{for all } x \in \Omega, \quad u \in B.$$

The set B is nondegenerate if the function $\varphi_0(x)$ above is in $C^1(\overline{\Omega})$ and $\varphi_0(x) > 0$ in Ω .

Note that φ_0 above can be a positive constant in the case of Neumann or Robin boundary conditions.

Then we have the following definitions for curves in the space of continuous functions.

DEFINITION 2.7. *A positive function with values in $C(\overline{\Omega})$ is nondegenerate at ∞ (respectively, $-\infty$) if there exists $t_0 \in \mathbf{R}$ such that u is defined in $[t_0, \infty)$ (respectively, $(-\infty, t_0]$) and*

$$\{u(t), t \geq t_0\} \text{ is a nondegenerate set}$$

(respectively, for $t \leq t_0$), that is, there exists a $C^1(\overline{\Omega})$ function $\varphi_0(x) > 0$ in Ω , (vanishing on $\partial\Omega$ in case of Dirichlet boundary conditions), such that

$$u(t, x) \geq \varphi_0(x) \quad \text{for all } x \in \Omega, \quad t \geq t_0$$

(respectively, for all $t \leq t_0$).

A family of curves in $C(\overline{\Omega})$, denoted $\{u_\sigma(t)\}_{\sigma \in \Sigma}$, is nondegenerate at ∞ if there exists $t_0 \in \mathbf{R}$ such that u_σ is defined in $[t_0, \infty)$ and

$$\{u_\sigma(t), t \geq t_0, \sigma \in \Sigma\} \text{ is a nondegenerate set.}$$

Finally, a family of curves in $C(\overline{\Omega})$, denoted $\{u_\sigma(t, s)\}_{\sigma \in \Sigma}$, defined in the intervals $[s, \infty)$ is nondegenerate as $s \rightarrow -\infty$ if there exists $s_0 \in \mathbf{R}$ such that for all $s \leq s_0$

$$\{u_\sigma(t), s \leq t \leq s_0, \sigma \in \Sigma\} \text{ is a nondegenerate set.}$$

For systems, analogously to Definition 2.6, a set $B \subset (C(\overline{\Omega}))^2$ is bounded away from zero if each projection of B is bounded away from zero in $C(\overline{\Omega})$. In a similar way, as in Definition 2.7, a family of curves $U_\sigma(x, \cdot) \in (C(\overline{\Omega}))^2$, $\sigma \in \Sigma$, is nondegenerate if both components are nondegenerate in $C(\overline{\Omega})$.

Now we can finally define when the system (1.1) or (1.3) is pullback permanent. Observe that we assume here that solutions are globally defined.

DEFINITION 2.8. *We say that system (1.1) is pullback permanent if for any bounded set of initial $B \subset (C(\overline{\Omega}))^2$ bounded away from zero, there exists $t_0 \in \mathbf{R}$ such that for any $t \leq t_0$ the family of solutions*

$$(2.2) \quad \{(u(t, s; u_0, v_0), v(t, s; u_0, v_0)), s \leq t, (u_0, v_0) \in B\}$$

is nondegenerate as $s \rightarrow -\infty$.

The system (1.1) is uniformly pullback permanent if it is pullback permanent and the functions φ_0 in Definition 2.7 are independent of B .

Note that using the regularizing properties of the solutions of (1.1) or (1.3), if the system is pullback permanent, as defined above, then the set (2.2) is nondegenerate as $s \rightarrow -\infty$ for any fixed $t \in \mathbf{R}$.

In an analogous although subtly different way, we can define when system (1.1) or (1.3) is forwards permanent.

DEFINITION 2.9. *We say that system (1.1) is forwards permanent if for any bounded set of initial $B \subset (C(\overline{\Omega}))^2$ bounded away from zero, and for any $s \in \mathbf{R}$, the family of solutions*

$$(2.3) \quad \{(u(t, s; u_0, v_0), v(t, s; u_0, v_0)), s \leq t, (u_0, v_0) \in B\}$$

is nondegenerate at ∞ .

The system (1.1) is uniformly forwards permanent if it is forwards permanent and the functions φ_0 in Definition 2.7 are independent of B .

Note that (1.1) always has the trivial solution $(0, 0)$ as well as semitrivial solutions $(u, 0)$ and $(0, v)$. Hence, if the system is permanent, as defined above, this implies that trivial and semitrivial solutions are unstable in the pullback or forwards sense; see, e.g., Langa, Robinson, and Suárez [24]. Also, note that permanence implies coexistence of the species, since the values of the solutions eventually remain far from zero in all points of the domain (except at the boundary in the case of Dirichlet boundary conditions).

In the next section we will give conditions on the coefficients of (1.3) for uniform permanence (both forwards and pullback), which will be moreover robust with respect to suitable perturbations on the coefficients.

3. Extinction and permanence for nonautonomous Lotka–Volterra equations: Competition, symbiosis, and prey–predator models. In this section we give results on extinction and pullback and forwards permanence for nonautonomous Lotka–Volterra systems of the type

$$(3.1) \quad \begin{cases} u_t - d_1 \Delta u = u(\lambda(t, x) - a(t, x)u - b(t, x)v), & x \in \Omega, t > s \\ v_t - d_2 \Delta v = v(\mu(t, x) - c(t, x)u - d(t, x)v), & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0, & x \in \partial\Omega, t > s \\ u(s) = u_s \geq 0, v(s) = v_s \geq 0, \end{cases}$$

with $d_1, d_2 > 0$; $\lambda, \mu, a, b, c, d \in C^\theta(\bar{Q})$, and $\bar{Q} = \mathbf{R} \times \bar{\Omega}$. Given a function $e \in C^\theta(\bar{Q})$, we define

$$e_L := \inf_{\bar{Q}} e(t, x) \quad e_M := \sup_{\bar{Q}} e(t, x).$$

We assume from now on that

$$(3.2) \quad a_L, d_L > 0,$$

and consider the three classical cases depending on the signs of b and c :

1. *Competition*: $b_L, c_L > 0$ in \bar{Q} .
2. *Symbiosis*: $b_M, c_M < 0$ in \bar{Q} .
3. *Prey–predator*: $b_L > 0, c_M < 0$ in \bar{Q} .

Also, note that we consider all nine possible choices for \mathcal{B}_i as in (1.2).

Using standard techniques, see for instance Pao [30], it can be shown that given $0 \leq u_s \in C(\bar{\Omega}), 0 \leq v_s \in C(\bar{\Omega})$ there exists, locally in time, a unique solution of (3.1) which is nonnegative, and which we will denote by

$$u = u(t, s, x; u_s, v_s) \geq 0, \quad v = v(t, s, x; u_s, v_s) \geq 0.$$

In fact, due to the strong maximum principle, if $u_s \geq 0$ and $v_s \geq 0$ are both nontrivial, then u and v are strictly positive in Ω . Furthermore, if we denote by \mathcal{C}_i and $\text{int}(\mathcal{C}_i)$ for $i = 1, 2$, respectively, the positive cones in $C^1_{\mathcal{B}_i}(\bar{\Omega})$ and their corresponding interior sets, we have

$$\text{int}(\mathcal{C}_i) := \{u \in \mathcal{C}_i : u > 0 \text{ in } \Omega, \text{ and } \frac{\partial u}{\partial \bar{n}} < 0 \text{ on } \partial\Omega\} \text{ if } \mathcal{B}_i u = u$$

and

$$\text{int}(\mathcal{C}_i) := \{u \in \mathcal{C}_i : u \geq \delta > 0, \text{ for some } \delta > 0 \text{ in } \bar{\Omega}\},$$

if $\mathcal{B}_i u = \frac{\partial u}{\partial \bar{n}}$ or $\mathcal{B}_i u = d_i \frac{\partial u}{\partial \bar{n}} + \sigma_i(x)u$.

Thus, if $u_s \geq 0$ and $v_s \geq 0$ are both nontrivial, then $(u, v) \in \text{int}(C_1) \times \text{int}(C_2)$ for $t > s$.

Note that (3.1) also admits semitrivial solutions of the form $(u, 0)$ or $(0, v)$. As indicated in the Introduction, the stability properties of semitrivial solutions play an important role in the global dynamics of (3.1). In fact, extinction requires some semitrivial solution is stable whereas permanence is only possible if semitrivial solutions are somehow unstable.

Thus, we first review some results on the solutions of scalar logistic equations that will be used further below. These results will be used to prove that the local solutions of (3.1) above are, in fact, globally defined. Also, they will be crucially used to prove the existence of a pullback attractor as in section 2.1, and to obtain our results on extinction and permanence as well.

3.1. On the nonautonomous logistic equation. Note that (3.1) always admits semitrivial solutions of the form $(u, 0)$ or $(0, v)$. In this case, when one species is not present, the other one satisfies the nonautonomous logistic equation

$$(3.3) \quad \begin{cases} u_t - d\Delta u = h(t, x)u - g(t, x)u^2 & \text{in } \Omega, t > s, \\ \mathcal{B}u = 0 & \text{on } \partial\Omega, \\ u(s) = u_s \geq 0 & \text{in } \Omega, \end{cases}$$

where $d > 0$ and \mathcal{B} as in (1.2), $u_s \in C(\overline{\Omega})$, $h, g \in C^\theta(\overline{Q})$, and

$$g_L > 0 \quad \text{in } \overline{Q}.$$

For $m \in L^\infty(\Omega)$ we denote by $\Lambda_{\mathcal{B}}(d, m)$ the first eigenvalue of

$$(3.4) \quad \begin{cases} -d\Delta u = \lambda u + m(x)u & \text{in } \Omega, \\ \mathcal{B}u = 0 & \text{on } \partial\Omega. \end{cases}$$

In particular, we denote by $\Lambda_{0, \mathcal{B}}(d) = \Lambda_{\mathcal{B}}(d, 0)$ the first eigenvalue of the operator $-d\Delta$ with boundary conditions \mathcal{B} . It is well known that $\Lambda_{\mathcal{B}}(d, m)$ is a simple eigenvalue and a continuous and decreasing function of m . Also note that if m_1 is constant, then

$$(3.5) \quad \Lambda_{\mathcal{B}}(d, m_1 + m_2) = \Lambda_{\mathcal{B}}(d, m_2) - m_1.$$

We write $\varphi_{1, \mathcal{B}}(d, m)$ for the positive eigenfunction associated to $\Lambda_{\mathcal{B}}(d, m)$, normalized such that $\|\varphi_{1, \mathcal{B}}(d, m)\|_{L^\infty(\Omega)} = 1$.

If there is no possible confusion we will suppress the dependence on d and \mathcal{B} in the notations above. When we need to distinguish these quantities with respect to \mathcal{B}_i , or d_i , $i = 1, 2$, we will employ superscripts as $\Lambda^i(m)$ or Λ_0^i .

Finally, for $h, g \in L^\infty(\Omega)$ with $g_L > 0$ consider the elliptic equation

$$(3.6) \quad \begin{cases} -d\Delta u = h(x)u - g(x)u^2 & \text{in } \Omega, \\ \mathcal{B}u = 0 & \text{on } \partial\Omega. \end{cases}$$

The following result is well known (Cantrell and Cosner [6]).

PROPOSITION 3.1. *If $\Lambda(h) \geq 0$, the unique nonnegative solution of (3.6) is the trivial one, i.e., $\omega_{[h, g]}(x) = 0$. On the other hand, if $\Lambda(h) < 0$ there exists a unique positive solution of (3.6), which we denote by $\omega_{[h, g]}(x)$. Moreover, $0 < \omega_{[h, g]}(x) \leq \Psi(x)$ in Ω , where*

$$\Psi(x) = \begin{cases} \frac{h_M}{g_L} & \text{for Dirichlet or Neumann BCs,} \\ -\frac{\Lambda(h)}{\varphi_L g_L} \varphi(x) & \text{for Robin BCs,} \end{cases}$$

with $\varphi = \varphi_{1, \mathcal{B}}(m)$.

The following result will be used in what follows.

LEMMA 3.2. Assume that $h_n \in L^\infty(\Omega)$ and that

$$h_n \rightarrow h_\infty \quad \text{in } L^\infty(\Omega),$$

with $\Lambda(h_\infty) < 0$. Then, there exist $n_0 \in \mathbb{N}$, and $\varphi \in \text{int}(\mathcal{C})$ such that

$$\varphi(x) \leq \omega_{[h_n, g]}(x) \quad \text{in } \Omega, \quad \text{for all } n \geq n_0,$$

where $\omega_{[h_n, g]}(x)$ is given by Proposition 3.1.

Proof. Since $\Lambda(h_\infty) < 0$, we can take $\varepsilon > 0$ such that $0 < \varepsilon < -\Lambda(h_\infty)$. For this $\varepsilon > 0$, there exists $n_0 \in \mathbb{N}$ such that for $n \geq n_0$

$$-\varepsilon < h_n - h_\infty < \varepsilon \quad \text{for all } x \in \Omega.$$

Consider $\varphi_\infty \in \text{int}(\mathcal{C})$ the eigenfunction associated to $\Lambda(h_\infty)$ with $\|\varphi_\infty\|_{L^\infty(\Omega)} = 1$. It is not hard to show that $\delta\varphi_\infty$ is a subsolution of (3.6) with $h = h_n$ provided that

$$\delta \leq -\frac{\varepsilon + \Lambda(h_\infty)}{g_M}.$$

So, $\delta\varphi_\infty(x) \leq \omega_{[h_n, g]}(x)$ in Ω . This completes the proof. \square

In [25] and [34] the following properties of (3.3) were proved.

THEOREM 3.3. Assume that in (3.3)

$$h_M < \infty \quad \text{and} \quad g_L > 0 \quad \text{in } \overline{\Omega}.$$

Then

1. For every nontrivial $u_s \in C(\overline{\Omega})$, $u_s \geq 0$, there exists a unique positive solution of (3.3) denoted by $\Theta_{[h, g]}(t, s, u_s)$. Moreover,

$$(3.7) \quad 0 \leq \Theta_{[h, g]}(t, s, u_s) \leq K,$$

where

$$K := \begin{cases} \max \left\{ (u_s)_M, \frac{h_M}{g_L} \right\} & \text{for Dirichlet or Neumann BCs,} \\ \max \left\{ \left(\frac{u_s}{\varphi} \right)_M, \frac{-\Lambda(h_M)}{\varphi_L g_L} \right\} & \text{for Robin BCs,} \end{cases}$$

and φ is the positive eigenfunction associated to $\Lambda(h_M)$ with $\|\varphi\|_{L^\infty(\Omega)} = 1$.

2. For fixed $t > s$, u_s , the map $h \mapsto \Theta_{[h, g]}(t, s, u_s)$ is increasing and $g \mapsto \Theta_{[h, g]}(t, s, u_s)$ is decreasing.

For fixed $t > s$, h and g , the map $u_s \mapsto \Theta_{[h, g]}(t, s, u_s)$ is increasing.

3. Define, for $x \in \Omega$,

$$h_0(x) := \inf_{t \in \mathbf{R}} h(t, x), \quad H_0(x) := \sup_{t \in \mathbf{R}} h(t, x)$$

and

$$g_0(x) := \inf_{t \in \mathbf{R}} g(t, x), \quad G_0(x) := \sup_{t \in \mathbf{R}} g(t, x).$$

Then, if $u_s \in \text{int}(\mathcal{C})$ and $\Lambda(h_0) < 0$ we have, for any $t > s$,

$$(3.8) \quad 0 < \varepsilon\varphi_1(x) \leq \Theta_{[h,g]}(t, s, x; u_s) \quad \text{in } \Omega,$$

where φ_1 is the positive eigenfunction associated to $\Lambda(h_0)$ and

$$\varepsilon = \varepsilon(u_s) := \min \left\{ \left(\frac{u_s}{\varphi_1} \right)_L, \frac{-\Lambda(h_0)}{g_M} \right\}.$$

4. If $\Lambda(H_0) > 0$, then for all initial data $u_s \geq 0$, $\Theta_{[h,g]}(t, s, u_s) \rightarrow 0$, in $C^1(\overline{\Omega})$, as $t - s \rightarrow \infty$. Moreover the convergence is exponential and uniform for bounded sets of initial data u_s .
5. If $\Lambda(h_0) < 0$, then there exists a unique bounded, complete, and nondegenerate trajectory at $\pm\infty$ of (3.3), $\varphi_{[h,g]}$, which moreover satisfies that for all s and any bounded set of nontrivial initial data $u_s \geq 0$, bounded away from 0,

$$\Theta_{[h,g]}(t, s, u_s) - \varphi_{[h,g]}(t) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

That is, $\varphi_{[h,g]}$ describes the forward behavior of all solutions. Also, $\varphi_{[h,g]}$ describes the pullback behavior of all nondegenerate solutions of (3.3), that is, for each t , if $s \mapsto u_s \geq 0$ is bounded and nondegenerate, then

$$\Theta_{[h,g]}(t, s, u_s) - \varphi_{[h,g]}(t) \rightarrow 0 \quad \text{as } s \rightarrow -\infty.$$

Both limits above are taken in $C^1(\overline{\Omega})$. Furthermore for all $t \in \mathbf{R}$, we have

$$\omega_{[h_0, G_0]}(x) \leq \varphi_{[h,g]}(t, x) \leq \omega_{[H_0, g_0]}(x) \quad \text{in } \Omega.$$

6. If h, g are independent of t and are in $L^\infty(\Omega)$ with $g_L > 0$ and $\Lambda(h) < 0$, then $\varphi_{[h,g]}(t, x) = \omega_{[h,g]}(x)$ is the unique positive solution of (3.6) and for all $t > s$ and u_s

$$\Theta_{[h,g]}(t, s, u_s) = \Theta_{[h,g]}(t - s, u_s) \rightarrow \omega_{[h,g]} \quad \text{in } C^1(\overline{\Omega}) \quad \text{as } t - s \rightarrow \infty$$

uniformly for bounded sets of initial data $u_s \geq 0$ bounded away from zero. In particular, there exist $m \leq 1 \leq M$ such that

$$m\omega_{[h,g]} \leq \Theta_{[h,g]}(t, s, u_s) \leq M\omega_{[h,g]},$$

for $t - s$ large.

Moreover in statements 4, 5, and 6 above the convergence as $t \rightarrow \infty$ is exponentially fast (see [33]).

3.2. Existence of the pullback attractor and complete trajectories for nonautonomous Lotka–Volterra systems. Our first purpose is to prove the existence of a nonautonomous pullback attractor for (3.1). To do this we will derive suitable estimates on the solutions of (3.1). In doing this we will use the following notation for the solutions of (3.3) with diffusion coefficients d_1 and d_2 and boundary conditions \mathcal{B}_1 and \mathcal{B}_2 , respectively,

$$\xi_{[\lambda,a]}(t, s) = \Theta_{[\lambda,a]}(t, s, u_s), \quad \eta_{[\mu,d]}(t, s) = \Theta_{[\mu,d]}(t, s, v_s),$$

where $u_s \geq 0$ and $v_s \geq 0$ in Ω .

THEOREM 3.4. *Provided that $a_L, d_L > 0$, for any solution (u, v) of (3.1), with initial data $u_s \geq 0, v_s \geq 0$, the following lower and upper bounds hold:*

1. *Competition*, $b_L > 0, c_L > 0$:

$$\xi_{[\lambda - b\eta_{[\mu, d], a}]} \leq u \leq \xi_{[\lambda, a]}, \quad \eta_{[\mu - c\xi_{[\lambda, a], d}]} \leq v \leq \eta_{[\mu, d]}.$$

2. *Symbiosis*, $b_M < 0, c_M < 0$: Assume that

$$(3.9) \quad b_L c_L < a_L d_L.$$

Then,

$$\xi_{[\lambda - b\eta_{[\mu, d], a}]} \leq u, \quad \eta_{[\mu - c\xi_{[\lambda, a], d}]} \leq v.$$

Assume furthermore that $d_1 = d_2$ and define

$$\gamma = \max\{\lambda_M, \mu_M\}, \quad M = \frac{a_L - c_L}{d_L - b_L} > 0, \quad K = \frac{a_L d_L - b_L c_L}{d_L - b_L} > 0,$$

and choose w_s such that $w_s \geq \max\{u_s, \frac{1}{M}v_s\}$. Denote by $\Theta_{[\gamma, K]}(t, s, w_s)$ the solution of (3.3) with $d = d_1$ and a certain boundary condition that depends on \mathcal{B}_1 and \mathcal{B}_2 and that will be specified in the proof. Then, we have the upper bounds

$$u \leq \Theta_{[\gamma, K]}(t, s, w_s), \quad v \leq M\Theta_{[\gamma, K]}(t, s, w_s).$$

3. *Prey-predator*, $b_L > 0, c_M < 0$:

$$\xi_{[\lambda - b\eta_{[\mu - c\xi_{[\lambda, a], d], a}]} \leq u \leq \xi_{[\lambda - b\eta_{[\mu, d], a}]} \leq \xi_{[\lambda, a]}, \quad \eta_{[\mu, d]} \leq v \leq \eta_{[\mu - c\xi_{[\lambda, a], d}]}$$

Proof. 1. Assume that $b_L, c_L > 0$. If we write the equation for u as

$$u_t - d_1 \Delta u = u(\lambda - bv) - au^2,$$

then using Theorem 3.3 we get

$$u = \xi_{[\lambda - bv, a]} \leq \xi_{[\lambda, a]},$$

and similarly,

$$v \leq \eta_{[\mu, d]}.$$

Hence, again by Theorem 3.3

$$u = \xi_{[\lambda - bv, a]} \geq \xi_{[\lambda - b\eta_{[\mu, d], a}]}$$

2. Assume now that $b_M, c_M < 0$. To have the lower bounds it is enough to check that in the equation for u one has

$$\xi_{[\lambda - b\eta_{[\mu, d], a}]}(\lambda - a\xi_{[\lambda - b\eta_{[\mu, d], a}]} - b\eta_{[\mu, d]}) \leq \xi_{[\lambda - b\eta_{[\mu, d], a}]}(\lambda - a\xi_{[\lambda - b\eta_{[\mu, d], a}]} - b\eta_{[\mu - c\xi_{[\lambda, a], d}]}),$$

or equivalently,

$$\eta_{[\mu - c\xi_{[\lambda, a], d}]} \geq \eta_{[\mu, d]},$$

which is true since $c < 0$. We treat the equation for v similarly.

On the other hand, assuming that $d_1 = d_2$, define

$$\bar{u} = \Theta_{[\gamma, K]}(t, s, w_s), \quad \bar{v} = M\Theta_{[\gamma, K]}(t, s, w_s)$$

with a suitable boundary condition, \mathcal{B} , to be described below. Then using the equations we get that \bar{u} and \bar{v} are supersolutions if

$$-K \geq -a - bM, \quad -K \geq -dM - c,$$

which is satisfied with the choice of M and K . To compare the solutions with the upper solutions on the boundary, if either u or v satisfies Dirichlet boundary conditions we take \mathcal{B} the boundary condition of the other component. If both u and v satisfy Robin or Neumann (i.e., $\sigma_i = 0$ in the latter case) boundary conditions we define

$$\sigma = \min\{\sigma_1, \sigma_2\},$$

and $\mathcal{B}u = d_1 \frac{\partial u}{\partial \bar{n}} + \sigma(x)u$.

3. Assume finally that $b_L > 0$, $c_M < 0$, then

$$u \leq \xi_{[\lambda, a]} \quad \text{and} \quad \eta_{[\mu, d]} \leq v.$$

Hence

$$v = \eta_{[\mu - c u, d]} \leq \eta_{[\mu - c \xi_{[\lambda, a]}, d]},$$

and then,

$$u = \xi_{[\lambda - b v, a]} \geq \xi_{[\lambda - b \eta_{[\mu - c \xi_{[\lambda, a]}, d]}, a]}. \quad \square$$

With the upper bounds in Theorem 3.4 and using the results for scalar logistic equations in Theorem 3.3, we get the following result.

THEOREM 3.5. *Under the assumptions in cases (1)–(3) of Theorem 3.4, all solutions of (3.1) are global in time and moreover there exists a pullback attractor $\mathcal{A}(t)$ of (3.1), which is bounded for all $t \in \mathbf{R}$. More precisely, we have*

$$\limsup_{t-s \rightarrow \infty} u(t, s; u_s, v_s) \leq M_\infty, \quad \limsup_{t-s \rightarrow \infty} v(t, s; u_s, v_s) \leq N_\infty,$$

uniformly in Ω and for bounded sets of initial data $u_s, v_s \geq 0$, for some constants $M_\infty \geq 0$ and $N_\infty \geq 0$ that depend on the coefficients of (3.1).

In particular, there exists at least one complete bounded trajectory $(u^*(t), v^*(t))$, $t \in \mathbf{R}$, for (3.1). Furthermore, all complete bounded trajectories of (3.1) are uniformly bounded by M_∞ and N_∞ and for all $t \in \mathbf{R}$.

Proof. Thanks to the upper bounds in Theorem 3.4, the positive solutions of (3.1) are always bounded by solutions of the logistic equation of the type (3.3). In particular, all solutions of (3.1) are globally defined.

Now we use that

$$0 \leq \Theta_{[\alpha, \beta]}(t, s; z) \leq \Theta_{[\alpha_M, \beta_L]}(t - s; z),$$

statements (4)–(6) in Theorem 3.3, and that $0 \leq \omega_{[\alpha_M, \beta_L]}(x) \leq \Psi_M$, with ω and $\Psi(x)$ as in Proposition 3.1, to get the estimates.

In particular, this implies the existence of bounded pullback absorbing sets for (3.1) in $C(\bar{\Omega}) \times C(\bar{\Omega})$.

Then following the proof of section 6 in Langa et al. [25], we can show the existence of a bounded pullback absorbing set in $C^1(\bar{\Omega}) \times C^1(\bar{\Omega})$, and so compact in $C(\bar{\Omega}) \times C(\bar{\Omega})$. Hence, we conclude using Theorem 2.3 as the existence of a bounded nonautonomous pullback attractor $\mathcal{A}(t)$, and thus the existence of at least one bounded complete trajectory $(u^*(t), v^*(t))$, $t \in \mathbf{R}$, follows. \square

3.3. Extinction for nonautonomous Lotka–Volterra systems. Note that with the arguments above there are some cases, when statement 4 in Theorem 3.3 can be used, in which one (or both) constants M_∞ and N_∞ are zero and we have then extinction of one of the species. This implies, in turn, that the semitrivial (or the trivial) solutions are stable in a forwards and pullback senses. More precisely, we have the following result. Observe that these sufficient conditions are far from optimal but qualitatively they describe the global stability of trivial or semitrivial solutions.

PROPOSITION 3.6. *With the notations in Theorems 3.4 and 3.5, we have*

1. *Competition, $b_L > 0, c_L > 0$. If*

$$\lambda_M < \Lambda_0^1, \quad \text{then } M_\infty = 0,$$

while if

$$\mu_M < \Lambda_0^2, \quad \text{then } N_\infty = 0.$$

2. *Symbiosis, $b_M < 0, c_M < 0, d_1 = d_2$ and (3.9), that is $b_L c_L < a_L d_L$. If*

$$\gamma < \Lambda_0^1, \quad \text{then } M_\infty = 0,$$

while if

$$\gamma < \Lambda_0^2, \quad \text{then } N_\infty = 0.$$

3. *Prey-predator, $b_L > 0, c_M < 0$. If*

$$\lambda_M < \Lambda_0^1, \quad \text{then } M_\infty = 0,$$

and in this case, if

$$\mu_M < \Lambda_0^2, \quad \text{then } M_\infty = 0.$$

On the other hand, if

$$\Lambda_0^1 < \lambda_M, \quad \text{and } \mu_M - c_L \frac{\lambda_M}{a_L} < \Lambda_0^2, \quad \text{then } N_\infty = 0.$$

In all the cases, when $M_\infty = 0$ the u component of the solutions of (3.1) extinguishes in pullback and forwards senses, while the v component of the solutions asymptotically follows the dynamics of the scalar logistic equation (3.3) with $h(t, x) = \mu(t, x)$ and $g(t, x) = d(t, x)$ as described in Theorem 3.3.

The case when $N_\infty = 0$ is analogous.

Proof. In fact, in the case of competition we have $0 \leq u \leq \xi_{[\lambda_M, a_L]}$ and $0 \leq v \leq \eta_{[\mu_M, d_L]}$. Hence, from statement 4 in Theorem 3.3 and using (3.5), if $\Lambda^1(\lambda_M) = \Lambda_0^1 - \lambda_M > 0$, then $M_\infty = 0$, while $N_\infty = 0$ if $\Lambda^2(\mu_M) = \Lambda_0^2 - \mu_M > 0$.

In the case of symbiosis, assuming $d_1 = d_2$, we have $0 \leq u \leq \Theta_{[\gamma, K]}(t, s, w_s)$, $0 \leq v \leq M\Theta_{[\gamma, K]}(t, s, w_s)$. Hence, if $\Lambda^1(\gamma) = \Lambda_0^1 - \gamma > 0$, then $M_\infty = 0$, while $N_\infty = 0$ if $\Lambda^2(\gamma) = \Lambda_0^2 - \gamma > 0$.

Finally, in the case of prey-predator, we have $0 \leq u \leq \xi_{[\lambda_M, a_L]}$, $0 \leq v \leq \eta_{[\mu_M - c_L \xi_{[\lambda_M, a_L]}, d_L]}$. Hence, if $\Lambda^1(\lambda_M) = \Lambda_0^1 - \lambda_M > 0$, then $M_\infty = 0$. In this case, $N_\infty = 0$ if $\Lambda^2(\mu_M) = \Lambda_0^2 - \mu_M > 0$.

On the other hand, if $\Lambda_0^1 < \lambda_M$, then for large values of $t - s$ we have $v \leq \eta_{[\mu_M - c_L(\omega_{[\lambda_M, a_L]} + \varepsilon), d_L]}$, and then $N_\infty = 0$ if $\Lambda^2(\mu_M - c_L \frac{\lambda_M}{a_L}) = \Lambda_0^2 - \mu_M + c_L \frac{\lambda_M}{a_L} > 0$.

The rest is immediate. \square

As we are interested in the “permanence” problem for (3.1), we will consider in what follows only the cases in which $M_\infty > 0$ and $N_\infty > 0$. In particular, note that for sufficiently large values of $\lambda_M > 0$ and $\mu_M > 0$ we can take, for the case of Dirichlet or Neumann boundary conditions in either one of the components u or v ,

$$M_\infty = \begin{cases} \frac{\lambda_M}{a_L} & \text{in the competition case,} \\ \frac{\gamma}{K} & \text{in the symbiosis case,} \\ \frac{\lambda_M}{a_L} & \text{in the prey-predator case,} \end{cases}$$

$$N_\infty = \begin{cases} \frac{\mu_M}{d_L} & \text{in the competition case,} \\ M \frac{\gamma}{K} & \text{in the symbiosis case,} \\ \frac{\mu_M - c_L \frac{\lambda_M}{a_L}}{d_L} & \text{in the prey-predator case,} \end{cases}$$

while for Robin boundary conditions we have

$$M_\infty = \begin{cases} \frac{\lambda_M - \Lambda_0^1}{(\varphi^1)_L a_L} & \text{in the competition case,} \\ \frac{\gamma - \Lambda_0^1}{(\varphi^1)_L K} & \text{in the symbiosis case,} \\ \frac{\lambda_M - \Lambda_0^1}{(\varphi^1)_L a_L} & \text{in the prey-predator case,} \end{cases}$$

$$N_\infty = \begin{cases} \frac{\mu_M - \Lambda_0^2}{(\varphi^2)_L d_L} & \text{in the competition case,} \\ M \frac{\gamma - \Lambda_0^2}{(\varphi^2)_L K} & \text{in the symbiosis case,} \\ \frac{\mu_M - c_L \frac{\lambda_M - \Lambda_0^1}{(\varphi^1)_L a_L} - \Lambda_0^2}{(\varphi^2)_L d_L} & \text{in the prey-predator case,} \end{cases}$$

where φ^i denotes the positive eigenfunction associated to Λ_0^i with $\|\varphi^i\|_{L^\infty(\Omega)} = 1$. Note that similar expressions can be given in the remaining five cases for the boundary conditions, although their explicit form becomes more cumbersome.

In fact, in the next section we will impose conditions on the coefficients to ensure that the pullback and forwards behavior of the solutions of (3.1), with nontrivial initial data, is far from the semitrivial and the trivial solutions.

3.4. Permanence for nonautonomous Lotka–Volterra systems: Nondegeneracy of solutions. Now, using the lower bounds in Theorem 3.4, we will give sufficient conditions for the system (3.1) to be uniformly permanent in pullback and forwards senses, as in section 2.2. For reasons that will become clear further below, we are interested in obtaining such nondegeneracy in a uniform way with respect to the coefficients λ, μ, a, b, c, d in the system. For this, recall the notations in (3.4) and that we always take nonnegative nontrivial initial data u_s, v_s .

Also note that in the results of this section we will use the quantities $\lambda_I \leq \lambda_S$, $\mu_I \leq \mu_S$, $a_I \leq a_S$, $b_I \leq b_S$, $c_I \leq c_S$, and $d_I \leq d_S$ to control the asymptotic sizes of

the coefficients λ, μ, a, b, c, d as $t \rightarrow \pm\infty$. As all the results will be given in terms of such quantities, the statements below show the robustness of the results with respect to perturbations in the coefficients of the system.

Finally, we stress here once again that the results below imply the instability of trivial and semitrivial solutions.

3.4.1. Competition.

PROPOSITION 3.7 (forwards permanence—competitive case). *Assume (3.2) and $b_L, c_L > 0$. Then:*

(i) *If $\lambda_I > \Lambda^1(-b_S\omega_{[\mu_S, d_I]})$ there exists $\psi_{11} \in \text{int}(C_1)$ such that whenever*

$$\lambda(t, x) \geq \lambda_I, \mu(t, x) \leq \mu_S, b(t, x) \leq b_S, a(t, x) \leq a_S, \text{ and } d(t, x) \geq d_I > 0$$

for all $x \in \Omega$ and $t \geq t_0$, for any $u_s, v_s > 0$, the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

(ii) *If $\mu_I > \Lambda^2(-c_S\omega_{[\lambda_S, a_I]})$ there exists $\psi_{22} \in \text{int}(C_2)$ such that whenever*

$$\lambda(t, x) \leq \lambda_S, \mu(t, x) \geq \mu_I, a(t, x) \geq a_I > 0, d(t, x) \leq d_S, c(t, x) \leq c_S$$

for all $x \in \Omega$ and $t \geq t_0$, for any $u_s, v_s > 0$, the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

Hence, if

$$(3.10) \quad \lambda_I > \Lambda^1(-b_S\omega_{[\mu_S, d_I]}) \quad \text{and} \quad \mu_I > \Lambda^2(-c_S\omega_{[\lambda_S, a_I]}),$$

then there exist $\psi_{11} \in \text{int}(C_1)$ and $\psi_{22} \in \text{int}(C_2)$ such that for any choice of coefficients that satisfy

$$\begin{aligned} \lambda_I \leq \lambda(t, x) \leq \lambda_S, \quad \mu_I \leq \mu(t, x) \leq \mu_S, \quad 0 < a_I \leq a(t, x) \leq a_S, \\ 0 < b_I \leq b(t, x) \leq b_S, \quad 0 < c_I \leq c(t, x) \leq c_S, \quad 0 < d_I \leq d(t, x) \leq d_S, \end{aligned}$$

for all $x \in \Omega$ and for all $t \geq t_0$, and for all nontrivial $u_s \geq 0, v_s \geq 0$ in a fixed bounded set of $C(\bar{\Omega})$ bounded away from 0, the solution (u, v) of (3.1) for $t > s \geq t_0$ is nondegenerate at ∞ and for all $t - s$ large enough,

$$u(t, s, x; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s, x; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly forwards permanent.

Proof. Since $\lambda_I > \Lambda^1(-b_S\omega_{[\mu_S, d_I]})$, by the continuity of $\Lambda^1(m)$ with respect to m , there exists $\varepsilon > 0$ such that

$$\lambda_I > \Lambda^1(-b_S(\omega_{[\mu_S, d_I]} + \varepsilon)) \quad \text{or equivalently by (3.5)} \quad \Lambda^1(\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon)) < 0.$$

Using Theorems 3.3 and 3.4, we get, for $t > s \geq t_0$,

$$u(t, s, u_s, v_s) \geq \xi_{[\lambda - b\eta_{[\mu, d]}, a]}(t, s, u_s) \geq \Theta_{[\lambda_I - b_S\eta_{[\mu_S, d_I]}, a_S]}(t - s, u_s).$$

Moreover, $\eta_{[\mu_S, d_I]}(t, s, v_s) \rightarrow \omega_{[\mu_S, d_I]}$ in $C^1(\bar{\Omega})$ and uniformly for v_s in bounded sets bounded away from zero, as $t - s \rightarrow \infty$, and so

$$(3.11) \quad u(t, s, u_s, v_s) \geq \Theta_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]}(t - s, u_s) \rightarrow \omega_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]}$$

in $C^1(\bar{\Omega})$ and uniformly for u_s in bounded sets bounded away from zero, as $t - s \rightarrow \infty$ by Theorem 3.3 and where we have used (3.10). Hence, the result follows for u .

On the other hand, we have analogously for the v component, for $t > s \geq t_0$,

$$v(t, s, u_s, v_s) \geq \eta_{[\mu - c\xi_{[\lambda, a]}, d]}(t, s, v_s) \geq \Theta_{[\mu_I - c_S \xi_{[\lambda_S, a_I]}, d_S]}(t - s, v_s).$$

Now, from (3.10), $\xi_{[\lambda_S, a_I]}(t, s, u_s) \rightarrow \omega_{[\lambda_S, a_I]}$ in $C^1(\bar{\Omega})$ and uniformly for u_s in bounded sets bounded away from zero, as $t - s \rightarrow \infty$, and so

$$(3.12) \quad v(t, s, u_s, v_s) \geq \Theta_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]}(t - s, v_s) \rightarrow \omega_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]}$$

in $C^1(\bar{\Omega})$ and uniformly for v_s in bounded sets bounded away from zero, as $t - s \rightarrow \infty$ by Theorem 3.3. \square

The same arguments as above, carried out in all pullback sense lead to the following result. Note that, in particular, this proposition guarantees the uniform nondegeneracy at $-\infty$ of complete nondegenerate trajectories with respect to the coefficients in the system.

PROPOSITION 3.8 (pullback permanence—competitive case). *Assume (3.2) and $b_L, c_L > 0$. Then:*

(i) *If $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_S, d_I]})$ there exists $\psi_{11} \in \text{int}(C_1)$ such that whenever*

$$\lambda(t, x) \geq \lambda_I, \quad \mu(t, x) \leq \mu_S, \quad b(t, x) \leq b_S, \quad a(t, x) \leq a_S, \quad d(t, x) \geq d_I > 0$$

for all $x \in \Omega$ and $t \leq t_0$ (for some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $u(t, x) \geq \psi_{11}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

(ii) *If $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_S, a_I]})$ there exists $\psi_{22} \in \text{int}(C_2)$ such that whenever*

$$\lambda(t, x) \leq \lambda_S, \quad \mu(t, x) \geq \mu_I, \quad a(t, x) \geq a_I > 0, \quad d(t, x) \leq d_S, \quad c(t, x) \leq c_S$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $v(t, x) \geq \psi_{22}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

Hence, if

$$(3.13) \quad \lambda_I > \Lambda^1(-b_S \omega_{[\mu_S, d_I]}) \quad \text{and} \quad \mu_I > \Lambda^2(-c_S \omega_{[\lambda_S, a_I]})$$

there exist functions $\psi_{11} \in \text{int}(C_1)$ and $\psi_{22} \in \text{int}(C_2)$ such that whenever

$$\lambda_I \leq \lambda(t, x) \leq \lambda_S, \quad \mu_I \leq \mu(t, x) \leq \mu_S, \quad 0 < a_I \leq a(t, x) \leq a_S, \\ 0 < b_I \leq b(t, x) \leq b_S, \quad 0 < c_I \leq c(t, x) \leq c_S, \quad 0 < d_I \leq d(t, x) \leq d_S,$$

for all $x \in \Omega$ and $t \leq t_0$ (for some $t_0 \in \mathbf{R}$), and for all nontrivial $u_s \geq 0, v_s \geq 0$ in a fixed bounded set, B , of $C(\bar{\Omega})$ bounded away from 0, the set of solutions of (3.1) $\{(u, v), s < t \leq t_0, (u_s, v_s) \in B\}$ is nondegenerate as $s \rightarrow -\infty$ and for all $t - s$ large enough

$$u(t, s, x; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s, x; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly pullback permanent and any bounded complete trajectory that is nondegenerate at $-\infty$ satisfies

$$u(t, x) \geq \psi_{11}(x) \quad \text{and} \quad v(t, x) \geq \psi_{22}(x) \quad \text{for all } x \in \Omega \text{ and } t \leq t_0.$$

Proof. The first part of the statements follow from (3.11) and (3.12), with $t - s \rightarrow \infty$ but now $s < t \leq t_0$.

For a complete solution, arguing as in Proposition 3.7 we get for any $t_0 \geq t > s$,

$$u(t) \geq \xi_{[\lambda - b\eta_{[\mu, d]}, a]}(t, s, u(s)) \geq \Theta_{[\lambda_I - b_S \eta_{[\mu_S, d_I]}, a_S]}(t - s, u(s)).$$

As v is nondegenerate at $-\infty$, 5 in Theorem 3.3 implies $\eta_{[\mu_S, d_I]}(t, s, v(s)) \rightarrow \omega_{[\mu_S, d_I]}$ in $C^1(\bar{\Omega})$ as $s \rightarrow -\infty$. Thus, for sufficiently negative s ,

$$(3.14) \quad u(t) \geq \Theta_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]}(t - s, u(s)) \rightarrow \omega_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]}$$

in $C^1(\bar{\Omega})$ as $s \rightarrow -\infty$, because u is nondegenerate at $-\infty$ and 5 in Theorem 3.3 again. Hence the result follows for u .

On the other hand, we have analogously for the v component for any $t_0 \geq t > s$,

$$v(t) \geq \eta_{[\mu - c\xi_{[\lambda, a]}, d]}(t, s, v(s)) \geq \Theta_{[\mu_I - c_S \xi_{[\lambda_S, a_I]}, d_S]}(t - s, v(s)).$$

Now, $\xi_{[\lambda_S, a_I]}(t, s, u(s)) \rightarrow \omega_{[\lambda_S, a_I]}$ in $C^1(\bar{\Omega})$ as $s \rightarrow -\infty$, because u is nondegenerate at $-\infty$, and so, for sufficiently negative s ,

$$(3.15) \quad v(t) \geq \Theta_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]}(t - s, v(s)) \rightarrow \omega_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]}$$

in $C^1(\bar{\Omega})$ as $s \rightarrow -\infty$ by Theorem 3.3, because v is nondegenerate at $-\infty$. \square

Results for the other cases can be proved analogously, as we now show.

3.4.2. Symbiosis. First for the case of symbiosis, we have the following result. Note that as we make no use here of the upper bound in Theorem 3.4, we do not assume below that $d_1 = d_2$.

PROPOSITION 3.9 (forwards permanence—symbiotic case). *Assume (3.2), $b_M, c_M < 0$, and (3.9), that is*

$$b_L c_L < a_L d_L.$$

Then:

- (i) *If $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_I, d_S]})$ there exists $\psi_{11} \in \text{int}(\mathcal{C}_1)$ such that whenever*

$$\lambda(t, x) \geq \lambda_I, \quad \mu(t, x) \geq \mu_I, \quad b(t, x) \leq b_S < 0, \quad a(t, x) \leq a_S, \quad d(t, x) \leq d_S$$

for all $x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$ the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

- (ii) *If $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]})$ there exists $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever*

$$\lambda(t, x) \geq \lambda_I, \quad \mu(t, x) \geq \mu_I, \quad a(t, x) \leq a_S, \quad d(t, x) \leq d_S, \quad c(t, x) \leq c_S < 0$$

for all $x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$ the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

Hence, if

$$(3.16) \quad \lambda_I > \Lambda^1(-b_S \omega_{[\mu_I, d_S]}) \quad \text{and} \quad \mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]}),$$

then there are functions $\psi_{11} \in \text{int}(\mathcal{C}_1)$ and $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever

$$\begin{aligned} \lambda_I &\leq \lambda(t, x), & \mu_I &\leq \mu(t, x), & a(t, x) &\leq a_S, \\ b(t, x) &\leq b_S < 0, & c(t, x) &\leq c_S < 0, & d(t, x) &\leq d_S \end{aligned}$$

$x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), and for all $u_s > 0, v_s > 0$ in a fixed bounded set of $C(\overline{\Omega})$ bounded away from 0, the solution (u, v) of (3.1) for $t > s \geq t_0$ is nondegenerate at ∞ , and for all $t - s$ large enough

$$u(t, s; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly forwards permanent.

Proof. We proceed as in the Proposition 3.7 using now that, as $t - s \rightarrow \infty$,

$$u \geq \xi_{[\lambda - b\eta_{[\mu, d], a}]} \geq \xi_{[\lambda_I - b_S \eta_{[\mu_I, d_S], a_S}]} \rightarrow \omega_{[\lambda_I - b_S \omega_{[\mu_I, d_S], a_S}]}$$

and

$$v \geq \eta_{[\mu - c\xi_{[\lambda, a], d}]} \geq \eta_{[\mu_I - c_S \xi_{[\lambda_I, a_S], d_S}]} \rightarrow \omega_{[\mu_I - c_S \omega_{[\lambda_I, a_S], d_S}]} \quad \square$$

On the other hand, for pullback permanence and for complete nondegenerate solutions, we have the following proposition along the same lines as that above.

PROPOSITION 3.10 (pullback permanence—symbiotic case). Assume (3.2), $b_M, c_M < 0$, and (3.9), that is

$$b_L c_L < a_L d_L.$$

Then

(i) If $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_I, d_S]})$ there exists $\psi_{11} \in \text{int}(C_1)$ such that whenever

$$\lambda(t, x) \geq \lambda_I, \quad \mu(t, x) \geq \mu_I, \quad b(t, x) \leq b_S < 0, \quad a(t, x) \leq a_S, \quad d(t, x) \leq d_S.$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $u(t, x) \geq \psi_{11}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

(ii) If $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]})$ there exists $\psi_{22} \in \text{int}(C_2)$ such that whenever

$$\lambda(t, x) \geq \lambda_I, \quad \mu(t, x) \geq \mu_I, \quad a(t, x) \leq a_S, \quad d(t, x) \leq d_S, \quad c(t, x) \leq c_S < 0,$$

$x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$) for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $v(t, x) \geq \psi_{22}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

Hence, if

$$(3.17) \quad \lambda_I > \Lambda^1(-b_S \omega_{[\mu_I, d_S]}) \quad \text{and} \quad \mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]}),$$

there exist $\psi_{11} \in \text{int}(C_1)$ and $\psi_{22} \in \text{int}(C_2)$ such that whenever

$$\begin{aligned} \lambda_I &\leq \lambda(t, x), & \mu_I &\leq \mu(t, x), & a(t, x) &\leq a_S, \\ b(t, x) &\leq b_S < 0, & c(t, x) &\leq c_S < 0, & d(t, x) &\leq d_S, \end{aligned}$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$), and for all nontrivial $u_s \geq 0, v_s \geq 0$ in a fixed bounded set, B , of $C(\overline{\Omega})$ bounded away from 0, the set of solutions of (3.1) $\{(u, v), s < t \leq t_0, (u_s, v_s) \in B\}$ is nondegenerate as $s \rightarrow -\infty$ and for all $t - s$ large enough

$$u(t, s, x; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s, x; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly pullback permanent and any bounded complete trajectory that is nondegenerate at $-\infty$ satisfies

$$u(t, x) \geq \psi_{11}(x) \quad \text{and} \quad v(t, x) \geq \psi_{22}(x) \quad \text{for all } x \in \Omega \text{ and } t \leq t_0.$$

3.4.3. Prey-predator. We also have for the prey-predator case the following result.

PROPOSITION 3.11 (forwards permanence—prey-predator case). *Assume (3.2) and $b_L > 0$ and $c_M < 0$. Then:*

(i) *If $\lambda_I > \Lambda^1(-b_S\omega_{[\mu_S-c_I\omega_{[\lambda_S,a_I],d_I]}})$, there exists $\psi_{11} \in \text{int}(\mathcal{C}_1)$ such that whenever*

$$\begin{aligned} \lambda_S \geq \lambda(t, x) \geq \lambda_I, & \quad \mu(t, x) \leq \mu_S, & \quad a_S \geq a(t, x) \geq a_I > 0, \\ b(t, x) \leq b_S, & \quad c(t, x) \geq c_I, & \quad d(t, x) \geq d_I > 0 \end{aligned}$$

for all $x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$ the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

(ii) *If $\mu_I > \Lambda_0^2$, there exists $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever*

$$\mu(t, x) \geq \mu_I, \quad d(x, t) \leq d_S$$

for all $x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$ the solution for $t > s \geq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

Hence, if

$$(3.18) \quad \lambda_I > \Lambda^1(-b_S\omega_{[\mu_S-c_I\omega_{[\lambda_S,a_I],d_I]}}) \quad \text{and} \quad \mu_I > \Lambda_0^2,$$

there are functions $\psi_{11} \in \text{int}(\mathcal{C}_1)$ and $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever

$$\begin{aligned} \lambda_I \leq \lambda(t, x) \leq \lambda_S, & \quad \mu_I \leq \mu(t, x) \leq \mu_S, & \quad a_S \geq a(t, x) \geq a_I > 0, \\ 0 < b_I \leq b(t, x) \leq b_S, & \quad c_I \leq c(t, x) \leq c_S < 0, & \quad d_S \geq d(t, x) \geq d_I > 0 \end{aligned}$$

for all $x \in \Omega$ and $t \geq t_0$ (some $t_0 \in \mathbf{R}$), and for all $u_s > 0, v_s > 0$ in a fixed bounded set of $C(\bar{\Omega})$ bounded away from 0, the solution (u, v) of (3.1) for $t > s \geq t_0$ is nondegenerate at ∞ and for all $t - s$ large enough

$$u(t, s, x; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s, x; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly forwards permanent.

Proof. As before, we use now that as $t - s \rightarrow \infty$,

$$u \geq \xi_{[\lambda-b\eta_{[\mu-c\xi_{[\lambda,a],d}],a}]} \geq \xi_{[\lambda_I-b_S\eta_{[\mu_S-c_I\xi_{[\lambda_S,a_I],d_I],a_I}]}] \rightarrow \omega_{[\lambda_I-b_S\eta_{[\mu_S-c_I\omega_{[\lambda_S,a_I],d_I],a_I}]}}$$

and

$$v \geq \eta_{[\mu,d]} \geq \eta_{[\mu_I,d_S]} \rightarrow \omega_{[\mu_I,d_S]}. \quad \square$$

PROPOSITION 3.12 (pullback permanence—prey-predator case). *Assume (3.2) and $b_L > 0$ and $c_M < 0$. Then:*

(i) *If $\lambda_I > \Lambda^1(-b_S\omega_{[\mu_S-c_I\omega_{[\lambda_S,a_I],d_I]}})$, there exists $\psi_{11} \in \text{int}(\mathcal{C}_1)$ such that whenever*

$$\begin{aligned} \lambda_S \geq \lambda(t, x) \geq \lambda_I, & \quad \mu(t, x) \leq \mu_S, & \quad a_S \geq a(t, x) \geq a_I > 0, \\ b(t, x) \leq b_S, & \quad c(t, x) \geq c_I, & \quad d(t, x) \geq d_I > 0 \end{aligned}$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t \in \mathbf{R}$), for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{11}(x) \leq u(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $u(t, x) \geq \psi_{11}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

(ii) If $\mu_I > \Lambda_0^2$, there exists $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever

$$\mu(t, x) \geq \mu_I, \quad d(x, t) \leq d_S$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$), for any $u_s, v_s > 0$, the solution for $s < t \leq t_0$ of (3.1) satisfies $\psi_{22}(x) \leq v(t, s, x; u_s, v_s)$ for $t - s$ large enough.

In particular, any complete trajectory of (3.1) that is nondegenerate at $-\infty$ satisfies $v(t, x) \geq \psi_{22}(x)$ for all $x \in \Omega$ and $t \leq t_0$.

Hence, if

$$(3.19) \quad \lambda_I > \Lambda^1(-b_S \omega_{[\mu_S - c_I \omega_{[\lambda_S, a_I], d_I]}}) \quad \text{and} \quad \mu_I > \Lambda_0^2,$$

there exist functions $\psi_{11} \in \text{int}(\mathcal{C}_1)$ and $\psi_{22} \in \text{int}(\mathcal{C}_2)$ such that whenever

$$\begin{aligned} \lambda_I \leq \lambda(t, x) \leq \lambda_S, \quad \mu_I \leq \mu(t, x) \leq \mu_S, \quad a_S \geq a(t, x) \geq a_I > 0, \\ 0 < b_I \leq b(t, x) \leq b_S, \quad c_I \leq c(t, x) \leq c_S < 0, \quad d_S \geq d(t, x) \geq d_I > 0, \end{aligned}$$

for all $x \in \Omega$ and $t \leq t_0$ (some $t_0 \in \mathbf{R}$), and for all nontrivial $u_s \geq 0, v_s \geq 0$ in a fixed bounded set, B , of $C(\overline{\Omega})$ bounded away from 0, the set of solutions of (3.1) $\{(u, v), s < t \leq t_0, (u_s, v_s) \in B\}$ is nondegenerate as $s \rightarrow -\infty$ and for all $t - s$ large enough

$$u(t, s, x; u_s, v_s) \geq \psi_{11}(x) \quad \text{and} \quad v(t, s, x; u_s, v_s) \geq \psi_{22}(x).$$

In particular, (3.1) is uniformly pullback permanent and any bounded complete trajectory that is nondegenerate at $-\infty$ satisfies

$$u(t, x) \geq \psi_{11}(x) \quad \text{and} \quad v(t, x) \geq \psi_{22}(x) \quad \text{for all } x \in \Omega \text{ and } t \leq t_0.$$

Remark 3.13. Note that in order to apply the previous results one has to check that the assumptions in Propositions 3.7–3.12 are meaningful. Indeed, conditions (3.10), (3.16), and (3.18) must define nonempty sets of coefficients. Here we analyze only Dirichlet or Neumann boundary conditions; Robin ones can be treated in a similar way although the estimates are a little more involved.

In fact, (3.16) includes all coefficients such that

$$\lambda_I > \Lambda_0^1, \quad \mu_I > \Lambda_0^2$$

since in this case $\lambda_I > \Lambda_0^1 > \Lambda^1(-b_S \omega_{[\mu_I, d_S]})$ and $\mu_I > \Lambda_0^2$ then $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]})$; see also [12].

However, in order to show that (3.10) defines a nonempty set we must impose some conditions on b or c . If, for example, $b_S \rightarrow 0$, then $\Lambda^1(-b_S \omega_{[\mu_S, d_I]}) \rightarrow \Lambda_0^1$. Also, if $c_S \rightarrow 0$, then $\Lambda^2(-c_S \omega_{[\lambda_S, a_I]}) \rightarrow \Lambda_0^2$. Hence if b_S or c_S are small, the conditions in (3.10) can be met; see also [27] and [29].

We analyze condition (3.18) for the prey-predator case in more detail. From Proposition 3.1, in the case of Dirichlet or Neumann boundary conditions, we have $\omega_{[h, g]} \leq h_M/g_L$, and so

$$\omega_{[\lambda_S, a_I]} \leq \frac{\lambda_S}{a_I} \quad \text{and then} \quad \omega_{[\mu_S - c_I \omega_{[\lambda_S, a_I], d_I]}} \leq \frac{\mu_S - c_I \left(\frac{\lambda_S}{a_I}\right)}{d_I},$$

and then using the monotonicity of $\Lambda(m)$ with respect to m and (3.5), we get

$$\Lambda^1(-b_S \omega_{[\mu_S - c_I \omega_{[\lambda_S, a_I], d_I]}}) \leq \Lambda^1(-b_S \frac{(a_I \mu_S - c_I \lambda_S)}{a_I d_I}) = \Lambda_0^1 + b_S \frac{(a_I \mu_S - c_I \lambda_S)}{a_I d_I}.$$

Hence, if λ_I and μ_I satisfy

$$\lambda_I > \Lambda_0^1 + \frac{b_S \mu_S}{d_I} + \frac{-b_S c_I}{a_I d_I} \lambda_S, \quad \mu_I > \Lambda_0^2,$$

then (3.18) defines a nonempty set of parameters.

Observe that the first condition above is a restriction on the oscillation of $\lambda(t, x)$ as $t \rightarrow \pm\infty$.

In particular, if $\Lambda_0^1 + \frac{b_S \mu_S}{d_I} > 0$, then a necessary condition is

$$a_I d_I + b_S c_I > 0.$$

In such a case the conditions above can be met.

Now, for reasons that will be apparent in the following sections, we are interested in some uniformity in the previous results with respect to the coefficients $b_I \leq b_S$ and $c_I \leq c_S$. More precisely, we are going to show that the functions $\psi_{11}(x)$ and $\psi_{22}(x)$ in all the previous propositions can be taken independent of $b(t, x)$ and $c(t, x)$, provided that one of the numbers $b_I \leq b_S$ or $c_I \leq c_S$ is sufficiently small. In fact we have the following:

THEOREM 3.14. (i) *The competitive case: $b_L, c_L > 0$. Assume either*

1. $\lambda_I > \Lambda_0^1$, $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_S, a_I]})$, and b_S is sufficiently small, or
2. $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_S, d_I]})$, $\mu_I > \Lambda_0^2$, and c_S is sufficiently small, for $t - s$ large enough.

Then the functions $\psi_{11}(x)$ and $\psi_{22}(x)$ in Propositions 3.7 and 3.8 can be taken also independent of b_S and c_S .

(ii) *The symbiotic case: $b_M, c_M < 0$ and $b_L c_L < a_L d_L$. Assume either*

1. $\lambda_I > \Lambda_0^1$, $\mu_I > \Lambda^2(-c_S \omega_{[\lambda_I, a_S]})$, or
2. $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_I, d_S]})$, $\mu_I > \Lambda_0^2$.

Then the functions $\psi_{11}(x)$ and $\psi_{22}(x)$ in Propositions 3.9 and 3.10 can be taken also independent of b_I and c_I .

(iii) *The prey-predator case: $b_L > 0$, $c_M < 0$. Assume either*

1. $\lambda_I > \Lambda_0^1$, $\mu_I > \Lambda_0^2$, and b_S is sufficiently small, or
2. $\lambda_I > \Lambda^1(-b_S \omega_{[\mu_S, d_I]})$, $\mu_I > \Lambda_0^2$, and $|c_I|$ is sufficiently small.

Then the functions $\psi_{11}(x)$ and $\psi_{22}(x)$ in Propositions 3.11 and 3.12 can be taken also independent of b_S and c_I .

Proof. We analyze only the competitive case. By the proof of Proposition 3.7 and Theorem 3.3, statement 6 we get

$$\begin{aligned} u(t, s, u_s, v_s) &\geq \Theta_{[\lambda_I - b_S \eta_{[\mu_S, d_I], a_S}]}(t - s, u_s) \geq \\ &\geq \omega_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]} \geq m \omega_{[\lambda_I - b_S(\omega_{[\mu_S, d_I]} + \varepsilon), a_S]}. \end{aligned}$$

It suffices to apply Lemma 3.2 as $b_S \rightarrow 0$ where $m < 1$.

On the other hand,

$$v(t, s, u_s, v_s) \geq \Theta_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]}(t - s, v_s) \rightarrow \omega_{[\mu_I - c_S(\omega_{[\lambda_S, a_I]} + \varepsilon), d_S]},$$

and so taking ε small, the result follows.

The other cases can be studied in an analogous way by Propositions 3.9, 3.10, 3.11, and 3.12. \square

4. Exponential decay for nonautonomous linear systems. Once the results on permanence of the previous section have been established, we turn now our

attention to determining ranges of parameters such that there exist some special asymptotically stable trajectories describing the asymptotic behavior of solutions of (3.1), either forwards or in a pullback sense. For this we have to develop some tools for linear systems.

Hence, in this section we give sufficient conditions for certain linear systems to have exponential decay. The results are of a perturbative nature and are based upon results in [33] for scalar equations.

4.1. Preliminary results for the scalar case. We start by recalling some results for the following scalar equation:

$$(4.1) \quad \begin{cases} u_t - d\Delta u = c(t, x)u & x \in \Omega, t > s \\ \mathcal{B}u = 0, & x \in \partial\Omega, t > s \\ u(s) = u_s. \end{cases}$$

Assume that $d > 0$, $c \in C^\theta(\mathbf{R}, L^p(\Omega))$, with $0 < \theta \leq 1$ and some $p > \max(N/2, 1)$. Then for any $u_s \in X$, where $X = L^q(\Omega)$ with $1 \leq q < \infty$, or $X = C(\overline{\Omega})$, (4.1) has a unique solution given by $u(t, s; u_s)$, which is a strong solution in $L^r(\Omega)$ for any $1 \leq r < p$. This solution can be used to define an order-preserving evolution operator T_c in X via the definition $T_c(t, s)u_s = u(t, s; u_s)$.

Moreover for each q and r with $1 \leq q \leq r \leq \infty$ and $R_0 > 0$ there exist $L_0 = L_0(R_0, r, q) > 0$ and $\delta_0 = \delta(R_0, r, q) > 0$ such that the evolution operator $T_c(t, s)$ satisfies

$$(4.2) \quad \|T_c(t, s)u_0\|_{L^r(\Omega)} \leq L_0 \frac{e^{\delta_0(t-s)}}{(t-s)^{\frac{N}{2}(\frac{1}{q}-\frac{1}{r})}} \|u_0\|_{L^q(\Omega)}, \quad t > s$$

for every $c \in C^\theta(\mathbf{R}, L^p(\Omega))$, with $0 < \theta \leq 1$ and some $p > N/2$, such that

$$\|c\|_{L^\infty(\mathbf{R}, L^p(\Omega))} \leq R_0.$$

Also, the evolution operator smooths the solutions. More precisely, for every $u_0 \in L^q(\Omega)$ and $t > s$, the map

$$(s, \infty) \ni t \longmapsto u(t, s; u_0) := T_c(t, s)u_0 \in \begin{cases} C_B^\nu(\overline{\Omega}) & \text{if } p > N/2, \\ C_B^{1,\nu}(\overline{\Omega}) & \text{if } p > N, \end{cases}$$

is continuous for some $\nu > 0$. Here

$$C_B^{j,\nu}(\overline{\Omega}) = \begin{cases} C_0^{j,\nu}(\overline{\Omega}) & \text{for Dirichlet BCs,} \\ C^{j,\nu}(\overline{\Omega}) & \text{for Neumann or Robin BCs,} \end{cases}$$

see, e.g., Rodríguez-Bernal [32].

The following proposition is taken from Lemma 4.1 in Robinson et al. [31] and Lemma 2.1 in Rodríguez-Bernal [33]:

PROPOSITION 4.1. *Suppose that for some q with $1 \leq q \leq \infty$ there exist $M > 0$ and $\beta \in \mathbf{R}$ such that*

$$(4.3) \quad \|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M e^{\beta(t-s)} \quad \text{for all } t > s.$$

Then for any $1 \leq r \leq \infty$ there exists a $K \geq 1$ such that

$$(4.4) \quad \|T_c(t, s)\|_{\mathcal{L}(L^r(\Omega))} \leq K e^{\beta(t-s)} \quad \text{for all } t > s.$$

The constant K can be taken as a continuous function of β, M .

Moreover, for each r with $1 \leq r \leq q \leq \infty$ and for any $\varepsilon > 0$, we have

$$(4.5) \quad \|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega), L^r(\Omega))} \leq M(\beta, \varepsilon) \frac{e^{(\beta+\varepsilon)(t-s)}}{(t-s)^\delta}, \quad t > s,$$

where $\delta = \frac{N}{2} \left(\frac{1}{r} - \frac{1}{q} \right)$,

$$(4.6) \quad M(\beta, \varepsilon) = \kappa(\beta, M) \begin{cases} \left(\frac{\delta}{e}\right)^\delta \varepsilon^{-\delta} & \text{if } 0 < \varepsilon < \varepsilon_0 = \frac{\delta}{e} \\ 1 & \text{if } \varepsilon \geq \varepsilon_0 = \frac{\delta}{e} \end{cases}$$

and

$$\kappa(\beta, M) = L_0 \varepsilon^{\delta_0} \max\{1, M \varepsilon^{-\beta}\}.$$

Note that the constants K and κ in the proposition also depend on q and r but we will not pay attention to this dependence.

Our main argument, further below in the paper, will rely on results of the following type. We start with an evolution operator $T_c(t, s)$ that satisfies the estimate

$$\|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1 \quad \text{for } t \geq s \text{ and } M_1 > 0$$

for either $s \geq s_0$ or for $t \leq t_0$. Then, we add to $c(t, x)$ a perturbation $p(t, x)$ in the class $C^\theta(\mathbf{R}, L^p(\Omega))$, with $0 < \theta \leq 1$ and some $p > \max(N/2, 1)$, and we want to guarantee that the solutions of the new evolution operator $T_{c+p}(t, s)$ decay exponentially. This means that we want to get estimates of the type

$$(4.7) \quad \|T_{c+p}(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M'_1 e^{\beta'(t-s)} \quad \text{for all } t > s \text{ and some } \beta' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$. Note also that we can always assume, without loss of generality, that the $L^\infty(\mathbf{R}, L^p(\Omega))$ norms of both $c(t, x)$ and $p(t, x)$ are bounded by R_0 , so (4.2) holds for $T_c(t, s)$ and $T_{c+p}(t, s)$.

In this direction, the following important result is a particular case of Corollary 3.3 in Rodríguez-Bernal [33], and it provides sufficient conditions on $p(t, x)$ to ensure that (4.7) holds.

PROPOSITION 4.2. *Assume that*

$$(4.8) \quad \|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1 \quad \text{for } t \geq s \text{ and } M_1 > 0$$

and for either $s \geq s_0$ or for $t \leq t_0$.

Let $p \in C^\theta(\mathbf{R}, L^p(\Omega))$, for some $0 < \theta \leq 1$ and $p > \max(N/2, 1)$, and assume that for $|t|$ sufficiently large, we have $p(t, x) \leq -\varphi(x)$ where

$$\varphi \in C^1(\overline{\Omega}), \quad \varphi \geq 0, \quad \text{and } \nabla\varphi \neq 0 \quad \text{at the points at which } \varphi = 0.$$

Then

$$(4.9) \quad \|T_{c+p}(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M'_1 e^{\beta'(t-s)} \quad \text{for all } t > s \text{ and some } \beta' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$, with $M'_1 = M'_1(M_1, \varphi)$ and $\beta' = \beta'(M_1, \varphi)$.

The constants $M'_1 = M'_1(M_1, \varphi)$ and $\beta' = \beta'(M_1, \varphi)$ depend continuously on M_1 and on $\varphi \in C^1(\overline{\Omega})$.

Note that the condition above holds, in particular if $p(t, x) \leq -\delta < 0$ (in which case the constants M'_1 and β' can be chosen so that they depend continuously on δ), or if $\varphi \in C^1_0(\bar{\Omega})$ is positive in Ω and $\frac{\partial \varphi}{\partial n} < 0$ on $\partial\Omega$. The former is a common situation in the case of Neumann or Robin boundary conditions and the latter in the case of Dirichlet boundary conditions.

In order to apply the above result, we need to show first that (4.8) holds. The next result gives conditions for an evolution operator to have bounds of the type (4.8); see [34], [33]. For this recall the definitions of complete trajectory and of nondegeneracy in section 2.2, which we apply here to solutions of (4.1). Hence, according to [33], we have the following proposition.

PROPOSITION 4.3. (i) *If there exists a positive nondegenerate solution $u(t, s; u_s)$ of (4.1) defined for all $t > s \geq s_0$ such that for some $M > 0$ and some q with $1 \leq q \leq \infty$*

$$\|u(t, s; u_s)\|_{L^q(\Omega)} \leq M,$$

then

$$(4.10) \quad 0 < M_0 \leq \|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1 \quad \text{for } t \geq s \geq s_0,$$

where M_0, M_1 are independent of t and s and depend continuously on M and on $\varphi_0 \in C^1(\bar{\Omega})$.

(ii) *If there exists a positive, complete nondegenerate solution $u(t)$ of (4.1) that is bounded as $t \rightarrow -\infty$, i.e.,*

$$\|u(t)\|_{L^q(\Omega)} \leq M \quad \text{for } t \leq t_0,$$

then

$$(4.11) \quad 0 < M_0 \leq \|T_c(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1 \quad \text{for } s \leq t \leq t_0,$$

where M_0, M_1 are independent of t and s and depend continuously on M and $\varphi_0 \in C^1(\bar{\Omega})$.

4.2. Perturbation and decay of linear systems. In this section we generalize the perturbation result in the previous section to the case of a system of linear equations. The main theorem in this section will be crucial in the analysis of Lotka–Volterra models in the following sections.

Consider the linear coupled nonautonomous system

$$(4.12) \quad \begin{cases} u_t - d_1 \Delta u = a_{11}(t, x)u + a_{12}(t, x)v, & x \in \Omega, t > s \\ v_t - d_2 \Delta v = a_{21}(t, x)u + a_{22}(t, x)v, & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 \\ u(s) = u_s, v(s) = v_s, \end{cases}$$

in $L^q(\Omega, \mathbf{R}^2) \doteq [L^q(\Omega)]^2$. Then define

$$D = \text{diag}(d_1, d_2) \quad \text{and} \quad A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

and note that setting $U = \begin{pmatrix} u \\ v \end{pmatrix}$, (4.12) can be written as

$$U_t - D \Delta U = A(t, x)U$$

with boundary conditions $\mathcal{B}U = \begin{pmatrix} \mathcal{B}_1 u \\ \mathcal{B}_2 v \end{pmatrix} = 0$ on the boundary of Ω .

If $A \in C^\theta(\mathbf{R}, L^p(\Omega, \mathbf{R}^4))$, with $0 < \theta \leq 1$, $p > N/2$, and $p > q \geq 1$, the existence of a unique solution $U(t, s; U_s)$ of (4.12), in $L^q(\Omega, \mathbf{R}^2)$, can be obtained from Theorems 11.2, 11.3, and 11.4 in Amann [1]. Thus, the time-dependent operator $-D\Delta - A(t, x)$ generates an evolution operator, $T_A(t, s)$, in $L^q(\Omega, \mathbf{R}^2)$ (Theorem 4.4.1 in Amann [2]) via the definition $T_A(t, s)U_s = U(t, s; U_s)$.

The following result, analogous to (4.2), can be proved along the lines of the scalar arguments in Rodríguez-Bernal [33], [32], and Robinson et al. [31].

PROPOSITION 4.4. *For any $1 \leq q \leq r \leq \infty$, and $R_0 > 0$ there exist $L_0 = L_0(R_0, r, q) > 0$ and $\delta_0 = \delta_0(R_0, r, q) > 0$ such that the evolution operator $T_A(t, s)$ satisfies*

$$(4.13) \quad \|T_A(t, s)U_s\|_{L^r(\Omega, \mathbf{R}^2)} \leq L_0 \frac{e^{\delta_0(t-s)}}{(t-s)^{\frac{N}{2}(\frac{1}{q}-\frac{1}{r})}} \|U_s\|_{L^q(\Omega, \mathbf{R}^2)},$$

for every $\|A\|_{L^\infty(\mathbf{R}, L^p(\Omega, \mathbf{R}^4))} \leq R_0$. In particular, $T_A(t, s)$ extends to an evolution operator in $L^q(\Omega, \mathbf{R}^2)$ for every $1 \leq q < \infty$.

Furthermore, the results of Proposition 4.1 for the scalar case remain true for system (4.12).

Along the same lines as for scalar equations, we consider the linear uncoupled system

$$(4.14) \quad \begin{cases} u_t - d_1\Delta u = q_{11}(t, x)u, & x \in \Omega, t > s \\ v_t - d_2\Delta v = q_{22}(t, x)v, & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 \\ u(s) = u_s, v(s) = v_s. \end{cases}$$

Observe that with the notations above and setting

$$Q = \text{diag}(q_{11}, q_{22}),$$

then the evolution operator $T_Q(t, s)$ is well defined in $L^q(\Omega, \mathbf{R}^2)$, $1 \leq q < \infty$.

Now we assume that each separate equation in (4.14) satisfies

$$\|T_{q_{ii}}(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1, \quad t > s,$$

with M_1 independent of t and s and for either $t \leq t_0$ or $s \geq s_0$. Therefore the evolution operator $T_Q(t, s)$ satisfies (4.8).

Our goal is to give conditions on the coupling perturbations such that the solutions of the perturbed system

$$\begin{cases} u_t - d_1\Delta u = q_{11}(t, x)u + p_{11}(t, x)u + p_{12}(t, x)v, & x \in \Omega, t > s \\ v_t - d_2\Delta v = q_{22}(t, x)v + p_{21}(t, x)u + p_{22}(t, x)v, & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 \\ u(s) = u_s, v(s) = v_s, \end{cases}$$

decay exponentially. Note that the perturbed system can be written as

$$(4.15) \quad U_t - D\Delta U = Q(t, x)U + P(t, x)U$$

with

$$Q = \text{diag}(q_{11}, q_{22}), \quad P = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}, \quad \text{and} \quad U = \begin{pmatrix} u \\ v \end{pmatrix},$$

with $Q, P \in C^\theta(\mathbf{R}, L^p(\Omega, \mathbf{R}^4))$, with $0 < \theta \leq 1$, $p > \max(N/2, 1)$.

Hence our goal is to obtain an estimate of the type

$$(4.16) \quad \|T_{Q+P}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M'_1 e^{\beta'(t-s)} \quad \text{for all } t > s \text{ and some } \beta' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$.

Note that again we will assume, without loss of generality, that all the evolution operators considered satisfy (4.13) with the same constants L_0 and δ_0 .

In what follows we will make use of the following singular Gronwall lemma (see Henry [16]):

LEMMA 4.5 (a singular Gronwall lemma). *Assume that $a \in L^\infty(\tau_0, \infty)$ with $\tau_0 \geq -\infty$ and that $z(t) \geq 0$ is a locally bounded function that for $t \geq s > \tau_0$ satisfies*

$$(4.17) \quad z(t) \leq A + \int_s^t \frac{a(\tau)}{(t-\tau)^\delta} z(\tau) \, d\tau$$

with $\delta < 1$. Then we have for $t \geq s > \tau_0$

$$0 \leq z(t) \leq A(\delta)e^{\gamma(t-s)}$$

with $\gamma = \gamma(a, s, \delta) = (\|a\|_{L^\infty(s, \infty)} \Gamma(1-\delta))^{1/(1-\delta)}$ and $A(\delta)$ depends only on the constants A and δ but not on the function $a(\cdot)$ or on s, γ , or τ_0 .

Our next result states that if the diagonal perturbing terms $p_{ii}(t, x)$ are sufficiently strong and the coupling terms $p_{ij}(t, x)$, $i \neq j$, are “small” at $\pm\infty$, then (4.16) is achieved.

THEOREM 4.6. *With the notations in (4.15), assume that the scalar evolution operators $T_{q_{ii}}(t, s)$ satisfy*

$$(4.18) \quad \|T_{q_{ii}}(t, s)\|_{\mathcal{L}(L^q(\Omega))} \leq M_1, \quad t > s,$$

with M_1 independent of t and s and for either $t \leq t_0$ or $s \geq s_0$.

Assume also that $p_{ii}(t, x)$ satisfies $p_{ii}(t, x) \leq -\varphi_{ii}(x)$ with $\varphi_{ii}(x)$ as in Proposition 4.2.

Then there exists a $\rho = \rho(M_1, \varphi_{11}, \varphi_{22}) > 0$ such that if

$$(4.19) \quad \limsup_{|t| \rightarrow \infty} \|p_{12}(t)\|_{L^p(\Omega)} \limsup_{|t| \rightarrow \infty} \|p_{21}(t)\|_{L^p(\Omega)} \leq \rho^2,$$

then

$$(4.20) \quad \|T_{Q+P}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M''_1 e^{\beta''(t-s)} \quad \text{for all } t > s \text{ and some } \beta'' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$, where $M''_1 = M''_1(M_1, \varphi_{11}, \varphi_{22})$ and $\beta'' = \beta''(M_1, \varphi_{11}, \varphi_{22})$.

The constants ρ , M''_1 , and β'' depend continuously on M_1 and $\varphi_{11}, \varphi_{22}$ as in Proposition 4.2.

Proof. Note that, using Proposition 4.4, we just need to prove the result for some suitably chosen $1 \leq q < \infty$. We proceed in several steps.

Step 1. If we define

$$P_1 = \begin{pmatrix} p_{11} & 0 \\ 0 & p_{22} \end{pmatrix},$$

then Proposition 4.2 applied to each separate equation gives the estimate

$$(4.21) \quad \|T_{Q+P_1}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M'_1 e^{\beta'(t-s)} \quad \text{for all } t > s \text{ and some } \beta' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$, with $M'_1 = M'_1(M_1, \varphi_{11}, \varphi_{22})$ and $\beta' = \beta'(M_1, \varphi_{11}, \varphi_{22})$.

Step 2. We will show that there exists a $\rho = \rho(M'_1, \beta')$, which depends continuously on M'_1, β' , such that if

$$\|p_{12}\|_{L^\infty(\mathbf{R}, L^p(\Omega))} \leq \rho \quad \text{and} \quad \|p_{21}\|_{L^\infty(\mathbf{R}, L^p(\Omega))} \leq \rho,$$

then

$$(4.22) \quad \|T_{Q+P}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M''_1 e^{\beta''(t-s)} \quad \text{for all } t > s \text{ and some } \beta'' < 0$$

and for either $s \geq s_0$ or for $t \leq t_0$, with

$$P = P_1 + P_2, \quad P_2 = \begin{pmatrix} 0 & p_{12} \\ p_{21} & 0 \end{pmatrix},$$

where $M''_1 = M''_1(M'_1, \beta', \rho)$ and $\beta'' = \beta''(M'_1, \beta', \rho)$, depend continuously on M'_1, β', ρ .

In fact, we have, by the variation of constants formula, that for every $U_0 \in L^q(\Omega, \mathbf{R}^2)$ the solution $U(t, s; U_0) = T_{Q+P}(t, s)U_0$ of (4.15) satisfies for $t \geq s$,

$$U(t, s; U_0) = T_{Q+P_1}(t, s)U_0 + \int_s^t T_{Q+P_1}(t, \tau)P_2(\tau)U(\tau, s; U_0) \, d\tau.$$

Now we choose q such that $p \geq q'$, so that $1/p + 1/q \leq 1$. In what follows we will apply (4.5) with $\frac{1}{r} = \frac{1}{p} + \frac{1}{q}$, and so with $\delta = N/2p$. With this choice, we have (4.21) and from (4.5)

$$\|T_{Q+P_1}(t, s)\|_{\mathcal{L}(L^r(\Omega, \mathbf{R}^2), L^q(\Omega, \mathbf{R}^2))} \leq M(\beta', \varepsilon) \frac{e^{(\beta'+\varepsilon)(t-s)}}{(t-s)^{\frac{N}{2p}}},$$

where $M(\beta', \varepsilon)$ is as in (4.6).

Since $P_2(\tau) \in L^p(\Omega, \mathbf{R}^2)$ and $U(\tau, s; U_0) \in L^q(\Omega, \mathbf{R}^2)$, then the term $P_2(\tau)U(\tau, s; U_0)$ can be estimated, using Hölder's inequality, in $L^r(\Omega, \mathbf{R}^2)$ with $\frac{1}{r} = \frac{1}{p} + \frac{1}{q}$. Thus,

$$\begin{aligned} \|U(t, s; U_0)\|_{L^q(\Omega, \mathbf{R}^2)} &\leq M'_1 e^{(\beta'+\varepsilon)(t-s)} \|U_0\|_{L^q(\Omega, \mathbf{R}^2)} \\ &+ M(\beta', \varepsilon) \int_s^t \frac{e^{(\beta'+\varepsilon)(t-\tau)}}{(t-\tau)^{\frac{N}{2p}}} \|P_2(\tau)\|_{L^p(\Omega, \mathbf{R}^2)} \|U(\tau, s; U_0)\|_{L^q(\Omega, \mathbf{R}^2)} \, d\tau. \end{aligned}$$

Then, multiplying by $e^{-(\beta'+\varepsilon)(t-s)}$, and setting $A = M'_1 \|U_0\|_{L^q(\Omega, \mathbf{R}^2)}$,

$$z(t) = e^{-(\beta'+\varepsilon)(t-s)} \|U(t, s; U_0)\|_{L^q(\Omega, \mathbf{R}^2)}, \quad \text{and} \quad a(\tau) = M(\beta', \varepsilon) \|P_2(\tau)\|_{L^p(\Omega, \mathbf{R}^2)}$$

we get, for all $t \geq s$,

$$z(t) \leq A + \int_s^t \frac{a(\tau)}{(\tau-s)^{\frac{N}{2p}}} z(\tau) \, d\tau.$$

We can apply the singular Gronwall lemma above with $\delta = \frac{N}{2p} < 1$, and we get

$$(4.23) \quad \|U(t, s; U_0)\|_{L^q(\Omega, \mathbf{R}^2)} \leq M''_1 e^{(\beta'+\mu(\varepsilon))(t-s)} \|U_0\|_{L^q(\Omega, \mathbf{R}^2)}, \quad t \geq s,$$

where

$$\mu(\varepsilon) = \varepsilon + \left(M(\beta', \varepsilon) \Gamma(1 - \delta) \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))} \right)^{\frac{1}{1-\delta}}.$$

Recalling (4.6), we get that

$$\mu(\varepsilon) = \begin{cases} \varepsilon + \varepsilon^{\frac{-\delta}{1-\delta}} A_0 \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}^{\frac{1}{1-\delta}} & \text{if } 0 < \varepsilon < \varepsilon_0 = \frac{\delta}{e} \\ \varepsilon + A_1 \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}^{\frac{1}{1-\delta}} & \text{if } \varepsilon \geq \varepsilon_0, \end{cases}$$

where

$$A_1 = \left(L_0 e^{\delta_0} \max\{1, M'_1 e^{-\beta'}\} \Gamma(1 - \delta) \right)^{1/(1-\delta)}, \quad A_0 = A_1 \left(\frac{\delta}{e} \right)^{\delta/(1-\delta)},$$

and L_0 and δ_0 are the constants in (4.2).

Thus $\mu(0) = \mu(\infty) = \infty$. But the function

$$h(\varepsilon) = \varepsilon + \varepsilon^{\frac{-\delta}{1-\delta}} A_0 \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}^{\frac{1}{1-\delta}}$$

has a unique minimum at

$$\varepsilon_1 = \left(A_0 \frac{\delta}{1-\delta} \right)^{1-\delta} \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))},$$

and

$$h(\varepsilon_1) = \frac{1}{\delta^\delta} \left(\frac{A_0}{1-\delta} \right)^{1-\delta} \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}.$$

Therefore, comparing ε_0 and ε_1 , and minimizing $\mu(\varepsilon)$ leads to

$$\|U(t, s; U_0)\|_{L^q(\Omega, \mathbf{R}^2)} \leq M_1'' e^{\beta''(t-s)} \|U_0\|_{L^q(\Omega, \mathbf{R}^2)}, \quad t \geq s$$

with

$$\begin{aligned} \beta'' &= \beta' + \min_{\{\varepsilon > 0\}} \mu(\varepsilon) \\ &= \beta' + \begin{cases} c_0 \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}, & \text{if } \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))} \leq s^*, \\ c_1 + c_2 \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))}^{\frac{1}{1-\delta}}, & \text{if } \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))} \geq s^*, \end{cases} \end{aligned}$$

where

$$c_0 = \frac{1}{\delta^\delta} \left(\frac{A_0}{1-\delta} \right)^{1-\delta}, \quad c_1 = \frac{\delta}{e}, \quad c_2 = A_1, \quad s^* = \frac{\delta}{e} \left(\frac{1-\delta}{A_0 \delta} \right)^{1-\delta}.$$

Thus, it is then clear that (4.22) follows, i.e., $\beta'' < 0$, provided that

$$\|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))} < \min \left\{ s^*, \frac{-\beta'}{c_0} \right\},$$

which reads

$$(4.24) \quad \|P_2\|_{L^\infty((s, \infty), L^p(\Omega, \mathbf{R}^2))} < \rho := \delta \left(\frac{1-\delta}{\delta A_0} \right)^{1-\delta} \min \left\{ -\beta', \frac{1}{e} \right\}.$$

Step 3. Now we show that the result in Step 2 above can be obtained only in terms of $\limsup_{|t| \rightarrow \infty} \|P_2(t)\|_{L^p(\Omega, \mathbf{R}^2)}$.

In fact, note that from (4.24), if we take $s \geq s_0$ sufficiently large, the conclusion with $\limsup_{t \rightarrow \infty} \|P_2(t)\|_{L^p(\Omega, \mathbf{R}^2)}$ is clear.

On the other hand, observe that we can set $P_2 = 0$ for $t \geq t_0$ and we still have (4.23) for $s \leq t \leq t_0$. Taking then t_0 very negative, (4.24) gives the result for $\limsup_{t \rightarrow -\infty} \|P_2(t)\|_{L^p(\Omega, \mathbf{R}^2)}$.

In particular, (4.22) follows, provided that

$$(4.25) \quad \limsup_{|t| \rightarrow \infty} \|P_2(t)\|_{L^p(\Omega, \mathbf{R}^2)} < \rho,$$

with ρ as in (4.24).

Step 4. The change of variables

$$U = \begin{pmatrix} u \\ v \end{pmatrix} \rightarrow V = \begin{pmatrix} \alpha u \\ \beta v \end{pmatrix}$$

with $\alpha, \beta > 0$, transforms the system (4.15) into

$$V_t - D\Delta V = Q(t, x)V + \tilde{P}(t, x)V$$

with

$$D = \text{diag}(d_1, d_2), \quad Q = \text{diag}(q_{11}, q_{22}), \quad \tilde{P} = \begin{pmatrix} p_{11} & \frac{\alpha}{\beta}p_{12} \\ \frac{\beta}{\alpha}p_{21} & p_{22} \end{pmatrix}.$$

Hence, we can apply Step 3 provided

$$\frac{\alpha}{\beta} \limsup_{|t| \rightarrow \infty} \|p_{12}(t)\|_{L^p(\Omega)} \leq \rho \quad \text{and} \quad \frac{\beta}{\alpha} \limsup_{|t| \rightarrow \infty} \|p_{21}(t)\|_{L^p(\Omega)} \leq \rho,$$

with $\rho > 0$ as in (4.24). We can choose α, β such that the above inequalities are satisfied if

$$\limsup_{|t| \rightarrow \infty} \|p_{12}(t)\|_{L^p(\Omega)} \limsup_{|t| \rightarrow \infty} \|p_{21}(t)\|_{L^p(\Omega)} \leq \rho^2$$

with $\rho > 0$ as in (4.24). □

Remark 4.7. Note that (4.24) gives a quantitative threshold for the size of the perturbation. In fact, from (4.24) and the expression of A_0 , it can be deduced that

$$\rho = \rho(M'_1, \beta') = \frac{e^\delta(1 - \delta)^{1-\delta}}{\Gamma(1 - \delta)} \frac{\min\{-\beta', \frac{1}{e}\}}{L_0 e^{\delta_0} \max\{1, M'_1 e^{-\beta'}\}},$$

where M'_1, β' are from Step 1.

Observe that Step 4 above is the only place where we used the fact that the system has only two components.

5. Attracting trajectories for general nonautonomous nonlinear systems. In this section we sketch out our approach to the existence of asymptotically stable complete trajectories for Lotka–Volterra systems. The key point is to write the equation satisfied by the difference of two solutions as a perturbation of an associated

linear system. Using then the permanence results in section 3, we can apply Theorem 4.6 to conclude that the difference of two solutions converges to zero as $t \rightarrow \infty$. A similar convergence result as the initial time $s \rightarrow -\infty$ will imply the uniqueness of complete nondegenerate solutions, which moreover describes the pullback behavior of the system.

First we treat the case of general nonautonomous nonlinear systems, before specializing to Lotka–Volterra models. Consider the general nonautonomous nonlinear system

$$(5.1) \quad \begin{cases} u_t - d_1 \Delta u = uf(t, x, u, v) & x \in \Omega, t > s \\ v_t - d_2 \Delta v = vg(t, x, u, v) & x \in \Omega, t > s \\ \mathcal{B}_1 u = 0, \mathcal{B}_2 v = 0 & x \in \partial\Omega, t > s \\ u(s) = u_s, v(s) = v_s. \end{cases}$$

We now sketch our strategy for analyzing the asymptotic behavior of solutions to (5.1). Consider two different pairs of nonnegative initial conditions (u_s^1, v_s^1) and (u_s^2, v_s^2) and consider the corresponding solutions of (5.1), $U_1 = \begin{pmatrix} u_1 \\ v_1 \end{pmatrix}$ and $U_2 = \begin{pmatrix} u_2 \\ v_2 \end{pmatrix}$, respectively. Write $y = u_2 - u_1$ and $z = v_2 - v_1$. Then, (y, z) satisfies

$$(5.2) \quad \begin{cases} y_t - d_1 \Delta y = q_{11}(t, x)y + p_{11}(t, x)y + p_{12}(t, x)z & x \in \Omega, t > s \\ z_t - d_2 \Delta z = q_{22}(t, x)z + p_{21}(t, x)y + p_{22}(t, x)z & x \in \Omega, t > s \\ \mathcal{B}_1 y = 0, \mathcal{B}_2 z = 0 & x \in \partial\Omega, t > s \\ y(s) = y_s, z(s) = z_s, \end{cases}$$

with $y_s = u_s^2 - u_s^1$, $z_s = v_s^2 - v_s^1$, and

$$\begin{aligned} q_{11}(t, x) &= f(t, x, u_2, v_2), & q_{22}(t, x) &= g(t, x, u_2, v_2) \\ p_{11}(t, x) &= u_1 \frac{f(t, x, u_2, v_1) - f(t, x, u_1, v_1)}{u_2 - u_1}, \\ p_{12}(t, x) &= u_1 \frac{f(t, x, u_2, v_2) - f(t, x, u_2, v_1)}{v_2 - v_1}, \\ p_{21}(t, x) &= v_1 \frac{g(t, x, u_2, v_1) - g(t, x, u_1, v_1)}{u_2 - u_1}, \\ p_{22}(t, x) &= v_1 \frac{g(t, x, u_2, v_2) - g(t, x, u_2, v_1)}{v_2 - v_1}. \end{aligned}$$

Most of the analysis that follows in the next section will be based on proving that the following results can be applied. The first one gives sufficient conditions to guarantee that two solutions have the same forwards asymptotic behavior, while the second gives a criterion to prove the coincidence of two complete trajectories and also describes the pullback behavior of solutions.

THEOREM 5.1 (forwards behavior). *Assume that both solutions of (5.1), $U_1 = \begin{pmatrix} u_1 \\ v_1 \end{pmatrix}$ and $U_2 = \begin{pmatrix} u_2 \\ v_2 \end{pmatrix}$, are globally defined and bounded in $L^\infty(\Omega, \mathbf{R}^2)$ for $t > s > t_0$. Moreover, suppose that u_1, v_1 are positive in Ω and $U_2(t)$ is positive, nondegenerate for $t > t_0$ and for some $p > \max(N/2, 1)$ and $0 < \theta \leq 1$, the coefficients in (5.2) satisfy $p_{ij}, q_{ii} \in C^\theta(\mathbf{R}, L^p(\Omega))$ for $i, j = 1, 2$ and $p_{ii}(t, x) \leq -\varphi_{ii}(x)$, for $t > t_0$, with $\varphi_{ii}(x)$ as in Proposition 4.2.*

Then there exists a $\rho > 0$ such that if

$$(5.3) \quad \limsup_{t \rightarrow \infty} \|p_{12}(t)\|_{L^p(\Omega)} \limsup_{t \rightarrow \infty} \|p_{21}(t)\|_{L^p(\Omega)} \leq \rho^2,$$

both solutions have the same forwards asymptotic behavior, i.e.,

$$U_1(t) - U_2(t) \rightarrow 0 \quad \text{exponentially in } C_{B_1}^1(\overline{\Omega}) \times C_{B_2}^1(\overline{\Omega}) \quad \text{as } t \rightarrow \infty.$$

In particular, $U_1(t)$ is also nondegenerate at $+\infty$.

Proof. Clearly, (5.2) can be written as

$$W_t - D\Delta W = QW + PW,$$

where

$$(5.4) \quad D = \text{diag}(d_1, d_2), \quad Q = \text{diag}(q_{11}, q_{22}), \quad P = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}, \quad W = \begin{pmatrix} y \\ z \end{pmatrix}.$$

Since $U_2 = \begin{pmatrix} u_2 \\ v_2 \end{pmatrix}$ is a positive, bounded, and nondegenerate solution, for $t > s > t_0$, of the diagonal system

$$W_t - D\Delta W = QW,$$

it follows from Propositions 4.3 and 4.4 that for any $1 \leq q < \infty$,

$$(5.5) \quad \|T_Q(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M_1, \quad t > s > t_0,$$

with M_1 independent of t and s , $t > s > t_0$.

Then, we apply Theorem 4.6 to obtain that there exists $\rho > 0$ such that if (5.3) holds, then

$$\|T_{Q+P}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M_1'' e^{\beta''(t-s)} \quad \text{for all } t > s > t_0 \quad \text{and some } \beta'' < 0.$$

Thus, from Proposition 4.4 (see also Proposition 4.1), we have, writing $W_s = (y_s, z_s)$ and for $t > s > t_0$,

$$(5.6) \quad \|W(t, s; W_s)\|_{L^\infty(\Omega, \mathbf{R}^2)} \leq M_2 e^{\beta''(t-s)} \|W_s\|_{L^\infty(\Omega, \mathbf{R}^2)} \rightarrow 0, \quad t \rightarrow \infty.$$

The uniform forwards convergence of trajectories follows. Standard parabolic regularization implies the $C_{B_1}^1(\overline{\Omega}) \times C_{B_2}^1(\overline{\Omega})$ convergence. \square

In the following result we use similar arguments to prove the coincidence of complete nondegenerate trajectories, and show that such a trajectory, when it exists, attracts (in the pullback sense) all bounded positive trajectories. In particular, the following results guarantee the uniqueness of complete nondegenerate solutions.

THEOREM 5.2 (coincidence of complete trajectories and pullback behavior). *Assume that $U_1 = \begin{pmatrix} u_1 \\ v_1 \end{pmatrix}$ is a complete trajectory that is bounded in $L^\infty(\Omega, \mathbf{R}^2)$ at $-\infty$, and nondegenerate for $t \leq t_0$. Suppose further that for some $p > \max(N/2, 1)$ and $0 < \theta \leq 1$, the coefficients in (5.2) satisfy $p_{ij}, q_{ii} \in C^\theta(\mathbf{R}, L^p(\Omega))$ for $i, j = 1, 2$ and $p_{ii}(t, x) \leq -\varphi_{ii}(x)$, for $t \leq t_0$, with $\varphi_{ii}(x)$ as in Proposition 4.2.*

Then there exists a $\rho > 0$ such that if

$$(5.7) \quad \limsup_{t \rightarrow -\infty} \|p_{12}(t)\|_{L^p(\Omega)} \limsup_{t \rightarrow -\infty} \|p_{21}(t)\|_{L^p(\Omega)} \leq \rho^2,$$

then

- (i) $U_1(t)$ is the unique complete trajectory that is bounded in $L^\infty(\Omega, \mathbf{R}^2)$ at $-\infty$, and
- (ii) if $U_2(s)$ is a family of positive initial data which is bounded in $L^\infty(\Omega, \mathbf{R}^2)$ as $s \rightarrow -\infty$, then $U_1(\cdot)$ pullback attracts $S(t, s)U_2(s)$, i.e., for any $t \in \mathbf{R}$

$$S(t, s)U_2(s) - U_1(t) \rightarrow 0 \quad \text{in } C_{B_1}^1(\overline{\Omega}) \times C_{B_2}^1(\overline{\Omega}) \quad \text{as } s \rightarrow -\infty.$$

Proof.

- (i) Let $U_2(t)$ be a complete trajectory bounded in $L^\infty(\Omega, \mathbf{R}^2)$ at $-\infty$. We write (5.2) as

$$W_t - D\Delta W = QW + PW, \quad W(s) = W_s = U_2(s) - U_1(s),$$

where Q, P , and W are defined as in (5.4). Since $U_1 = (u_1, v_1)$ is a complete, positive, bounded, and nondegenerate solution of the diagonal system

$$W_t - D\Delta W = QW,$$

it follows from Proposition 4.3 that for any $1 \leq q < \infty$, and sufficiently negative t_0 ,

$$(5.8) \quad \|T_Q(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M_1, \quad s < t \leq t_0,$$

with M_1 independent of t and s .

Then, we apply Theorem 4.6 to obtain that there exists $\rho > 0$ such that if (5.7) holds, then

$$\|T_{Q+P}(t, s)\|_{\mathcal{L}(L^q(\Omega, \mathbf{R}^2))} \leq M_1'' e^{\beta''(t-s)} \quad \text{for all } s < t \leq t_0 \quad \text{and some } \beta'' < 0.$$

Thus,

$$(5.9) \quad \|U_1(t) - U_2(t)\|_{L^q(\Omega, \mathbf{R}^2)} = \|W(t, s; W_s)\|_{L^q(\Omega, \mathbf{R}^2)} \leq M_1'' e^{\beta''(t-s)} \|W_s\|_{L^q(\Omega, \mathbf{R}^2)}.$$

The right-hand side tends to zero as $s \rightarrow -\infty$ since both complete trajectories are bounded, and the result follows.

- (ii) Proceeding as above, we obtain

$$\|U_1(t) - S(t, s)U_2(s)\|_{L^q(\Omega, \mathbf{R}^2)} \leq M_1'' e^{\beta''(t-s)} \|U_1(s) - U_2(s)\|_{L^q(\Omega, \mathbf{R}^2)}$$

for $s < t \leq t_0$ and some $\beta'' < 0$.

Thus, from Proposition 4.4 (see also Proposition 4.1), we get for $s < t \leq t_0$,

$$\|U_1(t) - S(t, s)U_2(s)\|_{L^\infty(\Omega, \mathbf{R}^2)} \leq M_2 e^{\beta''(t-s)} \|U_1(s) - U_2(s)\|_{L^q(\Omega, \mathbf{R}^2)} \rightarrow 0$$

as $s \rightarrow -\infty$. Standard parabolic regularization implies the convergence in $C_{B_1}^1(\overline{\Omega}) \times C_{B_2}^1(\overline{\Omega})$.

Now for every $\tau \geq t_0$, using the continuity of the nonlinear evolution process, we get, as $s \rightarrow -\infty$,

$$U_1(\tau, s) = S(\tau, t)U_1(t, s) \rightarrow S(\tau, t)U_2(t). \quad \square$$

The theorems above may perhaps appear more general than they really are. To verify the assumptions involved one must restrict the nonlinearities of the system and carefully choose the classes of solutions being considered. For example, the conditions $p_{ii}(t, x) \leq -\varphi_{ii}(x)$ and the smallness conditions on $p_{ij}(t, x)$, $i \neq j$ depend on the particular solutions considered.

Nevertheless, in the next section we will show that the assumptions required can be verified for our example of a general nonautonomous Lotka–Volterra system.

6. Attracting trajectories for nonautonomous Lotka–Volterra systems.

As (5.1) is far too general to apply Theorems 5.1 and 5.2 in a straightforward manner, in this section we apply these results to the solutions of (3.1). Note that we handle the three cases, competition, symbiosis, and prey-predator, in a unified way.

Then, for the difference of two solutions the coefficients in (5.2) are given by

$$(6.1) \quad \begin{aligned} q_{11}(t, x) &= \lambda(t, x) - a(t, x)u_2 - b(t, x)v_2, & q_{22}(t, x) &= \mu(t, x) - c(t, x)u_2 - d(t, x)v_2, \\ p_{11}(t, x) &= -a(t, x)u_1, & p_{12}(t, x) &= -b(t, x)u_1, \\ p_{21}(t, x) &= -c(t, x)v_1, & p_{22}(t, x) &= -d(t, x)v_1. \end{aligned}$$

Hence, to apply Theorems 5.1 or 5.2, since $a_L, d_L > 0$ and $u_1, v_1 \geq 0$, in order to find positive functions $\varphi_{ii}(x)$ such that $p_{ii}(t, x) \leq -\varphi_{ii}(x)$, $i = 1, 2$ we need positive functions $\psi_{ii}(x)$ such that $\psi_{11}(x) \leq u_1(t, x)$ and $\psi_{22}(x) \leq v_1(t, x)$, that is, we must consider nondegenerate solutions. The results in section 3.4 guarantee then that all solutions are nondegenerate.

On the other hand we must show that the product of the coupling terms

$$p_{12}(t, x)p_{21}(t, x)$$

is small at $\pm\infty$. Having obtained bounds on u_1, v_1 this will be achieved by a smallness condition on the coefficients $b(t, x)$ or $c(t, x)$.

But note that the nondegeneracy of solutions above depends on the functions $b(t, x)$ and $c(t, x)$ themselves. Therefore, we will use the results in section 3.4 which guarantee that solutions of (3.1) are nondegenerate for all sufficiently small “coupling” coefficients $b(t, x)$ or $c(t, x)$ and that the functions $\psi_{ii}(x)$, $i = 1, 2$ do not converge to zero as b or c vanish.

We first start with the forwards behavior in Theorem 5.1. Then we can prove the following theorem.

THEOREM 6.1. *There exists $\rho_0(M_\infty, N_\infty) > 0$, where M_∞ and N_∞ are given in Theorem 3.5, such that if*

$$\limsup_{t \rightarrow \infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow \infty} \|c\|_{L^\infty(\Omega)} < \rho_0(M_\infty, N_\infty),$$

and for some t_0 the coefficients of (3.1) satisfy for $t \geq t_0$ the assumptions of Theorem 3.14, then for any bounded set of positive initial data bounded away from zero, all solutions of (3.1) that start at a sufficiently large $s > t_0$ have the same asymptotic behavior as $t \rightarrow \infty$.

In particular, all complete positive trajectories in the pullback attractor have the same asymptotic behavior as $t \rightarrow \infty$.

Proof. Note that Theorem 3.14 implies that all forward solutions of (3.1) that start at $s \geq t_0$ are uniformly nondegenerate with respect to a bounded set of initial data $u_s > 0$, $v_s > 0$, bounded away from zero, and the coefficients. In particular, from Propositions 4.2 and 4.4 the constant M_1 in (5.5) can be taken independent of such $u_s > 0$, $v_s > 0$ and the coefficients.

Moreover, for such initial data and $t > s \geq t_0$, we have in (6.1)

$$\begin{aligned} p_{11}(t, x) &= -a(t, x)u_1 \leq -a_L\psi_{11}(x) = -\varphi_{11}(x), \\ p_{22}(t, x) &= -d(t, x)v_1 \leq -d_L\psi_{22}(x) = -\varphi_{22}(x) \end{aligned}$$

with $\psi_{11}(x)$ and $\psi_{22}(x)$ independent $u_s > 0$, $v_s > 0$, and of the coefficients. In particular $\varphi_{11}(x)$ and $\varphi_{22}(x)$ satisfy the assumptions in Proposition 4.2.

Hence the threshold value $\rho > 0$ in Theorem 5.1 is also uniform for $u_s > 0, v_s > 0$, and of the coefficients as in Theorem 3.14.

Now we have in (6.1) $p_{12}(t, x) = -b(t, x)u_1, p_{21}(t, x) = -c(t, x)v_1$, and hence (5.3) is satisfied if

$$\limsup_{t \rightarrow \infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow \infty} \|c\|_{L^\infty(\Omega)} < \rho^2(p, \Omega, M_\infty, N_\infty) = \rho_0,$$

where M_∞ and N_∞ are given in Theorem 3.5.

Therefore, from Theorem 5.1, all solutions have the same forwards behavior. \square

Our next result proves that if there is a complete trajectory that is nondegenerate at $-\infty$, then it must be unique and be pullback attracting, as in Theorem 5.2.

THEOREM 6.2. *Assume that there exists a complete, bounded solution of (3.1) that is nondegenerate at $-\infty, U^*(t), t \in \mathbf{R}$.*

Then there exists $\rho_0(M_\infty, N_\infty) > 0$, where M_∞ and N_∞ are given in Theorem 3.5, such that if

$$\limsup_{t \rightarrow -\infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow -\infty} \|c\|_{L^\infty(\Omega)} < \rho_0(M_\infty, N_\infty),$$

and for some t_0 the coefficients of (3.1) satisfy for $t \leq t_0$ the assumptions of Theorem 3.14, then $U^(t)$ is the unique bounded complete solution of (3.1) that is nondegenerate at $-\infty$. Moreover, for every $t \in \mathbf{R}, U^*(t)$ pullback attracts solutions $U_1(t, s)$ such that $U_1(s)$ are positive and bounded as $s \rightarrow -\infty$.*

If, in addition,

$$\limsup_{t \rightarrow \infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow \infty} \|c\|_{L^\infty(\Omega)} < \rho_0(M_\infty, N_\infty),$$

and for some t_1 the coefficients of (3.1) satisfy for $t \geq t_1$ the assumptions of Theorem 3.14, then for any $s \in \mathbf{R}$ and for any positive solution $U(t, s)$ of (3.1) we have

$$U(t, s) - U^*(t) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Proof. Assume there exists a complete, bounded nondegenerate solution at $-\infty$. Then Theorem 3.14 implies that all bounded nondegenerate solutions at $-\infty$ are uniformly nondegenerate with respect to the coefficients. In particular, from Propositions 4.2 and 4.4 the constant M_1 in (5.5) can be taken independent of the complete nondegenerate solution under consideration and of the coefficients. Moreover, we have in (6.1)

$$\begin{aligned} p_{11}(t, x) &= -a(t, x)u_1 \leq -a_L\psi_{11}(x) = -\varphi_{11}(x), \\ p_{22}(t, x) &= -d(t, x)v_1 \leq -d_L\psi_{22}(x) = -\varphi_{22}(x) \end{aligned}$$

with $\psi_{11}(x)$ and $\psi_{22}(x)$ independent of the complete nondegenerate solution and of the coefficients. In particular $\varphi_{11}(x), \varphi_{22}(x)$ satisfy the assumptions in Proposition 4.2.

Hence the threshold value $\rho > 0$ in Theorem 5.2 is also independent of the complete nondegenerate solution and of the coefficients.

Now we have in (6.1) $p_{12}(t, x) = -b(t, x)u_1, p_{21}(t, x) = -c(t, x)v_1$, and hence (5.7) is satisfied if

$$\limsup_{t \rightarrow -\infty} \|b\|_{L^\infty(\Omega)} \limsup_{t \rightarrow -\infty} \|c\|_{L^\infty(\Omega)} < \rho^2(p, \Omega, M_\infty, N_\infty) = \rho_0.$$

Therefore, from Theorem 5.2, there exists at most a complete nondegenerate solution at $-\infty$.

To show that $U^*(t)$ is pullback attracting, observe that for sufficiently negative t_0 we can proceed as in the proof of Theorem 5.2 to conclude that $U^*(t)$ pullback attracts solutions $U_1(t, s)$ such that $U_1(s)$ is positive and bounded as $s \rightarrow -\infty$.

The rest follows from Theorem 6.1. \square

7. Conclusions. We have obtained some results on permanence in nonautonomous Lotka–Volterra models without the assumption of any kind of periodicity. In particular we have found conditions under which there exists at least one complete trajectory, and for which all trajectories converge together as $t \rightarrow +\infty$. The key argument is a perturbation result for an associated linear system satisfied by the difference between two solutions, and using this we have been able to treat all the different classical cases – competition, symbiosis, and prey–predator – in a unified way. While this unified approach has its advantages, our method requires at least one of the coupling parameters in the system to be sufficiently small. Hence, we hope that a more detailed study of each particular situation could lead to some improvements in the conditions imposed on the nonautonomous terms while still using similar techniques.

It is a very interesting open problem to prove, for this Lotka–Volterra example, the existence of a complete trajectory that is nondegenerate at $-\infty$. Given this nondegeneracy one would get the uniqueness of such a trajectory, and its pullback attracting property. We believe that use of the concepts of sub- and super-trajectories (cf. Arnold and Chueshov [3] and Chueshov [10]), along with the sub- and super-solutions technique (cf. for example Pao [30]) should be able to provide this, and we intend to pursue this direction in a future paper.

However, it is certainly the case that the hypothesis that the time-dependent terms are bounded is important throughout the literature, as this assumption implies the existence of bounded global solutions, and in particular of bounded attracting trajectories. As the analysis in Langa et al. [23] shows, different kinds of forwards asymptotic behavior, such as the nonexistence of asymptotically stable trajectories, is possible if solutions are allowed to be unbounded.

REFERENCES

- [1] H. AMANN, *Nonhomogeneous linear and quasilinear elliptic and parabolic boundary value problems*, in Function Spaces, Differential Operators and Nonlinear Analysis (Friedrichroda, 1992), Teubner-Texte Math. 133, Teubner, Stuttgart, 1993, pp. 9–126.
- [2] H. AMANN, *Linear and Quasilinear Parabolic Problems*, Vol. I, Abstract Linear Theory. Monographs in Mathematics, 89, Birkhäuser Boston Inc., Boston, MA, 1995.
- [3] L. ARNOLD AND I. CHUESHOV, *Order-preserving random dynamical systems: Equilibria, attractors, applications*, Dynam. Stability Systems, 13 (1998), pp. 265–280.
- [4] R. S. CANTRELL AND C. COSNER, *Should a park be an island?*, SIAM J. Appl. Math., 53 (1993), pp. 219–252.
- [5] R. S. CANTRELL AND C. COSNER, *Practical persistence in ecological models via comparison methods*, Proc. R. Soc. Edin., 126A (1996), pp. 247–272.
- [6] R. S. CANTRELL AND C. COSNER, *Spatial Ecology via Reaction-Diffusion Equations*, John Wiley & Sons. Ltd., New York, 2003.
- [7] R. S. CANTRELL, C. COSNER, AND V. HUTSON, *Permanence in ecological systems with spatial heterogeneity*, Proc. R. Soc. Edin., 123A (1993), pp. 533–559.
- [8] R. S. CANTRELL, C. COSNER, AND V. HUTSON, *Ecological models, permanence and spatial heterogeneity*, Rocky Mountain J. Math., 26 (1996), pp. 1–35.
- [9] V. V. CHEPYZHOV AND M. I. VISHIK, *Attractors for Equations of Mathematical Physics*, AMS Colloquium Publications, Vol. 49, AMS, Providence, RI, 2002.
- [10] I. CHUESHOV, *Monotone Random Systems Theory and Applications*, Lecture Notes in Math. 1779, Springer-Verlag, Berlin, 2002.
- [11] H. CRAUEL, A. DEBUSSCHE, AND F. FLANDOLI, *Random attractors*, J. Dynam. Differential Equations, 9 (1997), pp. 397–341.

- [12] M. DELGADO, J. LÓPEZ-GÓMEZ, AND A. SUÁREZ, *On the symbiotic Lotka-Volterra model with diffusion and transport effects*, J. Differential Equations, 160 (2000), pp. 175–262.
- [13] J. E. FURTER AND J. LÓPEZ-GÓMEZ, *On the existence and uniqueness of coexistence states for the Lotka-Volterra competition model with diffusion and spatially dependent coefficients*, Nonlinear Anal., 25 (1995), pp. 363–398.
- [14] J. HALE, *Asymptotic Behavior of Dissipative Systems*, Math. Surveys and Monographs, AMS, Providence, RI, 1998.
- [15] J. HALE AND P. WALTMAN, *Persistence in infinite dimensional systems*, SIAM J. Math. Anal., 9 (1989), pp. 388–395.
- [16] D. HENRY, *Geometric Theory of Semilinear Parabolic Equations*, Lecture Notes Math. 840, Springer-Verlag, Berlin, 1981.
- [17] P. HESS, *Periodic-Parabolic Boundary Value Problems and Positivity*, Pitman Research Notes in Mathematics 247, Harlow Longman, 1991.
- [18] P. HESS AND A. C. LAZER, *On an abstract competition model and applications*, Nonlinear Anal., 16 (1991), pp. 917–940.
- [19] G. HETZER AND W. SHEN, *Uniform persistence, coexistence, and extinction in almost periodic/nonautonomous competition diffusion systems*, SIAM J. Math. Anal., 34 (2002), pp. 204–221.
- [20] V. HUTSON AND K. SCHMITT, *Permanence in dynamical systems*, Math. Biosci., 111 (1992), pp. 1–71.
- [21] P. E. KLOEDEN AND B. SCHMALFUSS, *Asymptotic behaviour of non-autonomous difference inclusions*, Systems Control Lett., 33 (1998), pp. 275–280.
- [22] N. LAKOS, *Existence of steady-state solutions for one-predator-two prey system*, SIAM J. Math. Anal., 21 (1990), pp. 647–659.
- [23] J. A. LANGA, J. C. ROBINSON, A. RODRÍGUEZ-BERNAL, A. SUÁREZ, AND A. VIDAL-LÓPEZ, *Existence and nonexistence of unbounded forward attractor for a class of non-autonomous reaction diffusion equations*, Discrete Contin. Dyn. Syst., 18 (2007), pp. 483–497.
- [24] J. A. LANGA, J. C. ROBINSON, AND A. SUÁREZ, *Stability, instability, and bifurcation phenomena in non-autonomous differential equations*, Nonlinearity, 15 (2002), pp. 887–903.
- [25] J. A. LANGA, J. C. ROBINSON, AND A. SUÁREZ, *Pullback permanence in the non-autonomous Lotka-Volterra competition model*, J. Differential Equations, 190 (2003), pp. 214–238.
- [26] A. W. LEUNG, *Monotone schemes for semilinear elliptic systems related to ecology*, Math. Methods Appl. Sci., 4 (1982), pp. 272–285.
- [27] J. LÓPEZ-GÓMEZ, *On the structure of the permanence region for competing species models with general diffusivities and transport effects*, Discrete Contin. Dyn. Syst., 2 (1996), pp. 525–542.
- [28] J. LÓPEZ-GÓMEZ AND R. PARDO, *Existence and uniqueness of coexistence states for the predator-prey model with diffusion: The scalar case*, Differential Int. Equations, 6 (1993), pp. 1025–1031.
- [29] J. LÓPEZ-GÓMEZ AND J. C. SABINA DE LIS, *Coexistence states and global attractivity for some convective diffusive competing species models*, Trans. Amer. Math. Soc., 347 (1995), pp. 3797–3833.
- [30] C. V. PAO, *Nonlinear parabolic and elliptic equations*, Plenum, New York, 1992.
- [31] J. C. ROBINSON, A. RODRÍGUEZ-BERNAL, AND A. VIDAL-LÓPEZ, *Pullback attractors and extremal complete trajectories for non-autonomous reaction-diffusion problems*, J. Differential Equations, 238 (2007), pp. 289–337.
- [32] A. RODRÍGUEZ-BERNAL, *On linear and nonlinear non-autonomous parabolic equations*, Departamento de Matemática Aplicada, UCM, Preprint Series MA-UCM-2006-15.
- [33] A. RODRÍGUEZ-BERNAL, *Perturbation of the exponential type of linear nonautonomous parabolic equations and applications to nonlinear equations*, Departamento de Matemática Aplicada, UCM, Preprint Series MA-UCM-2008-08.
- [34] A. RODRÍGUEZ-BERNAL AND A. VIDAL-LÓPEZ, *Existence, uniqueness and attractivity properties of positive complete trajectories for non-autonomous reaction-diffusion problems*, Discrete Contin. Dyn. Syst., 18 (2007), pp. 537–567.
- [35] B. SCHMALFUSS, *Attractors for the non-autonomous dynamical systems*, Proceedings of Equadiff 99 Berlin, B. Fiedler, K. Gröger, and J. Sprekels, eds., Singapore World Scientific, Singapore, 2000, pp. 684–689.
- [36] Y. YAMADA, *Stability of steady states for predator-prey diffusion equations with homogeneous Dirichlet conditions*, SIAM J. Math. Anal., 21 (1990), pp. 327–345.
- [37] R. TEMAM, *Infinite-Dimensional Dynamical Systems in Mechanics and Physics*, Springer-Verlag, New York, 1988.

AN ENTIRE SOLUTION TO THE LOTKA–VOLTERRA COMPETITION-DIFFUSION EQUATIONS*

YOSHIHISA MORITA[†] AND KOICHI TACHIBANA[‡]

Abstract. We deal with a system of Lotka–Volterra competition-diffusion equations on \mathbb{R} , which is a competing two species model with diffusion. It is known that the equations allow traveling waves with monotone profile. In this article we prove the existence of an entire solution which behaves as two monotone waves propagating from both sides of the x -axis, where an entire solution is meant by a classical solution defined for all space and time variables. The global dynamics for this entire solution exhibits the extinction of the inferior species by the superior one invading from both sides. The proof is carried out by applying the comparison principle for the competition-diffusion equations, that is, using an appropriate pair of a subsolution and a supersolution.

Key words. Lotka–Volterra, competition-diffusion, entire solution, traveling wave

AMS subject classifications. 35K57, 35B05, 35B40, 92B05

DOI. 10.1137/080723715

1. Introduction. In the field of population biology Lotka–Volterra competition equations are well accepted as a physiological model describing competing interaction of multiple species. Here we restrict our attention to a two species model. Taking random movement of the species into account, we get to the Lotka–Volterra competition-diffusion equations that are obtained by coupling diffusion terms to the Lotka–Volterra equations. In this article we are dealing with the following Lotka–Volterra competition-diffusion equations on \mathbb{R} :

$$(1.1) \quad \begin{cases} u_t = u_{xx} + (1 - u - k_1 v)u, \\ v_t = dv_{xx} + a(1 - v - k_2 u)v \end{cases} \quad (x \in \mathbb{R}),$$

where k_1, k_2, a , and d are positive constants. The variables $u(x, t)$ and $v(x, t)$ stand for the population density of two competing species. Thus we consider nonnegative $u(x, t), v(x, t)$. We note that the above system is normalized so that it has the equilibrium solutions $(u, v) = (1, 0), (0, 1)$, set as

$$e_u := (1, 0), \quad e_v := (0, 1).$$

In the diffusion-free case we can classify the asymptotic behavior of the solution as $t \rightarrow \infty$, depending on k_1, k_2 :

- (i) $(u(t), v(t)) \rightarrow e_u$ (u always wins) if $0 < k_1 < 1 < k_2$.
- (ii) Depending on the initial condition, $(u, v) \rightarrow e_u$ or $(u(t), v(t)) \rightarrow e_v$ in general if $1 < k_1, k_2$.
- (iii) $(u(t), v(t)) \rightarrow e_v$ (v always wins) if $0 < k_2 < 1 < k_1$.

*Received by the editors May 8, 2008; accepted for publication (in revised form) September 23, 2008; published electronically January 28, 2009. This research was supported in part by the Grant-in-Aid for Scientific Research (B) 19340026 and the Grant-in-Aid for Exploratory Research 19654030, Japan Society for the Promotion of Science.

<http://www.siam.org/journals/sima/40-6/72371.html>

[†]Department of Applied Mathematics and Informatics, Ryukoku University, Seta Otsu 520-2194, Japan (morita@rins.ryukoku.ac.jp).

[‡]Daiwa Technique Laboratory Ltd., 915 Mitsushima Kadoma, Osaka 571-0015, Japan (ham_one14@msn.com).

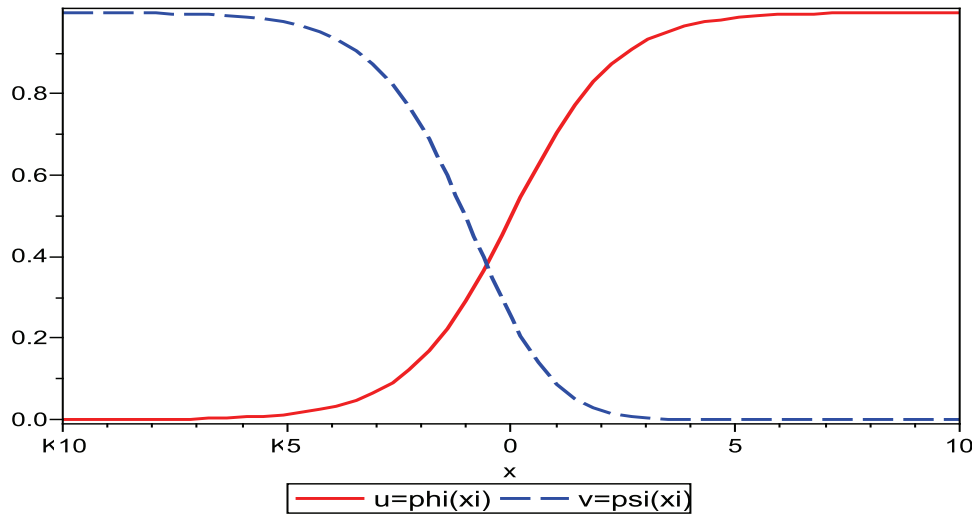


FIG. 1. The profile of a monotone traveling wave solution of (1.1). Monotone increasing and decreasing curves correspond to $u = \phi(\xi)$ and $v = \psi(\xi)$, respectively. Parameter values are given by $a = 3/2$, $d = 1/3$, $k_1 = 3/2$, and $k_2 = 2$.

- (iv) $(u(t), v(t))$ converges to a positive equilibrium (u and v coexist) if $0 < k_1, k_2 < 1$.

Here we do not treat case (iv). Then we may assume $k_2 > 1$ because the third case (iii) is obtained by exchanging the roles of u and v .

The above equations (1.1) are used for describing the invasion of the superior species. As a matter of fact, by virtue of the diffusion, we can observe behavior such as the superior species propagating from one side of the x -axis to the other side, pushing back the inferior species. This dynamics of the propagation is mathematically characterized by a traveling wave solution $(u, v) = (\phi(x + ct), \psi(x + ct))$ satisfying

$$(1.2) \quad \begin{cases} \phi'' - c\phi' + (1 - \phi - k_1\psi)\phi = 0, \\ d\psi'' - c\psi' + a(1 - \psi - k_2\phi)\psi = 0, \end{cases}$$

with

$$(1.3) \quad \lim_{\xi \rightarrow -\infty} (\phi(\xi), \psi(\xi)) = \mathbf{e}_v, \quad \lim_{\xi \rightarrow \infty} (\phi(\xi), \psi(\xi)) = \mathbf{e}_u,$$

$$(1.4) \quad \phi(\xi), \psi(\xi) > 0, \quad \phi'(\xi) > 0, \quad \psi'(\xi) < 0,$$

where $' = d/d\xi$, $'' = d^2/d\xi^2$. A solution to (1.2) with (1.3) and (1.4) gives a profile of the traveling wave. Here we assume the monotonicity of (1.4) which arises in most cases. A traveling wave with the monotonicity is called a monotone wave or a wave front.

We note that for a traveling wave solution $(\phi(\xi), \psi(\xi)) = (\phi(x + ct), \psi(x + ct))$ to (1.2)–(1.4), the reflected $(\phi(-x + ct), \psi(-x + ct))$ also gives a traveling wave with opposite speed. Then $(\tilde{\phi}(\xi), \tilde{\psi}(\xi)) = (\phi(-x + ct), \psi(-x + ct))$ satisfies

$$\lim_{\xi \rightarrow -\infty} (\tilde{\phi}(\xi), \tilde{\psi}(\xi)) = \mathbf{e}_u, \quad \lim_{\xi \rightarrow \infty} (\tilde{\phi}(\xi), \tilde{\psi}(\xi)) = \mathbf{e}_v.$$

Thus we obtain traveling wave solutions with opposite speeds simultaneously if we obtain a solution satisfying (1.3)–(1.4).

As for the study of traveling wave solutions, there are various works, including [3], [4], [15], and [7]. In particular, see [8], [9], [10], [11], [12], and the references therein.

Although the study for a traveling wave is a central issue of the Lotka–Volterra competition-diffusion equations, it is not enough for mathematical understanding of the dynamical structure of solutions to (1.1). In fact, as seen in the recent development of the study for entire solutions to a scalar reaction-diffusion equation, some combination of two traveling waves yields characteristic dynamics for the equation (see [6], [16], [2], [5], [1], [13]), where the entire solution is meant by a classical solution which is defined for every x and t . Those results revealed that the equation allows types of entire solutions other than equilibrium solutions and traveling wave solutions. For example, in [6], [16], [2], [5], it is shown that there exists an entire solution which behaves as two wave fronts propagating from both sides of the x -axis and then annihilating, while [13] shows entire solutions which behave as two wave fronts merging.

In this article, motivated by the work of [2] and [5], we consider the entire solution which behaves like two monotone traveling waves propagating from both sides of the x -axis. As t goes forward, this solution converges to the equilibrium e_u when u is superior. Namely, we are interested in the entire solution $(u(x, t), v(x, t))$ which behaves like $(\phi(x + ct), \psi(x + ct))$ for $x \in (0, \infty)$ and $(\phi(-x + ct), \psi(-x + ct))$ for $x \in (-\infty, 0)$ as $t \rightarrow -\infty$ and converges to a uniform state as $t \rightarrow \infty$. This solution corresponds to the phenomenon that the superior species invades from both sides and causes the extinction of the inferior species.

We prepare some notation for stating the main result. We assume that u is always superior so that the traveling wave solution of (1.2) with (1.3) and (1.4) has a positive speed $c > 0$. In addition, we allow the equations to have solutions with different speeds (corresponding to case (i)). As a matter of fact, we have traveling waves of the Fisher type for the case (i); that is, there is a minimum speed $c_{min} > 0$ so that the equations allow a family of monotone traveling wave solutions for $c \geq c_{min} > 0$ (see [10]). Thus we set a pair of traveling wave solutions

$$(\phi_j(\xi), \psi_j(\xi)) \quad (j = 1, 2)$$

satisfying

$$(1.5) \quad \begin{cases} \phi_j'' - c_j \phi_j' + (1 - \phi_j - k_1 \psi_j) \phi_j = 0, \\ d \psi_j'' - c_j \psi_j' + a(1 - \psi_j - k_2 \phi_j) \psi_j = 0 \end{cases} \quad (j = 1, 2),$$

with

$$(1.6) \quad \lim_{\xi \rightarrow -\infty} (\phi_j(\xi), \psi_j(\xi)) = e_v, \quad \lim_{\xi \rightarrow \infty} (\phi_j(\xi), \psi_j(\xi)) = e_u,$$

$$(1.7) \quad \phi_j(\xi), \psi_j(\xi) > 0, \quad \phi_j'(\xi) > 0, \quad \psi_j'(\xi) < 0,$$

where c_j ($j = 1, 2$) are positive. Note that we consider the case when they coincide, namely, $(\phi_1(\xi), \psi_1(\xi))$ identically equals $(\phi_2(\xi), \psi_2(\xi))$ up to a phase shift.

We assume that there is a positive constant η_0 such that

$$(1.8) \quad \eta_0 \leq \frac{\phi_j(\xi)}{1 - \psi_j(\xi)} \quad (\xi \leq 0)$$

holds. Then we obtain the following theorem, which is the main result in this article.

THEOREM 1.1. *Assume $k_1 \neq 1$ and $k_2 > 1$ in (1.1). Let (ϕ_j, ψ_j) be a solution to (1.5) with (1.6), (1.7), and $c_j > 0$ ($j = 1, 2$). If the condition (1.8) holds, then there exists a solution $(u(x, t), v(x, t))$ ($(x, t) \in \mathbb{R} \times \mathbb{R}$) of (1.1) satisfying*

$$\begin{aligned} \lim_{t \rightarrow -\infty} \sup_{x \geq (c_2 - c_1)t/2} \{|u(x, t) - \phi_1(x + c_1t)| + |v(x, t) - \psi_1(x + c_1t)|\} &= 0, \\ \lim_{t \rightarrow -\infty} \sup_{x \leq (c_2 - c_1)t/2} \{|u(x, t) - \phi_2(-x + c_2t)| + |v(x, t) - \psi_2(-x + c_2t)|\} &= 0, \end{aligned}$$

and

$$\lim_{t \rightarrow \infty} \sup_{x \in \mathbb{R}} \{|u(x, t) - 1| + |v(x, t)|\} = 0.$$

Remark 1.1. If $k_1 > 1$, the monotone traveling wave is unique up to translation (for instance, see [8]). Hence, when $k_1 > 1$, $(\phi_1, \psi_1) = (\phi_2, \psi_2)$ in the theorem.

Considering the fact that the equations (1.1) are invariant under the translation in x and the time shift, we easily see the following result from the above theorem.

COROLLARY 1.2. *Let $(u(x, t), v(x, t))$ be an entire solution to (1.1) given in Theorem 1.1. Given θ_1, θ_2 , there exist ξ and τ such that*

$$\begin{aligned} \lim_{t \rightarrow -\infty} \sup_{x \geq (c_2 - c_1)t/2} \{|u(x + \xi, t + \tau) - \phi_1(x + c_1t + \theta_1)| \\ + |v(x + \xi, t + \tau) - \psi_1(x + c_1t + \theta_1)|\} &= 0, \\ \lim_{t \rightarrow -\infty} \sup_{x \leq (c_2 - c_1)t/2} \{|u(x + \xi, t + \tau) - \phi_2(-x + c_2t + \theta_2)| \\ + |v(x + \xi, t + \tau) - \psi_2(-x + c_2t + \theta_2)|\} &= 0, \end{aligned}$$

and

$$\lim_{t \rightarrow \infty} \sup_{x \in \mathbb{R}} \{|u(x + \xi, t + \tau) - 1| + |v(x + \xi, t + \tau)|\} = 0.$$

The condition (1.8) is crucial in our argument, though it might be a technical condition. Here we give a sufficient condition for (1.8), which is easier to check. Let $\lambda_3^{(j)}$ be a positive root of

$$(1.9) \quad d\lambda^2 - c_j\lambda - a = 0.$$

For $k_1 > 1$, we have a unique positive root $\lambda_4^{(j)}$ of

$$(1.10) \quad \lambda^2 - c_j\lambda - (k_1 - 1) = 0,$$

while for $0 < k_1 < 1$, there is another positive root $\lambda_5^{(j)}$ ($\leq \lambda_4^{(j)}$) of (1.10) under the condition

$$(1.11) \quad c_j \geq c_{min} \geq 2\sqrt{1 - k_1},$$

where $\lambda_4^{(j)} = \lambda_5^{(j)}$ if $c_j = 2\sqrt{1 - k_1}$. Those roots $\lambda_k^{(j)}$ ($k = 3, 4, 5$) are related to the convergence rate of the traveling waves as $(\phi_j(\xi), \psi_j(\xi)) \rightarrow (0, 1)$ ($\xi \rightarrow -\infty$). We notice that

$$(1.12) \quad d \frac{\psi_j''}{1 - \psi_j} - c_j \frac{\psi_j'}{1 - \psi_j} + a\psi_j \left(1 - k_2 \frac{\phi_j}{1 - \psi_j} \right) = 0$$

by the second equation of (1.5). We can prove the limit of $\psi'_j(\xi)/(1 - \psi_j(\xi))$ as $\xi \rightarrow -\infty$ exists. In fact, we have the next lemma.

LEMMA 1.3. *Assume the same assumption of Theorem 1.1. Then (1.8) holds if*

$$(1.13) \quad \lim_{\xi \rightarrow -\infty} \frac{-\psi'_j(\xi)}{1 - \psi_j(\xi)} \neq \lambda_3^{(j)}.$$

Moreover, the condition (1.13) is realized for $\lambda_4^{(j)} < \lambda_3^{(j)}$ ($j = 1, 2$).

Remark 1.2. The condition $\lambda_4^{(j)} < \lambda_3^{(j)}$ ($j = 1, 2$) is clear and easily checked. However, there are cases which break this condition but allow (1.13). For example, (1.13) is verified by the condition $\lambda_5^{(j)} < \lambda_3^{(j)} \leq \lambda_4^{(j)}$ together with a generic condition for the asymptotic behavior of the traveling wave as $t \rightarrow -\infty$. We discuss this in section 2.

We briefly sketch an outline of the proof of Theorem 1.1. We will use the following functions to approximate the entire solution as $t \approx -\infty$:

$$(1.14) \quad \begin{cases} \underline{u} = \max\{\phi_1(x + c_1t), \phi_2(-x + c_2t)\}, \\ \bar{v} = \min\{\psi_1(x + c_1t), \psi_2(-x + c_2t)\}, \end{cases}$$

and

$$(1.15) \quad \begin{cases} \bar{u} = \phi_1(x + p_1(t)) + \phi_2(-x + p_2(t)) - \phi_1(x + p_1(t))\phi_2(-x + p_2(t)), \\ \underline{v} = \psi_1(x + p_1(t))\psi_2(-x + p_2(t)), \end{cases}$$

where

$$p_j(t) \approx c_jt \quad (t \approx -\infty).$$

Note that (\underline{u}, \bar{v}) and (\bar{u}, \underline{v}) certainly have the profile we are supposing as $t \rightarrow -\infty$. Since the comparison principle does work in the present competition-diffusion system, we prove that (1.14) and (1.15) are a subsolution and a supersolution with suitable $p_j(t)$, respectively, if $t < 0$. This yields a solution sandwiched by the subsolution and the supersolution, defined for every negative t . Then such a solution can be extended for any positive time. In what follows, we obtain the desired entire solution.

Although our idea for the proof of the main result is based on the previous work of [5], [1], and [13], we encounter a difficulty when we prove that the pairing of (1.15) is a supersolution. Fortunately, noticing the identity

$$(1.16) \quad \begin{aligned} &(u_1 + u_2 - u_1u_2)v_1v_2 - (1 - u_2)u_1v_1 - (1 - u_1)u_2v_2 \\ &= -u_1(1 - u_2)v_1(1 - v_2) - (1 - u_1)u_2(1 - v_1)v_2 + u_1u_2v_1v_2, \end{aligned}$$

we can overcome the difficulty (see section 4 for the details).

In the above theorem we obtain an entire solution characterized by the asymptotic behavior as $t \rightarrow -\infty$. Once we obtain such a solution, it is easy to prove the asymptotic behavior as $t \rightarrow \infty$ stated in the theorem, since the subsolution is defined globally in time. However, we have only the one side estimated by the subsolution; hence precise transient behavior for positive t is not described yet. Here, using an exact expression for the traveling wave solution, we show the entire behavior of the solution for specific parameter values. In the paper [14], traveling wave solutions are

exactly given in terms of the hyperbolic tangent function under a constraint of the parameters. We use one of them, which is given by

$$(1.17) \quad \phi(\xi) = \frac{1}{2} \left\{ 1 + \tanh \left(\frac{\sqrt{k_1} \xi}{2\sqrt{2}} \right) \right\}, \quad \psi(\xi) = \frac{1}{4} \left\{ 1 - \tanh \left(\frac{\sqrt{k_1} \xi}{2\sqrt{2}} \right) \right\}^2,$$

where

$$(1.18) \quad d = \frac{a}{3k_1}, \quad a = \frac{6 - 3k_1}{3k_2 - 5}, \quad c = \frac{2 - k_1}{\sqrt{2k_1}},$$

with the condition

$$(1.19) \quad 0 < k_1 < 2, \quad 5/3 < k_2.$$

It is a lengthy but simple computation to verify that (1.17) with (1.18) and (1.19) is a traveling wave of (1.1) satisfying (1.2)–(1.4). Moreover, this solution satisfies (1.13) as well as (1.8). Indeed, we can compute the corresponding positive roots of (1.9) and (1.10) as

$$\lambda_3 = \frac{1}{2} \sqrt{\frac{k_1}{2}} \{3k_2 - 5 + \sqrt{(3k_2 - 5)^2 + 24}\} > \sqrt{3k_1},$$

$$\lambda_4 = \begin{cases} \sqrt{k_1/2} & (2/3 \leq k_1), \\ \sqrt{2}(1 - k_1)/\sqrt{k_1} & (0 < k_1 < 2/3) \end{cases}$$

and

$$\lambda_5 = \begin{cases} \sqrt{2}(1 - k_1)/\sqrt{k_1} & (2/3 < k_1 < 1), \\ \sqrt{k_1/2} & (0 < k_1 \leq 2/3), \end{cases}$$

respectively.

We rewrite (1.17) as

$$\phi(\xi) = 1 - \frac{1}{1 + \exp(\xi\sqrt{k_1/2})}, \quad \psi(\xi) = \frac{1}{\{1 + \exp(\xi\sqrt{k_1/2})\}^2}.$$

Modify this expression as

$$(1.20) \quad u_*(x, t) := 1 - \frac{1}{1 + \Phi(x, t)}, \quad v_*(x, t) = \frac{1}{(1 + \Phi(x, t))^2}$$

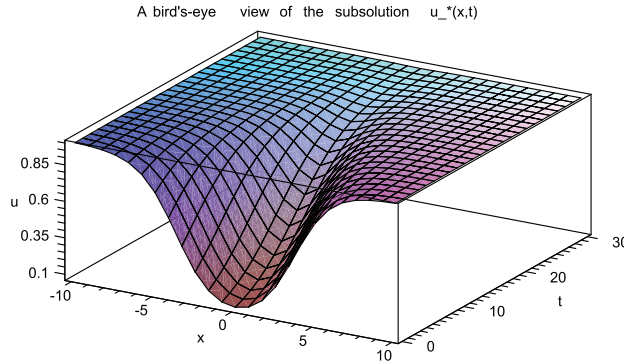
with

$$(1.21) \quad \Phi = \cosh \left(\frac{x\sqrt{k_1}}{\sqrt{2}} \right) \exp(\omega_0 t), \quad \omega_0 := 1 - \frac{k_1}{2}.$$

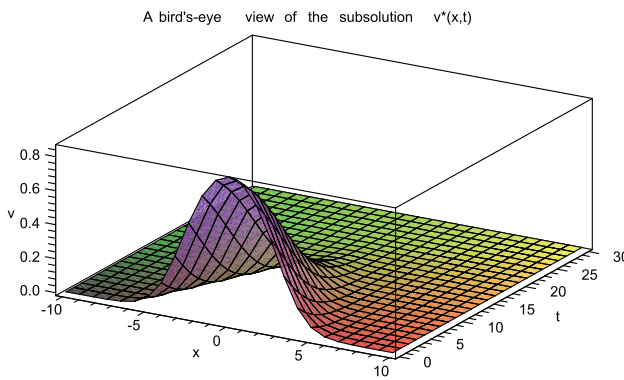
As for a profile of the functions of (1.20), see Figure 2.

We also define

$$(1.22) \quad u_*(x, t) := 1 - \frac{1}{1 + \Psi(x, t)}, \quad v_*(x, t) = \frac{1}{(1 + \Psi(x, t))^2}$$



(a) A bird's eye view of $u_*(x, t - 20)$



(b) A bird's eye view of $v^*(x, t - 20)$

FIG. 2. Bird's eye views of the functions $u_*(x, t - 20)$ and $v^*(x, t - 20)$ ($0 \leq t \leq 30$, $-10 \leq x \leq 10$), respectively, which are defined by (1.20) and (1.21). Parameter values are given by (1.18) with $k_1 = 3/2$ and $k_2 = 2$.

with

$$(1.23) \quad \Psi = \cosh\left(\frac{x\sqrt{k_1}}{\sqrt{2}}\right) \exp(p(t)),$$

where $p(t)$ is given by a solution to

$$(1.24) \quad \dot{p}(t) = \omega_0 + \frac{L_0 \exp(p(t))}{1 + \exp(p(t))}, \quad p(0) = p_0,$$

$$L_0 := \max\{k_1, a/2\}.$$

Note that the solution $p(t)$ is a monotone increasing function satisfying

$$\lim_{t \rightarrow -\infty} p(t) = -\infty, \quad \lim_{t \rightarrow \infty} p(t) = \infty.$$

The next result shows that the above functions give an approximation globally in time together with the existence of an entire solution.

PROPOSITION 1.4. *In addition to (1.18) and (1.19), assume that the solution $p(t)$ of (1.24) satisfies $p(t_0) = \omega_0 t_0$ for an arbitrarily given t_0 . Then a solution*

$(u(x, t), v(x, t))$ to (1.1) with the initial condition

$$(u(x, t_0), v(x, t_0)) = (u_*(x, t_0), v^*(x, t_0)) = (u^*(x, t_0), v_*(x, t_0))$$

satisfies

$$(1.25) \quad \begin{cases} u_*(x, t) \leq u(x, t) \leq u^*(x, t), \\ v_*(x, t) \leq v(x, t) \leq v^*(x, t) \end{cases} \quad ((t, x) \in (t_0, \infty) \times \mathbb{R}).$$

Moreover, if p_0 is a solution of

$$(1.26) \quad p_0 = \frac{L_0}{L_0 + \omega_0} \log \left(1 + \frac{L_0 + \omega_0}{\omega_0} \exp(p_0) \right),$$

then the asymptotics

$$\lim_{t \rightarrow -\infty} \sup_{x \in \mathbb{R}} [|u^*(x, t) - u_*(x, t)| + |v^*(x, t) - v_*(x, t)|] = 0$$

hold; hence there is an entire solution $(u(x, t), v(x, t))$ satisfying (1.25) with $t_0 = -\infty$.

We organize the paper as follows: In the next section we classify the asymptotic behavior of the traveling wave solutions as $\xi \rightarrow \pm\infty$. This result is given by investigating the linearized equations at the equilibrium points e_u and e_v . Then Lemma 1.3 is proved. Although the result is already established, we give a proof in the appendix for the reader’s convenience. In section 3 we give crucial estimates for the terms consisting of the traveling wave solutions. By virtue of the estimates we prove Theorem 1.1 together with Proposition 1.4 in section 4. The paper concludes with the appendix.

2. Asymptotic behavior of traveling front waves. All the results in Lemmas 2.1, 2.2, and 2.3 stated below were already obtained in [8] and [12]. For completeness of the present paper we give a proof in the appendix. We let $\lambda_1^{(j)}$ and $\lambda_2^{(j)}$ be negative roots of

$$(2.1) \quad \lambda^2 - c_j \lambda - 1 = 0$$

and

$$(2.2) \quad d\lambda^2 - c_j \lambda - a(k_2 - 1) = 0,$$

respectively. The first lemma is on the asymptotic behavior of the traveling wave solutions as $\xi \rightarrow \infty$.

LEMMA 2.1. *Let $(\phi_j(\xi), \psi_j(\xi))$ be traveling wave solutions to (1.5)–(1.7) with $c_j > 0$ ($j = 1, 2$), and let $\lambda_1^{(j)}$ and $\lambda_2^{(j)}$ be negative roots of (2.1) and (2.2), respectively. Then the solution exhibits the asymptotic behavior as $\xi \rightarrow \infty$ as follows:*

(i) *If $\lambda_2^{(j)} < \lambda_1^{(j)}$, then*

$$(2.3) \quad \begin{aligned} 1 - \phi_j(\xi) &= \alpha_j \exp(\lambda_1^{(j)} \xi) - \beta_j s_1^{(j)} \exp(\lambda_2^{(j)} \xi) + h.o.t., \\ \psi_j(\xi) &= \beta_j \exp(\lambda_2^{(j)} \xi) + h.o.t. \end{aligned}$$

hold for some numbers $\alpha_j \geq 0$ and $\beta_j > 0$, where

$$(2.4) \quad s_1^{(j)} := \frac{k_1}{(\lambda_2^{(j)})^2 - c_j \lambda_2^{(j)} - 1} = \begin{cases} > 0 & (\lambda_2^{(j)} < \lambda_1^{(j)}), \\ < 0 & (\lambda_2^{(j)} > \lambda_1^{(j)}). \end{cases}$$

(ii) If $\lambda_1^{(j)} < \lambda_2^{(j)}$, then

$$(2.5) \quad \begin{aligned} 1 - \phi_j(\xi) &= -\beta_j s_1^{(j)} \exp(\lambda_2^{(j)} \xi) + h.o.t., \\ \psi_j(\xi) &= \beta_j \exp(\lambda_2^{(j)} \xi) + h.o.t. \end{aligned}$$

for some $\beta_j > 0$, where $s_1^{(j)}$ is as in (2.4).

(iii) If $\lambda_1^{(j)} = \lambda_2^{(j)}$, then

$$(2.6) \quad \begin{aligned} 1 - \phi_j(\xi) &= \beta_j \xi \exp(\lambda_1^{(j)} \xi) + h.o.t., \\ \psi_j(\xi) &= -\beta_j \tau_1^{(j)} \exp(\lambda_1^{(j)} \xi) + h.o.t. \end{aligned}$$

for some $\beta_j > 0$, where

$$\tau_1^{(j)} = \frac{2\lambda_1^{(j)} - c_j}{k_1} < 0.$$

The second lemma is related to the asymptotic behavior of the traveling wave solution as $\xi \rightarrow -\infty$ for $k_1 > 1$, namely, the bistable case (ii) stated in the introduction.

LEMMA 2.2. *Let $(\phi_j(\xi), \psi_j(\xi))$ be traveling wave solutions to (1.5)–(1.7) with $c_j > 0$ ($j = 1, 2$) for $k_1 > 1$, and let $\lambda_3^{(j)}$ and $\lambda_4^{(j)}$ be positive roots of (1.9) and (1.10), respectively. Then the solution exhibits the asymptotic behavior as $\xi \rightarrow -\infty$ as follows:*

(i) If $\lambda_3^{(j)} < \lambda_4^{(j)}$, then

$$(2.7) \quad \begin{aligned} \phi_j(\xi) &= \beta_j \exp(\lambda_4^{(j)} \xi) + h.o.t., \\ 1 - \psi_j(\xi) &= \alpha_j \exp(\lambda_3^{(j)} \xi) - \beta_j s_2^{(j)} \exp(\lambda_4^{(j)} \xi) + h.o.t. \end{aligned}$$

hold for some numbers $\alpha_j \geq 0$ and $\beta_j > 0$, where

$$(2.8) \quad s_2^{(j)} := \frac{ak_2}{d(\lambda_4^{(j)})^2 - c_j \lambda_4^{(j)} - a} = \begin{cases} > 0 & (\lambda_3^{(j)} < \lambda_4^{(j)}), \\ < 0 & (\lambda_3^{(j)} > \lambda_4^{(j)}). \end{cases}$$

(ii) If $\lambda_4^{(j)} < \lambda_3^{(j)}$, then

$$(2.9) \quad \begin{aligned} \phi_j(\xi) &= \beta_j \exp(\lambda_4^{(j)} \xi) + h.o.t., \\ 1 - \psi_j(\xi) &= -\beta_j s_2^{(j)} \exp(\lambda_4^{(j)} \xi) + h.o.t. \end{aligned}$$

for some $\beta_j > 0$, where $s_2^{(j)}$ is defined as in (2.8).

(iii) If $\lambda_3^{(j)} = \lambda_4^{(j)}$, then

$$(2.10) \quad \begin{aligned} \phi_j(\xi) &= \beta_j \tau_2^{(j)} \exp(\lambda_3^{(j)} \xi) + h.o.t., \\ 1 - \psi_j(\xi) &= -\beta_j \xi \exp(\lambda_3^{(j)} \xi) + h.o.t. \end{aligned}$$

as $\xi \rightarrow -\infty$ for some β_j , where

$$(2.11) \quad \tau_2^{(j)} := \frac{2d\lambda_3^{(j)} - c_j}{k_1} > 0.$$

The next lemma is related to the asymptotic behavior as $\xi \rightarrow -\infty$ when $0 < k_1 < 1$, namely, the monostable case (i) stated in the introduction.

LEMMA 2.3. *Let $(\phi_j(\xi), \psi_j(\xi))$ be traveling wave solutions to (1.5)–(1.7) with $c_j > 0$ ($j = 1, 2$) for $k_1 \in (0, 1)$, and let $\lambda_5^{(j)}$ be a positive root of (1.10) with $\lambda_5^{(j)} \leq \lambda_4^{(j)}$, where $c_j \geq 2\sqrt{1-k_1}$. Then the solution behaves as $\xi \rightarrow -\infty$ in the following way:*

(i) $\lambda_5^{(j)} < \lambda_4^{(j)}$, $\lambda_4^{(j)}, \lambda_5^{(j)} \neq \lambda_3^{(j)}$: *there are numbers α_j, β_j , and nonnegative γ_j such that*

(2.12)

$$\phi_j(\xi) = \beta_j \exp(\lambda_4^{(j)} \xi) + \gamma_j \exp(\lambda_5^{(j)} \xi) + h.o.t.,$$

$$1 - \psi_j(\xi) = \alpha_j \exp(\lambda_3^{(j)} \xi) - \beta_j s_2^{(j)} \exp(\lambda_4^{(j)} \xi) - \gamma_j s_3^{(j)} \exp(\lambda_5^{(j)} \xi) + h.o.t.$$

with

$$\gamma_j \geq 0, \quad \beta_j > 0 \quad (\gamma_j = 0), \quad \alpha_j > 0 \quad (\lambda_3^{(j)} < \lambda_5^{(j)})$$

hold, where

$$(2.13) \quad s_3^{(j)} = \frac{ak_2}{d(\lambda_5^{(j)})^2 - c_j \lambda_5^{(j)} - a} = \begin{cases} > 0 & (\lambda_5^{(j)} > \lambda_3^{(j)}), \\ < 0 & (\lambda_5^{(j)} < \lambda_3^{(j)}). \end{cases}$$

(ii) $\lambda_5^{(j)} = \lambda_4^{(j)} < \lambda_3^{(j)}$: *there are β_j and γ_j such that*

(2.14)

$$\phi_j(\xi) = \beta_j \xi \exp(\lambda_4^{(j)} \xi) + \gamma_j \exp(\lambda_4^{(j)} \xi) + h.o.t.,$$

$$1 - \psi_j(\xi) = -\beta_j s_2^{(j)} \xi \exp(\lambda_4^{(j)} \xi) - \gamma_j s_2^{(j)} \exp(\lambda_4^{(j)} \xi) + h.o.t.$$

with $\beta_j \geq 0$ and $\gamma_j > 0$ ($\beta_j = 0$), where $s_2^{(j)}$ is defined as in (2.8).

(iii) $\lambda_5^{(j)} < \lambda_3^{(j)} = \lambda_4^{(j)}$: *there are β_j and γ_j such that*

(2.15)

$$\phi_j(\xi) = \beta_j \tau_2^{(j)} \exp(\lambda_3^{(j)} \xi) + \gamma_j \exp(\lambda_5^{(j)} \xi) + h.o.t.,$$

$$1 - \psi_j(\xi) = -\beta_j \xi \exp(\lambda_3^{(j)} \xi) - \gamma_j s_3^{(j)} \exp(\lambda_5^{(j)} \xi) + h.o.t.$$

with $\gamma_j \geq 0$ and $\beta_j > 0$ ($\gamma_j = 0$), where $\tau_2^{(j)}$ is defined as in (2.11).

(iv) $\lambda_5^{(j)} = \lambda_3^{(j)} < \lambda_4^{(j)}$: *there are α_j, β_j , and γ_j such that*

(2.16)

$$\phi_j(\xi) = \gamma_j \tau_2^{(j)} \exp(\lambda_3^{(j)} \xi) + \beta_j \exp(\lambda_4^{(j)} \xi) + h.o.t.,$$

$$1 - \psi_j(\xi) = \alpha_j \exp(\lambda_3^{(j)} \xi) - \gamma_j \xi \exp(\lambda_3^{(j)} \xi) - \beta_j s_2^{(j)} \exp(\lambda_4^{(j)} \xi) + h.o.t.$$

with $\gamma_j \geq 0$ and

$$\alpha_j \geq 0, \quad \beta_j > 0 \quad (\gamma_j = 0).$$

(v) $\lambda_5^{(j)} = \lambda_4^{(j)} = \lambda_3^{(j)}$: *there are β_j and γ_j such that*

(2.17)

$$\phi_j(\xi) = \beta_j \tau_2^{(j)} \exp(\lambda_3^{(j)} \xi) - \gamma_j \tau_2 \xi \exp(\lambda_3^{(j)} \xi) + h.o.t.,$$

$$1 - \psi_j(\xi) = -\beta_j \xi \exp(\lambda_3^{(j)} \xi) + \gamma_j (\xi^2/2) \exp(\lambda_3^{(j)} \xi) + h.o.t.$$

with $\gamma_j \geq 0$ and $\beta_j > 0$ ($\gamma_j = 0$).

The next lemma immediately follows from the above lemmas.

LEMMA 2.4. *There are positive constants $r_1, r_2, \eta_1, \eta_2, \omega_1, \omega_2$, and M such that*

$$(2.18) \quad r_1 \leq \frac{\phi'_j(\xi)}{1 - \phi_j(\xi)} \leq r_2 \quad (\xi \geq 0),$$

$$(2.19) \quad r_1 \leq \frac{1 - \phi_j(\xi)}{\psi_j(\xi)} \quad (\xi \geq 0),$$

$$(2.20) \quad r_1 \leq \frac{|\psi'_j(\xi)|}{\psi_j(\xi)} \quad (\xi \geq 0)$$

and

$$(2.21) \quad 0 < 1 - \psi_j(\xi) \leq M \exp(\omega_1 \xi) \quad (\xi \leq 0),$$

$$(2.22) \quad \eta_1 \leq \frac{|\psi'_j(\xi)|}{1 - \psi_j(\xi)} \leq \eta_2 \quad (\xi \leq 0),$$

$$(2.23) \quad 0 \leq \phi_j(\xi) \leq M \exp(\omega_2 \xi) \quad (\xi \leq 0),$$

$$(2.24) \quad \eta_1 \leq \frac{\phi'_j(\xi)}{\phi_j(\xi)} \leq \eta_2 \quad (\xi \leq 0).$$

The following lemma is clear from Lemmas 2.2 and 2.3.

LEMMA 2.5. *Inequality (1.8) holds provided that one of the following conditions is satisfied:*

- (1) $\lambda_4^{(j)} < \lambda_3^{(j)}$,
- (2) $\lambda_5^{(j)} < \lambda_3^{(j)} \leq \lambda_4^{(j)}$, $\gamma_j > 0$,
- (3) $\lambda_5^{(j)} \leq \lambda_3^{(j)} < \lambda_4^{(j)}$, $\alpha_j = \gamma_j = 0$,

where α_j, β_j , and γ_j are the numbers corresponding to each condition of (i)–(v) for $\lambda_k^{(j)}$ ($k = 3, 4, 5$) in Lemma 2.3.

Proof of Lemma 1.3. By Lemmas 2.2 and 2.3, the limit

$$\lim_{\xi \rightarrow -\infty} \frac{-\psi'_j(\xi)}{1 - \psi_j(\xi)} = \lim_{\xi \rightarrow -\infty} \frac{\psi''_j(\xi)}{\psi'_j(\xi)}$$

exists. Then the condition (1.13) implies that this limit is $\lambda_4^{(j)}$ or $\lambda_5^{(j)}$. Hence $\beta_j > 0$ or $\gamma_j > 0$. We assert that one of the conditions (1), (2), and (3) of Lemma 2.5 holds. This proves the first part of the assertion. The latter part is trivial by Lemma 2.5. \square

3. Key estimates. In this section we prove two lemmas which will play a crucial role in the argument for the proof of Theorem 1.1.

LEMMA 3.1. *Let (ϕ_j, ψ_j) be traveling wave solutions of (1.5)–(1.7) satisfying (1.8). Define*

$$(3.1) \quad A_1(z, p) := (1 - \phi_2(-z + p))\phi'_1(z + p) + (1 - \phi_1(z + p))\phi'_2(-z + p),$$

and put $\omega = \min\{\omega_1, \omega_2\}$. Then there exist $L_1, L_2 > 0$ such that, given $p < 0$,

$$(3.2) \quad \frac{\phi'_1(z + p)\phi'_2(-z + p)}{A_1(z, p)} \leq L_1 \exp(\omega p),$$

$$(3.3) \quad \frac{\phi_1(z + p)(1 - \phi_2(-z + p))\psi_1(z + p)(1 - \psi_2(-z + p))}{A_1(z, p)} \leq L_2 \exp(\omega p),$$

and

$$(3.4) \quad \frac{\phi_2(-z+p)(1-\phi_1(z+p))\psi_2(-z+p)(1-\psi_1(z+p))}{A_1(z,p)} \leq L_2 \exp(\omega p)$$

hold.

Proof. We first prove (3.2). With the aid of (2.23) and (2.24), for $z \leq 0$,

$$\begin{aligned} \frac{\phi_1'(z+p)\phi_2'(-z+p)}{A_1(z,p)} &\leq \frac{\phi_1'(z+p)}{1-\phi_1(z+p)} \\ &\leq \frac{\eta_2 M \exp(\omega_2(z+p))}{1-\phi_1(0)} \\ &\leq \frac{\eta_2 M \exp(\omega_2 p)}{1-\phi_1(0)}, \end{aligned}$$

while, for $z \geq 0$,

$$\begin{aligned} \frac{\phi_1'(z+p)\phi_2'(-z+p)}{A_1(z,p)} &\leq \frac{\phi_2'(-z+p)}{1-\phi_2(-z+p)} \\ &\leq \frac{\eta_2 M \exp(\omega_2(-z+p))}{1-\phi_2(0)} \\ &\leq \frac{\eta_2 M \exp(\omega_2 p)}{1-\phi_2(0)}. \end{aligned}$$

Thus we can take $L_1 = \eta_2 M \max\{1/(1-\phi_1(0)), 1/(1-\phi_2(0))\}$ to prove (3.2).

Next we show (3.3) and (3.4). Put

$$(3.5) \quad B_1(z,p) := \phi_1(z+p)\psi_1(z+p)(1-\phi_2(-z+p))(1-\psi_2(-z+p)),$$

$$(3.6) \quad \tilde{B}_1(z,p) := \phi_2(-z+p)\psi_2(-z+p)(1-\phi_1(z+p))(1-\psi_1(z+p)).$$

First we prove (3.3). Let $z \leq p$. By (2.18) we have

$$r_1 \leq \frac{\phi_2'(-z+p)}{1-\phi_2(-z+p)} \quad (z \leq p);$$

hence

$$(3.7) \quad \begin{aligned} \frac{B_1(z,p)}{A_1(z,p)} &\leq \frac{\phi_1(z+p)}{r_1(1-\phi_1(z+p))} \\ &\leq \frac{M \exp(\omega_2(z+p))}{r_1(1-\phi_1(0))} \\ &\leq \frac{M \exp(\omega_2 p)}{r_1(1-\phi_1(0))}. \end{aligned}$$

As for the case $p \leq z \leq 0$, noticing that

$$\eta_0 \eta_1 \leq \frac{\phi_2'(-z+p)}{1-\psi_2(-z+p)} \quad (p \leq z),$$

from (2.24) and (1.8) we have

$$\begin{aligned}
 \frac{B_1(z, p)}{A_1(z, p)} &\leq \frac{\phi_1(z + p)}{\eta_0 \eta_1 (1 - \phi_1(z + p))} \\
 &\leq \frac{M \exp(\omega_2(z + p))}{\eta_0 \eta_1 (1 - \phi_1(0))} \\
 (3.8) \qquad &\leq \frac{M \exp(\omega_2 p)}{\eta_0 \eta_1 (1 - \phi_1(0))}.
 \end{aligned}$$

We let $0 \leq z \leq -p$. With the aid of

$$\eta_1 \leq \frac{\phi_1'(z + p)}{\phi_1(z + p)} \quad (z \leq -p),$$

we obtain

$$\begin{aligned}
 \frac{B_1(z, p)}{A_1(z, p)} &\leq \frac{(1 - \phi_2(-z + p))(1 - \psi_2(-z + p))}{\eta_1 (1 - \phi_2(-z + p))} \\
 &\leq \frac{1 - \psi_2(-z + p)}{\eta_1} \\
 &\leq \frac{M \exp(\omega_1(-z + p))}{\eta_1} \\
 (3.9) \qquad &\leq \frac{M \exp(\omega_1 p)}{\eta_1}.
 \end{aligned}$$

Since

$$r_1^2 \leq \frac{\phi_1'(z + p)}{\psi_1(z + p)} \quad (-p \leq z),$$

by (2.18) and (2.19), we see that for $-p \leq z$,

$$\begin{aligned}
 \frac{B_1(z, p)}{A_1(z, p)} &\leq \frac{(1 - \phi_2(-z + p))(1 - \psi_2(-z + p))}{r_1^2 (1 - \phi_2(-z + p))} \\
 &\leq \frac{M \exp(\omega_1(-z + p))}{r_1^2} \\
 (3.10) \qquad &\leq \frac{M \exp(\omega_1 p)}{r_1^2}.
 \end{aligned}$$

Combining (3.7), (3.8), (3.9), and (3.10) leads to (3.3).

We go to the case (3.4). In the above argument for (3.5), we replace $\phi_1(z + p)$, $\phi_2(-z + p)$, $\psi_1(z + p)$, and $\psi_2(-z + p)$ by $\phi_2(-z + p)$, $\phi_1(z + p)$, $\psi_2(-z + p)$, and $\psi_1(z + p)$, respectively, and change z to $-z$. Then for (3.6) we similarly obtain the desired estimate. We will leave the details to the reader. \square

LEMMA 3.2. *Let*

$$(3.11) \qquad A_2(z, p) := \psi_2(-z + p)\psi_1'(z + p) + \psi_1(z + p)\psi_2'(-z + p) \quad (< 0).$$

Then there exist $L_3, L_4 > 0$ such that, given $p < 0$,

$$(3.12) \qquad \frac{\psi_1'(z + p)\psi_2'(-z + p)}{|A_2(z, p)|} \leq L_3 \exp(\omega p)$$

and

$$(3.13) \quad \frac{\psi_1(z+p)\psi_2(-z+p)(1-\psi_1(z+p))(1-\psi_2(-z+p))}{|A_2(z,p)|} \leq L_4 \exp(\omega p)$$

hold.

Proof. For $z \geq 0$, by (2.21) and (2.22) we have

$$\begin{aligned} \frac{\psi_1'(z+p)\psi_2'(-z+p)}{|A_2(z,p)|} &\leq \frac{|\psi_2'(-z+p)|}{\psi_2(-z+p)} \\ &\leq \frac{\eta_2 M \exp(\omega_1(-z+p))}{\psi_2(0)} \leq \frac{\eta_2 M \exp(\omega_1 p)}{\psi_2(0)}. \end{aligned}$$

Similarly, we obtain that, for $z \leq 0$,

$$\frac{\psi_1'(z+p)\psi_2'(-z+p)}{|A_2(z,p)|} \leq \frac{\eta_2 M \exp(\omega_1 p)}{\psi_1(0)}.$$

This proves (3.12).

We show (3.13). Put

$$(3.14) \quad B_2(z,p) := \psi_1(z+p)\psi_2(-z+p)(1-\psi_1(z+p))(1-\psi_2(-z+p)).$$

First let $0 \leq z \leq -p$. By (2.22),

$$\begin{aligned} \frac{B_2(z,p)}{|A_2(z,p)|} &\leq \frac{\psi_1(z+p)\psi_2(-z+p)(1-\psi_2(-z+p))}{\eta_1 \psi_2(-z+p)} \\ &\leq \frac{1-\psi_2(-z+p)}{\eta_1} \leq \frac{M \exp(\omega_1 p)}{\eta_1}. \end{aligned}$$

Next we consider the case $-p \leq z$. By (2.20)

$$\begin{aligned} \frac{B_2(z,p)}{|A_2(z,p)|} &\leq \frac{\psi_2(-z+p)(1-\psi_1(z+p))(1-\psi_2(-z+p))}{r_1 \psi_2(-z+p)} \\ &\leq \frac{1-\psi_2(-z+p)}{r_1} \leq \frac{M \exp(\omega_1 p)}{r_1}. \end{aligned}$$

As for $p \leq z \leq 0$ and $z \leq p$, we use the above argument similarly to obtain

$$\begin{aligned} \frac{B_2(z,p)}{|A_2(z,p)|} &\leq \frac{M \exp(\omega_1 p)}{\eta_1} \quad (p \leq z \leq 0), \\ \frac{B_2(z,p)}{|A_2(z,p)|} &\leq \frac{M \exp(\omega_1 p)}{r_1} \quad (z \leq p). \end{aligned}$$

This concludes the proof of (3.13). \square

4. Existence of the entire solution. Let $(\phi_j(\xi), \psi_j(\xi))$ ($j = 1, 2$) be the traveling wave solutions as in Theorem 1.1. To prove the theorem, we consider the combination of the solutions $(\phi_1(x+c_1t), \psi_1(x+c_1t))$ and $(\phi_2(-x+c_2t), \psi_2(-x+c_2t))$.

By the change of variables

$$(U(z,t), V(z,t)) = (u(x,t), v(x,t)), \quad z = x + \frac{c_1 - c_2}{2}t,$$

(1.1) is transformed into

$$(4.1) \quad \begin{cases} U_t = U_{zz} - \frac{c_1 - c_2}{2}U_z + U(1 - U - k_1V), \\ V_t = dV_{zz} - \frac{c_1 - c_2}{2}V_z + aV(1 - V - k_2U) \end{cases} \quad (x \in \mathbb{R}).$$

Then $(U, V) = (\phi_1(z + c_mt), \psi_1(z + c_mt))$ and $(\phi_2(-z + c_mt), \psi_2(-z + c_mt))$ give traveling wave solutions of (4.1) with speed $c_m = (c_1 + c_2)/2$ and the opposite speed $-c_m$, respectively. The latter has the profile such that the U component is monotone decreasing while the V component is monotone increasing.

We define supersolutions and subsolutions to (1.1) and (4.1). Put

$$\begin{cases} \mathcal{F}_1(u, v) := u_t - u_{xx} + \beta u_x - f(u, v), \\ \mathcal{F}_2(u, v) := v_t - dv_{xx} + \beta v_x - ag(u, v), \end{cases}$$

and consider the equations

$$(4.2) \quad \mathcal{F}_1(u, v) = 0, \quad \mathcal{F}_2(u, v) = 0,$$

where

$$\begin{aligned} f(u, v) &:= (1 - u - k_1v)u, & g(u, v) &:= (1 - v - k_2u)v, \\ \beta &:= \frac{c_1 - c_2}{2}. \end{aligned}$$

We call $(\underline{u}(x, t), \bar{v}(x, t)), (x, t) \in \mathbb{R} \times [T_1, T_2]$ a subsolution if

$$\mathcal{F}_1(\underline{u}, \bar{v}) \leq 0, \quad \mathcal{F}_2(\underline{u}, \bar{v}) \geq 0 \quad ((x, t) \in \mathbb{R} \times [T_1, T_2])$$

hold. We say that $(\bar{u}(x, t), \underline{v}(x, t)), (x, t) \in \mathbb{R} \times [T_1, T_2]$ is a supersolution to (4.2) if

$$\mathcal{F}_1(\bar{u}, \underline{v}) \geq 0, \quad \mathcal{F}_2(\bar{u}, \underline{v}) \leq 0 \quad ((x, t) \in \mathbb{R} \times [T_1, T_2]).$$

We note that any solutions $(u(x, t), v(x, t))$ of (4.2) are not only subsolutions but also supersolutions. Let $(u(x, t), v(x, t))$ be solutions to (4.2) in $t \in [T_1, T_2]$. The maximal principle for the competition-diffusion system yields that if

$$\begin{cases} \underline{u}(x, T_1) \leq u(x, T_1) \leq \bar{u}(x, T_1), \\ \underline{v}(x, T_1) \leq v(x, T_1) \leq \bar{v}(x, T_1) \end{cases}$$

hold, then these orders are preserved for all $t \in [T_1, T_2]$.

If $(u_1(x, t), v_1(x, t))$ and $(u_2(x, t), v_2(x, t))$ are subsolutions in $t \in (T_1, T_2)$, then the pairing of

$$\underline{u}(x, t) := \max_{x \in \mathbb{R}}\{u_1(x, t), u_2(x, t)\}, \quad \bar{v}(x, t) := \min_{x \in \mathbb{R}}\{v_1(x, t), v_2(x, t)\}$$

is a subsolution in $t \in (T_1, T_2)$. Similarly, if $(u_1(x, t), v_1(x, t))$ and $(u_2(x, t), v_2(x, t))$ are supersolutions, then the pairing of

$$\bar{u}(x, t) := \min_{x \in \mathbb{R}}\{u_1(x, t), u_2(x, t)\}, \quad \underline{v}(x, t) := \max_{x \in \mathbb{R}}\{v_1(x, t), v_2(x, t)\}$$

is a supersolution.

We define a subsolution of (4.1) by the pairing

$$(4.3) \quad \begin{cases} \underline{U}(z, t) := \max_{z \in \mathbb{R}} \{ \phi_1(z + c_m t + q), \phi_2(-z + c_m t + q) \}, \\ \underline{V}(z, t) := \min_{z \in \mathbb{R}} \{ \psi_1(z + c_m t + q), \psi_2(-z + c_m t + q) \}, \end{cases}$$

where q is an arbitrarily given number.

Next we construct a supersolution. We first notice the following identities:

$$\begin{aligned} & f(u_1 + u_2 - u_1 u_2, v_1 v_2) - (1 - u_2)f(u_1, v_1) - (1 - u_1)f(u_2, v_2) \\ &= -(1 - u_1)(1 - u_2)u_1 u_2 \\ &\quad - k_1 \{ (u_1 + u_2 - u_1 u_2)v_1 v_2 - (1 - u_2)u_1 v_1 - (1 - u_1)u_2 v_2 \}, \\ & g(u_1 + u_2 - u_1 u_2, v_1 v_2) - v_2 g(u_1, v_1) - v_1 g(u_2, v_2) \\ &= -v_1 v_2 (1 - v_1)(1 - v_2) + k_2 u_1 u_2 v_1 v_2. \end{aligned}$$

With the aid of (1.16), we obtain the inequalities

$$(4.4) \quad \begin{aligned} & f(u_1 + u_2 - u_1 u_2, v_1 v_2) - (1 - u_2)f(u_1, v_1) - (1 - u_1)f(u_2, v_2) \\ & \leq k_1 \{ u_1(1 - u_2)v_1(1 - v_2) + (1 - u_1)u_2(1 - v_1)v_2 \}, \end{aligned}$$

$$(4.5) \quad \begin{aligned} & g(u_1 + u_2 - u_1 u_2, v_1 v_2) - v_2 g(u_1, v_1) - v_1 g(u_2, v_2) \\ & \geq -v_1 v_2 (1 - v_1)(1 - v_2) \end{aligned}$$

for $0 \leq u_1, u_2 \leq 1$, and $0 \leq v_1, v_2$.

We introduce the following ordinary differential equation:

$$(4.6) \quad \begin{cases} \dot{p} = c_m + L e^{\omega p} & (t \leq 0), \\ p(0) = p_0 < 0. \end{cases}$$

A simple computation yields a solution to (4.6) as

$$(4.7) \quad p(t) = c_m t - \frac{1}{\omega} \log \left\{ e^{-\omega p_0} + \frac{L(1 - e^{\omega c_m t})}{c_m} \right\} < 0 \quad (t < 0)$$

having the asymptotics

$$(4.8) \quad \lim_{t \rightarrow -\infty} (p(t) - c_m t) = -\frac{1}{\omega} \log \left(e^{-\omega p_0} + \frac{L}{c_m} \right).$$

LEMMA 4.1. *Let $(\phi_j(x + c_j t), \psi_j(x + c_j t))$ ($j = 1, 2$) be traveling wave solutions to (1.1) satisfying (1.5), (1.6), (1.7), and (1.8). Let $p(t)$ be a solution to (4.6) with $L \geq \max\{2L_1 + 2k_1L_2, 2dL_3 + aL_4\}$. Then the pairing of*

$$(4.9) \quad \begin{cases} \overline{U}(z, t) := \phi_1(z + p(t)) + \phi_2(-z + p(t)) - \phi_1(z + p(t))\phi_2(-z + p(t)), \\ \underline{V}(z, t) := \psi_1(z + p(t))\psi_2(-z + p(t)) \end{cases}$$

is a supersolution of (4.1).

Proof. First note that $\phi_1 = \phi_1(\xi)$, $\psi_1 = \psi_1(\xi)$ ($\xi = z + c_m t$) satisfy

$$\begin{cases} \phi_1'' - c_1 \phi_1' + f(\phi_1, \psi_1) = 0, \\ d\psi_1'' - c_1 \psi_1' + ag(\phi_1, \psi_1) = 0, \end{cases}$$

while $\phi_2 = \phi_2(\xi)$, $\psi_2 = \psi_2(\xi)$ ($\xi = -z + c_m t$) satisfy

$$\begin{cases} \phi_2'' - c_2 \phi_2' + f(\phi_2, \psi_2) = 0, \\ d\psi_2'' - c_2 \psi_2' + ag(\phi_2, \psi_2) = 0, \end{cases}$$

where $' = d/d\xi$, $'' = d^2/d\xi^2$. Substituting (\bar{U}, \underline{V}) into

$$\begin{cases} \mathcal{F}_1(U, V) = U_t - U_{zz} + \frac{c_1 - c_2}{2} U_z - f(U, V), \\ \mathcal{F}_2(U, V) = V_t - dV_{zz} + \frac{c_1 - c_2}{2} V_z - ag(U, V) \end{cases}$$

yields

$$\begin{aligned} \mathcal{F}_1(\bar{U}, \underline{V}) &= \dot{p}\{(1 - \phi_2)\phi_1' + (1 - \phi_1)\phi_2'\} - (\phi_1'' + \phi_2'' - \phi_1''\phi_2 + 2\phi_1'\phi_2' - \phi_1\phi_2'') \\ &\quad + \frac{c_1 - c_2}{2}\{(1 - \phi_2)\phi_1' - (1 - \phi_1)\phi_2'\} - f(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) \\ &= \dot{p}\{(1 - \phi_2)\phi_1' + (1 - \phi_1)\phi_2'\} - (1 - \phi_2) \left(\phi_1'' - \frac{c_1 - c_2}{2} \phi_1' \right) \\ &\quad - (1 - \phi_1) \left(\phi_2'' + \frac{c_1 - c_2}{2} \phi_2' \right) - 2\phi_1'\phi_2' - f(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) \\ &= (\dot{p} - c_m)\{(1 - \phi_2)\phi_1' + (1 - \phi_1)\phi_2'\} - 2\phi_1'\phi_2' \\ &\quad - \{f(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) - (1 - \phi_2)f(\phi_1, \psi_1) - (1 - \phi_1)f(\phi_2, \psi_2)\}. \end{aligned}$$

Using (4.4), we obtain

$$\begin{aligned} \mathcal{F}_1(\bar{U}, \underline{V}) &\geq A_1(z, p)(\dot{p} - c_m) - 2\phi_1'\phi_2' - k_1(B_1(z, p) + \tilde{B}_1(z, p)) \\ (4.10) \quad &= A_1(z, p) \left\{ L \exp(\omega p) - \frac{2\phi_1'\phi_2'}{A_1(z, p)} - k_1 \frac{B_1(z, p) + \tilde{B}_1(z, p)}{A_1(z, p)} \right\}, \end{aligned}$$

where $A_1(z, p)$, $B_1(z, p)$, and $\tilde{B}_1(z, p)$ are defined in (3.1), (3.5), and (3.6), respectively. By Lemma 3.1 we obtain

$$\mathcal{F}_1(\bar{U}, \underline{V}) \geq 0 \quad ((z, t) \in \mathbb{R} \times (-\infty, 0])$$

for $L \geq 2L_1 + 2k_1L_2$.

Similarly,

$$\begin{aligned} \mathcal{F}_2(\bar{U}, \underline{V}) &= (\psi_1'\psi_2 + \psi_1\psi_2')\dot{p} - d(\psi_1''\psi_2 + \psi_1\psi_2'') + 2d\psi_1'\psi_2' \\ &\quad + \frac{c_1 - c_2}{2}(\psi_1'\psi_2 - \psi_1\psi_2') - ag(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) \\ &= (\psi_1'\psi_2 + \psi_1\psi_2')\dot{p} + \psi_2\{-c_1\psi_1' + ag(\phi_1, \psi_1)\} \\ &\quad + \psi_1\{-c_2\psi_2' + ag(\phi_2, \psi_2)\} + 2d\psi_1'\psi_2' \\ &\quad + \frac{c_1 - c_2}{2}(\psi_1'\psi_2 - \psi_1\psi_2') - ag(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) \\ &= (\psi_1'\psi_2 + \psi_1\psi_2')(\dot{p} - c_m) + 2d\psi_1'\psi_2' \\ &\quad - a\{g(\phi_1 + \phi_2 - \phi_1\phi_2, \psi_1\psi_2) - \psi_2g(\phi_1, \psi_1) - \psi_1g(\phi_2, \psi_2)\}. \end{aligned}$$

Applying (4.5) to the right-hand side of this equality yields

$$(4.11) \quad \mathcal{F}_2(\overline{U}, \underline{V}) \leq A_2(z, p) \left\{ L \exp(\omega p) + \frac{2d\psi'_1\psi'_2}{A_2(z, p)} + \frac{aB_2(z, p)}{A_2(z, p)} \right\},$$

where $A_2(z, p)$ and $B_2(z, p)$ are defined as in (3.11) and (3.14), respectively (recall $A_2(z, p) < 0$). By Lemma 3.2 we obtain

$$\mathcal{F}_2(\overline{U}, \underline{V}) \leq 0 \quad ((z, t) \in \mathbb{R} \times (-\infty, 0])$$

for $L \geq 2dL_3 + aL_4$. This implies that $(\overline{U}, \underline{V})$ is a supersolution. \square

Proof of Theorem 1.1. We use the subsolution (4.3) with

$$q = -\frac{1}{\omega} \log \left(e^{-\omega p_0} + \frac{L}{c_m} \right)$$

and the supersolution (4.9). For this q we have

$$\begin{aligned} \lim_{t \rightarrow -\infty} |p(t) - (c_m + q)t| &= 0, \\ \begin{cases} \underline{U}(z, t) < \overline{U}(z, t), \\ \underline{V}(z, t) < \overline{V}(z, t) \end{cases} & \quad (t < 0, z \in \mathbb{R}), \end{aligned}$$

and

$$\lim_{t \rightarrow -\infty} [\sup_{z \in \mathbb{R}} (\overline{U}(z, t) - \underline{U}(z, t)) + \sup_{z \in \mathbb{R}} (\overline{V}(z, t) - \underline{V}(z, t))] = 0.$$

Applying the same argument in [2] (see also [1]) with the aid of the comparison theorem yields that there is a solution $(U(z, t), V(z, t))$ satisfying

$$(4.12) \quad \begin{cases} \underline{U}(z, t) \leq U(z, t) \leq \overline{U}(z, t), \\ \underline{V}(z, t) \leq V(z, t) \leq \overline{V}(z, t) \end{cases} \quad ((z, t) \in \mathbb{R} \times (-\infty, 0]).$$

Since (1.1) (or (4.1)) has the invariance with time shift, we obtain the asymptotic behavior as $t \rightarrow -\infty$ in the statement in Theorem 1.1.

On the other hand, to prove the asymptotic behavior as $t \rightarrow \infty$ of the theorem, we recall that $\underline{U}(z, t), \overline{V}(z, t)$ in (4.3) are defined every t . By the asymptotic behavior

$$\lim_{t \rightarrow \infty} \sup_{z \in \mathbb{R}} (1 - \underline{U}(z, t)) = 0, \quad \lim_{t \rightarrow \infty} \sup_{z \in \mathbb{R}} \overline{V}(z, t) = 0$$

and the fact that the solution can be continued for positive t , we obtain the desired behavior as $t \rightarrow \infty$. \square

Proof of Corollary 1.2. Since the equation is invariant under phase shift and time shift, the desired result follows from Theorem 1.1. For the details of the argument, see [5]. \square

Proof of Proposition 1.4. A straightforward computation yields

$$\begin{aligned} \mathcal{F}_1(u_*, v^*) &= -\frac{k_1 \exp(2\omega_0 t)}{(1 + \Phi)^3} < 0, \\ \mathcal{F}_2(u_*, v^*) &= \frac{a \exp(2\omega_0 t)}{(1 + \Phi)^4} > 0. \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathcal{F}_1(u^*, v_*) &= \frac{\Psi}{(1 + \Psi)^2} \left\{ \dot{p} - \omega_0 - \frac{k_1 \exp(2p)}{\Psi(1 + \Psi)} \right\} \\ &\geq \frac{\Psi}{(1 + \Psi)^2} \left\{ \frac{L_0 \exp(p)}{1 + \exp(p)} - \frac{k_1 \exp(p)}{1 + \exp(p)} \right\} \geq 0, \\ \mathcal{F}_2(u^*, v_*) &= -\frac{2\Psi}{(1 + \Psi)^3} \left\{ \dot{p} - \omega_0 - \frac{a \exp(2p)}{2\Psi(1 + \Psi)} \right\} \\ &\leq -\frac{2\Psi}{(1 + \Psi)^3} \left\{ \frac{L_0 \exp(p)}{1 + \exp(p)} - \frac{a \exp(p)}{2(1 + \exp(p))} \right\} \leq 0, \end{aligned}$$

where we used $\Psi(x, t) \geq \Psi(0, t) = \exp(p)$. Hence the first assertion of the proposition is true.

By integration we obtain

$$p(t) - \omega_0 t = \frac{L_0}{L_0 + \omega_0} \log\{\omega_0 + (L_0 + \omega_0) \exp(p(t))\} + q,$$

where

$$q := p_0 - \frac{L_0}{L_0 + \omega_0} \log\{\omega_0 + (L_0 + \omega_0) \exp(p_0)\}.$$

As $t \rightarrow -\infty$,

$$p(t) - \omega_0 t \rightarrow \frac{L_0}{L_0 + \omega_0} \log \omega_0 + q.$$

Hence, if we take p_0 as a solution of

$$\frac{L_0}{L_0 + \omega_0} \log \omega_0 + q = 0,$$

namely, a solution of (1.26), then

$$\lim_{t \rightarrow -\infty} (p(t) - \omega_0 t) = 0.$$

This yields

$$u_*(x, t) < u^*(x, t), \quad v_*(x, t) < v^*(x, t) \quad ((x, t) \in \mathbb{R} \times \mathbb{R}),$$

and

$$\lim_{t \rightarrow -\infty} [\sup_{x \in \mathbb{R}} (u^*(x, t) - u_*(x, t)) + \sup_{x \in \mathbb{R}} (v^*(x, t) - v_*(x, t))] = 0.$$

The existence of the entire solution satisfying (1.25) is proven by the same argument found in [2] (see also [5] or [1]), which is left to the reader. \square

Appendix. In this section we prove Lemmas 2.1, 2.2, and 2.3. We investigate the asymptotic behavior as $\xi \rightarrow \pm\infty$. For simplicity we will not specify the dependence on j in the notation as long as there is no confusion.

We put

$$(u, w, v, y) = (\phi_j(\xi), \phi'_j(\xi), \psi_j(\xi), \psi'_j(\xi)).$$

Then the equations of (1.5) are written as

$$(A.1) \quad \begin{pmatrix} u' \\ w' \\ v' \\ y' \end{pmatrix} = \begin{pmatrix} w \\ c_j w - f(u, v) \\ y \\ (c_j/d)y - (a/d)g(u, v) \end{pmatrix},$$

where

$$f(u, v) := u(1 - u - k_1 v), \quad g(u, v) := v(1 - v - k_2 u).$$

Proof of Lemma 2.1. We first consider the linearized equations of (A.1) at $(u, w, v, y) = (1, 0, 0, 0)$:

$$(A.2) \quad \begin{pmatrix} W_1' \\ W_2' \\ Y_1' \\ Y_2' \end{pmatrix} = A_1 \begin{pmatrix} W_1 \\ W_2 \\ Y_1 \\ Y_2 \end{pmatrix},$$

$$(A.3) \quad A_1 := \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & c_j & k_1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -(a/d)(1 - k_2) & c_j/d \end{pmatrix}.$$

All the eigenvalues of A_1 consist of roots to (2.1) and (2.2). There is a negative eigenvalue satisfying (2.1) and a corresponding eigenvector, which are given by

$$(A.4) \quad \lambda_1 = -\frac{\sqrt{c_j^2 + 4} - c_j}{2}, \quad \begin{pmatrix} 1 \\ \lambda_1 \\ 0 \\ 0 \end{pmatrix}.$$

We recall $k_2 > 1$. Since (2.2) has a negative root, A_1 has another negative eigenvalue and a corresponding eigenvector

$$(A.5) \quad \lambda_2^{(j)} = -\frac{\sqrt{c_j^2 + 4ad(k_2 - 1)} - c_j}{2d}, \quad \begin{pmatrix} s_1 \\ s_1 \lambda_2 \\ 1 \\ \lambda_2 \end{pmatrix}, \quad s_1 := \frac{k_1}{\lambda_2^2 - c_j \lambda_2 - 1}.$$

When $\lambda_1 = \lambda_2$, A_1 has a generalized eigenvector

$$\begin{pmatrix} 0 \\ 1 \\ \tau_1 \\ \lambda_1 \tau_1 \end{pmatrix}, \quad \tau_1 := \frac{2\lambda_1 - c_j}{k_1}.$$

Thus every solution of (A.2), which converges to zeros as $\xi \rightarrow \infty$, is given by

$$\begin{pmatrix} W_1 \\ W_2 \\ Y_1 \\ Y_2 \end{pmatrix} = C_1 \begin{pmatrix} 1 \\ \lambda_1 \\ 0 \\ 0 \end{pmatrix} e^{\lambda_1 \xi} + C_2 \begin{pmatrix} s_1 \\ s_1 \lambda_2 \\ 1 \\ \lambda_2 \end{pmatrix} e^{\lambda_2 \xi} \quad (\lambda_1 \neq \lambda_2)$$

or

$$\begin{pmatrix} W_1 \\ W_2 \\ Y_1 \\ Y_2 \end{pmatrix} = C_1 \begin{pmatrix} 1 \\ \lambda_1 \\ 0 \\ 0 \end{pmatrix} e^{\lambda_1 \xi} + C_2 \left\{ \begin{pmatrix} 1 \\ \lambda_1 \\ 0 \\ 0 \end{pmatrix} \xi + \begin{pmatrix} 0 \\ 1 \\ \tau_1 \\ \lambda_1 \tau_1 \end{pmatrix} \right\} e^{\lambda_1 \xi} \quad (\lambda_1 = \lambda_2),$$

where C_1, C_2 are arbitrarily given real numbers.

Applying the stable manifold theorem to (A.1) yields that as $\xi \rightarrow \infty$, there are α and β such that

$$(A.6) \quad \begin{cases} \phi(\xi) = 1 - \alpha \exp(\lambda_1 \xi) + \beta s_1 \exp(\lambda_2 \xi) + h.o.t., \\ \psi(\xi) = \beta \exp(\lambda_2 \xi) + h.o.t. \end{cases}$$

for $\lambda_1 \neq \lambda_2$, while there are $\tilde{\alpha}$ and $\tilde{\beta}$ such that

$$(A.7) \quad \begin{cases} \phi(\xi) = 1 - \tilde{\alpha} \exp(\lambda_1 \xi) - \tilde{\beta} \xi \exp(\lambda_1 \xi) + h.o.t., \\ \psi(\xi) = -\tilde{\beta} \tau_1 \exp(\lambda_1 \xi) + h.o.t. \end{cases}$$

for $\lambda_1 = \lambda_2$.

We show $\beta \neq 0$. Equation (A.1) allows a solution of the form $(u, w, 0, 0)$, which corresponds to a solution to

$$(A.8) \quad u' = w, \quad w' = cw - f(u, 0) = cw - u(1 - u).$$

There is a stable manifold for the equilibrium $(u, w) = (1, 0)$ of (A.8). The stable manifold for $(u, w, v, y) = (1, 0, 0, 0)$ contains this dynamics that is obtained by putting $\beta = 0$ in (A.6). By this observation we can assert that $\beta \neq 0$ in (A.6) for the present traveling wave.

Since λ_1 is a negative root of (2.1), we notice that

$$s_1 > 0 \quad (\lambda_2 < \lambda_1 < 0), \quad s_1 < 0 \quad (\lambda_1 < \lambda_2 < 0).$$

Recall that $\phi(\xi)$ and $\psi(\xi)$ are monotone increasing and decreasing, respectively. Thus for $\lambda_2 < \lambda_1$, $\alpha > 0$ and $\beta > 0$ follow from (A.6), while for $\lambda_1 < \lambda_2$, $\beta > 0$ follows.

When $\lambda_1 = \lambda_2$, we can let $\tilde{\beta} \neq 0$ in the same argument above. Since $\tau_1 < 0$ holds, $\tilde{\beta}$ must be positive in (A.7). This concludes the proof of the lemma. \square

Next we consider the linearized equation at $(u, w, v, y) = (0, 0, 1, 0)$:

$$(A.9) \quad \mathbf{W}' = A_2 \mathbf{W},$$

$$(A.10) \quad A_2 := \begin{pmatrix} 0 & 1 & 0 & 0 \\ k_1 - 1 & c_j & 0 & 0 \\ 0 & 0 & 0 & 1 \\ ak_2/d & 0 & a/d & c_j/d \end{pmatrix}, \quad \mathbf{W} := \begin{pmatrix} W_1 \\ W_2 \\ Y_1 \\ Y_2 \end{pmatrix}.$$

In this case eigenvalues of A_2 are given by the roots of (1.9) and (1.10). We easily see from (1.9) that there is a positive eigenvalue and a corresponding eigenvector

$$(A.11) \quad \lambda_3 = \frac{c_j + \sqrt{c_j^2 + 4ad}}{2d}, \quad \begin{pmatrix} 0 \\ 0 \\ 1 \\ \lambda_3 \end{pmatrix}.$$

Proof of Lemma 2.2. This corresponds to the bistable case (ii) stated in the introduction. In addition to the positive eigenvalue λ_3 , there is a positive eigenvalue

$$(A.12) \quad \lambda_4 = \frac{c_j + \sqrt{c_j^2 + 4(k_1 - 1)}}{2}, \quad \begin{pmatrix} 1 \\ \lambda_4 \\ s_2 \\ s_2\lambda_4 \end{pmatrix}, \quad s_2 := \frac{ak_2}{d\lambda_4^2 - c_j\lambda_4 - a}$$

unless λ_4 satisfies (1.9). When $\lambda_3 = \lambda_4$, we have a generalized eigenvector

$$\begin{pmatrix} \tau_2 \\ \tau_2\lambda_3 \\ 0 \\ 1 \end{pmatrix}, \quad \tau_2 := \frac{2d\lambda_3 - c_j}{ak_2}.$$

Hence every solution which converges to zeros as $\xi \rightarrow -\infty$ is given by

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2\mathbf{p}_4e^{\lambda_4\xi}$$

or

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2(\mathbf{p}_1\xi + \mathbf{p}_2)e^{\lambda_3\xi},$$

where we put

$$\mathbf{p}_1 := \begin{pmatrix} 0 \\ 0 \\ 1 \\ \lambda_3 \end{pmatrix}, \quad \mathbf{p}_2 := \begin{pmatrix} \tau_2 \\ \tau_2\lambda_3 \\ 0 \\ 1 \end{pmatrix}, \quad \mathbf{p}_4 := \begin{pmatrix} 1 \\ \lambda_4 \\ s_2 \\ s_2\lambda_4 \end{pmatrix}.$$

Applying the unstable manifold theorem yields that as $\xi \rightarrow -\infty$, there exist α and β such that

$$(A.13) \quad \begin{cases} \phi_j(\xi) = \beta \exp(\lambda_4\xi) + h.o.t., \\ \psi_j(\xi) = 1 - \alpha \exp(\lambda_3\xi) + \beta s_2 \exp(\lambda_4\xi) + h.o.t. \end{cases}$$

for $\lambda_3 \neq \lambda_4$, while there are $\tilde{\alpha}$ and $\tilde{\beta}$ such that

$$(A.14) \quad \begin{cases} \phi_j(\xi) = \tilde{\beta}\tau_2 \exp(\lambda_3\xi) + h.o.t., \\ \psi_j(\xi) = 1 - \tilde{\alpha} \exp(\lambda_3\xi) + \tilde{\beta}\xi \exp(\lambda_3\xi) + h.o.t. \end{cases}$$

for $\lambda_3 = \lambda_4$. In this case we have an unstable manifold for the equilibrium $(v, y) = (1, 0)$ of

$$v' = y, \quad y' = c_jy - ag(0, v),$$

which corresponds to $\beta = 0$ and $\tilde{\beta} = 0$ in (A.13) and (A.14), respectively. Thus we may assume that $\beta \neq 0$, $\tilde{\beta} \neq 0$. By a similar argument applying to the behavior as $\xi \rightarrow \infty$, we obtain the assertion in Lemma 2.2. \square

Proof of Lemma 2.3. This corresponds to the monostable case (i) in the introduction. There is a family of traveling waves for $c_j \geq c_{min}^{(j)} \geq 2\sqrt{1 - k_1}$.

In addition to λ_3 and λ_4 , we have a positive eigenvalue and a corresponding eigenvector,

$$(A.15) \quad \lambda_5 = \frac{c_j - \sqrt{c_j^2 - 4(1 - k_1)}}{2}, \quad \begin{pmatrix} 1 \\ \lambda_5 \\ s_3 \\ s_3\lambda_5 \end{pmatrix}, \quad s_3 := \frac{ak_2}{d\lambda_5^2 - c_j\lambda_5 - a},$$

unless λ_5 satisfies (1.9). We note that $\lambda_5 \leq \lambda_4$. If $\lambda_5 < \lambda_4$ and $\lambda_4, \lambda_5 \neq \lambda_3$, any solution converging to zeros as $\xi \rightarrow -\infty$ is given as

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2\mathbf{p}_4e^{\lambda_4\xi} + C_3\mathbf{p}_5e^{\lambda_5\xi},$$

where

$$\mathbf{p}_5 := \begin{pmatrix} 1 \\ \lambda_5 \\ s_3 \\ s_3\lambda_5 \end{pmatrix}.$$

Applying the unstable manifold theorem, we obtain

$$\begin{cases} \phi_j(\xi) = \beta \exp(\lambda_4\xi) + \gamma \exp(\lambda_5\xi) + h.o.t., \\ \psi_j(\xi) = 1 - \alpha \exp(\lambda_3\xi) + \beta s_2 \exp(\lambda_4\xi) + \gamma s_3 \exp(\lambda_5\xi) + h.o.t. \end{cases}$$

By the same reason as in the previous case, we can assert $(\beta, \gamma) \neq (0, 0)$. This leads to (2.12).

When $\lambda_5 = \lambda_4 < \lambda_3$, we obtain

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2(\mathbf{p}_4\xi + \tilde{\mathbf{p}}_4)e^{\lambda_4\xi} + C_3\mathbf{p}_4e^{\lambda_4\xi},$$

where

$$\tilde{\mathbf{p}}_4 := \begin{pmatrix} 0 \\ 1 \\ \tau_4 \\ \tau_5 \end{pmatrix}, \quad \begin{pmatrix} \tau_4 \\ \tau_5 \end{pmatrix} := \frac{s_2}{d\lambda_4^2 - c_j\lambda_4 - a} \begin{pmatrix} c_j - 2d\lambda_4 \\ -d\lambda_4^2 - a \end{pmatrix} \quad (\lambda_4 = \lambda_5 = c_j/2).$$

Applying the unstable manifold theorem to this case yields

$$\begin{cases} \phi_j(\xi) = \beta \exp(\lambda_4\xi) + \gamma \exp(\lambda_4\xi) + h.o.t., \\ \psi_j(\xi) = 1 - \alpha \exp(\lambda_3\xi) + \beta(s_2\xi + \tau_4) \exp(\lambda_4\xi) + \gamma s_2 \exp(\lambda_4\xi) + h.o.t., \end{cases}$$

where we again see that $(\beta, \gamma) \neq (0, 0)$. This leads to (2.14).

For the cases $\lambda_5 < \lambda_3 = \lambda_4$ and $\lambda_5 = \lambda_3 < \lambda_4$ all the solutions converging to zeros as $\xi \rightarrow -\infty$ are given by

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2(\mathbf{p}_1\xi + \mathbf{p}_2)e^{\lambda_3\xi} + C_3\mathbf{p}_5e^{\lambda_5\xi}$$

and

$$\mathbf{W}(\xi) = C_1\mathbf{p}_1e^{\lambda_3\xi} + C_2(\mathbf{p}_1\xi + \mathbf{p}_2)e^{\lambda_3\xi} + C_3\mathbf{p}_4e^{\lambda_4\xi},$$

respectively. When $\lambda_3 = \lambda_4 = \lambda_5 = c_j/2$, we obtain one more generalized eigenvector so that the solution is given by

$$\mathbf{W}(\xi) = C_1 \mathbf{p}_1 e^{\lambda_3 \xi} + C_2 (\mathbf{p}_1 \xi + \mathbf{p}_2) e^{\lambda_3 \xi} + C_3 (\mathbf{p}_1 \xi^2/2 + \mathbf{p}_2 \xi + \mathbf{p}_3) e^{\lambda_3 \xi},$$

where

$$\mathbf{p}_3 := \begin{pmatrix} d/ak_2 \\ \lambda_3(3d-2)/ak_2 \\ 0 \\ 0 \end{pmatrix}.$$

By those solutions it is not so difficult to prove the remaining cases (iii), (iv), and (v) of Lemma 2.3. Completing the proof is left to the reader. \square

Acknowledgments. The authors would like to express their sincere thanks to Professor Yukio Kan-on for valuable comments on the asymptotic behavior of the monotone traveling wave to the Lotka–Volterra competition-diffusion equations. Their thanks also go to the referees for the useful comments on the revision of the paper.

REFERENCES

- [1] X. CHEN AND J.-S. GUO, *Existence and uniqueness of entire solutions for a reaction-diffusion equation*, J. Differential Equations, 212 (2005), pp. 62–84.
- [2] Y. FUKAO, Y. MORITA, AND H. NINOMIYA, *Some entire solutions of the Allen-Cahn equation*, Taiwanese J. Math., 8 (2004), pp. 15–32.
- [3] R. A. GARDNER, *Existence and stability of travelling wave solutions of competition models: A degree theoretic approach*, J. Differential Equations, 44 (1982), pp. 343–364.
- [4] R. A. GARDNER AND C. K. R. T. JONES, *Stability of travelling wave solutions of diffusive predator-prey systems*, Trans. Amer. Math. Soc., 327 (1991), pp. 465–524.
- [5] J.-S. GUO AND Y. MORITA, *Entire solutions of reaction-diffusion equations and an application to discrete diffusive equations*, Discrete Contin. Dynam. Systems, 12 (2005), pp. 193–212.
- [6] F. HAMEL AND N. NADIRASHVILI, *Entire solutions of the KPP equation*, Comm. Pure Appl. Math., 52 (1999), pp. 1255–1276.
- [7] Y. HOSONO, *Singular perturbation analysis of travelling waves for diffusive Lotka-Volterra competition models*, in Numerical and Applied Mathematics, Part II (Paris, 1988), IMACS Ann. Comput. Appl. Math. 1.2, Baltzer, Basel, 1989, pp. 687–692.
- [8] Y. KAN-ON, *Parameter dependence of propagation speed of travelling waves for competition-diffusion equations*, SIAM J. Math. Anal., 26 (1995), pp. 340–363.
- [9] Y. KAN-ON, *Existence of standing waves for competition-diffusion equations*, Japan J. Indust. Appl. Math., 13 (1996), pp. 117–133.
- [10] Y. KAN-ON, *Fisher wave fronts for the Lotka-Volterra competition model with diffusion*, Nonlinear Anal., 28 (1997), pp. 145–164.
- [11] Y. KAN-ON, *Instability of stationary solutions for a Lotka-Volterra competition model with diffusion*, J. Math. Anal. Appl., 208 (1997), pp. 158–170.
- [12] Y. KAN-ON AND Q. FANG, *Stability of monotone travelling waves for competition-diffusion equations*, Japan J. Indust. Appl. Math., 13 (1996), pp. 343–349.
- [13] Y. MORITA AND H. NINOMIYA, *Entire solutions with merging fronts to reaction-diffusion equations*, J. Dynam. Differential Equations, 18 (2006), pp. 841–861.
- [14] M. RODRIGO AND M. MIMURA, *Exact solutions of a competition-diffusion system*, Hiroshima Math. J., 30 (2000), pp. 257–270.
- [15] M. M. TANG AND P. C. FIFE, *Propagating fronts for competing species equations with diffusion*, Arch. Ration. Mech. Anal., 73 (1980), pp. 69–77.
- [16] H. YAGISITA, *Backward global solutions characterizing annihilation dynamics of travelling fronts*, Publ. Res. Inst. Math. Sci., 39 (2003), pp. 117–164.

HIGHER ORDER APPROXIMATIONS IN THE HEAT EQUATION AND THE TRUNCATED MOMENT PROBLEM*

YONG JUNG KIM[†] AND WEI-MING NI[‡]

Abstract. In this paper, we employ linear combinations of n heat kernels to approximate solutions to the heat equation. We show that such approximations are of order $O(t^{(\frac{1}{2p} - \frac{2n+1}{2})})$ in L^p -norm, $1 \leq p \leq \infty$, as $t \rightarrow \infty$. For positive solutions of the heat equation such approximations are achieved using the theory of truncated moment problems. For general sign-changing solutions these type of approximations are obtained by simply adding an auxiliary heat kernel. Furthermore, inspired by numerical computations, we conjecture that such approximations converge geometrically as $n \rightarrow \infty$ for any fixed $t > 0$.

Key words. heat equation, moments, asymptotics convergence rates, approximation of an integral formula, heat kernel

AMS subject classifications. 35K05, 78M05, 41A10

DOI. 10.1137/08071778X

1. Introduction. It is well known that

$$(1.1) \quad u(x, t) = \int \frac{u_0(c)}{\sqrt{4\pi t}} e^{-\frac{(x-c)^2}{4t}} dc$$

is the physically meaningful solution to the heat equation

$$(1.2) \quad u_t = u_{xx}, \quad u(x, 0) = u_0(x) \in L^1(\mathbf{R}), \quad x, u \in \mathbf{R}, \quad t > 0,$$

where, for simplicity, the initial value $u_0(x)$ is assumed to be continuous. In this paper, we shall refer to (1.1) as the solution of the heat equation (1.2) for the sake of brevity. If a general L^1 initial value is considered, no asymptotic convergence order to a fundamental solution is expected in L^1 -norm. Hence, the asymptotic convergence order is usually studied under suitable restrictions on its initial value $u_0(x)$ for $|x|$ large.

Since the analysis of this paper is based on the *moments* of the solution, the initial value $u_0(x)$ is required to have finite moments up to certain order, say, $2n$. We set $x^{2n}u_0(x) \in L^1(\mathbf{R})$ and the moments of the initial value $u_0(x)$ as

$$(1.3) \quad \gamma_k := \int x^k u_0(x) dx < \infty, \quad k = 0, 1, \dots, 2n.$$

For example, if the initial value has an algebraic decay order higher than $2n + 1$ for $|x|$ large, i.e., for $\varepsilon > 0$,

$$(1.4) \quad u_0(x) = O\left(|x|^{-(2n+1+\varepsilon)}\right) \quad \text{as } |x| \rightarrow \infty,$$

*Received by the editors March 4, 2008; accepted for publication (in revised form) September 23, 2008; published electronically February 20, 2009. This research was supported in part by the National Science Foundation.

<http://www.siam.org/journals/sima/40-6/71778.html>

[†]Department of Mathematical Sciences, KAIST, 335 Gwahangno, Yuseong-gu, Daejeon, 305-701, Republic of Korea (yongkim@kaist.edu). This author was supported by the Korea Science and Engineering Foundation (KOSEF) grant R01-2007-000-11307-0 funded by the Korean government (MOST).

[‡]School of Mathematics, University of Minnesota, 206 Church St. S.E., Minneapolis, MN 55455 (ni@math.umn.edu).

then the moments are well defined up to order $2n$. In the study of asymptotics the initial value is frequently assumed to have the order that a fundamental solution has for $|x|$ large. For the heat equation case the fundamental solution is the Gaussian and the corresponding decay order is $u_0(x) = O(e^{-x^2})$ as $|x| \rightarrow \infty$. Hence, the moment γ_k is defined for all order $k \geq 0$.

One may do the integration in the explicit formula (1.1) only approximately, even though the integration gives the exact value of the solution. In numerical computations finding an efficient way to compute such an integration has been an important issue. From this point of view, it seems useful to consider its approximation in a simpler form. Duoandikoetxea and Zuazua [9] showed that the following linear combination of derivatives of the Gaussian

$$(1.5) \quad \psi_{2n}(x, t) \equiv \sum_{i=0}^{2n-1} \frac{(-1)^i \gamma_i}{(i!) \sqrt{4\pi t}} \partial_x^i \left(e^{-\frac{x^2}{4t}} \right)$$

approaches to the solution u with a convergence order of

$$(1.6) \quad \|u(t) - \psi_{2n}(t)\|_p = O\left(t^{\left(\frac{1}{2p} - \frac{2n+1}{2}\right)}\right) \quad \text{as } t \rightarrow \infty \text{ for } 1 \leq p \leq \infty,$$

where $\|\cdot\|_p$ denotes the L^p -norm in the whole space \mathbf{R} and ∂_x^i the i th order partial differentiation with respect to x . Note that the original multidimensional result is written in a one-dimensional (1D) version for an easier comparison. This asymptotic convergence order indicates that ψ_{2n} is a good approximation of the solution $u(x, t)$ for t large. However, it does not necessarily mean that ψ_{2n} is a good approximation as $n \rightarrow \infty$ with a fixed $t > 0$. In fact, Table 7.3 shows that this L^p -norm difference may diverge geometrically as $n \rightarrow \infty$ if the fixed time $t > 0$ is not large enough. This is not surprising since the high order derivatives of the Gaussian in (1.5) diverge as their orders increase.

In this article we consider a linear combination of “ n ” heat kernels

$$(1.7) \quad \phi_n(x, t) \equiv \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-\frac{(x-c_i)^2}{4t}}$$

as an approximation to the solution $u(x, t)$. One may regard this summation as a discrete version of the integration in (1.1) by considering c_i 's as grid points and ρ_i 's as approximations of $u_0(c)dc$ in the interval (c_{i-1}, c_i) . However, we employ these $2n$ degrees of freedom, ρ_i 's and c_i 's, to match the first $2n$ initial moments, i.e., to satisfy the following $2n$ moment equations:

$$(1.8) \quad \lim_{t \rightarrow 0} \int x^k \phi_n(x, t) dx = \gamma_k, \quad k = 0, 1, \dots, 2n - 1.$$

If the initial value is positive, the theory of truncated moment problems [3] gives the solvability of this problem. Then ϕ_n and u share identical first $2n$ moments for all $t \geq 0$. Note that ψ_{2n} in (1.5) also has the same property, and Duoandikoetxea and Zuazua obtained the convergence order in (1.6) based on it. Hence, we may obtain the same convergence order for the approximation $\phi_n(x, t)$. In this paper, we actually go a little bit further and obtain the limit of $t^{\frac{2n+1}{2} - \frac{1}{2p}} \|u(t) - \phi_n(t)\|_p$ as $t \rightarrow \infty$. This convergence order is then improved in Lemma 2.3 for the case that this limit becomes zero. A multidimensional extension of this approach requires a theory of

multidimensional truncated moment problems. One may find one from a recent work by Curto and Fialkow [4].

From a practical point of view, it is desirable if the solution u can be approximated by ϕ_n as $n \rightarrow \infty$ for a fixed $t > 0$. Indeed, our numerical examples in section 7.2 indicate the following geometric convergence order:

$$(1.9) \quad \frac{\|u(t) - \phi_n(t)\|_\infty}{\|u(t) - \phi_{n+1}(t)\|_\infty} \rightarrow 1 + 4\frac{t}{v} \quad \text{as } n \rightarrow \infty,$$

where the constant $v > 0$ depends on the initial value $u_0(x)$. However, its proof has, thus far, eluded us; nevertheless, we will include a discussion of (1.9) in section 6.

This paper is organized as follows. First, in section 2, we compute the limit of $t^{(\frac{2n+1}{2} - \frac{1}{2p})} \|u(t) - \phi_n(t)\|_p$ as $t \rightarrow \infty$ under the assumption (1.8), which gives the convergence order in (1.6). A short introduction to the theory of truncated moment problems is given in section 3, which provides the existence and the uniqueness of ρ_i 's and c_i 's that solve (1.8). We remark that the theory is applicable for nonnegative initial values only (see [1, 3]). For general sign-changing solutions the existence and the uniqueness of such ρ_i 's and c_i 's do not hold. In section 4 we discuss this issue in detail for three cases with $n = 1, 2$, and 3. In section 5 we construct approximations for general sign-changing cases by adding an auxiliary heat kernel or by assigning c_i 's independently. The conjectured geometric convergence order for large $n > 0$ is discussed in section 6. The asymptotic convergence orders as $t \rightarrow \infty$ or $n \rightarrow \infty$ are numerically tested in sections 7.1 and 7.2. The convergence of the alternative approach using ψ_{2n} and the conjectured statements in section 6 are numerically tested in sections 7.3 and 7.4.

In the study of nonlinear diffusion or convection, fundamental solutions which have the Dirac measure as their initial value, i.e., $u_0(x) = \delta(x)$, often serve as canonical solutions. The Barenblatt solutions, the diffusion waves, and the N-waves are well-known examples (see [23]). In the study of porous medium equations, the Barenblatt solution is used as an asymptotic profile and the convergence order of general solutions to this special one has been studied in various cases (see [2, 6] and references therein). The diffusion wave and the Gaussian are the asymptotics of convection-diffusion equations for diffusion dominant cases (see [10, 11, 12, 16]). For convection dominant cases (see [14]) and inviscid convection equations (see [5, 8, 13, 20]) or hyperbolic systems (see [7, 17, 18]), N-waves represent the asymptotic behavior, where N-waves can be understood as a special solution with initial value $u_0(x) = \lim_{\varepsilon \rightarrow 0} [a\delta(x - \varepsilon) - b\delta(x + \varepsilon)]$ with $a, b > 0$. Placing the Dirac measure at the center of mass, the optimal convergence order of $O(t^{\frac{1}{2p} - \frac{3}{2}})$ in L^p -norm (or of $O(t^{-1})$ in L^1 -norm) has been obtained in several cases (see [2, 13, 15]). Therefore, the result of this paper can be viewed as an extreme case that exploits all of the moments of the initial value.

The approach in this paper can be directly employed to approximate the solutions to the Burgers equation via the Cole–Hopf transformation. To obtain the rigorous convergence order for the Burgers case it is required to check the well definedness of the transformed solutions as is done in Lemmas 3.1–3.2 and Theorem 3.3 in [15] for the special case $n = 1$. Considering that the Burgers equation has been used as a tool to study the asymptotic structure of the viscous systems of conservation laws (see, e.g., [19]), we hope the approach in this article may be useful for other general models.

2. Asymptotic convergence order. In this section, we show that the decay rate of a derivative of a solution is naturally transferred to the convergence order

of our approximation. This connection will be made by assigning the moments of a solution to its approximation. Let $\gamma_k(t)$ be the k th order moment of a solution $u(x, t)$ at time $t \geq 0$, i.e.,

$$\gamma_k(t) = \int x^k u(x, t) dx, \quad k = 0, 1, 2, \dots, t \geq 0.$$

(Notice that we are slightly abusing the notation γ_k in (1.3) in the following couple of paragraphs.) We can easily show how the moment $\gamma_k(t)$ evolves as $t \rightarrow \infty$.

LEMMA 2.1. *Let $u(x, t)$ be the solution to the heat equation and $\gamma_k(t)$ be its k th order moment at time $t \geq 0$. Then*

$$\frac{d}{dt} \gamma_k(t) = \begin{cases} 0, & k = 0 \text{ or } 1, \\ k(k-1)\gamma_{k-2}(t), & k \geq 2. \end{cases}$$

Proof. For $k = 0$, the lemma is equivalent to the conservation of mass. For $k = 1$, since $u_t = u_{xx}$, the integration by parts gives

$$\gamma'_1(t) = \int x u_t dx = \int x u_{xx} dx = [x u_x - u]_{-\infty}^{\infty} = 0.$$

Similarly, for $k \geq 2$, we obtain

$$\begin{aligned} \gamma'_k(t) &= \int x^k u_t dx = \int x^k u_{xx} dx \\ &= [x^k u_x - k x^{k-1} u]_{-\infty}^{\infty} + \int k(k-1) x^{k-2} u dx = k(k-1) \gamma_{k-2}(t). \quad \square \end{aligned}$$

This lemma shows that even numbered moments and odd numbered ones evolve independently. One may explicitly write

$$(2.1) \quad \begin{aligned} \gamma_{2n}(t) &= \sum_{k=0}^n \frac{(2n)!}{(n-k)!(2k)!} t^{n-k} \gamma_{2k}(0), \\ \gamma_{2n+1}(t) &= \sum_{k=0}^n \frac{(2n+1)!}{(n-k)!(2k+1)!} t^{n-k} \gamma_{2k+1}(0). \end{aligned}$$

If $\gamma_k(0) = 0$ for all $0 \leq k \leq n$, then $\gamma_k(t) = 0$ for all $0 \leq k \leq n$, $\gamma_k(t) = \gamma_k(0)$ for $k = n+1, n+2$, $\gamma_k(t)$ is linear for $k = n+3, n+4$, $\gamma_k(t)$ is quadratic for $k = n+5, n+6$, and so on.

Let $v(x, t)$ be an approximation solution of the exact one $u(x, t)$. Since the difference $E(x, t) = v(x, t) - u(x, t)$ is also a solution to the heat equation, the moments of $E(x, t)$ will be always zero up to certain order if they are initially zero. Hence, it is natural to expect a higher convergence order by matching the moments of the approximation solution to those of the exact one. We proceed with our discussion in this respect.

LEMMA 2.2. *If $x^m E_0(x) \in L^1(\mathbf{R})$ and*

$$(2.2) \quad \int_{-\infty}^{\infty} x^k E_0(x) dx = 0 \text{ for all } 0 \leq k < m,$$

then there exists $E_m \in W^{m,1}(\mathbf{R})$ such that

$$(2.3) \quad \partial_x^m E_m(x) = E_0(x).$$

Proof. The proof was given by Duoandikoetxea and Zuazua [9] for the multidimensional case. Here we provide its 1D counterpart. Consider a sequence of functions defined inductively by

$$(2.4) \quad E_k(x) = \int_{-\infty}^x E_{k-1}(y)dy, \quad 0 < k \leq m.$$

First, we show that E_k 's are well defined,

$$(2.5) \quad \int_{-\infty}^{\infty} E_k(x)dx = 0 \quad \text{and} \quad E_k(x) \rightarrow 0 \quad \text{as} \quad |x| \rightarrow \infty$$

for $k = 0, 1, \dots, m - 1$. It suffices to show (2.5) for $k = l < m$ under the assumption that (2.5) holds for all $k = 0, 1, \dots, l - 1$. Note that it is clearly satisfied for $k = 0$. Since $\int_{-\infty}^{\infty} E_{l-1}(x)dx = 0$, the integral $E_l(x)$ also decays to zero as $|x| \rightarrow \infty$. Using the fact that E_k decays to zero as $|x| \rightarrow \infty$ for all $0 \leq k \leq l$, we obtain

$$\int_{-\infty}^{\infty} E_l(x)dx = (-1)^l \int_{-\infty}^{\infty} \frac{x^l}{l!} E_0(x)dx = 0$$

using the integration by parts and then (2.2). Therefore, (2.5) holds for $k = l$ and, hence, for all $0 \leq k \leq m - 1$. Since $x^m E_0(x) \in L^1(\mathbf{R})$, $E_m \in W^{m,1}(\mathbf{R})$ and (2.3) is satisfied. \square

The existence of E_m satisfying (2.3) is the key observation to obtain the asymptotic convergence order. We now continue our discussion under the assumption that the initial value $E_0(x)$ satisfies (2.2) and $E_m(x)$ is its m th order antiderivative given in Lemma 2.2. However, the following discussions about the decay rate of derivatives of a solution can be considered independently. Let $E_m(x, t)$ be the solution to the heat equation with initial value $E_m(x, 0) = E_m(x) \in W^{m,1}(\mathbf{R})$, i.e.,

$$E_m(x, t) = \frac{1}{\sqrt{4\pi t}} \int e^{-(x-y)^2/(4t)} E_m(y)dy.$$

The dissipation of the solution can be easily shown by introducing similarity variables:

$$\xi = \frac{x}{\sqrt{t}}, \quad \zeta = \frac{y}{\sqrt{t}}, \quad \tilde{E}_m(\xi, t) = E_m(x, t).$$

Then $E_m(x, t)$ is transformed to

$$\tilde{E}_m(\xi, t) = \frac{1}{\sqrt{4\pi}} \int e^{-(\xi-\zeta)^2/4} E_m(\sqrt{t}\zeta)d\zeta,$$

and its m th order derivative is given by

$$\partial_{\xi}^m \tilde{E}_m(\xi, t) = \partial_x^m E_m(x, t)(\partial_{\xi} x)^m = \partial_x^m E_m(x, t)(\sqrt{t})^m.$$

Now consider the decay order of the m th order derivative of the solution $E_m(x, t)$. First, let $C_m := |\int E_m(y)dy|$ and consider the case $C_m \neq 0$. Then

$$(2.6) \quad (\sqrt{t})^{m+1} |\partial_x^m E_m(x, t)| = \sqrt{t} \left| \partial_{\xi}^m \tilde{E}_m(\xi, t) \right| = \frac{C_m}{\sqrt{4\pi}} \left| \int f(\zeta)g_t(\xi - \zeta)d\zeta \right|,$$

where

$$(2.7) \quad g_t(\xi) = \sqrt{t} E_m(\sqrt{t}\xi)/C_m, \quad f(\xi) = \partial_{\xi}^m \left(e^{-\xi^2/4} \right).$$

After taking the supremum on both sides of (2.6), one obtains that

$$(\sqrt{t})^{m+1} \|\partial_x^m E_m(t)\|_\infty \leq \frac{C_m}{\sqrt{4\pi}} \left\| \partial_\xi^m \left(e^{-\xi^2/4} \right) \right\|_\infty.$$

If one takes $t \rightarrow \infty$ limit to (2.6), then

$$\lim_{t \rightarrow \infty} (\sqrt{t})^{m+1} |\partial_x^m E_m(x, t)| = \frac{C_m}{\sqrt{4\pi}} |f(\xi)|.$$

Therefore, after taking the supremum on both sides again, we obtain

$$\lim_{t \rightarrow \infty} (\sqrt{t})^{m+1} \|\partial_x^m E_m(t)\|_\infty = \frac{C_m}{\sqrt{4\pi}} \left\| \partial_\xi^m \left(e^{-\xi^2/4} \right) \right\|_\infty.$$

On the other hand, if $1 \leq p < \infty$, then

$$\begin{aligned} & t^{(\frac{m+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m(t)\|_p \\ &= (\sqrt{t})^{m+1} \left(\frac{1}{\sqrt{t}} \right)^{1/p} \left(\int |\partial_x^m E_m(x, t)|^p dx \right)^{1/p} \\ (2.8) \quad &= \left(\int |(\sqrt{t})^{m+1} \partial_x^m E_m(x, t)|^p d\left(\frac{x}{\sqrt{t}}\right) \right)^{1/p} \\ &= \left(\int |\sqrt{t} \partial_\xi^m \tilde{E}_m(\xi, t)|^p d\xi \right)^{1/p} \\ &= \frac{C_m}{\sqrt{4\pi}} \left(\int \left| \int f(\zeta) g_t(\xi - \zeta) d\zeta \right|^p d\xi \right)^{1/p} = \frac{C_m}{\sqrt{4\pi}} \|f * g_t\|_p. \end{aligned}$$

Standard arguments imply that $\|f * g_t\|_p \rightarrow \|f\|_p$ as $t \rightarrow \infty$ (see [21, p. 62]). Therefore,

$$\lim_{t \rightarrow \infty} t^{(\frac{m+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m(t)\|_p = \frac{C_m}{\sqrt{4\pi}} \left\| \partial_\xi^m \left(e^{-\xi^2/4} \right) \right\|_p.$$

Now we consider the case that $C_m = 0$. Then one can easily show that this limit is zero. In fact, we will improve the convergence order by working with higher order antiderivatives. Let

$$(2.9) \quad E_k(x) = \int_{-\infty}^x E_{k-1}(y) dy, \quad k > m.$$

We can easily show that $\int_{-\infty}^\infty E_{k_0}(x) dx (= \lim_{x \rightarrow \infty} E_{k_0+1}(x)) \neq 0$ for some $k_0 > m$. Suppose that $\int_{-\infty}^\infty E_k(x) dx = 0$ for all $k > m$. Then $|E_k(x)|$ decays to zero for $|x|$ large, and, therefore, after integrating by parts k times with proper inductive arguments, one obtains

$$(-1)^k k! \int_{-\infty}^\infty E_k(x) dx = \int_{-\infty}^\infty x^k E_0(x) dx = 0.$$

On the other hand, by the Weierstrass approximation theorem, there exists a sequence of polynomials P_n such that

$$P_n(x) \rightarrow E_0(x) \quad \text{as} \quad n \rightarrow \infty$$

uniformly on any bounded domain $[-L, L]$. Therefore, we obtain

$$\|E_0\|_2^2 = \int_{-L}^L E_0^2(x)dx = \lim_{n \rightarrow \infty} \int_{-L}^L P_n(x)E_0(x)dx = 0.$$

Hence, if the initial value E_0 is not a trivial one, there exists $k_0 > m$ such that $\lim_{x \rightarrow \infty} E_k(x) = 0$ for all $0 \leq k \leq k_0$ and $C_{k_0} := |\lim_{x \rightarrow \infty} E_{k_0+1}(x)| \neq 0$. If $E_0(x)$ decays with an algebraic order $k_0 + 1 + \varepsilon$, $\varepsilon > 0$, for $|x|$ large, then $C_{k_0} < \infty$. However, C_{k_0} can be unbounded in general.

Let $E_{k_0}(x, t)$ be the solution with $E_{k_0}(x)$ as its initial value. Then, clearly, $\partial_x^m E_m = \partial_x^{k_0} E_{k_0} = E_0$ and, hence,

$$\lim_{t \rightarrow \infty} t^{(\frac{k_0+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m\|_p = \lim_{t \rightarrow \infty} t^{(\frac{k_0+1}{2} - \frac{1}{2p})} \|\partial_x^{k_0} E_{k_0}\|_p = \frac{C_{k_0}}{\sqrt{4\pi}} \left\| \partial_\xi^{k_0} e^{-\frac{\xi^2}{4}} \right\|_p.$$

Therefore, if $C_m := |\int E_m(y)dy| = 0$, one obtains a higher decay order. Summing up, we obtain the following lemma.

LEMMA 2.3. *Let $E_m(x, t)$ be the solution to the heat equation with a nontrivial initial value $E_m(x) \in W^{m,1}(\mathbf{R})$ and E_k 's be given inductively by (2.9). Then there exists $k_0 \geq m$ such that $\lim_{x \rightarrow \infty} E_k(x) = 0$ for $0 \leq k \leq k_0$ and $0 \neq |\lim_{x \rightarrow \infty} E_{k_0+1}(x)|$, and, for $m \leq k \leq k_0$,*

$$(2.10) \quad \lim_{t \rightarrow \infty} t^{(\frac{k+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m(t)\|_p = \frac{\left\| \partial_\xi^k e^{-\xi^2/4} \right\|_p}{\sqrt{4\pi}} \left| \int E_k(x)dx \right|, \quad 1 \leq p \leq \infty.$$

If $\int E_m(x)dx = 0$, then the limit in (2.10) implies that $\lim_{t \rightarrow \infty} t^{(\frac{m+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m\|_p = 0$. Hence, we may simply say that

$$(2.11) \quad \lim_{t \rightarrow \infty} t^{(\frac{m+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m(t)\|_p = \frac{\left\| \partial_\xi^m (e^{-\xi^2/4}) \right\|_p}{\sqrt{4\pi}} \left| \int E_m(x)dx \right|, \quad 1 \leq p \leq \infty,$$

which is a weaker statement than (2.10) is. Note that one may obtain the upper bound of the term $t^{(\frac{m+1}{2} - \frac{1}{2p})} \|\partial_x^m E_m(t)\|_p$ using Young's inequality. In fact, the corresponding upper bound for the estimate ψ_{2n} was obtained in [9].

In the following, we take the convergence order in (2.11) for simplicity. If an optimal convergence order is concerned and $\int E_m(x)dx = 0$, then one may refer to (2.10). It is well known that an L^1 solution to the heat equation decays to zero with order $O(t^{-1/2})$. Lemma 2.3 says that the decay order of its derivative is increased by $\frac{1}{2}$ after each differentiation. The asymptotic convergence order between two solutions is now obtained as a corollary of previous lemmas.

THEOREM 2.4. *Let $u(x, t)$ and $v(x, t)$ be solutions of the heat equation with initial values $u_0(x)$ and $v_0(x)$, respectively. Suppose that the initial difference $E_0(x) := u_0(x) - v_0(x)$ satisfies the assumptions in Lemma 2.2. Then, for $1 \leq p \leq \infty$,*

$$(2.12) \quad \lim_{t \rightarrow \infty} t^{(\frac{m+1}{2} - \frac{1}{2p})} \|u(t) - v(t)\|_p = \frac{\left\| \partial_\xi^m (e^{-\frac{1}{4}\xi^2}) \right\|_p}{\sqrt{4\pi}} \left| \int E_m(x)dx \right|,$$

where $E_m \in W^{m,1}(\mathbf{R})$ is the one that satisfies $\partial_x^m E_m(x) = E_0(x)$.

Proof. Let $E_m(x, t)$ be the solution to the heat equation with initial value $E_m(x)$. Then $\partial_x^m E_m(x, t)$ is the solution to the heat equation with initial value $\partial_x^m E_m(x) = E_0(x)$. Hence, $\partial_x^m E_m(x, t) (= E(x, t)) = u(x, t) - v(x, t)$ and (2.12) follows from (2.11). \square

Remark 2.5. In this section, we basically considered the convergence order of $\phi_n(x, t) (\cong v(x, t))$ to $u(x, t)$ as $t \rightarrow \infty$ with a fixed $n > 0$. However, the relation (2.6), for example, provides certain convergence information as $n \rightarrow \infty$ with a fixed $t > 0$, too. To obtain a convergence order as $n \rightarrow \infty$ we need to specify $E_m(x)$ corresponding to our approximation $\phi_n(x, t)$, which will be considered in section 6.

3. Positive solutions and truncated moment problems. Consider a linear combination of heat kernels

$$(3.1) \quad \phi_n(x, t) := \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-(x-c_i)^2/(4t)}.$$

The $2n$ freedom of choices in ρ_i 's and c_i 's are used to control the first $2n$ moments of the approximation. Remember that γ_k is to denote the initial k th moment, i.e.,

$$(3.2) \quad \gamma_k := \int x^k u_0(x) dx, \quad k = 0, 1, \dots, 2n - 1.$$

Let \mathbf{r}_k be a column n -vector and \mathbf{A} be the $n \times n$ Hankel matrix given by

$$(3.3) \quad \begin{aligned} \mathbf{r}_k &= (\gamma_k, \gamma_{k+1}, \dots, \gamma_{k+n-1})^t, & k &= 0, 1, \dots, n, \\ \mathbf{A} &\equiv (a_{ij}) = (\gamma_{i+j}), & i, j &= 0, 1, \dots, n - 1. \end{aligned}$$

Since $\phi_n(x, t) \rightarrow \sum_{i=1}^n \rho_i \delta_{c_i}(x)$ as $t \rightarrow 0$, the difference between the initial value and its approximation is

$$E_0(x) := u_0(x) - \sum_{i=1}^n \rho_i \delta_{c_i}(x),$$

where $\delta_{c_i}(x)$ is the Dirac measure centered at c_i , i.e., $\delta_{c_i}(x) = \delta(x - c_i)$. Hence, the zero moment conditions in (2.2) can be written as

$$(3.4) \quad \int \sum_{i=1}^n x^k \rho_i \delta_{c_i}(x) dx = \int x^k u_0(x) dx (\equiv \gamma_k), \quad 0 \leq k \leq 2n - 1,$$

or, in a matrix form, as

$$(3.5) \quad \begin{pmatrix} 1 & \cdots & 1 \\ c_1 & \cdots & c_n \\ \cdot & \cdots & \cdot \\ \cdot & \cdots & \cdot \\ \cdot & \cdots & \cdot \\ c_1^{2n-1} & \cdots & c_n^{2n-1} \end{pmatrix} \begin{pmatrix} \rho_1 \\ \cdot \\ \cdot \\ \rho_n \end{pmatrix} = \begin{pmatrix} \gamma_0 \\ \gamma_1 \\ \cdot \\ \cdot \\ \gamma_{2n-1} \end{pmatrix}.$$

After eliminating all ρ_i 's (see section 4.3), one may obtain n -equations involving c_i 's only:

$$(3.6) \quad \mathbf{A}\Psi = \mathbf{r}_n,$$

where the column vector $\Psi = (\psi_0, \dots, \psi_{n-1})^t$ is given by

$$(3.7) \quad \psi_0 = (-1)^{n+1} \prod_{i=1}^n c_i, \quad \psi_1 = (-1)^n \sum_{j=1}^n \prod_{i \neq j} c_i, \dots, \psi_{n-1} = \sum_{i=1}^n c_i.$$

Consequently, we set

$$(3.8) \quad g_n(x) := x^n - \sum_{j=0}^{n-1} \psi_j x^j = (x - c_1)(x - c_2) \cdots (x - c_n).$$

(Note that the coefficient of the leading order term is 1 and, hence, $g_n(x) \rightarrow \infty$ as $x \rightarrow \infty$.) Hence, if the initial moments in (3.4) are satisfied, then c_i 's are zero points of the polynomial $g_n(x)$, where its coefficients are given as a solution of (3.6).

To show the existence and the uniqueness of the approximation we should show that the Hankel matrix in (3.6) is nonsingular. Then there exists a unique column vector $\Psi = (\psi_0, \dots, \psi_{n-1})^t$ that satisfies (3.6). The next thing to show is that the polynomial $g_n(x)$ in (3.8) has n distinct real zeros $c_1 < \dots < c_n$. Then ρ_i 's are given by solving the Vandermonde given by the first n -equations in (3.5), i.e.,

$$(3.9) \quad \begin{pmatrix} 1 & 1 & \cdots & 1 \\ c_1 & c_2 & \cdots & c_n \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ c_1^{n-1} & c_2^{n-1} & \cdots & c_n^{n-1} \end{pmatrix} \begin{pmatrix} \rho_1 \\ \rho_2 \\ \cdot \\ \cdot \\ \cdot \\ \rho_n \end{pmatrix} = \begin{pmatrix} \gamma_0 \\ \gamma_1 \\ \cdot \\ \cdot \\ \cdot \\ \gamma_{n-1} \end{pmatrix}.$$

It is well known that the Vandermonde matrix is nonsingular if c_i 's are all different. Then we can easily check that c_i 's and ρ_i 's also satisfy the last n -equations in (3.5).

For a general sign-changing initial value $u_0(x)$ the Hankel matrix \mathbf{A} can be singular, and examples are given in section 4. However, if the initial value $u_0(x)$ is nonnegative, then the uniqueness and the existence are resolved by the theory for the moment problem (see [1, 3]). In the following, we assume $u_0(x) \geq 0$ and introduce this technique briefly for the completeness and the later use in this paper. Consider

$$\Psi^t \mathbf{A} \Psi = \sum_{i,j=0}^{n-1} \psi_i \psi_j \gamma_{i+j} = \int \sum_{i,j=0}^{n-1} \psi_i x^i \psi_j x^j u_0(x) dx = \int \left(\sum_{k=0}^{n-1} \psi_k x^k \right)^2 u_0(x) dx.$$

Since the integrand $(\sum_{k=0}^{n-1} \psi_k x^k)^2 u_0(x)$ is nonnegative, we have $\Psi^t \mathbf{A} \Psi \geq 0$. Furthermore, $\Psi^t \mathbf{A} \Psi = 0$ if and only if $(\sum_{k=0}^{n-1} \psi_k x^k)^2 u_0(x) = 0$ for all $x \in \mathbf{R}$. For $\Psi \neq 0$, the polynomial $\sum_{k=0}^{n-1} \psi_k x^k$ has at most $n - 1$ zeros and, therefore, $\Psi^t \mathbf{A} \Psi > 0$ if the support of the initial value u_0 consists of at least n points. Hence, we may conclude that the Hankel matrix $\mathbf{A} \equiv (\gamma_{i+j})_{i,j=0}^{n-1}$ is nonsingular. (The proof is originally done by Hamburger.)

To show that $g_n(x)$ has n -distinct real zeros, consider a linear functional S on the space of polynomials defined by

$$S(r) := \sum_{i=0}^l r_i \gamma_i = r_0 \gamma_0 + \cdots + r_l \gamma_l \quad \text{for } r(x) = \sum_{i=0}^l r_i x^i.$$

Then, the same statements used for the positivity of $\Psi^t \mathbf{A} \Psi$ also show that

$$S(r^2) = S \left(\sum_{i,j=0}^l r_i r_j x^{i+j} \right) = \sum_{i,j=0}^l r_i r_j \gamma_{i+j} > 0.$$

Suppose that $r(x) \geq 0$. Then the degree of the polynomial $r(x)$ is even and there exist two polynomials p, q such that $r(x) = p^2(x) + q^2(x)$ (see [1, p. 2]). So $S(r) = S(p^2) + S(q^2) > 0$.

Since \mathbf{A} is nonsingular, there exists an n -vector $\Psi = (\psi_0, \dots, \psi_{n-1})$ uniquely so that $\mathbf{A}\Psi = \mathbf{r}_n$, i.e.,

$$\sum_{j=0}^{n-1} \psi_j \mathbf{r}_j = \mathbf{r}_n$$

or

$$(3.10) \quad \gamma_{n+k} - \sum_{j=0}^{n-1} \psi_j \gamma_{j+k} = 0, \quad k = 0, 1, \dots, n-1.$$

Considering the polynomial $g_n(x)$ and the definition of the functional $S(r)$, we can easily check that (3.10) implies

$$(3.11) \quad S(g_n x^k) = 0, \quad k = 0, 1, \dots, n-1.$$

Suppose that $g_n(x)$ never changes its sign. Then $g_n(x) \geq 0$ and, hence, $S(g_n) > 0$, which contradicts (3.11) with $k = 0$. Suppose that $g_n(x)$ changes its sign at points c_1, \dots, c_l only. Then $g_n(x)(x - c_1) \cdots (x - c_l) \geq 0$ and $S(g_n(x)(x - c_1) \cdots (x - c_l)) > 0$. On the other hand, if $l < n$, then the linearity of the functional $S(r)$ together with (3.11) implies that $S(g_n(x)(x - c_1) \cdots (x - c_l)) = 0$. Hence, we obtain that $g_n(x)$ has n -distinct real roots, say, $c_1 < \dots < c_n$.

Now we show that there exist ρ_i 's that solve (3.5) in a unique way, i.e.,

$$(3.12) \quad \sum_{i=1}^n \rho_i c_i^l = \gamma_l, \quad l = 0, 1, \dots, 2n-1.$$

Since c_i 's are all different, there exists a unique solution for the Vandermonde (3.9); i.e., (3.12) is satisfied for all $0 \leq l < n$. Now we complete the proof using inductive arguments. Let $0 \leq k \leq n-1$. We will show that the identity in (3.12) holds for $l = n+k$ under the assumption that it holds for all $0 \leq l < n+k$. First, observe that, since c_i 's are zero points of $x^k g_n(x)$, $k \geq 0$,

$$c_i^{n+k} = \sum_{j=0}^{n-1} \psi_j c_i^{j+k} \quad \text{for any } 1 \leq i \leq n, k \geq 0.$$

Using the relations (3.10) and (3.12) for $l < n+k$, we obtain

$$\gamma_{n+k} = \sum_{j=0}^{n-1} \psi_j \gamma_{j+k} = \sum_{j=0}^{n-1} \psi_j \sum_{i=1}^n \rho_i c_i^{j+k} = \sum_{i=1}^n \rho_i \sum_{j=0}^{n-1} \psi_j c_i^{j+k} = \sum_{i=1}^n \rho_i c_i^{n+k}.$$

Hence, (3.12) holds by the induction.

In summary, the proof of the existence and the uniqueness of the solution to the problem (3.5) consists of three steps. The invertibility of the Hankel matrix \mathbf{A} in (3.6) and the existence of n -distinct real roots c_i 's of g_n are the first two. The latter depends on the positive definiteness of the matrix \mathbf{A} which is easily proved for positive initial value $u_0(x)$. On the other hand, after obtaining c_i 's, finding ρ_i 's that satisfy (3.5) does not require the positivity. It depends only on the recursive structure of the problem. The following theorem is now clear from Theorem 2.4.

THEOREM 3.1. *Let $u(x, t)$ be the solution to the heat equation with initial value $u_0(x)$. If $u_0(x)$ is nonnegative (or nonpositive) and $x^{2n}u_0(x) \in L^1(\mathbf{R})$, then there exist $\rho_i, c_i, i = 1, \dots, n$, such that, for $\phi_n(x, t) \equiv \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-(x-c_i)^2/(4t)}$,*

$$(3.13) \quad \lim_{t \rightarrow \infty} t^{\frac{2n+1}{2} - \frac{1}{2p}} \|u(t) - \phi_n(t)\|_p = \frac{\left\| \partial_\xi^{2n} \left(e^{-\frac{1}{4}\xi^2} \right) \right\|_p}{\sqrt{4\pi}} \left| \int E_{2n}(x) dx \right|,$$

where $1 \leq p \leq \infty$ and $E_{2n}(x) \in W^{2n,1}(\mathbf{R})$ is the $2n$ th order antiderivative of $E_0(x) = u(x, 0) - \phi_n(x, 0)$. Furthermore, such a function $\phi_n(x, t)$ is unique.

Remark 3.2. The system (3.5) can be solved by commercial software such as Maple. However, since the problem is highly nonlinear, it takes a very long time even for small n . Therefore, even for the computational purpose, one needs to follow the steps of the proof to construct $\phi_n(x)$.

4. General initial value. In this section, we consider a general initial value which may change its sign. Then the existence and the uniqueness theory of the previous section is not applicable since it is for positive solutions only. In this section, we observe that the existence and uniqueness may fail for a general solution.

4.1. Approximation with a single heat kernel. For the case $n = 1$, the approximation $\phi_1(x, t) = \frac{\rho_1}{\sqrt{4\pi t}} e^{-(x-c_1)^2/(4t)}$ is obtained by solving

$$(4.1) \quad \rho_1 = \gamma_0, \quad c_1 \rho_1 = \gamma_1.$$

If $\gamma_0 \neq 0$, c_1 is uniquely decided by $c_1 = \gamma_1/\gamma_0$; i.e., c_1 is the *center of the mass* of the initial mass distribution u_0 . The convergence order in Theorem 2.4 is written as

$$\lim_{t \rightarrow \infty} t^{\left(\frac{3}{2} - \frac{1}{2p}\right)} \|u(t) - \phi_1(t)\|_p = \frac{\left\| \partial_\xi^2 \left(e^{-\frac{1}{4}\xi^2} \right) \right\|_p}{\sqrt{4\pi}} \left| \int_{-\infty}^{\infty} E_2(x) dx \right|, \quad 1 \leq p \leq \infty,$$

where $E_2(x)$ is the second order antiderivative of the initial error $E_0(x) := u_0(x) - \rho_1 \delta_{c_1}(x)$ given by (2.4), i.e., $E_2(x) = \int_{-\infty}^x \int_{-\infty}^y (u_0(z) - \rho_1 \delta_{c_1}(z)) dz dy$.

Now consider the singular case $\gamma_0 = 0$. Then the approximation is simply $\phi_1 \equiv 0$. If $\gamma_1 = 0$, then the equation for the first moment is satisfied for any $c_1 \in \mathbf{R}$ and we obtain the above convergence order which is equivalent to the decay rate $u(x, t)$. If $\gamma_1 \neq 0$, (4.1) has no solution and we do not obtain a single heat kernel approximation ϕ_1 with the desirable convergence order $O(t^{\left(\frac{1}{2p} - \frac{3}{2}\right)})$ for t large.

4.2. Approximation with two heat kernels. The double heat kernel solution $\phi_2(x, t) = \sum_{i=1}^2 \frac{\rho_i}{\sqrt{4\pi t}} e^{-(x-c_i)^2/(4t)}$ that approximates the solution $u(x, t)$ is obtained by solving

$$(4.2) \quad \begin{aligned} \rho_1 + \rho_2 &= \gamma_0, & \rho_1 c_1 + \rho_2 c_2 &= \gamma_1, \\ \rho_1 c_1^2 + \rho_2 c_2^2 &= \gamma_2, & \rho_1 c_1^3 + \rho_2 c_2^3 &= \gamma_3. \end{aligned}$$

We may simplify the equation by eliminating ρ_i 's and obtain two equations of the form $\mathbf{A}\Psi = \mathbf{r}_2$, i.e.,

$$\begin{pmatrix} \gamma_0 & \gamma_1 \\ \gamma_1 & \gamma_2 \end{pmatrix} \begin{pmatrix} \psi_0 \\ \psi_1 \end{pmatrix} = \begin{pmatrix} \gamma_2 \\ \gamma_3 \end{pmatrix}.$$

First, we need to check the invertibility of the Hankel matrix. Its determinant is the variance of the initial value u_0 if it is a probability distribution, i.e.,

$$|\mathbf{A}| = \gamma_0\gamma_2 - \gamma_1^2.$$

If $|\mathbf{A}| \neq 0$, ψ_i 's can be solved using Cramer's rule, and c_i 's are zeros of a quadratic function

$$g_2(x) = x^2 + \frac{\gamma_1\gamma_2 - \gamma_0\gamma_3}{|\mathbf{A}|}x + \frac{\gamma_1\gamma_3 - \gamma_2^2}{|\mathbf{A}|}.$$

Hence, the centers c_1, c_2 are given by

$$(4.3) \quad c_{1,2} = \frac{(\gamma_0\gamma_3 - \gamma_1\gamma_2) \pm \sqrt{D}}{2|\mathbf{A}|}, \quad c_1 < c_2,$$

under two assumptions

$$(4.4) \quad |\mathbf{A}| = \gamma_0\gamma_2 - \gamma_1^2 \neq 0, \quad D := (\gamma_1\gamma_2 - \gamma_0\gamma_3)^2 - 4(\gamma_0\gamma_2 - \gamma_1^2)(\gamma_1\gamma_3 - \gamma_2^2) > 0.$$

After obtaining c_i 's, the problem (3.5) is easily solved and gives

$$(4.5) \quad \rho_1 = \frac{\gamma_0c_2 - \gamma_1}{c_2 - c_1}, \quad \rho_2 = \frac{\gamma_0c_1 - \gamma_1}{c_1 - c_2}.$$

From Theorem 2.4 we may conclude that if $D > 0$ and $|\mathbf{A}| \neq 0$, then

$$(4.6) \quad \lim_{t \rightarrow \infty} t^{\left(\frac{5}{2} - \frac{1}{2p}\right)} \|u(t) - \phi_2(t)\|_p \leq \frac{\left\| \partial_\xi^4 \left(e^{-\frac{1}{4}\xi^2} \right) \right\|_p}{\sqrt{4\pi}} \left| \int_{-\infty}^{\infty} E_4(x) dx \right|, \quad 1 \leq p \leq \infty,$$

where $E_4(x)$ is the fourth order antiderivative of the initial error $E_0(x) := u_0(x) - \sum_{i=1}^2 \rho_i \delta_{c_i}(x)$ given by (2.4).

Example 4.1. Consider an initial value

$$(4.7) \quad U_l(x) = \begin{cases} -1, & -2l - 0.5 < x < -l - 0.5, \quad l + 0.5 < x < 2l + 0.5, \\ 1, & -l - 0.5 \leq x \leq l + 0.5, \\ 0, & \text{otherwise,} \end{cases}$$

where $l > 0$. Let $\gamma_{k,l}$ be the k th moments of the function $U_l(x)$, i.e.,

$$\gamma_{k,l} := \int x^k U_l(x) dx, \quad k = 0, 1, \dots$$

Then $\gamma_{0,l} = 1$ for all $l > 0$ and, since U_l is an even function, $\gamma_{k,l} = 0$ for $k = 1, 3, 5, \dots$. Hence, $|A|$ and D in (4.4) are given by

$$|A| = \gamma_{2,l}, \quad D = 4(\gamma_{2,l})^3.$$

One may easily check that $\gamma_{2,l} = 0$, if and only if

$$l = l_2 := 0.5 \left(\sqrt[3]{2} - 1 \right) / \left(2 - \sqrt[3]{2} \right),$$

and $D = 4(\gamma_{2,l})^3 > 0$, if and only if $l < l_2$. Hence, the moment problem (4.2) with the initial value $U_l(x)$ is solvable only for $l < l_2$. This example says that, even if the Hankel matrix is nonsingular, $\phi_2(x, t)$ that satisfies convergence order in (4.6) may not exist.

4.3. Approximation with three heat kernels. The derivation of (3.6) from (3.5) is not clear without some calculations. In the following, such a derivation is given for an example. For the case $n = 3$ the system (3.5) reads

$$(4.8) \quad \begin{aligned} \rho_1 + \rho_2 + \rho_3 &= \gamma_0, \\ \rho_1 c_1 + \rho_2 c_2 + \rho_3 c_3 &= \gamma_1, \\ \rho_1 c_1^2 + \rho_2 c_2^2 + \rho_3 c_3^2 &= \gamma_2, \\ \rho_1 c_1^3 + \rho_2 c_2^3 + \rho_3 c_3^3 &= \gamma_3, \\ \rho_1 c_1^4 + \rho_2 c_2^4 + \rho_3 c_3^4 &= \gamma_4, \\ \rho_1 c_1^5 + \rho_2 c_2^5 + \rho_3 c_3^5 &= \gamma_5. \end{aligned}$$

Multiply c_1 to the k th equation and subtract $(k + 1)$ th one from it for $k = 1, \dots, 5$ and obtain five equations without ρ_1 , i.e.,

$$\begin{aligned} \rho_2(c_1 - c_2) + \rho_3(c_1 - c_3) &= \gamma_0 c_1 - \gamma_1, \\ \rho_2(c_1 - c_2)c_2 + \rho_3(c_1 - c_3)c_3 &= \gamma_1 c_1 - \gamma_2, \\ \rho_2(c_1 - c_2)c_2^2 + \rho_3(c_1 - c_3)c_3^2 &= \gamma_2 c_1 - \gamma_3, \\ \rho_2(c_1 - c_2)c_2^3 + \rho_3(c_1 - c_3)c_3^3 &= \gamma_3 c_1 - \gamma_4, \\ \rho_2(c_1 - c_2)c_2^4 + \rho_3(c_1 - c_3)c_3^4 &= \gamma_4 c_1 - \gamma_5. \end{aligned}$$

Do the similar process two more times and obtain three equations without ρ_i 's:

$$\begin{aligned} 0 &= \gamma_0 c_1 c_2 c_3 - \gamma_1(c_1 c_2 + c_2 c_3 + c_3 c_1) + \gamma_2(c_1 + c_2 + c_3) - \gamma_3, \\ 0 &= \gamma_1 c_1 c_2 c_3 - \gamma_2(c_1 c_2 + c_2 c_3 + c_3 c_1) + \gamma_3(c_1 + c_2 + c_3) - \gamma_4, \\ 0 &= \gamma_2 c_1 c_2 c_3 - \gamma_3(c_1 c_2 + c_2 c_3 + c_3 c_1) + \gamma_4(c_1 + c_2 + c_3) - \gamma_5, \end{aligned}$$

which are identical to (3.6)–(3.7) with $n = 3$, i.e.,

$$\begin{pmatrix} \gamma_0 & \gamma_1 & \gamma_2 \\ \gamma_1 & \gamma_2 & \gamma_3 \\ \gamma_2 & \gamma_3 & \gamma_4 \end{pmatrix} \begin{pmatrix} \psi_0 \\ \psi_1 \\ \psi_2 \end{pmatrix} = \begin{pmatrix} \gamma_3 \\ \gamma_4 \\ \gamma_5 \end{pmatrix},$$

where $\psi_0 = c_1 c_2 c_3$, $\psi_1 = -(c_1 c_2 + c_2 c_3 + c_3 c_1)$, and $\psi_2 = c_1 + c_2 + c_3$. The derivation is done for the case $n = 3$.

The determinant of the 3×3 Hankel matrix is given by

$$|\mathbf{A}| = \gamma_0 \gamma_2 \gamma_4 + 2\gamma_1 \gamma_2 \gamma_3 - \gamma_2^3 - \gamma_0 \gamma_3^2 - \gamma_1^2 \gamma_4.$$

If $|\mathbf{A}| \neq 0$, then ψ_i are given by Cramer's rule:

$$\begin{aligned} \psi_0 &= (2\gamma_3 \gamma_2 \gamma_4 + \gamma_3 \gamma_1 \gamma_5 - \gamma_3^3 - \gamma_2^2 \gamma_5 - \gamma_4^2 \gamma_1) / |\mathbf{A}|, \\ \psi_1 &= (\gamma_2 \gamma_5 \gamma_1 + \gamma_0 \gamma_4^2 + \gamma_3^2 \gamma_2 - \gamma_3 \gamma_1 \gamma_4 - \gamma_4 \gamma_2^2 - \gamma_0 \gamma_3 \gamma_5) / |\mathbf{A}|, \\ \psi_2 &= (\gamma_0 \gamma_2 \gamma_5 + \gamma_3^2 \gamma_1 + \gamma_2 \gamma_4 \gamma_1 - \gamma_0 \gamma_3 \gamma_4 - \gamma_3 \gamma_2^2 - \gamma_1^2 \gamma_5) / |\mathbf{A}|. \end{aligned}$$

The points c_i 's are zeros of third order polynomial

$$(4.9) \quad g_3(x) = x^3 - \psi_2 x^2 - \psi_1 x - \psi_0.$$

Hence, the solvability of the problem (4.8) is equivalent to the existence of three distinct real roots $c_1 < c_2 < c_3$ of (4.9). The convergence order in Theorem 2.4 gives the asymptotic convergence order:

$$(4.10) \quad \lim_{t \rightarrow \infty} t^{(\frac{7}{2} - \frac{1}{2p})} \|u(t) - \phi_3(t)\|_p \leq \frac{\left\| \partial_\xi^6 \left(e^{-\frac{1}{4}\xi^2} \right) \right\|_p}{\sqrt{4\pi}} \left| \int_{-\infty}^{\infty} E_6(x) dx \right|, \quad 1 \leq p \leq \infty,$$

where $E_6(x)$ is the sixth order antiderivative of the initial error $E_0(x) := u_0(x) - \sum_{i=1}^3 \rho_i \delta_{c_i}(x)$ given by (2.4).

Consider the initial value given in Example 4.1. Since $\gamma_{1,l} = \gamma_{3,l} = \gamma_{5,l} = 0$ and $\gamma_{0,l} = 1$, we obtain

$$|A| = \gamma_{2,l} (\gamma_{4,l} - \gamma_{2,l}^2), \quad \psi_1 = \gamma_{4,l} / \gamma_{2,l}, \quad \psi_0 = \psi_2 = 0.$$

Hence, if

$$|A| \neq 0 \quad \text{and} \quad \psi_1 > 0,$$

then $g_3(x)$ has three distinct real roots

$$c_1 = -\sqrt{\psi_1}, \quad c_2 = 0, \quad c_3 = \sqrt{\psi_1}.$$

One may show that $\gamma_{4,l} > 0$ if and only if $0 < l < l_4 := 0.5(\sqrt[5]{2} - 1)/(2 - \sqrt[5]{2})$. Therefore, if $l_4 < l < l_2$, then $\psi_1 < 0$ and the existence of $\phi_3(x, t)$ satisfying (4.10) is not guaranteed. This example shows that the solvability of (3.5) is not obvious for sign-changing initial values.

5. Approximation for sign-changing solutions. Now consider a general sign-changing initial value. First, consider the case that the initial value $u_0(x)$ decays for $|x|$ large with the order that the Gaussian has, i.e.,

$$(5.1) \quad u_0(x) = O\left(e^{-|x|^2}\right) \quad \text{as} \quad |x| \rightarrow \infty.$$

Then there exists $M > 0$ such that $v_0(x) := u_0(x) + \frac{M}{\sqrt{4\pi}} e^{-x^2/4} \geq 0$, and we may apply the theory in section 3 to the nonnegative function v_0 . Let ρ_i 's and c_i 's be the solutions of the moment problem with initial value $v_0(x)$. Then the solution $u(x, t)$ can be approximated by

$$(5.2) \quad u(x, t) \sim \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-(x-c_i)^2/(4t)} - \frac{M}{\sqrt{4\pi(t+1)}} e^{-x^2/4(t+1)}.$$

Since the auxiliary part of the approximation is the exact solution with the extra initial value added to $u_0(x)$, the convergence order of this approximation is the same as the one in Theorem 3.1. This example shows that we may obtain the same convergence order for general sign-changing solutions by simply adding an extra term.

On the other hand, if the grid points are preassigned, say, $c_i = \bar{c}_i$, then we have the freedom in choosing the weights ρ_i 's only. These ρ_i 's are simply obtained by solving

the first n equations in (3.5), where the corresponding matrix is the Vandermonde matrix, i.e.,

$$(5.3) \quad \begin{pmatrix} 1 & 1 & \cdots & 1 \\ \bar{c}_1 & \bar{c}_2 & \cdots & \bar{c}_n \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \bar{c}_1^{n-1} & \bar{c}_2^{n-1} & \cdots & \bar{c}_n^{n-1} \end{pmatrix} \begin{pmatrix} \rho_1 \\ \rho_2 \\ \cdot \\ \cdot \\ \rho_n \end{pmatrix} = \begin{pmatrix} \gamma_0 \\ \gamma_1 \\ \cdot \\ \cdot \\ \gamma_{n-1} \end{pmatrix}.$$

The Vandermonde determinant $\prod_{1 \leq i < j \leq n} (\bar{c}_j - \bar{c}_i)$ is not zero if \bar{c}_i are all different and, hence, (5.3) is solvable. Now construct a different kind of approximation:

$$(5.4) \quad \eta_n(x, t) := \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-(x-\bar{c}_i)^2/(4t)}.$$

Then $\eta_n(\cdot, t)$ converges to $u(\cdot, t)$ with the order

$$\|u(t) - \eta_n(t)\|_p = O\left(t^{\left(\frac{1}{2p} - \frac{n+1}{2}\right)}\right) \quad \text{as } t \rightarrow \infty,$$

since $\lim_{t \rightarrow 0} (\eta_n(x, t) - u_0(x))$ has zero moments up to $(n - 1)$ th order.

6. Convergence as $n \rightarrow \infty$ with fixed $t > 0$. In this section, we discuss the convergence of the approximation $\phi_n(x, t)$ to the solution $u(x, t)$ as $n \rightarrow \infty$ with a fixed $t > 0$. An interesting behavior of the approximation $\phi_n(x, t)$ that one may observe numerically is a geometric convergence order such as

$$(6.1) \quad \beta_n(t) := \frac{\|u(t) - \phi_n(t)\|_\infty}{\|u(t) - \phi_{n+1}(t)\|_\infty} \rightarrow 1 + 4\frac{t}{v} \left(\equiv \beta\left(\frac{t}{v}\right) \right) \quad \text{as } n \rightarrow \infty,$$

where $v > 0$ depends on the initial value $u_0(x)$. (We do not have a proof of it. Hence, the statements here are rather conjectures.) This convergence order implies that the error decays to zero very fast as $n \rightarrow \infty$ for any fixed time $t > 0$. This convergence order is somewhat extreme. For example, if $t > v/4$, then the approximation error is reduced into half whenever just a single heat kernel is added.

Set the approximation error as

$$e_n(x, t) = u(x, t) - \phi_n(x, t).$$

Consider a sequence of functions

$$E_k^n(x) = \int_{-\infty}^x E_{k-1}^n(y) dy, \quad k = 1, 2, \dots, 2n,$$

where

$$E_0^n(x) := e_n(x, 0) = u_0(x) - \sum_{i=1}^n \rho_i \delta(x - c_i).$$

Notice that the upper index n is to denote that E_k^n is related to the approximation $\phi_n(x, t)$ and the lower index k is to indicate that E_k^n is the k th order antiderivative of the initial approximation error $E_0^n(x)$. Then, from (2.6), one obtains

$$(6.2) \quad (\sqrt{t})^{2n+1} \|e_n(t)\|_\infty = \frac{C_{2n}}{\sqrt{4\pi}} \sup_{\xi} \left| \int \partial_{\xi}^{2n} \left(e^{-\zeta^2/4} \right) \sqrt{t} E_{2n}^n \frac{\sqrt{t}(\xi - \zeta)}{C_{2n}} d\zeta \right|,$$

where $C_{2n} := \int_{-\infty}^{\infty} E_{2n}^n(x) dx < \infty$.

An interesting observation is that, for n large, $E_{2n}^n(x)$ has a Gaussian-like structure. The following has been observed numerically.

CONJECTURE 6.1. *Suppose that the initial value $u_0(x)$ is nonnegative and has finite moments up to any order. Then there exist $c \in \mathbf{R}$ and $v > 0$ such that*

$$(6.3) \quad \left\| \frac{1}{\sqrt{v\pi}} e^{-\frac{(x-c)^2}{v}} - \frac{1}{C_{2n}} E_{2n}^n(x) \right\|_{\infty} \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

where $C_{2n} := \int_{-\infty}^{\infty} E_{2n}^n(x) dx < \infty$. Furthermore,

$$(6.4) \quad \frac{C_{2n}}{C_{2(n+1)}} \frac{\left\| D_x^{2n} e^{-\frac{x^2}{v}} \right\|_{\infty}}{\left\| D_x^{2n+2} e^{-\frac{x^2}{v}} \right\|_{\infty}} \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

The $(2n - 1)$ th order derivative of $E_{2n}^n(x)$ is $E_1^n(x)$ which is, at most, of order $O(1/n)$, which does not make any difference in the geometric convergence order such as (6.1). Hence, we may treat it as of order $O(1)$. Note that C_{2n} is obtained after integrating $E_0^n(x)$ $2n$ times and, hence, its order should be the reciprocal of the order of $\|D_x^{2n}(e^{-\frac{x^2}{v}})\|_{\infty}$, which is the $2n$ th derivative of the Gaussian. Hence, (6.4) is a natural conclusion if (6.3) is assumed. Furthermore, for $\xi = x/\sqrt{v}$,

$$\left\| D_x^{2n} \left(e^{-\frac{x^2}{v}} \right) \right\|_{\infty} = \left\| D_{\xi}^{2n} \left(e^{-\xi^2} \right) (\xi_x)^{2n} \right\|_{\infty} = \frac{1}{v^n} \left\| D_{\xi}^{2n} \left(e^{-\xi^2} \right) \right\|_{\infty}.$$

Under Conjecture 6.1, the right-hand side of (6.2) can be approximated using $A(n, v/t)$ given by

$$\begin{aligned} \sup_x \left| \int D_y^{2n} \left(e^{-y^2/4} \right) \sqrt{t} E_{2n}^n \frac{\sqrt{t}(x-y)}{C_{2n}} dy \right| \\ \cong \frac{\sqrt{t}}{\sqrt{v\pi}} \sup_x \left| \int D_y^{2n} \left(e^{-y^2/4} \right) e^{-\frac{(x-y-c/\sqrt{t})^2}{v/t}} dy \right| \\ = \frac{\sqrt{t}}{\sqrt{v\pi}} \left| \int D_y^{2n} \left(e^{-y^2/4} \right) e^{-\frac{y^2}{v/t}} dy \right| =: A(n, v/t). \end{aligned}$$

Notice that due to the symmetry of $D_x^{2n}(e^{-x^2/4})$ and $e^{-\frac{x^2}{v/t}}$ the supremum of the second line is obtained at $x - c/\sqrt{t} = 0$. Then we obtain from the relations (6.2) and (6.4) that

$$t^{-1} \frac{\|e_n(t)\|_{\infty}}{\|e_{n+1}(t)\|_{\infty}} \cong \frac{C_{2n}}{C_{2(n+1)}} \frac{A(n, v/t)}{A(n+1, v/t)} \cong \frac{v^n}{v^{n+1}} \frac{\left\| D_x^{2n+2} \left(e^{-x^2} \right) \right\|_{\infty}}{\left\| D_x^{2n} \left(e^{-x^2} \right) \right\|_{\infty}} \frac{A(n, v/t)}{A(n+1, v/t)}.$$

One can easily check that

$$\frac{\left\| D_x^{2n+2} \left(e^{-x^2} \right) \right\|_{\infty}}{\left\| D_x^{2n} \left(e^{-x^2} \right) \right\|_{\infty}} = 4n + 2, \quad \frac{A(n, v/t)}{A(n+1, v/t)} = \frac{4 + v/t}{4n + 2},$$

using a mathematical software such as Maple or by hand. Therefore, we obtain the convergence order in (6.1), i.e.,

$$\frac{\|e_n(t)\|_{\infty}}{\|e_{n+1}(t)\|_{\infty}} \cong t \frac{1}{v} (4n + 2) \frac{4 + v/t}{4n + 2} = 1 + 4 \frac{t}{v} \quad \text{for } n \text{ large.}$$

Notice that $c \in \mathbf{R}$ in (6.3) does not make any difference in the convergence order. The factor that decides the geometric convergence rate is the variance factor v of the limit function $\frac{1}{\sqrt{v\pi}}e^{-x^2/v}$. It seems that the variance factor v depends on the initial value $u_0(x)$, and another discussion about it will be included in section 7.2.

Remark 6.2. For $n > 0$ small, $E_{2n}^n(x)/C_{2n}$ is not close enough to the Gaussian and the arguments above do not apply. Then it is natural to ask how large n should be. The answer depends on the initial value. Clearly, if $u_0(x)$ itself is like a Gaussian, then such an $n > 0$ can be relatively small. In other cases, the corresponding $n > 0$ could be larger.

7. Numerical examples. In this section, we test the convergence orders numerically for $t > 0$ large and for $n > 0$ large. These tests confirm the convergence orders obtained in the previous sections. This section consists of four subsections. The first two are for $t \rightarrow \infty$ and for $n \rightarrow \infty$ limits of the approximation $\phi_n(x, t)$. In the third one, we test the behavior of the alternative approach $\psi_{2n}(x, t)$ as $n \rightarrow \infty$. In the last one, we do numerical tests for Conjecture 6.1.

There are two difficulties in observing the theoretical convergence order for $t > 0$ large. First, the convergence rate for small time $0 < t \ll 1$ is lower than the theoretical one for $t > 0$ large. So we need to wait a certain amount of time to observe the theoretical convergence order. On the other hand, since the convergence order is so high, the approximation error at the right moment can be as small as of order 10^{-36} or 10^{-64} (see Tables 7.1 and 7.2). So we should employ enough precisions in the computation to obtain meaningful numerical results.

The second difficulty, which is more restrictive, is in computing the solution $u(x, t)$. To compute the decay order of $\|u(x, t) - \phi_n(x, t)\|_\infty$ accurately, we should obtain the exact value $u(x, t)$ or compute it with a smaller error than the actual approximation error. However, it seems impossible to do the integration in (1.1) numerically with such a small error. (In this sense, one may say that the approximation $\phi_n(x, t)$ is more exact than the exact formula in (1.1).) To avoid such a difficulty, we consider the following two examples with explicit solutions. In the following numerical tests we employ these examples.

Example 7.1 (example with a single hump). Consider the solution of

$$(7.1) \quad u_t = u_{xx}, \quad u(x, 0) = K(x, t_0), \quad x \in \mathbf{R}, \quad t > 0,$$

where $K(x, t)$ is the heat kernel

$$K(x, t) = \frac{1}{\sqrt{4\pi t}}e^{-x^2/4t}.$$

Then the exact solution is simply $u(x, t) = K(x, t + t_0)$ and the variance of the initial value is $var = 2t_0$. This rather simple example illustrates certain convergence behavior very clearly.

Example 7.2 (example with double humps). Consider the solution of

$$(7.2) \quad u_t = u_{xx}, \quad u(x, 0) = \frac{1}{2}[K(x + 1, t_0) + K(x - 1, t_0)], \quad x \in \mathbf{R}, \quad t > 0.$$

Then the solution is simply $u(x, t) = \frac{1}{2}[K(x + 1, t + t_0) + K(x - 1, t + t_0)]$ and the variance of the initial value is $var = 1 + 2t_0$.

TABLE 7.1

The error $e_n(x, t) = u(x, t) - \phi_n(x, t)$ and the convergence order α_n in (7.5) have been computed for Examples 7.1 and 7.2 with $n = 4, 8$ and $t_0 = 1$. We observe that $\alpha_n(t) \rightarrow -(n + \frac{1}{2})$ as $t \rightarrow \infty$. (The norms in this and the following tables are L^∞ -norms.)

t	Example 7.1				Example 7.2			
	$\ e_4(t)\ $	$\alpha_4(t)$	$\ e_8(t)\ $	$\alpha_8(t)$	$\ e_4(t)\ $	$\alpha_4(t)$	$\ e_8(t)\ $	$\alpha_8(t)$
0.1	2.17e-01,	0.7	1.13e-01,	1.0	2.24e-01,	0.8	1.30e-01,	0.9
0.2	1.13e-01,	0.9	2.98e-02,	1.9	1.33e-01,	0.8	4.57e-02,	1.5
0.4	3.48e-02,	1.7	3.34e-03,	3.2	5.30e-02,	1.3	7.27e-03,	2.7
0.8	6.24e-03,	2.5	1.38e-04,	4.6	1.21e-02,	2.1	4.27e-04,	4.1
1.6	6.81e-04,	3.2	2.21e-06,	6.0	1.60e-03,	2.9	9.09e-06,	5.6
3.2	5.12e-05,	3.7	1.72e-08,	7.0	1.37e-04,	3.5	8.57e-08,	6.7
6.4	3.03e-06,	4.1	8.40e-11,	7.7	8.79e-06,	4.0	4.68e-10,	7.5
12.8	1.56e-07,	4.3	3.13e-13,	8.1	4.73e-07,	4.2	1.86e-12,	8.0
25.6	7.48e-09,	4.4	1.01e-15,	8.3	2.31e-08,	4.4	6.18e-15,	8.2
51.2	3.44e-10,	4.4	3.01e-18,	8.4	1.08e-09,	4.4	1.88e-17,	8.4
102.4	1.55e-11,	4.5	8.66e-21,	8.4	4.89e-11,	4.5	5.44e-20,	8.4
204.8	6.93e-13,	4.5	2.44e-23,	8.5	2.19e-12,	4.5	1.54e-22,	8.5

TABLE 7.2

The error $e_n(x, t) = u(x, t) - \phi_n(x, t)$ and the geometric convergence rate $\beta_n(t)$ in (6.1) have been computed for Examples 7.1 and 7.2 with $t = 1, 10$ and $t_0 = 1$. The ratio $\beta_n(t)$ converges to the limit in (6.1) quickly with $v = 2$ for Example 1 and slowly for Example 2.

n	Example 7.1				Example 7.2			
	$\ e_n(1)\ $	$\beta_n(1)$	$\ e_n(10)\ $	$\beta_n(10)$	$\ e_n(1)\ $	$\beta_n(1)$	$\ e_n(10)\ $	$\beta_n(10)$
2	2.8e-02	2.91	2.0e-04	20.84	4.28e-02	2.48	3.83e-04	15.83
3	9.6e-03	2.96	9.5e-06	20.92	1.72e-02	2.49	2.32e-05	16.53
4	3.2e-03	2.98	4.5e-07	20.95	6.7e-03	2.57	1.36e-06	17.06
7	1.2e-04	2.99	4.9e-11	20.98	3.67e-04	2.66	2.47e-10	17.89
10	4.5e-06	3.0	5.3e-15	20.99	1.89e-05	2.71	4.09e-14	18.35
13	1.7e-07	3.0	5.8e-19	21.0	9.35e-07	2.74	6.45e-18	18.65
16	6.1e-09	3.0	6.2e-23	21.	4.45e-08	2.76	9.65e-22	18.86
19	2.3e-10	3.0	6.7e-27	21.0	2.08e-09	2.78	1.41e-25	19.02
22	8.4e-12	3.0	7.3e-31	21.0	9.65e-11	2.79	2.02e-29	19.15
25	3.1e-13	3.0	7.9e-35	21.0	4.36e-12	2.8	2.85e-33	19.26
46	2.97e-23	3.0	1.35e-62	21.0	1.38e-21	2.85	2.25e-60	19.70
49	1.1e-24	3.0	1.46e-66	21.0	5.95e-21	2.86	2.94e-64	19.74

7.1. Numerical tests for the long time asymptotics. The approximation

$$(7.3) \quad \phi_n(x, t) \equiv \sum_{i=1}^n \frac{\rho_i}{\sqrt{4\pi t}} e^{-\frac{(x-c_i)^2}{4t}}$$

constructed in section 3 converges to the exact solution $u(x, t)$ with order

$$(7.4) \quad \|u(t) - \phi_n(t)\|_\infty = O\left(t^{-\frac{2n+1}{2}}\right) \quad \text{as } t \rightarrow \infty.$$

In Table 7.1 the error $e_n(x, t) = u(x, t) - \phi_n(x, t)$ and the convergence order α_n have been computed for $n = 4, 8$ as doubling the time from $t = 0.1$ to $t = 204.8$. The convergence order of the approximation has been measured by computing

$$(7.5) \quad \alpha_n(t) \sim \frac{\ln(\|e_n(t/2)\|_\infty / \|e_n(t)\|_\infty)}{\ln(1/2)}.$$

(Note that we measure the error in L^∞ -norm in the following numerical examples and denote it by $\|\cdot\|$ in the tables to get it fitted in the tables.) From Table 7.1, one

clearly observes that the convergence order $\alpha_n(t)$ approaches the optimal convergence order in (7.4) as $t \rightarrow \infty$. Notice that these numerical tests for $t \rightarrow \infty$ limits show similar patterns of the convergence for both Examples 7.1 and 7.2.

7.2. Numerical tests for $n \rightarrow \infty$ limits. Now we are going to check the convergence order for n large with a fixed $t > 0$. Consider the ratio

$$\beta_n(t) := \|e_{n-1}(t)\|_\infty / \|e_n(t)\|_\infty.$$

(The ratio $r = a_n/a_{n-1}$ is usually considered for a geometric sequence. Here we consider its reciprocal for easier comparison.) In Table 7.2, the error $\|e_n(t)\|_\infty$ and this ratio are computed for Examples 7.1 and 7.2 at two instances $t = 1$ and $t = 10$ with increasing n from $n = 2$ to $n = 25$. One can clearly observe a certain geometric convergence order as in (6.1). In both examples, we can clearly see that $10(\beta_n(1) - 1) \sim (\beta_n(10) - 1)$, which indicates that the corresponding constant $v > 0$ in (6.1) which decides the geometric convergence ratio does not depend on the time $t > 0$.

From the test for Example 7.1, one can clearly see that $\beta_n(1) \rightarrow 3$ and $\beta_n(10) \rightarrow 21$ as $n \rightarrow \infty$. In both cases, the corresponding v is $v = 2$ which is the variance of the initial value. The convergence pattern for Example 7.2 is different. First, the convergence speed of the ratio $\beta_n(t)$ is slow. It seems due to the complexity of the structure of the initial value. At the moment $n = 49$ the index v corresponding to the geometric convergence rate $\beta = 19.74$ is $v = 2.15$ and seems still decreasing. This value is already smaller than the variance of the initial value which is 3 and looks likely to converge to $v = 2$. It seems that the factor that decides the geometric convergence rate is not the variance but the tail of the initial value for $|x|$ large.

7.3. Approximation using derivatives of the Gaussian. We may write $\psi_{2n}(x, t)$ in (1.5) as

$$(7.6) \quad \psi_{2n}(x, t) = \sum_{i=0}^{2n-1} \frac{-\gamma_i}{2(i!)} \left(\frac{-1}{2\sqrt{t}}\right)^{n+1} H_i\left(\frac{x}{2\sqrt{t}}\right) e^{-x^2/4t},$$

where $H_i(x)$ is the Hermite polynomial of degree i . In Table 7.3, the approximation error and the geometric convergence ratio β_n are given for Example 7.1. To make the relation between the initial value and the convergence ratio more clear, we set

$$(7.7) \quad \beta_{2n}(t, t_0) := \|e_{2n-2}(t)\|_\infty / \|e_{2n}(t)\|_\infty,$$

where $e_{2n}(x, t) = u(x, t) - \psi_{2n}(x, t)$ and $u_0(x) = K(x, t_0)$ with $t_0 = 10$.

One may observe that $\beta_{2n}(1, 10) \rightarrow 0.1$, $\beta_{2n}(10, 10) \rightarrow 1$, and $\beta_{2n}(100, 10) \rightarrow 10$ as $n \rightarrow \infty$. This observation leads us to a conjecture

$$(7.8) \quad \lim_{n \rightarrow \infty} \frac{\|u(t) - \psi_{2n}(t)\|_\infty}{\|u(t) - \psi_{2n+2}(t)\|_\infty} = \frac{t}{t_0},$$

which indicates that the approximation error increases geometrically if $t < t_0$ and, hence, $\psi_{2n}(x, t)$ is meaningful for $t > t_0$ only. We may consider t_0 as the age of the initial value since $u_0(x) = K(x, t_0)$. This t_0 seems to be related to the time-shift t_* in [22]. Generally, one may call t_0 the age of a general initial value u_0 of the heat equation if $\lim_{n \rightarrow \infty} \frac{\|u(t_0) - \psi_{2n}(t_0)\|_\infty}{\|u(t_0) - \psi_{2n+2}(t_0)\|_\infty} = 1$.

TABLE 7.3

The error $e_n(x, t) = u(x, t) - \psi_{2n}(x, t)$ and the geometric convergence rate $\beta_n(t, t_0)$ in (7.7) have been computed numerically for Example 7.1 with $t_0 = 10$ and $t = 1, 10, 100$. We may observe the convergence rate in (7.8).

$2n$	$\ e_{2n}(1)\ $	$\beta_{2n}(1, 10)$	$\ e_{2n}(10)\ $	$\beta_{2n}(10, 10)$	$\ e_{2n}(100)\ $	$\beta_{2n}(100, 10)$
4	1.21e+00	0.1623821	1.85e-02	1.4142136	9.77e-05	13.4400610
6	9.37e+00	0.1295695	1.50e-02	1.2335625	8.11e-06	12.0475285
8	7.88e+01	0.1188625	1.29e-02	1.1610328	7.08e-07	11.4555359
14	5.74e+04	0.1089029	9.66e-03	1.0825322	5.40e-10	10.7781946
20	4.73e+07	0.1058095	8.05e-03	1.0553099	4.54e-13	10.5307485
26	4.12e+10	0.1043095	7.04e-03	1.0415615	3.99e-16	10.4026370
32	3.69e+13	0.1034247	6.34e-03	1.0332791	3.60e-19	10.3243276
38	3.38e+16	0.1028412	5.81e-03	1.0277462	3.31e-22	10.2715121
44	3.13e+19	0.1024275	5.39e-03	1.0237896	3.07e-25	10.2334861
50	2.93e+22	0.1021190	5.06e-03	1.0208199	2.88e-28	10.2048012

TABLE 7.4

We may observe conjectures in (6.3) and (6.4) numerically. In this table, those conjectures are tested using Examples 7.1 and 7.2 with $t_0 = 1$ and $v = 2t_0$.

n	$\left\ \frac{1}{\sqrt{2t_0}\pi} e^{-\frac{x^2}{2t_0}} - \frac{1}{C_{2n}} E_{2n}^n(x) \right\ $		$\frac{C_{2n-2}}{C_{2n}} \frac{\ D_x^{2n-2} e^{-x^2/2t_0}\ }{\ D_x^{2n} e^{-x^2/2t_0}\ }$	
	Example 7.1	Example 7.2	Example 7.1	Example 7.2
2	2.233e-02	4.378e-02	1.0	0.7500000
3	2.070e-02	3.675e-02	1.0	0.7826087
4	1.631e-02	3.547e-02	1.0	0.8070175
5	1.387e-02	3.353e-02	1.0	0.8237512
6	1.207e-02	3.194e-02	1.0	0.8369731
7	1.070e-02	3.054e-02	1.0	0.8474264
8	9.675e-03	2.932e-02	1.0	0.8561101
9	8.743e-03	2.825e-02	1.0	0.8634099
10	8.016e-03	2.729e-02	1.0	0.8696921

7.4. Numerical test for Conjecture 6.1. The geometric convergence rate (6.1) for n large has been obtained under Conjecture 6.1 and observed numerically in section 7.2. The conjecture itself is of an independent interest which has no direct relation with the heat equation. In this section, we test the limits (6.3) and (6.4) in the conjecture.

In Table 7.4, the uniform norm of the difference in (6.3) and the ratio in (6.4) are given for Examples 7.1 and 7.2. In both cases, we have taken $v = 2t_0$. The results for Example 7.1 show the convergence clearly. In particular, the test for (6.4) shows that the ratio is identically one if the initial value $u_0(x)$ is the Gaussian.

The columns for Example 7.2 also show similar convergence behavior. However, the speed of its convergence is a lot slower. In fact, it is not even clear that taking $v = 2t_0$ is the correct one for the case of Example 7.2. Since $2t_0$ is the variance of Example 7.1, one may also try $2t_0 + 1$, which is the variance of Example 7.2. However, one cannot obtain the convergence, and $v = 2t_0$ seems a better choice. The computation is done up to $n = 10$ and that was our best. If computations with higher n are performed, we conjecture that the convergence to the unity will be more clearly observed.

Acknowledgment. Authors would like to thank anonymous reviewers. Their comments and suggestions helped to improve this paper.

REFERENCES

- [1] N. I. AKHIEZER, *The Classical Moment Problem and Some Related Questions in Analysis*, Hafner, New York, 1965.
- [2] J. A. CARRILLO AND J. L. VÁZQUEZ, *Fine asymptotics for fast diffusion equations*, Comm. Partial Differential Equations, 28 (2003), pp. 1023–1056.
- [3] R. E. CURTO AND L. A. FIALKOW, *Recursiveness, positivity, and truncated moment problems*, Houston J. Math., 17 (1991), pp. 603–635.
- [4] R. E. CURTO AND L. A. FIALKOW, *Truncated K -moment problems in several variables*, J. Operator Theory, 54 (2005), pp. 189–226.
- [5] C. M. DAFERMOS, *Regularity and large time behaviour of solutions of a conservation law without convexity*, Proc. Roy. Soc. Edinburgh Sect. A, 99 (1985), pp. 201–239.
- [6] J. DENZLER AND R. MCCANN, *Fast diffusion to self-similarity: Complete spectrum, long time asymptotics, and numerology*, Arch. Ration. Mech. Anal., 175 (2005), pp. 301–342.
- [7] R. J. DiPERNA, *Decay and asymptotic behavior of solutions to nonlinear hyperbolic systems of conservation laws*, Indiana Univ. Math. J., 24 (1974/75), pp. 1047–1071.
- [8] J. DOLBEAULT AND M. ESCOBEDO, *L^1 and L^∞ intermediate asymptotics for scalar conservation laws*, Asymptot. Anal., 41 (2005), pp. 189–213.
- [9] J. DUOANDIKOETXEA AND E. ZUAZUA, *Moments, masses de Dirac et decomposition de fonctions*, C. R. Acad. Sci. Paris Ser. I Math., 315 (1992), pp. 693–698.
- [10] M. ESCOBEDO AND E. ZUAZUA, *Large time behavior for convection-diffusion equations in R^N* , J. Funct. Anal., 100 (1991), pp. 119–161.
- [11] M. ESCOBEDO, J. VÁZQUEZ, AND E. ZUAZUA, *Asymptotic behaviour and source-type solutions for a diffusion-convection equation*, Arch. Ration. Mech. Anal., 124 (1993), pp. 43–65.
- [12] E. HOPF, *The partial differential equation $u_t + uu_x = \mu u_{xx}$* , Comm. Pure Appl. Math., 3 (1950), pp. 201–230.
- [13] Y.-J. KIM, *Asymptotic behavior in scalar conservation laws and the optimal convergence order to N -waves*, J. Differential Equations, 192 (2003), pp. 202–224.
- [14] Y.-J. KIM, *An Oleinik type estimate for a convection-diffusion equation and the convergence to N -waves*, J. Differential Equations, 199 (2004), pp. 269–289.
- [15] Y.-J. KIM AND W.-M. NI, *On the rate of convergence and asymptotic profile of solutions to the viscous Burgers equation*, Indiana Univ. Math. J., 51 (2002), pp. 727–752.
- [16] Y.-J. KIM AND A. E. TZAVARAS, *Diffusive N -waves and metastability in the Burgers equation*, SIAM J. Math. Anal., 33 (2001), pp. 607–633.
- [17] P. D. LAX, *Hyperbolic systems of conservation laws. II*, Comm. Pure Appl. Math., 10 (1957), pp. 537–566.
- [18] T.-P. LIU, *Decay to N -waves of solutions of general systems of nonlinear hyperbolic conservation laws*, Comm. Pure Appl. Math., 30 (1977), pp. 586–611.
- [19] T.-P. LIU, *Nonlinear stability of shock waves for viscous conservation laws*, Mem. Amer. Math. Soc., 56 (1985), no. 328.
- [20] T.-P. LIU AND M. PIERRE, *Source-solutions and asymptotic behavior in conservation laws*, J. Differential Equations, 51 (1984), pp. 419–441.
- [21] E. M. STEIN, *Singular Integrals and Differentiability Properties of Functions*, Princeton Math. Ser. 30, Princeton University Press, Princeton, NJ, 1970.
- [22] T. P. WITELSKI AND A. J. BERNOFF, *Self-similar asymptotics for linear and nonlinear diffusion equations*, Stud. Appl. Math., 100 (1998), pp. 153–193.
- [23] G. WHITHAM, *Linear and Nonlinear Waves*, Pure Appl. Math., Wiley Interscience, New York, 1974.

ON DIVERGENCE FORM SPDES WITH VMO COEFFICIENTS*

N. V. KRYLOV†

Abstract. We present several results on solvability in Sobolev spaces W_p^1 of SPDEs in divergence form in the whole space.

Key words. stochastic partial differential equations, divergence equations, Sobolev spaces

AMS subject classifications. 60H15, 35R60

DOI. 10.1137/080726902

1. Introduction. The theory of (usual) partial differential equations has two rather different parts depending on whether the equations are written in divergence or nondivergence form. Quite often the starting point is the same: equations with constant coefficients, and then one uses different techniques to treat different types of equations.

By now, one can say that the L_p -theory of evolutionary second-order SPDEs is quite well developed. The most advanced results of this theory can be found in the following papers and references therein: [1] (nondivergence type equations), [2] and [3] (divergence type equations). The results of the present paper are close to the corresponding results of [2]. However, unlike [2] we do not assume that the leading coefficients are continuous in the space variable. Instead we assume that the leading coefficients of the “deterministic” part of the equation are in VMO (the space of functions with vanishing mean oscillation), which is a much wider class than C . Still the leading coefficients of the “stochastic” part are assumed to be continuous in x .

The exposition in [2] and [3] is based on the theory of solvability in spaces $H_p^\gamma = (1 - \Delta)^{-\gamma/2} L_p$ of SPDEs with coefficients independent of x . Then the method of “freezing” the coefficients is applied as in the general framework set out in [6]. This method does not work if the coefficients are only in VMO so we use a different technique based on recent results from [8] on deterministic parabolic equations with VMO coefficients. In addition, our technique allows us to avoid using the W_2^n theory of SPDEs, which is a starting point in the paper [6] and subsequent articles based on it.

One more difference of our approach from the one in [2] is that we represent the free term in the deterministic part in the form $D_i f^i + f^0$ with $f^j \in L_p$ (see (1.1) below). Of course, this is just a general form of a distribution from H_p^{-1} . However, the spaces H_p^γ are most appropriate for equations in nondivergence form. One general inconvenience of these spaces is that the space or space-time dilations affect the norms in a way which is hard to control. For divergence form equations with low regularity of coefficients the most important space is H_p^1 . This space coincides with the Sobolev space W_p^1 and the effect of dilations on the norm or on $D_i f^i + f^0$ can be easily taken into account.

The exposition here is self-contained apart from references to some very basic results of [6], [8], and [12], and is much more elementary than in [2], employing the

*Received by the editors June 11, 2008; accepted for publication (in revised form) October 30, 2008; published electronically February 20, 2009. This work was partially supported by NSF Grant DMS-0653121.

<http://www.siam.org/journals/sima/40-6/72690.html>

†School of Mathematics, University of Minnesota, 127 Vincent Hall, Minneapolis, MN 55455 (krylov@math.umn.edu).

derivatives instead of the powers of the Laplacian, and yet gives more information. In particular, the author intends to use Corollary 5.5 in order to largely simplify the theory in [2] of divergence form SPDEs in domains. It turns out that to develop this theory one need not first develop the theory of SPDEs in domains with coefficient independent of x , which in itself required quite a bit of work.

The author's interest in divergence-type equations and in simplifying the theory of them appeared after he realized that the corresponding results can be applied to filtering theory of partially observable diffusion processes, given by stochastic Itô equations. It turns out that, under Lipschitz and nondegeneracy conditions only, the filtering density is almost Lipschitz in x and almost Hölder 1/2 in time. This is proved in [10] on the basis of Theorems 2.2 through 2.6 of the present article. The filtering density satisfies an SPDE usually written in terms of the operators adjoint to operators in nondivergence form with Lipschitz continuous coefficients. Writing these adjoint operators in divergence form makes perfect sense and allows us to obtain the above-mentioned results (see [10]).

Our Theorem 2.2 is very close to Theorem 2.12 of [2]. Apart from weaker conditions on the coefficients, another important difference is the presence of the parameter λ in (2.10). One of differences in the proofs is that we avoid proving the solvability on small consecutive time intervals and then gluing together the results.

Let (Ω, \mathcal{F}, P) be a complete probability space with an increasing filtration $\{\mathcal{F}_t, t \geq 0\}$ of complete with respect to (\mathcal{F}, P) σ -fields $\mathcal{F}_t \subset \mathcal{F}$. Denote by \mathcal{P} the predictable σ -field in $\Omega \times (0, \infty)$ associated with $\{\mathcal{F}_t\}$. Let $w_t^k, k = 1, 2, \dots$, be independent one-dimensional Wiener processes with respect to $\{\mathcal{F}_t\}$.

We fix a stopping time τ and for $t \leq \tau$ in the Euclidean d -dimensional space \mathbb{R}^d of points $x = (x^1, \dots, x^d)$ we consider the following equation

$$(1.1) \quad du_t = (L_t u_t - \lambda u_t + D_i f_t^i + f_t^0) dt + (\Lambda_t^k u_t + g_t^k) dw_t^k,$$

where $u_t = u_t(x) = u_t(\omega, x)$ is an unknown function,

$$L_t \psi(x) = D_j \left(a_t^{ij}(x) D_i \psi(x) + a_t^j(x) \psi(x) \right) + b_t^i(x) D_i \psi(x) + c_t(x) \psi(x),$$

$$\Lambda_t^k \psi(x) = \sigma_t^{ik}(x) D_i \psi(x) + \nu_t^k(x) \psi(x),$$

the summation convention with respect to $i, j = 1, \dots, d$ and $k = 1, 2, \dots$ is enforced and detailed assumptions on the coefficients and the free terms will be given later.

One can rewrite (1.1) in the nondivergence form assuming that the coefficients a_t^{ij} and a_t^j are differentiable in x and then one could apply the results from [6]. It turns out that the differentiability of a_t^{ij} and a_t^j is not needed for the corresponding counterparts of the results in [6] to be true and showing this and generalizing the corresponding results of [2] is one of the main purposes of the present article.

The author is sincerely grateful to Kyeong-Hun Kim who kindly pointed out a serious error in the first draft of the article. Doyoon Kim and the referees of the paper made many valuable suggestions and their impact is greatly appreciated.

2. Main results. Fix a number

$$p \geq 2,$$

and denote $L_p = L_p(\mathbb{R}^d)$. We use the same notation L_p for vector- and matrix-valued or else ℓ_2 -valued functions such as $g_t = (g_t^k)$ in (1.1). For instance, if $u(x) =$

$(u^1(x), u^2(x), \dots)$ is an ℓ_2 -valued measurable function on \mathbb{R}^d , then

$$\|u\|_{L_p}^p = \int_{\mathbb{R}^d} |u(x)|_{\ell_2}^p dx = \int_{\mathbb{R}^d} \left(\sum_{k=1}^{\infty} |u^k(x)|^2 \right)^{p/2} dx.$$

Introduce

$$D_i = \frac{\partial}{\partial x^i}, \quad i = 1, \dots, d, \quad \Delta = D_1^2 + \dots + D_d^2.$$

By Du we mean the gradient with respect to x of a function u on \mathbb{R}^d .

As usual,

$$W_p^1 = \{u \in L_p : Du \in L_p\}, \quad \|u\|_{W_p^1} = \|u\|_{L_p} + \|Du\|_{L_p}.$$

Recall that τ is a stopping time and introduce

$$\mathbf{L}_p(\tau) := L_p((0, \tau], \mathcal{P}, L_p), \quad \mathbf{W}_p^1(\tau) := L_p((0, \tau], \mathcal{P}, W_p^1).$$

We also need the space $\mathcal{W}_p^1(\tau)$, which is the space of functions $u_t = u_t(\omega, \cdot)$ on $\{(\omega, t) : 0 \leq t \leq \tau, t < \infty\}$ with values in the space of generalized functions on \mathbb{R}^d and having the following properties:

- (i) We have $u_0 \in L_p(\Omega, \mathcal{F}_0, L_p)$;
- (ii) We have $u \in \mathbf{W}_p^1(\tau)$;
- (iii) There exist $f^i \in \mathbf{L}_p(\tau)$, $i = 0, \dots, d$, and $g = (g^1, g^2, \dots) \in \mathbf{L}_p(\tau)$ such that for any $\varphi \in C_0^\infty = C_0^\infty(\mathbb{R}^d)$ with probability 1 for all $t \in [0, \infty)$ we have

$$(2.1) \quad \begin{aligned} (u_{t \wedge \tau}, \varphi) &= (u_0, \varphi) + \sum_{k=1}^{\infty} \int_0^t I_{s \leq \tau} (g_s^k, \varphi) dw_s^k \\ &+ \int_0^t I_{s \leq \tau} ((f_s^0, \varphi) - (f_s^i, D_i \varphi)) ds. \end{aligned}$$

In particular, for any $\phi \in C_0^\infty$, the process $(u_{t \wedge \tau}, \phi)$ is \mathcal{F}_t -adapted and (a.s.) continuous.

The reader can find in [6] a discussion of (ii) and (iii), in particular, the fact that the series in (2.1) converges uniformly in probability on every finite subinterval of $[0, \tau]$. On the other hand, it is worth saying that the above introduced space \mathcal{W}_p^1 is not quite the same as $\mathcal{H}_p^1(\tau)$ in [6] or in [2]. There are three differences. One is that there is an additional restriction on u_0 in [6] and [2]. But in the main part of the article we are going to work with $\mathcal{W}_{p,0}^1(\tau)$ which is the subset of $\mathcal{W}_p^1(\tau)$ consisting of functions with $u_0 = 0$. Another issue is that in [6] and [2] in place of $D_i f^i + f^0$ we have just f such that

$$f \in \mathbf{H}_p^{-1}(\tau) = L_p((0, \tau], \mathcal{P}, H_p^{-1}).$$

Actually, this difference is fictitious because one knows that any $f \in H_p^{-1}$

- (a) has the form $D_i f^i + f^0$ with $f^j \in L_p$ and

$$\|f\|_{H_p^{-1}} \leq N \sum_{j=0}^d \|f^j\|_{L_p},$$

where N is independent of f, f^j , and on the other hand,

(b) for any $f \in H_p^{-1}$ there exist $f^j \in L_p$ such that $f = D_i f^i + f^0$ and

$$\sum_{j=0}^d \|f^j\|_{L_p} \leq N \|f\|_{H_p^{-1}},$$

where N is independent of f .

The third difference is that instead of (ii) the condition $D^2u \in \mathbf{H}_p^{-1}(\tau)$ is required in [6] and [2]. However, as it follows from Theorem 3.7 of [6] and the boundedness of the operator $D : L_p \rightarrow H_p^{-1}$, this difference disappears if τ is a bounded stopping time.

To summarize, the spaces $\mathcal{W}_{p,0}^1(\tau)$ introduced above coincide with $\mathcal{H}_{p,0}^1(\tau)$ from [6] if τ is bounded, and we choose a particular representation of the deterministic part of the stochastic differential just for convenience. In the remainder of the article the spaces $\mathcal{H}_{p,0}^1(\tau)$ do not appear and none of their properties is used.

In case that property (iii) holds, we write

$$(2.2) \quad du_t = (D_i f_t^i + f_t^0) dt + g_t^k dw_t^k$$

for $t \leq \tau$ and this explains the sense in which equation (1.1) is understood. Of course, we still need to specify appropriate assumptions on the coefficients and the free terms in (1.1).

Assumption 2.1. (i) The coefficients a_t^{ij} , a_t^i , b_t^i , σ_t^{ik} , c_t , and ν_t^k are measurable with respect to $\mathcal{P} \times \mathcal{B}(\mathbb{R}^d)$, where $\mathcal{B}(\mathbb{R}^d)$ is the Borel σ -field on \mathbb{R}^d .

(ii) There is a constant K such that for all values of indices and arguments

$$|a_t^i| + |b_t^i| + |c_t| + |\nu_t^k| \leq K, \quad c_t \leq 0.$$

(iii) There is a constant $\delta > 0$ such that for all values of the arguments and $\xi \in \mathbb{R}^d$

$$(2.3) \quad a_t^{ij} \xi^i \xi^j \leq \delta^{-1} |\xi|^2, \quad (a_t^{ij} - \alpha_t^{ij}) \xi^i \xi^j \geq \delta |\xi|^2,$$

where $\alpha_t^{ij} = (1/2)(\sigma^{i\cdot}, \sigma^{j\cdot})_{\ell_2}$. Finally, the constant $\lambda \geq 0$.

It is worth emphasizing that we do not require the matrix (a^{ij}) to be symmetric.

Assumption 2.1 guarantees that (1.1) makes perfect sense if $u \in \mathcal{W}_p^1(\tau)$. By the way, adding the term $-\lambda u_t$ with constant $\lambda \geq 0$ is one more technically convenient step. One can always introduce this term, if originally it is absent, by considering $v_t := u_t e^{\lambda t}$.

Let \mathbf{B} denote the set of balls $B \subset \mathbb{R}^d$ and let $\rho(B)$ be the radius of $B \in \mathbf{B}$. For functions $h_t(x)$ on $[0, \infty) \times \mathbb{R}^d$ and $B \in \mathbf{B}$ introduce

$$h_{t(B)} = \frac{1}{|B|} \int_B h_t(x) dx,$$

where $|B|$ is the volume of B . Also let \mathbf{Q} denote the set of all cylinders in $[0, \infty) \times \mathbb{R}^d$ of type $Q = (s, t) \times B$, where $B \in \mathbf{B}$ and $t - s = \rho^2(B)$. For such Q set $\rho(Q) = \rho(B)$. For $\rho \geq 0$, $s < t$, a continuous \mathbb{R}^d -valued function $x_r, r \in [s, t]$, and a $Q = (s, t) \times B \in \mathbf{Q}$, introduce

$$\begin{aligned} \text{osc}(h, Q, x.) &= \frac{1}{t-s} \int_s^t (|h_r - h_{r(B+x_r)}|)_{(B+x_r)} dr, \\ \text{Osc}(h, Q, \rho) &= \sup_{|x.|_C \leq \rho} \text{osc}(h, Q, x.), \quad \text{osc}(h, Q) = \text{osc}(h, Q, 0), \end{aligned}$$

where $|x.|_C$ is the sup norm of $|x.|$.

Observe that $\text{osc}(h, Q, x) = 0$ if $h_t(x)$ is independent of x .

Denote by B_ρ the open ball with radius $\rho > 0$ centered at the origin, define $Q_\rho = (0, \rho^2) \times B_\rho$ and for $t \geq 0$ and $x \in \mathbb{R}^d$ set $B_\rho(x) = B_\rho + x$, $Q_\rho(t, x) = Q_\rho + (t, x)$.

In the remaining two assumptions we use constants $\beta > 0$ and $\beta_1 > 0$, the values of which will be specified later.

Let $t_0 \geq 0$, $x_0 \in \mathbb{R}^d$, and constants $\varepsilon \geq \varepsilon_1 > 0$. We say that the couple (a, σ) is $(\varepsilon, \varepsilon_1)$ -regular at point (t_0, x_0) if (for any ω) either

(i) we have $\sigma_t^{nm}(x_0) = 0$ for $t \in (t_0, t_0 + \varepsilon_1^2)$ and all n, m and

$$(2.4) \quad \text{osc}(a^{ij}, Q) \leq \beta, \quad \forall i, j,$$

for all $Q \in \mathbf{Q}$ such that $Q \subset Q_\varepsilon(t_0, x_0)$, or

(ii) for all $Q \in \mathbf{Q}$ such that $Q \subset Q_\varepsilon(t_0, x_0)$ we have

$$(2.5) \quad \text{Osc}(a^{ij}, Q, \varepsilon) \leq \beta, \quad \forall i, j.$$

Note that (a, σ) is $(\varepsilon, \varepsilon_1)$ -regular at any point (t_0, x_0) for any $\beta > 0$ if, for instance, a^{ij} depend only on x and are of class VMO.

Assumption 2.2. There exist $\varepsilon \geq \varepsilon_1 > 0$ such that (a, σ) is $(\varepsilon, \varepsilon_1)$ -regular at any point (t_0, x_0) and

$$\left(a_t^{jk}(x) - a_t^{jk}(y) \right) \xi^j \xi^k \geq \delta |\xi|^2$$

for all t, ξ, x , and y satisfying $|x - y| \leq \varepsilon$.

Assumption 2.3. There exists an $\varepsilon_2 > 0$ such that

$$(2.6) \quad \left| \sigma_t^{i\cdot}(x) - \sigma_t^{i\cdot}(y) \right|_{\ell_2} \leq \beta_1$$

for all i, t, x , and y satisfying $|x - y| \leq \varepsilon_2$.

Needless to say that Assumptions 2.2 and 2.3 are satisfied with any $\beta, \beta_1 > 0$ and slightly reduced δ if (2.3) holds and $a_t^{ij}(x)$ and $\sigma_t^{i\cdot}(x)$ are uniformly continuous in x uniformly with respect to (ω, t) .

Finally, we describe the space of initial data. Recall that for $p \geq 2$ the Slobodetskii space $W_p^{1-2/p} = W_p^{1-2/p}(\mathbb{R}^d)$ of functions $u_0(x)$ can be introduced as the space of traces on $t = 0$ of (deterministic) functions u such that

$$u \in L_p(\mathbb{R}_+, H_p^1), \quad \partial u / \partial t \in L_p(\mathbb{R}_+, H_p^{-1}),$$

where $\mathbb{R}_+ = (0, \infty)$. For such functions there is a (unique) modification denoted again u such that u_t is a continuous L_p -valued function on $[0, \infty)$ so that u_0 is well defined. Any such u_t is called an extension of u_0 .

The norm in $W_p^{1-2/p}$ can be defined as the infimum of

$$\|u\|_{L_p(\mathbb{R}_+, H_p^1)} + \|\partial u / \partial t\|_{L_p(\mathbb{R}_+, H_p^{-1})}$$

over all extensions u_t of elements u_0 . It is also well known that an equivalent norm of u_0 can be introduced as

$$\|u\|_{L_p((0,1), W_p^1)},$$

where $u = u_t$ is defined as the (unique) solution of the heat equation $\partial u_t(x) / \partial t = \Delta u_t(x)$ with initial condition $u_0(x)$.

For $s \geq 0$ we introduce

$$\text{tr}_s \mathcal{W}_p^1 = L_p \left(\Omega, \mathcal{F}_s, W_p^{1-2/p} \right).$$

The following auxiliary result helps understand the role of $\text{tr}_s \mathcal{W}_p^1$. We use spaces $\mathcal{W}_p^1([S, T])$ and $\mathbf{W}_p^1((S, T))$, which are introduced in the same way as $\mathcal{W}_p^1(\tau)$ and $\mathbf{W}_p^1(\tau)$, but the functions are considered only on $[S, T]$ and (S, T) , respectively.

LEMMA 2.1. *Let $s \geq 0$ be a fixed number and let u_s be an \mathcal{F}_s -measurable function with values in the set of distributions over \mathbb{R}^d .*

(i) *We have $u_s \in \text{tr}_s \mathcal{W}_p^1$ if and only if there exists a $v \in \mathcal{W}_p^1([s, \infty))$ satisfying the equation*

$$(2.7) \quad \partial v / \partial t = \Delta v - v, \quad t \geq s,$$

(which is a particular case of (1.1) and is understood in the same sense) with initial data u_s . This v is unique and satisfies

$$(2.8) \quad \|v\|_{\mathbf{W}_p^1((s, \infty))} \leq N \|u_s\|_{\text{tr}_s \mathcal{W}_p^1}, \quad \|u_s\|_{\text{tr}_s \mathcal{W}_p^1} \leq N \|v\|_{\mathbf{W}_p^1((s, \infty))},$$

where the constants N are independent of s, u_s , and v .

(ii) *We have $u_s \in \text{tr}_s \mathcal{W}_p^1$ if and only if there exists a $v \in \mathcal{W}_p^1([s, s+1])$ such that $v_s = u_s$.*

(iii) *If such a v exists and $dv_t = (D_i f_t^i + f_t^0) dt + g_t^k dw_t^k, t \geq s$, then*

$$(2.9) \quad \|u_s\|_{\text{tr}_s \mathcal{W}_p^1} \leq N \left(\|v\|_{\mathbf{W}_p^1((s, s+1))} + \sum_{j=0}^d \|f^j\|_{\mathbf{L}_p((s, s+1))} + \|g\|_{\mathbf{L}_p((s, s+1))} \right),$$

where the constant N is independent of s, u_s , and v .

(iv) *If $s > 0$ and we have a $u \in \mathcal{W}_p^1(s)$, then $u_s \in \text{tr}_s \mathcal{W}_p^1$ and*

$$\|u_s\|_{\text{tr}_s \mathcal{W}_p^1} \leq N \left(\|u\|_{\mathbf{W}_p^1(s)} + \sum_{j=0}^d \|f^j\|_{\mathbf{L}_p(s)} + \|g\|_{\mathbf{L}_p(s)} \right),$$

where N is independent of u , and f^j and g^k are taken from (2.2).

We prove this lemma in section 5.

Here are our main results concerning (1.1). The following theorem is very close to Theorem 2.12 of [2]. Important differences are the presence of the parameter λ in (2.10) and weaker assumptions on the coefficients of the deterministic part of the equation.

THEOREM 2.2. *Let the above assumptions be satisfied with $\beta = \beta(d, p, \delta) = \beta_0/3$, where β_0 is the constant from Lemma 5.1, and $\beta_1 = \beta_1(d, p, \delta, \varepsilon) > 0$ taken from the proof of Lemma 5.2. Let $\lambda \geq 0$, let $f^j, g \in \mathbf{L}_p(\tau)$, and let $u_0 \in \text{tr}_0 \mathcal{W}_p^1$.*

(i) *Then equation (1.1) for $t \leq \tau \wedge T$ has a unique solution $u \in \mathcal{W}_p^1(\tau \wedge T)$ with initial data u_0 and any $T \in (0, \infty)$. Moreover, if*

$$\lambda \geq \lambda_0(d, p, \delta, K, \varepsilon, \varepsilon_1, \varepsilon_2) \geq 1,$$

then (1.1) for $t \leq \tau$ has a unique solution $u \in \mathcal{W}_p^1(\tau)$ with initial data u_0 .

(ii) Furthermore, if a $v \in \mathcal{W}_p^1(\infty)$ is defined by (2.7) with initial condition u_0 , then the above solution u satisfies

$$(2.10) \quad \lambda^{1/2} \|u\|_{\mathbf{L}_p(\tau)} + \|Du\|_{\mathbf{L}_p(\tau)} \leq N \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} + \|Dv\|_{\mathbf{L}_p(\tau)} \right) + N\lambda^{-1/2} \|f^0\|_{\mathbf{L}_p(\tau)} + N\lambda^{1/2} \|v\|_{\mathbf{L}_p(\tau)},$$

provided that $\lambda \geq \lambda_0$, where the constants $N, \lambda_0 \geq 1$ depend only on $d, p, \delta, K, \varepsilon, \varepsilon_1$, and ε_2 .

(iii) Finally, there exists a set $\Omega' \subset \Omega$ of full probability such that $u_{t \wedge \tau} I_{\Omega'}$ is a continuous \mathcal{F}_t -adapted L_p -valued function of $t \in [0, \infty)$.

Observe that estimate (2.10) shows one of good reasons for writing the free term in (1.1) in the form $D_i f^i + f^0$, because $f^i, i = 1, \dots, d$, and f^0 enter (2.10) differently.

Remark 2.3. As it follows from our proofs, if $p = 2$, Assumptions 2.2 and 2.3 are not needed for Theorem 2.2 to be true and mentioning $\varepsilon, \varepsilon_1$, and ε_2 can be dropped in the statement. Thus we provide a new way to prove the classical result on Hilbert space solvability of SPDEs (cf., for instance, [13]).

We prove Theorem 2.2 in section 6 after we prepare necessary tools in sections 3–5. In section 3 we prove uniqueness part of Theorem 2.2 on the basis of Itô’s formula from [12]. Here Assumptions 2.2 and 2.3 are not used. In section 4 we treat the case of the heat equation with the random right-hand side and present a simplified version of the corresponding result from [6]. In section 5 we prove an auxiliary existence theorem and derive some a priori estimates.

Here is a result about continuous dependence of solutions on the data.

THEOREM 2.4. *Assume that for each $n = 1, 2, \dots$ we are given functions $a_{nt}^{ij}, a_{nt}^i, b_{nt}^i, c_{nt}, \sigma_{nt}^{ik}, \nu_{nt}^k, f_{nt}^j, g_{nt}^k$, and u_{n0} having the same meaning as the original ones and satisfying the same assumptions as those imposed on the original ones in Theorem 2.2 (with the same δ, K, β, \dots). Assume that for $i, j = 1, \dots, d$ and almost all (ω, t, x) we have*

$$\begin{aligned} \left(a_{nt}^{ij}, a_{nt}^i, b_{nt}^i, c_{nt} \right) &\rightarrow \left(a_t^{ij}, a_t^i, b_t^i, c_t \right), \\ \left| \sigma_{nt}^{i\cdot} - \sigma_t^{i\cdot} \right|_{\ell_2} + \left| \nu_{nt} - \nu_t \right|_{\ell_2} &\rightarrow 0, \end{aligned}$$

as $n \rightarrow \infty$. Also assume that

$$\sum_{j=0}^d \left\| f_n^j - f^j \right\|_{\mathbf{L}_p(\tau)} + \|g_n - g\|_{\mathbf{L}_p(\tau)} + \|u_{n0} - u_0\|_{\text{tr}_0 \mathcal{W}_p^1} \rightarrow 0$$

as $n \rightarrow \infty$. Take $\lambda \geq \lambda_0$, take the function u from Theorem 2.2, and let $u_n \in \mathcal{W}_p^1(\tau)$ be the unique solutions of (1.1) for $t \leq \tau$ constructed from $a_{nt}^{ij}, a_{nt}^i, b_{nt}^i, c_{nt}, \sigma_{nt}^{ik}, \nu_{nt}^k, f_{nt}^j$, and g_{nt}^k and having initial values u_{n0} .

Then, as $n \rightarrow \infty$, we have $\|u_n - u\|_{\mathbf{W}_p^1(\tau)} \rightarrow 0$ and for any finite $T \in [0, \infty)$

$$(2.11) \quad E \sup_{t \leq \tau \wedge T} \|u_{nt} - u_t\|_{L_p}^p \rightarrow 0.$$

Proof. Set $v_{nt} = u_{nt} - u_t$. Then

$$dv_{nt} = \left(L_{nt} v_{nt} - \lambda v_{nt} + D_i \tilde{f}_{nt}^i + \tilde{f}_{nt}^0 \right) dt + \left(\Lambda_{nt}^k v_{nt} + \tilde{g}_{nt}^k \right) dw_t^k,$$

where L_{nt} and Λ_{nt}^k are the operators constructed from a_{nt}^{ij} , a_{nt}^i , b_{nt}^i , c_{nt} , and σ_{nt}^{ik} , ν_{nt}^k , respectively, and

$$\begin{aligned} \tilde{f}_{nt}^i &= f_{nt}^i - f_t^i + (a_{nt}^{ji} - a_t^{ji}) D_j u_t + (a_{nt}^i - a_t^i) u_t, \\ \tilde{f}_{nt}^0 &= f_{nt}^0 - f_t^0 + (b_{nt}^i - b_t^i) D_i u_t + (c_{nt} - c_t) u_t, \\ \tilde{g}_{nt}^k &= g_{nt}^k - g_t^k + (\sigma_{nt}^{ik} - \sigma_t^{ik}) D_i u_t + (\nu_{nt}^k - \nu_t^k) u_t. \end{aligned}$$

By Theorem 2.2 we know that $u \in \mathbf{W}_p^1(\tau)$. This, along with our assumptions and the dominated convergence theorem, implies that

$$\sum_{j=0}^d \left\| \tilde{f}_n^j \right\|_{\mathbf{L}_p(\tau)} + \|\tilde{g}_n\|_{\mathbf{L}_p(\tau)} \rightarrow 0$$

as $n \rightarrow \infty$. After that by applying (2.10) to v_{nt} we immediately see that $\|u_n - u\|_{\mathbf{W}_p^1(\tau)} \rightarrow 0$.

Assertion (2.11) is, actually, a simple corollary of the above. Indeed, by introducing \hat{f}_n^j and \hat{g}_n^k in an obvious way, we can write

$$(2.12) \quad dv_{nt} = \left(D_i \hat{f}_{nt}^i + \hat{f}_{nt}^0 \right) dt + \hat{g}_{nt}^k dw_t^k,$$

and

$$\sum_{j=0}^d \left\| \hat{f}_n^j \right\|_{\mathbf{L}_p(\tau)} + \|\hat{g}_n\|_{\mathbf{L}_p(\tau)} \rightarrow 0.$$

It is standard (see, for instance, our Theorem 3.1) to derive from here the estimate

$$E \sup_{t \leq \tau \wedge T} \|u_{nt} - u_t\|_{L_p}^p \leq N \left(\sum_{j=0}^d \left\| \hat{f}_n^j \right\|_{\mathbf{L}_p(\tau \wedge T)} + \|\hat{g}_n\|_{\mathbf{L}_p(\tau \wedge T)} + E \|u_{n0} - u_0\|_{L_p}^p \right),$$

where N is independent of n . It is also well known that $W_p^{1-2/p} \subset L_p$, that is

$$\|u_{n0} - u_0\|_{L_p} \leq N \|u_{n0} - u_0\|_{W_p^{1-2/p}}.$$

By combining all this together we obtain (2.11) and the theorem is proved. \square

The following result could be proved on the basis of Theorem 2.4 in the same way as Corollary 5.11 of [6], where the solutions are approximated by solutions of equations with smooth coefficients and then a stopping time technique was used. We give here a shorter proof based on a different idea.

THEOREM 2.5. *Let $p_1, p_2 \in [2, \infty)$, $p_1 < p_2$, and let the above assumptions be satisfied with $\beta \leq \beta(d, p, \delta)$ for all $p \in [p_1, p_2]$ and $\beta_1 \leq \beta_1(d, p, \delta, \varepsilon)$ for all $p \in [p_1, p_2]$. Let $\lambda \geq 0$, and suppose that for $p \in [p_1, p_2]$ we have $f^j, g \in \mathbf{L}_p(\tau)$, and $u_0 \in \text{tr}_0 \mathcal{W}_p^1$.*

Then the solutions corresponding to $p = p_1$ and $p = p_2$ coincide, that is, there is a unique solution $u \in \mathcal{W}_{p_1}^1(\tau) \cap \mathcal{W}_{p_2}^1(\tau)$ of (1.1) with initial data u_0 .

Proof. Obviously, it suffices to concentrate on bounded τ . As is explained above in that case we may assume that λ is as large as we like. We take it so large that one could use assertion (ii) of Theorem 2.2 with any $p \in [p_1, p_2]$.

Denote by u the solution corresponding to $p = p_2$ and observe that, owing to uniqueness of solutions in $\mathcal{W}_{p_1}^1(\tau)$, we need only show that $u \in \mathcal{W}_{p_1}^1(\tau)$.

Take a nonnegative $\zeta \in C_0^\infty$ such that $\zeta(0) = 1$, set $\zeta_n(x) = \zeta(x/n)$, and notice that $u^n := u\zeta_n$ satisfies

$$du_t^n = (L_t u_t^n - \lambda u_t^n + D_i f_{nt}^i + f_{nt}^0) dt + (\Lambda_t^k u_t^n + g_{nt}^k) dw_t^k,$$

where

$$\begin{aligned} f_{nt}^i &= f_t^i \zeta_n - u a_t^{ji} D_j \zeta_n, \quad i \geq 1, \\ f_{nt}^0 &= f_t^0 \zeta_n - f_t^i D_i \zeta_n - \left(a_t^{ij} D_i u_t + a_t^j u \right) D_j \zeta_n - b_t^i u_t D_i \zeta_n, \\ g_{nt}^k &= g_t^k \zeta_n - \sigma_t^{ik} u_t D_i \zeta_n. \end{aligned}$$

It follows that for $p_1 \leq p \leq p_2$ we have

$$(2.13) \quad \|u^n\|_{\mathcal{W}_p^1(\tau)} \leq N \left(\sum_{i=0}^d \|f_n^i\|_{\mathbf{L}_p(\tau)} + \|g_n\|_{\mathbf{L}_p(\tau)} + \|u_0 \zeta_n\|_{\text{tr}_0 \mathcal{W}_p^1} \right).$$

One knows that with constants N independent of n

$$\|u_0 \zeta_n\|_{\text{tr}_0 \mathcal{W}_p^1} \leq N \left(\|u_0 \zeta_n\|_{\text{tr}_0 \mathcal{W}_{p_1}^1} + \|u_0 \zeta_n\|_{\text{tr}_0 \mathcal{W}_{p_2}^1} \right) \leq N \left(\|u_0\|_{\text{tr}_0 \mathcal{W}_{p_1}^1} + \|u_0\|_{\text{tr}_0 \mathcal{W}_{p_2}^1} \right).$$

Similarly, and by Hölder's inequality,

$$\|f_n^i\|_{\mathbf{L}_p(\tau)} \leq N + N \|u D \zeta_n\|_{\mathbf{L}_p(\tau)} \leq N + \|u\|_{\mathbf{L}_{p_2}(\tau)} \|D \zeta_n\|_{\mathbf{L}_q(\tau)},$$

where

$$q = \frac{pp_2}{p_2 - p}.$$

Similar estimates are available for other terms in the right-hand side of (2.13). Since

$$\|D \zeta_n\|_{\mathbf{L}_q(\tau)} = N n^{-1+(p_2-p)d/(p_2 p)} \rightarrow 0$$

as $n \rightarrow \infty$ if

$$(2.14) \quad \frac{1}{p} - \frac{1}{p_2} < \frac{1}{d},$$

estimate (2.13) implies that $u \in \mathcal{W}_p^1(\tau)$.

Thus, knowing that $u \in \mathcal{W}_{p_2}^1(\tau)$ allowed us to conclude that $u \in \mathcal{W}_p^1(\tau)$ as long as $p \in [p_1, p_2]$ and (2.14) holds. We can now replace p_2 with a smaller p and keep going in the same way each time increasing $1/p$ by the same amount until p reaches p_1 . Then we get that $u \in \mathcal{W}_{p_1}^1(\tau)$. The theorem is proved. \square

In many situations the following maximum principle is useful.

THEOREM 2.6. *Let the above assumptions be satisfied with $\beta \leq \beta(d, q, \delta)$ for all $q \in [2, p]$ and $\beta_1 \leq \beta_1(d, q, \delta, \varepsilon)$ for all $q \in [2, p]$. Let $\lambda \geq 0$ and $f^0 \in \mathbf{L}_p(\tau)$, $u_0 \in \text{tr}_0 \mathcal{W}_p^1$, $f^i = 0$, $i = 1, \dots, d$, $g = 0$ be such that $u_0 \geq 0$ and $f^0 \geq 0$. Then for the solution u almost surely we have $u_t \geq 0$ for all finite $t \leq \tau$.*

Proof. If $p = 2$ the result is proved in [9]. For general $p \geq 2$ take the same function ζ_n as in the preceding proof, introduce $f^{ni} = f^i \zeta_n$, $g_n^k = 0$, and call u^n the solution of (1.1) with so modified free terms and the initial data $u_0 \zeta_n$. By Theorem 2.5 we have $u^n \in \mathcal{W}_p^1(\tau) \cap \mathcal{W}_2^1(\tau)$. By the above, $u^n \geq 0$ and it only remains to use Theorem 2.4. The theorem is proved. \square

3. Itô’s formula and uniqueness. The following two “standard” results are taken from [12].

THEOREM 3.1. *Let $u \in \mathcal{W}_p^1(\tau)$, $f^j \in \mathbf{L}_p(\tau)$, $g = (g^k) \in \mathbf{L}_p(\tau)$, and assume that (2.2) holds for $t \leq \tau$ in the sense of generalized functions. Then there is a set $\Omega' \subset \Omega$ of full probability such that*

- (i) $u_{t \wedge \tau} I_{\Omega'}$ is a continuous L_p -valued \mathcal{F}_t -adapted function on $[0, \infty)$;
- (ii) for all $t \in [0, \infty)$ and $\omega \in \Omega'$ Itô’s formula holds:

$$(3.1) \quad \begin{aligned} \int_{\mathbb{R}^d} |u_{t \wedge \tau}|^p dx &= \int_{\mathbb{R}^d} |u_0|^p dx + p \int_0^{t \wedge \tau} \int_{\mathbb{R}^d} |u_s|^{p-2} u_s g_s^k dx dw_s^k \\ &+ \int_0^{t \wedge \tau} \left(\int_{\mathbb{R}^d} [p|u_t|^{p-2} u_t f_t^0 - p(p-1)|u_t|^{p-2} f_t^i D_i u_t \right. \\ &\left. + (1/2)p(p-1)|u_t|^{p-2} |g_t|_{\ell_2}^2] dx \right) dt. \end{aligned}$$

Furthermore, for any $T \in [0, \infty)$

$$(3.2) \quad \begin{aligned} E \sup_{t \leq \tau \wedge T} \|u_t\|_{L_p}^p &\leq 2E\|u_0\|_{L_p}^p + NT^{p-1} \|f^0\|_{\mathbf{L}_p(\tau)}^p \\ &+ NT^{(p-2)/2} \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)}^p + \|g\|_{\mathbf{L}_p(\tau)}^p + \|Du\|_{\mathbf{L}_p(\tau)}^p \right), \end{aligned}$$

where $N = N(d, p)$.

Here is an “energy” estimate.

COROLLARY 3.2. *Under the conditions of Theorem 3.1 assume that $\tau < \infty$ (a.s.).*

Then

$$(3.3) \quad \begin{aligned} E \int_{\mathbb{R}^d} |u_0|^p dx + E \int_0^\tau \left(\int_{\mathbb{R}^d} [p|u_t|^{p-2} u_t f_t^0 - p(p-1)|u_t|^{p-2} f_t^i D_i u_t \right. \\ \left. + (1/2)p(p-1)|u_t|^{p-2} |g_t|_{\ell_2}^2] dx \right) dt \geq EI_{\tau < \infty} \int_{\mathbb{R}^d} |u_\tau|^p dx. \end{aligned}$$

Furthermore, if τ is bounded, then there is an equality instead of inequality in (3.3).

The next result implies, in particular, uniqueness in Theorem 2.2.

LEMMA 3.3. *Under Assumption 2.1 there exist $\lambda_0 \geq 0$ and N depending only on d, p, K , and δ such that, for any strictly positive $\lambda \geq \lambda_0$ and any solution $u \in \mathcal{W}_{p,0}^1(\tau)$ of (1.1) for $t \leq \tau$, we have*

$$(3.4) \quad \lambda \|u\|_{\mathbf{L}_p(\tau)} \leq N \lambda^{1/2} \left(\sum_{j=1}^d \|f^j\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} \right) + N \|f^0\|_{\mathbf{L}_p(\tau)}.$$

Furthermore, if $a^i = b^i = \nu^k \equiv 0$, then one can take $\lambda_0 = 0$.

Proof. We may assume that $f^j \in \mathbf{L}_p(\tau)$, $g = (g^k) \in \mathbf{L}_p(\tau)$, since otherwise the right-hand side of (3.4) is infinite.

If (3.4) is true for $\tau \wedge T$ in place of τ and any $T \in (0, \infty)$, then it is obviously also true as is. Therefore, we may assume that τ is finite. An advantage of this assumption is that we can use Corollary 3.2. Write (3.3) with \hat{f}_t^i, \hat{f}_t^0 , and \hat{g}_t^k in place of f_t^i, f_t^0 , and g_t^k , respectively, where

$$\begin{aligned} \hat{f}_t^i &= a_t^{j_i} D_j u_t + a_t^i u_t + f_t^i, \\ \hat{f}_t^0 &= b_t^i D_i u_t + (c_t - \lambda) u_t + f_t^0, \quad \hat{g}_t^k = \sigma_t^{ik} D_i u_t + \nu_t^k u_t + g_t^k. \end{aligned}$$

Then observe that inequalities like $(a + b)^2 \leq (1 + \varepsilon)a^2 + (1 + \varepsilon^{-1})b^2$ show that for any $\varepsilon \in (0, 1]$ we have

$$\begin{aligned} |\hat{g}_t|_{\ell_2}^2 &\leq (1 + \varepsilon) \left| \sum_{i=1}^d \sigma_t^{i \cdot} D_i u_t \right|_{\ell_2}^2 + 2\varepsilon^{-1} |\nu_t u_t + g_t|_{\ell_2}^2 \\ &\leq 2(1 + \varepsilon) \alpha_t^{ij} (D_i u_t) D_j u_t + N\varepsilon^{-1} (|u_t|^2 + |g_t|_{\ell_2}^2). \end{aligned}$$

Owing to (2.3), for $\varepsilon = \varepsilon(\delta) > 0$ small enough

$$\begin{aligned} (3.5) \quad I_t &:= (1/2)|u_t|^{p-2} |\hat{g}_t|_{\ell_2}^2 - |u_t|^{p-2} \hat{f}_t^i D_i u_t + (p - 1)^{-1} |u_t|^{p-2} u_t b_t^i D_i u_t \\ &\leq -(\delta/2)|u_t|^{p-2} |Du_t|^2 \\ &\quad + N|u_t|^{p-2} \left(|u_t|^2 + |g_t|_{\ell_2}^2 + |Du_t| |u_t| + |Du_t| \sum_{i=1}^d |f_t^i| \right). \end{aligned}$$

Next, we use that for any $\gamma > 0$

$$\begin{aligned} |u_t|^{p-1} |Du_t| &= \left(|u_t|^{(p-2)/2} |Du_t| \right) |u_t|^{p/2} \leq \gamma |u_t|^{p-2} |Du_t|^2 + \gamma^{-1} |u_t|^p, \\ |u_t|^{p-2} |Du_t| |f_t^i| &\leq \gamma |u_t|^{p-2} |Du_t|^2 + \gamma^{-1} |u_t|^{p-2} |f_t^i|^2, \end{aligned}$$

and by choosing γ appropriately find from (3.5) that

$$(3.6) \quad I_t \leq N|u_t|^p + N|u_t|^{p-2} \left(\sum_{i=1}^d |f_t^i|^2 + |g_t|_{\ell_2}^2 \right).$$

After that Hölder’s inequality and (3.3), where the right-hand side is nonnegative, immediately lead to

$$(\lambda - N_1) \|u\|_{\mathbf{L}_p(\tau)}^p \leq N \|u\|_{\mathbf{L}_p(\tau)}^{p-2} \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)}^2 + \|g\|_{\mathbf{L}_p(\tau)}^2 \right) + N \|u\|_{\mathbf{L}_p(\tau)}^{p-1} \|f^0\|_{\mathbf{L}_p(\tau)}.$$

Furthermore, simple inspection of the above argument shows that, if $a^i = b^i = \nu^k \equiv 0$, then the terms with $|u_t|^2$ and $|u_t| |Du_t|$ in (3.5) and the term with $|u_t|^p$ in (3.6) disappear, so that we can take $N_1 = 0$ in this case (recall that $c \leq 0$). Generally, for $\lambda \geq 2N_1$ we have $\lambda - N_1 \geq (1/2)\lambda$ and

$$\bar{U}^p \leq N\bar{U}^{p-2} \bar{G}^2 + N\bar{U}^{p-1} \bar{F},$$

where

$$\bar{U} = \lambda \|u\|_{\mathbf{L}_p(\tau)}, \quad \bar{G} = \lambda^{1/2} \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} \right), \quad \bar{F} = \|f^0\|_{\mathbf{L}_p(\tau)}.$$

It follows that $\bar{U} \leq N(\bar{G} + \bar{F})$, which is (3.4) and the lemma is proved. \square

4. Case of the heat equation. To move further we need the following analytic fact established in [4] (see also [7] for a complete proof).

LEMMA 4.1. *Denote by T_t the heat semigroup in \mathbb{R}^d and let $p \geq 2$, $-\infty \leq a < b \leq \infty$, $g \in L_p((a, b) \times \mathbb{R}^d, \ell_2)$. Then*

$$\int_{\mathbb{R}^d} \int_a^b \left[\int_a^t |DT_{t-s} g_s(x)|_{\ell_2}^2 ds \right]^{p/2} dt dx \leq N(d, p) \int_{\mathbb{R}^d} \int_a^b |g_t(x)|_{\ell_2}^p dt dx.$$

In this section we deal with the following model equation

$$(4.1) \quad du_t = \Delta u_t dt + g_t^k dw_t^k.$$

LEMMA 4.2. *Assume that $\tau \leq T$, where the constant $T \in [0, \infty)$. Then for any $g = (g^1, g^2, \dots) \in \mathbf{L}_p(\tau)$ there exists a unique $u \in \mathcal{W}_{p,0}^1(\tau)$ satisfying (4.1) for $t \leq \tau$. Furthermore, for this solution we have*

$$(4.2) \quad E \sup_{t \leq \tau} \|u_t\|_{L_p}^p \leq N(d, p) T^{(p-2)/2} \|g\|_{\mathbf{L}_p(\tau)}^p,$$

$$(4.3) \quad \|Du\|_{\mathbf{L}_p(\tau)} \leq N(d, p) \|g\|_{\mathbf{L}_p(\tau)}.$$

Proof. By replacing the unknown function u_t with $v_t e^{\lambda t}$ we see that v_t satisfies

$$dv_t = (\Delta v_t - \lambda v_t) dt + e^{-\lambda t} g_t^k dw_t^k.$$

Since τ is bounded, the inclusions $u \in \mathcal{W}_{p,0}^1(\tau)$ and $v \in \mathcal{W}_{p,0}^1(\tau)$ are equivalent and our assertion about uniqueness follows from Lemma 3.3.

In the proof of existence we borrow part of the proof of Theorem 4.2 of [6]. As we have pointed out in the Introduction, the beginning of the theory of divergence and nondivergence type equations is the same. The main difference with that proof is that here we take $f \equiv 0$.

We take an integer $m \geq 1$, some stopping times $\tau_0 \leq \tau_1 \leq \dots \leq \tau_m \leq T$, and some (nonrandom) functions $g^{ij} \in C_0^\infty$, $i, j = 1, \dots, m$. Then we define

$$g_t^k(x) = \sum_{i=1}^m g^{ik}(x) I_{(\tau_{i-1}, \tau_i]}(t),$$

$$v_t(x) = \sum_{k=1}^m \int_0^t g_s^k(x) dw_s^k = \sum_{i,k=1}^m g^{ik}(x) \left(w_{t \wedge \tau_i}^k - w_{t \wedge \tau_{i-1}}^k \right), \quad t \geq 0.$$

Obviously, for any ω , the function $v_t(x)$ is continuous and bounded in (t, x) along with any derivative in x . Furthermore, the function and its derivative in x are Hölder 1/3 continuous in t uniformly with respect to x (for almost any ω). Also, $v_t(x)$ has compact support in x .

These properties of $v_t(x)$ imply that for any ω there exists a unique classical solution of the heat equation

$$\frac{\partial}{\partial t} \bar{u}_t = \Delta \bar{u}_t + \Delta v_t, \quad t > 0,$$

with zero initial data. Furthermore,

$$(4.4) \quad \bar{u}_t(x) = \int_0^t T_{t-s} \Delta v_s(x) ds.$$

This formula shows, in particular, that $\bar{u}_t(x)$ is \mathcal{F}_t -adapted. Adding the fact that \bar{u}_t is continuous in t proves that $\bar{u}_t(x)$ is predictable. The same holds for

$$(\bar{u}_t, \phi) = \int_0^t (T_{t-s} \Delta v_s, \phi) ds$$

with any $\phi \in C_0^\infty$. The following corollary of Minkowski's inequality

$$(4.5) \quad \|\bar{u}_t\|_{L_p} \leq \int_0^t \|\Delta v_s\|_{L_p} ds$$

shows that \bar{u}_t is L_p -valued. Since (\bar{u}_t, ϕ) is predictable for any $\phi \in C_0^\infty$, \bar{u}_t is weakly and, hence, strongly predictable as an L_p -valued process.

One can differentiate (4.4) with respect to x as many times as one wants and get similar statements about the derivatives of \bar{u}_t . In particular, (4.5) implies that for any multi-index α

$$E \int_0^T \int_{\mathbb{R}^d} |D^\alpha \bar{u}_t|^p dx dt \leq T^p E \int_0^T \int_{\mathbb{R}^d} |D^\alpha \Delta v_t|^p dx dt < \infty,$$

so that $\bar{u}_t \in \mathcal{W}_{p,0}^1(T)$.

Now, it is easily seen that

$$u_t(x) := \bar{u}_t(x) + v_t(x)$$

satisfies (4.1) pointwisely and by the above $u_t \in \mathcal{W}_{p,0}^1(T)$. The (deterministic) Fubini's theorem also shows that u_t satisfies (4.1) in the sense of distributions.

Next, we use the same simple transformation as in the proof of Lemma 4.1 of [6] and conclude that for any t and x almost surely

$$Du_t(x) = \sum_{k=1}^m \int_0^t T_{t-s} Dg_s^k(x) dw_s^k.$$

Hence, by Burkholder–Davis–Gundy inequality

$$E|Du_t(x)|^p \leq NE \left[\int_0^t |T_{t-s} Dg_s(x)|_{\ell_2}^2 ds \right]^{p/2},$$

which along with Lemma 4.1 proves (4.3) for our particular g . Theorem 3.1 shows that (4.2) follows from (4.3) and (4.1).

The rest is trivial since the set of g 's like the one above is dense in $\mathbf{L}_p(T)$ by Theorem 3.10 of [6]. The lemma is proved. \square

Next we introduce the parameter λ into (4.1).

LEMMA 4.3. *Assume that $\tau \leq T$, where the constant $T \in [0, \infty)$. Let $\lambda > 0$. Then for any $g = (g^1, g^2, \dots) \in \mathbf{L}_p(\tau)$ there exists a unique $u \in \mathcal{W}_{p,0}^1(\tau)$ satisfying*

$$(4.6) \quad du_t = (\Delta u_t - \lambda u_t) dt + g_t^k dw_t^k$$

for $t \leq \tau$. Furthermore, for this solution we have

$$(4.7) \quad \lambda^{p/2} \|u\|_{\mathbf{L}_p(\tau)}^p \leq N(d, p) \|g\|_{\mathbf{L}_p(\tau)}^p,$$

$$(4.8) \quad \|Du\|_{\mathbf{L}_p(\tau)} \leq N(d, p) \|g\|_{\mathbf{L}_p(\tau)}.$$

Proof. Uniqueness and estimate (4.7) follow from Lemma 3.3. The existence immediately follows from Lemma 4.2 and the result of transformation described in the beginning of its proof. To establish (4.8) consider the heat equation

$$(4.9) \quad \frac{\partial}{\partial t} v_t = \Delta v_t - \lambda v_t.$$

Since $u \in \mathbf{L}_p(\tau)$, for almost any ω we have $u \in L_p((0, \tau) \times \mathbb{R}^d)$ and by a classical result (see, for instance, [11]) for almost any ω equation (4.9) with zero initial data has a unique solution in the class of functions such that along with derivatives in x up to the second order they belong to $L_p((0, \tau) \times \mathbb{R}^d)$. Furthermore, after writing the term $-\lambda u_t$ as $-\lambda v_t - \lambda z_t$, where $z_t = u_t - v_t$, differentiating the equation once with respect to x and using the contraction property of the heat semigroup we find

$$(4.10) \quad \|Dv\|_{L_p((0,\tau)\times\mathbb{R}^d)}^p \leq \|Dz\|_{L_p((0,\tau)\times\mathbb{R}^d)}^p.$$

The solution v_t can be given by an integral formula, which implies that v_t is \mathcal{F}_t -adapted. It is also continuous as an L_p -valued process, and hence, is a predictable L_p -valued process. Taking expectations of both parts of (4.10) shows that $v \in \mathcal{W}_p^1(\tau)$ if $z \in \mathcal{W}_p^1(\tau)$.

Now observe that

$$dz_t = \Delta z_t dt + g_t^k dw_t^k,$$

which by Lemma 4.2 implies that $z \in \mathcal{W}_p^1(\tau)$ and

$$\|Dz\|_{\mathbf{L}_p(\tau)}^p \leq N \|g\|_{\mathbf{L}_p(\tau)}^p.$$

Upon combining this with (4.10) and using the fact that $u = v + z$, we come to (4.8). The lemma is proved. \square

5. A priori estimates in the general case. First we deal with the case when $\sigma = \nu = 0$.

LEMMA 5.1. *Suppose that $\sigma^{ik} \equiv \nu^k \equiv 0$. Also suppose that Assumptions 2.1 and 2.2 are satisfied with $\beta \leq \beta_0$, where the way to estimate the constant $\beta_0(d, p, \delta) > 0$ is described in the proof. Let $f^j \in \mathbf{L}_p(\tau)$ and $g \in \mathbf{L}_p(\tau)$.*

Then there exist constants $\lambda_0 \geq 1$ and N , depending only on d, p, δ, K , and ε , such that for any $\lambda \geq \lambda_0$ there exists a unique $u \in \mathcal{W}_{p,0}^1(\tau)$ satisfying (1.1) for $t \leq \tau$. Furthermore, this solution satisfies the estimate

$$(5.1) \quad \lambda^{1/2} \|u\|_{\mathbf{L}_p(\tau)} + \|Du\|_{\mathbf{L}_p(\tau)} \leq N \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} \right) + N \lambda^{-1/2} \|f^0\|_{\mathbf{L}_p(\tau)}.$$

Proof. Uniqueness and part of estimate (5.1) follow from Lemma 3.3. In the rest of the proof we may assume that τ is bounded and split our argument into two parts.

Case $g^k \equiv 0$. First assume that the coefficients and f^j are nonrandom. We extend the coefficients of L following the example $a_t^{ij}(x) = \delta^{ij}$, $t < 0$, and extend f_t^j beyond $(0, \tau)$ arbitrary only requiring $f^j \in L_p(\mathbb{R}^{d+1})$.

Then by Theorem 4.5 and Remark 2.4 of [8] the equation

$$(5.2) \quad \frac{\partial}{\partial t} u_t = L_t u_t - \lambda u_t + D_i f_t^i + f_t^0$$

in \mathbb{R}^{d+1} has a unique solution with finite norms

$$\|u\|_{L_p(\mathbb{R}^{d+1})} \quad \text{and} \quad \|Du\|_{L_p(\mathbb{R}^{d+1})}$$

provided that $\lambda \geq \lambda_0$. By Theorem 4.4 of [8]

(5.3)

$$\lambda^{1/2}\|u\|_{L_p(\mathbb{R}^{d+1})} + \|Du\|_{L_p(\mathbb{R}^{d+1})} \leq N \left(\sum_{i=1}^d \|f^i\|_{L_p(\mathbb{R}^{d+1})} + \lambda^{-1/2} \|f^0\|_{L_p(\mathbb{R}^{d+1})} \right).$$

By Theorem 3.1 the function u_t is a continuous L_p -valued function.

The proof of Theorem 4.4 of [8] is achieved on the basis of the a priori estimate (5.3) and the method of continuity by considering the family of equations

$$(5.4) \quad \frac{\partial}{\partial t} u_t = (\theta L_t + (1 - \theta)\Delta)u_t - \lambda u_t + D_i f_t^i + f_t^0,$$

where the parameter θ changes in $[0, 1]$. We remind briefly the method of continuity because we want to show that certain properties of (5.4) which we know for $\theta = 0$ propagate from $\theta = 0$ to $\theta = 1$.

We fix a $\theta_0 \in [0, 1]$ and to solve (5.4) for given f^j define a sequence of $u^n \in L_p(\mathbb{R}, W_p^1)$ by solving the equation

$$(5.5) \quad \begin{aligned} \frac{\partial}{\partial t} u_t^{n+1} &= (\theta_0 L_t + (1 - \theta_0)\Delta)u_t^{n+1} - \lambda u_t^{n+1} \\ &+ D_i f_t^i + f_t^0 + (\theta - \theta_0)(L_t - \Delta)u^n, \quad n \geq 1, \quad u^0 = 0. \end{aligned}$$

If we know that (5.4) is uniquely solvable with θ_0 in place of θ for arbitrary $f^j \in L_p(\mathbb{R}^{d+1})$, then the sequence u^n is well defined. Furthermore, estimate (5.3) easily shows that for θ sufficiently close to θ_0 the $L_p(\mathbb{R}, W_p^1)$ norm of $u^{n+1} - u^n$ goes to zero geometrically as $n \rightarrow \infty$. In this way, passing to the limit in (5.5), we obtain the solution of (5.4) for θ close to θ_0 . Then we can repeat the procedure and starting from $\theta = 0$ and moving step by step eventually reach $\theta = 1$.

For $\theta = 0$ we are dealing with solvability of the heat equation which is proved by giving the solution explicitly by means of the heat semigroup. This representation formula has two important implications:

- (i) For any constant $T \in \mathbb{R}$, changing f_t^j for $t \geq T$ does not affect u_t for $t \leq T$;
- (ii) If f^j are $L_p(\mathbb{R}^{d+1})$ -valued measurable functions of a parameter, say ω from a measurable space, say (Ω, \mathcal{F}_T) , then the solution $u \in L_p(\mathbb{R}, W_p^1)$, which now depends on ω is also \mathcal{F}_T -measurable.

Property (i) is obtained by inspecting the representation formula. Property (ii) is true because the mapping $L_p(\mathbb{R}^{d+1}) \ni f^j \rightarrow u \in L_p(\mathbb{R}, W_p^1)$ is continuous and, hence, Borel measurable.

Obviously, both properties propagate from $\theta = 0$ to $\theta = 1$ by the above method of continuity. In particular, solutions of (5.2) on the time interval $(-\infty, T]$ depend only on the values of f_t^j for $t \in (-\infty, T]$. It follows that with the same λ and N , for any $T \in \mathbb{R}$,

$$(5.6) \quad \begin{aligned} \lambda^{1/2}\|u\|_{L_p((-\infty, T), L_p)} + \|Du\|_{L_p((-\infty, T), L_p)} \\ \leq N \left(\sum_{i=1}^d \|f^i\|_{L_p((-\infty, T), L_p)} + \lambda^{-1/2} \|f^0\|_{L_p((-\infty, T), L_p)} \right). \end{aligned}$$

From now on, we allow the coefficients and f^j to be random, continue f^j as zero for $t < 0$, and solve (5.2) for each ω . By (5.6) with $T = 0$ we have that $u_t = 0$ for $t \leq 0$, and it makes sense considering (5.2) on $(0, T)$ for each $T \in (0, \infty)$ with zero initial condition. In such situation properties (i) and (ii) still hold.

In particular, if f^j are measurable $L_p((0, T), L_p)$ -valued functions of a parameter, say ω from a measurable space, say (Ω, \mathcal{F}_T) , then the solution $u \in L_p((0, T), W_p^1)$ is also \mathcal{F}_T -measurable. Then, from the equation itself, it follows that (u_T, ϕ) is \mathcal{F}_T -measurable for any $\phi \in C_0^\infty$. Since u_T takes values in L_p , it is an L_p -valued \mathcal{F}_T -measurable function.

If f_t^i are predictable L_p -valued functions, the above conclusions are valid for any $T \in [0, \infty)$. In particular, u_t is \mathcal{F}_t -adapted as an L_p -valued function and since it is continuous, u_t is a predictable L_p -valued function.

These properties and the fact that (5.6) holds for any $T \in (0, \infty)$ and ω prove the lemma in the particular case under consideration.

General case. By Lemma 4.3 there is a unique solution $v \in \mathcal{W}_{p,0}^1(\tau)$ of (4.6). Observe that

$$(L_t - \Delta)v_t = D_i \hat{f}_t^i + \hat{f}_t^0,$$

where \hat{f}_t^j are functions of class $\mathbf{L}_p(\tau)$ defined by

$$\begin{aligned} \hat{f}_t^j &= \left(a_t^{ij} - \delta^{ij} \right) D_i v_t + a_t^j v_t, \quad j = 1, \dots, d, \\ \hat{f}_t^0 &= b_t^i D_i v_t + c_t v_t. \end{aligned}$$

By the above there is a unique solution $u \in \mathcal{W}_{p,0}^1(\tau)$ of

$$\frac{\partial}{\partial t} u_t = L_t u_t - \lambda u_t + (L_t - \Delta)v_t + D_i f_t^i + f_t^0.$$

Obviously, $v_t + u_t$ is a solution of class $\mathcal{W}_{p,0}^1(\tau)$ of (1.1). By the particular case

$$\begin{aligned} \lambda^{1/2} \|u\|_{\mathbf{L}_p(\tau)} + \|Du\|_{\mathbf{L}_p(\tau)} &\leq N \sum_{i=1}^d \left(\|f^i\|_{\mathbf{L}_p(\tau)} + \|\hat{f}^i\|_{\mathbf{L}_p(\tau)} \right) \\ &\quad + N \lambda^{-1/2} \left(\|f^0\|_{\mathbf{L}_p(\tau)} + \|\hat{f}^0\|_{\mathbf{L}_p(\tau)} \right) \end{aligned}$$

and to obtain (5.1), it remains only to use the estimates of v_t provided by Lemma 4.3. The lemma is proved. \square

Now we allow $\sigma \neq 0$.

LEMMA 5.2. (i) *Suppose that Assumption 2.1 is satisfied with $K = 0$ and take $\varepsilon \geq \varepsilon_1 > 0$, $\varepsilon_2 \in (0, \varepsilon/4]$, $t_0 \geq 0$, and $x_0 \in \mathbb{R}^d$.*

(ii) *Let $f^j \in \mathbf{L}_p(\tau)$, $g \in \mathbf{L}_p(\tau)$, and $u \in \mathcal{W}_{p,0}^1(\tau)$ be such that (1.1) holds for $t \leq \tau$. Assume that $u_t(x) = 0$ if*

$$(t, x) \notin \Gamma := (t_0, t_0 + \varepsilon_1^2) \times B_{\varepsilon_2}(x_0).$$

(iii) *Assume that the couple (a, σ) is $(\varepsilon, \varepsilon_1)$ -regular at (t_0, x_0) with $\beta = \beta_0/3$ in (2.4) and (2.5), where β_0 is the constant from Lemma 5.1. Also assume that*

$$|\sigma_t^i(x) - \sigma_t^i(x_0)|_{\ell_2} \leq \beta_1, \quad (a_t^{jk}(y) - \alpha_t^{jk}(x_0)) \xi^j \xi^k \geq \delta |\xi|^2$$

for all values of indices and arguments such that $(t, x) \in \Gamma$ and $(t, y) \in Q_\varepsilon(t_0, x_0)$, where $\beta_1 = \beta_1(d, \delta, p, \varepsilon) > 0$ is a constant an estimate from below for which can be obtained from the proof.

Then there exist constants $\lambda_0 \geq 1$ and N , depending only on d, p, δ , and ε , such that estimate (5.1) holds provided that $\lambda \geq \lambda_0$.

Proof. Without loss of generality we may and will assume that $x_0 = 0$. Also we modify, if necessary, a and σ in such a way that $\sigma_t^{ik}(x) = 0$ if $t \notin (t_0, t_0 + \varepsilon_1^2)$, and $a_t^{ij}(x) = \delta^{-1}\delta^{ij}$ if $t \notin (t_0, t_0 + \varepsilon_1^2)$. Obviously, under this modification assumption (iii) is preserved and (1.1) remains unaffected due to assumption (ii). The rest of the proof we split into two cases.

Case $\sigma_t^{ik}(x) = \sigma_t^{ik}(0)$ for $|x| \leq \varepsilon_2$ and $t \geq 0$. We want to apply Lemma 5.1 and for that, even if $\sigma \equiv 0$, we need a^{ij} to satisfy at least the condition $\text{osc}(a^{ij}, Q) \leq \beta$ for all $Q \in \mathbf{Q}$ with $\rho(Q) \leq \varepsilon$. To achieve this we modify $a_t^{ij}(x)$ for $|x| \geq \varepsilon/4$ using the fact that such modifications have no effect on (1.1) since $u_t(x) = 0$ for $|x| \geq \varepsilon_2$ and $\varepsilon_2 \leq \varepsilon/4$.

Take a $\xi \in C_0^\infty(\mathbf{R}^d)$ with support lying in the ball of radius $\varepsilon/2$ centered at the origin and such that $\xi(x) = 1$ for $|x| \leq \varepsilon/4$ and $0 \leq \xi \leq 1$. Set

$$\hat{a}_t^{ij} := \xi a_t^{ij} + \delta^{-1}(1 - \xi)\delta^{ij}.$$

We can use \hat{a} in place of a in (1.1). It follows by Lemma 4.7 of [6] (Itô–Wentzell formula) that the function $v_t(x) := u_t(x + x_t)$ satisfies the equation

$$(5.7) \quad dv_t(x) = (\bar{L}_t v_t(x) - \lambda v_t + D_i \bar{f}_t^i + \bar{f}_t^0) dt + \bar{g}_t^k(x + x_t) dw_t^k,$$

where

$$\begin{aligned} \bar{L}_t \phi &= D_j \left(\bar{a}_t^{ij} D_i \phi \right), \quad \bar{a}_t^{ij}(x) = \hat{a}_t^{ij}(x + x_t) - \alpha_t^{ij}(0), \\ \bar{f}_t^i(x) &:= f_t^i(x + x_t) - \sigma_t^{ik}(0)g_t^k(x + x_t), \quad i = 1, \dots, d, \\ \bar{f}_t^0(x) &:= f_t^0(x + x_t), \quad \bar{g}_t^k(x) = g_t^k(x + x_t), \end{aligned}$$

and the process $x_t = (x_t^1, \dots, x_t^d)$ is defined by

$$x_t^i = - \int_0^t \sigma_s^{ik}(0) dw_s^k.$$

This fact shows that the assertion of the present lemma is a direct consequence of Lemma 5.1 in case the latter is applicable to (5.7).

As is easy to see we will be able to apply Lemma 5.1 to (5.7) if we can find $\varepsilon' = \varepsilon'(d, \delta, \varepsilon, p) > 0$ such that

$$(5.8) \quad \frac{1}{t-s} \int_s^t \left(\left| \bar{a}_r^{ij} - \bar{a}_{r(B)}^{ij} \right| \right)_{(B)} dr \leq \beta_0,$$

whenever $(s, t) \times B \in \mathbf{Q}$ and $\rho(B) \leq \varepsilon'$.

Denote by N , with or without subscripts, various (large) constants depending only on d, δ , and ε and observe that $|D\xi| \leq N$. It follows easily that for $B \in \mathbf{B}$ we have

$$(5.9) \quad \begin{aligned} \left(\left| \bar{a}_r^{ij} - \bar{a}_{r(B)}^{ij} \right| \right)_{(B)} &\leq \left(\left| \xi a_r^{ij} - (\xi a_r^{ij})_{(B+x_r)} \right| \right)_{(B+x_r)} + \delta^{-1}\delta^{ij}(|\xi - \xi_{(B+x_r)}|)_{(B+x_r)} \\ &\leq \left(\left| \xi a_r^{ij} - (\xi a_r^{ij})_{(B+x_r)} \right| \right)_{(B+x_r)} + N_1 \rho =: I_r + N_1 \rho, \end{aligned}$$

where and below $\rho = \rho(B)$.

Let z be the center of B and set

$$y_r = (z + x_r)(\rho + \varepsilon/2)|z + x_r|^{-1}$$

if $|z + x_r| \geq \rho + \varepsilon/2$ and $y_r = z + x_r$ otherwise. Observe that y_r is continuous in r and

$$(5.10) \quad |y_r| \leq \rho + \varepsilon/2.$$

Next we claim that

$$(5.11) \quad I_r \leq 2 \left(\left| a_r^{ij} - a_{r(B_\rho+y_r)}^{ij} \right| \right)_{(B_\rho+y_r)} + N_2\rho.$$

If (5.11) is true, then by combining it with (5.9) and using (5.10) we find that the left-hand side of (5.8) is less than

$$(N_1 + N_2)\rho + 2 \sup_{|y| \leq \rho + \varepsilon/2} \text{osc}(a^{ij}, Q_\rho + (0, y), 0)$$

if $\sigma_t^{nm}(0) = 0$ for all t, n, m or, in general, less than

$$(N_1 + N_2)\rho + 2\text{Osc}(a^{ij}, Q_\rho, \rho + \varepsilon/2),$$

where $Q_\rho = (s, t) \times B_\rho$. Now (2.4) and (2.5) imply that (5.8) is satisfied for $\rho \leq \varepsilon'$ if we choose $\varepsilon' > 0$ so that

$$(N_1 + N_2)\varepsilon' \leq \beta_0/3, \quad \varepsilon' \leq \varepsilon/4.$$

Therefore, it remains only to prove the claim. Obviously, if $|z + x_r| \geq \rho + \varepsilon/2$, then $I_r = 0$ and (5.11) holds.

In case $|z + x_r| < \rho + \varepsilon/2$ the estimates

$$\begin{aligned} (|h_r - h_{r(B')}|)_{(B')} &\leq \frac{1}{|B'|^2} \int_{B'} \int_{B'} |h_r(y) - h_r(z)| \, dydz \leq 2(|h_r - h_{r(B')}|)_{(B')}, \\ |\xi(y)a_r^{ij}(y) - \xi(z)a_r^{ij}(z)| &\leq \xi(y) |a_r^{ij}(y) - a_r^{ij}(z)| + N|\xi(y) - \xi(z)| \end{aligned}$$

show that

$$I_r \leq 2 \left(\left| a_r^{ij} - a_{r(B+x_r)}^{ij} \right| \right)_{(B+x_r)} + N\rho,$$

which is equivalent to (5.11). This proves the lemma in the particular case under consideration.

General case. We rewrite the term $\Lambda_t^k u_t + g_t^k$ in (1.1) as $\sigma_t^{ik}(0)D_i u_t + \bar{g}_t^k$ with $\bar{g}_t^k = g_t^k + (\sigma_t^{ik} - \sigma_t^{ik}(0))D_i u_t$ and use the above result to conclude that estimate (5.1) holds with $N = N_1 = N_1(d, p, \delta, \varepsilon)$ if we add to its right-hand side

$$N_2(d, p, \delta, \varepsilon)\beta_1 \|Du\|_{\mathbf{L}_p(\tau)}.$$

By choosing $\beta_1 = \beta_1(d, p, \delta, \varepsilon)$ so that $N_2\beta_1 \leq 1/2$, we get (5.1) with $2N_1$ in place of N_1 . The lemma is proved. \square

Remark 5.3. If Assumption 2.1 is satisfied with $K = 0$ and a_t^{ij} and σ_t^{ik} depend only on ω and t , then the assertion of Lemma 5.2 is true with $\lambda_0 = 0$ and $N = N(d, p, \delta)$ and without requiring u to have compact support. This fact can be obtained

by following the arguments in section 4.3 of [6]. Even though those arguments are much longer, they allow one to prove a very general result that, roughly speaking, “whatever estimate can be established for solutions of the heat equation in Banach function spaces with norms that are invariant under time dependent shifting of the x coordinate, the same estimate with the same constant also holds for solutions of the parabolic equations with no lower order terms and with the matrix of the second order coefficients depending only on t and dominating (in the matrix sense) the unit matrix” (see [5]).

The next step is to consider equations with lower order terms. The following lemma and its corollary are stated in a slightly more general form than it is needed in the present article. The point is that we intend to use them in a subsequent article about equations in half spaces.

LEMMA 5.4. *Let $G \subset \mathbb{R}^d$ be a domain (perhaps, $G = \mathbb{R}^d$) and take $\varepsilon \geq \varepsilon_1 > 0$ and $\varepsilon_2 \in (0, \varepsilon/4]$.*

(i) *Let $f^j, g \in \mathbf{L}_p(\tau)$ and let $u \in \mathcal{W}_{p,0}^1(\tau)$ satisfy (1.1) for $t \leq \tau$ and be such that $u_t(x) = 0$ if $x \notin G$.*

(ii) *Suppose that Assumption 2.1 is satisfied.*

(iii) *Suppose that assumption (iii) of Lemma 5.2 is satisfied for any $t_0 \geq 0$ and x_0 such that $\text{dist}(x_0, G) \leq \varepsilon_2$.*

Then there exist constants $N, \lambda_0 \geq 0$, depending only on $d, p, K, \delta, \varepsilon, \varepsilon_1$, and ε_2 , such that estimate (5.1) holds true whenever $\lambda \geq \lambda_0$.

Proof. As usual we will use partitions of unity. Take a nonnegative $\xi \in C_0^\infty(B_{\varepsilon_2})$ with unit L_p -norm and take a nonnegative $\eta \in C_0^\infty((0, \varepsilon_1^2))$ with unit L_p -norm. For $s \in \mathbb{R}$ and $y \in \mathbb{R}^d$ introduce

$$\zeta(t, x) = \xi(x)\eta(t), \quad \zeta^{s,y}(t, x) = \zeta(t - s, x - y), \quad u_t^{s,y}(x) = \zeta^{s,y}(t, x)u_t(x)$$

so that, in particular,

$$(5.12) \quad |u_t(x)|^p = \int_{\mathbb{R}^{d+1}} |u_t^{s,y}(x)|^p dy ds.$$

Observe that, for each s, y ,

$$(5.13) \quad \begin{aligned} du_t^{s,y} = & \left(\sigma_t^{ik} D_i u_t^{s,y} + \hat{g}_t^{s,y,k} \right) dw_t^k \\ & + \left(D_j \left(a_t^{ij} D_i u_t^{s,y} \right) - \lambda u_t^{s,y} + D_j \hat{f}_t^{s,y,j} + \hat{f}_t^{s,y,0} \right) dt \end{aligned}$$

for $t \leq \tau$, where we dropped the argument x (and ω) and

$$\hat{g}_t^{s,y,k} = \zeta^{s,y} (\nu_t^k u_t + g_t^k) - u_t \sigma_t^{ik} D_i \zeta^{s,y},$$

$$\hat{f}_t^{s,y,j} = \zeta^{s,y} (a_t^j u_t + f_t^j) - a_t^{ij} u_t D_i \zeta^{s,y}, \quad j = 1, \dots, d,$$

$$\hat{f}_t^{s,y,0} = \zeta^{s,y} (f_t^0 + b_t^i D_i u_t + c_t u_t) - f_t^j D_j \zeta^{s,y} - \left(a_t^{ij} D_i u_t + a_t^j u_t \right) D_j \zeta^{s,y} + \zeta_t^{s,y} u_t,$$

and $\zeta_t^{s,y}(t, x) = \xi(x - y)\eta'(t - s)$.

As is easy to see $u_t^{s,y}(t, x) = 0$ for $(t, x) \notin (s_+, s_+ + \varepsilon_1^2) \times B_{\varepsilon_2}(y)$. Therefore, by Lemma 5.2 if $\text{dist}(y, G) \leq \varepsilon_2$, then

$$(5.14) \quad \begin{aligned} & \lambda^{p/2} \|u^{s,y}\|_{\mathbf{L}_p(\tau)}^p + \|Du^{s,y}\|_{\mathbf{L}_p(\tau)}^p \\ & \leq N \left(\sum_{j=1}^d \|\hat{f}^{s,y,j}\|_{\mathbf{L}_p(\tau)}^p + \|\hat{g}^{s,y}\|_{\mathbf{L}_p(\tau)}^p \right) + N\lambda^{-p/2} \|\hat{f}^{s,y,0}\|_{\mathbf{L}_p(\tau)}^p \end{aligned}$$

provided that $\lambda \geq \lambda_0$, where N and λ_0 depend only on d, δ, p , and ε . This estimate also, obviously, holds if $\text{dist}(y, G) > \varepsilon_2$ since then $u_t^{s,y} \equiv 0$.

Next,

$$\begin{aligned} |\hat{f}_t^{s,y,j}| &\leq N\bar{\zeta}^{s,y}|u_t| + \zeta^{s,y}|f_t^j|, \quad j = 1, \dots, d, \\ |\hat{f}_t^{s,y,0}| &\leq N\bar{\zeta}^{s,y}(|Du_t| + |u_t|) + N\bar{\zeta}^{s,y} \sum_{j=0}^d |f_t^j|, \\ |\hat{g}_t^{s,y}|_{\ell_2} &\leq N\bar{\zeta}^{s,y}|u_t| + \zeta^{s,y}|g_t|_{\ell_2}, \end{aligned}$$

where $\bar{\zeta} = \zeta + |D\zeta| + |\zeta_t|$, $\bar{\zeta}^{s,y}(t, x) = \bar{\zeta}(t - s, x - y)$, and here and below we allow the constants N to depend only on $d, p, \delta, K, \varepsilon, \varepsilon_1$, and ε_2 .

We also notice that $|\zeta^{s,y}Du_t| \leq |D(\zeta^{s,y}u_t)| + \bar{\zeta}^{s,y}|u_t|$. Then we find that

$$\begin{aligned} &\lambda^{p/2} \|\zeta^{s,y}u\|_{\mathbf{L}_p(\tau)}^p + \|\zeta^{s,y}Du\|_{\mathbf{L}_p(\tau)}^p \\ &\leq N \left(\sum_{i=1}^d \|\bar{\zeta}^{s,y}f^i\|_{\mathbf{L}_p(\tau)}^p + \|\zeta^{s,y}g\|_{\mathbf{L}_p(\tau)}^p + \|\bar{\zeta}^{s,y}u\|_{\mathbf{L}_p(\tau)}^p \right) \\ &\quad + N\lambda^{-p/2} \left(\|\bar{\zeta}^{s,y}f^0\|_{\mathbf{L}_p(\tau)}^p + \|\bar{\zeta}^{s,y}Du\|_{\mathbf{L}_p(\tau)}^p \right). \end{aligned}$$

We integrate through this estimate and use formulas like (5.12). Then we obtain

$$\begin{aligned} &\lambda^{p/2} \|u\|_{\mathbf{L}_p(\tau)}^p + \|Du\|_{\mathbf{L}_p(\tau)}^p \\ &\leq N_1 \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)}^p + \|g\|_{\mathbf{L}_p(\tau)}^p + \|u\|_{\mathbf{L}_p(\tau)}^p \right) + N_1\lambda^{-p/2} \left(\|f^0\|_{\mathbf{L}_p(\tau)}^p + \|Du\|_{\mathbf{L}_p(\tau)}^p \right). \end{aligned}$$

Finally, we increase $\lambda_0 \geq 0$, if necessary, in such a way that $N_1\lambda^{-p/2} \leq 1/2$ for $\lambda \geq \lambda_0$. Then we obviously arrive at (5.1) with $N = 2N_1$. The lemma is proved. \square

To the best of the author’s knowledge, the following multiplicative estimate is new even in the deterministic case.

COROLLARY 5.5. *Let $\lambda = 0$. Then under the assumptions of Lemma 5.4 we have*

$$\|Du\|_{\mathbf{L}_p(\tau)} \leq N \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} + \|f^0\|_{\mathbf{L}_p(\tau)}^{1/2} \|u\|_{\mathbf{L}_p(\tau)}^{1/2} + \|u\|_{\mathbf{L}_p(\tau)} \right),$$

where N depends only on $d, p, K, \delta, \varepsilon, \varepsilon_1$, and ε_2 .

Indeed, take $\lambda > 0$ and add and subtract the term $(\lambda_0 + \lambda)u_t dt$ on the right in (1.1), thus introducing λ into the equation and modifying f_t^0 by including into it one of $(\lambda_0 + \lambda)u_t$. Then after applying (5.1), we see that

$$\begin{aligned} \|Du\|_{\mathbf{L}_p(\tau)} &\leq N \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} \right. \\ &\quad \left. + (\lambda_0 + \lambda)^{-1/2} \|f^0\|_{\mathbf{L}_p(\tau)} + (\lambda_0 + \lambda)^{1/2} \|u\|_{\mathbf{L}_p(\tau)} \right). \end{aligned}$$

Now it remains only to take the inf with respect to $\lambda > 0$.

Proof of Lemma 2.1. By bearing in mind an obvious shifting of time, we see that in the proof of assertions (i)–(iii) we may assume that $s = 0$.

(i) First of all observe that uniqueness of solution of (2.7) is well known even in a much wider class than $\mathcal{W}_p^1(\infty)$.

Let $u_0 \in \text{tr}_0 \mathcal{W}_p^1$, then $u_0 \in W^{1-2/p}$ for almost each ω and there is a unique solution of the heat equation

$$dv_t = \Delta v_t dt$$

of class $L_p((0, 1), W_p^1)$ with initial condition u_0 . Furthermore,

$$\|v\|_{L_p((0,1),W_p^1)} \sim \|u_0\|_{W_p^{1-2/p}}.$$

Next, take a $\zeta \in C_0^\infty(\mathbb{R})$ such that $\zeta_0 = 1$ and $\zeta_t = 0$ for $t \geq 1/2$ and define $\psi_t(x) = e^{-t} v_t(x) \zeta_t$ for $t \in [0, 1]$ and as zero if $t \geq 1/2$. Notice that (a.s.)

$$\psi \in L_p(\mathbb{R}_+, W_p^1),$$

and

$$\frac{\partial}{\partial t} \psi_t = \Delta \psi_t - \psi_t + e^{-t} \zeta'_t v_t.$$

Then it is a classical result that there exists a unique $\phi \in L_p(\mathbb{R}_+, W_p^2)$ which solves the equation

$$d\phi_t = (\Delta \phi_t - \phi_t + e^{-t} \zeta'_t v_t) dt$$

with zero initial condition. In addition,

$$\|\phi\|_{L_p(\mathbb{R}_+, W_p^2)} \leq N \|\zeta' v\|_{L_p(\mathbb{R}_+, L_p)} \leq N \|u_0\|_{W_p^{1-2/p}},$$

where the constants N depend only on d and p . Owing to these estimates and uniqueness, the operators mapping u_0 into v and ϕ are continuous (and nonrandom). Since u_0 is \mathcal{F}_0 -measurable, the same is true for ψ , ϕ , and $u = \psi - \phi$, which is of class $L_p(\mathbb{R}_+, W_p^1)$, satisfies (2.7), and equals u_0 for $t = 0$. Also, for each ω

$$\|u\|_{L_p(\mathbb{R}_+, W_p^1)} \leq \|\psi\|_{L_p(\mathbb{R}_+, W_p^1)} + \|\phi\|_{L_p(\mathbb{R}_+, W_p^1)} \leq N \|u_0\|_{W_p^{1-2/p}},$$

where N depends only on d and p . By raising the extreme terms to the p th power and taking expectations we get the first inequality in (2.8) and also finish proving the “only if” part of (i).

To prove the “if” part assume that we have a $v \in \mathcal{W}_p^1(\infty)$ satisfying (2.7) and equal u_0 at $t = 0$. Then $u_t = v_t e^t$ satisfies $\partial u_t / \partial t = \Delta u_t$ and is of class $\mathcal{W}_p^1(1)$. It follows that almost all ω we have $u \in L_p((0, 1), W_p^1)$, $u_0 \in W_p^{1-2/p}$, and

$$\|u_0\|_{W_p^{1-2/p}} \leq N \|u\|_{L_p((0,1),W_p^1)} \leq N \|v\|_{L_p(\mathbb{R}_+, W_p^1)}.$$

By raising all expressions to the power p and taking expectations we arrive at the second estimate in (2.8). Assertion (i) is proved.

The “only if” part in (ii) is, actually, proved above. To prove the “if” part write

$$dv_t = (D_i f_t^i + f_t^0) dt + g_t^k dw_t^k = \left(\Delta v_t - \lambda v_t + D_i \hat{f}_t^i + \hat{f}_t^0 \right) dt + g_t^k dw_t^k,$$

where the constant $\lambda > 0$ will be chosen later, $\hat{f}_t^i = f_t^i - D_i v_t$, $i = 1, \dots, d$, $\hat{f}_t^0 = f_t^0 + \lambda v_t$, and $\hat{f}^j, g \in \mathbf{L}_p(1)$. Next, take the function ζ as above, set $u = v\zeta$, and observe that

$$(5.15) \quad du_t = (\Delta u_t - \lambda u_t + D_i \check{f}_t^i + \check{f}_t^0) dt + \check{g}_t^k dw_t^k,$$

where $\check{f}^0 = \zeta \hat{f}^0 + v\zeta'$, $\check{f}_t^i = \zeta \hat{f}_t^i$, $i = 1, \dots, d$, $\check{g}^k = \zeta g^k$, and $\check{f}^j, \check{g} \in \mathbf{L}_p(\infty)$ and $u \in \mathcal{W}_p^1(\infty)$.

By Lemma 5.1, for λ fixed and large enough (actually, one can take $\lambda = 1$, which is shown by using dilations), (5.15) with zero initial condition admits a unique solution $\psi \in \mathcal{W}_p^1(\infty)$ and

$$\begin{aligned} \|\psi\|_{\mathbf{W}_p^1(\infty)} &\leq N \left(\sum_{j=0}^d \|\check{f}^j\|_{\mathbf{L}_p(\infty)} + \|\check{g}\|_{\mathbf{L}_p(\infty)} \right) \\ &\leq N \left(\sum_{j=0}^d \|f^j\|_{\mathbf{L}_p(1)} + \|g\|_{\mathbf{L}_p(1)} + \|v\|_{\mathbf{W}_p^1(1)} \right). \end{aligned}$$

Then the difference $\phi = u - \psi$ satisfies (2.7), is of class $\mathcal{W}_p^1(\infty)$, and $\phi_0 = u_0$. By assertion (i) we have $u_0 \in \text{tr}_0 \mathcal{W}_p^1$, which proves the “if” part in (ii). Furthermore,

$$\begin{aligned} \|u_0\|_{\text{tr}_0 \mathcal{W}_p^1} &\leq N \|\phi\|_{\mathbf{W}_p^1(\infty)} \leq N \|u\|_{\mathbf{W}_p^1(\infty)} + N \|\psi\|_{\mathbf{W}_p^1(\infty)} \\ &\leq N \|v\|_{\mathbf{W}_p^1(1)} + N \|\psi\|_{\mathbf{W}_p^1(\infty)} \\ &\leq N \left(\sum_{j=0}^d \|f^j\|_{\mathbf{L}_p(1)} + \|g\|_{\mathbf{L}_p(1)} + \|v\|_{\mathbf{W}_p^1(1)} \right). \end{aligned}$$

This proves assertion (iii).

To prove (iv) observe that obvious dilations of the t axis allow us to assume that $s = 1$. Then write (2.2) for $t \in [0, 1]$ and notice that tu_t admits representation (2.2) with new f^j and g^k having simple relations with u_t and the original f^j and g^k . It follows that in the rest of the proof we may assume that $u_0 = 0$.

In that case, take a sufficiently large $\lambda > 0$ and consider the equation

$$dv_t = (\Delta v_t - \lambda v_t + D_i \bar{f}_t^i + \bar{f}_t^0) dt + \bar{g}_t^k dw_t^k$$

for $t \geq 0$ with zero initial condition, where

$$\begin{aligned} \bar{f}_t^i &= f_t^i I_{(0,1)}(t) - D_i u_t I_{(0,1)}(t), \quad i = 1, \dots, d, \\ \bar{f}_t^0 &= (f_t^0 + \lambda u_t) I_{(0,1)}(t), \quad \bar{g}_t^k = g_t^k I_{(0,1)}(t). \end{aligned}$$

By uniqueness, $v_t = u_t$ for $t \in [0, 1]$ and by assertion (iii) we have $v_1 \in \text{tr}_1 \mathcal{W}_p^1$. This fact combined with already known estimates of v proves assertion (iv). The lemma is proved. \square

6. Proof of Theorem 2.2. Owing to Lemma 2.1 we may assume that we are given a v as in assertion (i) of the lemma. By introducing a new unknown function $\bar{u} = u - v$ we see that u satisfies (1.1) and $u_0 = v_0$ if and only if $\bar{u}_0 = 0$ and

$$d\bar{u}_t = \left(L_t \bar{u}_t - \lambda \bar{u}_t + D_j \bar{f}_t^j + \bar{f}_t^0 \right) dt + \left(\Lambda_t^k \bar{u}_t + \bar{g}_t^k \right) dw_t^k,$$

where

$$\begin{aligned} \bar{f}_t^j &= f_t^j - D_j v_t + a_t^{ij} D_i v_t + a_t^j v_t, \quad j = 1, \dots, d, \\ \bar{f}_t^0 &= f_t^0 + b_t^i D_i v_t + (c_t - \lambda + 1)v_t, \\ \bar{g}_t^k &= g_t^k + \sigma_t^{ik} D_i v_t + \nu_t^k v_t. \end{aligned}$$

By Lemma 2.1 we have $\bar{f}^j, \bar{g} \in \mathbf{L}_p(\tau)$ and the problem of finding solutions of (1.1) with initial data u_0 is thus reduced to the same problem but with zero initial data.

Furthermore, if estimate (2.10) holds for solutions with zero initial condition, then (for $\lambda \geq \lambda_0$)

$$\begin{aligned} \lambda^{1/2} \|u\|_{\mathbf{L}_p(\tau)} + \|Du\|_{\mathbf{L}_p(\tau)} - \lambda^{1/2} \|v\|_{\mathbf{L}_p(\tau)} - \|Dv\|_{\mathbf{L}_p(\tau)} & \\ \leq \lambda^{1/2} \|\bar{u}\|_{\mathbf{L}_p(\tau)} + \|D\bar{u}\|_{\mathbf{L}_p(\tau)} & \\ \leq N \left(\sum_{i=1}^d \|\bar{f}^i\|_{\mathbf{L}_p(\tau)} + \|\bar{g}\|_{\mathbf{L}_p(\tau)} \right) + N\lambda^{-1/2} \|\bar{f}^0\|_{\mathbf{L}_p(\tau)} & \\ \leq N \left(\sum_{i=1}^d \|f^i\|_{\mathbf{L}_p(\tau)} + \|g\|_{\mathbf{L}_p(\tau)} + \|v\|_{\mathbf{W}_p^1(\tau)} \right) & \\ + N\lambda^{-1/2} \left(\|f^0\|_{\mathbf{L}_p(\tau)} + \|v\|_{\mathbf{W}_p^1(\tau)} \right) + N\lambda^{1/2} \|v\|_{\mathbf{L}_p(\tau)}, & \end{aligned}$$

which yields (2.10) in full generality.

It follows that while proving (2.10) we may also assume that $u_0 = 0$. Therefore, in the rest of the proof of assertions (i) and (ii) we assume that $u_0 = 0$. Having in mind the substitution $u_t = v_t e^{-\mu t}$, we see that while proving assertion (i) it suffices to concentrate on large λ and prove only the second part of the assertion.

We recall that we suppose that Assumption 2.2 is satisfied with $\beta = \beta_0/3$ and β_0 from Lemma 5.1 and Assumption 2.3 is satisfied with β_1 defined in Lemma 5.2. It follows that assumption (iii) of Lemma 5.2 is satisfied for any (t_0, x_0) .

Now we take λ_0 larger than the one in Lemma 3.3 and the one in Lemma 5.4. In that case uniqueness follows from Lemma 3.3. In the proof of existence we will rely on the method of continuity and the a priori estimate (5.1) which is established in Lemma 5.4. For $\lambda \geq \lambda_0$ and $\theta \in [0, 1]$ we consider the equation

$$(6.1) \quad du_t = [(\theta L_t + (1 - \theta)\Delta)u_t - \lambda u_t + D_i f_t^i + f_t^0] dt + (\theta \Lambda_t^k u_t + g_t^k) dw_t^k.$$

We call a $\theta \in [0, 1]$ “good” if the assertions of the theorem hold for (6.1). Observe that 0 is a “good” point by Lemma 5.1. Now to prove the theorem it suffices to show that there exists a $\gamma > 0$ such that if θ_0 is a good point, then all points of the interval $[\theta_0 - \gamma, \theta_0 + \gamma] \cap [0, 1]$ are “good”. So fix a “good” θ_0 and for any $v \in \mathbf{W}_p^1(\tau)$ consider the equation

$$(6.2) \quad \begin{aligned} du_t &= [(\theta_0 L_t + (1 - \theta_0)\Delta)u_t - \lambda u_t + (\theta - \theta_0)(L_t - \Delta)v_t + D_i f_t^i + f_t^0] dt \\ &\quad + (\theta_0 \Lambda_t^k u_t + (\theta - \theta_0)\Lambda^k v_t + g_t^k) dw_t^k. \end{aligned}$$

Observe that

$$(L_t - \Delta)v_t = D_j \left((a^{ij} - \delta^{ij}) D_i v_t + a_t^j v_t \right) + b_t^i D_i v_t + cv_t$$

and recall that $v \in \mathbf{W}_p^1(\tau)$. It follows by assumption that (6.2) has a unique solution $u \in \mathcal{W}_{p,0}^1(\tau) (\subset \mathbf{W}_p^1(\tau))$.

In this way, for f^j and g being fixed, we define a mapping $v \rightarrow u$ in the space $\mathbf{W}_p^1(\tau)$. It is important to keep in mind that the image u of $v \in \mathbf{W}_p^1(\tau)$ is always in $\mathcal{W}_{p,0}^1(\tau)$. Take $v', v'' \in \mathbf{W}_p^1(\tau)$ and let u', u'' be their corresponding images. Then $u := u' - u''$ satisfies

$$\begin{aligned} du_t = & [(\theta_0 L_t + (1 - \theta_0)\Delta)u_t - \lambda u_t + (\theta - \theta_0)(L_t - \Delta)v_t] dt \\ & + (\theta_0 \Lambda_t^k u_t + (\theta - \theta_0)\Lambda^k v_t) dw_t^k, \end{aligned}$$

where $v = v' - v''$. It follows by Lemma 5.4 that

$$\|u\|_{\mathbf{W}_p^1(\tau)} \leq N|\theta - \theta_0| \|v\|_{\mathbf{W}_p^1(\tau)}$$

with a constant N independent of f, g, v', v'', θ_0 , and θ . For θ sufficiently close to θ_0 , our mapping is a contraction and, since $\mathbf{W}_p^1(\tau)$ is a Banach space, it has a fixed point. This fixed point is in $\mathcal{W}_{p,0}^1(\tau)$ and, obviously, satisfies (6.1). This proves assertion (i) of the theorem.

Estimate (2.10) is proved above in Lemma 5.4 and assertion (iii) follows from Theorem 3.1. The theorem is proved. \square

REFERENCES

- [1] KYEONG-HUN KIM, *On stochastic partial differential equations with variable coefficients in C^1 domains*, Stochastic Process. Appl., 112 (2004), pp. 261–283.
- [2] KYEONG-HUN KIM, *On L_p -theory of stochastic partial differential equations of divergence form in C^1 domains*, Probab. Theory Related Fields, 130 (2004), pp. 473–492.
- [3] KYEONG-HUN KIM, *L_p estimates for SPDE with discontinuous coefficients in domains*, Electron. J. Probab., 10 (2005), pp. 1–20.
- [4] N.V. KRYLOV, *A generalization of the Littlewood-Paley inequality and some other results related to stochastic partial differential equations*, Ulam Quarterly, 2 (1994), pp. 16–26, <http://www.ulam.usm.edu/VIEW2.4/krylov.ps>
- [5] N.V. KRYLOV, *A parabolic Littlewood-Paley inequality with applications to parabolic equations*, Topological Methods in Nonlinear Analysis, Journal of the Juliusz Schauder Center, 4 (1994), pp. 355–364.
- [6] N.V. KRYLOV, *An analytic approach to SPDEs*, in Stochastic Partial Differential Equations: Six Perspectives, Mathematical Surveys and Monographs, Vol. 64, AMS, Providence, RI, 1999, pp. 185–242.
- [7] N.V. KRYLOV, *On the foundation of the L_p -theory of SPDEs*, in Stochastic Partial Differential Equations and Applications-VII, G. Da Prato, L. Tubaro eds., A Series of Lecture Notes in Pure and Applied Math., Chapman & Hall/CRC, London, 2006, pp. 179–191.
- [8] N.V. KRYLOV, *Parabolic equations with VMO coefficients in Sobolev spaces with mixed norms*, J. Function. Anal., 250 (2007), pp. 521–558.
- [9] N.V. KRYLOV, *Maximum principle for SPDEs and its applications*, in Stochastic Differential Equations: Theory and Applications, A Volume in Honor of Professor Boris L. Rozovskii, P.H. Baxendale, and S.V. Lototsky, eds., Interdisciplinary Mathematical Sciences, Vol. 2, World Scientific, 2007, pp. 311–338.
- [10] N.V. KRYLOV, *Filtering equations for partially observable diffusion processes with Lipschitz continuous coefficients*, The Oxford Handbook of Nonlinear Filtering, Oxford University Press, to appear.
- [11] N.V. KRYLOV, *Lectures on elliptic and parabolic equations in Sobolev spaces*, Amer. Math. Soc., Graduate Studies in Math. 96, Providence, RI, 2008.
- [12] N.V. KRYLOV, *Itô's formula for the L_p -norm of stochastic W_p^1 -valued processes*, <http://arxiv.org/abs/0806.1557>
- [13] B.L. ROZOVSKII, *Stochastic Evolution Systems*, Kluwer, Dordrecht, 1990.

SEMILINEAR STOCHASTIC EQUATIONS IN A HILBERT SPACE WITH A FRACTIONAL BROWNIAN MOTION*

T. E. DUNCAN[†], B. MASLOWSKI[‡], AND B. PASIK-DUNCAN[†]

Abstract. The solutions of a family of semilinear stochastic equations in a Hilbert space with a fractional Brownian motion are investigated. The nonlinear term in these equations has primarily only a growth condition assumption. An arbitrary member of the family of fractional Brownian motions can be used in these equations. Existence and uniqueness for both weak and mild solutions are obtained for some of these semilinear equations. The weak solutions are obtained by a measure transformation that verifies absolute continuity with respect to the measure for the solution of the associated linear equation. Some examples of stochastic differential and partial differential equations are given that satisfy the assumptions for the solutions of the semilinear equations.

Key words. semilinear stochastic equations, fractional Brownian motion, stochastic partial differential equations, absolute continuity of measures

AMS subject classifications. 60H15, 60G18, 60G15

DOI. 10.1137/08071764X

1. Introduction. Fractional Brownian motion denotes a family of Gaussian processes with continuous sample paths that are indexed by the Hurst parameter $H \in (0, 1)$ and that have properties that appear empirically in a wide variety of physical phenomena, such as hydrology, economic data, telecommunications, and medicine. Since some physical phenomena are naturally modeled by stochastic partial differential equations and the randomness can be described by a fractional Gaussian noise, it is important to study the problems of the solutions of stochastic differential equations in a Hilbert space with a fractional Brownian motion. A significant family of these stochastic equations is the set of semilinear equations, so it is important to investigate the existence and the uniqueness of the solutions of the equations and the sample path properties of the solutions. If primarily only some growth assumptions are made on the nonlinear terms in the semilinear equations, then it is natural to investigate weak solutions, especially those that arise by an absolutely continuous transformation of the measure of the solution of the associated linear stochastic equation.

The study of the solutions of stochastic equations in an infinite-dimensional space with a (cylindrical) fractional Brownian motion (for example, stochastic partial differential equations) has been relatively limited. For the Hurst parameter $H \in (1/2, 1)$, linear and semilinear equations with an additive fractional Gaussian noise, the formal derivative of a fractional Brownian motion, are considered in [8, 13, 15, 28]. Random dynamical systems described by such stochastic equations and their fixed points are studied in [22]. A pathwise (or nonprobabilistic) approach is used in [21] to study a parabolic equation with a fractional Gaussian noise where the stochastic term is a nonlinear function of the solution. Strong solutions of bilinear evolution equations

*Received by the editors March 4, 2008; accepted for publication (in revised form) November 11, 2008; published electronically February 20, 2009. Research supported in part by NSF grants DMS 0204669, DMS 0505706, ANI 0124510, and GACR 201/04/0750.

<http://www.siam.org/journals/sima/40-6/71764.html>

[†]Department of Mathematics, University of Kansas, Lawrence, KS 66045 (duncan@math.ku.edu, bozena@math.ku.edu).

[‡]Institute of Mathematics, Czech Academy of Sciences, Prague, Czech Republic (maslow@math.cas.cz).

with a fractional Brownian motion are considered in [11, 12], and the same type of equation is studied in [33], where a fractional Feynman–Kac formula is obtained. A stochastic wave equation with a fractional Gaussian noise is considered in [2], and a stochastic heat equation with a multiparameter fractional Gaussian noise is studied in [16, 18].

One facet of the motivation for the study of weak solutions in an infinite-dimensional space follows from some results [7, 27] for weak solutions in finite-dimensional spaces that use an absolutely continuous transformation of measures which generalize the result of Girsanov [14] for Brownian motion.

In this paper, a similar analysis is made for infinite-dimensional state spaces. While the structure of the infinite-dimensional Girsanov theorem is analogous to the finite-dimensional case, significant distinct difficulties arise when the application of this theorem is used for stochastic equations in infinite-dimensional spaces. First the driving process is only cylindrical, so the Girsanov theorem can only be used to transform the semilinear equation to a linear equation that is a fractional Ornstein–Uhlenbeck process. Since there is no classical strong solution to the linear equation, the mild solution must be used, making the analysis of the transformation of the measures by the Radon–Nikodým derivative more difficult because a suitable sample path regularity of the Ornstein–Uhlenbeck process must be verified. Unlike the finite-dimensional case, this regularity is not immediate and some assumptions on the coefficients in the linear equation must be made which are known for the case of Brownian motion. The sample regularity requirement increases as the Hurst parameter H increases. Dually the operators that appear in the Radon–Nikodým derivative are less regular as H increases. Thus the applicability of the Girsanov theorem is not immediate in this case and some conditions must be determined for the whole procedure to succeed. Furthermore, for $H > \frac{1}{2}$ in the finite-dimensional case it is assumed that the nonlinear term in the semilinear equation satisfies a global Hölder condition, but this assumption is not satisfied in many typical examples in infinite-dimensional spaces, such as reaction-diffusion equations. Thus this Hölder condition is relaxed here as well.

In section 2, some results from fractional calculus are given, and these results are used to describe a kernel function for an integral operator that provides an isometry of the second moment of Wiener-type stochastic integrals with respect to a fractional Brownian motion and the Lebesgue space of square integrable functions. Furthermore, some recent results for the solution of a linear stochastic equation in a Hilbert space [28] are described. In section 3, semilinear stochastic equations in a Hilbert space are studied. Initially, an absolute continuity of measures result for transforming the solution of a linear stochastic equation is verified that can be viewed as an analogue of the result of Girsanov [14] for a transformation of a finite-dimensional standard Brownian motion. For a semilinear stochastic equation where the nonlinear term satisfies a linear growth condition and some additional conditions are satisfied, it is shown that there is one and only one weak solution. The weak solution is obtained by verifying an absolute continuity of the measure of the solution with respect to the measure of the solution of the associated linear equation. The cases $H \in (0, 1/2)$ and $H \in (1/2, 1)$ are treated separately. Absolute continuity of the above measures is verified when the nonlinearity satisfies a power growth condition and some additional assumptions are made. In section 4, some examples of stochastic differential and partial differential equations are given that satisfy the assumptions of the theorems.

2. Preliminaries. In this section, a cylindrical fractional Brownian motion in a separable Hilbert space is introduced, a Wiener-type stochastic integral with respect to this process is defined, and some basic properties of this integral are noted. Initially, some facts from the theory of fractional integration (cf. [31]) are described. Let $(V, \|\cdot\|, \langle \cdot, \cdot \rangle)$ be a separable Hilbert space. If $\varphi \in L^1([0, T], V)$, then for $\alpha > 0$ the left-side and the right-side fractional (Riemann–Liouville) integrals of φ are defined (for almost all $t \in [0, T]$) by

$$(I_{0+}^\alpha \varphi)(t) = \frac{1}{\Gamma(\alpha)} \int_0^t (t-s)^{\alpha-1} \varphi(s) ds$$

and

$$(I_{T-}^\alpha \varphi)(t) = \frac{1}{\Gamma(\alpha)} \int_t^T (s-t)^{\alpha-1} \varphi(s) ds,$$

respectively, where $\Gamma(\cdot)$ is the gamma function. For $\alpha \in (0, 1)$ the inverse operators of these fractional integrals are called fractional derivatives and can be given by their respective Weyl representations

$$(D_{0+}^\alpha \psi)(t) = \frac{1}{\Gamma(1-\alpha)} \left(\frac{\psi(t)}{t^\alpha} + \alpha \int_0^t \frac{\psi(t) - \psi(s)}{(t-s)^{\alpha+1}} ds \right)$$

and

$$(D_{T-}^\alpha \psi)(t) = \frac{1}{\Gamma(1-\alpha)} \left(\frac{\psi(t)}{(T-t)^\alpha} + \alpha \int_t^T \frac{\psi(s) - \psi(t)}{(s-t)^{\alpha+1}} ds \right),$$

where $\psi \in I_{0+}^\alpha (L^1([0, T], V))$ and $\psi \in I_{T-}^\alpha (L^1([0, T], V))$, respectively.

Let $K_H(t, s)$ for $0 \leq s \leq t \leq T$ be the real-valued kernel function

$$(2.1) \quad K_H(t, s) = \frac{\tilde{c}_H (t-s)^{H-\frac{1}{2}}}{\Gamma(H+\frac{1}{2})} + \frac{\tilde{c}_H (\frac{1}{2}-H)}{\Gamma(H+\frac{1}{2})} \int_s^t (u-s)^{H-\frac{3}{2}} \left(1 - \left(\frac{s}{u}\right)^{\frac{1}{2}-H}\right) du$$

for $H \in (0, 1/2)$. If $H \in (1/2, 1)$, then K_H has a simpler form as

$$(2.2) \quad K_H(t, s) = \frac{\hat{c}_H}{\Gamma(H-\frac{1}{2})} s^{\frac{1}{2}-H} \int_s^t (u-s)^{H-\frac{3}{2}} u^{H-\frac{1}{2}} du.$$

The terms \tilde{c}_H and \hat{c}_H are constants that depend only on H .

Define the integral operator \mathbb{K}_H induced from the kernel K_H by

$$(2.3) \quad \mathbb{K}_H \varphi(t) = \int_0^t K_H(t, s) h(s) ds$$

for $h \in L^2([0, T], V)$. It is well known [31] that

$$\mathbb{K}_H: L^2([0, T], V) \rightarrow I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V))$$

is a bijection and \mathbb{K}_H can be described as

$$(2.4) \quad \mathbb{K}_H h(s) = \bar{c}_H I_{0+}^{2H} \left(u_{\frac{1}{2}-H} I_{0+}^{\frac{1}{2}-H} \left(u_{H-\frac{1}{2}} h \right) \right) (s)$$

for $H \in (0, 1/2]$ and

$$(2.5) \quad \mathbb{K}_H h(s) = c_H I_{0+}^1 \left(u_{H-\frac{1}{2}} I_{0+}^{H-\frac{1}{2}} \left(u_{\frac{1}{2}-H} h \right) \right) (s)$$

for $H \in [1/2, 1)$, where

$$(2.6) \quad c_H = \left[\frac{2H\Gamma(H + \frac{1}{2})\Gamma(\frac{3}{2} - H)}{\Gamma(2 - 2H)} \right]^{\frac{1}{2}},$$

$$\bar{c}_H = c_H \Gamma(2H),$$

and

$$u_a(s) = s^a I$$

for $s \geq 0$ and $a \in \mathbb{R}$. The inverse operator

$$\mathbb{K}_H^{-1} : I_{0+}^{H+\frac{1}{2}} (L^2([0, T], V)) \rightarrow L^2([0, T], V)$$

is given by

$$(2.7) \quad \mathbb{K}_H^{-1} \varphi(s) = \bar{c}_H^{-1} s^{\frac{1}{2}-H} D_{0+}^{\frac{1}{2}-H} \left(u_{H-\frac{1}{2}} D_{0+}^{2H} \varphi \right) (s)$$

for $H \in (0, 1/2]$ and

$$(2.8) \quad \mathbb{K}_H^{-1} \varphi(s) = c_H^{-1} s^{H-\frac{1}{2}} D_{0+}^{H-\frac{1}{2}} \left(u_{\frac{1}{2}-H} D \varphi \right) (s)$$

for $H \in [1/2, 1)$ and $\varphi \in I_{0+}^{H+\frac{1}{2}} (L^2([0, T], V))$. Note that if $\varphi \in H^1([0, T], V)$, the Sobolev space, then

$$(2.9) \quad \mathbb{K}_H^{-1} \varphi(s) = \bar{c}_H^{-1} s^{H-\frac{1}{2}} I_{0+}^{\frac{1}{2}-H} \left(u_{\frac{1}{2}-H} \varphi' \right) (s)$$

for $H \in (0, 1/2]$.

Since the operator \mathbb{K}_H^{-1} plays an important role in what follows, it is desirable to have some information about its domain $I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V))$. It is straightforward that $I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)) \subset C^\beta([0, T], V)$ for $\beta > \frac{1}{2} - H$ and $H \in (0, 1/2)$. However, in section 3, a more refined result is needed. If $H \in (1/2, 1)$, then $I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)) \subset L^2([0, T], V)$.

A definition of the stochastic integral of a deterministic V -valued function with respect to a scalar fractional Brownian motion $(\beta(t), t \geq 0)$ is given. This definition uses the methods in [1, 6, 11, 30]. An alternative, equivalent method is given in [10].

A family of linear operators $(\mathcal{K}_H^*, H \in (0, 1))$ is defined which provides an isometry between Wiener-type integrals of a fractional Brownian motion and $L^2([0, T], V)$. It is written as an adjoint because the linear operator \mathcal{K}_H occurs naturally in the factorization of the covariance for a fractional Brownian motion in L^2 .

Let $\mathcal{K}_H^* : \mathcal{E} \rightarrow L^2([0, T], V)$ be the linear map given by

$$(2.10) \quad \mathcal{K}_H^* \varphi(t) = \varphi(t) K_H(T, t) + \int_t^T (\varphi(s) - \varphi(t)) \frac{\partial K_H}{\partial s}(s, t) ds$$

for $\varphi \in \mathcal{E}$ and K_H given by (2.1), where \mathcal{E} is the linear space of V -valued step functions on $[0, T]$, that is, $\varphi \in \mathcal{E}$ if

$$\varphi(t) = \sum_{i=1}^{n-1} x_i \mathbb{1}_{[t_i, t_{i+1})}(t),$$

where $x_i \in V$ for $i \in \{1, \dots, n-1\}$ and $0 = t_1 < t_2 < \dots < t_n = T$.

Define the stochastic integral as

$$\int_0^T \varphi d\beta := \sum_{i=1}^n x_i (\beta(t_{i+1}) - \beta(t_i)).$$

It follows directly that

$$(2.11) \quad \mathbb{E} \left\| \int_0^T \varphi d\beta \right\|^2 = |\mathcal{K}_H^* \varphi|_{L^2([0, T], V)}^2,$$

where $|\cdot|_{L^2([0, T], V)}$ is the norm in $L^2([0, T], V)$ induced by the inner product. Let $(\mathcal{H}, |\cdot|_{\mathcal{H}}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ be the Hilbert space obtained by the completion of the pre-Hilbert space \mathcal{E} with the inner product

$$(2.12) \quad \langle \varphi, \psi \rangle_{\mathcal{H}} := \langle \mathcal{K}_H^* \varphi, \mathcal{K}_H^* \psi \rangle_{L^2([0, T], V)}$$

for $\varphi, \psi \in \mathcal{E}$. The stochastic integral is extended to an arbitrary $\varphi \in \mathcal{H}$ by the isometry (2.11). Thus \mathcal{H} is a linear space of integrable functions, and it is useful to obtain some more specific information about \mathcal{H} . If $H \in (1/2, 1)$, then it is easily verified that $\mathcal{H} \supset \tilde{\mathcal{H}}$, where $\tilde{\mathcal{H}}$ is the Banach space of Borel measurable functions with the norm $|\cdot|_{\tilde{\mathcal{H}}}$ given by

$$(2.13) \quad |\varphi|_{\tilde{\mathcal{H}}}^2 := \int_0^T \int_0^T \|\varphi(u)\| \|\varphi(v)\| \phi_H(u-v) du dv,$$

where $\phi_H(u) = H(2H-1)|u|^{2H-2}$, and it can be verified that $\tilde{\mathcal{H}} \supset L^{\frac{1}{H}}([0, T], V)$ and, in particular, $\tilde{\mathcal{H}} \supset L^2([0, T], V)$ (cf. [12]). If $\varphi \in \tilde{\mathcal{H}}$ and $H > 1/2$, then

$$(2.14) \quad \mathbb{E} \left\| \int_0^T \varphi d\beta \right\|^2 = \int_0^T \int_0^T \langle \varphi(u), \varphi(v) \rangle \phi_H(u-v) du dv.$$

If $H \in (0, 1/2)$, then the space of integrable functions is smaller than for $H \in (1/2, 1)$. For $H \in (0, 1/2)$ it is known that $\mathcal{H} \supset H^1([0, T], V)$ (cf. [17, Theorem 5.20]) and $\mathcal{H} \supset C^\beta([0, T], V)$ for each $\beta > 1/2 - H$ (a more specific result is given in the next section). If $H \in (0, 1/2)$, then the linear operator \mathcal{K}_H^* can be described by a fractional derivative

$$(2.15) \quad \mathcal{K}_H^* \varphi(t) = c_H t^{\frac{1}{2}-H} D_{T-}^{\frac{1}{2}-H} \left(u_{H-\frac{1}{2}} \varphi \right) (t),$$

where its domain is $\mathcal{H} = I_{T-}^{\frac{1}{2}-H} (L^2([0, T], V))$ (cf. [1, Proposition 6]). If $H \in (1/2, 1)$, then

$$(2.16) \quad \mathcal{K}_H^* \varphi(t) = c_H t^{\frac{1}{2}-H} I_{T-}^{H-\frac{1}{2}} \left(u_{H-\frac{1}{2}} \varphi \right) (t).$$

A standard cylindrical fractional Brownian motion is defined now.

DEFINITION 2.1. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a complete probability space. A cylindrical process $\langle B, \cdot \rangle: \Omega \times \mathbb{R}_+ \times V \rightarrow \mathbb{R}$ on $(\Omega, \mathcal{F}, \mathbb{P})$ is called a standard cylindrical fractional Brownian motion with the Hurst parameter $H \in (0, 1)$ if

- (1) for each $x \in V \setminus \{0\}$, $\frac{1}{\|x\|} \langle B(\cdot), x \rangle$ is a standard scalar fractional Brownian motion with the Hurst parameter H ;
- (2) for $\alpha, \beta \in \mathbb{R}$ and $x, y \in V$

$$\langle B(t), \alpha x + \beta y \rangle = \alpha \langle B(t), x \rangle + \beta \langle B(t), y \rangle \quad \text{a.s. } \mathbb{P}.$$

Note that $\langle B(t), x \rangle$ has the interpretation of the evaluation of the functional $B(t)$ at x though the process $B(\cdot)$ does not take values in V .

For $H = 1/2$, this definition is the usual one for a standard cylindrical Wiener process in V . For a complete orthonormal basis $(e_n, n \in \mathbb{N})$ of V , letting $\beta_n(t) = \langle B(t), e_n \rangle$ for $n \in \mathbb{N}$, the sequence of scalar processes $(\beta_n, n \in \mathbb{N})$ is independent and B can be represented by the formal series

$$(2.17) \quad B(t) = \sum_{n=1}^{\infty} \beta_n(t) e_n$$

that does not converge a.s. in V .

Naturally associated with a standard cylindrical fractional Brownian motion is a standard cylindrical Wiener process $(W(t), t \geq 0)$ in V such that, formally, $\dot{B}(t) = \mathcal{K}_H \dot{W}(t)$. For $x \in V \setminus \{0\}$, let $\beta_x(t) = \langle B(t), x \rangle$. It is elementary to verify from (2.1) that there is a scalar Wiener process $(w_x(t), t \geq 0)$ such that

$$(2.18) \quad \beta_x(t) = \int_0^t K_H(t, s) dw_x(s)$$

for $t \in \mathbb{R}_+$. Dually, $w_x(t) = \beta_x((\mathcal{K}_H^*)^{-1} \mathbb{1}_{[0,t]})$, where \mathcal{K}_H^* is given by (2.15) or (2.16) and $V = \mathbb{R}$. Thus there is a formal expansion of W ,

$$(2.19) \quad W(t) = \sum_{n=1}^{\infty} w_n(t) e_n,$$

where $(e_n, n \in \mathbb{N})$ is a complete orthonormal basis for V and $w_n = \langle W, e_n \rangle$ for $n \in \mathbb{N}$.

Now, the stochastic integral $\int_0^T G dB$ is defined for a suitable operator-valued function $G: [0, T] \rightarrow \mathcal{L}(V)$ so that the integral is a V -valued random variable.

DEFINITION 2.2. Let $G: [0, T] \rightarrow \mathcal{L}(V)$ be Borel measurable, let $(e_n, n \in \mathbb{N})$ be a complete orthonormal basis in V , let $G(\cdot)e_n \in \mathcal{H}$ for each $n \in \mathbb{N}$, and let B be a standard cylindrical fractional Brownian motion for some fixed $H \in (0, 1)$. The stochastic integral $\int_0^T G dB$ is defined as

$$(2.20) \quad \int_0^T G dB := \sum_{n=1}^{\infty} \int_0^T G e_n d\beta_n,$$

provided the infinite series converges in $L^2(\Omega, V)$.

It is elementary to verify that this definition does not depend on the complete orthonormal basis that is used.

The following proposition describes some $\mathcal{L}(V)$ -valued functions G that can be used as integrands in Definition 2.2.

PROPOSITION 2.3. *Let $G: [0, T] \rightarrow \mathcal{L}(V)$ be Borel measurable and let $G(\cdot)x \in \mathcal{H}$ for each $x \in V$. Let $\Gamma_T: V \rightarrow L^2([0, T], V)$ be given as*

$$(2.21) \quad (\Gamma_T x)(t) = (\mathcal{K}_H^* G x)(t)$$

for $t \in [0, T]$ and $x \in V$. If $\Gamma_T \in \mathcal{L}_2(V, L^2([0, T], V))$, the linear space of Hilbert–Schmidt operators, then the stochastic integral (2.20) is a centered Gaussian V -valued random variable with the covariance operator \tilde{Q}_T given by

$$(2.22) \quad \tilde{Q}_T x = \int_0^T \sum_{n=1}^{\infty} \langle (\Gamma_T e_n)(s), x \rangle (\Gamma_T e_n)(s) ds.$$

This integral does not depend on the choice of the complete orthonormal basis $(e_n, n \in \mathbb{N})$ in V .

Proof. Substituting G in the definition of the stochastic integral (2.20), it is clear that the terms of the summation on the right-hand side are V -valued Gaussian random variables by the construction of the integral for a scalar fractional Brownian motion, and the sequence of random variables $(\int_0^T G e_n d\beta_n, n \in \mathbb{N})$ is independent. Computing the second moment of the tail of the series in (2.20) yields

$$\begin{aligned} \mathbb{E} \left\| \sum_{k=m}^{\infty} G e_k d\beta_k \right\|^2 &= \sum_{k=m}^{\infty} \mathbb{E} \left\| \int_0^T G e_k d\beta_k \right\|^2 = \sum_{k=m}^{\infty} \int_0^T \|(\mathcal{K}_H^* G e_k)(s)\|^2 ds \\ &= \sum_{k=m}^{\infty} \int_0^T \|(\Gamma_T e_k)(s)\|^2 ds = \sum_{k=m}^{\infty} \|\Gamma_T e_k\|_{L^2([0, T], V)}^2. \end{aligned}$$

It is clear that this final series tends to zero as m tends to infinity. Thus there is convergence in $L^2(\Omega, V)$ of the partial sums of the infinite series in (2.20).

To verify that (2.20) is a Gaussian random variable and the form of the covariance \tilde{Q}_T , initially note that for any $\varphi \in \mathcal{H}$ and $x \in V$, there is the equality

$$(2.23) \quad \int_0^T \varphi d\beta_x = \int_0^T \mathcal{K}_H^* \varphi dw_x,$$

where w_x is the Wiener process given by (2.18). The terms in the infinite series on the right-hand side of (2.20) are V -valued, independent centered Gaussian random variables with the sequence of covariance operators $(\tilde{Q}_T^{(n)}, n \in \mathbb{N})$

$$(2.24) \quad \tilde{Q}_T^{(n)} x = \int_0^T \langle (\mathcal{K}_H^* G e_n)(s), x \rangle (\mathcal{K}_H^* G e_n)(s) ds$$

for each $n \in \mathbb{N}$ and $x \in V$. Thus

$$\begin{aligned} (2.25) \quad \tilde{Q}_T x &= \sum_{n=1}^{\infty} \int_0^T \langle (\mathcal{K}_H^* G e_n)(s), x \rangle (\mathcal{K}_H^* G e_n)(s) ds \\ &= \int_0^T \sum_{n=1}^{\infty} \langle (\Gamma_T e_n)(s), x \rangle (\Gamma_T e_n)(s) ds. \end{aligned}$$

The summability of the infinite series on the right-hand side follows from the Hilbert–Schmidt property of Γ_T . The independence of the stochastic integral from the choice of the complete orthonormal basis follows from (2.23) and the analogous property for stochastic integrals with respect to a standard cylindrical Wiener process. \square

Since for almost all $t \in [0, T]$ the linear operator $\Gamma_T(\cdot)(t) : V \rightarrow V$ is Hilbert–Schmidt, so we denote for almost all $t \in [0, T]$ the adjoint of $\Gamma_T(\cdot)(t)$ as $\Gamma_T^*(\cdot)(t) : V \rightarrow V$. It follows by (2.25) that for $x, y \in V$,

$$\begin{aligned} \langle \tilde{Q}_T x, y \rangle &= \int_0^T \sum_{n=1}^\infty \langle (\Gamma_T e_n)(s), x \rangle \langle (\Gamma_T e_n)(s), y \rangle ds \\ &= \int_0^T \sum_{n=1}^\infty \langle e_n, (\Gamma_T^* x)(s) \rangle \langle e_n, (\Gamma_T^* y)(s) \rangle ds \\ &= \int_0^T \langle (\Gamma_T^* x)(s), (\Gamma_T^* y)(s) \rangle ds \\ &= \int_0^T \langle \Gamma_T \Gamma_T^* x(s), y \rangle ds. \end{aligned}$$

If $H \in (1/2, 1)$, then \tilde{Q}_T satisfies

$$\tilde{Q}_T = \int_0^T \int_0^T G(u) G^*(v) \phi_H(u - v) du dv,$$

where $\phi_H(u) = H(2H - 1)|u|^{2H-2}$ and G is assumed to satisfy

$$\int_0^T \int_0^T |G(u)|_{\mathcal{L}_2(V)} |G(v)|_{\mathcal{L}_2(V)} \phi_H(u - v) du dv < \infty$$

(cf. [13, Proposition 2.2]).

The next proposition shows that some densely defined linear operators commute with the stochastic integration.

PROPOSITION 2.4. *If $\tilde{A} : \text{Dom}(\tilde{A}) \rightarrow V$ is a closed linear operator, $\text{Dom}(\tilde{A}) \subset V$, and $G : [0, T] \rightarrow \mathcal{L}(V)$ is Borel measurable such that $G([0, T]) \subset \text{Dom}(\tilde{A})$ and both G and $\tilde{A}G$ satisfy the conditions for G in Proposition 2.3, then*

$$\int_0^T G dB \subset \text{Dom}(\tilde{A}) \quad \text{a.s. } \mathbb{P}$$

and

$$(2.26) \quad \tilde{A} \int_0^T G dB = \int_0^T \tilde{A}G dB \quad \text{a.s. } \mathbb{P}.$$

Proof. By the assumptions on G and $\tilde{A}G$, it follows that $Ge_n \in \mathcal{H}$ and $\tilde{A}Ge_n \in \mathcal{H}$ for $n \in \mathbb{N}$, so by a standard argument using a sequence of step function integrands, the following equality is satisfied:

$$\tilde{A} \int_0^T Ge_n d\beta_n = \int_0^T \tilde{A}Ge_n d\beta_n.$$

Since the sequence of integrals that are obtained from a complete orthonormal basis $(e_n, n \in \mathbb{N})$ are Gaussian random variables it follows that

$$(2.27) \quad \lim_{m \rightarrow \infty} \sum_{n=1}^m \int_0^T G e_n d\beta_n = \int_0^T G dB \quad \text{a.s. } \mathbb{P}$$

and

$$\lim_{m \rightarrow \infty} \tilde{A} \left(\sum_{n=1}^m \int_0^T G e_n d\beta_n \right) = \lim_{m \rightarrow \infty} \sum_{n=1}^m \int_0^T \tilde{A} G e_n d\beta_n = \int_0^T \tilde{A} G dB \quad \text{a.s. } \mathbb{P}.$$

Since \tilde{A} is a closed linear operator it follows that $\int_0^T G dB \in \text{Dom}(\tilde{A})$ a.s. \mathbb{P} and equality (2.26) is satisfied. \square

Some results are reviewed for a linear stochastic differential equation with a cylindrical fractional Brownian motion whose solution is often called a fractional Ornstein–Uhlenbeck process. This process is a mild solution of the linear stochastic equation

$$(2.28) \quad \begin{aligned} dZ(t) &= AZ(t) dt + \Phi dB(t), \\ Z(0) &= x, \end{aligned}$$

where $Z(t), x \in V, (B(t), t \geq 0)$ is a standard cylindrical fractional Brownian with $H \in (0, 1), \Phi \in \mathcal{L}(V), A: \text{Dom}(A) \rightarrow V, \text{Dom}(A) \subset V,$ and A is the infinitesimal generator of a strongly continuous semigroup $(S(t), t \geq 0)$ on V . A mild solution of (2.28) is

$$(2.29) \quad \begin{aligned} Z(t) &= S(t)x + \int_0^t S(t-r)\Phi dB(r) \\ &= S(t)x + \hat{Z}(t), \end{aligned}$$

where the stochastic integral in (2.29) is given by Definition 2.2.

Typically it is assumed that $(S(t), t \geq 0)$ is an analytic semigroup. In this case, there is a $\hat{\beta} \in \mathbb{R}$ such that the operator $\hat{\beta}I - A$ is uniformly positive on V ; that is, the resolvent set contains $\{\lambda \in \mathbb{C}; |\arg \lambda| < \pi/2 + \delta\} \cup U,$ where $\delta > 0$ and U is a neighborhood of zero.

For each $\delta \geq 0, (V_\delta, \|\cdot\|_\delta)$ is a Hilbert space where $V_\delta = \text{Dom}((\hat{\beta}I - A)^\delta)$ with the graph norm topology so that

$$\|x\|_\delta = \left\| (\hat{\beta}I - A)^\delta x \right\|.$$

For the mild solution of (2.28), the cases $H \in (0, 1/2)$ and $H \in (1/2, 1)$ have been treated separately [13, 28] because the conditions for similar results are somewhat different. The case $H = 1/2$ (Brownian motion) has been studied extensively (cf. [4]).

For $H \in (1/2, 1),$ the following sample path property of the solution is described in [13].

PROPOSITION 2.5. *If $H \in (1/2, 1), S(t)\Phi \in \mathcal{L}_2(V)$ for each $t > 0$ and*

$$(2.30) \quad \int_0^{T_0} \int_0^{T_0} u^{-\alpha} v^{-\alpha} |S(u)\Phi|_{\mathcal{L}_2(V)} |S(v)\Phi|_{\mathcal{L}_2(V)} \phi_H(u-v) du dv < \infty$$

for some $T_0 > 0$ and $\alpha > 0$, then there is a Hölder continuous V -valued version of the process $(\hat{Z}(t), t \geq 0)$ with Hölder exponent $\beta < \alpha$, where \hat{Z} is the stochastic convolution in (2.29) and ϕ_H is given in (2.13). If $(S(t), t \geq 0)$ is an analytic semigroup, then there is a version of the process $(\hat{Z}(t), t \in [0, T])$ with $C^\beta([0, T], V_\delta)$ sample paths for each $T > 0$ and $\beta + \delta < \alpha$.

For each $H \in (0, 1)$, there are the following results for the sample path behavior of the mild solution [28].

PROPOSITION 2.6. *Let $(S(t), t \geq 0)$ be an analytic semigroup, let $H \in (0, 1)$, and let*

$$(2.31) \quad |S(t)\Phi|_{\mathcal{L}_2(V)} \leq ct^{-\gamma}$$

for $t \in [0, T]$, some $c > 0$, and $\gamma \in [0, H)$. Let $\alpha \geq 0$ and $\delta \geq 0$ satisfy

$$(2.32) \quad \alpha + \delta + \gamma < H.$$

Then there is a version of the process $(\hat{Z}(t), t \in [0, T])$ with $C^\alpha([0, T], V_\delta)$ sample paths. If it is assumed instead of (2.31) and (2.32) that $\Phi \in \mathcal{L}_2(V)$ and $\alpha + \delta < H$, then the process $(\hat{Z}(t), t \in [0, T])$ has a $C^\alpha([0, T], V_\delta)$ version. In particular, there is a $C^\alpha([0, T], V)$ version for $0 < \alpha < H$.

3. Semilinear stochastic equations. In this section, both weak and mild solutions are obtained for various semilinear stochastic equations with a fractional Brownian motion. The cases $H \in (0, 1/2)$ and $H \in (1/2, 1)$ are treated separately as in the case of the linear stochastic equations (Propositions 2.5 and 2.6). The weak solution of a semilinear equation is obtained by an absolutely continuous transformation of the measure for the solution of the associated linear equation. The absolute continuity methods given here are an analogue of the results for the measure of a finite-dimensional fractional Brownian motion [7, 9, 25, 26] and the results for Wiener measure [3, 14]. For a fixed $H \in (0, 1)$ and $T > 0$, let $(\mathcal{F}_t, t \in [0, T])$ be the filtration for the standard cylindrical fractional Brownian motion $(B(t), t \in [0, T])$ with the Hurst parameter H . The sub- σ -algebra $\mathcal{F}_t \subset \mathcal{F}$ can be generated by $\sigma(\beta_n(s), s \in [0, t], n \in \mathbb{N})$, where $(\beta_n, n \in \mathbb{N})$ is a sequence of independent scalar fractional Brownian motions with the Hurst parameter H that is given in the definition of a standard cylindrical fractional Brownian motion (Definition 2.1).

The following result describes an absolute continuity for a transformation of a standard cylindrical fractional Brownian motion.

THEOREM 3.1. *Let $H \in (0, 1)$ and $T > 0$ be fixed and let $(u(t), t \in [0, T])$ be a V -valued, (\mathcal{F}_t) -adapted process such that*

1.

$$\int_0^T \|u(t)\| dt < \infty \quad \text{a.s. } \mathbb{P}$$

and

2.

$$U(\cdot) := \int_0^\cdot u(s) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)) \quad \text{a.s. } \mathbb{P}.$$

Furthermore, it is assumed that

$$\mathbb{E}\xi(T) = 1,$$

where

$$(3.1) \quad \xi(T) = \exp \left[\int_0^T \langle \mathbb{K}_H^{-1}(U)(t), dW(t) \rangle - \frac{1}{2} \int_0^T \|\mathbb{K}_H^{-1}(U)(t)\|^2 dt \right],$$

where $(W(t), t \in [0, T])$ is a standard cylindrical Wiener process in V given by (2.19) and \mathbb{K}_H^{-1} is the inverse of the integral operator \mathbb{K}_H in (2.3). Then the process $(\tilde{B}(t), t \in [0, T])$ given by

$$\tilde{B}(t) := B(t) - U(t)$$

is a standard cylindrical fractional Brownian motion in V with the Hurst parameter H on the probability space $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$, where

$$(3.2) \quad \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} = \xi(T) \quad \text{a.s.}$$

Proof. Initially, it is noted that for an (\mathcal{F}_t) -adapted process, $(\eta(t), t \in [0, T])$ with $\eta \in L^2([0, T], V)$ a.s. \mathbb{P} , $\int_0^T \langle \eta, dW \rangle$ is defined by

$$\int_0^T \langle \eta, dW \rangle = \sum_{n=1}^{\infty} \int_0^T \langle \eta, e_n \rangle dw_n,$$

where the sequences $(\beta_n, n \in \mathbb{N})$ and $(w_n, n \in \mathbb{N})$ are related by (2.18). It is shown that $\mathbb{K}_H^{-1}U$ satisfies the conditions of η so that the stochastic integral in (2.20) is well-defined. Recall that the linear operator \mathbb{K}_H given in (2.3) is a bijection

$$\mathbb{K}_H: L^2([0, T], V) \rightarrow I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)),$$

so by assumption 1 in Theorem 3.1, $\mathbb{K}_H^{-1}(U) \in L^2([0, T], V)$ a.s. \mathbb{P} . From the definition of \mathbb{K}_H , it follows that $(\mathbb{K}_H^{-1}(U)(t), t \in [0, T])$ is an (\mathcal{F}_t) -adapted process because U is (\mathcal{F}_t) -adapted. By the construction of the standard cylindrical Wiener process W , it is a Wiener process with respect to (\mathcal{F}_t) so ξ_T is a well-defined random variable. By a Girsanov theorem for Wiener processes in infinite dimensions (cf. [4, 24]), equality (3.2) defines a probability $\tilde{\mathbb{P}}$ on (Ω, \mathcal{F}) such that

$$\tilde{W}(t) := W(t) - \int_0^t \mathbb{K}_H^{-1}(U)(s) ds$$

is a standard cylindrical Wiener process in V . Let

$$\tilde{\beta}_n(t) := \langle B(t), e_n \rangle - \langle U(t), e_n \rangle$$

and

$$\tilde{w}_n(t) = \langle W(t), e_n \rangle - \left\langle \int_0^t \mathbb{K}_H^{-1}(U)(s) ds, e_n \right\rangle.$$

It follows that

$$(3.3) \quad \begin{aligned} \int_0^t K_H(t, s) d\tilde{w}_n(s) &= \int_0^t K_H(t, s) dw_n(s) - \int_0^t K_H(t, s) \langle \mathbb{K}_H^{-1}(U)(s), e_n \rangle ds \\ &= \beta_n(t) - \left\langle \int_0^t K_H(t, s) (\mathbb{K}_H^{-1}(U)(s)) ds, e_n \right\rangle \\ &= \beta_n(t) - \langle \mathbb{K}_H \mathbb{K}_H^{-1}(U)(t), e_n \rangle \\ &= \beta_n(t) - \langle U(t), e_n \rangle = \tilde{\beta}_n(t). \end{aligned}$$

Thus $(\tilde{B}(t), t \in [0, T])$ is a standard cylindrical fractional Brownian motion in V with the Hurst parameter H on $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$. \square

In this section, the following semilinear stochastic equation is considered:

$$(3.4) \quad dX(t) = (AX(t) + F(X(t))) dt + \Phi dB(t),$$

where $t \in \mathbb{R}_+$, $X(t)$, $X_0 \in V$ is nonrandom, $(B(t), t \geq 0)$ is a standard cylindrical fractional Brownian motion with the Hurst parameter $H \in (0, 1)$, $\Phi \in \mathcal{L}(V)$, $A: \text{Dom}(A) \rightarrow V$, $\text{Dom}(A) \subset V$, and A is the infinitesimal generator of a strongly continuous semigroup $(S(t), t \geq 0)$ on V . The function $F: V \rightarrow V$ is nonlinear, and for the applications to stochastic partial differential equations it is more useful to assume that F is defined only on a (dense) subspace of V . So, let $(E, \|\cdot\|_E)$ be a separable Banach space that is continuously embedded in V and $F: E \rightarrow E$ with $X_0 \in E$. Subsequently, it is assumed that $F: E \rightarrow E$ is Borel measurable, $\text{Im}(F) \subset \text{Im}(\Phi)$, for $G := \Phi^{-1}F$, $G: E \rightarrow V$, and

$$(3.5) \quad \|G(x)\| \leq \hat{k}(1 + \|x\|_E^\rho)$$

and

$$(3.6) \quad \|F(x)\|_E \leq \hat{k}(1 + \|x\|_E^\rho)$$

for each $x \in E$ and some $\rho \geq 1$. Furthermore, it is assumed that there is a constant \bar{K} such that for each pair (x, y) in $\text{Dom}(A)$, there is a $z^* \in \partial\|z\|_E$ such that

$$(3.7) \quad \langle Ax - Ay + F(x) - F(y), z^* \rangle_{E, E^*} \leq \bar{K}\|x - y\|_E,$$

where $\partial\|z\|_E$ is the subdifferential of the norm $\|z\|_E$ at the point $z = x - y$ and $\langle \cdot, \cdot \rangle_{E, E^*}$ is the pairing between E and E^* . The basic results on subdifferentials can be found in [32]. Inequality (3.7) is a one-sided growth condition that ensures the absence of explosions of solutions of (3.4) in a finite time. Some subsequent examples should clarify its interpretation.

The notions of a weak and a mild solution of (3.4) are given now.

DEFINITION 3.2. *A weak solution of (3.4) is a triple $(X(t), B(t), (\tilde{\Omega}, \tilde{\mathcal{F}}, (\tilde{\mathcal{F}}_t), \tilde{\mathbb{P}}), t \geq 0)$, where $(B(t), t \geq 0)$ is a standard cylindrical fractional Brownian motion in V that is defined on the probability space $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$, $(B(t), t \geq 0)$ and $(X(t), t \geq 0)$ are adapted to the filtration $(\tilde{\mathcal{F}}_t)$, and $(X(t), t \geq 0)$ is an E -valued process satisfying*

$$(3.8) \quad X(t) = S(t)X_0 + \int_0^t S(t-r)F(X(r)) dr + \int_0^t S(t-r)\Phi dB(r).$$

A mild solution, $(X(t), t \geq 0)$ of (3.4), is an E -valued process on a fixed probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t), \mathbb{P})$ with a given standard cylindrical fractional Brownian motion that is the fractional Brownian motion in (3.8), $(B(t), t \geq 0)$ and $(X(t), t \geq 0)$ are adapted to the filtration (\mathcal{F}_t) , and the process $(X(t), t \geq 0)$ satisfies (3.8).

Equation (3.8) has a unique weak solution if, for any two weak solutions $(X(t), B(t), (\Omega, \mathcal{F}, (\mathcal{F}_t), \mathbb{P}), t \geq 0)$ and $(\tilde{X}(t), \tilde{B}(t), (\tilde{\Omega}, \tilde{\mathcal{F}}, (\tilde{\mathcal{F}}_t), \tilde{\mathbb{P}}), t \geq 0)$, the processes $(X(t), t \geq 0)$ and $(\tilde{X}(t), t \geq 0)$ have the same probability law.

The equation has a unique mild solution if, for any two processes $(X_1(t), t \geq 0)$ and $(X_2(t), t \geq 0)$ that satisfy (3.8) on the same probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t), \mathbb{P})$ with the same standard cylindrical fractional Brownian motion, $\mathbb{P}(X_1(t) = X_2(t), t \geq 0) = 1$.

A primary goal in this section is to verify weak existence and weak uniqueness of a solution of (3.4). Note that (H1) alone is not sufficient to ensure that the stochastic convolution has values in E . While the assumption (H2) is given in a rather general form, it is verified for particular examples in section 4. Since the cases $H \in (0, 1/2)$ and $H \in (1/2, 1)$ require different methods, they are treated separately.

The following three assumptions are made to construct a solution of (3.4):

- (H1) The semigroup $(S(t), t \geq 0)$ generated by A is analytic on V and for each $t \geq 0$, $S(t)|_E \in \mathcal{L}(E)$ and $\|S(t)|_E\|_{\mathcal{L}(E)}$ is bounded on compact time intervals.
- (H2) $\Phi \in \mathcal{L}(V)$ is injective and for $T > 0$ the stochastic convolution process

$$\left(\int_0^t S(t-r)\Phi dB(r), t \in [0, T] \right)$$

has a version with $C([0, T], E)$ sample paths.

- (H3) The function $F: E \rightarrow V$ in (3.4) is Borel measurable, $\text{Im}(F) \subset \text{Im}(\Phi)$, and the function $G = \Phi^{-1}F: E \rightarrow V$ satisfies

$$(3.9) \quad \|G(x)\| \leq k(1 + \|x\|_E)$$

for some $k > 0$ and all $x \in E$.

The following result verifies a weak solution for $H \in (0, 1/2)$.

THEOREM 3.3. *If $H \in (0, 1/2)$ and the conditions (H1)–(H3) are satisfied, then the semilinear equation (3.4) has a weak solution. If additionally $F: E \rightarrow E$ and*

$$(3.10) \quad \|F(x)\|_E \leq k_1(1 + \|x\|_E)$$

for some $k_1 > 0$ and all $x \in E$, then the weak solution is unique.

Proof. Initially, existence of a weak solution is verified. By a standard method that has been used for equations of the form (3.4) with a standard cylindrical Brownian motion (cf. [4, 24]), it suffices to verify that the cylindrical process

$$\tilde{B}(t) = B(t) - \int_0^t G(Z(s)) ds$$

is a standard cylindrical fractional Brownian motion in a suitable probability space where

$$Z(t) = S(t)X_0 + \tilde{Z}(t)$$

satisfies the associated linear equation. To use Theorem 3.1 it is necessary to verify that $G = \Phi^{-1}F$ satisfies the conditions of U in this theorem, that is,

$$(3.11) \quad \int_0^\cdot G(Z(s)) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V))$$

and

$$(3.12) \quad \mathbb{E} \exp[\rho(Z)] = 1,$$

where

$$(3.13) \quad \rho(Z) = \int_0^T \left\langle \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) (t), dW(t) \right\rangle - \frac{1}{2} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) (t) \right\|^2 dt,$$

\mathbb{K}_H^{-1} is the inverse of \mathbb{K}_H in (2.3), and $(W(t), t \geq 0)$ is a standard cylindrical Wiener process in V by (2.19).

From (2.9), it follows that

$$\begin{aligned}
 (3.14) \quad & \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)}^2 \\
 &= \bar{c}_H^{-2} \left| u_{H-\frac{1}{2}} I_{0+}^{\frac{1}{2}-H} \left(u_{\frac{1}{2}-H} G(Z) \right) \right|_{L^2([0,T],V)}^2 \\
 &= \hat{c}_H \int_0^T \left(s^{H-\frac{1}{2}} \left\| \int_0^s r^{\frac{1}{2}-H} (s-r)^{-\frac{1}{2}-H} G(Z(r)) dr \right\| \right)^2 ds \\
 &\leq \hat{c}_H k^2 \left(1 + |\tilde{Z}|_{C([0,T],E)} + \sup_{t \in [0,T]} \|S(t)X_0\|_E \right)^2 \int_0^T s^{2H-1} \\
 &\quad \cdot \left(\int_0^s r^{\frac{1}{2}-H} (s-r)^{-\frac{1}{2}-H} dr \right)^2 ds \\
 &\leq c_T \left(1 + |\tilde{Z}|_{C([0,T],E)}^2 \right)
 \end{aligned}$$

for some $c_T > 0$ that depends only on T . This inequality verifies (3.11). By (3.14) it follows directly that

$$(3.15) \quad \mathbb{E} \exp \left[\hat{k} \int_0^T \left\| K_H^{-1} \left(\int_0^\cdot G(Z) \right) (t) \right\|^2 dt \right] \leq c \mathbb{E} \exp \left[\hat{k} c_T |\tilde{Z}|_{C([0,T],E)}^2 \right]$$

for some $c > 0$. Substituting $v = \frac{t}{s}$ in the integral with respect to r on the right-hand side of (3.14), it easily follows that $c_T \downarrow 0$ as $T \downarrow 0$. Since \tilde{Z} is a $C([0, T], E)$ -valued Gaussian random variable, it follows that

$$(3.16) \quad \mathbb{E} \exp \left[\hat{k} c_T |\tilde{Z}|_{C([0,T],E)}^2 \right] < \infty$$

is satisfied for $T > 0$ sufficiently small by the Fernique inequality. Clearly, (3.16) is the Novikov condition [19] which implies the equality (3.12) for $T > 0$ sufficiently small. For arbitrary $T > 0$, a simple iteration verifies the result, that is,

$$(3.17) \quad \mathbb{E} \exp \left[\hat{k} \int_{T_{m-1}}^{T_m} \left\| K_H^{-1} \left(\int_0^\cdot G(Z) \right) (t) \right\|^2 dt \right] < \infty$$

for a sufficiently fine partition $0 = T_0 < T_1 < \dots < T_n = T$. Using a downward induction procedure from the well-known proofs of the martingale property for the Radon–Nikodým derivative in (3.12) for an arbitrary $T > 0$ (see, e.g., [20, Example 6.2.3]), the verification of the equality in (3.12) is obtained.

Now, uniqueness of the weak solution is verified. Uniqueness in law can be proved in a standard way by removing the term F in (3.4) by absolute continuity of measures, which is a suitable inverse of the above construction of a weak solution.

Let $(\tilde{X}(t), t \in [0, T])$ be a solution to the equation

$$(3.18) \quad \tilde{X}(t) = S(t)x_0 + \int_0^t S(t-r)F(X(r)) dr + \tilde{Z}(t),$$

where $\tilde{Z}(t) = \int_0^t S(t-r)\Phi dB(r)$ and $(B(t), t \in [0, T])$ is some standard cylindrical fractional Brownian motion on a probability space $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$.

The process $(\tilde{X}(t), t \in [0, T])$ is defined on the same probability space as $(B(t), t \in [0, T])$. Let $(W(t), t \in [0, T])$ be the Wiener process associated with $(B(t), t \in [0, T])$ by (2.18). It suffices to show that

(3.19)

$$\begin{aligned} & \exp \left[\tilde{\rho}(\tilde{X}) \right] \\ & := \exp \left[- \int_0^T \left\langle \mathbb{K}_H^{-1} \left(\int_0^\cdot G(\tilde{X}) \right) (t), dW(t) \right\rangle - \frac{1}{2} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(\tilde{X}) \right) (t) \right\|^2 dt \right] \end{aligned}$$

is a Radon–Nikodým derivative on $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$, so $\tilde{\mathbb{P}}$ is the measure for a fractional Ornstein–Uhlenbeck process and uniqueness in law follows. Thus it is necessary to show that

(3.20)
$$\int_0^\cdot G(\tilde{X}(s)) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V))$$

and

(3.21)
$$\tilde{\mathbb{E}} \exp \left[\tilde{\rho}(\tilde{X}) \right] = 1,$$

where \tilde{E} is integration with respect to $\tilde{\mathbb{P}}$. The verifications of (3.20) and (3.21) are analogous to the verifications of (3.11) and (3.12), respectively. However, since \tilde{X} is not a Gaussian process, the Fernique inequality cannot be used directly. Initially, it is verified that there is a $c > 0$ such that

(3.22)
$$|\tilde{X}|_{C([0, T], E)} \leq c \left(1 + \|X_0\|_E + |\tilde{Z}|_{C([0, T], E)} \right),$$

where \tilde{Z} is the stochastic process described in (H2). Let

$$\begin{aligned} u(t) &= \tilde{X}(t) - \tilde{Z}(t) \\ &= S(t)X_0 + \int_0^t S(t-r)F(u(r) + \tilde{Z}(r)) dr. \end{aligned}$$

Thus

(3.23)
$$\|u(t)\|_E \leq c_1 \|X_0\| + c_2 \int_0^t \left(1 + \|u(r)\|_E + \|\tilde{Z}(r)\|_E \right) dr$$

for some positive constants c_1 and c_2 . By the Gronwall lemma it follows that

(3.24)
$$\|u(t)\|_E \leq c_1 \left(1 + \|X_0\|_E + |\tilde{Z}|_{C([0, T], E)} \right)$$

for $t \in [0, T]$, so the inequality (3.22) is verified. The exponential that usually occurs in the Gronwall inequality is bounded by $e^{c_2 T}$. Making the analogous computations in (3.14), it follows that

(3.25)

$$\left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) (s) \right|_{L^2([0, T], V)}^2 \leq c_T \left(1 + |X|_{C([0, T], E)}^2 \right) \leq \tilde{c}_T \left(1 + |\tilde{Z}|_{C([0, T], E)}^2 \right),$$

where $\tilde{c}_T \downarrow 0$ as $T \downarrow 0$, so (3.20) is satisfied. Thus the method in (3.15)–(3.17) can be used to verify (3.21).

The random variable $\exp(\tilde{\rho}(\tilde{X}))$ in (3.19) is a Radon–Nikodým derivative and it defines a probability measure \mathbb{Q} on $\tilde{\Omega}$. By this Girsanov-type theorem the process defined by $\tilde{B}(t) = B(t) + \int_0^t u(s)ds$, where $u(s) = G(\tilde{X}(s))$, is a standard cylindrical fractional Brownian motion with respect to the measure \mathbb{Q} . Let $(\tilde{W}(t), t \in [0, T])$ be the Wiener process associated with $(\tilde{B}(t), t \in [0, T])$. Let $U(s) = \mathbb{K}_H^{-1}(\int_0^s G(\tilde{X})(s))$ for $s \in [0, T]$ and let \mathbb{E} and $\mathbb{E}_\mathbb{Q}$ denote the expectations with respect to the measures $\tilde{\mathbb{P}}$ and \mathbb{Q} , respectively. For a bounded measurable function Ψ on $C([0, T], V)$ it follows that

$$\begin{aligned} \mathbb{E}[\Psi(\tilde{X})] &= \int_{\tilde{\Omega}} \Psi \frac{d\tilde{\mathbb{P}}}{d\mathbb{Q}} d\mathbb{Q} = \mathbb{E}_\mathbb{Q}[\Psi(\tilde{X}) \exp(-\tilde{\rho}(\tilde{X}))] \\ &= \mathbb{E}_\mathbb{Q} \left[\Psi(\tilde{X}) \exp \left(\int_0^T \langle U(r), d\tilde{W}(r) \rangle - \frac{1}{2} \int_0^T \|U(r)\|^2 dr \right) \right] \\ &= \mathbb{E}_\mathbb{Q} \left[\Psi \left(S(\cdot)X_0 + \int_0^\cdot S(\cdot - r) d\tilde{B}(r) \right) \exp \left(\int_0^T \left\langle \mathbb{K}_H^{-1} \left(\int_0^\cdot G \left(S(\cdot)X_0 \right. \right. \right. \right. \right. \right. \\ &\quad \left. \left. \left. \left. \left. + \int_0^\cdot S(\cdot - r) \Phi d\tilde{B}(r) \right) \right) (s), d\tilde{W}(s) \right\rangle \right. \\ &\quad \left. \left. - \frac{1}{2} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G \left(S(\cdot)X_0 + \int_0^\cdot S(\cdot - r) \Phi d\tilde{B}(r) \right) \right) (s) \right\|^2 ds \right) \right]. \end{aligned}$$

Since the processes \tilde{W} and \tilde{B} are standard cylindrical Brownian motions and standard cylindrical fractional Brownian motions, respectively, the final expectation on the right-hand side above does not depend on the realization of \tilde{X} , so the uniqueness in law is verified. \square

Now the existence and the uniqueness of a weak solution of (3.4) is verified for $H \in (1/2, 1)$.

THEOREM 3.4. *If $H \in (1/2, 1)$, (H1)–(H3) are satisfied, and*

$$(3.26) \quad \|G(x) - G(y)\| \leq k_G \|x - y\|^\gamma$$

for all $x, y \in E$, some $\gamma \in (0, 1]$, $k_G > 0$, and $\tilde{Z} \in C^\beta([0, T], V)$ for some β satisfying

$$(3.27) \quad \beta > \frac{H - \frac{1}{2}}{\gamma},$$

where \tilde{Z} is the stochastic convolution process in (H2), then (3.4) has a weak solution. If, additionally, (3.10) is satisfied, then the weak solution is unique.

Proof. Initially, the existence of a solution is verified as in the proof of Theorem 3.3. It is shown that

$$(3.28) \quad \int_0^\cdot G(Z(s)) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)) \quad \text{a.s.}$$

and

$$(3.29) \quad \mathbb{E} \exp[\rho(Z)] = 1,$$

where ρ is given by (3.13). By (2.8) it follows that

$$\begin{aligned}
 (3.30) \quad & \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)}^2 \\
 &= c_H^{-2} \left| u_{H-\frac{1}{2}} D_{0+}^{H-\frac{1}{2}} \left(u_{\frac{1}{2}-H} G(Z) \right) \right|_{L^2([0,T],V)}^2 \\
 &= c_H^{-2} \int_0^T \left\| \frac{s^{H-\frac{1}{2}}}{\Gamma(\frac{3}{2}-H)} \left(\frac{s^{\frac{1}{2}-H} G(Z(s))}{s^{H-\frac{1}{2}}} \right) \right. \\
 &\quad \left. + \left(H - \frac{1}{2} \right) \int_0^s \frac{s^{\frac{1}{2}-H} G(Z(s)) - r^{\frac{1}{2}-H} G(Z(r))}{(s-r)^{H+\frac{1}{2}}} dr \right\|^2 ds \\
 &\leq c \int_0^T \left(s^{\frac{1}{2}-H} \|G(Z(s))\| + s^{H-\frac{1}{2}} \int_0^s \frac{s^{\frac{1}{2}-H} - r^{\frac{1}{2}-H}}{(s-r)^{H+\frac{1}{2}}} \|G(Z(r))\| dr \right. \\
 &\quad \left. + \int_0^s \frac{\|G(Z(s)) - G(Z(r))\|}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 ds.
 \end{aligned}$$

Using (3.9) and (3.26), the analyticity of the semigroup $S(\cdot)$ on V , and the inequality

$$\int_0^s \frac{s^{\frac{1}{2}-H} - r^{\frac{1}{2}-H}}{(s-r)^{H+\frac{1}{2}}} dr \leq cs^{1-2H},$$

where c is a generic constant, it follows that

$$\begin{aligned}
 (3.31) \quad & \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)} \\
 &\leq c_T \left(1 + \|X_0\|_E^2 + |\tilde{Z}|_{C([0,T],E)}^2 \right) \\
 &\quad + c_T \int_0^T \left[\left(\int_0^s \frac{\|S(s)X_0 - S(r)X_0\|^\gamma}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 + \left(\int_0^s \frac{\|\tilde{Z}(s) - \tilde{Z}(r)\|^\gamma}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 \right] ds \\
 &\leq c_T \left(1 + \|X_0\|_E^2 + |\tilde{Z}|_{C([0,T],E)}^2 \right) + c_T \int_0^T \left(\|X_0\|^\gamma \int_0^s \frac{(s-r)^{\gamma\lambda}}{r^{\gamma\lambda}(s-r)^{H+\frac{1}{2}}} dr \right)^2 \\
 &\quad + c_T \int_0^T |\tilde{Z}|_{C^\beta([0,T],V)}^2 \left(\int_0^s \frac{(s-r)^{\gamma\beta}}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 ds,
 \end{aligned}$$

where $\lambda > 0$ satisfies $\gamma\lambda < 1$ and $H + 1/2 - \gamma\lambda < 1$. The first integral term in the initial inequality in (3.31) is obtained by the analyticity of the semigroup $S(\cdot)$ on V , which implies that

$$\|(S(s) - S(r))x\| \leq c(s-r)^\lambda r^\lambda \|x\|, \quad x \in V, \lambda \geq 0, 0 < r < s \leq T,$$

for some $c = c(\lambda)$. It follows that

$$(3.32) \quad \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)}^2 \leq c_T \left(1 + \|X_0\|_E^2 + \|\tilde{Z}\|_{C([0,T],E)}^2 + |\tilde{Z}|_{C^\beta([0,T],V)}^2 \right),$$

where $c_T \downarrow 0$ as $T \downarrow 0$, so (3.28) is verified and by the Fernique inequality (3.29) is also verified.

Now the uniqueness of the weak solution is verified. Let $(\tilde{X}(t), t \in [0, T])$ be the solution to (3.4) on a probability space $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$. As in the proof of Theorem 3.3, it is shown that

$$(3.33) \quad \int_0^\cdot G(\tilde{X}) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V)) \quad \text{a.s.}$$

and

$$(3.34) \quad \tilde{\mathbb{E}} \exp \left[\tilde{\rho}(\tilde{X}) \right] = 1,$$

where $\tilde{\rho}$ is given by (3.13). It is necessary to obtain inequality (3.22), used in the proof of Theorem 3.3. Inequality (3.22) is verified by verifying the inequality

$$(3.35) \quad \|\hat{X}\|_{C^\beta([0,T],V)} \leq L \left(1 + \|X_0\|_E + |\tilde{Z}|_{C([0,T],E)} + |\tilde{Z}|_{C^\beta([0,T],V)} \right),$$

where $\hat{X}(t) = \tilde{X}(t) - S(t)X_0$ and $L > 0$. Let $w(t) = \tilde{X}(t) - S(t)X_0 - \tilde{Z}(t)$ for $t \geq 0$. The process w satisfies

$$(3.36) \quad w(t) = \int_0^t S(t-r)\psi(r) dr$$

for $t \in [0, T]$, where

$$\psi(t) = F(w(t) + S(t)X_0 + \tilde{Z}(t)).$$

By inequalities (3.10) and (3.35) it follows that $\psi \in L^\infty([0, T], V)$ a.s. \mathbb{P} . Since the semigroup $S(\cdot)$ is analytic on V , w is α -Hölder continuous for each $\alpha \in (0, 1)$, and using the method of proof of [29, Theorem 4.3.1] there are constants $c_i > 0$ for $i = 1, 2$ such that

$$(3.37) \quad |w|_{C^\beta([0,T],V)} \leq c_1 |\psi|_{L^\infty([0,T],V)} \leq c_2 \left(|w|_{L^\infty([0,T],E)} + \|X_0\|_E + |\tilde{Z}|_{C([0,T],E)} \right).$$

Thus

$$(3.38) \quad \begin{aligned} |\tilde{X}|_{C^\beta([0,T],V)} &\leq |w|_{C^\beta([0,T],V)} + |\tilde{Z}|_{C^\beta([0,T],V)} \\ &\leq c_2 \left(|w|_{L^\infty([0,T],E)} + \|X_0\|_E + |\tilde{Z}|_{C([0,T],E)} + |\tilde{Z}|_{C^\beta([0,T],V)} \right). \end{aligned}$$

Using (3.22) again to bound $|w|_{L^\infty([0,T],E)}$, inequality (3.35) follows.

Now, using the methods for inequalities (3.30)–(3.32), where $Z(t)$ is replaced by $\tilde{X}(t) = S(t)X_0 + \hat{X}(t)$, it follows that

(3.39)

$$\left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)}^2 \leq c_T \left(1 + \|X_0\|_E^2 + |\tilde{X}|_{C([0,T],V)}^2 + |\tilde{X}|_{C^\beta([0,T],V)}^2 \right),$$

where $c_T \downarrow 0$ as $T \downarrow 0$, which by (3.22) and (3.35) verifies (3.33). Equality (3.34) is obtained from (3.39) by the Fernique inequality as in the proof of Theorem 3.3. \square

Remark 3.5. The proofs of Theorems 3.3 and 3.4 have verified, in addition to weak existence and uniqueness of a solution to (3.4), the mutual absolute continuity (equivalence) of the probability laws of the solution to (3.4) and the solution of (3.4) with $F \equiv 0$ (the fractional Ornstein–Uhlenbeck process) in the path space.

The next objective is to relax the linear growth conditions (3.9) and (3.10) and the Hölder continuity (3.26). The linear growth condition is replaced by a dissipativity condition of the drift term of (3.4), but some other conditions are also imposed so that there is existence and (strong) uniqueness of a mild solution. The main contribution of the following two theorems is a mutual absolute continuity of the probability laws of the solutions of (3.4) with a nonzero F and (3.4) with $F \equiv 0$.

Initially, the case $H \in (0, 1/2)$ is considered.

THEOREM 3.6. *Let $H \in (0, 1/2)$ and let (H1) and (H2) be satisfied. Let $\Phi \in \mathcal{L}(V)$ be injective, let $\Phi^{-1} \in \mathcal{L}(E, V)$, and let $(S(t)|_E, t \geq 0)$ be a strongly continuous semigroup on E such that*

$$(3.40) \quad |S(t)|_E|_{\mathcal{L}(E)} \leq e^{\tilde{w}t}$$

for $t \geq 0$ and some $\tilde{w} \in \mathbb{R}$. Let $F: E \rightarrow E$ be continuous and satisfy

$$(3.41) \quad \|F(x)\|_E \leq k_1 (1 + \|x\|_E^\rho)$$

for $x \in E$ for some $k_1 \geq 0$ and $\rho \geq 1$, and for each pair $x, y \in E$, there is a $z^* \in \partial\|x - y\|_E$ where $\partial\|z\|_E$ is the subdifferential of the norm $\|\cdot\|_E$ at $z \in E$ such that

$$(3.42) \quad \langle F(x) - F(y), z^* \rangle_{E, E^*} \leq k_2 \|x - y\|_E$$

for some $k_2 \in \mathbb{R}$; that is, $F - k_2I$ is dissipative on E . Then there is one and only one mild solution of (3.4), and its probability law on the Borel σ -algebra of $\tilde{\Omega} = C([0, T], E)$ is mutually absolutely continuous with respect to the probability law of the fractional Ornstein–Uhlenbeck process (3.24) on Ω .

Proof. Let $(F_\lambda, \lambda > 0)$ be a family of Lipschitz continuous functions from E to E such that each F_λ satisfies inequalities (3.41) and (3.42) for F with the same constants ρ, k_1, k_2 . It is shown that there is a $\bar{k} > 0$ depending only on \tilde{w}, k_1 , and k_2 such that

$$(3.43) \quad \|v_\lambda(t)\|_E \leq \bar{k} \left(1 + \|X_0\|_E + \|\phi\|_{C([0,T],E)}^\rho \right)$$

for $t \in [0, T]$ is satisfied for each $\lambda > 0$ and $\phi \in C([0, T], E)$, where v_λ is a solution of the equation

$$(3.44) \quad v_\lambda(t) = S(t)X_0 + \int_0^t S(t-r)F_\lambda(v_\lambda(r) + \phi(r)) \, dr$$

for $t \in [0, T]$.

To verify inequality (3.43), it can be assumed by translation that $k_2 = 0$ in (3.42) (replace F_λ and A by $F_\lambda - k_2I$ and $A + k_2I$, respectively). Thus F_λ is dissipative on E for each $\lambda > 0$ and by the assumptions

$$(3.45) \quad \langle A_E z, z^* \rangle_{E, E^*} \leq \tilde{w} \|z\|_E^2$$

for each $z \in \text{Dom}(A_E)$ and $z^* \in \partial\|z\|_E$, where A_E is the restriction of A to E that generates the semigroup $S(\cdot)|_E$. For each pair $x, y \in \text{Dom}(A_E)$ and $\lambda > 0$, there is a $z_\lambda^* \in \partial\|x - y\|_E$ such that

$$\langle A_E(x - y) + F_\lambda(x) - F_\lambda(y), z_\lambda^* \rangle_{E, E^*} \leq \tilde{w} \|x - y\|_E.$$

By [5, Proposition 5.5.6], there is a sequence $(v_\lambda^n, n \in \mathbb{N})$ such that $v_\lambda^n \in C^1([0, T], E) \cap C([0, T], \text{Dom}(A_E))$ such that $v_\lambda^n \rightarrow v_\lambda$ and $\delta_\lambda^n = \frac{d}{dt} v_\lambda^n - A_E v_\lambda^n - F_\lambda(v_\lambda^n + \phi) \rightarrow 0$ in $C([0, T], E)$ as $n \rightarrow \infty$. It follows that

$$(3.46) \quad \begin{aligned} \frac{d^-}{dt} \|v_\lambda^n(t)\|_E &\leq \langle A_E v_\lambda^n(t) + F_\lambda(v_\lambda^n(t) + \phi(t)), (v_\lambda^n(t))^* \rangle_{E, E^*} + \|\delta_\lambda^n(t)\|_E \\ &= \langle A_E v_\lambda^n(t) + F_\lambda(v_\lambda^n(t) + \phi(t)) - F_\lambda(\phi(t)), (v_\lambda^n(t))^* \rangle \\ &\quad + \langle F_\lambda(\phi(t)), (v_\lambda^n(t))^* \rangle_{E, E^*} + \|\delta_\lambda^n(t)\|_E \\ &\leq \bar{w} \|v_\lambda^n(t)\|_E + k_2 \left(1 + |\phi|_{C([0, T], E)}^\rho + \|\delta_\lambda^n(t)\|_E \right) \end{aligned}$$

for $t \in [0, T]$. Using the Gronwall lemma, and letting $n \rightarrow \infty$, verifies inequality (3.43).

The mild solution to (3.4) can be expressed as $X(t) = v(t) + \tilde{Z}(t)$, where v satisfies the equation

$$(3.47) \quad v(t) = S(t)X_0 + \int_0^t S(t-r)F(v(r) + \tilde{Z}(r)) dr$$

for $t \in [0, T]$. Thus the existence and the uniqueness of a mild solution follows from the corresponding pathwise deterministic result (cf. [5, Proposition 5.5.6]).

The equivalence of the probability laws is shown by application of Theorem 3.1. As in the proof of Theorem 3.3, it suffices to show that

$$(3.48) \quad \int_0^\cdot G(Z(s)) ds \in I_{0+}^{H+\frac{1}{2}} (L^2([0, T], V))$$

and

$$(3.49) \quad \mathbb{E} \exp [\rho(Z)] = 1,$$

where ρ is given by (3.13). While G is not assumed to have at most linear growth as in Theorem 3.3, there is the growth condition

$$(3.50) \quad \|G(x)\| \leq \hat{k} (1 + \|x\|_E^\rho)$$

for all $x \in E$ and a constant \hat{k} . Proceeding as in (3.14), it follows that

(3.51)

$$\begin{aligned} & \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,t],V)}^2 \\ & \leq c_1 \int_0^T \left(s^{H-\frac{1}{2}} \left\| \int_0^s r^{\frac{1}{2}-H} (s-r)^{-\frac{1}{2}-H} G(Z(r)) dr \right\| \right)^2 ds \\ & \leq c_2 \left(1 + |\tilde{Z}|_{C([0,T],E)}^\rho + \sup_{t \in [0,T]} \|S(t)X_0\|_E^\rho \right) \int_0^T s^{2H-1} \\ & \quad \cdot \left(\int_0^s r^{\frac{1}{2}-H} (s-r)^{-\frac{1}{2}-H} dr \right)^2 ds \\ & \leq c_3 \left(1 + \|X_0\|_E^{2\rho} + |\tilde{Z}|_{C([0,T],E)}^{2\rho} \right) \end{aligned}$$

for suitable constants c_1, c_2, c_3 . This inequality verifies (3.48). To verify equality (3.49), it suffices to assume that F is dissipative (that is, $k_2 = 0$ in (3.42)). Recall that m -dissipative mapping F is defined as a dissipative mapping satisfying $\text{Range}(I - \lambda F) = E$ for each $\lambda > 0$; cf. [22]. Since F is continuous, it is m -dissipative (cf. [23]), so the family $(F_\lambda, \lambda > 0)$ of Yosida approximations of F is defined as

$$(3.52) \quad F_\lambda(x) = F(R_\lambda(x)) = \frac{1}{\lambda}(R_\lambda(x) - x)$$

for $x \in E$, where

$$(3.53) \quad R_\lambda(x) = (I - \lambda F)^{-1}(x).$$

It is well known that $F_\lambda: E \rightarrow E$ for $\lambda > 0$ is Lipschitz continuous, so by Theorem 3.3, there is the equality

$$(3.54) \quad \mathbb{E} \exp[\rho_\lambda(Z)] = 1$$

for $\lambda > 0$, where

(3.55)

$$\rho_\lambda(Z) = \int_0^T \left\langle \mathbb{K}_H^{-1} \left(\int_0^\cdot G_\lambda(Z) \right) (t), dW(t) \right\rangle - \frac{1}{2} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_\lambda(Z) \right) (t) \right\|^2 dt$$

and $G_\lambda := \Phi^{-1}F_\lambda$. As in (3.51), it follows that

$$(3.56) \quad \begin{aligned} & \mathbb{E} \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot (G_\lambda(Z) - G(Z)) \right) \right|_{L^2([0,T],V)} \\ & \leq c_T \mathbb{E} \int_0^T \left(s^{H-\frac{1}{2}} \int_0^s r^{\frac{1}{2}-H} (s-r)^{-\frac{1}{2}-H} \|G_\lambda(Z(r)) - G(Z(r))\| dr \right)^2 ds. \end{aligned}$$

By some well-known properties of the Yosida approximations and for $x \in E$,

$$(3.57) \quad \|G_\lambda(x) - G(x)\| \leq |\Phi^{-1}|_{\mathcal{L}(E,V)} \|F_\lambda(x) - F(x)\|,$$

it follows that $F_\lambda \rightarrow F$ as $\lambda \rightarrow 0$ and the right-hand side of (3.57) tends to zero as $\lambda \downarrow 0$, and

$$\begin{aligned}
 (3.58) \quad \|G_\lambda(x)\| &\leq |\Phi^{-1}|_{\mathcal{L}(E,V)} \|F_\lambda(x)\|_E \\
 &\leq |\Phi^{-1}|_{\mathcal{L}(E,V)} \|F(x)\|_E \\
 &\leq |\Phi^{-1}|_{\mathcal{L}(E,V)} k_1 (1 + \|x\|_E^\rho),
 \end{aligned}$$

so the right-hand side of (3.56) tends to zero as $\lambda \downarrow 0$. For a sequence $(\lambda_n, n \in \mathbb{N})$ that decreases to zero, it follows that

$$(3.59) \quad \lim_{n \rightarrow \infty} \exp[\rho_{\lambda_n}(Z)] = \exp[\rho(Z)] \quad \text{a.s. } \mathbb{P}.$$

To obtain equality (3.54) from equality (3.59) for $\lambda_n, n \in \mathbb{N}$, it is necessary and sufficient to show that the sequence $(\exp[\rho_{\lambda_n}(Z)], n \in \mathbb{N})$ is uniformly integrable. A sufficient condition for this uniform integrability is to verify that

$$(3.60) \quad \sup_n \mathbb{E}[(\exp[\rho_{\lambda_n}(Z)]) |\log(\exp[\rho_{\lambda_n}(Z)])|] = \sup_n \mathbb{E}[(\exp[\rho_{\lambda_n}(Z)]) \rho_{\lambda_n}(Z)] < \infty.$$

By Theorem 3.3,

$$(3.61) \quad \mathbb{E}[\rho_{\lambda_n}(Z) \exp[\rho_{\lambda_n}(Z)]] \leq \tilde{\mathbb{E}}_{\lambda_n} \left[2 \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(Z) \right) (t) \right\|^2 dt \right],$$

where $\tilde{\mathbb{E}}_{\lambda_n}$ is the expectation with respect to $\tilde{\mathbb{P}}_{\lambda_n}$ and

$$\frac{d\tilde{\mathbb{P}}_{\lambda_n}}{d\mathbb{P}} = \exp[\rho_{\lambda_n}(Z)],$$

and $Z(\cdot)$ satisfies (2.28). On the probability space with the measure \mathbb{P}_{λ_n} , $Z(\cdot)$ satisfies the following semilinear equation, where $B(\cdot)$ is a fractional Brownian motion with respect to \mathbb{P}_{λ_n} :

$$\begin{aligned}
 (3.62) \quad dX_{\lambda_n}(t) &= (AX(t) + F_{\lambda_n}(X(t))) dt + \Phi dB(t), \\
 X_{\lambda_n}(0) &= X_0.
 \end{aligned}$$

Since F_{λ_n} is Lipschitz continuous, there is a unique mild solution on a given probability space, so it suffices to show

$$(3.63) \quad \mathbb{E} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(X_{\lambda_n}) \right) (t) \right\|^2 dt \leq c$$

for some $c \in \mathbb{R}_+$ that does not depend on λ_n . Repeating inequalities (3.51), where G and Z are replaced by G_{λ_n} and X_{λ_n} , respectively, and using inequality (3.58), it follows that

$$(3.64) \quad \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(X_{\lambda_n}) \right) (t) \right\|^2 dt \leq c_5 \left(1 + \|X_0\|_E^{2\rho} + |\tilde{X}_{\lambda_n}|_{C([0,T],E)}^{2\rho} \right)$$

for a constant c_5 that does not depend on $n \in \mathbb{N}$ where $\tilde{X}_{\lambda_n}(t) = X_{\lambda_n}(t) - S(t)X_0$. By inequality (3.43) there is a constant c_6 that does not depend on n such that

(3.65)

$$\mathbb{E} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(X_{\lambda_n}) \right) (t) \right\|^2 dt \leq c_6 \left(1 + \|X_0\|_E^{2\rho} + \mathbb{E} |Z|_{C([0,T],E)}^{4\rho^2} \right) = C < \infty.$$

This inequality verifies (3.60). Thus the sequence $(\exp[\rho_{\lambda_n}(Z)], n \in \mathbb{N})$ converges in L^1 and equality (3.54) is satisfied. \square

Now the case $H \in (1/2, 1)$ is considered.

THEOREM 3.7. *Let $H \in (1/2, 1)$ and the other assumptions in Theorem 3.6 be satisfied. Let $\Phi^{-1} \in \mathcal{L}(V)$, $\tilde{Z} \in C^\beta([0, T], V)$ for some $\beta \in (0, 1)$,*

$$(3.66) \quad \langle F(x) - F(y), x - y \rangle \leq k_2 \|x - y\|^2$$

for each pair $x, y \in E$ and a $k_2 \in \mathbb{R}_+$ (that is, $F - k_2I$ is dissipative on E with respect to the norm on V) and

$$(3.67) \quad \|F(x) - F(y)\| \leq k_3 (1 + \|x\|_E^q + \|y\|_E^q) \|x - y\|^\gamma$$

for each $x, y \in E$, with some $k_3 > 0$, $q \geq 1$, and $\gamma \in (0, 1]$ such that

$$(3.68) \quad \gamma\beta > H - \frac{1}{2}.$$

Then there is one and only one mild solution to (3.4), and its probability law is mutually absolutely continuous with respect to the probability law of the fractional Ornstein-Uhlenbeck process (2.28) on Ω .

Proof. As in the proof of Theorem 3.6, it is shown that

$$(3.69) \quad \int_0^\cdot G(Z(s)) ds \in I_{0+}^{H+\frac{1}{2}}(L^2([0, T], V))$$

and

$$(3.70) \quad \mathbb{E} \exp[\rho(Z)] = 1.$$

The methods to verify (3.69) and (3.70) are similar to those used in the proof of Theorem 3.6, but now the operator \mathbb{K}_H^{-1} has a different form. Using inequality (3.41) and the Hölder continuity condition (3.67), it follows that

$$(3.71) \quad \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right\|_{L^2([0,T],V)}^2 \leq c_1 \int_0^T \left(s^{\frac{1}{2}-H} \|G(Z(s))\| + s^{H-\frac{1}{2}} \int_0^s \frac{s^{\frac{1}{2}-H} - r^{\frac{1}{2}-H}}{(s-r)^{H+\frac{1}{2}}} \|G(Z(r))\| dr + \int_0^s \frac{\|G(Z(s)) - G(Z(r))\|}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 ds \leq c_2 \left[1 + |Z|_{C([0,T],E)}^{2\rho} + c_3 \left(1 + |Z|_{C([0,T],E)}^{2q} \right) \cdot \int_0^T \left(\int_0^s \frac{\|S(s)X_0 - S(r)X_0\|^\gamma + \|\tilde{Z}(s) - \tilde{Z}(r)\|^\gamma}{(s-r)^{H+\frac{1}{2}}} dr \right)^2 ds \right]$$

for some constants c_1, c_2, c_3 . By the analyticity of the semigroup $S(\cdot)$ on V , it follows that

$$\begin{aligned}
 (3.72) \quad & \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G(Z) \right) \right|_{L^2([0,T],V)}^2 \\
 & \leq c_4 \left[1 + \|X_0\|_E^{2\rho} \right. \\
 & \quad \left. + \left(\|X_0\|_E^{2q} + |Z|_{C([0,T],V)}^{2q} + 1 \right) \left(\|X_0\|_E^{2\gamma} + |Z|_{C^\beta([0,T],V)}^{2\gamma} \right) \right] \\
 & \leq c_5 \left(1 + \|X_0\|_E^m + |Z|_{C([0,T],V)}^m + |Z|_{C^\beta([0,T],V)}^m \right)
 \end{aligned}$$

for some constants c_4 and c_5 and m sufficiently large. Thus (3.69) is verified. To verify equality (3.70) consider the family of Yosida approximations $(F_\lambda, \lambda > 0)$ of F as in the proof of Theorem 3.6. By the dissipativity of F in the norm on V , $F_\lambda: V \rightarrow V$ is Lipschitz continuous for each $\lambda > 0$ and has at most polynomial growth, so F_λ satisfies the assumptions of Theorem 3.4 so that

$$(3.73) \quad \mathbb{E} \exp [\rho_\lambda(Z)] = 1,$$

where ρ_λ is given by (3.55). By the method used to obtain inequality (3.71), it follows that

$$\begin{aligned}
 (3.74) \quad & \mathbb{E} \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot (G_\lambda(Z) - G(Z)) \right) \right|_{L^2([0,T],V)}^2 \\
 & \leq c_6 \mathbb{E} \int_0^T \left[s^{\frac{1}{2}-H} \|G_\lambda(Z(s)) - G(Z(s))\| \right. \\
 & \quad \left. + s^{H-\frac{1}{2}} \int_0^s \frac{s^{\frac{1}{2}-H} - r^{\frac{1}{2}-H}}{(s-r)^{H+\frac{1}{2}}} \|G_\lambda(Z(r)) - G(Z(r))\| \, dr \right. \\
 & \quad \left. + \int_0^s \frac{\|G_\lambda(Z(s)) - G(Z(s)) - G_\lambda(Z(r)) + G(Z(r))\|}{(s-r)^{H+\frac{1}{2}}} \, dr \right]^2 ds.
 \end{aligned}$$

By inequalities (3.57) and (3.58), it follows that $\|G_\lambda(x) - G(x)\| \rightarrow 0$ as $\lambda \downarrow 0$ for each $x \in E$ and the family $(G_\lambda, \lambda > 0)$ satisfies the growth condition

$$(3.75) \quad \|G_\lambda(x)\| \leq c_7 (1 + \|x\|_E^\rho)$$

for $x \in E$ and some $c_7 > 0$. From the V -dissipativity of F , it follows by [5, Proposition 5.5.3] that

$$\|R_\lambda(x) - R_\lambda(y)\| \leq \|x - y\|$$

for $x, y \in E$, so that

$$\begin{aligned}
 (3.76) \quad & \|F_\lambda(x) - F_\lambda(y)\| = \|F(R_\lambda(x)) - F(R_\lambda(y))\| \\
 & \leq k_3 (1 + \|R_\lambda(x)\|_E^q + \|R_\lambda(y)\|_E^q) \|x - y\|^\gamma
 \end{aligned}$$

for $x, y \in E$. Since

$$\|R_\lambda(x)\|_E \leq \|x\|_E + \lambda \|F(x)\|_E \leq c_8 (1 + \|x\|_E^\rho)$$

for $x, y \in E$, $c_8 \in \mathbb{R}_+$, and $\lambda \in (0, 1]$, there is the inequality

$$(3.77) \quad \|F_\lambda(x) - F_\lambda(y)\| \leq c_9 (1 + \|x\|_E^m + \|y\|_E^m) \|x - y\|^\gamma$$

for $x, y \in E$, $c_0 \in \mathbb{R}_+$, $m \geq 1$, and $\lambda \in (0, 1]$. So F_λ and G_λ satisfy inequality (3.67) uniformly in $\lambda \in (0, 1]$. Thus the right-hand side of inequality (3.74) tends to zero as $\lambda \downarrow 0$ by the dominated convergence theorem where a majorizing function is provided by the estimates (3.75) and (3.77), whose integrability is shown in (3.71) and (3.72), and there is a decreasing sequence $(\lambda_n, n \in \mathbb{N})$ whose limit is zero such that

$$(3.78) \quad \lim_{n \rightarrow \infty} \exp[\rho_{\lambda_n}(Z)] = \exp[\rho(Z)] \quad \text{a.s. } \mathbb{P}.$$

The uniform integrability of the sequence $(\exp[\rho_{\lambda_n}(Z)], n \in \mathbb{N})$ is shown by verifying the analogue of (3.60). Equivalently,

$$(3.79) \quad \sup_n \mathbb{E} \int_0^T \left\| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(X_{\lambda_n}) \right) (t) \right\|^2 dt \leq c < \infty,$$

where $X_{\lambda_n}(\cdot)$ is the unique mild solution to (3.62). The analogous inequalities (3.71)–(3.74) are obtained by replacing G by G_λ using the polynomial growth bound and the local Hölder continuity that are uniform in $(\lambda_n, n \in \mathbb{N})$, and $Z(\cdot)$ is replaced by $X_{\lambda_n}(\cdot)$. For some constants c_{10} and $m \geq 1$,

$$(3.80) \quad \left| \mathbb{K}_H^{-1} \left(\int_0^\cdot G_{\lambda_n}(X_{\lambda_n}) \right) \right|_{L^2([0,T],V)}^2 \leq c_{10} \left(1 + \|X_0\|_E^m + |\tilde{X}_{\lambda_n}|_{C([0,T],E)}^m + |\tilde{X}_{\lambda_n}|_{C^\beta([0,T],V)}^m \right),$$

where $\tilde{X}_{\lambda_n}(t) = X_{\lambda_n}(t) - S(t)X_0$. By inequality (3.43), it follows that

$$(3.81) \quad |\tilde{X}_{\lambda_n}|_{C([0,T],E)} \leq c_{11} \left(1 + \|X_0\| + |Z|_{C([0,T],E)}^\rho \right)$$

for some $c_{11} > 0$. Let $w_{\lambda_n}(t) = \tilde{X}_{\lambda_n}(t) - \tilde{Z}(t)$ so that

$$(3.82) \quad w_{\lambda_n}(t) = \int_0^t S(t-r)F_{\lambda_n} \left(w_{\lambda_n}(s) + S(s)X_0 + \tilde{Z}(s) \right) ds$$

for $t \in [0, T]$. Inequality (3.81) provides a uniform bound on $|w_{\lambda_n}|_{C([0,T],E)}$, so by repeating the arguments for inequalities (3.37) and (3.38), it follows that

$$(3.83) \quad |X_{\lambda_n}|_{C^\beta([0,T],V)} \leq c_{12} \left(1 + \|X_0\| + |\tilde{Z}|_{C([0,T],E)}^\rho + |\tilde{Z}|_{C^\beta([0,T],V)} \right)$$

for some $c_{12} > 0$. Inequalities (3.80) and (3.81) verify inequality (3.79), so the sequence $(\exp[\rho_{\lambda_n}(Z)], n \in \mathbb{N})$ is uniformly integrable and equality (3.73) is verified. \square

4. Some examples. The first example is a finite-dimensional stochastic equation with a nonlinear drift. Consider the equation

$$(4.1) \quad dX(t) = f(X(t)) dt + \Phi dB(t),$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\Phi \in \mathcal{L}(\mathbb{R}^n)$, and $(B(t), t \geq 0)$ is an \mathbb{R}^n -valued standard fractional Brownian motion with Hurst parameter $H \in (0, 1)$. This case can be subsumed in the infinite-dimensional results given here, though some of the assumptions and the results simplify significantly. Let $E = V = \mathbb{R}^n$, $S(t) = I$ for $t \in \mathbb{R}_+$, and assume that $Q = \Phi\Phi^*$ is positive definite. The process

$$\left(\int_0^t \Phi dB, t \in [0, T] \right)$$

has sample paths in $C^\beta([0, T], \mathbb{R}^n)$ for $0 < \beta < H$. If $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Borel measurable and

$$(4.2) \quad \|f(x)\| \leq k_1(1 + \|x\|)$$

for some $k_1 > 0$ and all $x \in \mathbb{R}^n$, then for $H \in (0, 1/2)$ there is one and only one weak solution of (4.1) by Theorem 3.3. If, additionally, it is assumed that

$$(4.3) \quad \|f(x) - f(y)\| \leq k\|x - y\|^\gamma$$

for all $x, y \in \mathbb{R}^n$ and some $\gamma > 1 - \frac{1}{2H}$, then for $H \in (1/2, 1)$ there is one and only one weak solution. In each of these cases, the probability measure of the solution is mutually absolutely continuous with respect to the probability measure of the process $(\Phi B(t), t \in [0, T])$.

Now, replace the inequality in (4.2) by

$$(4.4) \quad \|f(x)\| \leq k_1(1 + \|x\|^\rho)$$

for some $\rho \geq 1$ and $k_1 > 0$. Assume that $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and satisfies

$$(4.5) \quad \langle f(x) - f(y), x - y \rangle \leq k_3\|x - y\|^2$$

for some $k_3 > 0$ and all $x, y \in \mathbb{R}^n$. If $H \in (1/2, 1)$, then assume that

$$(4.6) \quad \|f(x) - f(y)\| \leq k_4(1 + \|x\|^q + \|y\|^q)\|x - y\|^\gamma$$

for some $q \geq 1$, $k_4 > 0$, $\gamma > 1 - \frac{1}{2H}$. For $H \in (0, \frac{1}{2})$ Theorem 3.6 can be used to verify that the probability law of the solution of (4.1) is mutually absolutely continuous with respect to the probability law of $(\Phi B(t), t \in [0, T])$. Furthermore, there is one and only one mild solution of (4.1); in fact, since the state space is finite-dimensional, the mild solution is a strong solution. For $H \in (\frac{1}{2}, 1)$ Theorem 3.7 can be used to verify mutual absolute continuity and one and only one mild solution as for the case $H \in (0, \frac{1}{2})$. Note that inequalities (4.4)–(4.6) are satisfied for the important case of models where f is a polynomial of odd degree with a negative leading coefficient.

The second example is a stochastic parabolic equation of $2m$ th order:

$$(4.7) \quad \frac{\partial u}{\partial t}(t, \xi) = [L_{2m}u](t, \xi) + f(u(t, \xi)) + \eta(t, \xi)$$

for $(t, \xi) \in [0, T] \times \mathcal{O}$ with the initial condition

$$(4.8) \quad u(0, \xi) = x(\xi)$$

for $\xi \in \mathcal{O}$ and the Dirichlet boundary condition

$$(4.9) \quad \frac{\partial^k u}{\partial \nu^k}(t, \xi) = 0$$

for $(t, \xi) \in [0, T] \times \partial\mathcal{O}$, $k \in \{0, \dots, m - 1\}$, with $\frac{\partial}{\partial v}$ denoting the conormal derivative, \mathcal{O} a bounded domain in \mathbb{R}^d with a smooth boundary, and L_{2m} a $2m$ th order uniformly elliptic operator

$$(4.10) \quad L_{2m} = \sum_{|\alpha| \leq 2m} a_\alpha(\xi) D^\alpha$$

with $a_\alpha \in C_b^\infty(\mathcal{O})$. For example, if $m = 1$, then this equation is called the stochastic heat equation. The process η denotes a space-dependent noise process that is fractional in time with the Hurst parameter $H \in (0, 1)$ and, possibly, in space. The system (4.7)–(4.9) is modeled as

$$(4.11) \quad \begin{aligned} dX(t) &= AX(t)dt + F(X(t))dt + \Phi dB(t), \\ X(0) &= x \end{aligned}$$

in the space $V = L^2(\mathcal{O})$, where $A = L_{2m}$,

$$\text{Dom}(A) = \left\{ \varphi \in H^{2m}(\mathcal{O}) \mid \frac{\partial^k}{\partial v^k} \varphi = 0 \text{ on } \partial D \text{ for } k \in \{0, \dots, m - 1\} \right\},$$

$F : V \rightarrow V$ is the operator, $F(x)(\xi) = f(x(\xi))$, $x \in V, \xi \in \mathcal{O}$, $\Phi \in \mathcal{L}(V)$ defines the space correlation of the noise process, and $(B(t), t \geq 0)$ is a cylindrical standard fractional Brownian motion in V (formally, $\eta(t, \cdot) = \Phi(\partial/\partial t)B(t, \cdot)$). For $\Phi = I$, the noise process is uncorrelated in space. It is well known that A generates an analytic semigroup $(S(t), t \geq 0)$. Furthermore

$$(4.12) \quad |S(t)\Phi|_{\mathcal{L}_2(V)} \leq |S(t)|_{\mathcal{L}_2(V)} |\Phi|_{\mathcal{L}(V)} \leq ct^{-\frac{d}{4m}}$$

for $t \in [0, T]$, so if

$$(4.13) \quad H > \frac{d}{4m},$$

then the conditions of Proposition 2.6 are satisfied with $\gamma = \frac{d}{4m}$. Therefore, for any $\Phi \in \mathcal{L}(V)$, the stochastic convolution process

$$\left(\int_0^t S(t-r)\Phi dB(r), t \in [0, T] \right)$$

is well-defined and has a version with $C^\beta([0, T], V)$ sample paths for $\beta \geq 0$ satisfying

$$(4.14) \quad \beta < H - \frac{d}{4m}.$$

Note that the condition (4.13) extends the well-known result for a standard Wiener process ($H = \frac{1}{2}$).

Theorems 3.3 and 3.4 are applied to the present example. Assume inequality (4.13) and let Φ be boundedly invertible on V . Furthermore, let $f : \mathbb{R} \rightarrow \mathbb{R}$ be measurable and satisfy

$$(4.15) \quad |f(\xi)| \leq k_1(1 + |\xi|), \quad \xi \in \mathbb{R}.$$

By the preceding part of this example, conditions (H1)–(H3) are satisfied for $E = V = L^2(\mathcal{O})$ and the map $F : V \rightarrow V$ has at most linear growth. Thus by Theorem 3.3 if $H < \frac{1}{2}$, then there exists a unique weak solution to (4.11).

If $H > \frac{1}{2}$, some additional conditions are required. Assume that

$$(4.16) \quad \frac{d}{4m} < \frac{1}{2}$$

(which is more restrictive than (4.13)) and suppose that

$$(4.17) \quad |f(\xi) - f(\lambda)| \leq k|\xi - \lambda|^\gamma, \quad \xi, \lambda \in \mathbb{R},$$

for some $k > 0$ and $\gamma > 0$,

$$(4.18) \quad \frac{H - 1/2}{H - d/4m} < \gamma \leq 1.$$

Then, letting β be such that $\beta < H - \frac{d}{4m}$ and $\gamma\beta > H - \frac{1}{2}$, it is clear that all of the conditions of Theorem 3.4 are verified so there is a unique weak solution to (4.11).

The third example is a one-dimensional stochastic equation of reaction-diffusion type. Consider the equation

$$(4.19) \quad \frac{\partial u}{\partial t}(t, \xi) = \frac{\partial^2 u}{\partial \xi^2}(t, \xi) + f(u(t, \xi)) + \eta(t, \xi)$$

for $(t, \xi) \in (0, T) \times (0, 1)$ and

$$\begin{aligned} u(0, \xi) &= x_0(\xi), \\ \frac{\partial u}{\partial \xi}(t, 0) &= \frac{\partial u}{\partial \xi}(t, 1) = 0 \end{aligned}$$

for $(t, \xi) \in (0, T) \times (0, 1)$, where f and η are given in the previous example (with $\mathcal{O} = (0, 1)$). The above formal equation can be rewritten in the form (4.11) with $V = L^2([0, 1])$, $A = \frac{\partial^2}{\partial \xi^2}$,

$$\text{Dom}(A) = \left\{ \phi \in H^2([0, 1]) : \frac{\partial}{\partial \xi} \phi(0) = \frac{\partial}{\partial \xi} \phi(1) = 0 \right\},$$

where $\Phi \in \mathcal{L}(V)$ and F is as given in the preceding example. The semigroup generated by A satisfies the estimate (4.12) (with $m = d = 1$), so if f satisfies the conditions of the previous example, the same conclusions on existence and uniqueness of the weak solution are obtained.

However, it is desirable to relax condition (4.15) of the linear growth of the function f , which is very restrictive in view of reaction-diffusion models, where f is often a polynomial. Let $H > 1/2$ and assume that

$$(4.20) \quad |f(\xi)| \leq k(1 + |\xi|^\rho),$$

$$(4.21) \quad (f(\xi) - f(\lambda)) \text{sgn}(\xi - \lambda) \leq k(\xi - \lambda),$$

$$(4.22) \quad |f(\xi) - f(\lambda)| \leq k(1 + |\xi|^q + |\lambda|^q)|\xi - \lambda|^\gamma$$

for all $\xi, \lambda \in \mathbb{R}$ and some universal constants $\rho > 0, q > 0, k > 0$ and γ satisfying

$$(4.23) \quad \frac{H - 1/2}{H - 1/4} < \gamma \leq 1.$$

Note that these conditions are satisfied if f is Lipschitz or if f is a polynomial of odd degree with a negative leading coefficient.

The conditions of Theorem 3.7 are verified now. Take the state space $E = C([0, 1])$. It is well known that the restriction of A to E generates a strongly continuous semigroup of contractions on E . By Proposition 2.6 the stochastic convolution

$$(4.24) \quad \left(\int_0^t S(t-r) \Phi dB(r), t \in [0, T] \right)$$

has $C^\beta([0, T], V_\delta)$ sample paths for $\beta + \delta < H - 1/4$, and, hence, by the Sobolev embedding theorem, in the space $C([0, T], E) \cap C^\beta([0, T], V)$ for $0 < \beta < H - 1/4$ (by (4.23) β can be chosen such that $\beta\gamma > H - 1/2$). It remains to verify the conditions imposed on F . The polynomial growth condition (3.41), the “dissipativity of $F - kI$ on V ” (3.66), and the local Hölder continuity of the form (3.67) follow easily from the corresponding conditions on f , that is, (4.20), (4.21), and (4.22). The dissipativity of $F - kI$ on E (3.42) is a well-known consequence of (4.21) by the characterization of the subdifferential of the norm on $E = C([0, 1])$ (cf. [32]). The characterization of the subdifferentials for this example is as follows: Given $x \in E$, let $M_x = \{\xi \in [0, 1]; |x(\xi)| = \|x\|_E\}$. Then $\mu \in \delta\|x\|_E$ if and only if the following three conditions are satisfied: (i) μ is a Radon measure on $[0, 1]$ with $\|\mu\| = 1$; (ii) the support of μ is contained in M_x ; and (iii) $\int_\Gamma \operatorname{sgn} x(\xi) \mu(d\xi) \geq 0$ for each Borel set Γ in $[0, 1]$. In particular if $x \in E$ has the property that $M_x = \{\xi_0\}$, then $\delta\|x\|_E = \delta_{\xi_0}$ for $x(\xi_0) = \|x\|_E$ and $\delta\|x\|_E = -\delta_{\xi_0}$ for $x(\xi_0) = -\|x\|_E$, where δ_{ξ_0} is the Dirac distribution at ξ_0 . The family of $x \in E$ with this latter property is dense in E . Therefore, all of the conditions of Theorem 3.7 are satisfied, and it follows that there is a unique weak solution in the present case.

REFERENCES

- [1] E. ALÒS, O. MAZET, AND D. NUALART, *Stochastic calculus with respect to Gaussian processes*, Ann. Probab., 29 (2001), pp. 766–801.
- [2] P. CAITHAMER, *The stochastic wave equation driven by fractional Brownian noise and temporally correlated smooth noise*, Stoch. Dyn., 5 (2005), pp. 45–64.
- [3] R. H. CAMERON AND W. T. MARTIN, *Transformations of Wiener integrals under translations*, Ann. of Math. (2), 2 (1944), pp. 386–396.
- [4] G. DA PRATO AND J. ZABCZYK, *Stochastic Equations in Infinite Dimensions*, Encyclopedia Math. Appl. 44, Cambridge University Press, Cambridge, UK, 1992.
- [5] G. DA PRATO AND J. ZABCZYK, *Ergodicity for Infinite-Dimensional Systems*, London Math. Soc. Lecture Note Ser. 229, Cambridge University Press, Cambridge, UK, 1996.
- [6] L. DECREUSEFOND AND A. S. ÜSTÜNEL, *Stochastic analysis of the fractional Brownian motion*, Potential Anal., 10 (1999), pp. 177–214.
- [7] L. DENIS, M. ERRAOUI, AND Y. OUKNINE, *Existence and uniqueness for solutions of one dimensional SDE's driven by an additive fractional noise*, Stoch. Stoch. Rep., 76 (2004), pp. 409–427.
- [8] T. E. DUNCAN, *Some stochastic semilinear equations in Hilbert space with fractional Brownian motion*, in Optimal Control and Partial Differential Equations, J. Menaldi, E. Rofman, and A. Sulem, eds., IOS Press, Amsterdam, 2000, pp. 241–247.
- [9] T. E. DUNCAN, *Some processes associated with a fractional Brownian motion*, in Mathematics of Finance, Contemp. Math. 351, AMS, Providence, RI, 2004, pp. 93–101.
- [10] T. E. DUNCAN, Y. HU, AND B. PASIK-DUNCAN, *Stochastic calculus for fractional Brownian motion I. Theory*, SIAM J. Control Optim., 38 (2000), pp. 582–612.

- [11] T. E. DUNCAN, J. JAKUBOWSKI, AND B. PASIK-DUNCAN, *Stochastic integration for fractional Brownian motion in a Hilbert space*, Stoch. Dyn., 6 (2006), pp. 53–75.
- [12] T. E. DUNCAN, B. MASLOWSKI, AND B. PASIK-DUNCAN, *Stochastic equations in Hilbert space with a multiplicative fractional Gaussian noise*, Stochastic Process. Appl., 115 (2005), pp. 1357–1383.
- [13] T. E. DUNCAN, B. PASIK-DUNCAN, AND B. MASLOWSKI, *Fractional Brownian motion and stochastic equations in Hilbert spaces*, Stoch. Dyn., 2 (2002), pp. 225–250.
- [14] I. V. GIRSANOV, *On transforming a class of stochastic processes by absolutely continuous substitution of measures*, Teor. Veroyatnost. i Primenen., 5 (1960), pp. 314–330.
- [15] W. GRECKSCH AND V. V. ANH, *A parabolic stochastic differential equation with fractional Brownian motion input*, Statist. Probab. Lett., 41 (1999), pp. 337–346.
- [16] Y. HU, *Heat equations with fractional white noise potentials*, Appl. Math. Optim., 43 (2001), pp. 221–243.
- [17] Y. HU, *Integral transformations and anticipative calculus for fractional Brownian motions*, Mem. Amer. Math. Soc., 175 (2005).
- [18] Y. HU, B. ØKSENDAL, AND T. ZHANG, *General fractional multiparameter white noise theory and stochastic partial differential equations*, Comm. Partial Differential Equations, 29 (2004), pp. 1–23.
- [19] I. KARATZAS AND S. E. SHREVE, *Brownian Motion and Stochastic Calculus*, Grad. Texts in Math. 113, Springer-Verlag, New York, 1988.
- [20] R. S. LIPTSER AND A. N. SHIRYAEV, *Statistics of Random Processes I: General Theory*, 2nd ed., Springer-Verlag, Berlin, 2001.
- [21] B. MASLOWSKI AND D. NUALART, *Evolution equations driven by a fractional Brownian motion*, J. Funct. Anal., 202 (2003), pp. 277–305.
- [22] B. MASLOWSKI AND B. SCHMALFUSS, *Random dynamical systems and stationary solutions of differential equations driven by the fractional Brownian motion*, Stochastic Anal. Appl., 22 (2004), pp. 1577–1607.
- [23] R. H. MARTIN, JR., *A global existence theorem for autonomous differential equations in a Banach space*, Proc. Amer. Math. Soc., 26 (1970), pp. 307–314.
- [24] R. MIKULEVICIUS AND B. L. ROZOVSKY, *Martingale paths for stochastic PDE's*, in Stochastic Partial Differential Equations: Six Perspectives, Math. Survey Monogr. 64, AMS, Providence, RI, 1999, pp. 243–325.
- [25] S. MORET AND D. NUALART, *Onsager-Machlup functional for the fractional Brownian motion*, Probab. Theory Related Fields, 124 (2002), pp. 227–260.
- [26] I. NORROS, E. VALKEILA, AND J. VIRTAMO, *An elementary approach to a Girsanov formula and other analytical results on fractional Brownian motions*, Bernoulli, 5 (1999), pp. 571–587.
- [27] D. NUALART AND Y. OUKNINE, *Regularization of differential equations by fractional noises*, Stoch. Process. Appl., 102, (2002) pp. 103–116.
- [28] B. PASIK-DUNCAN, T. E. DUNCAN, AND B. MASLOWSKI, *Linear stochastic equations in a Hilbert space with a fractional Brownian motion*, in Stochastic Processes, Optimization, and Control Theory: Applications in Financial Engineering, Queueing Networks, and Manufacturing Systems, H. Yan, G. Yin, and Q. Zhang, eds., Springer-Verlag, New York, 2006, pp. 201–221.
- [29] A. PAZY, *Semigroups of Linear Operators and Applications to Partial Differential Equations*, Appl. Math. Sci. 44, Springer-Verlag, New York, 1983.
- [30] V. PIPIRAS AND M. S. TAQQU, *Integration questions related to fractional Brownian motion*, Probab. Theory Related Fields, 118 (2000), pp. 251–291.
- [31] S. G. SAMKO, A. A. KILBAS, AND O. I. MARICHEV, *Fractional Integrals and Derivatives*, Gordon and Breach, Yverdon, 1993.
- [32] E. SINISTRARI, *Accretive differential operators*, Boll. Un. Mat. Ital. B (5), 13 (1976), pp. 19–31.
- [33] S. TINDEL, C. A. TUDOR, AND F. VIENS, *Stochastic evolution equations with fractional Brownian motion*, Probab. Theory Related Fields, 127 (2003), pp. 186–204.

ON THE EXISTENCE OF SOLUTIONS OF EQUILIBRIA IN LUBRICATED JOURNAL BEARINGS*

I. CIUPERCA[†], M. JAI[‡], AND J. I. TELLO[§]

Abstract. In this paper we study a system of equations concerning equilibrium positions of journal bearings. The problem consists of two surfaces in relative motion separated by a small distance filled with a lubricant. The shape of the inlet surface is circular, while the other surface has a more general shape. Our result shows the existence of at least one equilibrium by using degree theory.

Key words. Reynolds variational inequality, inverse problem, existence of solutions, degree theory

AMS subject classifications. 35J20, 47H11, 49J10

DOI. 10.1137/080724228

1. Introduction. We consider in this paper a lubricated system called a journal bearing consisting of two cylinders in relative motion. An incompressible fluid, the lubricant, is introduced in the narrow space between the cylinders. An exterior force $F = (F_1, F_2) \in \mathbb{R}^2$ is applied on the inner cylinder (shaft) which turns with a given velocity ω .

The wedge between the two cylinders is assumed to satisfy the thin-film hypothesis, so that the pressure (assumed time-independent) does not depend on the normal coordinate to the bodies and obeys the Reynolds equation.

In order to introduce the Reynolds equation we need to describe the geometry and the dynamic of the system. We suppose that the interior cylinder has a circular form of constant radius (assumed 1) which rotates with known velocity. The transversal axis of the shaft is assumed to have only two degrees of freedom in the transversal plane, i.e., parallel to the exterior cylinder (bush) which is fixed and not necessarily of constant radius.

Let us consider (O, y_1, y_2) a reference system in the transversal plane to the cylinders, and suppose that the distance between O and the surface of the bush is larger than the radius of the shaft (see Figure 1.1). We also assume that the representation in polar coordinates (r, θ) of the bush is given by

$$(1.1) \quad r = 1 + \delta\rho(\theta),$$

where $\rho : [0, 2\pi] \rightarrow [1, +\infty[$ is a known function and $\delta > 0$ (the clearance) is a small parameter, representing the distance between the two cylinders when O and the center of the shaft coincide.

*Received by the editors May 14, 2008; accepted for publication (in revised form) November 12, 2008; published electronically February 20, 2009.

<http://www.siam.org/journals/sima/40-6/72422.html>

[†]Université de Lyon, Université Lyon 1, CNRS, UMR 5208, Institut Camille Jordan, Bat. Braconnier, 43, blvd du 11 novembre 1918, F-69622 Villeurbanne Cedex, France (ciuperca@math.univ-lyon1.fr).

[‡]ICJ CNRS-UMR 5208, INSA de Lyon, Bât. Léonard de Vinci, 20 A. A. Einstein, 69621 Villeurbanne Cedex, France (Mohammed.Jai@insa-lyon.fr).

[§]Matemática Aplicada, E.U.I. Informática, Universidad Politécnica de Madrid, 28031 Madrid, Spain (jtello@eui.upm.es). This author was partially supported by project MTM2005-03463 of DG-ISGPI of Spain and projects of CAM: CCG07-UPM/000-3199 at UPM and CCG07-UCM/ESP-2787 at UCM.

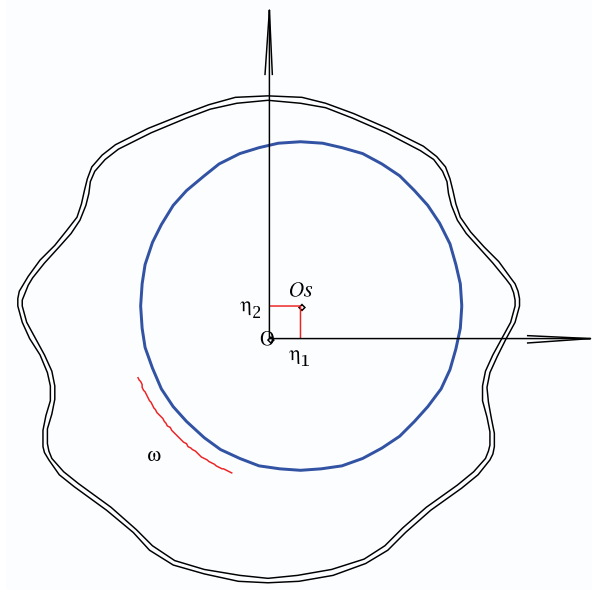


FIG. 1.1. Scheme of the journal bearing.

Remark 1.1. The particular case $\rho \equiv 1$ corresponds to a circular bush with radius $1 + \delta$ and O the center of the bush.

Let us now denote by O_s the center of the shaft. The position of O_s is given in cartesian coordinates by $(\delta\eta_1, \delta\eta_2)$ and in polar coordinates by $(\delta\eta, \alpha)$, that is,

$$(1.2) \quad \begin{aligned} \eta_1 &= \eta \cos \alpha, \\ \eta_2 &= \eta \sin \alpha. \end{aligned}$$

It is well known (see, for instance, [5]) that the distance between the two cylinders is given by $\delta h(\theta, \eta, \alpha) + O(\delta^2)$ with

$$(1.3) \quad h(\theta, \eta, \alpha) = \rho(\theta) - \eta \cos(\theta - \alpha) = \rho(\theta) - \eta_1 \cos \theta - \eta_2 \sin \theta.$$

The formulation of the problem is complete when the distance between the surfaces is of order δ (i.e., $h \in O(1)$) and second order terms are neglected. Admissible forces in that case are of order $\frac{1}{\delta^2}$ or smaller. If the distance becomes smaller (i.e., $h \ll 1$) or forces are large (larger than $O(\frac{1}{\delta^2})$), second order terms cannot be neglected and the formulation loses its physical meaning.

Now, the problem will be posed in a fixed domain $\Omega =]0, 2\pi[\times]0, 1[$ which parametrizes the space between shaft and bush. In fact the gap is approximated by the following domain given in cylindrical coordinates by (r, θ, x) with respect to O_s :

$$1 \leq r \leq 1 + \delta h(\theta, \eta, \alpha), \quad \theta \in [0, 2\pi[, \quad x \in [0, 1],$$

where $O_s = (\delta\eta \cos(\alpha), \delta\eta \sin(\alpha))$.

Since the wedge between the two cylinders satisfies the thin-film hypothesis, the pressure of the lubricant fluid (assumed time-independent) does not depend on the

normal coordinate to the bodies and obeys the Reynolds equation (see [9]). We also consider that there is an alimentation region along the circles $\{x = 0\}$ and $\{x = 1\}$, respectively, where the pressure in the fluid equals the atmospheric pressure supposed to be 0 by translation. Then the pressure $p : (\theta, x) \in \Omega \rightarrow \mathbb{R}$ satisfies the following problem written in nondimensional form:

$$(1.4) \quad \begin{cases} \nabla \cdot (h^3 \nabla p) = \frac{\partial h}{\partial \theta} & \text{on } \Omega, \\ p & \text{is } 2\pi\text{-periodic in } \theta, \\ p = 0 & \text{on }]0, 2\pi[\times \{0\} \cup]0, 2\pi[\times \{1\}. \end{cases}$$

In general the solution of (1.4) is not always nonnegative and we must replace (1.4) by the corresponding variational inequality.

In this work we are interested in an equilibrium problem which entails finding the position (η_1, η_2) of the shaft such that the hydrodynamic force (load) created by the pressure film equilibrates the exterior force F . Thus the problem is formulated as follows:

Find $p \in K, (\eta_1, \eta_2) \in A$ such that

$$(1.5) \quad \int_{\Omega} h^3 \nabla p \cdot \nabla (\varphi - p) \geq \int_{\Omega} h \frac{\partial}{\partial \theta} (\varphi - p) \quad \forall \varphi \in K,$$

$$(1.6) \quad \int_{\Omega} p \cos \theta d\theta dx = F_1,$$

$$(1.7) \quad \int_{\Omega} p \sin \theta d\theta dx = F_2,$$

where h is given in (1.3),

$$K = \{\varphi \in H_0^1(\Omega) : \varphi \geq 0\},$$

and $A \subset \mathbb{R}^2$. The set of admissible positions of the shaft (to be defined later) is such that $h(\theta, \eta, \alpha) > 0 \quad \forall \theta \in [0, 2\pi]$.

As far as we know, very few works can be found in the literature concerning existence of equilibrium in the lubricated devices in spite of a larger number of references to numerical simulations; see, for instance, [3] or [4].

Some results exist in the case of sliders, that is, mechanisms consisting of an almost plane surface sliding above a horizontal plane surface.

Exact solutions for the case in which the upper surface is an inclined plane of angle θ and infinite width have been known for a long time, since the problem becomes one-dimensional and is easily integrated [7].

For more general shapes of the upper surfaces an existence result is obtained in [1] for the equation case and in [2] for the variational inequality in the one-dimensional case.

The content of the paper is as follows: In section 2 we give the main result. In section 3 we recall some elements of the degree theory which will be used. In section 4 some preliminary results are given, and section 5 is devoted to the proof of the main result.

2. Main results and assumptions. For the rest of the paper we assume that

$$(2.1) \quad \rho \in C^3(\mathbb{R}), \quad \rho, \rho', \text{ and } \rho'' \text{ are } 2\pi\text{-periodic,}$$

$$(2.2) \quad \rho''(\theta) + \rho(\theta) > 0 \quad \forall \theta \in \mathbb{R},$$

$$(2.3) \quad \min_{0 \leq \theta \leq 2\pi} \rho(\theta) = 1.$$

We denote the following:

$$(2.4) \quad \begin{aligned} \rho_M &= \max_{0 \leq \theta \leq 2\pi} \rho(\theta), \\ m &= \min_{0 \leq \theta \leq 2\pi} (\rho''(\theta) + \rho(\theta)) > 0, \\ M &= \max_{0 \leq \theta \leq 2\pi} (\rho''(\theta) + \rho(\theta)) > 0. \end{aligned}$$

Remark 2.1. All these assumptions are clearly satisfied in the particular case of a circular bush, corresponding to $\rho \equiv 1$.

Remark 2.2. Assumption (2.2) is the most restrictive of the assumptions. It is introduced for technical reasons and guarantees that in the limit case the set where $h = 0$ is a single line.

Let us introduce the function $a(\alpha) : \mathbb{R} \rightarrow \mathbb{R}_+$ given by

$$(2.5) \quad a(\alpha) := \min_{\alpha - \frac{\pi}{2} < \theta < \alpha + \frac{\pi}{2}} \left\{ \frac{\rho(\theta)}{\cos(\theta - \alpha)} \right\}.$$

The fact that

$$\lim_{\theta \rightarrow \alpha \pm \frac{\pi}{2}} \frac{\rho(\theta)}{\cos(\theta - \alpha)} = \infty$$

and the continuity and boundedness of ρ guarantee the existence of at least one minimum.

Now we define the set A by

$$A = \left\{ (\eta_1, \eta_2) \in \mathbb{R}^2 : 0 \leq \eta < a(\alpha) \right\},$$

where (η, α) are given by (1.2).

It is clear that $h(\theta) > 0 \forall \theta \in [0, 2\pi]$ if and only if $(\eta_1, \eta_2) \in A$.

For any fixed $(\eta_1, \eta_2) \in A$, problem (1.5) has been studied by several authors. The existence and uniqueness of solutions can be obtained by using the direct methods in the calculus of variations along with strict convexity of the associated functional. Then (1.5) admits a unique classical solution $p \in K$; see, for instance, Kinderlehrer and Stampacchia [6].

Thus problem (1.5)–(1.7) is equivalent to the following problem:

$$(2.6) \quad \begin{cases} \text{Find } (\eta_1, \eta_2) \in A & \text{such that} \\ G(\eta_1, \eta_2) = (0, 0), \end{cases}$$

where $G = (G_1, G_2) : A \rightarrow \mathbb{R}^2$ is given by

$$(2.7) \quad G_1(\eta_1, \eta_2) = \int_{\Omega} p \cos \theta d\theta dx - F_1,$$

$$(2.8) \quad G_2(\eta_1, \eta_2) = \int_{\Omega} p \sin \theta d\theta dx - F_2,$$

where p (depending on η_1, η_2) is the unique solution of (1.5).

The main result of the paper is presented in the following theorem.

THEOREM 2.1. *Under assumptions (2.1)–(2.3) and for any $F \in \mathbb{R}^2$ there exists at least one solution $(\eta_1, \eta_2) \in A$ of (2.6).*

3. Known results on degree theory. In order to prove Theorem 2.1 we use the topological degree theory, which is rapidly recalled in the following.

The topological degree for continuous mappings between n -dimensional Euclidean spaces was first introduced by L. E. J. Brouwer in 1912. We first introduce the definition of degree for C^1 maps.

Let S be a bounded open subset of \mathbb{R}^n and a $C^1(S)$ function $f : S \rightarrow \mathbb{R}^n$. Let $y_0 \in \mathbb{R}^n$ such that $y_0 \notin f(\partial S)$, and suppose that $f \in C^1(S)$ and that $Df(x)$ is invertible for all $x \in f^{-1}(y_0)$. Then $f(x) = y_0$ has either no solutions in S or a finite number r of solutions, say x_1, x_2, \dots, x_r and $\det(Df(x_i)) \neq 0$ for $i = 1, 2, \dots, r$.

DEFINITION 3.1. *We define the degree of f in S at y_0 as follows:*

If $r = 0$, then $d(f, S, y_0) := 0$; else

$$d(f, S, y_0) := \sum_{i=1}^r \text{sign}(\det(Df(x_i))).$$

Remark 3.1. In the particular case where f is a linear function, i.e., $f(x) = Ax$, where A is an invertible matrix, the general formula for calculating the degree of linear functions is the following:

$$\text{deg}(Ax, S, 0) = \text{sgn}(\det A) \quad \text{if } 0 \text{ is contained in } S.$$

The definition of degree can be extended to continuous maps; see [8, Extension Lemma, p. 60]. For the reader's convenience we state the following result (see, for instance, [8, Corollary 4 to Theorem 1.12, p. 81]) adapted to the finite-dimensional case.

THEOREM 3.1. *Let S be a bounded open set in \mathbb{R}^n . Let f_0 and f_1 be two continuous functions from \bar{S} to \mathbb{R}^n . We assume moreover that S is a star domain with respect to the point $y_0 \in S$ and that*

$$[f_0(x) - y_0] \cdot [f_1(x) - y_0] > 0 \text{ for any } x \in \partial S.$$

Then

$$d(f_0, S, y_0) = d(f_1, S, y_0).$$

Remark 3.2. It is clear that if $d(f, S, y_0) \neq 0$, then there exists at least a solution $x \in S$ of the equation $f(x) = y_0$.

4. Preliminary results.

LEMMA 4.1. *Let $a(\alpha)$ be the function defined in (2.5). Then*

$$1 = \min_{0 \leq \theta \leq 2\pi} \rho \leq a(\alpha) \leq \max_{0 \leq \theta \leq 2\pi} \rho = \rho_M.$$

Proof. Since

$$\min_{0 \leq \theta \leq 2\pi} \rho \leq \rho \leq \max_{0 \leq \theta \leq 2\pi} \rho$$

we have that

$$\min_{0 \leq \theta \leq 2\pi} \{\rho\} \min_{\alpha - \frac{\pi}{2} < \theta < \alpha + \frac{\pi}{2}} \left\{ \frac{1}{\cos(\theta - \alpha)} \right\} \leq a(\alpha)$$

and

$$a(\alpha) \leq \max_{0 \leq \theta \leq 2\pi} \{\rho\} \min_{\alpha - \frac{\pi}{2} < \theta < \alpha + \frac{\pi}{2}} \left\{ \frac{1}{\cos(\theta - \alpha)} \right\}.$$

The fact that

$$\min_{\alpha - \frac{\pi}{2} < \theta < \alpha + \frac{\pi}{2}} \left\{ \frac{1}{\cos(\theta - \alpha)} \right\} = 1$$

implies

$$\min\{\rho\} \leq a(\alpha) \leq \max\{\rho\},$$

and the proof is complete. \square

LEMMA 4.2. *There exists s with $0 < s < \frac{\pi}{2}$ such that*

$$|\tilde{\theta}_\alpha - \alpha| < s$$

for any $\alpha \in \mathbb{R}$ and for any $\tilde{\theta}_\alpha \in]\alpha - \frac{\pi}{2}, \alpha + \frac{\pi}{2}[$ satisfying

$$\frac{\rho(\tilde{\theta}_\alpha)}{\cos(\tilde{\theta}_\alpha - \alpha)} = a(\alpha).$$

Proof. From Lemma 4.1 we deduce

$$\cos(\tilde{\theta}_\alpha - \alpha) \geq \frac{1}{\rho_M},$$

which implies

$$(4.1) \quad |\tilde{\theta}_\alpha - \alpha| \leq \arccos \frac{1}{\rho_M} < \frac{\pi}{2}.$$

Now taking $s = \frac{1}{2}(\frac{\pi}{2} + \arccos \frac{1}{\rho_M})$ we obtain the result. \square

LEMMA 4.3. *Under assumptions (2.1)–(2.3), for any $\alpha \in \mathbb{R}$, the set of $\theta \in]\alpha - \frac{\pi}{2}, \alpha + \frac{\pi}{2}[$, which minimizes the function $\theta \mapsto \frac{\rho(\theta)}{\cos(\theta - \alpha)}$, is a single point which we denote by θ_α satisfying*

$$\frac{\rho(\theta_\alpha)}{\cos(\theta_\alpha - \alpha)} = a(\alpha)$$

and

$$\rho'(\theta_\alpha) \cos(\theta_\alpha - \alpha) + \rho(\theta_\alpha) \sin(\theta_\alpha - \alpha) = 0.$$

Moreover the function $\alpha \mapsto \theta_\alpha$ belongs to $C^1(\mathbb{R})$.

Proof. Let α be a fixed value of \mathbb{R} and consider the critical points of the function $\theta \in]\alpha - \frac{\pi}{2}, \alpha + \frac{\pi}{2}[\mapsto \frac{\rho(\theta)}{\cos(\theta - \alpha)}$, which satisfies

$$(4.2) \quad \left(\frac{\rho(\theta)}{\cos(\theta - \alpha)} \right)' = \frac{\rho'(\theta)}{\cos(\theta - \alpha)} + \frac{\rho(\theta) \sin(\theta - \alpha)}{\cos^2(\theta - \alpha)} = 0.$$

Consider

$$f(\theta, \alpha) := \rho' \cos(\theta - \alpha) + \rho \sin(\theta - \alpha).$$

Notice that if θ is a minimum of

$$\frac{\rho(\theta)}{\cos(\theta - \alpha)},$$

then $f(\theta, \alpha) = 0$. We now have that

$$\frac{\partial f}{\partial \theta} = (\rho'' + \rho) \cos(\theta - \alpha),$$

which is positive in $|\theta - \alpha| < \frac{\pi}{2}$ (by hypothesis (2.2)). The assertion follows using also the implicit function theorem on the open set $(\theta, \alpha) \in]-s, s[\times \mathbb{R}$ with s given in Lemma 4.2. \square

We now introduce the function $h_0 : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$h_0(\theta, \alpha) = h(\theta, a(\alpha), \alpha) = \rho(\theta) - a(\alpha) \cos(\theta - \alpha).$$

Remark 4.1. We have from Lemma 4.3

$$(4.3) \quad h_0(\theta_\alpha, \alpha) = \frac{\partial h_0}{\partial \theta}(\theta_\alpha, \alpha) = 0,$$

$$(4.4) \quad \frac{\partial^2 h_0}{\partial \theta^2}(\theta_\alpha, \alpha) = \rho''(\theta_\alpha) + \rho(\theta_\alpha).$$

LEMMA 4.4. *Let us denote for any $\alpha \in [0, 2\pi]$ and $\epsilon > 0$ small enough*

$$I_{\alpha, \epsilon} := [\theta_\alpha - 2\sqrt{\epsilon}, \theta_\alpha + \sqrt{\epsilon}].$$

Then for any $\alpha \in [0, 2\pi]$ and $\theta \in I_{\alpha, \epsilon}$ we have

(i) $h(\theta, a(\alpha) - \epsilon, \alpha) \leq (2M + 2)\epsilon,$

(ii) $-\frac{\partial h}{\partial \theta}(\theta, a(\alpha) - \epsilon, \alpha) \geq \frac{m}{2}\sqrt{\epsilon},$

with m and M defined in (2.4).

Proof.

(i) We have

$$h(\theta, a(\alpha) - \epsilon, \alpha) = h_0(\theta, \alpha) + \epsilon \cos(\theta - \alpha).$$

From Taylor development of h_0 at $\theta = \theta_\alpha$, using (4.3) and (4.4) we have (i).

(ii) We have

$$-\frac{\partial h}{\partial \theta}(\theta, a(\alpha) - \epsilon, \alpha) = -\frac{\partial h_0}{\partial \theta}(\theta, \alpha) + \epsilon \sin(\theta - \alpha).$$

Using the Taylor development of $\frac{\partial h_0}{\partial \theta}$ at $\theta = \theta_\alpha$ we obtain

$$-\frac{\partial h}{\partial \theta}(\theta, a(\alpha) - \epsilon, \alpha) = \frac{\partial^2 h_0}{\partial \theta^2}(\hat{\theta}, \alpha)(\theta_\alpha - \theta) + \epsilon \sin(\theta - \alpha)$$

with $\hat{\theta} \in I_{\alpha, \epsilon}$.

From (4.4), the uniform continuity of $\frac{\partial^2 h_0}{\partial \theta^2}$, and (2.4) we deduce

$$\frac{\partial^2 h_0}{\partial \theta^2}(\hat{\theta}, \alpha) \geq \frac{2}{3}m.$$

Since $\theta_\alpha - \theta \geq \sqrt{\epsilon}$, we obtain for ϵ small enough

$$-\frac{\partial h}{\partial \theta}(\theta, a(\alpha) - \epsilon, \alpha) \geq \frac{2}{3}m\sqrt{\epsilon} - \epsilon,$$

which ends the proof. \square

LEMMA 4.5. *Let ϵ be small enough and $\eta := a(\alpha) - \epsilon$. There exists a constant $c > 0$ independent of ϵ and α such that*

$$\inf_{0 \leq \alpha \leq 2\pi} \int_{\Omega} h^3(\theta, \eta, \alpha) |\nabla p|^2 d\theta dx > c\epsilon^{-\frac{1}{2}} \quad \text{as } \epsilon \rightarrow 0.$$

Proof. We take in (1.5) $\phi := p + \varphi$ with $\varphi \in K$ arbitrary. Then

$$\int_{\Omega} h \frac{\partial \varphi}{\partial \theta} d\theta dx \leq \int_{\Omega} h^3 \nabla p \nabla \varphi d\theta dx.$$

By the Cauchy–Schwarz inequality we have

$$\int_{\Omega} h \frac{\partial \varphi}{\partial \theta} d\theta dx \leq \left| \int_{\Omega} h^3 |\nabla p|^2 d\theta dx \right|^{\frac{1}{2}} \left| \int_{\Omega} h^3 |\nabla \varphi|^2 d\theta dx \right|^{\frac{1}{2}}$$

and we deduce

$$(4.5) \quad \left| \int_{\Omega} h^3 |\nabla p|^2 d\theta dx \right|^{\frac{1}{2}} \geq \sup_{\varphi \in K, \varphi \neq 0} \frac{-\int_{\Omega} \varphi \frac{\partial h}{\partial \theta} d\theta dx}{\left| \int_{\Omega} h^3 |\nabla \varphi|^2 d\theta dx \right|^{\frac{1}{2}}}.$$

Let m be given by (2.4) and $\psi \in C^2(\mathbb{R})$ such that

- (i) $\text{supp}(\psi) \subset [-2, -1]$;
- (ii) $\psi \geq 0$;
- (iii) $\int_{\mathbb{R}} \psi > 0$.

As a consequence of (i)–(iii) we have

$$(4.6) \quad \frac{\int_{\mathbb{R}} \psi}{\left| \int_{\mathbb{R}} |\psi'|^2 \right|^{\frac{1}{2}}} = c_1 > 0.$$

Let

$$\varphi_\epsilon(\theta, x) := \psi \left(\frac{\theta - \theta_\alpha}{\sqrt{\epsilon}} \right) x(1 - x)$$

with θ_α defined in Lemma 4.3.

It is clear that $\text{supp}(\varphi_\epsilon) = I_{\alpha, \epsilon} \times [0, 1]$ with $I_{\alpha, \epsilon}$ as in Lemma 4.4. Using Lemma 4.4 (ii) we deduce

$$(4.7) \quad \begin{aligned} -\int_{\Omega} \frac{\partial h}{\partial \theta} \varphi_\epsilon d\theta dx &\geq \min_{\theta \in I_{\alpha, \epsilon}} \left(-\frac{\partial h}{\partial \theta} \right) \int_{\Omega} \varphi_\epsilon d\theta dx \\ &\geq \frac{m}{2}\epsilon \int_{\mathbb{R}} \psi dy \int_0^1 x(1 - x) dx = \frac{m}{12}\epsilon \int_{\mathbb{R}} \psi dy. \end{aligned}$$

On the other hand, using Lemma 4.4 (i) we have

$$\begin{aligned} \left| \int_{\Omega} h^3 |\nabla \varphi_{\epsilon}|^2 d\theta dx \right|^{\frac{1}{2}} &\leq \max_{I_{\alpha, \epsilon}} \{h^3\}^{\frac{1}{2}} \left| \int_{\Omega} |\nabla \varphi_{\epsilon}|^2 d\theta dx \right|^{\frac{1}{2}} \\ &\leq ((2M + 2)\epsilon)^{3/2} \left| \int_{\Omega} |\nabla \varphi_{\epsilon}|^2 d\theta dx \right|^{\frac{1}{2}}. \end{aligned}$$

A simple calculation gives

$$\int_{\Omega} |\nabla \varphi_{\epsilon}|^2 d\theta dx \leq c_3 \epsilon^{-1/2}.$$

We then have

$$(4.8) \quad \left| \int_{\Omega} h^3 |\nabla \varphi_{\epsilon}|^2 d\theta dx \right|^{\frac{1}{2}} \leq c_4 \epsilon^{5/4},$$

which, also using (4.7), gives us

$$(4.9) \quad \frac{-\int_{\Omega} \varphi_{\epsilon} \frac{\partial h}{\partial \theta} d\theta dx}{\left| \int_{\Omega} h^3 |\nabla \varphi_{\epsilon}|^2 d\theta dx \right|^{1/2}} \geq c_5 \epsilon^{-\frac{1}{4}}.$$

Finally from (4.5) and (4.9) we have the result. \square

We now introduce for any ϵ small the function $\tilde{h} : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$(4.10) \quad \tilde{h}(\theta, \alpha) := h(\theta_{\alpha}, a(\alpha) - \epsilon, \alpha) + [\rho''(\theta_{\alpha}) + \rho(\theta_{\alpha})](1 - \cos(\theta - \theta_{\alpha})).$$

Notice that

$$(4.11) \quad \begin{cases} h(\theta_{\alpha}, a(\alpha) - \epsilon, \alpha) = \tilde{h}(\theta_{\alpha}, \alpha) = \epsilon \cos(\theta_{\alpha} - \alpha), \\ \frac{\partial h}{\partial \theta}(\theta_{\alpha}, a(\alpha) - \epsilon, \alpha) - \frac{\partial \tilde{h}}{\partial \theta}(\theta_{\alpha}, \alpha) = \frac{\partial h}{\partial \theta}(\theta_{\alpha}, a(\alpha) - \epsilon, \alpha) = -\epsilon \sin(\theta_{\alpha} - \alpha), \\ \frac{\partial^2 h}{\partial \theta^2}(\theta_{\alpha}, a(\alpha) - \epsilon, \alpha) - \frac{\partial^2 \tilde{h}}{\partial \theta^2}(\theta_{\alpha}, \alpha) = -\epsilon \cos(\theta_{\alpha} - \alpha). \end{cases}$$

LEMMA 4.6. *There exists a constant c independent of ϵ and α such that*

$$\int_0^{2\pi} \frac{\left(h(\theta, a(\alpha) - \epsilon, \alpha) - \tilde{h}(\theta, \alpha) \right)^2}{h^3(\theta, a(\alpha) - \epsilon, \alpha)} d\theta \leq c.$$

Proof. By the Taylor polynomial of $h - \tilde{h}$ and h at θ_{α} and (4.11) we have

$$(4.12) \quad |h(\theta, a(\alpha) - \epsilon, \alpha) - \tilde{h}(\theta, \alpha)| \leq \epsilon |\theta - \theta_{\alpha}| + \epsilon |\theta - \theta_{\alpha}|^2 + c |\theta - \theta_{\alpha}|^3$$

and

$$(4.13) \quad \begin{aligned} h(\theta, a(\alpha) - \epsilon, \alpha) &= \epsilon \cos(\theta_{\alpha} - \alpha) - \epsilon \sin(\theta_{\alpha} - \alpha)(\theta - \theta_{\alpha}) \\ &+ \frac{1}{2} \left(\rho''(\theta_{\alpha}) + \rho(\theta_{\alpha}) - \epsilon \cos(\theta_{\alpha} - \alpha) \right) (\theta - \theta_{\alpha})^2 \\ &+ \frac{1}{6} \frac{\partial^3 h}{\partial \theta^3}(\hat{\theta}, a(\alpha) - \epsilon, \alpha) (\theta - \theta_{\alpha})^3 \end{aligned}$$

with $\hat{\theta} \in \mathbb{R}$.

From Lemma 4.2, (4.13), and (2.2) there exists m_1, m_2 , and m_3 positive and independent of ϵ and α such that

$$(4.14) \quad h(\theta, a(\alpha) - \epsilon, \alpha) \geq m_1\epsilon + m_2(\theta - \theta_\alpha)^2$$

$\forall \theta$ such that $|\theta - \theta_\alpha| \leq m_3$.

We have from (4.12)

$$\int_0^{2\pi} \frac{(h - \tilde{h})^2}{h^3} d\theta \leq 2(\epsilon^2 I_1 + \epsilon^2 I_2 + c^2 I_3)$$

with $I_k = \int_0^{2\pi} \frac{(\theta - \theta_\alpha)^{2k}}{h^3} d\theta, k = 1, 2, 3$.

Now from (4.14) it is clear that I_3 is bounded uniformly in ϵ and α .

We will prove the uniform estimate for $\epsilon^2 I_1$. For $\epsilon^2 I_2$ the proof is similar.

We have

$$I_1 = I_1^1 + I_1^2 + I_1^3,$$

where I_1^1, I_1^2, I_1^3 are the subintegrals, respectively, in the intervals $|\theta - \theta_\alpha| \geq m_3, \epsilon^{1/3} \leq |\theta - \theta_\alpha| \leq m_3$, and $0 \leq |\theta - \theta_\alpha| \leq \epsilon^{1/3}$.

It is clear that $\epsilon^2 I_1^1$ is bounded uniformly in ϵ and α since h is lower bounded by a positive constant on the interval $|\theta - \theta_\alpha| \geq m_3$.

From (4.14) we have

$$I_1^2 \leq m_2^{-3} \int_{\epsilon^{1/3} \leq |\theta - \theta_\alpha| \leq m_3} \frac{d\theta}{|\theta - \theta_\alpha|^4} \leq c\epsilon^{-4/3}$$

and

$$I_1^3 \leq m_1^{-3} \epsilon^{-3} \int_{0 \leq |\theta - \theta_\alpha| \leq \epsilon^{1/3}} |\theta - \theta_\alpha|^2 d\theta \leq m_1^{-3} \epsilon^{-2},$$

which proves the lemma. \square

5. Proof of Theorem 2.1. We apply the degree theory recalled in section 3 to the function G .

Since G is not defined on \bar{A} we introduce for any $\epsilon > 0$ small enough the domain

$$A_\epsilon = \left\{ (\eta_1, \eta_2) : 0 \leq \eta \leq a(\alpha) - \epsilon \right\}.$$

Let us now introduce the vector field

$$W : (\eta_1, \eta_2) \in A \rightarrow (\eta \sin \theta_\alpha, -\eta \cos \theta_\alpha) \in \mathbb{R}^2$$

with (η, α) defined in (1.2) and θ_α as in Lemma 4.3.

We observe that

$$W \cdot (-\eta_2, \eta_1) = -\eta^2 \cos(\theta_\alpha - \alpha).$$

From Lemmas 4.2 and 4.3 we have that

$$W \cdot (-\eta_2, \eta_1) < 0 \quad \forall (\eta_1, \eta_2) \in A.$$

From Theorem 3.1 we deduce that

$$\text{deg}(W, A_\epsilon, 0) = \text{deg}((-\eta_2, \eta_1), A_\epsilon, 0) = 1.$$

Notice that $\text{deg}(f, S, 0) = (-1)^n \text{deg}(-f, S, 0)$, where n is the dimension of the domain (for the application here, $n = 2$, so the sign does not change).

It suffices now to prove the inequality

$$(5.1) \quad G(\eta_1, \eta_2) \cdot W(\eta_1, \eta_2) > 0 \quad \forall (\eta_1, \eta_2) \in \partial A_\epsilon,$$

and the proof is finished using again Theorem 3.1.

We have, $\forall (\eta_1, \eta_2) \in \partial A_\epsilon, \eta = a(\alpha) - \epsilon$. Then

$$(5.2) \quad G(\eta_1, \eta_2) \cdot W(\eta_1, \eta_2) = -\eta \int_{\Omega} p \sin(\theta - \theta_\alpha) d\theta dx + \eta(-F_1 \sin \theta_\alpha + F_2 \cos \theta_\alpha).$$

Now taking $\varphi = 0$ and $\varphi = 2p$, respectively, in (1.5) we obtain

$$(5.3) \quad \int_{\Omega} h^3 |\nabla p|^2 d\theta dx = - \int_{\Omega} \frac{\partial h}{\partial \theta} p d\theta dx.$$

Notice that

$$(5.4) \quad \int_{\Omega} \frac{\partial h}{\partial \theta} p d\theta dx = \int_{\Omega} \frac{\partial}{\partial \theta} (h - \tilde{h}) p d\theta dx + \int_{\Omega} \frac{\partial \tilde{h}}{\partial \theta} p d\theta dx$$

with \tilde{h} given by (4.10).

We now have

$$(5.5) \quad \begin{aligned} \int_{\Omega} \frac{\partial}{\partial \theta} (h - \tilde{h}) p d\theta dx &= - \int_{\Omega} (h - \tilde{h}) \frac{\partial p}{\partial \theta} d\theta dx \\ &\geq - \int_{\Omega} \frac{|h - \tilde{h}|}{h^{3/2}} h^{3/2} |\nabla p| d\theta dx \\ &\geq - \frac{1}{2} \int_{\Omega} \frac{|h - \tilde{h}|^2}{h^3} d\theta dx - \frac{1}{2} \int_{\Omega} h^3 |\nabla p|^2 d\theta dx. \end{aligned}$$

Since

$$\frac{\partial \tilde{h}}{\partial \theta} = (\rho''(\theta_\alpha) + \rho(\theta_\alpha)) \sin(\theta - \theta_\alpha)$$

we obtain from (5.3) and (5.5)

$$\begin{aligned} \int_{\Omega} h^3 |\nabla p|^2 d\theta dx &\leq \frac{1}{2} \int_{\Omega} \frac{|h - \tilde{h}|^2}{h^3} + \frac{1}{2} \int_{\Omega} h^3 |\nabla p|^2 d\theta dx \\ &\quad - (\rho''(\theta_\alpha) + \rho(\theta_\alpha)) \int_{\Omega} \sin(\theta - \theta_\alpha) p d\theta dx, \end{aligned}$$

which implies, using also hypotheses (2.2),

$$- \int_{\Omega} \sin(\theta - \theta_\alpha) p d\theta dx \geq \frac{\int_{\Omega} h^3 |\nabla p|^2 d\theta dx - \int_{\Omega} \frac{|h - \tilde{h}|^2}{h^3} d\theta dx}{2(\rho''(\theta_\alpha) + \rho(\theta_\alpha))}.$$

Since ρ is in C^2 and from Lemmas 4.5 and 4.6, we obtain for ϵ small enough

$$-\int_{\Omega} \sin(\theta - \theta_{\alpha})p \geq c(\epsilon^{-1/2} - 1)$$

with $c > 0$ a constant independent of ϵ and α .

Since $\eta = a(\alpha) - \epsilon$ and thanks to (5.2) we deduce

$$G(\eta_1, \eta_2) \cdot W(\eta_1, \eta_2) \geq \eta \left(c(\epsilon^{-1/2} - 1) - \|F\| \right).$$

Taking ϵ small enough, we prove the theorem.

Remark 5.1. In the same manner we can prove the existence of at least one equilibrium solution for some other similar problems in this context. For instance, the case of Dirichlet boundary conditions on every boundary of Ω , or the case of a one-dimensional domain with Dirichlet boundary conditions.

Acknowledgments. The third author thanks INSA-Lyon and Institut Camille Jordan at University of Lyon for their hospitality.

REFERENCES

- [1] G. BUSCAGLIA, I. CIUPERCA, I. HAFIDI, AND M. JAI, *Existence of equilibria in articulated bearings*, J. Math. Anal. Appl., 328 (2007), pp. 24–45.
- [2] G. BUSCAGLIA, I. CIUPERCA, I. HAFIDI, AND M. JAI, *Existence of equilibria in articulated bearings in presence of cavity*, J. Math. Anal. Appl., 335 (2007), pp. 841–859.
- [3] J. DURANY, G. GARCÍA, AND C. VÁZQUEZ, *Numerical simulation of a lubricated Hertzian contact problem under imposed load*, Finite Elem. Anal. Des., 38 (2002), pp. 645–658.
- [4] J. DURANY, J. PEREIRA, AND F. VARAS, *Numerical solution of steady and transient problems in thermohydrodynamic lubrication using a combination of finite element, finite volume and boundary element methods*, Finite Elem. Anal. Des., 44 (2008), pp. 686–695.
- [5] J. FRÈNE, D. NICOLAS, B. DEGUERCE, D. BERTHE, AND M. GODET, *Lubrification Hydrodynamique*, Eyrolles, Paris, 1990.
- [6] D. KINDERLEHRER AND G. STAMPACCHIA, *An Introduction to Variational Inequalities and Their Applications*, Academic Press, New York, 1980.
- [7] L. LELOUP, *Etude de la Lubrification et Calcul des Paliers*, Dunod, Paris, 1962.
- [8] E. H. ROTHE, *Introduction to Various Aspects of Degree Theory in Banach Spaces*, Math. Surveys Monogr. 23, AMS, Providence, RI, 1986.
- [9] O. REYNOLDS, *On the theory of lubrication and its application to Mr Beauchamp Tower's experiments, including an experimental determination of the viscosity of olive oil*, Phil. Trans. Roy. Soc. A, 117 (1886), pp. 157–234.

POINTWISE GREEN FUNCTION BOUNDS AND LONG-TIME STABILITY OF LARGE-AMPLITUDE NONCHARACTERISTIC BOUNDARY LAYERS*

SHANTIA YARAHMADIAN† AND KEVIN ZUMBRUN†

Abstract. Using pointwise semigroup techniques of Zumbrun–Howard and Mascia–Zumbrun, we obtain sharp global pointwise Green function bounds for noncharacteristic boundary layers of arbitrary amplitude. These estimates allow us to analyze linearized and nonlinearized stability of noncharacteristic boundary layers of one-dimensional systems of conservation laws, showing that both follow from (and linearized stability is equivalent to) a numerically checkable Evans function condition. Our results extend to the large-amplitude case results obtained for small amplitudes by Matsumura, Nishihara, and others using energy estimates.

Key words. boundary layer stability, Evans function, pointwise Green function bounds

AMS subject classifications. Primary, 35B35; Secondary, 76N20

DOI. 10.1137/080714804

1. Introduction. Boundary layers appear in many physical settings, such as gas dynamics, MHD, and rotating fluids; see, for example, the physical discussion in [34]. In this paper, we study the stability of boundary layers assuming that the boundary layer solution is noncharacteristic, which means that signals are transmitted into or out of, but not along the boundary. Specifically, we consider a boundary layer, or stationary solution,

$$(1.1) \quad u = \bar{u}(x), \quad \lim_{x \rightarrow +\infty} \bar{u}(x) = u_+, \quad \bar{u}(0) = u_0$$

of a system of conservation laws on the quarter-plane

$$(1.2) \quad u_t + f(u)_x = (B(u)u_x)_x, \quad x, t > 0,$$

$u, f \in \mathbb{R}^n$, $B \in \mathbb{R}^{n \times n}$, with initial data $u(x, 0) = g(x)$ and Dirichlet boundary condition

$$(1.3) \quad u(0, t) = h(t).$$

A fundamental question is whether or not such boundary layer solutions are *stable* in the sense of PDE, i.e., whether or not a sufficiently small perturbation of \bar{u} remains close to \bar{u} , or converges time-asymptotically to \bar{u} , under the evolution of (1.2).

Long-time stability of boundary layers has been considered for scalar equations in [24, 25] and for the equations of isentropic gas dynamics in [26, 23]. The latter results, obtained by energy estimates, apply to arbitrary amplitude layers of “expansive inflow” type analogous to rarefaction waves, but only to small-amplitude layers of “compressive inflow or outflow” type analogous to shock waves or “expansive outflow” type. For general symmetric hyperbolic-parabolic systems, stability of small-amplitude noncharacteristic boundary layers has been shown in multidimensions for

*Received by the editors January 31, 2008; accepted for publication (in revised form) October 9, 2008; published electronically February 25, 2009. This work was supported in part by the National Science Foundation grant number DMS-0300487.

<http://www.siam.org/journals/sima/40-6/71480.html>

†Department of Mathematics, Indiana University, Bloomington, IN 47405 (syarahma@indiana.edu, kzumbrun@indiana.edu).

strictly parabolic systems in [11], and in one dimension for partially parabolic (“real viscosity”) systems in [30].

Here, in the spirit of results obtained for shock waves in [37, 27, 28], we show for general strictly parabolic systems of conservation laws that linearized and nonlinear stability follow from and linearized stability is equivalent to a generalized spectral stability condition phrased in terms of the Evans function associated with the linearized equations about the wave, *independent of the amplitude of the boundary layer in question*. The Evans function, introduced for boundary layers in [32], is described further in section 2; for its origins in the study of stability of traveling waves, see for example, [1, 10] and references therein.

The Evans condition is readily checkable numerically, and in some cases analytically; see [3, 4, 5, 6, 20, 2, 17]. In particular, stability of small-amplitude uniformly noncharacteristic boundary layers has been shown for general hyperbolic–parabolic systems in multidimensions in [13] using elementary Evans function arguments (convergence to the constant layer). An exhaustive numerical study for isentropic gas layers in one dimension has been carried out in [7], with the conclusion of stability for arbitrary amplitudes. Results of instability in some cases have been shown in [32, 33] using a mod two stability index in the spirit of [10].

Our method of analysis [35] is by pointwise Green function methods like those used in [37, 27, 28], and especially [16], to analyze the stability of viscous shock layers. Similar results have been obtained for the related small-viscosity-limit problem in [12, 29, 14]. In particular, we point to the analysis of Grenier and Rousset [12] as using pointwise Green function estimates very similar to those that we use here, though adapted for different purposes.

1.1. Equations and assumptions. Consider a *viscous boundary layer*, a standing-wave solution (1.1) of a general parabolic system of conservation laws (1.2).

Assume, similarly, as in the treatment of the viscous shock case in [16]:

(H0) $f, B \in C^3$.

(H1) $Re \sigma(B) > 0$.

(H2) $\sigma(f'(u_+))$ real, distinct, and nonzero.

(H3) $Re \sigma(-ikf'(u_+) - k^2B(u_+)) < -\theta k^2$ for all real k and $\theta > 0$.

LEMMA 1.1 ([27, 36]). *Assuming (H0)–(H3), a solution \bar{u} of (1.1) if it exists is also unique; moreover,*

$$(1.4) \quad |(d/dx)^k(\bar{u} - u_+)| \leq C e^{-\theta x}, \quad k = 0, \dots, 4,$$

as $x \rightarrow +\infty$, for some $\theta > 0$.

Proof. As in the shock case [28, 36], (1.4) follows by the observation that, under hypotheses (H0)–(H3), u_+ is a hyperbolic rest point of the layer profile ODE. Uniqueness follows by the observation [27] that the standing-wave ODE may be integrated from x to $+\infty$ and rearranged to yield

$$(1.5) \quad B(u)u' = f(u) - f(u_+)$$

together with the boundary conditions at $x = 0$, thus determining a unique solution for all $x \geq 0$. \square

1.2. Linearized stability and the Evans function. After linearizing (1.2) about the stationary solution \bar{u} , we obtain the linearized equation

$$(1.6) \quad u_t = Lu := -(Au)_x + (Bu_x)_x, \quad A, B \in C^2,$$

where

$$(1.7) \quad B := B(\bar{u})$$

and

$$(1.8) \quad Au := dF(\bar{u})u - dB(\bar{u})(u, \bar{u}_x).$$

DEFINITION 1.2. *The boundary layer \bar{u} is said to be linearly asymptotically stable if $u(\cdot, t)$ approaches 0 as $t \rightarrow \infty$, for any solution u of (1.6) with initial data bounded in some specified norm.*

We define the following *stability criterion*, where $D(\lambda)$, described below, denotes the Evans function associated with the linearized operator L about the layer, an analytic function analogous to the characteristic polynomial of a finite-dimensional operator, whose zeroes away from the essential spectrum agree in location and multiplicity with the eigenvalues of L :

$$(1.9) \quad \text{There exist no zeroes of } D(\cdot) \text{ in the nonstable half-plane } \text{Re}\lambda \geq 0.$$

As discussed, e.g., in [31], under assumptions (H0)–(H3), this is equivalent to *strong spectral stability*, $\sigma(L) \subset \{\text{Re}\lambda < 0\}$, *transversality* of \bar{u} as a solution of the connection problem in the associated standing-wave ODE, and *hyperbolic stability* of an associated boundary value problem obtained by formal matched asymptotics. Here, and elsewhere, σ denotes spectrum of a linearized operator or matrix.

Our first main result is as follows.

THEOREM 1.3. *Assuming (H0)–(H3), linearized asymptotic $L^1 \cap L^p \rightarrow L^p$ stability, $p > 1$, is equivalent to (1.9).*

Theorem 1.3 is obtained as a consequence of the following detailed, pointwise bounds on the Green function $G(x, t; y)$ of the linearized evolution equations (1.6) with homogeneous boundary conditions (more properly speaking, a distribution), defined by:

- (i) $(\partial_t - L_x)G = 0$ in the distributional sense, for all $x, y, t > 0$;
- (ii) $G(x, t; y) \rightarrow \delta(x - y)$ as $t \rightarrow 0$;
- (iii) $G(0, t; y) \equiv 0$, for all $y, t > 0$.

Denote by

$$(1.10) \quad a_1^+ < a_2^+ < \dots < a_n^+$$

the eigenvalues of the limiting convection matrix $A_+ := df(u_+)$.

Then, our second main result is as follows.

THEOREM 1.4. *Assuming (H0)–(H3) and stability condition (1.9),*

$$(1.11) \quad \begin{aligned} & |\partial_x^\gamma \partial_y^\alpha G(x, t; y)| \leq C e^{-\eta(|x-y|+t)} \\ & + C \left(t^{-|\alpha|/2} + |\alpha| e^{-\eta|y|} + |\gamma| e^{-\eta|x|} \right) \left(\sum_{k=1}^n t^{-1/2} e^{-(x-y-a_k^- t)^2/Mt} \right. \\ & \left. + \sum_{a_k^+ < 0, a_j^+ > 0} \chi_{\{|a_k^+ t| \geq |y|\}} t^{-1/2} e^{-(x-a_j^+(t-|y/a_k^+|))^2/Mt} \right), \end{aligned}$$

$0 \leq |\alpha|, |\gamma| \leq 1$, for some $\eta, C, M > 0$, where x^\pm denotes the positive/negative part of x , indicator function $\chi_{\{|a_k^- t| \geq |y|\}}$ is 1 for $|a_k^- t| \geq |y|$ and 0 otherwise.

1.3. Nonlinear stability.

DEFINITION 1.5. *The boundary layer \bar{u} is said to be nonlinearly asymptotically stable if $\tilde{u}(\cdot, t)$ exists for all $t \geq 0$ and approaches \bar{u} as $t \rightarrow \infty$, for any solution \tilde{u} of (1.2) with initial data sufficiently close in some norm to the original layer \bar{u} .*

Denoting by

$$(1.12) \quad a_1^+ < a_2^+ < \dots < a_n^+$$

the eigenvalues of of the limiting convection matrix $A_+ := df(u_+)$, define

$$(1.13) \quad \theta(x, t) := \sum_{a_j^+ > 0} (1+t)^{-1/2} e^{-|x-a_j^+t|^2/Lt},$$

$$(1.14) \quad \psi_1(x, t) := \chi(x, t) \sum_{a_j^+ > 0} (1+|x|+t)^{-1/2} (1+|x-a_j^+t|)^{-1/2},$$

and

$$(1.15) \quad \psi_2(x, t) := (1-\chi(x, t)) \left(1+|x-a_n^+t|+t^{1/2}\right)^{-3/2},$$

where $\chi(x, t) = 1$ for $x \in [0, a_n^+t]$ and $\chi(x, t) = 0$ otherwise and $L > 0$ is a sufficiently large constant. For simplicity, take B identically constant. Then, our third and final main result is as follows.

THEOREM 1.6. *Assuming (H0)–(H3), $B \equiv \text{constant}$, and the linear stability condition (1.9), the profile \bar{u} is nonlinearly asymptotically stable with respect to perturbations g, h in initial and boundary data satisfying*

$$|g(x)| \leq E_0(1+|x|)^{-3/2}, \quad |h(t)| \leq E_0(1+|t|)^{-3/2}, \quad |h'(t)| \leq E_0(1+|t|)^{-1}$$

for E_0 sufficiently small. More precisely,

$$(1.16) \quad |\tilde{u}(x, t) - \bar{u}(x)| \leq CE_0(\theta + \psi_1 + \psi_2)(x, t),$$

where \tilde{u} denotes the solution of (1.2) with initial data $\tilde{g} = \bar{u} + g$ and boundary data $\tilde{h} = u_0 + h$.

Remark 1.7. Pointwise bound (1.16) yields as a corollary the sharp L^p decay rate

$$(1.17) \quad |\tilde{u}(x, t) - \bar{u}(x)|_{L^p} \leq CE_0(1+t)^{-\frac{1}{2}(1-\frac{1}{p})}, \quad 1 \leq p \leq \infty.$$

1.4. Discussion and open problems. The case of boundary layers is quite analogous to the undercompressive shock case; in particular, pointwise estimates as in [16] appear to be necessary to close the one-dimensional analysis by a linearized semigroup approach suitable for large-amplitude layers. (On the other hand, small-amplitude stability has been established using energy estimates in, e.g., [26, 11].) A new feature of the present analysis as compared to those of [16, 18, 19] is the admission of perturbations in boundary as well as initial data. Open problems are extensions to systems with physical (partial) or quasilinear viscosity and to multidimensional boundary layers.

2. The Evans function. Before starting the analysis, we review the basic Evans function methods and gap/conjugation lemma.

2.1. The gap/conjugation lemma. Consider a family of first-order ODE systems on the half-line:

$$(2.1) \quad \begin{aligned} W' &= \mathbb{A}(x, \lambda)W, & \lambda \in \Omega & \text{ and } x > 0, \\ \mathbb{B}(\lambda)W &= 0, & \lambda \in \Omega & \text{ and } x = 0. \end{aligned}$$

These systems of ODEs should be considered as a generalized eigenvalue equation, with λ representing frequency. We assume that the boundary matrix \mathbb{B} is analytic in λ and that the coefficient matrix \mathbb{A} is analytic in λ as a function from Ω into $L^\infty(x)$, C^K in x , and approaches exponentially to a limit $\mathbb{A}_+(\lambda)$ as $x \rightarrow \infty$, with uniform exponentially decay estimates

$$(2.2) \quad |(\partial/\partial x)^k(\mathbb{A} - \mathbb{A}_+)| \leq C_1 e^{-\theta|x|/C_2}, \quad \text{for } x > 0, 0 \leq k \leq K,$$

$C_j, \theta > 0$, on compact subsets of Ω . Now we can state a refinement of the ‘‘Gap Lemma’’ of [10, 21], relating solutions of the variable-coefficient ODE to the solutions of its constant-coefficient limiting equations

$$(2.3) \quad Z' = \mathbb{A}_+(\lambda)Z$$

as $x \rightarrow +\infty$.

LEMMA 2.1 (conjugation Lemma [29]). *Under assumption (2.2), there exists locally to any given $\lambda_0 \in \Omega$ a linear transformation $P_+(x, \lambda) = I + \Theta_+(x, \lambda)$ on $x \geq 0$, Φ_+ analytic in λ as functions from Ω to $L^\infty[0, +\infty)$, such that:*

(i) $|P_+|$ and their inverses are uniformly bounded, with

$$(2.4) \quad |(\partial/\partial \lambda)^j(\partial/\partial x)^k \Theta_+| \leq C(j)C_1C_2 e^{-\theta|x|/C_2} \quad \text{for } x > 0, 0 \leq k \leq K + 1,$$

$j \geq 0$, where $0 < \theta < 1$ is an arbitrary fixed parameter, and $C > 0$ and the size of the neighborhood of definition depend only on θ, j , the modulus of the entries of \mathbb{A} at λ_0 , and the modulus of continuity of \mathbb{A} on some neighborhood of $\lambda_0 \in \Omega$.

(ii) *The change of coordinates $W := P_+Z$ reduces (2.1) on $x \geq 0$ to the asymptotic constant-coefficient equations (2.3). Equivalently, solutions of (2.1) may be conveniently factorized as*

$$(2.5) \quad W = (I + \Theta_+)Z_+,$$

where Z_+ are solutions of the constant-coefficient equations, and Θ_+ satisfy (2.4).

Proof. As described in [27], for $j = k = 0$ this is a straightforward corollary of the gap lemma as stated in [Z.3], applied to the ‘‘lifted’’ matrix-valued ODE

$$P' = \mathbb{A}_+P - P\mathbb{A} + (\mathbb{A} - \mathbb{A}_+)P$$

for the conjugating matrices P_+ . The x -derivative bounds $0 < k \leq K + 1$ then follow from the ODE and its first K derivatives. Finally, the λ -derivative bounds follow from standard interior estimates for analytic functions. \square

DEFINITION 2.2. *Following [1], we define the domain of consistent splitting for the ODE system $W' = \mathbb{A}(x, \lambda)W$ as the (open) set of λ such that the limiting matrix \mathbb{A}_+ is hyperbolic (has no center subspace) and the boundary matrix \mathbb{B} is full rank, with $\dim S_+ = \text{rank } \mathbb{B}$.*

LEMMA 2.3. *On any simply connected subset of the domain of consistent splitting, there exists an analytic basis $\{v_1, \dots, v_k\}$ for the subspace S_+ defined in Definition 2.2.*

Proof. By spectral separation of S_+ from the complementary unstable subspace U_+ , the associated (group) eigenprojection is analytic. The existence of an analytic eigenbasis then follows by a standard result of Kato; see [22], pp. 99–102. \square

COROLLARY 2.4. *By the Conjugation Lemma, on the domain of consistent splitting, the stable manifold of solutions decaying as $x \rightarrow +\infty$ of (2.1) is*

$$(2.6) \quad \mathcal{S}^+ := \text{span} \{P_+v_1^+, \dots, P_+v_k^+\},$$

where $W_+^j := P_+v_j^+$ are analytic in λ and C^{K+1} in x for $\mathbb{A} \in C^K$.

2.2. Definition of the Evans function. On any simply connected subset of the domain of consistent splitting, let $W_1^+, \dots, W_k^+ = P_+v_1^+, \dots, P_+v_k^+$ be the analytic basis described in Corollary 2.4 of the subspace \mathcal{S}^+ of solutions W of (2.1) satisfying the boundary condition $W \rightarrow 0$ at $+\infty$. Then, the *Evans function* for the ODE systems $W' = \mathbb{A}(x, \lambda)W$ associated with this choice of limiting bases is defined as the $k \times k$ Gramian determinant

$$(2.7) \quad \begin{aligned} D(\lambda) &:= \det (\mathbb{B}W_1^+, \dots, \mathbb{B}W_k^+) |_{x=0, \lambda} \\ &= \det (\mathbb{B}P_+v_1^+, \dots, \mathbb{B}P_+v_k^+) |_{x=0, \lambda}. \end{aligned}$$

Remark 2.5. Note that D is independent of the choice of P_+ as, by uniqueness of stable manifolds, the exterior products (minors) $P_+v_1^+ \wedge \dots \wedge P_+v_k^+$ are uniquely determined by their behavior as $x \rightarrow +\infty$.

PROPOSITION 2.6. *Both the Evans function and the subspace \mathcal{S}^+ are analytic on the entire simply connected subset of the domain of consistent splitting on which they are defined. Moreover, for λ within this region, equation (2.1) admits a nontrivial solution $W \in L^2(x > 0)$ if and only if $D(\lambda) = 0$.*

Proof. Analyticity follows by uniqueness, and local analyticity of P_+, v_k^+ . Noting that the first $P_+v_j^+$ are a basis for the stable manifold of (2.1) at $x \rightarrow +\infty$, we find that the determinant of $\mathbb{B}P_+v_j^+$ vanishes if and only if $\mathbb{B}(\lambda)$ has nontrivial kernel on $\mathcal{S}_+(\lambda, 0)$, whence, the second assertion follows. \square

Remark 2.7. In the case that the ODE system describes an eigenvalue equation associated with an ordinary differential operator L , Proposition 2.6 implies that eigenvalues of L agree in location with zeroes of D . (Indeed, they agree also in multiplicity; see [8, 9]; Lemma 6.1, [37]; or Proposition 6.15 of [27].)

When $\ker \mathbb{B}$ has an analytic basis $v_{k+1}^0, \dots, v_{N-k}^0$, for example, in the commonly occurring case, as here, that $\mathbb{B} \equiv \text{constant}$, we have the following useful alternative formulation. This is the version that we will use in our analysis of the Green function and resolvent kernel.

PROPOSITION 2.8. *Let $v_{k+1}^0, \dots, v_{N-k}^0$ be an analytic basis of $\ker \mathbb{B}$, normalized so that $\det(\mathbb{B}^*, v_{k+1}^0, \dots, v_N^0) \equiv 1$. Then, the solutions W_j^0 of (2.1) determined by initial data $W_j^0(\lambda, 0) = v_j^0$ are analytic in λ and C^{K+1} in x , and*

$$(2.8) \quad D(\lambda) := \det (W_1^+, \dots, W_k^+, W_{k+1}^0, \dots, W_N^0) |_{x=0, \lambda}.$$

Proof. Analyticity/smoothness follow by analytic/smooth dependence on initial data/parameters. By the chosen normalization, and standard properties of Gramian determinants, $D(\lambda) = \det(W_1^+, \dots, W_k^+, v_{k+1}^0, \dots, v_N^0) |_{x=0, \lambda}$, yielding (2.8). \square

3. Construction of the resolvent kernel. In this section we construct the explicit form of the resolvent kernel, which is nothing more than the Green function

$G_\lambda(x, y)$ associated with the elliptic operator $(L - \lambda I)$, where

$$(3.1) \quad (L - \lambda I)G_\lambda(\cdot, y) = \delta_y I, \quad G_\lambda(0, y) \equiv 0.$$

Let Λ be the region of consistent splitting for L . It is an established fact (see [15]) that the resolvent $(L - \lambda I)^{-1}$ and the Green function $G_\lambda(x, y)$ are meromorphic in λ on Λ , with isolated poles of finite order. G_λ , in fact, admits a meromorphic extension to a sector

$$(3.2) \quad \Omega_\theta = \{\lambda : \operatorname{Re}(\lambda) \geq -\theta_1 - \theta_2 |\operatorname{Im}(\lambda)|\}, \quad \theta_1, \theta_2 > 0.$$

Writing the associated eigenvalue equation in the form of a first-order system (2.1), we obtain

$$(3.3) \quad W' = A(\lambda, x)W, \quad \mathbb{B}W(0) = 0,$$

where

$$W = \begin{pmatrix} w \\ w' \end{pmatrix} \in \mathbb{C}^{2n}, \quad A = \begin{pmatrix} 0 & I \\ \lambda B^{-1} + A'B^{-1} & AB^{-1} - B'B^{-1} \end{pmatrix},$$

and $\mathbb{B} \equiv$ constant is the rank- n projection onto the first coordinate w of W , with kernel spanned by the constant basis $v_{n+j}^0 = e_{n+j}$, $j = 1, \dots, n$ and e_j the j th standard basis element.

Denote by

$$(3.4) \quad \Phi^0 = (\phi_1^0(x; \lambda) \ \cdots \ \phi_n^0(x; \lambda)) = (W_1^0 \ \cdots \ W_n^0)$$

and

$$(3.5) \quad \Phi^+ = (\phi_1^+(x; \lambda) \ \cdots \ \phi_n^+(x; \lambda)) = (W_{n+1}^+ \ \cdots \ W_{2n}^+) = (P_+ v_1^+ \ \cdots \ P_+ v_k^+)$$

the matrices whose columns span the subspaces of solutions of (2.1) decaying at $x = 0, +\infty$, respectively, denoting (analytically chosen) complementary subspaces by

$$(3.6) \quad \Psi^0 = (\psi_1^0(x; \lambda) \ \cdots \ \psi_n^0(x; \lambda)) = (W_{n+1}^0 \ \cdots \ W_{2n}^0)$$

and

$$(3.7) \quad \Psi^+ = (\psi_1^+(x; \lambda) \ \cdots \ \psi_n^+(x; \lambda)) = (W_1^+ \ \cdots \ W_n^+).$$

As described in the previous subsection, eigenfunctions decaying at both $0, +\infty$ occur precisely when the subspaces $\operatorname{span} \Phi^0$ and $\operatorname{span} \Phi^+$ intersect, i.e., at zeros of the Evans function defined in (2.8):

$$(3.8) \quad D_L(\lambda) := \det(\Phi^0, \Phi^+)_{|x=0} = (\phi_1^0 \wedge \cdots \wedge \phi_n^0 \wedge \phi_1^+ \wedge \cdots \wedge \phi_n^+)_{|x=0}.$$

LEMMA 3.1 ([10, 37]). *For $\theta_1, \theta_2 > 0$ sufficiently small, D_L is locally analytic on sector Ω_θ as defined in (3.2).*

Proof. Direct calculation showing that the domain Λ of consistent splitting is contained in $\Omega_\theta - B(0, r)$ for $r > 0$ arbitrary and θ sufficiently small, with v_j^\pm extending analytically to $B(0, r)$. \square

LEMMA 3.2. Let $H_\lambda(x, y)$ denote the Green function for the adjoint operator $(L - \lambda I)^*$ on the half-plane $x \geq 0$. Then $G_\lambda(x, y) = H_\lambda^*(x, y)$. In particular, for $x \neq y$, the matrix $z = G_\lambda(x, \cdot)$ satisfies

$$(3.9) \quad (z'B)' = -z'A + z\lambda.$$

Proof. Standard duality argument; see [37] for operators on the whole line. \square
 Considering (3.9) as an ODE system for the vector $Z = (z, z')$, it becomes

$$(3.10) \quad Z' = Z\tilde{A}(\lambda, x),$$

where

$$(3.11) \quad \tilde{A} = \begin{pmatrix} 0 & \lambda B^{-1} - A'B^{-1} \\ I & -AB^{-1} - B'B^{-1} \end{pmatrix}.$$

LEMMA 3.3 ([37]). Z is a solution of (3.11) if and only if $ZSW \equiv \text{constant}$ for any solution W of (2.1), where $\mathcal{S} = \begin{pmatrix} -A & B \\ -B & 0 \end{pmatrix}$.

Proof. Direct computation/comparison with 0 of $(ZSW)'$; see [37]. \square
 Using Lemma 3.3, we can define dual bases \tilde{W}_j^0 and \tilde{W}_j^+ by the relations

$$(3.12) \quad \tilde{W}_j^{0,+} \mathcal{S} W_k^{0,+} = \delta_k^j.$$

Likewise, $\tilde{A}_{0,+}$ can be defined as

$$(3.13) \quad \tilde{A}_{0,+} = \begin{pmatrix} 0 & \lambda B_{0,+}^{-1} \\ I & -A_{0,+} B_{0,+}^{-1} \end{pmatrix}.$$

We define also the dual subspaces

$$(3.14) \quad \tilde{\Phi}^0 = (\tilde{\phi}_1^0(x; \lambda) \ \cdots \ \tilde{\phi}_n^0(x; \lambda)) = (\tilde{W}_{n+1}^0 \ \cdots \ \tilde{W}_{2n}^0),$$

$$(3.15) \quad \tilde{\Phi}^+ = (\tilde{\phi}_1^+(x; \lambda) \ \cdots \ \tilde{\phi}_n^+(x; \lambda)) = (\tilde{W}_1^+ \ \cdots \ \tilde{W}_n^+),$$

$$(3.16) \quad \tilde{\Psi}^0 = (\tilde{\psi}_1^0(x; \lambda) \ \cdots \ \tilde{\psi}_n^0(x; \lambda)) = (\tilde{W}_1^0 \ \cdots \ \tilde{W}_n^0),$$

$$(3.17) \quad \tilde{\Psi}^+ = (\tilde{\psi}_1^+(x; \lambda) \ \cdots \ \tilde{\psi}_n^+(x; \lambda)) = (\tilde{W}_{n+1}^+ \ \cdots \ \tilde{W}_{2n}^+).$$

With these preparations, the construction of the resolvent kernel goes exactly as in the construction performed in [37, 27] on the whole line.

LEMMA 3.4. We have the the representation

$$(3.18) \quad \begin{pmatrix} G_\lambda & G_{\lambda_y} \\ G_{\lambda_x} & G_{\lambda_{xy}} \end{pmatrix} = \begin{cases} \Phi^+(\lambda, x) M^+(\lambda) \tilde{\Psi}^0(\lambda, y) & \text{for } x > y, \\ \Phi^0(\lambda, x) M^0(\lambda) \tilde{\Psi}^+(\lambda, y) & \text{for } x < y, \end{cases}$$

where $M^{0,+}$ are to be determined.

Proof. See [37] Lemma 4.6. \square

Using Lemma 3.4, we find the explicit coordinate-free representation for $x > y$:

$$(3.19) \quad \begin{pmatrix} G_\lambda & G_{\lambda_y} \\ G_{\lambda_x} & G_{\lambda_{xy}} \end{pmatrix} = \mathcal{F}^{z \rightarrow x} \Pi_+(z) \mathcal{S}^{-1}(z) \tilde{\Pi}_0(z) \tilde{\mathcal{F}}^{z \rightarrow y},$$

where

$$(3.20) \quad \Pi_+(y) = (\Phi^+(y), 0) (\Phi^+(y), \Phi^-(y))^{-1},$$

$$(3.21) \quad \tilde{\Pi}_0(y) = \begin{pmatrix} \tilde{\Psi}^0(y) \\ \tilde{\Psi}^+(y) \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\Psi}^0(y) \\ 0 \end{pmatrix},$$

$$(3.22) \quad \mathcal{F}^{z \rightarrow x} = (\Phi^+(x), \Phi^0(x)) (\Phi^+(z), \Phi^0(z))^{-1},$$

$$(3.23) \quad \tilde{\mathcal{F}}^{z \rightarrow y} = \begin{pmatrix} \tilde{\Psi}^0(z) \\ \Phi^+(z) \end{pmatrix} \begin{pmatrix} \Psi^0(y) \\ \Psi^+(y) \end{pmatrix}^{-1},$$

and similarly for $x < y$.

COROLLARY 3.5. *The resolvent kernel may be expressed as*

$$(3.24) \quad G_\lambda(x, y) = \begin{cases} (I_n, 0)\Phi^+(x; \lambda)M^+(\lambda)\tilde{\Psi}^{0*}(y; \lambda)(I_n, 0)^{tr} & x > y, \\ -(I_n, 0)\Phi^0(x; \lambda)M^0(\lambda)\tilde{\Psi}^{+*}(y; \lambda)(I_n, 0)^{tr} & x < y, \end{cases}$$

where

$$(3.25) \quad M(\lambda) := \text{diag}(M^+(\lambda), M^0(\lambda)) = \Phi^{-1}(z; \lambda)\bar{S}^{-1}(z)\tilde{\Psi}^{-1*}(z; \lambda).$$

4. Low-frequency bounds. Our goal in this section is the estimation of the resolvent kernel in the critical regime $|\lambda| \rightarrow 0$, i.e., the large time behavior of the Green function G , or global behavior in space and time. We are basically following the same treatment as that carried out for viscous shock waves of strictly parabolic conservation laws in [37, 27]; we refer to those references for details. In the low frequency case the behavior is essentially governed by the equation

$$(4.1) \quad U_t = L_+ U := -A_+ U_x + B_+ U_{xx}.$$

PROPOSITION 4.1. *Assuming (H0)–(H3), let K be the order of the pole of G_λ at $\lambda = 0$ and r be sufficiently small that there are no other poles in $B(0, r)$. Then for $\lambda \in \Omega_\theta$ such that $|\lambda| \leq r$ and for $x > y > 0$ we have*

$$(4.2) \quad \begin{pmatrix} G_\lambda & G_{\lambda_y} \\ G_{\lambda_x} & G_{\lambda_{xy}} \end{pmatrix} = \sum_{j,k} d_{jk}(\lambda)\phi_j^+(x)\tilde{\psi}_k^+(y) + \sum_k \phi_k^+(x)\tilde{\phi}_k^+(y),$$

where $d_{jk}(\lambda) = \mathcal{O}(\lambda^{-K})$ is a scalar meromorphic function, moreover, $K \leq$ order of vanishing of the Evans function $D(\lambda)$ at $\lambda = 0$.

Proof. See [37], Proposition 7.1 for the first statement and Theorem 6.3 for the second statement linking order K of the pole to multiplicity of the zero of the Evans function. \square

LEMMA 4.2. *Assuming (H0)–(H3), for $|\lambda|$ sufficiently small, the eigenvalue equation $(L_+ - \lambda)W = 0$ associated with the limiting, constant-coefficient operator L_+ has a basis of $2n$ solutions $\bar{W}_j^+ = e^{\mu_j^+(\lambda)x}V_j^+(\lambda)$ where μ_j^+ and V_j^+ are analytic in λ , consisting of n fast modes*

$$(4.3) \quad \begin{aligned} \mu_j^+ &= \gamma_j^+ + \mathcal{O}(\lambda), \\ V_j^+ &= S_j^+ + \mathcal{O}(\lambda), \end{aligned}$$

where γ_j^+, S_j^+ are eigenvalues and associated right eigenvectors of $B_+^{-1}A_+$, and n slow modes

$$(4.4) \quad \begin{aligned} \mu_{r+j}^+(\lambda) &:= -\lambda/a_j^+ + \lambda^2\beta_j^+/a_j^{+3} + \mathcal{O}(\lambda^3), \\ V_{r+j}^+(\lambda) &:= r_j^+ + \mathcal{O}(\lambda), \end{aligned}$$

where a_j^+, l_j^+, r_j^+ are eigenvalues and left and right eigenvectors of $A_+ := dF(u_+)$, and $\beta_j^+ := l_j^+ B_+ r_j^+ > 0$ with $B_+ := B(u_+)$. The same is true for the adjoint eigenvalue equation

$$(L^+ - \lambda)^* Z = 0;$$

i.e., it has a basis of solutions

$$\bar{W}_j^+ = e^{-\mu_j^+(\lambda)x} \tilde{V}_j(\lambda)$$

with

$$(4.5) \quad \tilde{V}_j^+(\lambda) = \tilde{T}_j^+ + \mathcal{O}(\lambda),$$

$$(4.6) \quad \tilde{V}_{r+j}^+(\lambda) = l_j^+ + \mathcal{O}(\lambda),$$

\tilde{V}^+ analytic in λ .

Proof. See [27]. \square

PROPOSITION 4.3. Assume (H0)–(H3) and (1.9), then, for $r > 0$ sufficiently small, the resolvent kernel G_λ associated with the linearized evolution equation

$$(4.7) \quad U_t = L_+ U := -A_+ U_x + B_+ U_{xx}$$

satisfies, for $0 \leq y \leq x$:

$$(4.8) \quad \begin{aligned} |\partial_x^\gamma \partial_y^\alpha G_\lambda(x, t; y)| \leq C & \left(|\lambda|^\gamma + e^{-\theta|x|} \right) \left(|\lambda|^\alpha + e^{-\theta|y|} \right) \left(\sum_{a_k^+ > 0} \left| e^{(-\lambda/a_k^+ + \lambda^2\beta_k^+/a_k^{+3})(x-y)} \right| \right. \\ & \left. + \sum_{a_k^+ < 0, a_j^+ > 0} \left| e^{(-\lambda/a_j^+ + \lambda^2\beta_j^+/a_j^{+3})x + (\lambda/a_k^+ - \lambda^2\beta_k^+/a_k^{+3})y} \right| \right), \end{aligned}$$

$0 \leq |\alpha|, |\gamma| \leq 1, \theta > 0$, with similar bounds for $0 \leq x \leq y$. Moreover, each term in the summation on the right-hand side of (4.8) bounds a separately analytic function.

Proof. By 1.8 D does not vanish on $Re(\lambda) \geq 0$, hence, by continuity, on $|\lambda| \leq r$. Thus, according to (4.2), all $|d_{jk}(\lambda)|$ are uniformly bounded on $|\lambda| \leq r$, and so it is enough to find estimates for fast and slow modes $\phi_j^+, \tilde{\phi}_j^+, \psi_j^+$, and $\tilde{\psi}_j^+$. By using (3.5) we find:

$$(4.9) \quad \begin{pmatrix} \phi_j^+ \\ \partial_x \phi_j^+ \end{pmatrix} = e^{\mu_j(\lambda)x} P^+ \begin{pmatrix} V_j \\ \mu_j V_j \end{pmatrix} = e^{\mu_j(\lambda)x} (I + \Theta) \begin{pmatrix} V_j \\ \mu_j V_j \end{pmatrix}$$

and similarly for $\tilde{\phi}_j^+, \psi_j^+$, and $\tilde{\psi}_j^+$. Now using (2.4) and the fact, by (4.4), that $e^{\mu_j(\lambda)x}$ is of order $e^{-|\theta x|}$ for fast modes and order $e^{-\lambda/a_j^+ + \lambda^2\beta_j^+/a_j^{+3} + \mathcal{O}(\lambda^3)}$ for slow modes, substituting this and corresponding dual estimates in (4.9) and grouping terms, we obtain the result. \square

5. High frequency bounds. To analyze the high frequency behavior of the Green function of the boundary layer, we first establish some bounds for the projection terms in the Green function, using the symmetric formula

$$(5.1) \quad \begin{pmatrix} G_\lambda(x, y) & \partial_y G_\lambda(x, y) \\ \partial_x G_\lambda(x, y) & \partial_x \partial_y G_\lambda(x, y) \end{pmatrix} = \begin{cases} \mathcal{F}^{y \rightarrow x} \Pi_+(x) \mathcal{S}^{-1}(y) & \text{if } x > y, \\ \tilde{\mathcal{F}}^{x \rightarrow y} \tilde{\Pi}_+(x) \mathcal{S}^{-1}(x) & \text{if } x < y. \end{cases}$$

By setting $\bar{x} = |\lambda^{\frac{1}{2}}|x$, $\bar{\lambda} = \frac{\lambda}{|\lambda|}$, $\bar{B}(\bar{x}) = B(\frac{\bar{x}}{\lambda^{\frac{1}{2}}})$, $\bar{w}(\bar{x}) = w(\frac{x}{\lambda^{\frac{1}{2}}})$ in the eigenvalue equation $Lw = \lambda w$ associated with (1.6) we obtain

$$(5.2) \quad \bar{W}' = \mathbb{B}\bar{W} + \mathcal{O}\left(|\lambda^{-\frac{1}{2}}|\right) \bar{W}$$

where

$$(5.3) \quad \mathbb{B} = \begin{pmatrix} 0 & I \\ \bar{\lambda}\bar{B} & 0 \end{pmatrix}$$

and $\mathbb{B}' = \mathcal{O}(|\lambda^{-\frac{1}{2}}|)$ and $|\bar{\lambda}| = 1$. Since $\mathbb{B}(\lambda, \bar{x})$ varies within a compact set, then there are C^1 eigenprojections P_0 and P_+ with property $|P'_+| = \mathcal{O}(|\lambda^{-\frac{1}{2}}|)$ and $|P'_0| = \mathcal{O}(|\lambda^{-\frac{1}{2}}|)$ taking \bar{W} onto the stable and unstable subspaces. By using the two new coordinates $Y_+ = P_+ \bar{W}$ and $Y_0 = P_0 \bar{W}$, we obtain

$$(5.4) \quad \begin{pmatrix} Y_+ \\ Y_0 \end{pmatrix}' = \begin{pmatrix} A_+ & 0 \\ 0 & A_0 \end{pmatrix} \begin{pmatrix} Y_+ \\ Y_0 \end{pmatrix} + \mathcal{O}\left(|\lambda^{-\frac{1}{2}}|\right) \begin{pmatrix} y \\ y \end{pmatrix}.$$

Equivalently, we can find continuous invertible transformations Q_+ , Q_0 such that $E_+ = Q_+ A_+ Q_+^{-1}$ and $E_0 = Q_0 A_0 Q_0^{-1}$, where

$$(5.5) \quad \text{Re}(E) := \frac{1}{2}(E_+ + E_+^*) < -\beta^{-\frac{1}{2}}I$$

in the sense of quadratic forms.

Again, by coordinate change $Z_+ = Q_+ Y_+$, $Z_0 = Q_0 Y_0$, we find

$$(5.6) \quad \begin{pmatrix} z_+ \\ z_0 \end{pmatrix}' = \begin{pmatrix} E_+ & 0 \\ 0 & E_0 \end{pmatrix} \begin{pmatrix} z_+ \\ z_0 \end{pmatrix} + \mathcal{O}\left(|\lambda^{-\frac{1}{2}}|\right) \begin{pmatrix} z \\ z \end{pmatrix},$$

where

$$(5.7) \quad \frac{|\bar{w}|}{C} \leq |z| \leq C|\bar{w}|.$$

From this we find by energy estimate that

$$(5.8) \quad (|z_+|^2)' < -2\beta^{-\frac{1}{2}}|z_+|^2$$

and, hence,

$$(5.9) \quad \frac{|z_+(x)|}{|z_+(y)|} \leq e^{-\tilde{\beta}^{-\frac{1}{2}}|x-y|}$$

for any solution z_+ decaying at ∞ , where $\tilde{\beta} < \beta$ and, thus,

$$(5.10) \quad \frac{|z(x)|}{|z(y)|} \leq e^{-\beta^{-\frac{1}{2}}|x-y|}$$

for $x > y$, provided that $|\lambda|$ is sufficiently large. From this we obtain

$$(5.11) \quad \frac{|\tilde{W}(x)|}{|\tilde{W}(y)|} \leq C^2 e^{-\beta^{-\frac{1}{2}}|x-y|}$$

where C is as in (5.7).

Applying a symmetric argument for the adjoint equation, we obtain the following lemma.

LEMMA 5.1. *On the manifolds Φ_+ and $\tilde{\Psi}_+$ defined in (3.5) and (3.17), for λ sufficiently large, within the sector $\Omega_\theta = \{\lambda : \text{Re}(\lambda) \geq -\theta_1 - \theta_2 |\text{Im}(\lambda)|\}$, $\theta_1, \theta_2 > 0$, we have in rescaled coordinates \bar{x} , for some uniform $C > 0$,*

$$(5.12) \quad |\mathcal{F}^{y \rightarrow x}|, |\tilde{\mathcal{F}}^{x \rightarrow y}| \leq C e^{-\frac{|y-x|}{c}}$$

for $x > y$ and $x < y$, respectively.

LEMMA 5.2. *In rescaled coordinates \bar{x} , $\bar{\lambda}$, for the projection terms $\Pi_+(y)$ and $\tilde{\Pi}_+(x)$, the projection along Φ_0 onto Φ_+ , for λ sufficiently large,*

$$(5.13) \quad |\Pi_+(y)|, |\tilde{\Pi}_+(x)| < C$$

for some uniform $C > 0$.

Proof. Choosing the coordinates $\begin{pmatrix} W_1 \\ W_2 \end{pmatrix} \in \mathbb{C}^{2n}$ where $W^j = \begin{pmatrix} W_1^j \\ W_2^j \end{pmatrix}$, we show for small enough ϵ and fixed $c > 0$ such that $\frac{|W_2^j|}{|W_1^j|} \leq c\epsilon$ and $\frac{|W_1^j|}{|W_2^j|} < c$ the projection along $E = \text{span}(W^1, \dots, W^n)$ onto $F = \text{span}(W^{n+1}, \dots, W^{2n})$

$$(5.14) \quad \Pi := (w^1, \dots, w^{2n}) (O_n, I_n) (w^1, \dots, w^{2n})^{-1}$$

satisfies

$$(5.15) \quad |\Pi| \leq C_2(c, \epsilon).$$

To show this without loss of generality we assume that

$$(5.16) \quad (w^{n+1}, \dots, w^{2n}) = \begin{pmatrix} I_n \\ \mathcal{O}(\epsilon) \end{pmatrix} (w^1, \dots, w^n) = \begin{pmatrix} \mathcal{O}(1) \\ I_n \end{pmatrix}.$$

Now it is sufficient to show that

$$(5.17) \quad |(w^1, \dots, w^{2n})| \leq C_2(c, \epsilon).$$

However, this amounts to showing that

$$(5.18) \quad \left| \begin{pmatrix} M & I_n \\ I_n & O \end{pmatrix} + \mathcal{O}(\epsilon)^{-1} \right| \leq C,$$

which amounts to showing that

$$(5.19) \quad \left| \begin{pmatrix} M & I_n \\ I_n & O \end{pmatrix}^{-1} \right| \leq C_2(c),$$

where $|M| \leq c$. However, this is easy to show because

$$\begin{pmatrix} M & I_n \\ I_n & O \end{pmatrix}^{-1} = \begin{pmatrix} I_n & -M \\ O & I_n \end{pmatrix},$$

and so

$$\left| \begin{pmatrix} M & I_n \\ I_n & O \end{pmatrix}^{-1} \right| \leq 1 + |M| \leq C. \quad \square$$

PROPOSITION 5.3. *Assume (H0)–(H3) and (1.9). Then, for $R > 0$ sufficiently large, the resolvent kernel G_λ associated with the linearized evolution equation (4.7) satisfies, for $c, C > 0$ and $0 \leq |\alpha|, |\gamma| \leq 1$:*

$$(5.20) \quad |\partial_x^\gamma \partial_y^\alpha G_\lambda(x, y)| \leq C |\lambda|^{\left(\frac{|\alpha|+|\gamma|-1}{2}\right)} e^{-\sqrt{\lambda} \frac{|y-x|}{c}}.$$

Proof. Recalling the coordinate-free representation (5.1) and combining with (5.12) and (5.13), we find that the Green function \bar{G}_λ in rescaled coordinates $\bar{x}, \bar{\lambda}$ satisfies

$$(5.21) \quad |\partial_x^\gamma \partial_y^\alpha \bar{G}_\lambda(\bar{x}, \bar{y})| \leq C e^{-\frac{|\bar{y}-\bar{x}|}{c}},$$

whence, (5.20) follows in the original coordinates. \square

Remark 5.4. The argument of Lemma 5.2 is the key new ingredient in the resolvent estimates for the boundary layer case as compared to the analysis on the whole line carried out for viscous shock layers in [37], making essential use of compatibility of the boundary condition with high-frequency behavior. On the whole line, there is no such requirement and high-frequency stability is automatic.

6. Pointwise Green function bounds. With the pointwise bounds established on the resolvent kernel G_λ , we obtain pointwise bounds on the Green function through the inverse Laplace transform formula by a simplified version of the stationary-phase arguments used in [37] for the shock case, repeated here for completeness.

Proof of Theorem 1.4. By sectoriality of L , we have the inverse Laplace transform representation (see [37]):

$$(6.1) \quad G(t; x, y) = \int_\Gamma e^{\lambda t} G_\lambda(x, y) d\lambda.$$

Let $\theta_1 > 0, \theta_2 > 0$ be chosen sufficiently small, in particular so small as to satisfy the hypotheses of all previous assertions. By assumption (1.9), the large- $|\lambda|$ bounds on the resolvent kernel, and analyticity of the Evans function $D_L(\lambda)$, it follows that G_λ has finitely many poles in Ω_θ (corresponding to roots of D_L), each with strictly negative real part. Choosing θ_1, θ_2 still smaller, if necessary, we can, thus, arrange that G_λ is analytic on Ω_θ . It follows from Cauchy’s Theorem that

$$(6.2) \quad G(x, t; y) = \int_\Gamma e^{\lambda t} G_\lambda(x, y) d\lambda$$

for any contour Γ that can be expressed as $\Gamma = \partial(\Omega_\theta \setminus \mathcal{S})$ for $\mathcal{S} \subset \mathbb{C}$ open.

Case I. $|x - y|/t$ large. We first treat the trivial case that $|x - y|/t \geq S, S$ sufficiently large, the regime in which standard short-time parabolic theory applies. Set

$$(6.3) \quad \bar{\alpha} := \frac{|x - y|}{2\beta t}, \quad R := \beta \bar{\alpha}^2,$$

where β is as in (5.5), and consider again the representation of G , that is

$$(6.4) \quad G(x, t; y) = \int_{\Gamma_1 \cup \Gamma_2} e^{\lambda t} G_\lambda(x, y) d\lambda,$$

where $\Gamma_1 := \partial B(0, R) \cap \bar{\Omega}_\theta$ and $\Gamma_2 := \partial\Omega_\theta \setminus B(0, R)$. Note that the intersection of Γ with the real axis is $\lambda_{min} = R = \beta\bar{\alpha}^2$. \square

By the large $|\lambda|$ estimates of Proposition 5.3, we have for all $\lambda \in \Gamma_1 \cup \Gamma_2$ that

$$|G_\lambda(x, y)| \leq C \frac{e^{-\sqrt{|\lambda|} \frac{|y-x|}{c}}}{\sqrt{|\lambda|}}.$$

Further, we have

$$(6.5) \quad Re\lambda \leq R(1 - \eta\omega^2), \quad \lambda \in \Gamma_1, Re\lambda \leq Re\lambda_0 - \eta(|Im\lambda| - |Im\lambda_0|), \quad \lambda \in \Gamma_2$$

for R sufficiently large, where ω is the argument of λ and λ_0 and λ_0^* are the two points of intersection of Γ_1 and Γ_2 , for some $\eta > 0$ independent of $\bar{\alpha}$. Combining these estimates, we obtain

$$(6.6) \quad \begin{aligned} \left| \int_{\Gamma_1} e^{\lambda t} G_\lambda d\lambda \right| &\leq \int_{\Gamma_1} C|\lambda|^{-\frac{1}{2}} e^{Re\lambda t - \beta^{-\frac{1}{2}}|\lambda|^{-\frac{1}{2}}|x-y|} d\lambda \\ &\leq C e^{-\beta\bar{\alpha}^2 t} \int_{-L}^{+L} R^{-\frac{1}{2}} e^{-\beta R\eta\omega^2 t} R d\omega \leq C t^{-\frac{1}{2}} e^{-\beta\bar{\alpha}^2 t}. \end{aligned}$$

Likewise,

$$(6.7) \quad \begin{aligned} \left| \int_{\Gamma_2} e^{\lambda t} G_\lambda d\lambda \right| &\leq \int_{\Gamma_2} C|\lambda|^{-\frac{1}{2}} C e^{Re\lambda t - \beta^{-\frac{1}{2}}|\lambda|^{-\frac{1}{2}}|x-y|} d\lambda \\ &\leq C e^{Re(\lambda_0)t - |\beta|^{-\frac{1}{2}}|\lambda_0|^{-\frac{1}{2}}|x-y|} \int_{\Gamma_2} |\lambda|^{-\frac{1}{2}} e^{(Re\lambda - Re\lambda_0)t} |d\lambda| \\ &\leq C e^{-\beta\bar{\alpha}^2 t} \int_{\Gamma_2} |Im\lambda|^{-\frac{1}{2}} e^{-\eta|Im\lambda - Im\lambda_0|t} |dIm\lambda| \\ &\leq C t^{-\frac{1}{2}} e^{-\beta\bar{\alpha}^2 t}. \end{aligned}$$

Combining these last two estimates, we have

$$(6.8) \quad |G(x, t; y)| \leq C t^{-\frac{1}{2}} e^{-\frac{\beta\bar{\alpha}^2 t}{2}} e^{-\frac{(x-y)^2}{8\beta t}} \leq C t^{-\frac{1}{2}} e^{-\eta t} e^{-\frac{(x-y)^2}{8\beta t}},$$

for $\eta > 0$ independent of $\bar{\alpha}$. Observing that $\frac{|x-at|}{2t} \leq \frac{|x-y|}{t} \leq \frac{2|x-at|}{t}$ for any bounded a , for $\frac{|x-y|}{t}$ sufficiently large, we find that this contribution may be absorbed in any summand $t^{-\frac{1}{2}} e^{-\frac{(x-y-a_k^+ t)^2}{Mt}}$.

Case II. $|x - y|/t$ bounded. We now turn to the critical case that $|x - y|/t \leq S$. A few remarks are in order at the outset. Our goal is to bound $|G|$ by terms of form $Ct^{-1/2}e^{-\bar{\alpha}^2 t/M}$, where $\bar{\alpha} := (x - a_j^+(t - |y/a_k^+|))/2t$ or $\bar{\alpha} := (x - y - a_k^+ t)/2t$ are now uniformly bounded, by

$$(6.9) \quad |x - y|/2t + \max_j \{ |a_j^+| \} / 2 \leq S/2 + \max |a_j^+| / 2.$$

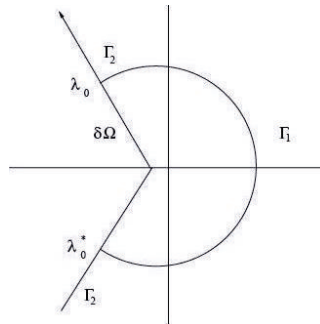


FIG. 1.

Thus, in particular, contributions of order $t^{-1/2}e^{-\eta t}$, $\eta > 0$, can be absorbed in any summand $t^{-1/2}e^{-(x-y-a_k^+t)^2/Mt}$ if we take M sufficiently large. Likewise, for G_x and G_y , contributions of order $t^{-1}e^{-\eta t}$ can be absorbed. We will use this observation repeatedly.

In contrast to the previous case of large characteristic speed $|x - y|/t \geq S$, we are not trying to show rapid time-exponential decay. Rather, we are trying to show that the rate of exponential decay of the solution does not degrade too rapidly as $\bar{\alpha} \rightarrow 0$: precisely, that it vanishes to order $\bar{\alpha}^2$ and no more. Thus, the crucial part of our analysis will be for small $\bar{\alpha}$. All other situations can be estimated crudely as described just above.

Let r be sufficiently small that the small- $|\lambda|$ bounds hold on $B(0, r)$. Next, choose θ_1 and θ_2 still smaller than before, if necessary, so that $\Omega_\theta \setminus B(0, r) \subset \Lambda$. This implies that $\partial\Omega_\theta \cap B(0, r) \neq \emptyset$, giving the configuration pictured in Figure 1. Similarly as in the previous case, define $\Gamma = \Gamma_1 \cup \Gamma_2$, where Γ_1 is the portion of the circle $\partial B(0, r)$ contained in $\bar{\Omega}_\theta$, and Γ_2 is the portion of $\partial\Omega_\theta$ outside $B(0, r)$.

$$(6.10) \quad G(x, t; y) = \int_{\Gamma_1} e^{\lambda t} G_\lambda(x, y) d\lambda + \int_{\Gamma_2} e^{\lambda t} G_\lambda(x, y) d\lambda.$$

We separately estimate the terms \int_{Γ_1} and \int_{Γ_2} .

Large- and medium- λ estimates. The \int_{Γ_2} term is straightforward. The points λ_0, λ_0^* where Γ_1 meets Γ_2 satisfy $Re(\lambda_0) = -\eta < 0$. Moreover, combining the results low-frequency case, we have the bound $|G_\lambda| \leq C|\lambda|^{-\frac{1}{2}}$ for $\lambda \in \Gamma_2$. Thus, we have

$$(6.11) \quad \left| \int_{\Gamma_2} e^{\lambda t} G_\lambda d\lambda \right| \leq C e^{-Re \lambda_0 t} \int_{\Gamma_2} |Im \lambda|^{-\frac{1}{2}} e^{-\eta |Im \lambda - Im \lambda_0| t} |d Im \lambda| \leq C t^{-\frac{1}{2}} e^{-\eta t}.$$

This contribution can be absorbed as described above. An analogous computation using $|G_{\lambda_x}|, |G_{\lambda_x}| \leq C|\lambda|^{-1}$ shows that the Γ_2 contribution to G_x and G_y is $O(t^{-1}e^{-\eta t})$, and can likewise be absorbed.

Small $|\lambda|$ estimates. It remains to estimate the critical term $\int_{\Gamma_1} e^{\lambda t} G_\lambda d\lambda$. This we will estimate in different ways, depending on the size of t .

Bounded time. For t bounded, we can use the medium- λ bounds $|G_\lambda|, |G_{\lambda_x}|, |G_{\lambda_y}| \leq C$ to obtain $|\int_{\Gamma_1} e^{\lambda t} G_\lambda d\lambda| \leq C_2 |\Gamma_1|$. This contribution is order $Ce^{-\eta t}$ for bounded time, and hence can be absorbed.

Large time. For t large, we must instead estimate $\int_{\Gamma_1} e^{\lambda t} G_\lambda d\lambda$ using the small- $|\lambda|$ expansions. First, observe that all coefficient functions $d_{jk}(\lambda)$ are uniformly bounded (since $|\lambda|$ is bounded in this case).

Expanding $G = \int_{\Gamma} e^{\lambda t} G_\lambda(x, y) d\lambda$ as

$$\begin{pmatrix} G & G_x \\ G_y & G_{xy} \end{pmatrix} = \int_{\Gamma} e^{\lambda t} \begin{pmatrix} G_\lambda & G_{\lambda_x} \\ G_{\lambda_y} & G_{\lambda_{xy}} \end{pmatrix} d\lambda,$$

we estimate the \int_{Γ_1} contributions to G , G_x , and G_y simultaneously.

Case II(i). ($0 < y < x$). By our low-frequency estimates, we have

$$(6.12) \quad \int_{\Gamma} e^{\lambda t} \begin{pmatrix} G_\lambda & G_{\lambda_x} \\ G_{\lambda_y} & G_{\lambda_{xy}} \end{pmatrix} d\lambda = \int_{\Gamma} \sum_{j,k} e^{\lambda t} \phi_j^+(x) d_{jk} \tilde{\psi}_k^+(y) d\lambda + \int_{\Gamma} \sum_{j,k} e^{\lambda t} \psi_k^+(x) \tilde{\psi}_k^+(y) d\lambda,$$

where each d_{jk} is analytic and, hence, bounded. We estimate separately each of the terms

$$\int_{\Gamma_1} e^{\lambda t} \phi_j^+(x) d_{jk} \tilde{\psi}_k^+(y) d\lambda$$

on the right-hand side of (6.12). Estimates for terms

$$\int_{\Gamma} \sum_{j,k} e^{\lambda t} \psi_k^+(x) \tilde{\psi}_k^+(y) d\lambda$$

go similarly.

Case II(ia). First, consider the critical case $a_j^+ > 0$, $a_k^+ < 0$. For this case,

$$\left| \phi_{j(x)}^+ d_{jk} \tilde{\psi}_k^+(y) \right| \leq C e^{Re(\rho_j^+ x - \nu_k^+ y)},$$

where

$$\begin{cases} \nu_k^+(\lambda) = -\lambda/a_k^+ + \lambda^2 \beta_k^+ / (a_k^+)^3 + \mathcal{O}(\lambda^3) \\ \rho_j^+(\lambda) = -\lambda/a_j^+ + \lambda^2 \beta_j^+ / (a_j^+)^3 + \mathcal{O}(\lambda^3). \end{cases}$$

Set

$$\bar{\alpha} = \frac{a_k^+ x/a_j^+ - y - a_k^+ t}{2t}, \quad p := \frac{\beta_j^+ a_k^+ x / (a_j^+)^3 - \beta_k^+ y / (a_k^+)^2}{t} > 0.$$

Define Γ'_{1a} to be the portion contained in Ω_θ of the hyperbola

(6.13)

$$\begin{aligned} & Re(\rho_j^+ x - \nu_k^+ y) + \mathcal{O}(\lambda^3)(|x| + |y|) \\ &= (1/a_k^+) Re \left[\lambda(-a_k^+ x/a_j^+ + y) + \lambda^2 \left(x\beta_j^+ a_k^+ / (a_j^+)^3 - y\beta_k^- / (a_k^+)^2 \right) \right] \\ &\equiv \text{constant} \\ &= (1/a_k^-) \left[\lambda_{min}(-a_k^- x/a_j^+ + y) + \lambda_{min}^2 \left(x\beta_j^+ a_k^+ / (a_j^+)^3 - y\beta_k^+ / (a_k^+)^2 \right) \right], \end{aligned}$$

where

$$(6.14) \quad \lambda_{min} := \begin{cases} \frac{\bar{\alpha}}{p} & \text{if } \left| \frac{\bar{\alpha}}{p} \right| \leq \epsilon, \\ \pm\epsilon & \text{if } \frac{\bar{\alpha}}{p} \gtrless \epsilon. \end{cases}$$

Denoting by λ_1, λ_1^* the intersections of this hyperbola with $\partial\Omega_\theta$, define Γ'_{1_b} to be the union of $\lambda_1\lambda_0$ and $\lambda_0^*\lambda_1^*$, and define $\Gamma'_1 = \Gamma'_{1_a} \cup \Gamma'_{1_b}$. Note that $\lambda = \bar{\alpha}/p$ minimizes the left-hand side of (6.13) for λ real. Note also that p is bounded for $\bar{\alpha}$ sufficiently small, since $\bar{\alpha} \leq \epsilon$ implies that

$$(|a_k^+x/a_j^+| + |y|)/t \leq 2|a_k^+| + 2\epsilon;$$

i.e., $(|x| + |y|)/t$ is controlled by $\bar{\alpha}$.

With these definitions, we readily obtain that

$$(6.15) \quad \begin{aligned} \operatorname{Re}(\lambda t + \rho_j^+x - \nu_k^+y) &\leq -(t/a_k^-) (\bar{\alpha}^2/4p) - \eta \operatorname{Im}(\lambda)^2t \\ &\leq -\bar{\alpha}^2t/M - \eta \operatorname{Im}(\lambda)^2t, \end{aligned}$$

for $\lambda \in \Gamma'_{1_a}$ (note: here, we have used the crucial fact that $\bar{\alpha}$ controls $(|x| + |y|)/t$, in bounding the error term $\mathcal{O}(\lambda^3)(|x| + |y|)/t$ arising from expansion). Likewise, we obtain for any q that

$$(6.16) \quad \int_{\Gamma'_{1_a}} |\lambda|^q e^{\operatorname{Re}(\lambda t + \rho_j^+x - \nu_k^-y)} d\lambda \leq C t^{-\frac{1}{2} - \frac{q}{2}} e^{-\bar{\alpha}^2t/M},$$

for suitably large $C, M > 0$ (depending on q). Observing that

$$\bar{\alpha} = (a_k^+/a_j^+) (x - a_j^+ (t - |y/a_k^+|)) / 2t,$$

we find that the contribution of (6.16) can be absorbed in the described bounds for $t \geq |y/a_k^-|$. At the same time, we find that $\bar{\alpha} \geq x > 0$ for $t \leq |y/a_k^+|$, whence,

$$\bar{\alpha} \geq (x - y - a_j^+t) / Mt + |x|/M,$$

for some $\epsilon > 0$ sufficiently small and $M > 0$ sufficiently large.

This gives

$$e^{-\bar{\alpha}^2/p} \leq e^{-(x-y-a_k^+t)^2/Mt} e^{-\eta|x|}$$

provided $|x|/t > a_j^+$, a contribution which can again be absorbed. On the other hand, if $t \leq |x/a_j^+|$, we can use the dual estimate

$$(6.17) \quad \begin{aligned} \bar{\alpha} &= (-y - a_k^+ (t - |x/a_j^+|)) / 2t \\ &\geq (x - y - a_k^+t) / Mt + |y|/M, \end{aligned}$$

together with $|y| \geq |a_k^-t|$, to obtain

$$e^{-\bar{\alpha}^2/p} \leq e^{-(x-y-a_j^+t)^2/Mt} e^{-\eta|y|},$$

a contribution that can likewise be absorbed.

Case II(ib). In case $a_j^+ < 0$ or $a_k^+ > 0$, terms $|\varphi_j^+| \leq C e^{-\eta|x|}$ and $|\tilde{\psi}_j^+| \leq C e^{-\eta|y|}$ are strictly smaller than those already treated in Case II(ia), so may be absorbed in previous terms.

Case II(ii) ($0 < x < y$). The case $0 < x < y$ can be treated very similarly to the previous one; see [37] for details. This completes the proof of Case II, and the theorem.

7. Nonlinear analysis. Introducing the perturbation variable

$$(7.1) \quad u(x, t) := \tilde{u}(x, t) - \bar{u}(x),$$

we obtain

$$(7.2) \quad u_t - Lu = Q(u)_x,$$

where the second-order Taylor remainder satisfies

$$(7.3) \quad Q(u) := f(\bar{u} + u) - f(\bar{u}) - df(\bar{u})u = \mathcal{O}(|u|^2)$$

so long as $|u|$ remains bounded.

LEMMA 7.1 (integral formulation). *Under the assumptions of Theorem 1.6, there exists a classical solution of (7.2) for $0 < t \leq T$, $T > 0$, continuous in $L^\infty(x)$ at $t = 0$, extending for all $t > 0$ such that $u(\cdot, t)$ remains sufficiently small in $L^1 \cap L^\infty$, given by*

$$(7.4) \quad \begin{aligned} u(x, t) = & \int_0^\infty G(x, t; y)g(y) dy + \int_0^t G_y(x, t - s; 0)Bh(s) ds \\ & - \int_0^t \int_0^\infty G_y(x, t - s; y)Q(u)(y, s) dy ds. \end{aligned}$$

Proof. From Lemma 3.2 and the inverse Laplace representation (6.1) we find that $G(x, t - s; y)$ considered as a function of y, s satisfies the adjoint equation

$$(7.5) \quad (\partial_s - L_y)^* G^*(x, t - \cdot; \cdot) = 0,$$

or

$$(7.6) \quad -G_s - (GA)_y + GA_y = (G_y B)_y.$$

Likewise, reviewing the construction of the resolvent, we find $G_\lambda(x, 0) \equiv 0$, yielding

$$(7.7) \quad G(x, t - s; 0) \equiv 0.$$

That is, $G^*(x, t - \cdot; \cdot)$ is the Green function for the adjoint equation, as may alternatively be seen directly by a duality argument analogous to the proof of Lemma 3.2.

Thus, integrating G against (7.2), integrating by parts, and using the fact that $G = 0$ and $u = h$ on the boundary $y = 0$, we obtain for any classical solution of (7.2) that

$$(7.8) \quad \begin{aligned} & \int_0^t \int_0^\infty G(x, t - s; y)Q(u(y, s))_y dy ds = \\ & \int_0^t \int_0^\infty G(x, t - s; y)(\partial_s - L_y)u(y, s) dy ds \\ & = \int_0^t \int_0^\infty ((\partial_s - L_y)^* G^*)^*(x, t - s; y)u(y, s) dy ds \\ & + u(x, t) - \int_0^\infty G(x, t; y)g(y) dy - \int_0^t G_y(x, t - s; 0)Bh(s) ds, \end{aligned}$$

from which we obtain (7.4) by rearranging and integrating by parts the term $\int_0^t \int_0^\infty G(x, t - s; y) Q(u(y, s))_y dy ds$.

Indeed, (7.4) may be taken as the definition of a weak solution in $L^\infty(x, t)$. (One can see using convolution identities that this agrees with the usual definition in terms of integration against test functions $\phi \in C_0^\infty(\mathbb{R} \times \mathbb{R})$.) Existence of weak solutions can be obtained by a standard contraction mapping/continuation argument using the convolution bounds of Lemmas 7.2–7.4 below; we omit the details, since we shall carry out quite similar but more difficult estimates in the proof of stability. Smoothness of solutions may then be obtained by a bootstrapping argument as sketched in Appendix A. \square

To establish stability, we use the following lemmas proved in [16].

LEMMA 7.2 (linear estimates [16]). *Under the assumptions of Theorem 1.6,*

$$(7.9) \quad \int_0^{+\infty} |G(x, t; y)|(1 + |y|)^{-3/2} dy \leq C(\theta + \psi_1 + \psi_2)(x, t),$$

for $0 \leq t \leq +\infty$, some $C > 0$.

LEMMA 7.3 (nonlinear estimates [16]). *Under the assumptions of Theorem 1.6,*

$$(7.10) \quad \int_0^t \int_0^{+\infty} |G_y(x, t - s; y)|\Psi(y, s) dy ds \leq C(\theta + \psi_1 + \psi_2)(x, t),$$

for $0 \leq t \leq +\infty$, some $C > 0$, where

$$(7.11) \quad \begin{aligned} \Psi(y, s) := & (1 + s)^{1/2} s^{-1/2} (\theta + \psi_1 + \psi_2)^2(y, s) \\ & + (1 + s)^{-1} (\theta + \psi_1 + \psi_2)(y, s). \end{aligned}$$

We require also the following estimate accounting boundary effects.

LEMMA 7.4 (boundary estimate). *Under the assumptions of Theorem 1.6,*

$$(7.12) \quad \left| \int_0^t G_y(x, t - s; 0) B h(s) ds \right| \leq C E_0 (\theta + \psi_1 + \psi_2)(x, t),$$

for $0 \leq t \leq +\infty$, some $C > 0$.

Proof. The estimate on \int_0^{t-1} , where $G_y(x, t - s; 0)$ is nonsingular, follows readily by estimates similar to but somewhat simpler than those of Lemma (7.3), which we, therefore, omit.

To bound the singular part \int_{t-1}^t , we integrate (7.6) in y from 0 to $+\infty$, recalling that $G(x, t - s; 0) \equiv 0$, to obtain

$$(7.13) \quad G_y B = - \int_0^{+\infty} A_y(y) G(x, t - s; y) dy - \int_0^{+\infty} G_s(x, t - s; y) dy.$$

Substituting in the left-hand side of (7.12), and integrating by parts in s , we obtain

$$(7.14) \quad \begin{aligned} \int_{t-1}^t G_y B h(s) ds = & \int_0^1 \left(\int_0^{+\infty} A_y(y) G(x, \tau; y) dy \right) h(t - \tau) d\tau \\ & - \int_0^1 \left(\int_0^{+\infty} G(x, \tau; y) dy \right) h'(t - \tau) d\tau \\ & + \left(\int_0^{+\infty} G(x, 1; y) dy \right) h(t - 1), \end{aligned}$$

which, by $\int |G| dy \leq C$, has norm bounded by $\max_{0 \leq \tau \leq 1} (|h| + |h'|)(t - \tau)$.

Combining this with the more straightforward estimate

$$\begin{aligned}
 \left| \int_{t-1}^t G_y(x, t; 0) Bh(s) ds \right| &\leq \int_0^1 |G_y(x, \tau; 0)| Bh(s) ds \\
 &\leq C \max_{0 \leq \tau \leq 1} |h(t - \tau)| \int_0^1 \tau^{-1} e^{-|x|^2/C\tau} d\tau \\
 (7.15) \qquad &= C|x|^{-2} \max_{0 \leq \tau \leq 1} |h(t - \tau)| \\
 &\quad \times \int_0^1 (|x|^2/\tau) e^{-|x|^2/C\tau} d\tau \\
 &\leq C \max_{0 \leq \tau \leq 1} |h(t - \tau)| |x|^{-2},
 \end{aligned}$$

we find that the contribution from \int_{t-1}^t has norm bounded by

$$\max_{0 \leq \tau \leq 1} (|h| + |h'|)(t - \tau)(1 + |x|)^{-2}.$$

Combining this estimate with the one for \int_0^{t-1} , we obtain (7.12). \square

With these preparations, the proof of stability is straightforward.

Proof of Theorem 1.6. Define

$$(7.16) \qquad \zeta(t) := \sup_{y, 0 \leq s \leq t} |u|(\theta + \psi_1 + \psi_2)^{-1}(y, t).$$

We will establish:

Claim. For all $t \geq 0$ for which a solution exists with ζ uniformly bounded by some fixed, sufficiently small constant, there holds

$$(7.17) \qquad \zeta(t) \leq C_2 (E_0 + \zeta(t)^2).$$

From this result, provided $E_0 < 1/4C_2^2$, we have that $\zeta(t) \leq 2C_2E_0$ implies $\zeta(t) < 2C_2E_0$, and so we may conclude by continuous induction that

$$(7.18) \qquad \zeta(t) < 2C_2E_0$$

for all $t \geq 0$. (By Lemma 7.1 and standard short-time estimates, $u \in C^0(x)$ exists and ζ remains continuous so long as ζ remains bounded by some uniform constant; hence (7.18) is an open condition.) From (7.18) and the definition of ζ in (7.16) we then obtain the bounds of (1.16). Thus, it remains only to establish the claim above.

Proof of Claim. We must show that $u(\theta + \psi_1 + \psi_2)^{-1}$ is bounded by $C(E_0 + \zeta(t)^2)$, for some $C > 0$, all $0 \leq s \leq t$, so long as ζ remains sufficiently small. By (7.16), we have for all $t \geq 0$ and some $C > 0$ that

$$(7.19) \qquad |u(x, t)| \leq \zeta(t)(\theta + \psi_1 + \psi_2)(x, t),$$

and, therefore,

$$(7.20) \qquad |Q(u)(y, s)| \leq C\zeta(t)^2\Psi(y, s)$$

with Ψ as defined in (7.11), for $0 \leq s \leq t$. Combining (7.20) with representation (7.4)

and applying Lemmas 7.2–7.4, we obtain

$$\begin{aligned}
 |u(x, t)| &\leq \int_0^\infty \left| \tilde{G}(x, t; y) \right| |g(y)| dy + \left| \int_0^t G_y(x, t - s; 0) Bh(s) ds \right| \\
 &\quad + \int_0^t \int_0^\infty \left| \tilde{G}_y(x, t - s; y) \right| |(Q(u))(y, s)| dy ds \\
 (7.21) \quad &\leq E_0 \int_0^\infty \left| \tilde{G}(x, t; y) \right| (1 + |y|)^{-3/2} dy \\
 &\quad + \left| \int_0^t G_y(x, t - s; 0) Bh(s) ds \right| \\
 &\quad + C\zeta(t)^2 \int_0^t \int_0^\infty \left| \tilde{G}_y(x, t - s; y) \right| \Psi(y, s) dy ds \\
 &\leq C (E_0 + \zeta(t)^2) (\theta + \psi_1 + \psi_2)(x, t).
 \end{aligned}$$

Dividing by $(\theta + \psi_1 + \psi_2)(x, t)$, we obtain (7.17) as claimed. This completes the proof of the claim, and the theorem. \square

Appendix A. Smoothness of solutions.

In this appendix, we briefly sketch the proof that weak solutions defined by (7.4) are necessarily smooth, classical solutions as well, by indicating how to get the necessary derivative bounds.

Time-derivative. Rewriting the second boundary term on the right-hand side of (7.4) using its convolution structure, as

$$\int_0^t G_y(x, \tau; 0) Bh(t - \tau) d\tau,$$

and differentiating in t , we obtain

$$G_y(x, t; 0) Bh(0) + \int_0^t G_y(x, \tau; 0) Bh'(t - \tau) d\tau,$$

for which the first term is bounded and smooth for $x, t > 0$, and the second by the same estimate as in (7.15) is bounded by

$$C|x|^{-2} \int_0^t |h'(t - \tau)| d\tau \leq C|x|^{-2} \log(1 + t).$$

Differentiating the first and third terms, with respect to t and integrating the third term by parts in y , yields

$$\begin{aligned}
 (A.1) \quad &\int_0^\infty G(x, t; y) g(y) dy - \int_{t/2}^t \int_0^\infty G_y(x, t - s; y) Q(u)_s(y, s) dy ds \\
 &\quad - \int_0^{t/2} \int_0^\infty G_{yt}(x, t - s; y) Q(u)(y, s) dy ds,
 \end{aligned}$$

from which, in combination with the boundary estimate already performed, we may readily obtain a short-time bound $|u_t| \leq C|x|^{-2}t^{-1}$ by Picard iteration.

Spatial-derivatives. Likewise, differentiating (7.14) with respect to x , we may bound the x -derivative of the boundary term $\int_0^t G_y B ds$ by

$$C \int_0^t \tau^{-1/2} (|h'| + |h|)(t - \tau) d\tau \leq C \log(1 + t).$$

Differentiating the first and third terms of the right-hand side of (7.4), with respect to x and integrating the third term by parts in y , yields

$$(A.2) \quad \int_0^\infty G(x, t; y)g(y) dy - \int_{t/2}^t \int_0^\infty G_x(x, t-s; y)Q(u)_y(y, s) dy ds \\ - \int_0^{t/2} \int_0^\infty G_{yx}(x, t-s; y)Q(u)(y, s) dy ds,$$

from which, in combination with the boundary estimate already performed, we obtain a short-time bound $|u_x| \leq Ct^{-1/2}$ by Picard iteration. From the bounds on $|u_t|$ and $|u_x|$, finally, we obtain bounds on $|u_{xx}|$ by the equation satisfied by u .

REFERENCES

- [1] J. ALEXANDER, R. GARDNER, AND C.K.R.T. JONES, *A topological invariant arising in the analysis of traveling waves*, J. Reine Angew. Math., 410 (1990), pp. 167–212.
- [2] B. BARKER, J. HUMPHERYS, K. RUDD, AND K. ZUMBRUN, *Stability of viscous shocks in isentropic gas dynamics*, Comm. Math. Phys., 281 (2008), pp. 231–249.
- [3] L.Q. BRIN, *Numerical testing of the stability of viscous shock waves*, Doctoral thesis, Indiana University (1998).
- [4] L.Q. BRIN, *Numerical testing of the stability of viscous shock waves*, Math. Comput., 70 (2001), pp. 1071–1088.
- [5] L. BRIN AND K. ZUMBRUN, *Analytically varying eigenvectors and the stability of viscous shock waves*, Seventh Workshop on Partial Differential Equations, Part I (Rio de Janeiro, 2001). Mat. Contemp., 22 (2002), pp. 19–32.
- [6] T. BRIDGES, G. DERKS, AND G. GOTTWALD, *Stability and instability of solitary waves of the fifth-order KdV equation: A numerical framework*, Phys. D, 172 (2002), pp. 190–216.
- [7] N. COSTANZINO, J. HUMPHERYS, T. NGUYEN, AND K. ZUMBRUN, *Spectral stability of noncharacteristic boundary layers of isentropic Navier–Stokes equations*, Arch. Ration. Mech. Anal., to appear.
- [8] R. GARDNER AND C.K.R.T. JONES, *A stability index for steady state solutions of boundary value problems for parabolic systems*, J. Differential Equations, 91 (1991), pp. 181–203.
- [9] R. GARDNER AND C.K.R.T. JONES, *Traveling waves of a perturbed diffusion equation arising in a phase field model*, Indiana Univ. Math. J., 38 (1989), pp. 1197–1222.
- [10] R. GARDNER AND K. ZUMBRUN, *The Gap Lemma and geometric criteria for instability of viscous shock profiles*, Commun. Pure Appl. Math., 51 (1998), pp. 797–855.
- [11] E. GRENIER AND O. GUES, *Boundary layers for viscous perturbations of noncharacteristic quasilinear hyperbolic problems*, J. Differential Equations, 143 (1998), pp. 110–146.
- [12] E. GRENIER AND F. ROUSSET, *Stability of one-dimensional boundary layers by using Green’s function*, Commun. Pure Appl. Math., 54 (2001), pp. 1343–1385.
- [13] O. GUÈS, G. MÉTIVIER, M. WILLIAMS, AND K. ZUMBRUN, *Multidimensional stability of small-amplitude noncharacteristic boundary layers*, preprint (2007).
- [14] O. GUÈS, G. MÉTIVIER, M. WILLIAMS, AND K. ZUMBRUN, *Viscous boundary value problems for symmetric systems with various multiplicities*, J. Differential Equations, 244 (2008), pp. 309–387.
- [15] D. HENRY, *Geometric theory of semilinear parabolic equations*, Lecture Notes in Mathematics, Springer-Verlag, Berlin, 1981.
- [16] P. HOWARD AND K. ZUMBRUN, *Stability of undercompressive viscous shock waves*, in press, J. Differential Equations, 225 (2006), pp. 308–360.
- [17] J. HUMPHERYS, O. LAFITTE, AND K. ZUMBRUN, *Stability of viscous shock profiles in the high Mach number limit*, preprint (2007).
- [18] P. HOWARD AND M. RAOOFI, *Pointwise asymptotic behavior of perturbed viscous shock profiles*, Adv. Differential Equations, 11 (2006), pp. 1031–1080.
- [19] P. HOWARD, M. RAOOFI, AND K. ZUMBRUN, *Sharp pointwise bounds for perturbed shock waves*, J. Hyperbolic Differential Equations, 3 (2006), pp. 297–374.
- [20] J. HUMPHERYS AND K. ZUMBRUN, *Spectral stability of small amplitude shock profiles for dissipative symmetric hyperbolic–parabolic systems*, Z. Angew. ed., Math. Phys., 53 (2002), pp. 20–34.

- [21] T. KAPITULA AND B. SANDSTED, *Stability of bright solitary-wave solutions to perturbed nonlinear Schrödinger equations*, Phys. D, 124 (1998), pp. 58–103.
- [22] T. KATO, *Perturbation Theory for Linear Operators*, Springer–Verlag, Berlin, Heidelberg (1985).
- [23] S. KAWASHIMA, S. NISHIBATA, AND P. ZHU, *Asymptotic stability of the stationary solution to the compressible Navier-Stokes equations in the half space*, Commun. Math. Phys., 240 (2003), pp. 483–500.
- [24] T.P. LIU AND K. NISHIHARA, *Asymptotic behavior for scalar viscous conservation laws with boundary effect*, J. Differential Equations, 133 (1997), pp. 296–320.
- [25] T.P. LIU AND S.-H. YU, *Propagation of stationary viscous Burgers shock under the effect of boundary*, Arch. Ration. Mech. Anal., 139 (1997), pp. 57–82.
- [26] A. MATSUMURA AND K. NISHIHARA, *Large-time behaviors of solutions to an inflow problem in the half space for a one-dimensional system of compressible viscous gas*, Commun. Math. Phys., 222 (2001), pp. 449–474.
- [27] C. MASCIA AND K. ZUMBRUN, *Pointwise Green’s function bounds for shock profiles with degenerate viscosity*, Arch. Ration. Mech. Anal., 169 (2003), pp. 177–263.
- [28] C. MASCIA AND K. ZUMBRUN, *Stability of large-amplitude shock profiles of hyperbolic–parabolic systems*, Arch. Ration. Mech. Anal., 172 (2004), pp. 93–131.
- [29] G. MÉTIVIER AND K. ZUMBRUN, *Large viscous boundary layers for noncharacteristic nonlinear hyperbolic problems*, Mem. Amer. Math. Soc., 175(826):vi+107 (2005).
- [30] F. ROUSSET, *Stability of small amplitude boundary layers for mixed hyperbolic-parabolic systems*, Trans. Amer. Math. Soc., 355 (2003), pp. 2291–3008.
- [31] F. ROUSSET, *Inviscid boundary conditions and stability of viscous boundary layers*, Asymptot. Anal., 26 (2001), pp. 285–306.
- [32] D. SERRE, *Sur la stabilité des couches limites de viscosité*, Annales de l’institut Fourier, 51 (2001), pp. 109–130.
- [33] D. SERRE AND K. ZUMBRUN, *Boundary layer stability in real vanishing-viscosity limit*, Commun. Math. Phys., 221 (2001), pp. 267–292.
- [34] H. SCHLICHTING, K. GERSTEN, E. KRAUSE, AND H. OERTEL, JR., *Boundary-Layer Theory*, Springer; 8th ed. 2000. Corr. 2nd printing edition (March 22, 2004).
- [35] S. YARAHMADIAN, *Pointwise Green function bounds and long-time stability of large-amplitude noncharacteristic boundary layers*, Doctoral thesis, Indiana University (2008).
- [36] K. ZUMBRUN, *Stability of large-amplitude shock waves of compressible Navier-Stokes equations*, In Handbook of mathematical fluid dynamics. Vol. III, pp. 311–533, North-Holland, Amsterdam, 2004. With an appendix by Helge Kristian Jenssen and Gregory Lyng.
- [37] K. ZUMBRUN AND P. HOWARD, *Pointwise semigroup methods and stability of viscous shock waves*, Indiana Math. J., 47 (1998), pp. 741–871.

A HIGHER ORDER MODEL FOR IMAGE RESTORATION: THE ONE-DIMENSIONAL CASE*

G. DAL MASO[†], I. FONSECA[‡], G. LEONI[‡], AND M. MORINI[†]

Abstract. The higher order total variation-based model for image restoration proposed by Chan, Marquina, and Mulet in [*SIAM J. Sci. Comput.*, 22 (2000), pp. 503–516] is analyzed in one dimension. A suitable functional framework in which the minimization problem is well posed is being proposed, and it is proved analytically that the higher order regularizing term prevents the occurrence of the *staircase effect*. The generalized version of the model considered here includes, as particular cases, some curvature dependent functionals.

Key words. image segmentation, total variation models, staircase effect, higher order regularization, relaxation, curvature dependent functionals

AMS subject classifications. 49J45, 26A45, 65K10, 68U10

DOI. 10.1137/070697823

1. Introduction. Deblurring and denoising of images are fundamental problems in image processing and gave rise in the past few years to a vast variety of techniques and methods touching different fields of mathematics. Among them, variational methods based on the minimization of some energy functional have been successfully employed to treat a fairly general class of image restoration problems. Typically, such functionals present a *fidelity term*, which penalizes the distance between the reconstructed image u and the noisy image g with respect to a suitable metric, and a regularizing term, which makes high frequency noise energetically unfavorable.

When the fidelity term is given by the squared L^2 distance multiplied by a parameter $\lambda > 0$ and the regularizing term is represented by the total variation, we are led to the following minimization problem:

$$(1.1) \quad \min \left\{ |Du|(\Omega) + \lambda \int_{\Omega} |u - g|^2 dx : u \in BV(\Omega) \right\},$$

which was proposed by Rudin, Osher, and Fatemi in [13]. Here Ω is an open bounded domain in one or two dimensions, $BV(\Omega)$ denotes the space of functions of bounded variations in Ω , and $|Du|(\Omega)$ stands for the total variation of u in Ω . The main feature of the total variation-based image restoration is perhaps represented by the tendency to yield (almost) piecewise constant solutions or, in other words, “blocky” images. Typically, one observes that *ramps* (i.e., affine regions) in the original image give rise to staircase-like structures in the reconstructed image, a phenomenon which is often referred to as the *staircase effect*. This means that the original edges are well preserved by this method but also that many artificial discontinuities can be generated by the presence of noise, while the finer details of the objects contained in the image may not be properly recovered.

Several variants of (1.1) have been subsequently proposed in order to fix these drawbacks. In this paper we follow the approach of Chan, Marquina, and Mulet [6]:

*Received by the editors July 9, 2007; accepted for publication (in revised form) October 10, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/69782.html>

[†]SISSA, Via Beirut 2, 34014 Trieste, Italy (dalmaso@sissa.it, morini@sissa.it).

[‡]Carnegie Mellon University, Pittsburgh, PA 15213-3890 (fonseca@andrew.cmu.edu, giovanni@andrew.cmu.edu).

Since the total variation does not distinguish between jumps and smooth transitions, their idea is to consider an additional penalization of the discontinuities by taking second derivatives into account. More precisely, they propose a regularizing term of the form

$$(1.2) \quad \int_{\Omega} |\nabla u| dx + \int_{\Omega} \psi(|\nabla u|)h(\Delta u) dx,$$

where ψ is a function that must satisfy suitable conditions at infinity in order to allow jumps.

In this paper we consider the following one-dimensional (1D) version of (1.2):

$$(1.3) \quad \mathcal{F}_p(u) := \int_a^b |u'| dx + \int_a^b \psi(|u'|)|u''|^p dx,$$

where $a < b$ are real numbers, $p \in [1, +\infty)$, and $\psi : \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying suitable integrability assumptions (see (2.1) and (3.1)). We remark that by taking

$$\psi(t) := \frac{1}{(1+t^2)^{\frac{1}{2}(3p-1)}}$$

the functional (1.3) takes the form

$$\int_a^b |u'| dx + \int_{\text{Graph } u} |k|^p d\mathcal{H}^1,$$

where k denotes the curvature of the graph of u and \mathcal{H}^1 stands for the 1D Hausdorff measure. This seems to suggest a broader applicability of our results: Functionals of this kind are often encountered in many computer vision and graphics applications, such as, among the others, corner preserving geometry denoising and segmentation with depth (see, e.g., [3], [15], and the references contained therein).

Our main analytical objective is twofold:

- (i) to set up a proper functional framework where the minimization problem corresponding to

$$\mathcal{F}_p(u) + \lambda \int_a^b |u - g|^2 dx$$

is well posed;

- (ii) to give an analytical proof of the fact that the higher order regularizing term eliminates the staircase effect.

We carry out the first part of this program in sections 2 and 3 by using the theory of relaxation (see [7] for a general introduction): We regard \mathcal{F}_p as defined for all functions in the Sobolev space $W^{2,p}(]a, b[)$, we extend it to $L^1(]a, b[)$ by setting $\mathcal{F}_p(u) := +\infty$ if $u \in L^1(]a, b[) \setminus W^{2,p}(]a, b[)$, and then we identify its lower semicontinuous envelope with respect to the strong L^1 convergence. The two cases $p = 1$ and $p > 1$ require a different treatment, and, in fact, the analysis turns out to be considerably more delicate in the case $p = 1$. Moreover, the domains of the relaxed functionals are quite peculiar (see Definitions 2.1 and 3.1) and display properties that are qualitatively different in the two cases. In particular, it turns out that piecewise constant functions corresponding to images with genuine edges are approximable by sequences with bounded energy

only for $p = 1$. The extension of these relaxation results to higher dimensions will be the subject of a subsequent paper.

The second part of the program is carried out in section 4. We start by exhibiting an analytical example of staircasing for the Rudin–Osher–Fatemi model (Theorem 4.3). More precisely, we show that if g has the form $g = g_1 + h$, with g_1 an affine function (which represents the original “clean” signal) and h a highly oscillating noise, then, for certain choices of h , the reconstructed signal u may display a stairlike or piecewise constant structure, even if h is uniformly small. In particular, we note that

- (a) the discontinuity set of u is much larger than the discontinuity set of the original signal g_1 ,
- (b) the derivatives of u and u_1 are very far apart,

where u_1 denotes the solution to the Rudin–Osher–Fatemi minimization problem corresponding to the datum g_1 . Although we do not attempt to give a formal definition of staircasing, we consider the presence/absence of (a) and (b) as a way to detect the presence/absence of the staircase effect. In Theorems 4.5 and 4.8 we prove that the new model eliminates the effect by showing that (a) and (b) do not occur. More precisely, we show that whenever the datum g is of the form $g = g_1 + h$, with g_1 a regular signal and h a highly oscillating noise close to 0 in the L^∞ -weak* topology, the reconstructed image u is regular as well (in particular, no new artificial discontinuities are created) and u and u_1 are close in some strong norm ($W^{1,q}$ -norm for every $q > 1$ or C^1 -norm, depending on whether $p = 1$ or $p > 1$), where u_1 is the solution corresponding to the clean signal g_1 .

For completeness we conclude by mentioning that other approaches have been considered to avoid staircasing: The works by Geman and Reynolds [9] and Chambolle and Lions [5] contain a different use of higher order derivatives as regularizing terms; in [2], Blomgren, Chan, and Mulet propose a $BV-H^1$ interpolation approach, while Kindermann, Osher, and Jones avoid in [11] the use of second derivatives by considering a sort of nonlocal total variation.

2. The case $p = 1$. We start by studying the compactness properties and the relaxation of (1.3) in the case $p = 1$. Throughout this section $\psi: \mathbb{R} \rightarrow]0, +\infty[$ will be a bounded Borel function such that

$$(2.1) \quad M := \int_{-\infty}^{+\infty} \psi(t) dt < +\infty$$

and

$$(2.2) \quad \inf_{t \in K} \psi(t) > 0 \quad \text{for every compact set } K \subset \mathbb{R}.$$

Let $\Psi_1: \overline{\mathbb{R}} \rightarrow [0, M]$ be the increasing function defined by

$$\Psi_1(t) := \int_{-\infty}^t \psi(s) ds,$$

and let $\Psi_1^{-1}: [0, M] \rightarrow \overline{\mathbb{R}}$ be its inverse function.

Given a bounded open interval $]a, b[$ in \mathbb{R} , we let $\mathcal{F}_1: L^1(]a, b[) \rightarrow [0, +\infty[$ be the functional defined by

$$(2.3) \quad \mathcal{F}_1(u) := \begin{cases} \int_a^b |u'| dx + \int_a^b \psi(u')|u''| dx & \text{if } u \in W^{2,1}(]a, b[), \\ +\infty & \text{otherwise.} \end{cases}$$

The first step in the study of (2.3) will consist in identifying the subspace of L^1 functions which can be approximated by energy bounded sequences. In order to do so, we need to introduce some notation and recall some basic facts about BV functions of one variable. This will be the content of the next subsection.

2.1. BV functions of one variable. We recall that a function $u \in L^1(]a, b[)$ belongs to $BV(]a, b[)$ if and only if

$$(2.4) \quad \sup \left\{ \int_a^b u \varphi' dx : \varphi \in C_c^1(]a, b[), |\varphi| \leq 1 \right\} < +\infty.$$

Note that this implies that the distributional derivative u' of u is a bounded Radon measure in $]a, b[$. We will often consider the Lebesgue decomposition

$$u' = (u')^a \mathcal{L}^1 + (u')^s,$$

where $(u')^a$ is the density of the absolutely continuous part of u' with respect to the Lebesgue measure \mathcal{L}^1 on $]a, b[$ while $(u')^s$ is its singular part. We will denote the total variation measure of u' by $|u'|$. In particular, $|u'| (]a, b[)$ equals the value of the supremum in (2.4). For every function $u \in BV(]a, b[)$ the following left and right approximate limits

$$u_-(y) := \lim_{\varepsilon \rightarrow 0^+} \frac{1}{\varepsilon} \int_{y-\varepsilon}^y u(x) dx, \quad u_+(y) := \lim_{\varepsilon \rightarrow 0^+} \frac{1}{\varepsilon} \int_y^{y+\varepsilon} u(x) dx,$$

respectively, are well defined at every point $y \in]a, b[$. In fact, $u_-(y)$ is well defined also at $y = b$ while $u_+(y)$ exists also at $y = a$. The functions u_- and u_+ coincide \mathcal{L}^1 -a.e. with u and are left and right continuous, respectively. Moreover, it turns out that the set $S_u := \{y \in]a, b[: u_-(y) \neq u_+(y)\}$ is at most countable. The set S_u is often referred to as the set of essential discontinuities or *jump points* of u .

It is well known that, in turn, the singular part $(u')^s$ splits into the sum of an atomic measure concentrated on S_u and a singular diffuse measure $(u')^c$, called the *Cantor part* of u' :

$$(u')^s = [u] \mathcal{H}^0 \llcorner S_u + (u')^c,$$

where we set $[u] := u_+ - u_-$ and \mathcal{H}^0 stands for the counting measure. Finally, we recall that every $u \in BV(]a, b[)$ is differentiable at \mathcal{L}^1 -a.e. y in $]a, b[$ with derivative given by $(u')^a(y)$. In this case, we will often write, with a slight abuse of notation, $u'(y)$ instead of $(u')^a(y)$.

We say that a sequence $\{u_k\}$ of functions in $BV(]a, b[)$ *weakly star converges* in $BV(]a, b[)$ to a function $u \in BV(]a, b[)$ if $u_n \rightarrow u$ in $L^1(]a, b[)$ and $u'_k \rightarrow u'$ weakly* in $M_b(]a, b[)$, where $M_b(]a, b[)$ is the space of bounded Radon measures.

We will also need sometimes the notion of total variation for a function defined everywhere. We recall that $u:]a, b[\rightarrow \mathbb{R}$ has bounded *pointwise total variation* over the interval $]c, d[\subset]a, b[$ if

$$\text{Var}(u;]c, d[) := \sup \sum_{i=1}^k |u(y_i) - u(y_{i-1})| < +\infty,$$

where the supremum is taken over all finite families y_0, y_1, \dots, y_k such that $c < y_0 < y_1 < \dots < y_k < d$, $k \in \mathbb{N}$. It is easy to see that if u has bounded pointwise total

variation in $]a, b[$, then it admits left and right limits at every point, it belongs to $BV(]a, b[)$, and $|u'|(|c, d|) \leq \text{Var}(u;]c, d|)$ for every interval $]c, d| \subset]a, b[$. Conversely, if $u \in BV(]a, b[)$, the *precise representatives* u_- and u_+ have bounded pointwise total variation and satisfy

$$|u'|(|c, d|) = \text{Var}(u_-;]c, d|) = \text{Var}(u_+;]c, d|)$$

for every interval $]c, d| \subset]a, b[$.

Finally, we recall the Helly theorem: For every bounded sequence of functions $u_k :]a, b[\rightarrow \mathbb{R}$ such that $\sup_k \text{Var}(u_k;]a, b|) < +\infty$, there exist u , with pointwise total variation in $]a, b[$, and a subsequence (not relabeled) such that $u_k \rightarrow u$ pointwise.

We refer to [14] and [10] for an exhaustive exposition of the properties of BV functions of one variable.

2.2. Compactness. To define the subspace of L^1 functions that can be approximated by energy bounded sequences, for every function $u \in BV(]a, b|)$ we consider the sets

$$(2.5) \quad Z^+[(u')^a] := \left\{ x \in]a, b[: \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{x-\varepsilon}^{x+\varepsilon} (u')^a dx = +\infty \right\},$$

$$(2.6) \quad Z^-[(u')^a] := \left\{ x \in]a, b[: \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{x-\varepsilon}^{x+\varepsilon} (u')^a dx = -\infty \right\}.$$

It is also convenient to define

$$Z[(u')^a] := Z^+[(u')^a] \cup Z^-[(u')^a].$$

DEFINITION 2.1. Let $X_\psi^1(]a, b|)$ be the set of all functions $u \in BV(]a, b|)$ such that $v := \Psi_1 \circ (u')^a$ belongs to $BV(]a, b|)$ and the positive part $((u')^c)^+$ and the negative part $((u')^c)^-$ of the measure $(u')^c$ are concentrated on $Z^+[(u')^a]$ and $Z^-[(u')^a]$, respectively.

Remark 2.2. Note that if $u \in X_\psi^1(]a, b|)$, then the limits

$$(2.7) \quad (u')^a_-(y) := \lim_{\varepsilon \rightarrow 0^+} \frac{1}{\varepsilon} \int_{y-\varepsilon}^y (u')^a dx, \quad (u')^a_+(y) := \lim_{\varepsilon \rightarrow 0^+} \frac{1}{\varepsilon} \int_y^{y+\varepsilon} (u')^a dx$$

exist in $\overline{\mathbb{R}}$ for every y . More precisely, $(u')^a_-$ exists also at $y = b$, while $(u')^a_+$ is well defined also at $y = a$. Indeed, since $v = \Psi_1 \circ (u')^a$ is a BV function, it admits a precise representative \tilde{v} such that the right and left limits exist at every point, and the same property holds for $\Psi_1^{-1}(\tilde{v})$. As $\Psi_1^{-1}(\tilde{v}) = (u')^a$ \mathcal{L}^1 -a.e. in $]a, b[$, the limits considered in (2.7) are everywhere well defined. Moreover, the set $S_{(u')^a} := S_v$ is at most countable and

$$(2.8) \quad (u')^a_- = (u')^a_+ \quad \text{on }]a, b[\setminus S_{(u')^a}.$$

We also remark that $(u')^a_-$ and $(u')^a_+$ are left and right continuous, respectively, which, in turn, implies that the functions defined by

$$(u')^a_{\vee}(x) := \max \{ (u')^a_+(x), (u')^a_-(x) \}, \quad (u')^a_{\wedge}(x) := \min \{ (u')^a_+(x), (u')^a_-(x) \}$$

if $x \in]a, b[$ and by

$$(u')^a_{\vee}(a) = (u')^a_{\wedge}(a) := (u')^a_+(a), \quad (u')^a_{\vee}(b) = (u')^a_{\wedge}(b) := (u')^a_-(b)$$

are upper and lower semicontinuous in $[a, b]$, respectively. By (2.8) we have

$$\begin{aligned} Z^+[(u')^a] \setminus S_{(u')^a} &= \{x \in]a, b[: (u')^a_\wedge(x) = +\infty\} \setminus S_{(u')^a}, \\ Z^-[(u')^a] \setminus S_{(u')^a} &= \{x \in]a, b[: (u')^a_\vee(x) = -\infty\} \setminus S_{(u')^a}. \end{aligned}$$

Therefore, $((u')^c)^+$ is concentrated on the set $\{x \in]a, b[: (u')^a_\wedge(x) = +\infty\}$ and $((u')^c)^-$ is concentrated on the set $\{x \in]a, b[: (u')^a_\vee(x) = -\infty\}$.

Before we proceed we show that the space $X^1_\psi(]a, b[)$ contains functions with nontrivial Cantor part when ψ satisfies suitable decay estimates at infinity.

PROPOSITION 2.3. *Assume that $\psi : \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying (2.1), (2.2), and*

$$(2.9) \quad \psi(t) \leq \frac{c}{t^\alpha}$$

for all $t \geq 1$ and for some $c > 0, \alpha > 1$. Then there exists $u \in X^1_\psi(]a, b[)$ with $(u')^c \neq 0$.

Proof. For simplicity we take $]a, b[=]0, 1[$.

Step 1. We start by recalling the definition of the generalized Cantor set \mathbb{D}_δ , where $\delta \in]0, \frac{1}{2}[$ (see, for instance, [8, Chapter 1, section 2.4]). The construction is entirely similar to the one of the (ternary) Cantor set with the only difference that the middle intervals removed at each step have length $1 - 2\delta$ times the length of the intervals remaining from the previous step. To be more precise, remove from $[0, 1]$ the interval $I_{11} := (\delta, 1 - \delta)$. At the second step remove from each of the remaining closed intervals $[0, \delta]$ and $[1 - \delta, 1]$ the middle intervals, denoted by I_{12} and I_{22} , of length $\delta(1 - 2\delta)$. Continuing in this fashion at each step n we remove 2^{n-1} middle intervals $I_{1n}, \dots, I_{2^{n-1}n}$, each of length $\delta^{n-1}(1 - 2\delta)$. The *generalized Cantor set* \mathbb{D}_δ is defined as

$$\mathbb{D}_\delta := [0, 1] \setminus \bigcup_{n=1}^\infty \bigcup_{k=1}^{2^{n-1}} I_{kn}.$$

The set \mathbb{D}_δ is closed (since its complement is given by a family of open intervals), and

$$\begin{aligned} \mathcal{L}^1(\mathbb{D}_\delta) &= 1 - \sum_{n=1}^\infty \sum_{k=1}^{2^{n-1}} \mathcal{L}^1(I_{kn}) = 1 - \sum_{n=1}^\infty \sum_{k=1}^{2^{n-1}} \delta^{n-1}(1 - 2\delta) \\ &= 1 - (1 - 2\delta) \sum_{n=1}^\infty (2\delta)^{n-1} = 0. \end{aligned}$$

Next, we recall the definition of the corresponding Cantor function f_δ . Set

$$g_n := \frac{1}{(2\delta)^n} \left(1 - \sum_{j=1}^n \sum_{k=1}^{2^{j-1}} \chi_{I_{kj}} \right),$$

and define $f_n(x) := \int_0^x g_n(t) dt$. It can be shown that $\{f_n\}$ converges uniformly to a continuous nondecreasing function f_δ such that $f_\delta(0) = 0, f_\delta(1) = 1$, and $f'_\delta = (f'_\delta)^c$ is supported on \mathbb{D}_δ .

Step 2. We claim that it is enough to find a constant $\delta \in]0, \frac{1}{2}[$ for which it is possible to construct a continuous integrable function $w_\delta :]0, 1[\rightarrow [0, +\infty[$ such

that $\Psi_1 \circ w_\delta \in BV(]0, 1[)$ and $w_\delta(x) = +\infty$ if and only if $x \in \mathbb{D}_\delta$. Indeed, setting $u_\delta(x) := \int_0^x w_\delta(t) dt + f_\delta(x)$, we have that $u_\delta \in BV(]0, 1[)$, u_δ is continuous, and $(u'_\delta)^a = w_\delta$ so that $Z^+[(u'_\delta)^a] = Z[(u'_\delta)^a] = \mathbb{D}_\delta$ and $\Psi_1 \circ (u'_\delta)^a \in BV(]0, 1[)$. Moreover, $(u'_\delta)^c = (f'_\delta)^c$ is supported on $\mathbb{D}_\delta = Z^+[(u'_\delta)^a]$. Hence, u_δ belongs to $X^1_\psi(]a, b[)$.

Step 3. It remains to construct w_δ for a suitable $\delta \in]0, \frac{1}{2}[$. Consider a convex function $\phi :]0, 1[\rightarrow [0, +\infty)$ such that

$$(2.10) \quad \lim_{x \rightarrow 0^+} \phi(x) = \lim_{x \rightarrow 1^-} \phi(x) = +\infty, \quad \phi\left(\frac{1}{2}\right) = 0,$$

and

$$(2.11) \quad \int_0^1 \phi(x) dx = 1.$$

Choose $s > 0$ so large that

$$(2.12) \quad \alpha > \frac{s+1}{s}.$$

For $x \in I_{kn}$ (see Step 1) define

$$(2.13) \quad \phi_{kn}(x) := 2^{sn} + \phi\left(\frac{x - a_{kn}}{\delta^{n-1}(1-2\delta)} + \frac{1}{2}\right),$$

where a_{kn} is the midpoint of the interval I_{kn} . Finally, set

$$w_\delta := \sum_{n=1}^\infty \sum_{k=1}^{2^{n-1}} \phi_{kn} \chi_{I_{kn}} + I_{\mathbb{D}_\delta},$$

where $I_{\mathbb{D}_\delta}$ is the indicator function of the set \mathbb{D}_δ , that is,

$$I_{\mathbb{D}_\delta}(x) := \begin{cases} +\infty & \text{if } x \in \mathbb{D}_\delta, \\ 0 & \text{otherwise.} \end{cases}$$

Using the fact that

$$\int_{I_{kn}} \phi_{kn} dx = (2^{sn} + 1) \delta^{n-1} (1 - 2\delta),$$

which follows from (2.11) and a change of variables, we have

$$\int_0^1 w_\delta dx = \sum_{n=1}^\infty \sum_{k=1}^{2^{n-1}} (2^{sn} + 1) \delta^{n-1} (1 - 2\delta) < \infty$$

for $\delta < \frac{1}{2^{s+1}}$. To estimate the total variation of $v := \Psi_1 \circ w_\delta$, we consider the approximating sequence

$$v_m(x) := \begin{cases} \Psi_1 \circ \phi_{kn}(x) & \text{if } x \in I_{kn}, 1 \leq k \leq 2^{n-1}, 1 \leq n \leq m, \\ M & \text{otherwise.} \end{cases}$$

By (2.9), (2.10), (2.13), and the convexity of ϕ , it can be seen that

$$\text{Var}(v_m; I_{kn}) = 2(M - \Psi_1(2^{sn})) = 2 \int_{2^{sn}}^{+\infty} \psi(t) dt \leq \frac{2c}{\alpha - 1} \frac{1}{2^{sn(\alpha-1)}}.$$

It follows that

$$\text{Var}(v_m;]0, 1[) \leq \frac{2c}{\alpha - 1} \sum_{n=1}^m \sum_{k=1}^{2^{n-1}} \frac{1}{2^{sn(\alpha-1)}} \leq \frac{2c}{\alpha - 1} \sum_{n=1}^{\infty} \frac{1}{2^{sn(\alpha-1)-n+1}}.$$

The last series is finite, thanks to (2.12). Therefore, the v_m 's have equibounded total variations, and, since $v_m \rightarrow v$ in $L^1(]0, 1[)$, we conclude that $v \in BV(]0, 1[)$. \square

Energy bounded sequences are compact in $X_\psi^1(]a, b[)$, as made precise by the following theorem.

THEOREM 2.4. *Let $\{u_k\}$ be a sequence of functions bounded in $L^1(]a, b[)$ such that*

$$(2.14) \quad C := \sup_k \mathcal{F}_1(u_k) < +\infty.$$

Then there exist a subsequence (not relabeled) $\{u_k\}$ and a function $u \in X_\psi^1(]a, b[)$ such that

$$(2.15) \quad u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(]a, b[),$$

$$(2.16) \quad \Psi_1 \circ u'_k \rightharpoonup \Psi_1 \circ (u')^a \quad \text{weakly}^* \text{ in } BV(]a, b[),$$

$$(2.17) \quad u'_k \rightarrow (u')^a \quad \text{pointwise } \mathcal{L}^1\text{-a.e. in }]a, b[.$$

Proof. By (2.3) and (2.14) we have that each u_k belongs to $W^{2,1}(]a, b[)$ and

$$(2.18) \quad C_1 := \sup_k \int_a^b [|u_k| + |u'_k| + \psi(u'_k)|u''_k|] dx < +\infty.$$

Let us define

$$(2.19) \quad v_k := \Psi_1 \circ u'_k.$$

As Ψ_1 is Lipschitz in \mathbb{R} , the functions v_k belong to $W^{1,1}(]a, b[)$ and

$$(2.20) \quad v'_k = \psi(u'_k)u''_k \quad \mathcal{L}^1\text{-a.e. on }]a, b[.$$

It follows from (2.1) and (2.14) that

$$(2.21) \quad \int_a^b [|v_k| + |v'_k|] dx \leq M(b - a) + C.$$

By (2.18) and (2.21) and the Helly theorem, passing to a subsequence if necessary, we may assume that

$$u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(]a, b[)$$

and

$$(2.22) \quad v_k(x) \rightarrow v(x) \quad \text{for all } x \in]a, b[$$

for some $u \in BV(]a, b[)$ and $v:]a, b[\rightarrow [0, M]$ with pointwise bounded variation. Note that (2.22) determines the values of v at every $x \in]a, b[$.

Since Ψ_1^{-1} is continuous, we obtain

$$(2.23) \quad u'_k \rightarrow w := \Psi_1^{-1}(v) \quad \text{pointwise in }]a, b[.$$

Since by (2.18) we have $\sup_k \|u'_k\|_1 < +\infty$, it follows from Fatou's lemma and (2.23) that w is integrable, and so

$$(2.24) \quad w \text{ is finite } \mathcal{L}^1\text{-a.e. in }]a, b[.$$

Moreover, w has left and right limits in $\overline{\mathbb{R}}$ at each point $x \in]a, b[$, denoted by $w_-(x)$ and $w_+(x)$, respectively, and

$$(2.25) \quad w(x) = w_-(x) = w_+(x) \quad \text{except for a countable set of points } x.$$

We now split the remaining part of the proof into two steps.

Step 1. We prove that

$$(2.26) \quad w = (u')^a \quad \mathcal{L}^1\text{-a.e. in }]a, b[.$$

If not, we have $\mathcal{L}^1(\{w \neq (u')^a\}) > 0$. By (2.24) there exists $t_0 > 0$ such that

$$\mathcal{L}^1(\{w \neq (u')^a\} \cap \{|w| < t_0\}) > 0,$$

and, in particular, we may find an infinite number of disjoint open intervals I such that

$$(2.27) \quad \mathcal{L}^1(\{w \neq (u')^a\} \cap \{|w| < t_0\} \cap I) > 0.$$

By a change of variables we obtain

$$(2.28) \quad \int_I \psi(u'_k) |u''_k| dx \geq \int_{m_k}^{M_k} \psi(t) dt,$$

where

$$m_k := \inf_I u'_k \quad \text{and} \quad M_k := \sup_I u'_k.$$

We claim that at least one of the two sequences $\{m_k\}$ and $\{M_k\}$ is divergent. Indeed, if not, a subsequence of $\{u'_k\}$ would be bounded in $L^\infty(I)$. This implies that $u' \in L^\infty(I)$ and that $u'_k \rightharpoonup u'$ weakly* in $L^\infty(I)$. As $u'_k \rightarrow w$ pointwise \mathcal{L}^1 -a.e. in I , we deduce that $u' = w$ \mathcal{L}^1 -a.e. in I , which contradicts (2.27). Hence, the claim holds. If

$$(2.29) \quad \lim_{k \rightarrow \infty} M_k = +\infty,$$

then by (2.23) and (2.27)

$$(2.30) \quad \limsup_{k \rightarrow \infty} m_k < t_0.$$

From (2.28), (2.30), and (2.29) we obtain

$$\liminf_{k \rightarrow \infty} \int_I \psi(u'_k) |u''_k| dx \geq \int_{t_0}^{+\infty} \psi(t) dt > 0.$$

Analogously, if $\lim_k m_k = -\infty$, then

$$\liminf_{k \rightarrow \infty} \int_I \psi(u'_k) |u''_k| dx \geq \int_{-\infty}^{-t_0} \psi(t) dt > 0.$$

In any case, we can choose an arbitrarily large number m of disjoint intervals I satisfying (2.27). Adding the contributions of each interval we obtain

$$\liminf_{k \rightarrow \infty} \int_a^b \psi(u'_k) |u''_k| dx \geq m \min \left\{ \int_{t_0}^{+\infty} \psi(t) dx, \int_{-\infty}^{-t_0} \psi(t) dt \right\},$$

which contradicts (2.18) for m large enough. This concludes the proof of (2.26).

Step 2. To prove that $u \in X^1_\psi(]a, b[)$, it remains to show that the positive part $((u')^c)^+$ and the negative part $((u')^c)^-$ of the measure $(u')^c$ are concentrated on $Z^+[(u')^a]$ and $Z^-[(u')^a]$, respectively, that is,

$$(2.31) \quad ((u')^c)^\pm(]a, b[\setminus Z^\pm[(u')^a]) = 0.$$

To this purpose we introduce the sets

$$(2.32) \quad E^+[u'] := \left\{ x \in]a, b[: \lim_{\varepsilon \rightarrow 0^+} \frac{(u')^+(]x - \varepsilon, x + \varepsilon])}{2\varepsilon} = +\infty \right\},$$

$$(2.33) \quad E^-[u'] := \left\{ x \in]a, b[: \lim_{\varepsilon \rightarrow 0^+} \frac{(u')^-(]x - \varepsilon, x + \varepsilon])}{2\varepsilon} = +\infty \right\},$$

$$E[u'] := \left\{ x \in]a, b[: \lim_{\varepsilon \rightarrow 0^+} \frac{|u'|(|x - \varepsilon, x + \varepsilon])}{2\varepsilon} = +\infty \right\}.$$

Since $((u')^s)^+ = ((u')^+)^s$ is concentrated on $E^+[u']$ and $((u')^s)^- = ((u')^-)^s$ is concentrated on $E^-[u']$ (see, e.g., [1, Theorem 2.22]), to prove (2.31) it is enough to show that

$$(2.34) \quad E^+[u'] \setminus Z^+[(u')^a] \text{ and } E^-[u'] \setminus Z^-[(u')^a] \text{ are at most countable.}$$

We show only that $E^+[u'] \setminus Z^+[(u')^a]$ is at most countable, since the other property can be proved in a similar way. Assume, by contradiction, that $E^+[u'] \setminus Z^+[(u')^a]$ is not countable. Since by (2.5) and (2.26)

$$Z^+[(u')^a] \subset \{x \in]a, b[: \max\{w_-(x), w_+(x)\} = +\infty\},$$

by (2.25) there exists $t_0 > 0$ such that

$$(E^+[u'] \setminus Z^+[(u')^a]) \cap \{w < t_0\} \text{ is uncountable.}$$

Fix $t_1 > t_0$, and let x_1, \dots, x_m be m distinct points in $(E^+[u'] \setminus Z^+[(u')^a]) \cap \{w < t_0\}$. By (2.32) there exists $\varepsilon > 0$ such that the intervals $I_j :=]x_j - \varepsilon, x_j + \varepsilon[$ are pairwise disjoint and

$$(2.35) \quad \frac{(u')^+(]x_j - \varepsilon, x_j + \varepsilon])}{2\varepsilon} > t_1 \quad \text{for } i = 1, \dots, m.$$

By a change of variables we obtain

$$(2.36) \quad \int_{I_j} \psi(u'_k) |u''_k| dx \geq \int_{m_{kj}}^{M_{kj}} \psi(t) dt,$$

where

$$m_{kj} := \inf_{I_j} u'_k \quad \text{and} \quad M_{kj} := \sup_{I_j} u'_k.$$

By (2.23) and the fact that $w(x_j) < t_0$, we deduce that

$$(2.37) \quad \limsup_{k \rightarrow \infty} m_{kj} < t_0$$

for $j = 1, \dots, m$. On the other hand, (2.15) and (2.35) yield

$$\liminf_{k \rightarrow \infty} \frac{1}{2\varepsilon} \int_{x_j - \varepsilon}^{x_j + \varepsilon} (u'_k)^+ dx \geq \frac{(u')^+(]x_j - \varepsilon, x_j + \varepsilon])}{2\varepsilon} > t_1$$

(this can be seen as a particular case of the Reshetnyak lower semicontinuity theorem, with $f = (\cdot)^+$). This implies that $\liminf_{k \rightarrow \infty} M_{kj} > t_1$ for $j = 1, \dots, m$. Hence, also by (2.36) and (2.37) we obtain

$$\liminf_{k \rightarrow \infty} \sum_{j=1}^m \int_{I_j} \psi(u'_k) |u''_k| dx \geq \sum_{j=1}^m \liminf_{k \rightarrow \infty} \int_{I_j} \psi(u'_k) |u''_k| dx \geq m \int_{t_0}^{t_1} \psi(t) dt,$$

which contradicts (2.18) for m large enough. This shows (2.34) and concludes the proof of the theorem. \square

2.3. Relaxation. The following theorem, which is the main result of the section, is devoted to the characterization of the relaxation of \mathcal{F}_1 with respect to strong convergence in $L^1(]a, b[)$.

THEOREM 2.5. *Let $\overline{\mathcal{F}}_1 : L^1(]a, b[) \rightarrow [0, +\infty]$ be defined by*

$$\overline{\mathcal{F}}_1(u) := \inf \left\{ \liminf_{k \rightarrow \infty} \mathcal{F}_1(u_k) : u_k \rightarrow u \text{ in } L^1(]a, b[) \right\}$$

for every $u \in L^1(]a, b[)$. Then

$$(2.38) \quad \overline{\mathcal{F}}_1(u) = \begin{cases} |u'|(|a, b|) + |v'|(|a, b[\setminus S_u) + \sum_{x \in S_u} \Phi(\nu_u, (u')_-^a, (u')_+^a) & \text{if } u \in X_\psi^1(]a, b[), \\ +\infty & \text{otherwise,} \end{cases}$$

where $v := \Psi_1 \circ (u')^a$, $\nu_u := \text{sign}(u_+ - u_-)$, and

$$(2.39) \quad \begin{aligned} \Phi(1, t_1, t_2) &:= \int_{t_1}^{+\infty} \psi(t) dt + \int_{t_2}^{+\infty} \psi(t) dt, \\ \Phi(-1, t_1, t_2) &:= \int_{-\infty}^{t_1} \psi(t) dt + \int_{-\infty}^{t_2} \psi(t) dt. \end{aligned}$$

Remark 2.6. For every $x \in S_u$ we have

$$\Phi(\nu_u(x), (u')_-^a(x), (u')_+^a(x)) = |v'|(\{x\}) + \hat{\Phi}(\nu_u(x), (u')_-^a(x), (u')_+^a(x)),$$

where

$$\hat{\Phi}(1, t_1, t_2) := \int_{\max\{t_1, t_2\}}^{+\infty} \psi(t) dt \quad \text{and} \quad \hat{\Phi}(-1, t_1, t_2) := \int_{-\infty}^{\min\{t_1, t_2\}} \psi(t) dt.$$

In particular, for every Borel set $B \subset]a, b[$

$$\begin{aligned} &|v'| (B \setminus S_u) + \sum_{x \in S_u \cap B} \Phi(\nu_u, (u')_-^a, (u')_+^a) \\ &= |v'| (B) + \sum_{x \in S_u \cap B} \hat{\Phi}(\nu_u, (u')_-^a, (u')_+^a) \geq |v'| (B). \end{aligned}$$

Proof of Theorem 2.5. Let \mathcal{G} be the functional defined by the right-hand side of (2.38). We prove that, for every $u_k \rightarrow u$ in $L^1(]a, b[)$, we have

$$(2.40) \quad \mathcal{G}(u) \leq \liminf_{k \rightarrow \infty} \mathcal{F}_1(u_k).$$

It is enough to consider sequences $\{u_k\}$ for which the liminf is a limit and has a finite value and $u_k \rightarrow u$ pointwise \mathcal{L}^1 -a.e. in $]a, b[$. Then u_k belongs to $W^{2,1}(]a, b[)$ and (2.14) is satisfied. This implies that

$$(2.41) \quad |u'| (]a, b[) \leq \liminf_{k \rightarrow \infty} \int_a^b |u'_k| dx.$$

Moreover, it follows from Theorem 2.4 that $u \in X_{\psi}^1(]a, b[)$ and that, up to a subsequence, $\{u'_k\}$ converges to $(u')^a$ pointwise \mathcal{L}^1 -a.e. in $]a, b[$.

Let F be a finite subset of S_u . We want to prove that

$$(2.42) \quad |v'| (]a, b[\setminus F) + \sum_{x \in F} \Phi(\nu_u, (u')^a_-, (u')^a_+) \leq \liminf_{k \rightarrow \infty} \int_a^b \psi(u'_k) |u''_k| dx.$$

We write F as $\{x_1, \dots, x_m\}$, with $a < x_1 < \dots < x_m < b$. For every $\varepsilon > 0$ there exists $\delta = \delta(\varepsilon) \in]0, \varepsilon[$ such that $a < x_1 - \delta < x_1 + \delta < x_2 - \delta < x_2 + \delta < \dots < x_{m-1} - \delta < x_{m-1} + \delta < x_m - \delta < x_m + \delta < b$ and

$$(2.43) \quad |u(x_j - \delta) - u_-(x_j)| < \varepsilon, \quad |u(x_j + \delta) - u_+(x_j)| < \varepsilon,$$

$$(2.44) \quad |(u')^a(x_j - \delta) - (u')^a_-(x_j)| < \varepsilon, \quad |(u')^a(x_j + \delta) - (u')^a_+(x_j)| < \varepsilon,$$

$$(2.45) \quad u_k(x_j - \delta) \rightarrow u(x_j - \delta), \quad u_k(x_j + \delta) \rightarrow u(x_j + \delta) \quad \text{as } k \rightarrow \infty,$$

$$(2.46) \quad u'_k(x_j - \delta) \rightarrow (u')^a(x_j - \delta), \quad u'_k(x_j + \delta) \rightarrow (u')^a(x_j + \delta) \quad \text{as } k \rightarrow \infty,$$

$$|(u')^a(x_j - \delta)| + |(u')^a(x_j + \delta)| + \varepsilon < \frac{|[u](x_j)| - 4\varepsilon}{2\delta}$$

for $j = 1, \dots, m$.

Since $v_k \rightarrow v$ pointwise \mathcal{L}^1 -a.e. in $]a, b[$ and $v'_k = \psi(u'_k)u''_k$ \mathcal{L}^1 -a.e. in $]a, b[$, we obtain

$$|v'| (]x_j + \delta, x_{j+1} - \delta[) \leq \liminf_{k \rightarrow \infty} \int_{x_j + \delta}^{x_{j+1} - \delta} \psi(u'_k) |u''_k| dx$$

for $j = 1, \dots, m - 1$. A similar result holds for the intervals $]a, x_1 - \delta[$ and $]x_m + \delta, b[$. Let F_δ be the union of the intervals $[x_j - \delta, x_j + \delta]$ for $j = 1, \dots, m$. Summing with respect to j and adding the contributions of the intervals $]a, x_1 - \delta[$ and $]x_m + \delta, b[$, we obtain

$$(2.47) \quad |v'| (]a, b[\setminus F_\delta) \leq \liminf_{k \rightarrow \infty} \int_{]a, b[\setminus F_\delta} \psi(u'_k) |u''_k| dx.$$

We consider now the interval $I_j^\delta := [x_j - \delta, x_j + \delta]$, assuming that $[u](x_j) = u_+(x_j) - u_-(x_j) > 0$. By the mean value theorem there exists $y_{kj}^\delta \in]x_j - \delta, x_j + \delta[$ such that

$$(2.48) \quad u'_k(y_{kj}^\delta) = \frac{u_k(x_j + \delta) - u_k(x_j - \delta)}{2\delta} \geq \frac{[u](x_j) - 4\varepsilon}{2\delta},$$

where the last inequality follows from (2.43) and (2.45) for k sufficiently large. By a change of variables we obtain

$$\int_{u'_k(x_j-\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt \leq \int_{x_j-\delta}^{y_{kj}^\delta} \psi(u'_k)|u''_k| dx,$$

$$\int_{u'_k(x_j+\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt \leq \int_{y_{kj}^\delta}^{x_j+\delta} \psi(u'_k)|u''_k| dx.$$

Adding these inequalities and taking the limit as $k \rightarrow \infty$ we obtain, thanks to (2.46),

$$(2.49) \quad \int_{(u')^a(x_j-\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt + \int_{(u')^a(x_j+\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt \leq \liminf_{k \rightarrow \infty} \int_{x_j-\delta}^{x_j+\delta} \psi(u'_k)|u''_k| dx.$$

Similarly, if $[u](x_j) < 0$, then we have

$$(2.50) \quad \int_{\frac{[u](x_j)+4\epsilon}{2\delta}}^{(u')^a(x_j-\delta)} \psi(t) dt + \int_{\frac{[u](x_j)+4\epsilon}{2\delta}}^{(u')^a(x_j+\delta)} \psi(t) dt \leq \liminf_{k \rightarrow \infty} \int_{x_j-\delta}^{x_j+\delta} \psi(u'_k)|u''_k| dx.$$

From (2.47), (2.49), and (2.50) we deduce that

$$\begin{aligned} |v'|(\]a, b[\setminus F_\delta) + \sum_{[u](x_j) > 0} & \left(\int_{(u')^a(x_j-\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt + \int_{(u')^a(x_j+\delta)}^{\frac{[u](x_j)-4\epsilon}{2\delta}} \psi(t) dt \right) \\ + \sum_{[u](x_j) < 0} & \left(\int_{\frac{[u](x_j)+4\epsilon}{2\delta}}^{(u')^a(x_j-\delta)} \psi(t) dt + \int_{\frac{[u](x_j)+4\epsilon}{2\delta}}^{(u')^a(x_j+\delta)} \psi(t) dt \right) \\ & \leq \liminf_{k \rightarrow \infty} \int_a^b \psi(u'_k)|u''_k| dx. \end{aligned}$$

Taking the limit as $\epsilon \rightarrow 0$ (which implies $\delta(\epsilon) \rightarrow 0$) we obtain (2.42), thanks to (2.44).

Since S_u is at most countable, (2.40) can be obtained from (2.42) by taking the supremum over all finite sets F contained in S_u .

Conversely, let $u \in X^1_\psi(\]a, b[)$. We claim that there exists a sequence $\{u_k\}$ in $W^{2,1}(\]a, b[)$ such that $u_k \rightarrow u$ in $L^1(\]a, b[)$ and

$$(2.51) \quad \mathcal{G}(u) \geq \limsup_{k \rightarrow \infty} \mathcal{F}_1(u_k).$$

It is clearly enough to consider the case $\mathcal{G}(u) < +\infty$.

We divide the proof into three steps.

Step 1. We prove (2.51) under the additional assumptions that $(u')^a$ is bounded and that $S_u = \{x_1, \dots, x_m\}$, with $x_1 < \dots < x_m$. Note that, in this case, $Z[(u')^a] = \emptyset$ and, hence, $(u')^c = 0$.

Construct a sequence $\{v_k\}$ in $W^{1,1}(\]a, b[)$ such that $v_k \rightarrow v = \Psi_1 \circ (u')^a$ pointwise \mathcal{L}^1 -a.e. in $\]a, b[$, $\Psi_1(-\|(u')^a\|_\infty) \leq v_k \leq \Psi_1(\|(u')^a\|_\infty)$, and

$$\int_a^b |v'_k(x)| dx \rightarrow |v'|(\]a, b[).$$

Setting $w_k := \Psi_1^{-1}(v_k)$, we have $w_k \in W^{1,1}(\]a, b[)$, thanks to (2.2),

$$(2.52) \quad w_k \rightarrow (u')^a \quad \text{pointwise } \mathcal{L}^1\text{-a.e. in } \]a, b[,$$

and $\|w_k\|_\infty \leq \|(u')^a\|_\infty$. Find $\delta_k \rightarrow 0^+$ such that

$$(2.53) \quad w_k(x_j - \delta_k) \rightarrow (u')_-^a(x_j), \quad w_k(x_j + \delta_k) \rightarrow (u')_+^a(x_j) \quad \text{for } j = 1, \dots, m$$

and

$$(2.54) \quad \int_{x_{j-1}+\delta_k}^{x_j-\delta_k} |v'_k| dx \rightarrow |v'|([x_{j-1}, x_j]) \quad \text{for } j = 2, \dots, m,$$

$$\int_a^{x_1-\delta_k} |v'_k| dx \rightarrow |v'|([a, x_1]), \quad \int_{x_m+\delta_k}^b |v'_k| dx \rightarrow |v'|([x_m + \delta_k, b]).$$

By (2.52) and by the dominated convergence theorem we have

$$(2.55) \quad u_+(x_{j-1}) + \int_{x_{j-1}+\delta_k}^{x_j-\delta_k} w_k(s) ds \longrightarrow u_+(x_{j-1}) + \int_{x_{j-1}}^{x_j} (u')^a ds = u_-(x_j)$$

for $j = 2, \dots, m$, with the obvious changes for $j = 1$ and $j = m + 1$.

To deal with the jump point x_j , assume first that

$$(2.56) \quad u_+(x_j) - u_-(x_j) > 0.$$

In this case, we need to construct functions $f_{kj} \in C^2([x_j - \delta_k, x_j + \delta_k])$ that satisfy the following properties: There exists $y_{kj} \in]x_j - \delta_k, x_j + \delta_k[$ such that

$$(2.57) \quad f_{kj}(x_j - \delta_k) = u_+(x_{j-1}) + \int_{x_{j-1}+\delta_k}^{x_j-\delta_k} w_k(s) ds, \quad f_{kj}(x_j + \delta_k) = u_+(x_j),$$

$$(2.58) \quad f'_{kj}(x_j - \delta_k) = w_k(x_j - \delta_k), \quad f'_{kj}(x_j + \delta_k) = w_k(x_j + \delta_k),$$

$$(2.59) \quad f''_{kj}(x) > 0 \text{ for } x \in]x_j - \delta_k, y_{kj}[, \quad f''_{kj}(x) < 0 \text{ for } x \in]y_{kj}, x_j + \delta_k[,$$

$$(2.60) \quad \left| f_{kj}(x_j - \delta_k) - \min_{[x_j-\delta_k, y_{kj}]} f_{kj} \right| \leq \frac{1}{k}, \quad \left| f_{kj}(x_j + \delta_k) - \max_{[y_{kj}, x_j+\delta_k]} f_{kj} \right| \leq \frac{1}{k},$$

where we replace x_{j-1} and $x_{j-1} - \delta_k$ by a in the case $j = 1$.

We now discuss briefly the existence of such functions. We observe that the latter conditions in (2.57)–(2.59) imply that the graph of f_{kj} in the interval $[y_{kj}, x_j + \delta_k]$ lies below the straight line passing through the point $(x_j + \delta_k, u_+(x_j))$ with slope $w_k(x_j + \delta_k)$, i.e.,

$$f_{kj}(x) \leq u_+(x_j) + w_k(x_j + \delta_k)(x - x_j - \delta_k)$$

for $x \in [y_{kj}, x_j + \delta_k]$. It is then easy to see that the inequality

$$(2.61) \quad u_+(x_j) - 2w_k(x_j + \delta_k)\delta_k - u_+(x_{j-1}) - \int_{x_{j-1}+\delta_k}^{x_j-\delta_k} w_k(s) ds > 0$$

allows us to fulfill also the former conditions in (2.57)–(2.59), as well as (2.60). By (2.53), (2.55), and (2.56), inequality (2.61) is satisfied when δ_k is small enough.

If the left-hand side of (2.56) is negative, then we choose f_{kj} so that (2.57) and (2.58) hold, and there exists $y_{kj} \in]x_j - \delta_k, x_j + \delta_k[$ such that

$$f''_{kj}(x) < 0 \text{ for } x \in]x_j - \delta_k, y_{kj}[, \quad f''_{kj}(x) > 0 \text{ for } x \in]y_{kj}, x_j + \delta_k[,$$

$$\left| f_{kj}(x_j - \delta_k) - \max_{[x_j-\delta_k, x_j+\delta_k]} f_{kj} \right| \leq \frac{1}{k}, \quad \left| f_{kj}(x_j + \delta_k) - \min_{[x_j-\delta_k, x_j+\delta_k]} f_{kj} \right| \leq \frac{1}{k}.$$

In the same way the construction is possible if δ_k is small enough.

We are now ready to define the approximating sequence

$$u_k(x) := \begin{cases} u_+(a) + \int_a^x w_k(s) ds & \text{if } a \leq x < x_1 - \delta_k, \\ f_{kj}(x) & \text{if } x_j - \delta_k \leq x < x_j + \delta_k, j = 1, \dots, m, \\ u_+(x_{j-1}) + \int_{x_{j-1} + \delta_k}^x w_k(s) ds & \text{if } x_{j-1} + \delta_k \leq x < x_j - \delta_k, j = 2, \dots, m, \\ u_+(x_m) + \int_{x_m + \delta_k}^x w_k(s) ds & \text{if } x_m + \delta_k \leq x < b. \end{cases}$$

Let us define $x_0 := a$ and $x_{m+1} := b$. Since $w_k \rightarrow (u')^a$ in $L^1(]a, b[)$, we have

$$u_k(x) \rightarrow u_+(x_{j-1}) + \int_{x_{j-1}}^x (u')^a(s) ds = u(x)$$

for every $x \in]x_{j-1}, x_j[$ and $j = 1, \dots, m + 1$ and, in turn, $u_k \rightarrow u$ in $L^1(]a, b[)$. As

$$\begin{aligned} \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} |u'_k| dx + \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} \psi(u'_k) |u''_k| dx &= \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} |w_k| dx + \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} \psi(w_k) |w'_k| dx \\ &\leq \int_{x_{j-1}}^{x_j} |w_k| dx + \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} |v'_k| dx, \end{aligned}$$

by (2.54) and the fact that $w_k \rightarrow (u')^a$ in $L^1(]a, b[)$ we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} \left(\int_{x_{j-1} + \delta_k}^{x_j - \delta_k} |u'_k| dx + \int_{x_{j-1} + \delta_k}^{x_j - \delta_k} \psi(u'_k) |u''_k| dx \right) \\ (2.62) \qquad \qquad \qquad \leq \int_{x_{j-1}}^{x_j} |(u')^a| dx + |v'| (]x_{j-1}, x_j]). \end{aligned}$$

Similarly,

$$\begin{aligned} \limsup_{k \rightarrow \infty} \left(\int_a^{x_1 - \delta_k} |u'_k| dx + \int_a^{x_1 - \delta_k} \psi(u'_k) |u''_k| dx \right) \\ (2.63) \qquad \qquad \qquad \leq \int_a^{x_1} |(u')^a| dx + |v'| (]a, x_1]), \end{aligned}$$

$$\begin{aligned} \limsup_{k \rightarrow \infty} \left(\int_{x_m + \delta_k}^b |u'_k| dx + \int_{x_m + \delta_k}^b \psi(u'_k) |u''_k| dx \right) \\ (2.64) \qquad \qquad \qquad \leq \int_{x_m}^b |(u')^a| dx + |v'| (]x_m, b]). \end{aligned}$$

Assume that $[u](x_j) = u_+(x_j) - u_-(x_j) > 0$. Then (2.56) holds for the functions u_k

if k is sufficiently large. By (2.58), (2.59), (2.60), and a change of variables we obtain

$$\begin{aligned}
 & \int_{x_j - \delta_k}^{x_j + \delta_k} |u'_k| \, dx + \int_{x_j - \delta_k}^{x_j + \delta_k} \psi(u'_k) |u''_k| \, dx \\
 &= \int_{x_j - \delta_k}^{x_j + \delta_k} |f'_{kj}| \, dx + \int_{x_j - \delta_k}^{x_j + \delta_k} \psi(f'_{kj}) |f''_{kj}| \, dx \\
 &\leq f_{kj}(x_j + \delta_k) - f_{kj}(x_j - \delta_k) + \int_{w_k(x_j - \delta_k)}^{f'_{kj}(y_{kj})} \psi(t) \, dt \\
 &\quad + \int_{w_k(x_j + \delta_k)}^{f'_{kj}(y_{kj})} \psi(t) \, dt + \frac{2}{k}.
 \end{aligned}
 \tag{2.65}$$

By (2.59) we have

$$f'_{kj}(y_{kj}) = \max_{[x_j - \delta_k, x_j + \delta_k]} f'_{kj} \geq \frac{1}{2\delta_k} [f_{kj}(x_j + \delta_k) - f_{kj}(x_j - \delta_k)].
 \tag{2.66}$$

By (2.57) and the fact that $w_k \rightarrow (u')^a$ in $L^1([a, b])$ we obtain

$$f_{kj}(x_j + \delta_k) - f_{kj}(x_j - \delta_k) \rightarrow u_+(x_j) - \left(u_+(x_{j-1}) + \int_{x_{j-1}}^{x_j} (u')^a \, ds \right) = [u](x_j).$$

In turn, using (2.66), we get that $f'_{kj}(y_{kj}) \rightarrow \infty$. Thus, letting $k \rightarrow \infty$ in (2.65) and using (2.53), we infer

$$\begin{aligned}
 & \limsup_{k \rightarrow \infty} \left(\int_{x_j - \delta_k}^{x_j + \delta_k} |u'_k| \, dx + \int_{x_j - \delta_k}^{x_j + \delta_k} \psi(u'_k) |u''_k| \, dx \right) \\
 &\leq [u](x_j) + \int_{(u')^-_+(x_j)}^{+\infty} \psi(t) \, dt + \int_{(u')^+_+(x_j)}^{+\infty} \psi(t) \, dt.
 \end{aligned}
 \tag{2.67}$$

Similarly, if $[u](x_j) = u_+(x_j) - u_-(x_j) < 0$, we find

$$\begin{aligned}
 & \limsup_{k \rightarrow \infty} \left(\int_{x_j - \delta_k}^{x_j + \delta_k} |u'_k| \, dx + \int_{x_j - \delta_k}^{x_j + \delta_k} \psi(u'_k) |u''_k| \, dx \right) \\
 &\leq |[u](x_j)| + \int_{-\infty}^{(u')^-_-(x_j)} \psi(t) \, dt + \int_{-\infty}^{(u')^+_-(x_j)} \psi(t) \, dt.
 \end{aligned}
 \tag{2.68}$$

Summing over j in (2.62), (2.67), and (2.68) and combining with (2.63) and (2.64), inequality (2.51) follows.

Step 2. Assume only that $u \in X^1_\psi([a, b])$ and that S_u is finite. We claim that there exists a sequence $\{u_k\}$ such that $u_k \rightarrow u$ in $L^1([a, b])$, each u_k satisfies the hypotheses of Step 1, and

$$\mathcal{G}(u) \geq \limsup_{k \rightarrow \infty} \mathcal{G}(u_k).
 \tag{2.69}$$

Note that, if (2.69) holds, then, by applying Step 1 to each u_k , we may find a sequence $u_{km} \in W^{2,1}([a, b])$ converging to u_k in $L^1([a, b])$ and satisfying

$$\mathcal{G}(u_k) \geq \limsup_{m \rightarrow \infty} \mathcal{F}_1(u_{km}).$$

By (2.69) we then have

$$\mathcal{G}(u) \geq \limsup_{k \rightarrow \infty} \limsup_{m \rightarrow \infty} \mathcal{F}_1(u_{km}),$$

and a standard diagonalization argument now yields the existence of a sequence $m_k \rightarrow \infty$ such that $u_{km_k} \rightarrow u$ in $L^1(a, b]$ and

$$\mathcal{G}(u) \geq \limsup_{k \rightarrow \infty} \mathcal{F}_1(u_{km_k}).$$

In the construction of the sequence satisfying (2.69), we need to consider the precise representatives $(u')^\nabla$ and $(u')^\wedge$ defined in Remark 2.2. We recall that $(u')^\nabla$ is upper semicontinuous while $(u')^\wedge$ is lower semicontinuous, and so for each $k \in \mathbb{N}$ we may decompose the open sets $\{(u')^\wedge > k\}$ and $\{(u')^\nabla < -k\}$ into the union of two finite sequences of pairwise disjoint open sets U_{kj}^+ and U_{kj}^- , that is,

$$\bigcup_j U_{kj}^+ = \{(u')^\wedge > k\}, \quad \bigcup_j U_{kj}^- = \{(u')^\nabla < -k\},$$

such that

$$(2.70) \quad \text{diam}(U_{kj}^+) \leq \mathcal{L}^1(\{(u')^\wedge > k\}), \quad \text{diam}(U_{kj}^-) \leq \mathcal{L}^1(\{(u')^\nabla < -k\})$$

for every j . Note that, setting $v_\nabla := \Psi_1 \circ (u')^\nabla$ and $v_\wedge := \Psi_1 \circ (u')^\wedge$, we have

$$(2.71) \quad |v'|]c, d[= \text{Var}(v_\nabla;]c, d[) = \text{Var}(v_\wedge;]c, d[)$$

for every interval $]c, d[\subset]a, b[$.

For every set U_{kj}^\pm , we fix a nonnegative function $g_{kj}^\pm \in C_c^1(U_{kj}^\pm)$ such that

$$(2.72) \quad \int_{U_{kj}^\pm} g_{kj}^\pm(x) \, dx = ((u')^c)^\pm(U_{kj}^\pm),$$

and $(g_{kj}^\pm)'$ has only one zero in the interior of the support of g_{kj}^\pm . Then we define

$$(2.73) \quad g_k^+ := \sum_j g_{kj}^+, \quad g_k^- := \sum_j g_{kj}^-, \quad g_k := g_k^+ - g_k^-, \quad w_k := T_{-k}^k \circ (u')^a + g_k,$$

where, for any pair of constants $h < k$, the truncation function T_h^k is defined by

$$T_h^k(t) := \begin{cases} h & \text{for } t \leq h, \\ t & \text{for } h \leq t \leq k, \\ k & \text{for } t \geq k. \end{cases}$$

We claim that

$$(2.74) \quad w_k \mathcal{L}^1 \rightharpoonup (u')^a \mathcal{L}^1 + (u')^c \text{ weakly* in } M_b(a, b].$$

Define

$$A_k := \{(u')^\wedge > k\} \cup \{(u')^\nabla < -k\}.$$

Since by the Chebychev inequality

$$(2.75) \quad k\mathcal{L}^1(A_k) \rightarrow 0,$$

it suffices to show that

$$(2.76) \quad \left(\sum_j g_{kj}^\pm \right) \mathcal{L}^1 \rightharpoonup ((u')^c)^\pm \text{ weakly* in } M_b(]a, b[).$$

Let $\varphi \in C_0(]a, b[)$ and $\varepsilon > 0$. By uniform continuity, there exists $\delta = \delta(\varepsilon) > 0$ such that $|\varphi(x) - \varphi(y)| \leq \varepsilon$ for all $x, y \in]a, b[$ with $|x - y| \leq \delta$. In view of (2.70) and (2.75), for all k sufficiently large and for all j we have that $\text{diam}(U_{kj}^\pm) \leq \delta$. Let us fix $y_{kj}^\pm \in U_{kj}^\pm$. Then by (2.72)

$$\begin{aligned} & \left| \int_{U_{kj}^\pm} \varphi(x) g_{kj}^\pm(x) dx - \int_{U_{kj}^\pm} \varphi(x) d((u')^c)^\pm(x) \right| \\ &= \left| \int_{U_{kj}^\pm} [\varphi(x) - \varphi(y_{kj}^\pm)] g_{kj}^\pm(x) dx - \int_{U_{kj}^\pm} [\varphi(x) - \varphi(y_{kj}^\pm)] d((u')^c)^\pm(x) \right| \\ &\leq \varepsilon \left(\int_{U_{kj}^\pm} g_{kj}^\pm(x) dx + ((u')^c)^\pm(U_{kj}^\pm) \right) \leq 2\varepsilon ((u')^c)^\pm(U_{kj}^\pm). \end{aligned}$$

Summing over j and using the fact that the measures $(\sum_j g_{kj}^+) \mathcal{L}^1$ and $((u')^c)^+$ are concentrated on $\{(u')_\wedge^a > k\}$, while the measures $(\sum_j g_{kj}^-) \mathcal{L}^1$ and $((u')^c)^-$ are concentrated on $\{(u')_\vee^a < -k\}$ (see Remark 2.2), we obtain (2.76).

Moreover, we claim that

$$(2.77) \quad \lim_{k \rightarrow \infty} \int_a^b |w_k| dx = \int_a^b |(u')^a| dx + |(u')^c|(]a, b[).$$

Indeed, using (2.72), (2.73), and Remark 2.2, we deduce that

$$\begin{aligned} \int_a^b |w_k| dx &\leq \int_{\{|(u')^a| \leq k\}} |(u')^a| dx + k\mathcal{L}^1(A_k) + \sum_j \int_{U_{kj}^+} g_{kj}^+ dx + \sum_j \int_{U_{kj}^-} g_{kj}^- dx \\ &\leq \int_a^b |(u')^a| dx + k\mathcal{L}^1(A_k) + \sum_j ((u')^c)^+(U_{kj}^+) + \sum_j ((u')^c)^-(U_{kj}^-) \\ &\leq \int_a^b |(u')^a| dx + k\mathcal{L}^1(A_k) + |(u')^c|(]a, b[), \end{aligned}$$

and the limit superior inequality follows from (2.75). The limit inferior inequality follows from (2.74) and the lower semicontinuity of the total variation.

Set

$$(2.78) \quad u_k(x) := u_+(a) + \int_a^x w_k(s) ds + \sum_{x_j < x, x_j \in S_u} [u](x_j)$$

and $v_k := \Psi_1 \circ (u'_k)^a = \Psi_1 \circ w_k$.

We claim that $u_k \rightarrow u$ in $L^1]a, b[$. For $x \in]a, b[$ by (2.74) and (2.77) it follows that

$$\int_a^x w_k dy \rightarrow \int_a^x (u')^a dy + (u')^c]a, x[,$$

and so u_k converges to u pointwise \mathcal{L}^1 -a.e. and, in turn, in $L^1]a, b[$.

Next, we show that

$$(2.79) \quad \limsup_{k \rightarrow \infty} \mathcal{G}(u_k) \leq \mathcal{G}(u) .$$

From (2.77) we get

$$(2.80) \quad |u'_k|]a, b[\rightarrow |u'|]a, b[.$$

Moreover, as $v_k = T_{\Psi_1(-k)}^{\Psi_1(k)} \circ v$ \mathcal{L}^1 -a.e. in the open set $V_k :=]a, b[\setminus \text{supp } g_k$, we have $|v'_k| \leq |v'|$ as measures in V_k . In particular, this yields $|v'_k|]a, b[\setminus (A_k \cup S_u) \leq |v'|]a, b[\setminus (A_k \cup S_u)$ and, hence,

$$(2.81) \quad |v'_k|]a, b[\setminus (A_k \cup S_u) \leq |v'|]a, b[\setminus (A_\infty \cup S_u) ,$$

where

$$A_\infty := \bigcap_k A_k = \{(u')^a_\wedge = +\infty\} \cup \{(u')^a_\vee = -\infty\} .$$

Using the properties of g_{kj}^+ , we have

$$(2.82) \quad \begin{aligned} |v'_k| (\{(u')^a_\wedge > k\} \setminus S_u) &= \sum_j \int_{U_{kj}^+} \psi(k + g_{kj}^+) \left| (g_{kj}^+)' \right| dx \\ &= 2 \sum_j \int_k^{k + \sup g_{kj}^+} \psi(t) dt \leq 2\mathcal{H}^0 \left(\{j : ((u')^c)^+ (U_{kj}^+) > 0\} \right) \int_k^\infty \psi(t) dt . \end{aligned}$$

We claim that

$$\begin{aligned} 2\mathcal{H}^0 \left(\{j : ((u')^c)^+ (U_{kj}^+) > 0\} \right) \int_k^\infty \psi(t) dt \\ \leq |v'| (\{(u')^a_\wedge > k\} \setminus S_u) + 4 \int_k^\infty \psi(t) dt . \end{aligned}$$

Indeed, if $((u')^c)^+(U_{kj}^+) > 0$, then there exists a connected component $I_{kj}^+ =]a_{kj}, b_{kj}[$ of $U_{kj}^+ \setminus S_u$ such that $((u')^c)^+(I_{kj}^+) > 0$. Assume that $I_{kj}^+ \subset\subset]a, b[$. Then by Remark 2.2 we may find $c_{kj} \in I_{kj}^+$ such that $(u')^a_\wedge(c_{kj}) = +\infty$, while $(u')^a_\wedge(a_{kj}), (u')^a_\wedge(b_{kj}) \leq k$. Hence, by (2.71)

$$|v'| (U_{kj}^+ \setminus S_u) \geq |v'| (I_{kj}^+) \geq 2 \int_k^\infty \psi(t) dt .$$

Summing over all such intervals and adding the possible contribution of the intervals I_{kj}^+ with at least one end point in $\{a, b\}$, we obtain the claim. In turn, by (2.82) we have

$$|v'_k| (\{(u')^a_\wedge > k\} \setminus S_u) \leq |v'| (\{(u')^a_\wedge > k\} \setminus S_u) + 4 \int_k^\infty \psi(t) dt .$$

A similar estimate holds for the set $\{(u')^a_{\nabla} < -k\} \setminus S_u$, thus, yielding

$$(2.83) \quad \limsup_{k \rightarrow \infty} |v'_k| (A_k \setminus S_u) \leq |v'| (A_\infty \setminus S_u) .$$

Combining (2.81) with (2.83) we obtain

$$\limsup_{k \rightarrow \infty} |v'_k| (]a, b[\setminus S_u) \leq |v'| (]a, b[\setminus S_u) .$$

Next, we show that

$$(2.84) \quad \lim_{k \rightarrow \infty} \sum_{x \in S_{u_k}} \Phi (\nu_{u_k}, (u'_k)^a_-, (u'_k)^a_+) = \sum_{x \in S_u} \Phi (\nu_u, (u')^a_-, (u')^a_+) .$$

Note that $S_{u_k} = S_u$ and $\nu_{u_k}(x) = \nu_u(x)$ for all k by (2.78). Moreover, for every $x \in S_u$ if $(u')^a_+(x) \in \mathbb{R}$, then $|(u')^a_+(y)| \leq k_0$ for all y in a right neighborhood of x and for some integer k_0 . Thus, by (2.73) and (2.78) we have that $(u'_k)^a(y) = (u')^a(y)$ for $k \geq k_0$ and for \mathcal{L}^1 -a.e. y in the same right neighborhood. In turn, by (2.7) we infer $(u'_k)^a_+(x) = (u')^a_+(x)$ for all $k \geq k_0$. If $(u')^a_+(x) = \infty$, then for all k we have $(u')^a_+ > k$ in a right neighborhood of x by right continuity (see Remark 2.2). By construction this implies that $(u'_k)^a_+ = w_k \geq k$ \mathcal{L}^1 -a.e. in the same right neighborhood. Thus, $(u'_k)^a_+(x) \geq k \rightarrow (u')^a_+(x)$. Similarly, $(u'_k)^a_-(x) \rightarrow (u')^a_-(x)$ so that

$$\Phi (\nu_{u_k}(x), (u'_k)^a_-(x), (u'_k)^a_+(x)) \rightarrow \Phi (\nu_u(x), (u')^a_-(x), (u')^a_+(x)) .$$

Hence, (2.84) follows. This, together with (2.80) and (2.83), yields (2.79).

Step 3. Let now u be an arbitrary function in $X^1_\psi(]a, b[)$ such that $\mathcal{G}(u) < +\infty$. As in the previous step it suffices to construct $u_k \in X^1_\psi(]a, b[)$ satisfying the hypotheses of Step 2, converging to u in $L^1(]a, b[)$ and such that (2.79) holds. Write $S_u = \{x_j\}$, and for each k define $S^k_u := \{x_j : j \leq k\}$ and

$$u_k(x) = u_+(a) + \int_a^x (u')^a dt + (u')^c(]a, x]) + \sum_{x_j < x, x_j \in S^k_u} [u](x_j) .$$

It is clear that $\{u_k\}$ converges to u in $L^1(]a, b[)$ and that $|u'_k|(]a, b[) \rightarrow |u'|(]a, b[)$. Moreover, $|v'_k|(]a, b[\setminus S_u) = |v'| (]a, b[\setminus S_u)$ and

$$\begin{aligned} \lim_{k \rightarrow \infty} \sum_{x \in S_{u_k}} \Phi (\nu_{u_k}, (u'_k)^a_-, (u'_k)^a_+) &= \lim_{k \rightarrow \infty} \sum_{x \in S^k_u} \Phi (\nu_u, (u')^a_-, (u')^a_+) \\ &= \sum_{x \in S_u} \Phi (\nu_u, (u')^a_-, (u')^a_+) . \end{aligned}$$

This concludes the proof of the theorem. \square

We end the section with a compactness result for energy bounded sequences in $X^1_\psi(]a, b[)$.

COROLLARY 2.7. *Let $\{u_k\}$ be a sequence of functions in $X^1_\psi(]a, b[)$ bounded in $L^1(]a, b[)$ and such that*

$$(2.85) \quad C := \sup_k \overline{\mathcal{F}}_1(u_k) < +\infty .$$

Then there exist a subsequence (not relabeled) $\{u_k\}$ and a function $u \in X^1_\psi(]a, b[)$ such that

$$(2.86) \quad u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(]a, b[),$$

$$(2.87) \quad \Psi_1 \circ (u'_k)^a \rightharpoonup \Psi_1 \circ (u')^a \quad \text{weakly}^* \text{ in } BV(]a, b[),$$

$$(u'_k)^a \rightarrow (u')^a \quad \text{pointwise } \mathcal{L}^1\text{-a.e. in }]a, b[.$$

Proof. It is well known that convergence in measure is metrizable with the following metric:

$$d(u_1, u_2) := \int_a^b \frac{|u_1 - u_2|}{1 + |u_1 - u_2|} dx,$$

where u_1 and u_2 are (equivalent classes of) measurable functions.

By Theorems 2.4 and 2.5, for every $k \in \mathbb{N}$ we may find $w_k \in W^{2,1}(]a, b[)$ such that

$$(2.88) \quad \int_a^b |u_k - w_k| dx \leq \frac{1}{k}, \quad d((u'_k)^a, w'_k) \leq \frac{1}{k},$$

and

$$\mathcal{F}_1(w_k) \leq C + 1.$$

By Theorem 2.4 we may find a subsequence (not relabeled) of $\{w_k\}$ and a function $u \in X^1_\psi(]a, b[)$ such that (2.15), (2.16), and (2.17) hold (with w_k in place of u_k). It now follows from (2.88) that $u_k \rightarrow u$ in $L^1(]a, b[)$ and $(u'_k)^a \rightarrow (u')^a$ in measure and, hence, pointwise \mathcal{L}^1 -a.e. in $]a, b[$, up to a further subsequence. From the bound (2.85), the uniqueness of the limit, and the invertibility of Ψ_1 , we deduce (2.86) and (2.87). \square

3. The case $p > 1$. In this section we analyze the functional (1.3) in the case $p > 1$.

Let us state precisely the standing assumptions. Throughout this section p denotes any exponent in $]1, +\infty[$, $\psi: \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying

$$(3.1) \quad M := \int_{-\infty}^{+\infty} (\psi(t))^{1/p} dt < +\infty$$

in addition to (2.2), and $\Psi_p: \overline{\mathbb{R}} \rightarrow [0, M]$ denotes the antiderivative of $\psi^{1/p}$ defined by

$$(3.2) \quad \Psi_p(t) := \int_{-\infty}^t (\psi(s))^{1/p} ds.$$

The function $\Psi_p^{-1}: [0, M] \rightarrow \overline{\mathbb{R}}$ stands for the inverse function of Ψ_p .

We now consider the functional $\mathcal{F}_p: L^1(]a, b[) \rightarrow [0, +\infty]$ defined by

$$(3.3) \quad \mathcal{F}_p(u) := \begin{cases} \int_a^b |u'| dx + \int_a^b \psi(u') |u''|^p dx & \text{if } u \in W^{2,p}(]a, b[), \\ +\infty & \text{otherwise.} \end{cases}$$

It turns out that piecewise smooth functions with bounded derivative and nonempty discontinuity set cannot be approximated by sequences with equibounded energy. This is a consequence of Remark 3.2(i) and Theorem 3.3 below, and to this end we introduce a suitable space of functions. Recall that $Z^\pm[(u')^a]$ are the sets defined in (2.5) and (2.6), while $(u')^s$ denotes the singular part of the gradient measure u' .

DEFINITION 3.1. Let $X_\psi^p(]a, b[)$ be the set of all functions $u \in BV(]a, b[)$ such that $v := \Psi_p \circ (u')^a$ belongs to $W^{1,p}(]a, b[)$, and the positive part $((u')^s)^+$ and the negative part $((u')^s)^-$ of the measure $(u')^s$ are concentrated on $Z^+[(u')^a]$ and $Z^-[(u')^a]$, respectively.

Remark 3.2. (i) It follows immediately from the definition that if $u \in X_\psi^p(]a, b[)$, then $(u')^a = \Psi_p^{-1}(v)$ is continuous on $[a, b]$ with values in \mathbb{R} . In particular, it turns out that

$$Z^\pm[(u')^a] = \{x \in]a, b[: (u')^a = \pm\infty\}.$$

By the hypothesis on the support of the singular part $(u')^s$, we have $\lim_{x \rightarrow x_0} (u')^a(x) = +\infty$ for every jump point x_0 with $u_+(x_0) - u_-(x_0) > 0$ and $\lim_{x \rightarrow x_0} (u')^a(x) = -\infty$ for every jump point x_0 with $u_+(x_0) - u_-(x_0) < 0$. This means that if S_u is nonempty, then u cannot have bounded derivative outside the discontinuity set. In particular, piecewise constant functions are not included in the class $X_\psi^p(]a, b[)$.

(ii) We observe that the function $(u')^a$ is differentiable \mathcal{L}^1 -a.e. in $]a, b[$ with

$$(3.4) \quad v' = \psi^{\frac{1}{p}}((u')^a)((u')^a)'$$

To see this, we consider the open set

$$A_k := \{x \in]a, b[: -k < (u')^a < k\}.$$

Since by (2.2) the function Ψ_p^{-1} is Lipschitz continuous in the interval $[\Psi_p(-k), \Psi_p(k)]$ and $v \in W^{1,p}(]a, b[)$, by the chain rule, we have that $(u')^a = \Psi_p^{-1} \circ v \in W^{1,p}(A_k)$ and, in particular, it is differentiable \mathcal{L}^1 -a.e. in A_k , and (3.4) holds. Since $(u')^a$ is integrable, we have that

$$\mathcal{L}^1\left(]a, b[\setminus \bigcup_k A_k\right) = 0,$$

and the conclusion follows.

(iii) It is easy to check that $X_\psi^p(]a, b[)$ may contain discontinuous functions. An example is given by the following construction: Let $\psi: \mathbb{R} \rightarrow]0, +\infty[$ be defined by

$$\psi(t) := \begin{cases} 1 & \text{if } |t| \leq 1, \\ \frac{1}{|t|^\alpha} & \text{if } |t| > 1, \end{cases}$$

where α is any number in $]1, +\infty[$, and let $p \in]1, \frac{\alpha+1}{2}[$. Consider now the discontinuous functions $u:]-1, 1[\rightarrow \mathbb{R}$ given by

$$u(x) := \begin{cases} -|x|^\beta & \text{if } x \leq 0, \\ 1 + x^\beta & \text{if } x > 0, \end{cases}$$

with

$$0 < \beta < 1 - \frac{p-1}{\alpha-p}.$$

A simple computation shows that the function $\Psi_p \circ (u')^\alpha$ belongs to $W^{1,p}(-1, 1)$, which in turn implies that $u \in X_\psi^p(-1, 1)$.

(iv) Finally, the same construction of Proposition 2.3 shows that, for every admissible ψ satisfying (2.9), the space $X_\psi^p(a, b)$ contains a function with nontrivial Cantor part if p is sufficiently close to 1. We omit the details of this fact, which can be easily checked following step by step the proof of Proposition 2.3.

The next theorem is the counterpart of Theorem 2.4 for the case $p > 1$. It establishes that energy bounded sequences are relatively compact in $X_\psi^p(a, b)$. The proof is similar to the one of Theorem 2.4; nevertheless, since this is the main result of this section we reproduce it in full detail for the reader's convenience.

THEOREM 3.3. *Let $\{u_k\}$ be a sequence of functions bounded in $L^1(a, b)$ and such that*

$$(3.5) \quad C := \sup_k \mathcal{F}_p(u_k) < +\infty.$$

Then there exist a subsequence (not relabeled) $\{u_k\}$ and a function $u \in X_\psi^p(a, b)$ such that

$$(3.6) \quad u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(a, b),$$

$$(3.7) \quad \begin{aligned} \Psi_p \circ u'_k &\rightharpoonup \Psi_p \circ (u')^\alpha \quad \text{weakly in } W^{1,p}(a, b), \\ u'_k &\rightarrow (u')^\alpha \quad \text{pointwise in }]a, b[. \end{aligned}$$

Proof. By (3.3) and (3.5) we may assume that each u_k belongs to $W^{2,p}(a, b)$ and that

$$(3.8) \quad C_1 := \sup_k \int_a^b [|u_k| + |u'_k| + \psi(u'_k) |u''_k|^p] dx < +\infty.$$

Let us define

$$(3.9) \quad v_k := \Psi_p \circ u'_k.$$

As Ψ_p is Lipschitz in \mathbb{R} , the functions v_k belong to $W^{1,p}(a, b)$ and

$$(3.10) \quad v'_k = (\psi(u'_k))^{1/p} u''_k \quad \mathcal{L}^1\text{-a.e. on }]a, b[.$$

It follows from (3.1) and (3.5) that

$$(3.11) \quad \int_a^b [|v_k|^p + |v'_k|^p] dx \leq M^p(b-a) + C_1.$$

By (3.8) and (3.11), passing to a subsequence (not relabeled), we may assume that

$$u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(a, b)$$

and

$$(3.12) \quad v_k \rightharpoonup v \quad \text{weakly in } W^{1,p}(a, b)$$

for some functions $u \in BV(a, b)$ and $v \in W^{1,p}(a, b; [0, M])$.

Since Ψ_p^{-1} is continuous, we obtain

$$(3.13) \quad u'_k = \Psi_p^{-1} \circ v_k \rightarrow w := \Psi_p^{-1} \circ v \quad \text{pointwise in }]a, b[.$$

Note, also, that w is continuous with values in $\overline{\mathbb{R}}$.

We now split the remaining part of the proof into two steps.

Step 1. We prove that

$$(3.14) \quad w = (u')^a \quad \mathcal{L}^1\text{-a.e. on }]a, b[.$$

If not, arguing as for (2.27), we may find $t_0 > 0$ and an infinite number of disjoint open intervals I such that

$$(3.15) \quad \mathcal{L}^1(\{w \neq (u')^a\} \cap \{|w| < t_0\} \cap I) > 0.$$

By Hölder’s inequality and a change of variables, we obtain

$$(3.16) \quad \begin{aligned} \int_I \psi(u'_k) |u''_k|^p dx &\geq \frac{1}{\mathcal{L}^1(I)^{p-1}} \left(\int_I (\psi(u'_k))^{1/p} |u''_k| dx \right)^p \\ &\geq \frac{1}{(b-a)^{p-1}} \left(\int_{m_k}^{M_k} (\psi(t))^{1/p} dt \right)^p, \end{aligned}$$

where $m_k := \inf_I u'_k$ and $M_k := \sup_I u'_k$.

Reasoning as in the first step of the proof of Theorem 2.4, we can show that at least one of the two sequences $\{m_k\}$ and $\{M_k\}$ is divergent. If $\lim_k M_k = +\infty$, then by (3.13) $\limsup_k m_k < t_0$ and, in turn, from (3.16) we obtain

$$\liminf_{k \rightarrow \infty} \int_I \psi(u'_k) |u''_k|^p dx \geq \frac{1}{(b-a)^{p-1}} \left(\int_{t_0}^{+\infty} (\psi(t))^{1/p} dt \right)^p > 0.$$

Analogously, if $\lim_k m_k = -\infty$, then

$$\liminf_{k \rightarrow \infty} \int_I \psi(u'_k) |u''_k|^p dx \geq \frac{1}{(b-a)^{p-1}} \left(\int_{-\infty}^{t_0} (\psi(t))^{1/p} dt \right)^p > 0.$$

In any case, for an arbitrarily large number m of disjoint intervals I satisfying (3.15), adding the contributions of each interval we obtain

$$\begin{aligned} \liminf_{k \rightarrow \infty} \int_a^b \psi(u'_k) |u''_k|^p dx \\ \geq \frac{m}{(b-a)^{p-1}} \min \left\{ \left(\int_{t_0}^{+\infty} (\psi(t))^{1/p} dt \right)^p, \left(\int_{-\infty}^{t_0} (\psi(t))^{1/p} dt \right)^p \right\}, \end{aligned}$$

which contradicts (3.8) for m large enough. This concludes the proof of (3.14) and, in turn, of (3.7).

Step 2. To prove that $u \in X^p_\psi(]a, b[)$, it remains to show that the positive part $((u')^s)^+$ and the negative part $((u')^s)^-$ of the measure $(u')^s$ are concentrated on $Z^+[(u')^a]$ and $Z^-[(u')^a]$, respectively.

Arguing as in Step 2 of the proof of Theorem 2.4, one can see that it is enough to show

$$(3.17) \quad E^+[u'] \setminus Z^+[(u')^a] \text{ and } E^-[u'] \setminus Z^-[(u')^a] \text{ are empty,}$$

where $E^+[u']$ and $E^-[u']$ are the sets introduced in (2.32) and (2.33). We show only that $E^+[u'] \setminus Z^+[(u')^a]$ is empty, since the other property can be proved in the same way.

Assume, by contradiction, that $E^+[u'] \setminus Z^+[(u')^a]$ contains a point x_0 . Denote $t_0 := 2|w(x_0)|$, fix any $t_1 > t_0$, and choose $\varepsilon_0 > 0$ such that

$$(3.18) \quad \frac{1}{(2\varepsilon_0)^{p-1}} \left(\int_{t_0}^{t_1} (\psi(t))^{1/p} dt \right)^p > C,$$

where C is the constant appearing in (3.5). By (2.32) there exists $0 < \varepsilon < \varepsilon_0$ such that

$$(3.19) \quad \frac{(u')^+(\text{]}x_0 - \varepsilon, x_0 + \varepsilon[})}{2\varepsilon} > t_1.$$

Set $I := \text{]}x_0 - \varepsilon, x_0 + \varepsilon[$. By Hölder’s inequality and a change of variables (see (3.16)), we obtain

$$(3.20) \quad \int_I \psi(u'_k) |u''_k|^p dx \geq \frac{1}{(2\varepsilon_0)^{p-1}} \left(\int_{m_k}^{M_k} (\psi(t))^{1/p} dt \right)^p,$$

where $m_k := \inf_I u'_k$ and $M_k := \sup_I u'_k$. By (3.13) and the fact that $w(x_0) < t_0$, we deduce that

$$(3.21) \quad \limsup_{k \rightarrow \infty} m_k < t_0.$$

On the other hand, reasoning as at the end of the proof of Theorem 2.4, we deduce from (3.6) and (3.19) that

$$\liminf_{k \rightarrow \infty} \frac{1}{2\varepsilon} \int_{x_0-\varepsilon}^{x_0+\varepsilon} (u'_k)^+ dx \geq \frac{(u')^+(\text{]}x_0 - \varepsilon, x_0 + \varepsilon[})}{2\varepsilon} > t_1,$$

which implies that

$$(3.22) \quad \liminf_{k \rightarrow \infty} M_k > t_1.$$

From (3.18), (3.20), (3.21), and (3.22) we obtain

$$\liminf_{k \rightarrow \infty} \int_I \psi(u'_k) |u''_k|^p dx \geq \frac{1}{(2\varepsilon_0)^{p-1}} \left(\int_{t_0}^{t_1} (\psi(t))^{1/p} dt \right)^p > C,$$

which contradicts (3.8). This shows (3.17) and concludes the proof of the theorem. \square

We next identify the relaxation of \mathcal{F}_p with respect to strong convergence in $L^1(\text{]}a, b[)$.

THEOREM 3.4. *Let $\overline{\mathcal{F}}_p : L^1(\text{]}a, b[) \rightarrow [0, +\infty]$ be defined by*

$$(3.23) \quad \overline{\mathcal{F}}_p(u) := \inf \left\{ \liminf_{k \rightarrow \infty} \mathcal{F}_p(u_k) : u_k \rightarrow u \text{ in } L^1(\text{]}a, b[) \right\}$$

for every $u \in L^1(]a, b[)$. Then

$$(3.24) \quad \overline{\mathcal{F}}_p(u) = \begin{cases} |u'|(|a, b[) + \int_a^b |v'|^p dx & \text{if } u \in X_\psi^p(]a, b[), \\ +\infty & \text{otherwise,} \end{cases}$$

where $v := \Psi_p \circ (u')^a$.

Proof. We sketch the proof focusing only on the main changes with respect to the proof of Theorem 2.5. Let \mathcal{G}_p be the functional defined by the right-hand side of (3.24).

We start by showing that

$$(3.25) \quad \mathcal{G}_p(u) \leq \liminf_{k \rightarrow \infty} \mathcal{F}_p(u_k)$$

whenever $u_k \rightarrow u$ in $L^1(]a, b[)$. It is enough to consider sequences $\{u_k\}$ for which the liminf is a limit and has a finite value. Then u_k belongs to $W^{2,1}(]a, b[)$ and (3.5) is satisfied. Setting $v_k := \Psi_p \circ u'_k$, by Theorem 3.3, we have $v_k \rightarrow v$ weakly in $W^{1,p}(]a, b[)$. Using the fact that $|v'_k|^p = \psi(u'_k)|u''_k|^p$, we deduce that

$$(3.26) \quad \int_a^b |v'|^p dx \leq \liminf_{k \rightarrow \infty} \int_a^b \psi(u'_k)|u''_k|^p dx.$$

Inequality (3.25) follows now from (3.26) and the lower semicontinuity of the total variation.

We split the proof of the limsup inequality into several steps.

Step 1. Let $u \in X_\psi^p(]a, b[)$ be such that $(u')^s = 0$. We claim that there exists a sequence $\{u_k\}$ in $W^{2,p}(]a, b[)$ such that $u_k \rightarrow u$ in $L^1(]a, b[)$ and

$$(3.27) \quad \limsup_{k \rightarrow \infty} \mathcal{F}_p(u_k) \leq \mathcal{G}_p(u).$$

Define $w_k := ((u')^a \vee -k) \wedge k$. Using the fact that $(u')^a \in W^{1,p}(A_{2k})$, where

$$A_{2k} := \{x \in]a, b[: -2k < (u')^a < 2k\}$$

as observed in Remark 3.2(ii), one sees that $w_k \in W^{1,p}(]a, b[)$. Define

$$u_k(x) := u_+(a) + \int_a^x w_k(y) dy.$$

It is easy to see that $u_k \rightarrow u$ in $L^1(]a, b[)$ and (3.27) holds.

Step 2. Assume that $u \in X_\psi^p(]a, b[)$, $(u')^c = 0$, and S_u is finite. We claim that there exists a sequence $\{u_k\}$ of functions in $X_\psi^p(]a, b[)$, with $(u'_k)^s = 0$, such that $u_k \rightarrow u$ in $L^1(]a, b[)$ and

$$(3.28) \quad \limsup_{k \rightarrow \infty} \mathcal{G}_p(u_k) \leq \mathcal{G}_p(u).$$

Since the construction is local, it is enough to consider the case $S_u = \{x_0\}$ for some $x_0 \in]a, b[$ with $[u](x_0) > 0$. By the properties of $X_\psi^p(]a, b[)$, we can find two sequences $x_k \nearrow x_0$ and $y_k \searrow x_0$ such that

$$u(x_k) \rightarrow u_-(x_0), \quad u(y_k) \rightarrow u_+(x_0), \quad \text{and } (u')^a(x_k) = (u')^a(y_k) \rightarrow (u')^a(x_0) = +\infty.$$

Consider the affine functions $h_k(x) := u(x_k) + (u')^a(x_k)(x - x_k)$. For every k sufficiently large there exists $z_k \in]x_k, b[$ such that $h_k(z_k) = u_+(x_0)$. Since $(u')^a(x_k) \rightarrow +\infty$ and $x_k \rightarrow x_0$, we have that $z_k \rightarrow x_0$ as $k \rightarrow \infty$. Define

$$u_k(x) := \begin{cases} u(x) & \text{if } a < x \leq x_k, \\ h_k(x) & \text{if } x_k < x \leq z_k, \\ u(x + y_k - z_k) + u_+(x_0) - u(y_k) & \text{if } z_k < x < b. \end{cases}$$

Using the fact that $(u')^a(x_k) = (u')^a(y_k)$, it is easy to check that $u_k \in X_\psi^p(]a, b[)$, with $(u'_k)^s = 0$, $u_k \rightarrow u$ in $L^1(]a, b[)$, and (3.28) holds.

Step 3. Assume that $u \in X_\psi^p(]a, b[)$ and $(u')^c = 0$. We claim that there exists a sequence of functions u_k in $X_\psi^p(]a, b[)$, with $(u'_k)^c = 0$ and S_{u_k} finite, such that $u_k \rightarrow u$ in $L^1(]a, b[)$ and (3.28) holds.

To see this, it is enough to consider the same approximation constructed in Step 3 of the proof of Theorem 2.5.

Step 4. Assume that $u \in X_\psi^p(]a, b[)$. We claim that there exists a sequence of functions u_k in $X_\psi^p(]a, b[)$, with $(u'_k)^c = 0$, such that $u_k \rightarrow u$ in $L^1(]a, b[)$ and (3.28) holds.

Since $(u')^a$ is continuous from $]a, b[$ into $\overline{\mathbb{R}}$ and integrable (see Remark 3.2), we have that $K := \{x \in]a, b[: |(u')^a| = +\infty\}$ is relatively closed in $]a, b[$ with zero \mathcal{L}^1 measure. Hence, we may find a sequence of open sets $A_k \subset]a, b[$ such that $A_k \searrow K$. Let $\{I_j^k\}_j$ be the collection of all connected components of A_k intersecting K . Let $c_j^k := (u')^s(I_j^k) > 0$. By the properties of $X_\psi^p(]a, b[)$, for every j we may choose $x_j^k \in I_j^k \cap K$ such that $(u')^a(x_j^k) = +\infty$ if $c_j^k > 0$ and $(u')^a(x_j^k) = -\infty$ if $c_j^k < 0$. Define

$$u_k(x) := u_+(a) + \int_a^x (u')^a(y) dy + \sum_{j: x_j^k \leq x} c_j^k.$$

Using the definition of $X_\psi^p(]a, b[)$, one can check that

$$\sum_j c_j^k \delta_{x_j^k} \rightharpoonup (u')^s \quad \text{weakly}^* \text{ in } M_b(]a, b[)$$

and $|\sum_j c_j^k \delta_{x_j^k}|(]a, b[) \rightarrow |(u')^s|(]a, b[)$ as $k \rightarrow \infty$. Using this fact it is easy to see that the sequence $\{u_k\}$ meets all of the requirements.

By combining Steps 1–4 with a diagonal argument, one can finally prove that (3.27) holds for every u in $X_\psi^p(]a, b[)$. \square

COROLLARY 3.5. *Let $\{u_k\}$ be a sequence of functions in $X_\psi^p(]a, b[)$ bounded in $L^1(]a, b[)$ and such that*

$$(3.29) \quad C := \sup_k \overline{\mathcal{F}}_p(u_k) < +\infty.$$

Then there exists a subsequence (not relabeled) $\{u_k\}$ and a function $u \in X_\psi^p(]a, b[)$ such that

$$(3.30) \quad u_k \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(]a, b[),$$

$$(3.31) \quad \Psi_p \circ (u'_k)^a \rightharpoonup \Psi_p \circ (u')^a \quad \text{weakly in } W^{1,p}(]a, b[),$$

$$(3.32) \quad (u'_k)^a \rightarrow (u')^a \quad \text{pointwise in }]a, b[.$$

Proof. With an argument entirely similar to the one used in the proof of Corollary 2.7 we can extract a subsequence $\{u_k\}$ which satisfies (3.30) and (3.31). In turn, (3.31) and the continuity of Ψ_p^{-1} in $\overline{\mathbb{R}}$ imply (3.32). \square

4. The staircase effect. The purpose of this section is to show analytically that the presence of the higher order term in the functional $\overline{\mathcal{F}}$ prevents the occurrence of the so-called *staircase effect*, as opposed to what happens in image reconstructions based on the total variation functional.

4.1. The Rudin–Osher–Fatemi model. We start by showing that staircase-like structures do appear in solutions to the Rudin–Osher–Fatemi problem, i.e., in minimizers for the functional $\text{ROF}_{\lambda,g} : BV([a, b]) \rightarrow \mathbb{R}$ defined by

$$\text{ROF}_{\lambda,g}(w) := |w'|([a, b]) + \lambda \int_a^b (w - g)^2 dx,$$

where $\lambda > 0$ is the *fidelity parameter* and $g \in L^2([a, b])$ is the given “signal” to be processed. This fact is well known and numerically observed in many situations. We provide here a simple analytical example. A different example is provided in [12], and a more detailed analysis on the nature of the staircase effect in the Rudin–Osher–Fatemi model can be found in a recent paper of Caselles, Chambolle, and Novaga [4]. Our example will be constructed by means of the following proposition which deals with minimizers of $\text{ROF}_{\lambda,g}$ when g is a monotone function.

PROPOSITION 4.1. *Let $g : [a, b] \rightarrow [0, 1]$ be a nondecreasing function such that $g_+(a) = 0$ and $g_-(b) = 1$. Let g^{-1} denote the left-continuous generalized inverse of g , defined by*

$$(4.1) \quad g^{-1}(c) := \inf\{x \in [a, b] : g(x) \geq c\}$$

for every $c \in [0, 1]$, and assume that there exist $0 < c_1 < c_2 < 1$ such that

$$(4.2) \quad 2\lambda \int_a^{g^{-1}(c_1)} (c_1 - g(x)) dx = 1 \quad \text{and} \quad 2\lambda \int_{g^{-1}(c_2)}^b (g(x) - c_2) dx = 1.$$

Then the function u , defined by

$$u(x) := \begin{cases} c_1 & \text{if } a \leq x \leq g^{-1}(c_1), \\ g(x) & \text{if } g^{-1}(c_1) < x \leq g^{-1}(c_2), \\ c_2 & \text{if } g^{-1}(c_2) < x \leq b, \end{cases}$$

is the unique minimizer of $\text{ROF}_{\lambda,g}$ in $BV([a, b])$.

Remark 4.2. Since

$$\int_a^{g^{-1}(c)} (c - g(x)) dx = \int_0^c g^{-1}(y) dy, \quad \int_{g^{-1}(c)}^b (g(x) - c) dx = \int_c^1 g^{-1}(y) dy$$

for all $c \in [0, 1]$, the continuity of the integral implies that condition (4.2) is satisfied for every λ sufficiently large.

Proof of Proposition 4.1. We split the proof into two steps.

Step 1. We assume first that u is absolutely continuous. In order to prove the minimality of u , by density it suffices to show that $\text{ROF}_{\lambda,g}(u + \varphi) \geq \text{ROF}_{\lambda,g}(u)$ for

every $\varphi \in C^1([a, b])$, which, in turn, due to the convexity of $\text{ROF}_{\lambda, g}$, is equivalent to proving that

$$(4.3) \quad \left. \frac{d^+}{d\varepsilon} \text{ROF}_{\lambda, g}(u + \varepsilon\varphi) \right|_{\varepsilon=0} \geq 0 \quad \text{for every } \varphi \in C^1([a, b]),$$

where $\frac{d^+}{d\varepsilon}$ denotes the right derivative. By a straightforward computation, we have

$$(4.4) \quad \left. \frac{d^+}{d\varepsilon} \text{ROF}_{\lambda, g}(u + \varepsilon\varphi) \right|_{\varepsilon=0} = \int_{\{u'=0\}} |\varphi'| dx + \int_{\{u'>0\}} \varphi' dx + 2\lambda \int_a^b (u - g)\varphi dx.$$

Consider now the function $\theta : [a, b] \rightarrow [0, 1]$ defined by $\theta(x) := 2\lambda \int_a^x (u - g) dt$. Using (4.2) and the definition of u , one can check that $\theta(a) = \theta(b) = 0$, $0 \leq \theta \leq 1$, and $\theta \equiv 1$ in $[g^{-1}(c_1), g^{-1}(c_2)]$. In particular, $\{u' > 0\} \subset [g^{-1}(c_1), g^{-1}(c_2)] \subset \{\theta = 1\}$ so that by (4.4)

$$\left. \frac{d^+}{d\varepsilon} \text{ROF}_{\lambda, g}(u + \varepsilon\varphi) \right|_{\varepsilon=0} \geq \int_a^b \varphi' \theta dx + 2\lambda \int_a^b (u - g)\varphi dx = 0,$$

where the last equality is obtained by integrating by parts and by using the fact that $\theta' = 2\lambda(u - g)$ and $\theta(a) = \theta(b) = 0$. This shows (4.3) and concludes the proof of Step 1.

Step 2. In the general case, we construct a sequence $\{g_k\} \subset AC([g^{-1}(c_1), g^{-1}(c_2)])$ of nondecreasing functions such that $g_k(g^{-1}(c_1)) = c_1$, $g_k(g^{-1}(c_2)) = c_2$, and $g_k \rightarrow g$ in $L^2([g^{-1}(c_1), g^{-1}(c_2)])$. Let \tilde{g}_k be the function coinciding with g_k in $[g^{-1}(c_1), g^{-1}(c_2)]$ and with g elsewhere in $[a, b]$, and, analogously, let u_k coincide with g_k in the interval $[g^{-1}(c_1), g^{-1}(c_2)]$ and with u elsewhere. For any $v \in BV([a, b])$, by applying the previous step we obtain

$$\text{ROF}_{\lambda, \tilde{g}_k}(v) \geq \text{ROF}_{\lambda, \tilde{g}_k}(u_k) = \text{ROF}_{\lambda, g}(u).$$

The minimality of u follows by letting $k \rightarrow \infty$. Finally, uniqueness is a consequence of the strict convexity of $\text{ROF}_{\lambda, g}$. \square

As a corollary of the previous result, we can prove analytically the occurrence of the staircase effect in a very simple case. Let $g(x) := x$, $x \in [0, 1]$, be the original 1D image to which we add the “noise”

$$h_n(x) := \frac{i}{n} - x \quad \text{if } \frac{i-1}{n} \leq x < \frac{i}{n}, \quad i = 1, \dots, n,$$

where $n \in \mathbb{N}$, so that the resulting degraded 1D image is given by the staircase function

$$(4.5) \quad g_n(x) := \frac{i}{n} \quad \text{if } \frac{i-1}{n} \leq x < \frac{i}{n}, \quad i = 1, \dots, n.$$

Note that, even though $h_n \rightarrow 0$ uniformly, the reconstructed image u_n preserves the staircase structure of g_n . Indeed, we show that there exists a nondegenerate interval $I \subset [0, 1]$ such that each u_n coincides with the degraded 1D image g_n in I for all $n \in \mathbb{N}$. More precisely, we have the following theorem.

THEOREM 4.3 (staircase effect). *Let $\lambda > 4$, let g_n be as in (4.5), and let u_n be the unique minimizer of $\text{ROF}_{\lambda, g_n}$ in $BV([0, 1])$. Then for all n sufficiently large there exist $0 < a_n < b_n < 1$, with*

$$a_n \rightarrow \frac{1}{\sqrt{\lambda}}, \quad b_n \rightarrow 1 - \frac{1}{\sqrt{\lambda}}$$

as $n \rightarrow \infty$, such that $u_n = g_n$ on $[a_n, b_n]$ and u_n is constant on each interval $[0, a_n]$ and $(b_n, 1]$.

Proof. Let g_n^{-1} denote the generalized inverse function of g_n defined by (4.1) with g replaced by g_n . As both $\{g_n\}$ and $\{g_n^{-1}\}$ converge uniformly to $g(x) = x$ and since $\lambda > 4$, one can check that for n large enough there exist $0 < c_1^{(n)} < c_2^{(n)} < 1$ satisfying

$$2\lambda \int_0^{g_n^{-1}(c_1^{(n)})} (c_1^{(n)} - g_n) \, dx = 1 \quad \text{and} \quad 2\lambda \int_{g_n^{-1}(c_2^{(n)})}^1 (g_n - c_2^{(n)}) \, dx = 1$$

with $c_1^{(n)} \rightarrow c_1$ and $c_2^{(n)} \rightarrow c_2$ as $n \rightarrow \infty$, where c_1 and c_2 are defined by

$$(4.6) \quad 2\lambda \int_0^{c_1} (c_1 - x) \, dx = 1 \quad \text{and} \quad 2\lambda \int_{c_2}^1 (x - c_2) \, dx = 1.$$

By Proposition 4.1, the unique minimizer u_n of $\text{ROF}_{\lambda, g_n}$ in $BV(]0, 1[)$ takes the form

$$u_n(x) = \begin{cases} c_1^{(n)} & \text{if } 0 \leq x \leq g_n^{-1}(c_1^{(n)}), \\ g_n(x) & \text{if } g_n^{-1}(c_1^{(n)}) < x \leq g_n^{-1}(c_2^{(n)}), \\ c_2^{(n)} & \text{if } g_n^{-1}(c_2^{(n)}) < x \leq 1. \end{cases}$$

The conclusion follows by observing that $a_n := g_n^{-1}(c_1^{(n)}) \rightarrow c_1$ and $b_n := g_n^{-1}(c_2^{(n)}) \rightarrow c_2$ and that $c_1 = \frac{1}{\sqrt{\lambda}}$ and $c_2 = 1 - \frac{1}{\sqrt{\lambda}}$, thanks to (4.6). \square

4.2. Absence of the staircase effect: The case $p = 1$. Next, we show that the presence of the higher order term in the functional $\overline{\mathcal{F}}_1$ prevents the occurrence of the staircase effect. We begin with the case $p = 1$. We consider the minimization problem

$$(4.7) \quad \min \left\{ \overline{\mathcal{F}}_1(u) + \lambda \int_a^b (u - g)^2 \, dx : u \in X_\psi^1(]a, b[) \right\},$$

where $\overline{\mathcal{F}}_1$ is the relaxed functional given in (2.38). To prove the absence of the staircase effect we need the following auxiliary result that is of independent interest.

PROPOSITION 4.4. *Assume that $\psi: \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying (2.1) and (2.2). Let $g : [a, b] \rightarrow \mathbb{R}$ be Lipschitz continuous, and let $u \in X_\psi^1(]a, b[)$ be a solution of the minimization problem (4.7). Then u is Lipschitz continuous and $u' \in BV(]a, b[)$.*

Proof. The plan of the proof is the following. We will show that the discontinuity set S_u is empty and that the left and right limits $(u')_-^a$ and $(u')_+^a$, respectively, defined in (2.7), are finite everywhere on $]a, b[$ and on $[a, b[$, respectively. Note that this will imply that the sets $Z^\pm[(u')^a]$ (see (2.5) and (2.6)) are empty and, in turn, that $u \in W^{1,1}(]a, b[)$ by the properties of the space $X_\psi^1(]a, b[)$. Moreover, recalling that the functions $(u')_-^a$ and $(u')_+^a$ defined in Remark 2.2 are upper and lower semicontinuous on $[a, b]$, it will also follow that both $(u')_-^a$ and $(u')_+^a$ are bounded, yielding the Lipschitz continuity of u . In turn, the fact that $u' \in BV(]a, b[)$ is a consequence of the local Lipschitz continuity of Ψ_1^{-1} .

Step 1. We start by showing that S_u is empty. We argue, by contradiction, assuming that S_u contains a point x_0 . Without loss of generality, we may suppose that

$\nu_u(x_0) = 1$; i.e., $u_+(x_0) > u_-(x_0)$. We also assume that $\frac{1}{2}(u_+(x_0) + u_-(x_0)) \geq g(x_0)$. In the following, it is convenient to think of u as coinciding everywhere with its lower semicontinuous representative $u_\wedge := \min\{u_-, u_+\}$.

Find $\varepsilon > 0$ so small that

$$(4.8) \quad \sum_{\substack{x \in S_u \\ x \in]x_0, x_0 + \varepsilon[}} |[u](x)| < \frac{[u](x_0)}{4},$$

and let $C > 0$ satisfy

$$(4.9) \quad C > 2\|g'\|_\infty \quad \text{and} \quad \frac{1}{2}(u_+(x_0) + u_-(x_0)) + C\varepsilon > u_-(x_0 + \varepsilon).$$

For $t \in [0, 1]$ consider the affine function

$$h^t(x) := \frac{(1-t)}{2}(u_+(x_0) + u_-(x_0)) + t \left(\frac{1}{4}u_-(x_0) + \frac{3}{4}u_+(x_0) \right) + C(x - x_0)$$

and note that, by (4.9), there exists $x^t \in]x_0, x_0 + \varepsilon[$ such that

$$(4.10) \quad (x^t, h^t(x^t)) \in \Gamma_u \quad \text{and} \quad g < h^t < u \quad \text{in} \quad]x_0, x^t[,$$

where Γ_u stands for the extended graph of u defined by

$$\Gamma_u := \{(x, t) \in]a, b[\times \mathbb{R} : \min\{u_-(x), u_+(x)\} \leq t \leq \max\{u_-(x), u_+(x)\}\}.$$

Let u^t be the function defined by

$$(4.11) \quad u^t(x) := \begin{cases} h^t(x) & \text{if } x \in]x_0, x^t[, \\ u(x) & \text{otherwise,} \end{cases}$$

and note that

$$(4.12) \quad \lambda \left(\int_a^b |u - g|^2 dx - \int_a^b |u^t - g|^2 dx \right) \geq \lambda \left(\int_a^b |u - g|^2 dx - \int_a^b |u^1 - g|^2 dx \right) =: \eta > 0$$

for every $t \in [0, 1]$. Now it is convenient to approximate u with functions having only finitely many jump points. Hence, the following approximation procedure is needed only when S_u is infinite. In this case, write $S_u = \{x_0, x_1, \dots, x_j, \dots\}$, for each k define $S_u^k := \{x_j : 0 \leq j \leq k\}$, and for $x \in]a, b[$ set

$$u_k(x) = u_+(a) + \int_a^x (u')^a dt + (u')^c]a, x[+ \sum_{x_j < x, x_j \in S_u^k} [u](x_j).$$

Note that, since $u_k \rightarrow u$ in $L^\infty]a, b[$, for k large enough it follows from (4.9) and (4.10) that for every $t \in [0, 1]$ there exists $x_k^t \in]x_0, x_0 + \varepsilon[$ such that

$$(x_k^t, h^t(x_k^t)) \in \Gamma_{u_k} \quad \text{and} \quad g < h^t < u_k \quad \text{in} \quad]x_0, x_k^t[,$$

where Γ_{u_k} denotes the extended graph of u_k . For all such k , we consider the comparison function u_k^t defined as in (4.11), with u and x^t replaced by u_k and x_k^t , respectively. Using the uniform convergence of $\{u_k\}$ to u and (4.10), we have that $x^t \leq \liminf_k x_k^t$, which yields $u^t \geq \limsup_k u_k^t$ \mathcal{L}^1 -a.e. on $]a, b[$. Moreover, $u_k \rightarrow u$ in $\overline{\mathcal{F}}_1$ energy. Hence, also by (4.12) we may find k so large that for $t \in [0, 1]$

$$(4.13) \quad \lambda \left(\int_a^b |u_k - g|^2 dx - \int_a^b |u_k^t - g|^2 dx \right) \geq \lambda \left(\int_a^b |u_k - g|^2 dx - \int_a^b |u_k^1 - g|^2 dx \right) \geq \frac{\eta}{2},$$

$$(4.14) \quad \overline{\mathcal{F}}_1(u_k) + \lambda \int_a^b |u_k - g|^2 dx \leq \overline{\mathcal{F}}_1(u) + \lambda \int_a^b |u - g|^2 dx + \frac{\eta}{4}.$$

Let us fix k satisfying (4.13) and (4.14). We claim that there exists $\bar{t} \in [0, 1]$ such that $x_k^{\bar{t}}$ is a continuity point for u_k . Indeed, if not, then for every $t \in [0, 1]$ there exists a jump point x_j , with $1 \leq j \leq k$, such that $x_k^t = x_j$ and the point $(x_k^t, h^t(x_k^t))$ belongs to the corresponding vertical segment of the extended graph of u_k . Setting $I_j := \{t \in [0, 1] : x_k^t = x_j\}$ and $\sigma_j := \{(x_j, h^t(x_j)) : t \in I_j\}$, it is clear that $[0, 1] = \cup_{j=1}^k I_j$ and $\mathcal{H}^1(\sigma_j) = \mathcal{H}^1(\{(x_0, h^t(x_0)) : t \in I_j\})$. Thus,

$$\sum_{\substack{x \in S_u \\ x \in]x_0, x_0 + \varepsilon[}} |u(x)| \geq \sum_{j=1}^k \mathcal{H}^1(\sigma_j) = \mathcal{H}^1(\{(x_0, h^t(x_0)) : t \in [0, 1]\}) = \frac{[u](x_0)}{4},$$

in contradiction with (4.8).

Since from now on \bar{t} and k are fixed, to simplify the notation we set $\hat{x} := x_k^{\bar{t}}$, $\hat{u} := u_k^{\bar{t}}$, $\hat{h} := h^{\bar{t}}$, and $\hat{v} := \Psi_1 \circ (\hat{u}')^a$. By construction (see (4.11)) we have

$$(4.15) \quad |\hat{u}'|(\cdot]a, b[) \leq |u_k'|(\cdot]a, b[).$$

Next, we claim that

$$(4.16) \quad (u')^a_-(\hat{x}) \leq \hat{h}'(\hat{x}) = C.$$

If $(u')^a_-(\hat{x}) \leq 0$, there is nothing to prove. If $(u')^a_-(\hat{x}) > 0$, then by left continuity $(u')^a_-(y) > 0$ for y sufficiently close to \hat{x} , which, in turn, implies $(u')^c(\cdot]y, \hat{x}[) \geq 0$ by the properties of $X^1_\psi(\cdot]a, b[)$. Since S_{u_k} is finite and \hat{x} is a continuity point, for y in a left neighborhood of \hat{x} we can write

$$\hat{h}(\hat{x}) = u_k(\hat{x}) = u_k(y) + \int_y^{\hat{x}} (u')^a(s) ds + (u')^c(\cdot]y, \hat{x}[) > \hat{h}(y) + \int_y^{\hat{x}} (u')^a(s) ds,$$

where we have used the fact that $u_k(\hat{x}) = \hat{h}(\hat{x})$ and $\hat{h} < u_k$ in a left neighborhood of \hat{x} . Claim (4.16) follows.

Now, recalling that $\Phi(1, t_1, t_2) = 2\Psi_1(+\infty) - \Psi_1(t_1) - \Psi_1(t_2)$ for every $t_1, t_2 \in \overline{\mathbb{R}}$

by (2.39) and using Remark 2.6, we estimate

$$\begin{aligned}
 & |v'|([x_0, \hat{x}] \setminus S_u) + \sum_{x \in S_u \cap [x_0, \hat{x}]} \Phi(\nu_u, (u')_-^a, (u')_+^a) \\
 & \geq |v'|([x_0, \hat{x}]) + \Phi(1, (u')_-^a(x_0), (u')_+^a(x_0)) \\
 & \geq |\Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(\hat{x}))| + |\Psi_1((u')_+^a(\hat{x})) - \Psi_1((u')_-^a(\hat{x}))| \\
 & \quad + \Phi(1, (u')_-^a(x_0), (u')_+^a(x_0)) \\
 & = |\Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(\hat{x}))| + |\Psi_1((u')_+^a(\hat{x})) - \Psi_1((u')_-^a(\hat{x}))| \\
 (4.17) \quad & \quad + 2\Psi_1(+\infty) - \Psi_1((u')_-^a(x_0)) - \Psi_1((u')_+^a(x_0)) \\
 & \geq -\Psi_1((u')_-^a(\hat{x})) + 2\Psi_1(+\infty) - \Psi_1((u')_-^a(x_0)) \\
 & \quad + |\Psi_1((u')_+^a(\hat{x})) - \Psi_1((u')_-^a(\hat{x}))| \\
 & = \Psi_1(C) - \Psi_1((u')_-^a(\hat{x})) + 2\Psi_1(+\infty) - \Psi_1((u')_-^a(x_0)) - \Psi_1(C) \\
 & \quad + |\Psi_1((u')_+^a(\hat{x})) - \Psi_1((u')_-^a(\hat{x}))| \\
 & \geq |\Psi_1(C) - \Psi_1((u')_+^a(\hat{x}))| + \Phi(1, (\hat{u}')_-^a(x_0), (\hat{u}')_+^a(x_0)) \\
 & = |\hat{v}'|([x_0, \hat{x}] \setminus S_{\hat{u}}) + \sum_{x \in S_{\hat{u}} \cap [x_0, \hat{x}]} \Phi(\nu_{\hat{u}}, (\hat{u}')_-^a, (\hat{u}')_+^a),
 \end{aligned}$$

where in the last inequality we have used (4.11) and (4.16). Collecting (4.13), (4.15), and (4.17) we deduce that

$$\overline{\mathcal{F}}_1(\hat{u}) + \lambda \int_a^b |\hat{u} - g|^2 + \frac{\eta}{2} \leq \overline{\mathcal{F}}_1(u_k) + \lambda \int_a^b |u_k - g|^2$$

and, in turn, by (4.14)

$$(4.18) \quad \overline{\mathcal{F}}_1(\hat{u}) + \lambda \int_a^b |\hat{u} - g|^2 dx < \overline{\mathcal{F}}_1(u) + \lambda \int_a^b |u - g|^2 dx,$$

which contradicts the minimality of u .

If $\frac{1}{2}(u_+(x_0) + u_-(x_0)) < g(x_0)$, then we proceed in a similar manner: The comparison function \hat{u} is now constructed by replacing u_k with an affine function (defined as before and with C and t properly chosen) in a left neighborhood of x_0 . The argument is completely analogous to the previous one, and we omit the details.

Step 2. We finally show that $(u')_-^a$ and $(u')_+^a$ are finite everywhere in $]a, b[$ and in $]a, b[$, respectively. We give the details only for $(u')_-^a$, since one can argue for $(u')_+^a$ in an entirely similar way.

Recall that, by the previous step, u is continuous. Once again, we reason, by contradiction, by assuming that there exists $\bar{x} \in]a, b[$ such that $|(u')_-^a(\bar{x})| = +\infty$. Without loss of generality, we may suppose that $(u')_-^a(\bar{x}) = +\infty$. Using Remark 2.2 and the differentiability properties of BV functions we may choose a point $x_1 \in]a, \bar{x}[$ such that u is differentiable at x_1 and

$$u(x_1) \neq g(x_1), \quad u'(x_1) = (u')_-^a(x_1) = (u')_+^a(x_1), \quad u'(x_1) > 2\|g'\|_\infty, \quad |v'|([a, x_1]) > 0.$$

The first condition is a consequence of the fact that g is Lipschitz and u cannot be Lipschitz in any left neighborhood of \bar{x} , since $|(u')_-^a(\bar{x})| = +\infty$. The last condition follows easily from the fact that $(u')^a$ cannot be constant \mathcal{L}^1 -a.e. on $]a, \bar{x}[$. Assume

that $u(x_1) > g(x_1)$. Then, by our choice of x_1 and by the previous step, we can find $\varepsilon \in [0, \frac{1}{2}[$, with

$$(4.19) \quad \Psi_1(u'(x_1)) - \Psi_1((1 - \varepsilon)u'(x_1)) < |v'|([x_1, b]),$$

such that the affine function $h(x) := u(x_1) + (1 - \varepsilon)u'(x_1)(x - x_1)$ satisfies one of the following conditions: Either there exists a point $x_2 \in]x_1, b[$ for u such that

$$(4.20) \quad h(x_2) = u(x_2) \text{ and } g < h < u \text{ in }]x_1, x_2[$$

or

$$(4.21) \quad g < h < u \text{ in }]x_1, b[.$$

In the latter case, we set $x_2 := b$. We now consider the comparison function

$$\hat{u}(x) := \begin{cases} h(x) & \text{if } x \in]x_1, x_2[, \\ u(x) & \text{otherwise,} \end{cases}$$

and we denote $\hat{v} := \Psi_1 \circ (\hat{u}')^\alpha$. We claim that (4.18) holds, contradicting the minimality of u . By (4.20) and (4.21) in any case we have

$$\lambda \int_a^b |\hat{u} - g|^2 dx < \lambda \int_a^b |u - g|^2 dx.$$

Moreover, if $x_2 < b$, we have $|\hat{u}'|([x_1, x_2]) = u(x_2) - u(x_1) \leq |u'|([x_1, x_2])$, while if $x_2 = b$, we have $|\hat{u}'|([x_1, b]) = u_-(b) - u(x_1) \leq |u'|([x_1, b])$ so that in both cases $|\hat{u}'|([a, b]) \leq |u'|([a, b])$. Hence, (4.18) will follow if we show that $|\hat{v}'|([x_1, x_2]) \leq |v'|([x_1, x_2])$, where $[x_1, x_2]$ is replaced by $[x_1, b]$ if $x_2 = b$. To see this, we first assume that (4.20) holds. Arguing as for (4.16), we deduce $(u')_-^\alpha(x_2) \leq h'(x_2) = (1 - \varepsilon)u'(x_1)$. Therefore, using the properties of x_1 , we have

$$\begin{aligned} |v'|([x_1, x_2]) &= |v'|([x_1, x_2]) + |v'|(\{x_2\}) \\ &\geq \Psi_1(u'(x_1)) - \Psi_1((u')_-^\alpha(x_2)) + |\Psi_1((u')_-^\alpha(x_2)) - \Psi_1((u')_+^\alpha(x_2))| \\ &= \Psi_1(u'(x_1)) - \Psi_1((1 - \varepsilon)u'(x_1)) + \Psi_1((1 - \varepsilon)u'(x_1)) - \Psi_1((u')_-^\alpha(x_2)) \\ &\quad + |\Psi_1((u')_-^\alpha(x_2)) - \Psi_1((u')_+^\alpha(x_2))| \\ &\geq \Psi_1(u'(x_1)) - \Psi_1((1 - \varepsilon)u'(x_1)) + |\Psi_1((1 - \varepsilon)u'(x_1)) - \Psi_1((u')_+^\alpha(x_2))| \\ &= |\hat{v}'|([x_1, x_2]). \end{aligned}$$

If (4.21) holds, then by (4.19) we obtain

$$|v'|([x_1, b]) > \Psi_1(u'(x_1)) - \Psi_1((1 - \varepsilon)u'(x_1)) = |\hat{v}'|([x_1, b]).$$

If $u(x_1) < g(x_1)$, we modify the previous argument in the following way. We now choose $\varepsilon \in [0, \frac{1}{2}[$ satisfying (4.19) with $|v'|([x_1, b])$ replaced by $|v'|([a, x_1])$ and such that the affine function $h(x)$ defined before satisfies one of the following conditions: Either there exists a point $x_2 \in]a, x_1[$ such that $h(x_2) = u(x_2)$ and $u < h < g$ in $]x_2, x_1[$ or $u < h < g$ in $]a, x_1[$. In the latter case, we set $x_2 := a$. We now consider the comparison function

$$\hat{u}(x) := \begin{cases} h(x) & \text{if } x \in]x_2, x_1[, \\ u(x) & \text{otherwise,} \end{cases}$$

and we proceed exactly as before to show (4.18). \square

We now turn to the main theorem of this subsection.

THEOREM 4.5. Assume that $\psi: \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying (2.1) and (2.2), let $g: [a, b] \rightarrow \mathbb{R}$ be Lipschitz continuous, and let $\{h_n\}$ satisfy

$$(4.22) \quad h_n \rightharpoonup 0 \quad \text{weakly}^* \text{ in } L^\infty(]a, b[).$$

Define \mathcal{A}_n as the class of all solutions to (4.7), with g replaced by $g_n := g + h_n$. Then for n large enough every solution $u_n \in \mathcal{A}_n$ is Lipschitz continuous. Moreover,

$$(4.23) \quad \limsup_{n \rightarrow \infty} \sup_{w \in \mathcal{A}_n} \|w\|_{1, \infty} < +\infty$$

and for every sequence $\{u_n\} \subset \mathcal{A}_n$ there exists a subsequence (not relabeled) and a solution u to (4.7) such that $u_n \rightarrow u$ in $W^{1,p}(]a, b[)$ for all $p \in [1, +\infty[$.

Proof. It will be enough to prove that for any (sub)sequence $\{u_n\} \subset \mathcal{A}_n$ we may extract a further subsequence (not relabeled) and find a solution u to (4.7) such that u_n is Lipschitz continuous for n large enough,

$$(4.24) \quad \limsup_{n \rightarrow \infty} \|u_n\|_{1, \infty} < +\infty,$$

and $u_n \rightarrow u$ in $W^{1,p}(]a, b[)$ for all $p \in [1, +\infty[$. Since the sequence h_n is bounded in $L^\infty(]a, b[)$, for any $w \in X_\psi^1(]a, b[)$ we have

$$\sup_n \left(\overline{\mathcal{F}}_1(u_n) + \lambda \int_a^b (u_n - g_n)^2 dx \right) \leq \overline{\mathcal{F}}_1(w) + \lambda \int_a^b (w - g_n)^2 dx \leq C < \infty$$

for a suitable constant $C > 0$ independent of n . By Corollary 2.7 there exist a subsequence not relabeled and a function $u \in X_\psi^1(]a, b[)$ such that

$$(4.25) \quad u_n \rightharpoonup u \quad \text{weakly}^* \text{ in } BV(]a, b[)$$

and

$$(4.26) \quad u'_n \rightarrow (u')^a \quad \text{pointwise } \mathcal{L}^1\text{-a.e. in }]a, b[.$$

Moreover, since also the functions h_n^2 are equibounded, upon extracting a further subsequence we may find $f \in L^\infty(]a, b[)$ such that

$$(4.27) \quad h_n^2 \rightharpoonup f \quad \text{weakly}^* \text{ in } L^\infty(]a, b[).$$

It is convenient to “localize” the functional $\overline{\mathcal{F}}_1$: For every Borel set $B \subset]a, b[$ and for $w \in X_\psi^1(]a, b[)$ we set

$$(4.28) \quad \overline{\mathcal{F}}_1(w; B) := |w'|_-(B) + |v'|_-(B \setminus S_w) + \sum_{x \in S_w \cap B} \Phi(\nu_w, (w')_-^a, (w')_+^a),$$

where $v := \Psi_1 \circ (w')^a$. We divide the remaining part of the proof into two steps.

Step 1. We claim that u is a solution of the minimization problem (4.7) and that, for every open interval $I =]c, d[$, with $a \leq c < d \leq b$ and $c, d \in [a, b] \setminus S_{(u')^a}$,

$$(4.29) \quad \lim_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; I) = \overline{\mathcal{F}}_1(u; I).$$

To see this, note that for each $n \in \mathbb{N}$

$$\lambda \int_I (u_n - g_n)^2 dx = \lambda \int_I (u_n - g)^2 dx - 2\lambda \int_I (u_n - g) h_n dx + \lambda \int_I h_n^2 dx.$$

By (4.22), (4.25), and (4.27) it follows that

$$(4.30) \quad \lim_{n \rightarrow \infty} \int_I (u_n - g_n)^2 dx = \int_I (u - g)^2 dx + \int_I f dx.$$

Recall, also, that by lower semicontinuity

$$(4.31) \quad \liminf_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; A) \geq \overline{\mathcal{F}}_1(u; A)$$

for every open set $A \subset]a, b[$.

By the minimality of u_n , for every $w \in X_\psi^1(]a, b[)$ we have

$$\begin{aligned} \overline{\mathcal{F}}_1(w) + \lambda \int_a^b (w - g)^2 dx - 2\lambda \int_a^b (w - g) h_n dx + \lambda \int_a^b h_n^2 dx \\ = \overline{\mathcal{F}}_1(w) + \lambda \int_a^b (w - g_n)^2 dx \geq \overline{\mathcal{F}}_1(u_n) + \lambda \int_a^b (u_n - g_n)^2 dx. \end{aligned}$$

Using (4.31) (with $A =]a, b[$) and once again (4.22) and (4.27), we get

$$\begin{aligned} \overline{\mathcal{F}}_1(w) + \lambda \int_a^b (w - g)^2 dx + \lambda \int_a^b f dx &\geq \limsup_{n \rightarrow \infty} \left(\overline{\mathcal{F}}_1(u_n) + \lambda \int_a^b (u_n - g_n)^2 dx \right) \\ &\geq \liminf_{n \rightarrow \infty} \left(\overline{\mathcal{F}}_1(u_n) + \lambda \int_a^b (u_n - g_n)^2 dx \right) \\ &\geq \overline{\mathcal{F}}_1(u) + \lambda \int_a^b (u - g)^2 dx + \lambda \int_a^b f dx. \end{aligned}$$

Given the arbitrariness of $w \in X_\psi^1(]a, b[)$ this implies that u is a solution of the minimization problem (4.7). Moreover, taking $w = u$ in the previous inequalities and using (4.30) we deduce (4.29) for $I =]a, b[$; i.e.,

$$(4.32) \quad \lim_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n) = \overline{\mathcal{F}}_1(u).$$

It remains to prove (4.29) for every open interval of the form $I =]c, d[$, with $c, d \in [a, b] \setminus S_{(u)^c}$. To this end, fix one such interval and assume, by contradiction, that

$$(4.33) \quad \limsup_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; I) > \overline{\mathcal{F}}_1(u; I).$$

As u is continuous by Proposition 4.4, our assumption on I implies that the end points c and d do not charge $\overline{\mathcal{F}}_1(u; \cdot)$ so that $\overline{\mathcal{F}}_1(u; I) = \overline{\mathcal{F}}_1(u; \overline{I} \cap]a, b[)$. Therefore, combining (4.31), (4.32), and (4.33) we obtain

$$\begin{aligned} \overline{\mathcal{F}}_1(u) &= \overline{\mathcal{F}}_1(u; \overline{I} \cap]a, b[) + \overline{\mathcal{F}}_1(u;]a, b[\setminus \overline{I}) = \overline{\mathcal{F}}_1(u; I) + \overline{\mathcal{F}}_1(u;]a, b[\setminus \overline{I}) \\ &< \limsup_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; I) + \liminf_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n;]a, b[\setminus \overline{I}) \leq \lim_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n) = \overline{\mathcal{F}}_1(u), \end{aligned}$$

which is a contradiction. This concludes the proof of (4.29).

Step 2. We now show that u_n is Lipschitz continuous for n large enough and that (4.24) holds. Note that the convergence of u_n to u in $W^{1,p}(]a, b[)$ for all $p \in [1, +\infty[$ will then easily follow from (4.24) and (4.26). Assume, by contradiction, that the conclusion is false. Then, arguing as at the beginning of the proof of Proposition 4.4, we may find a subsequence (not relabeled) and points $x_n \in]a, b[$ such that one of the following two cases holds:

- (i) $x_n \notin S_{(u'_n)^a}$ and $|(u'_n)^a(x_n)| \rightarrow +\infty$;
- (ii) $x_n \in S_{u_n}$ for every $n \in \mathbb{N}$.

Assume that (i) holds and, without loss of generality, that $(u'_n)^a(x_n) \rightarrow +\infty$. Upon extracting a further subsequence, we may also assume that $x_n \rightarrow x_0 \in [a, b]$. Recall that, by Proposition 4.4 and by the previous step, the function u is Lipschitz continuous. Hence, there are two cases: Either

$$(4.34) \quad \overline{\mathcal{F}}_1(u; \{x_0\} \cap]a, b]) = 0$$

or

$$(4.35) \quad x_0 \in S_{(u')^a}, \quad (u')^a_{\pm}(x_0) \in \mathbb{R}, \quad \overline{\mathcal{F}}_1(u; \{x_0\}) = |\Psi_1((u')^a_{+}(x_0)) - \Psi_1((u')^a_{-}(x_0))|.$$

Assume first that (4.34) holds. Set $L := \|u'\|_{\infty}$ and fix ε so small that

$$\overline{\mathcal{F}}_1(u; I_{\varepsilon}) < \int_{L+1}^{+\infty} \psi(t) dt,$$

where $I_{\varepsilon} :=]x_0 - \varepsilon, x_0 + \varepsilon[\cap]a, b[$. By (4.29) we also have

$$(4.36) \quad \overline{\mathcal{F}}_1(u_n; I_{\varepsilon}) < \int_{L+1}^{+\infty} \psi(t) dt$$

for n large enough. On the other hand, by (4.26) there exists $y \in I_{\varepsilon}$ such that $(u'_n)^a(y) < L + 1$ for n large. Moreover, taking into account (i), we also have $(u'_n)^a(x_n) > L + 1$ for n large enough. Thus,

$$\overline{\mathcal{F}}_1(u_n; I_{\varepsilon}) \geq |v'_n|(I_{\varepsilon}) \geq |\Psi_1((u'_n)^a(x_n)) - \Psi_1((u'_n)^a(y))| \geq \Psi_1((u'_n)^a(x_n)) - \Psi_1(L + 1).$$

Passing to the limit as $n \rightarrow \infty$, we then obtain

$$\liminf_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; I_{\varepsilon}) \geq \Psi_1(+\infty) - \Psi_1(L + 1) = \int_{L+1}^{+\infty} \psi(t) dt,$$

which contradicts (4.36).

In case (4.35) holds, then $x_0 \in]a, b[$. Set

$$(4.37) \quad \eta := 2\Psi_1(+\infty) - \Psi_1((u')^a_{+}(x_0)) - \Psi_1((u')^a_{-}(x_0)) - |\Psi_1((u')^a_{+}(x_0)) - \Psi_1((u')^a_{-}(x_0))| > 0,$$

and choose ε such that both $x_0 - \varepsilon$ and $x_0 + \varepsilon$ belong to $]a, b[\setminus S_{(u')^a}$ and

$$(4.38) \quad \overline{\mathcal{F}}_1(u; I_{\varepsilon}) < |\Psi_1((u')^a_{+}(x_0)) - \Psi_1((u')^a_{-}(x_0))| + \frac{\eta}{3},$$

$$(4.39) \quad |\Psi_1((u')^a_{\pm}(y)) - \Psi_1((u')^a_{\pm}(x_0))| < \frac{\eta}{4} \quad \text{for } y \in I_{\varepsilon}^{\pm},$$

where $I_\varepsilon :=]x_0 - \varepsilon, x_0 + \varepsilon[$, $I_\varepsilon^+ :=]x_0, x_0 + \varepsilon[$, and $I_\varepsilon^- :=]x_0 - \varepsilon, x_0[$. Note that by (4.29) and (4.38) we have

$$(4.40) \quad \overline{\mathcal{F}}_1(u_n; I_\varepsilon) < |\Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(x_0))| + \frac{\eta}{3}$$

for n large enough. Moreover, by (4.26) and (4.39) we may find $y^-, y^+ \in I_\varepsilon$, with $y^- < x_0 < y^+$, such that

$$(4.41) \quad y^\pm \notin S_{(u'_n)^a} \quad \text{and} \quad |\Psi_1((u'_n)^a(y^\pm)) - \Psi_1((u')_\pm^a(x_0))| < \frac{\eta}{4}$$

for n large enough. As $y^- < x_n < y^+$ for n sufficiently large, we have

$$(4.42) \quad \begin{aligned} \overline{\mathcal{F}}_1(u_n; I_\varepsilon) &\geq |v'_n|(I_\varepsilon) \geq |\Psi_1((u'_n)^a(x_n)) - \Psi_1((u'_n)^a(y^-))| \\ &\quad + |\Psi_1((u'_n)^a(x_n)) - \Psi_1((u'_n)^a(y^+))| \\ &\geq |\Psi_1((u'_n)^a(x_n)) - \Psi_1((u')_-^a(x_0))| + |\Psi_1((u'_n)^a(x_n)) - \Psi_1((u')_+^a(x_0))| - \frac{\eta}{2}, \end{aligned}$$

where the last inequality follows from (4.41). Letting $n \rightarrow \infty$ in (4.42) and recalling (4.37) we deduce

$$\begin{aligned} \liminf_{n \rightarrow \infty} \overline{\mathcal{F}}_1(u_n; I_\varepsilon) &\geq 2\Psi_1(+\infty) - \Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(x_0)) - \frac{\eta}{2} \\ &= |\Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(x_0))| + \frac{\eta}{2}, \end{aligned}$$

which contradicts (4.40). This concludes the proof of (4.24) if (i) holds.

Assume now that (ii) holds, and let $x_0 \in [a, b]$ be the limit of a subsequence (not relabeled) of $\{x_n\}$. Passing to a further subsequence, we may also assume that $\nu_{u_n}(x_n)$ is constant, say, 1 (the other case is analogous). Again, either (4.34) or (4.35) holds. If (4.34) holds, we may choose, as before, ε so small that (4.36) holds. On the other hand, by (4.26) there exists $y \in I_\varepsilon$ such that $(u'_n)^a(y) < L + 1$ for n large, and, thus, taking into account (ii), we have

$$\begin{aligned} \overline{\mathcal{F}}_1(u_n; I_\varepsilon) &\geq |v'_n|(I_\varepsilon \setminus \{x_n\}) + \Phi(1, (u'_n)_-^a(x_n), (u'_n)_+^a(x_n)) \\ &\geq |\Psi_1((u'_n)^a(y)) - \Psi_1(+\infty)| = \int_{(u'_n)^a(y)}^{+\infty} \psi(t) dt > \int_{L+1}^{+\infty} \psi(t) dt, \end{aligned}$$

which contradicts (4.36). If (4.35) holds, then we may argue exactly as for case (i), with the only difference that (4.42) must be replaced by

$$\begin{aligned} \overline{\mathcal{F}}_1(u_n; I_\varepsilon) &\geq |v'_n|(I_\varepsilon \setminus \{x_n\}) + \Phi(1, (u'_n)_-^a(x_n), (u'_n)_+^a(x_n)) \\ &\geq |\Psi_1((u'_n)_-^a(x_n)) - \Psi_1((u'_n)^a(y^-))| + |\Psi_1((u'_n)_+^a(x_n)) - \Psi_1((u'_n)^a(y^+))| \\ &\quad + \Phi(1, (u'_n)_-^a(x_n), (u'_n)_+^a(x_n)) \\ &\geq 2\Psi_1(+\infty) - \Psi_1((u')_-^a(x_0)) - \Psi_1((u')_+^a(x_0)) - \frac{\eta}{2} \\ &= |\Psi_1((u')_+^a(x_0)) - \Psi_1((u')_-^a(x_0))| + \frac{\eta}{2}, \end{aligned}$$

where we have used (4.37) and (4.41). The last chain of inequalities contradicts (4.40) and concludes the proof of the theorem. \square

4.3. Absence of the staircase effect: The case $p > 1$. We now turn to the case $p > 1$. We consider the minimization problem

$$(4.43) \quad \min \left\{ \overline{\mathcal{F}}_p(u) + \lambda \int_a^b |u - g|^2 dx : u \in X_\psi^p(]a, b[) \right\},$$

where $\overline{\mathcal{F}}_p$ is the relaxed functional given in (3.24). We start with two auxiliary results.

PROPOSITION 4.6. *Let $p > 1$, and assume that $\psi: \mathbb{R} \rightarrow]0, +\infty[$ is a bounded Borel function satisfying (2.2) and (3.1). Let g be Lipschitz continuous, and let u_n be a sequence in $X_\psi^p(]a, b[)$ such that $\sup_n \overline{\mathcal{F}}_p(u_n) < +\infty$ and $u_n \rightarrow g$ in $L^2(]a, b[)$. Then $g \in C^1([a, b]) \cap X_\psi^p(]a, b[)$. Moreover, $u_n \in C^1([a, b])$ for n large enough and $u_n \rightarrow g$ in $C^1([a, b])$.*

Proof. By the assumptions and by Corollary 3.5 we deduce that $g \in X_\psi^p(]a, b[)$. The fact that $g \in C^1([a, b])$ now follows from Remark 3.2(i). To prove the last part of the statement we start by showing that $(u'_n)^a \rightarrow g'$ uniformly in $]a, b[$. Again, by Corollary 3.5, the whole sequence u_n satisfies

$$(4.44) \quad \Psi_p \circ (u'_n)^a \rightharpoonup \Psi_p \circ g' \quad \text{weakly in } W^{1,p}(]a, b[),$$

which implies, in particular, that

$$(4.45) \quad (\Psi_p \circ (u'_n)^a)([a, b]) \subset [\Psi_p(-2\|g'\|_\infty), \Psi_p(2\|g'\|_\infty)] \quad \text{for } n \text{ large enough.}$$

Since by (2.2) Ψ_p^{-1} is Lipschitz continuous on $[\Psi_p(-2\|g'\|_\infty), \Psi_p(2\|g'\|_\infty)]$, it follows from (4.44) and (4.45) that $(u'_n)^a \rightarrow g'$ uniformly in $]a, b[$. In turn, by Definition 3.1 we have that $u'_n = (u'_n)^a$ in $]a, b[$. In particular, $u_n \in C^1([a, b])$ by Remark 3.2(i) and $u_n \rightarrow g$ in $C^1([a, b])$. \square

PROPOSITION 4.7. *Let p and ψ be as in the previous proposition. Then, for every $C > 0$, there exists $\bar{\lambda} = \bar{\lambda}(C)$ with the following property: For all $g \in C^1([a, b]) \cap X_\psi^p(]a, b[)$, with $\|g\|_{C^1([a, b])} \leq C$ and $\overline{\mathcal{F}}_p(g) \leq C$, and for all $\lambda \geq \bar{\lambda}$, every solution u to (4.43) belongs to $C^1([a, b])$.*

Proof. Assume, by contradiction, that for every $n \in \mathbb{N}$ there exist $g_n \in C^1([a, b]) \cap X_\psi^p(]a, b[)$, with $\|g'_n\|_\infty \leq C$ and $\overline{\mathcal{F}}_p(g_n) \leq C$, and a solution u_n to

$$\min \left\{ \overline{\mathcal{F}}_p(u) + n \int_a^b |u - g_n|^2 dx : u \in X_\psi^p(]a, b[) \right\}$$

which does not belong to $C^1([a, b])$. Owing to Proposition 4.6 we may assume, without loss of generality, that $g_n \rightarrow g$ in $C^1([a, b])$ for a suitable function $g \in C^1([a, b]) \cap X_\psi^p(]a, b[)$. Moreover, by minimality, we have

$$\overline{\mathcal{F}}_p(u_n) + n \int_a^b |u_n - g_n|^2 dx \leq \overline{\mathcal{F}}_p(g_n) \leq C.$$

It follows, in particular, that $\sup_n \overline{\mathcal{F}}_p(u_n) < +\infty$ and $u_n \rightarrow g$ in $L^2(]a, b[)$. By Proposition 4.6 we conclude that $u_n \in C^1([a, b])$ for n large enough, which gives a contradiction. \square

The next theorem shows that also in the case $p > 1$ the staircase effect does not occur.

THEOREM 4.8. *Let ψ and p be as in Proposition 4.6, let $g \in C^1([a, b]) \cap X_\psi^p(]a, b[)$, and let h_n satisfy (4.22). For $\lambda > 0$ and $n \in \mathbb{N}$ let $\mathcal{A}_{\lambda,n} \subset X_\psi^p(]a, b[)$ be the class of the solutions to the minimization problem (4.43), with g replaced by $g_n := g + h_n$. Let $\bar{\lambda}$ be as in Proposition 4.7, with $C := \max\{\|g\|_{C^1([a,b])}, \bar{\mathcal{F}}_p(g)\}$. Then for all $\lambda \geq \bar{\lambda}$ we have $\mathcal{A}_{\lambda,n} \subset C^1([a, b])$ for n sufficiently large. Moreover,*

$$(4.46) \quad \lim_{\lambda \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{u \in \mathcal{A}_{\lambda,n}} \|u - g\|_{C^1([a,b])} = 0.$$

Proof. We start by showing the second part of the statement. Assume, by contradiction, that (4.46) does not hold. Then there exist $\delta > 0$, a sequence of real numbers $\lambda_k \rightarrow +\infty$, and, for every k , a sequence of integers $n_j^k \rightarrow \infty$ as $j \rightarrow \infty$ such that for every k, j

$$(4.47) \quad \|u_{\lambda_k, n_j^k} - g\|_{C^1([a,b])} \geq \delta$$

for a suitable function $u_{\lambda_k, n_j^k} \in \mathcal{A}_{\lambda_k, n_j^k}$, with the understanding that

$$\|u_{\lambda_k, n_j^k} - g\|_{C^1([a,b])} = +\infty \text{ if } u_{\lambda_k, n_j^k} \notin C^1([a, b]).$$

Arguing exactly as in Step 1 of the proof of Theorem 4.5 we can show that for every k there exist a subsequence (still denoted by n_j^k) and a solution u_k to (4.43), with λ replaced by λ_k , such that

$$(4.48) \quad u_{\lambda_k, n_j^k} \rightharpoonup u_k \text{ weakly}^* \text{ in } BV(]a, b[) \quad \text{and} \quad \bar{\mathcal{F}}_p(u_{\lambda_k, n_j^k}) \rightarrow \bar{\mathcal{F}}_p(u_k)$$

as $j \rightarrow \infty$. Moreover, since $g \in C^1([a, b]) \cap X_\psi^p(]a, b[)$, we have, by minimality, that

$$(4.49) \quad \bar{\mathcal{F}}_p(u_k) + \lambda_k \int_a^b |u_k - g|^2 dx \leq \bar{\mathcal{F}}_p(g),$$

which shows, in particular, that $u_k \rightarrow g$ in $L^2(]a, b[)$. Combining (4.48) and (4.49) and using a diagonal argument, we may find a subsequence $n_{j_k}^k$ such that

$$\sup_k \bar{\mathcal{F}}_p(u_{\lambda_k, n_{j_k}^k}) < +\infty \quad \text{and} \quad u_{\lambda_k, n_{j_k}^k} \rightarrow g \text{ in } L^2(]a, b[).$$

Proposition 4.6 then implies that $u_{\lambda_k, n_{j_k}^k} \rightarrow g$ in $C^1([a, b])$, which contradicts (4.47).

Finally, the first part of the statement follows from a similar argument, by contradiction, as a consequence of Propositions 4.6 and 4.7 and from the fact that if $u_n \in \mathcal{A}_{\lambda,n}$ then, up to subsequences, u_n converges to a solution of (4.43). \square

Acknowledgments. The authors thank the Center for Nonlinear Analysis (NSF grants DMS-0405343 and DMS-0635983) for its support during the preparation of this paper. The research of G. Dal Maso and M. Morini was partially supported by the projects ‘‘Calculus of Variations’’ 2004 and ‘‘Problemi di Calcolo delle Variazioni in Meccanica e in Scienza dei Materiali’’ 2006–2008, supported by the Italian Ministry of Education, University, and Research, and by the project ‘‘Variational Problems with Multiple Scales’’ 2006, supported by the Italian Ministry of University and Research. The research of I. Fonseca was partially supported by the NSF under grant DMS-040171 and that of G. Leoni under grants DMS-0405423 and DMS-0708039.

REFERENCES

- [1] L. AMBROSIO, N. FUSCO, AND D. PALLARA, *Functions of Bounded Variation and Free Discontinuity Problems*, Oxford Math. Monogr., Clarendon Press, Oxford University Press, New York, 2000.
- [2] P. BLOMGREN, T. F. CHAN, AND P. MULET, *Extensions to total variation denoising*, in Proceedings of the SPIE 1997, San Diego, International Society for Optical Engineering, 1997.
- [3] A. BRAIDES AND R. MARCH, *Approximation by Γ -convergence of a curvature-depending functional in visual reconstruction*, Commun. Pure Appl. Math., 59 (2006), pp. 71–121.
- [4] V. CASELLES, A. CHAMBOLLE, AND M. NOVAGA, *The discontinuity set of solutions of the TV denoising problem and some extensions*, Multiscale Model. Simul., 6 (2007), pp. 879–894.
- [5] A. CHAMBOLLE AND P. L. LIONS, *Image recovery via total variation minimization and related problems*, Numer. Math., 76 (1997), pp. 167–188.
- [6] T. CHAN, A. MARQUINA, AND P. MULET, *High-order total variation-based image restoration*, SIAM J. Sci. Comput., 22 (2000), pp. 503–516.
- [7] G. DAL MASO, *An Introduction to Γ -Convergence*, Birkhäuser, Boston, 1993.
- [8] E. DiBENEDETTO, *Real Analysis*, Birkhäuser, Boston, 2002.
- [9] D. GEMAN AND G. REYNOLDS, *Constrained restoration and the recovery of discontinuities*, IEEE Trans. on Pattern Anal. and Mach. Intell., 14 (1992), pp. 367–383.
- [10] E. HEWITT AND K. STROMBERG, *Real and Abstract Analysis. A Modern Treatment of the Theory of Functions of a Real Variable*, Springer-Verlag, New York, Heidelberg, 1975.
- [11] S. KINDERMANN, S. OSHER, AND P. W. JONES, *Deblurring and denoising of images by nonlocal functionals*, Multiscale Model. Simul., 4 (2005), pp. 1091–1115.
- [12] W. RING, *Structural properties of solutions to total variation regularization problems*, M2AN Math. Model. Numer. Anal., 34 (2000), pp. 799–810.
- [13] L. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Phys. D, 60 (1992), pp. 259–268.
- [14] W. RUDIN, *Real and Complex Analysis*, 2nd ed., McGraw–Hill, New York, Düsseldorf, Johannesburg, 1974.
- [15] W. ZHU, T. CHAN, AND S. ESEDOGLU, *Segmentation with depth: A level set approach*, SIAM J. Sci. Comput., 28 (2006), pp. 1957–1973.

ENTIRE SOLUTIONS IN DELAYED LATTICE DIFFERENTIAL EQUATIONS WITH MONOSTABLE NONLINEARITY*

ZHI-CHENG WANG[†], WAN-TONG LI[‡], AND JIANHONG WU[§]

Abstract. We construct new types of entire solutions for a class of monostable delayed lattice differential equations with global interaction by mixing a heteroclinic orbit of the spatially averaged ordinary differential equations with traveling wave fronts with different speeds. We also establish the uniqueness of entire solutions and the continuous dependence of such an entire solution on parameters, such as wave speeds, for the spatially discrete Fisher-KPP equation.

Key words. entire solution, traveling wave front, heteroclinic orbit, delayed lattice differential equation, monostable nonlinearity

AMS subject classifications. 35B40, 35R10, 37L60, 58D25

DOI. 10.1137/080727312

1. Introduction. We consider the following delayed lattice differential equations:

$$(1.1) \quad u'_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - du_n(t) + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau)),$$

where $D > 0$ is a given constant, $\tau \geq 0$, $I(i) = I(-i) \geq 0$, $J(i) = J(-i) \geq 0$, $\sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) = 1$, $\sum_{i \in \mathbb{Z}} J(i) = 1$, $\sum_{i \in \mathbb{Z} \setminus \{0\}} e^{\lambda|i|} I(i) < \infty$, and $\sum_{i \in \mathbb{Z}} e^{\lambda|i|} J(i) < \infty$ for every $\lambda \geq 0$. The birth function $b \in C^2(\mathbb{R})$, and we assume that there exists a constant $K > 0$ such that

$$b(0) = dK - b(K) = 0$$

and that

(H1) for $u \in (0, K)$, there hold $b(u) > du$, $b'(u) \geq 0$, and $b(u) \leq b'(0)u$;

(H2) $b'(K) < d < b'(0)$.

A specific function $b(u) = pue^{-\alpha u}$ with $p > 0$ and $\alpha > 0$, which has been widely used in the mathematical biology literature, satisfies the above conditions for a wide range of parameters p and α .

A special case when $I(i) = 0$ for $|i| \neq 1$ and $I(1) = \frac{1}{2}$ is

$$(1.2) \quad u'_n = \frac{D}{2} [u_{n+1} + u_{n-1} - 2u_n] - du_n + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau)).$$

*Received by the editors June 14, 2008; accepted for publication (in revised form) October 24, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/72731.html>

[†]School of Mathematics and Statistics, Lanzhou University, Lanzhou, Gansu 730000, People's Republic of China (wangzhch@lzu.edu.cn); Department of Mathematics and Statistics, York University, Toronto, Ontario, M3J 1P3, Canada. This work was partially supported by the NSF of Gansu Province of China (0710RJZA020) and The Fundamental Research Fund for Physics and Mathematics of Lanzhou University (LZULL200807).

[‡]School of Mathematics and Statistics, Lanzhou University, Lanzhou, Gansu, 730000, People's Republic of China (wtli@lzu.edu.cn). This work was partially supported by the NSFC (10871085).

[§]Laboratory for Industrial and Applied Mathematics, Department of Mathematics and Statistics, York University, Toronto, Ontario, M3J 1P3, Canada (wujh@mathstat.yorku.ca). This work was partially supported by Canada Research Chairs Program, by Natural Sciences and Engineering Research Council of Canada, and by Mathematics for Information Technology and Complex Systems.

This system was derived by Weng, Huang, and Wu [40] for the dynamics of growth of a single species population with two age classes distributed over a patchy environment consisting of all integer nodes of a one-dimensional (1-D) lattice. Another special case when $\tau = 0$, $J(0) = 1$, and $J(i) = 0$ for $|i| > 0$ is

$$(1.3) \quad u'_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] + f(u_n(t)),$$

which was derived by Bates and Chmaj [1] as an l_2 -gradient flow for a Helmholtz-free energy functional with general long range interactions. Both lattice systems include, as a special example, the following spatially discrete Fisher-KPP equation:

$$(1.4) \quad u'_n(t) = \frac{D}{2} [u_{n+1} + u_{n-1} - 2u_n] + f(u_n(t)).$$

It is shown in [25] that (1.1) admits a nondecreasing traveling wave front $\phi_c(n+ct)$ satisfying $\phi_c(-\infty) = 0$ and $\phi_c(+\infty) = K$ for every $c \geq c^* > 0$. Furthermore, $\lim_{\xi \rightarrow -\infty} \phi_c(\xi)e^{-\lambda_1(c)\xi} = 1$ and $\lim_{\xi \rightarrow -\infty} \phi'_c(\xi)e^{-\lambda_1(c)\xi} = \lambda_1(c)$ for $c > c^*$, where c and $\lambda_1(c)$ satisfy

$$(1.5) \quad \Delta(\lambda, c) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) (e^{-\lambda i} - 1) - c\lambda - d + b'(0) e^{-\lambda c\tau} \sum_{i \in \mathbb{Z}} J(i) e^{-\lambda i} = 0$$

and c^* is determined by $\Delta(\lambda, c) = 0$ and $\frac{\partial}{\partial \lambda} \Delta(\lambda, c) = 0$. More precisely, there exist $c^* > 0$ and $\lambda^* > 0$ such that

(D1) if $0 < c < c^*$ and $\lambda > 0$, then $\Delta(\lambda, c) > 0$;

(D2) if $c = c^*$, then the equation $\Delta(\lambda, c^*) = 0$ has a double real root $\lambda_1(c^*) = \lambda_2(c^*)$ with $0 < \lambda_1(c^*) = \lambda_2(c^*) = \lambda^*$ such that $\Delta(\lambda, c^*) > 0$ for $\lambda \neq \lambda^*$;

(D3) if $c > c^*$, then the equation $\Delta(\lambda, c) = 0$ has two positive real roots $\lambda_1(c)$ and $\lambda_2(c)$ with $0 < \lambda_1(c) < \lambda^* < \lambda_2(c)$ such that $\lambda'_1(c) < 0$, $\lambda'_2(c) > 0$, $\frac{d}{dc} \{c\lambda_1(c)\} < 0$, and

$$\Delta(\lambda, c) = \begin{cases} > 0 & \text{for } \lambda < \lambda_1(c), \\ < 0 & \text{for } \lambda \in (\lambda_1(c), \lambda_2(c)), \\ > 0 & \text{for } \lambda > \lambda_2(c). \end{cases}$$

We note that in [25] (see also [1], where bistable waves were considered) there is a further assumption on the kernel I ; that is, *the support of I contains either $i = 1$ or two relatively prime integers*, to ensure $\phi'_c(\xi) > 0$ for every $c \geq c^*$. It is interesting to note that, for $c > c^*$, we can confirm $\phi'_c(\xi) > 0$ for any $\xi \in \mathbb{R}$ without this assumption. In fact, for a fixed $c > c^*$, their proof of the existence of nondecreasing traveling wave fronts ϕ_c with $\phi_c(-\infty) = 0$, $\phi_c(+\infty) = K$, and $\lim_{\xi \rightarrow -\infty} \phi_c(\xi)e^{-\lambda_1(c)\xi} = 1$ is independent of this assumption; so is the proof of $\lim_{\xi \rightarrow -\infty} \phi'_c(\xi)e^{-\lambda_1(c)\xi} = \lambda_1(c)$. We now note that $\phi'_c(\xi) \geq 0$ for $\xi \in \mathbb{R}$. Assume that there exists $\xi_0 \in \mathbb{R}$ such that $\phi'_c(\xi_0) = 0$. Then there must be $\phi''_c(\xi_0) = 0$. It is obvious that ϕ_c satisfies

$$c\phi''_c(\xi) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [\phi'_c(\xi - i) - \phi'_c(\xi)] - d\phi'_c(\xi) + \sum_{i \in \mathbb{Z}} J(i)b'(\phi_c(\xi - i - c\tau))\phi'_c(\xi - i - c\tau),$$

which implies that $\sum_{i \in \mathbb{Z} \setminus \{0\}} I(i)\phi'_c(\xi_0 - i) = 0$. By $\sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) = 1$, there exists a $i_0 \in \mathbb{N}$ such that $I(i_0) > 0$ and $\phi'_c(\xi_0 - i_0) = 0$. Let $\xi_1 = \xi_0 - i_0$. A similar argument yields $\phi'_c(\xi_0 - 2i_0) = \phi'_c(\xi_1 - i_0) = 0$. Continuing this procedure, we have

$\phi'_c(\xi_0 - mi_0) = 0$ for all $m \in \mathbb{N}$, a contradiction to the fact $\lim_{\xi \rightarrow -\infty} \phi'_c(\xi)e^{-\lambda_1(c)\xi} = \lambda_1(c)$. Therefore, we have $\phi'_c(\xi) > 0$ for $\xi \in \mathbb{R}$. Of course, to prove $\phi'_c(\xi) > 0$ for the case $c = c^*$ the above assumption on the support of I seems necessary. Since in what follows in this paper we use only the traveling wave fronts ϕ_c with $c > c^*$, we shall not require this assumption.

In the remainder of this paper, we always normalize the traveling wave front $\phi_c(n + ct)$ so that $\phi_c(0) = \frac{K}{2}$. Then, for each $c \in (c^*, +\infty)$, we set

$$(1.6) \quad \alpha_c = \lim_{z \rightarrow -\infty} \phi_c(z)e^{-\lambda_1(c)z}.$$

Furthermore, we define $A_c > 0$ for each $c \in (c^*, +\infty)$ by

$$(1.7) \quad A_c = \inf \left\{ A > 0 : A \geq \phi_c(z)e^{-\lambda_1(c)z} \text{ for any } z \in \mathbb{R} \right\}.$$

It is easy to see that $A_c \geq \alpha_c$.

Our focus is on the so-called entire solutions; here an entire solution of (1.1) is a solution defined for all $(n, t) \in \mathbb{Z} \times \mathbb{R}$. In what follows, we say that a sequence of functions $\Phi_p(t) = \{\Phi_{n,p}(t)\}_{n \in \mathbb{Z}}$ converges to a function $\Phi_{p_0}(t) = \{\Phi_{n,p_0}(t)\}_{n \in \mathbb{Z}}$ in \mathcal{T} if, for every compact set $S \subset \mathbb{Z} \times \mathbb{R}$, the functions $\Phi_{n,p}(t)$ and $\frac{d}{dt}\Phi_{n,p}(t)$ converge uniformly in $(n, t) \in S$ to $\Phi_{n,p_0}(t)$ and $\frac{d}{dt}\Phi_{n,p_0}(t)$ as $p \rightarrow p_0$.

One of our main results can be stated as follows.

THEOREM 1.1. *Let $\Gamma(t)$ be a heteroclinic orbit of the following functional differential equation:*

$$\frac{d}{dt}u(t) = -du(t) + b(u(t - \tau)),$$

which is increasing and satisfies $\Gamma(-\infty) = 0, \Gamma(+\infty) = K, \lim_{t \rightarrow -\infty} e^{-\lambda_* t}\Gamma(t) = K$, and $\Gamma(t) \leq Ke^{\lambda_* t}$ for all $t \in \mathbb{R}$, where $\lambda_* > 0$ is the unique real root of the equation $\lambda + d - b'(0)e^{-\lambda\tau} = 0$. Then for every $c_1, \dots, c_m, c'_1, \dots, c'_l > c^*, \theta_0, \theta_1, \dots, \theta_m, \theta'_1, \dots, \theta'_l \in \mathbb{R}$, and $\chi \in \{0, 1\}$, there exists an entire solution $\Phi(t) = \{\Phi_n(t)\}_{n \in \mathbb{Z}}$ of (1.1) such that

$$(1.8) \quad \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i t + \theta_i), \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j t + \theta'_j), \chi \Gamma(t + \theta_0) \right\} \\ \leq \Phi_n(t) \leq \min \left\{ \vartheta_m^+(n, t), \vartheta_l^-(n, t), \vartheta^0(n, t) \right\}$$

on $(n, t) \in \mathbb{Z} \times \mathbb{R}$, where

$$\vartheta_m^+(n, t) = \min_{1 \leq i \leq m} \left\{ \phi_{c_i}(n + c_i t + \theta_i) + \chi K e^{\lambda_*(t + \theta_0)} \right. \\ \left. + \sum_{1 \leq j \leq m, j \neq i} A_{c_j} e^{\lambda_1(c_j)(n + c_j t + \theta_j)} + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n + c'_j t + \theta'_j)} \right\},$$

$$\vartheta_l^-(n, t) = \min_{1 \leq i \leq l} \left\{ \phi_{c'_i}(-n + c'_i t + \theta'_i) + \chi K e^{\lambda_*(t + \theta_0)} \right. \\ \left. + \sum_{1 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n + c_j t + \theta_j)} + \sum_{1 \leq j \leq l, j \neq i} A_{c'_j} e^{\lambda_1(c'_j)(-n + c'_j t + \theta'_j)} \right\},$$

$$\vartheta^0(n, t) = \chi \Gamma(t + \theta_0) + \sum_{1 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n + c_j t + \theta_j)} + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n + c'_j t + \theta'_j)},$$

and $m, l \in \mathbb{N} \cup \{0\}$ with $\chi + m + l \geq 2$. Moreover, the following statements hold:

- (i) For any $n \in \mathbb{Z}$, $\Phi'_n(t) > 0$ for $t \in \mathbb{R}$.
- (ii) $\lim_{t \rightarrow \infty} \sup_{n \in \mathbb{Z}} |\Phi_n(t) - K| = 0$ and $\lim_{t \rightarrow -\infty} \sup_{|n| \leq N_0} |\Phi_n(t)| = 0$ for every given $N_0 \in \mathbb{N}$.
- (iii) If $m \geq 1$, then $\lim_{n \rightarrow \infty} \|\Phi_n(\cdot) - K\|_{L^\infty[a, +\infty)} = 0$ for every $a \in \mathbb{R}$; if $l \geq 1$, then $\lim_{n \rightarrow -\infty} \|\Phi_n(\cdot) - K\|_{L^\infty[a, +\infty)} = 0$ for every $a \in \mathbb{R}$.
- (iv) If $\chi = 1$ and $m = 0$ ($l = 0$, respectively), then $\Phi_n(t)$ converges uniformly on $t \in [a, b]$ to $\Gamma(t + \theta_0)$ as $n \rightarrow +\infty$ ($n \rightarrow -\infty$, respectively) for any $a, b \in \mathbb{R}$ with $a < b$.
- (v) If $\chi = 1$, then $\Phi_n(t) \sim Ke^{\lambda_*(t+\theta_0)}$ as $t \rightarrow -\infty$ for every $n \in \mathbb{Z}$.
- (vi) If $\chi = 0$, then, for every $n \in \mathbb{Z}$, there exist $B_2(n) > B_1(n) > 0$ such that

$$B_1(n)e^{c_{\max}\lambda_1(c_{\max})t} < \Phi_n(t) < B_2(n)e^{c_{\max}\lambda_1(c_{\max})t} \quad \text{for every } t \ll -1,$$

where $c_{\max} = \max\{\max_{1 \leq i \leq m} c_i, \max_{1 \leq j \leq l} c'_j\}$.

- (vii) If we denote $\Phi(t)$ by $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ when $\chi = 1$ and denote $\Phi(t)$ by $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l}(t)$ when $\chi = 0$, then

$$\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$$

converges to $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l}(t)$ as $\theta_0 \rightarrow -\infty$ in \mathcal{T} and uniformly on $(n, t) \in \mathbb{Z} \times (-\infty, a]$ for every $a \in \mathbb{R}$; $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ converges to K as $\theta_0 \rightarrow +\infty$ in \mathcal{T} and uniformly on $(n, t) \in \mathbb{Z} \times [a, +\infty)$ for every $a \in \mathbb{R}$.

- (viii) $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ converges to

$$\Phi_{c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$$

as $\theta_i \rightarrow -\infty$ in \mathcal{T} and uniformly on $(n, t) \in \{n : n \leq N_0, n \in \mathbb{Z}\} \times (-\infty, a]$ for every $N_0 \in \mathbb{Z}$ and $a \in \mathbb{R}$. $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ converges to

$$\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_{j-1}, c'_{j+1}, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_{j-1}, \theta'_{j+1}, \dots, \theta'_l; \theta_0}(t)$$

as $\theta'_j \rightarrow -\infty$ in \mathcal{T} and uniformly on $(n, t) \in \{n : n \geq N_0, n \in \mathbb{Z}\} \times (-\infty, a]$ for every $N_0 \in \mathbb{Z}$ and $a \in \mathbb{R}$. Similar results hold for $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l}(t)$.

- (ix) $\Phi(t)$ converges to K as $\theta_i \rightarrow +\infty$ in \mathcal{T} and uniformly on $(n, t) \in \{n : n \geq N_0, n \in \mathbb{Z}\} \times [a, +\infty)$ for every $N_0 \in \mathbb{Z}$ and $a \in \mathbb{R}$; $\Phi(t)$ converges to K as $\theta'_j \rightarrow +\infty$ in \mathcal{T} and uniformly on $(n, t) \in \{n : n \leq N_0, n \in \mathbb{Z}\} \times [a, +\infty)$ for every $N_0 \in \mathbb{Z}$ and $a \in \mathbb{R}$.

From (iv) and (v) of Theorem 1.1 and the fact $\lambda_* < c_{\max}\lambda_1(c_{\max})$, it follows that $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ are completely different from

$$\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l}(t).$$

Theorem 1.1 applies to the spatially discrete Fisher-KPP equation (1.4), where $f \in C^2$ satisfies $f(0) = f(1) = 0$, $f'(0) > 0$, $f'(1) < 0$, $f(u) > 0$, and $f(u) \leq f'(0)u$ for $u \in (0, 1)$. In this case, $K = 1$, $d = \max_{u \in [0, 1]} |f'(u)|$, and $b(u) = du + f(u)$. The existence, uniqueness, and stability of traveling wave fronts of (1.4) were studied in Chen, Fu, and Guo [3], Chen and Guo [4, 5], and Zinner [44]; the entire solutions of (1.4) were studied by Guo and Morita [18] and Guo [19]. However, there seem to be no results on the uniqueness of entire solutions of (1.4) and the continuous dependence on parameters $c_1, \dots, c_m, c'_1, \dots, c'_l, \theta_0, \theta_1, \dots, \theta_m, \theta'_1, \dots, \theta'_l$ given by Theorem 1.1. The following theorem is devoted to this topic and is a spatially discrete version of results of Hamel and Nadirashvili [20], where the reaction-diffusion Fisher-KPP equation was considered.

THEOREM 1.2. *For any $c, c' > c^*$, $\theta_0, \theta, \theta' \in \mathbb{R}$, and $\varrho, \varrho', \chi \in \{0, 1\}$ with $\varrho + \varrho' + \chi \geq 2$, there exists a unique entire solution $\Phi(t) = \{\Phi_n(t)\}_{n \in \mathbb{Z}}$ of (1.4) such that (i)–(ix) of Theorem 1.1 hold and*

$$\begin{aligned}
 (1.9) \quad & \max \{ \varrho \phi_c(n + ct + \theta), \varrho' \phi_{c'}(-n + c't + \theta'), \chi \Gamma(t + \theta_0) \} \\
 & \leq \Phi_n(t) \\
 & \leq \min \left\{ \varrho \phi_c(n + ct + \theta) + \chi e^{f'(0)(t+\theta_0)} + \varrho' A_{c'} e^{\lambda_1(c')(-n+c't+\theta')}, \right. \\
 & \quad \varrho' \phi_{c'}(-n + c't + \theta') + \chi e^{f'(0)(t+\theta_0)} + \varrho A_c e^{\lambda_1(c)(n+ct+\theta)}, \\
 & \quad \left. \chi \Gamma(t + \theta_0) + \varrho A_c e^{\lambda_1(c)(n+ct+\theta)} + \varrho' A_{c'} e^{\lambda_1(c')(-n+c't+\theta')} \right\}
 \end{aligned}$$

on $(n, t) \in \mathbb{Z} \times \mathbb{R}$. In particular, when $\varrho = \varrho' = \chi = 1$, the entire solutions $\Phi = \Phi_{c,c',\theta,\theta',\theta_0}$ depend continuously on $(c, c', \theta, \theta', \theta_0) \in (c^*, +\infty)^2 \times \mathbb{R}^3$ in \mathcal{T} ; when $\varrho = \varrho' = 1$ and $\chi = 0$, the entire solutions $\Phi = \Phi_{c,c',\theta,\theta'}$ depend continuously on $(c, c', \theta, \theta') \in (c^*, +\infty)^2 \times \mathbb{R}^2$ in \mathcal{T} ; when $\varrho = \chi = 1$ and $\varrho' = 0$, the entire solutions $\Phi = \Phi_{c,\theta,\theta_0}$ depend continuously on $(c, \theta, \theta_0) \in (c^*, +\infty) \times \mathbb{R}^2$ in \mathcal{T} ; when $\varrho' = \chi = 1$ and $\varrho = 0$, the entire solutions $\Phi = \Phi_{c',\theta',\theta_0}$ depend continuously on $(c', \theta', \theta_0) \in (c^*, +\infty) \times \mathbb{R}^2$ in \mathcal{T} .

We note that when $\varrho' = \chi = 0$ and $\varrho = 1$, $\Phi_n(t) = \Phi_{n;c,\theta}(t) = \phi(n + ct + \theta)$ for $(n, t) \in \mathbb{Z} \times \mathbb{R}$; when $\varrho = \chi = 0$ and $\varrho' = 1$, $\Phi_n(t) = \Phi_{n;c',\theta'}(t) = \phi(-n + c't + \theta')$ for $(n, t) \in \mathbb{Z} \times \mathbb{R}$; and when $\varrho = \varrho' = 0$ and $\chi = 1$, $\Phi_n(t) = \Phi_{n;\theta_0}(t) = \Gamma(t + \theta_0)$ for $(n, t) \in \mathbb{Z} \times \mathbb{R}$. Therefore, similarly to the discussions in Hamel and Nadirashvili [20], it follows from Theorems 1.1 and 1.2 that the functions $\Phi_{c,c',\theta,\theta',\theta_0}(t)$ ($\Phi_{c,c',\theta,\theta'}(t)$, $\Phi_{c,\theta,\theta_0}(t)$, $\Phi_{c',\theta',\theta_0}(t)$, respectively) established by Theorem 1.2 are the 5-D (4-D, 3-D, and 3-D, respectively) manifold of entire solutions of (1.4). In addition, (1.4) possesses two 2-D manifolds of entire solutions of traveling wave type, namely, $\Phi_{c,\theta}^+(t) = \{\phi_c(n + ct + \theta)\}_{n \in \mathbb{Z}}$ and $\Phi_{c',\theta'}^-(t) = \{\phi_{c'}(-n + c't + \theta')\}_{n \in \mathbb{Z}}$, and a 1-D manifold of spatially homogeneous entire solutions, namely, $\Gamma(t + \theta_0)$. Let \mathcal{M}_5 (\mathcal{M}_4 , \mathcal{M}_3^+ , \mathcal{M}_3^- , \mathcal{M}_2^+ , \mathcal{M}_2^- , and \mathcal{M}_1 , respectively) be the above 5-D (4-D, 3-D, 3-D, 2-D, 2-D, and 1-D, respectively) manifold of entire solutions. Then, from Theorems 1.1 and 1.2, it follows that \mathcal{M}_4 is on the boundary of \mathcal{M}_5 (via taking the limit $\theta_0 \rightarrow -\infty$) and \mathcal{M}_3^+ (or \mathcal{M}_3^-) is on the boundary of \mathcal{M}_5 (via taking the limit $\theta \rightarrow -\infty$) (or $\theta' \rightarrow -\infty$). \mathcal{M}_2^+ (or \mathcal{M}_2^-) is on the boundary of \mathcal{M}_4 (via taking the limit $\theta' \rightarrow -\infty$) (or $\theta \rightarrow -\infty$) and is also on the boundary of \mathcal{M}_3^+ (or \mathcal{M}_3^-) (via taking the limit $\theta_0 \rightarrow -\infty$). \mathcal{M}_1 is on the boundary of \mathcal{M}_3^+ (or \mathcal{M}_3^-) (via taking the limit $\theta \rightarrow -\infty$) (or $\theta' \rightarrow -\infty$). In particular, \mathcal{M}_2^+ (or \mathcal{M}_2^-) is on the boundary of \mathcal{M}_5 (via taking the limits $\theta' \rightarrow -\infty$ and $\theta_0 \rightarrow -\infty$) (or $\theta \rightarrow -\infty$ and $\theta_0 \rightarrow -\infty$), and \mathcal{M}_1 is on the boundary of \mathcal{M}_5 (via taking the limits $\theta \rightarrow -\infty$ and $\theta' \rightarrow -\infty$). We can also easily show that the functions $\Phi_{c,c',\theta,\theta',\theta_0}$ converge to $\Phi_{c,\theta}^+$ as $\theta' \rightarrow -\infty$ and $\theta_0 \rightarrow -\infty$ in \mathcal{T} and to $\Phi_{c',\theta'}^-$ as $\theta \rightarrow -\infty$ and $\theta_0 \rightarrow -\infty$ in \mathcal{T} , and that $\Phi_{c,c',\theta,\theta',\theta_0}$ converge to Φ_{θ_0} as $\theta \rightarrow -\infty$ and $\theta' \rightarrow -\infty$ in \mathcal{T} .

Contrasting to [18, 20], we require only $f(u) \leq f'(0)u$ for any $u \in (0, 1)$ other than $f'(u) \leq f'(0)$. We also note some differences on the uniqueness of entire solutions up to a spatial-temporal translation between a reaction-diffusion equation and its spatially discrete analogue (see a similar remark for the bistable nonlinearity reported by Wang, Li, and Ruan [39]). Namely, consider the reaction-diffusion KPP equation

$$(1.10) \quad \frac{d}{dt}u(x, t) = D\Delta u(x, t) + f(u),$$

for which the existence of entire solutions was established by Hamel and Nadirashvili ([20, Theorems 1.1, 1.3, and 1.4 and Corollary 1.5]). For comparison, we consider only the entire solutions established by [20, Theorem 1.3], corresponding to the case $\chi = 0$ and $\varrho = \varrho' = 1$ in our Theorem 1.2. The entire solution $v_{c,c',h,h'}(x, t)$ of (1.10) established by [20, Theorem 1.3] and satisfying (1.4) of [20] is unique for each given $(c, c', h, h') \in (c^*, +\infty)^2 \times \mathbb{R}^2$. Consequently, it is easy to see that for any $(\bar{h}, \bar{h}') \neq (h, h')$,

$$v_{c,c',\bar{h},\bar{h}'}(x, t) = v_{c,c',h,h'}(x + x_0, t + t_0) \quad \text{for } (x, t) \in \mathbb{R}^2,$$

where

$$x_0 = \frac{c(\bar{h}' - h') - c'(\bar{h} - h)}{c + c'}, \quad t_0 = \frac{(\bar{h}' - h') + (\bar{h} - h)}{c + c'}.$$

But for (1.4) if $(\bar{\theta}, \bar{\theta}') \neq (\theta, \theta')$, then $\Phi_{n;c,c',\bar{\theta},\bar{\theta}'}(t) = \Phi_{n+n_0;c,c',\theta,\theta'}(t + t_0)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$ if and only if

$$\frac{c(\bar{h}' - h') - c'(\bar{h} - h)}{c + c'} \in \mathbb{Z},$$

and, hence, $n_0 = \frac{c(\bar{h}' - h') - c'(\bar{h} - h)}{c + c'}$, $t_0 = \frac{(\bar{h}' - h') + (\bar{h} - h)}{c + c'}$. When $(\bar{c}, \bar{c}') \neq (c, c')$, as proved by [24, Theorem 1.1], there exists no $(x_0, t_0) \in \mathbb{R}^2$ such that $v_{\bar{c},\bar{c}',h,h'}(\cdot, \cdot) = v_{c,c',h,h'}(\cdot + x_0, \cdot + t_0)$ on \mathbb{R}^2 for (1.10). Similarly, for (1.4), there exists no $(n_0, t_0) \in \mathbb{Z} \times \mathbb{R}$ such that $\Phi_{n;\bar{c},\bar{c}',h,h'}(t) = v_{n+n_0;c,c',h,h'}(t)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$.

There have been extensive studies about the dynamics of lattice delay systems (1.1), as reported in a recent survey by Gourley and Wu [16]. In particular, the asymptotic speed of propagation and the existence of monotone traveling waves were studied in [25, 40]. The existence, uniqueness, and stability of traveling wave solutions of (1.1) and (1.2) with monostable and bistable nonlinearities have henceforth been studied; see Ma and Zou [26] for the bistable case and Ma and coworkers [25, 27] for the monostable case. Also, Gourley and Wu [17] proved for (1.2) that if the birth rate is so small that a patch alone cannot sustain a positive equilibrium, then the whole population in the patchy environment will become extinct; and if the birth rate is large enough that each patch can sustain a positive equilibrium and if the maturation time is moderate, then the model exhibits nonlinear oscillations characterized by the occurrence of multiple periodic traveling waves. A stage-structured model for a single species on a finite 1-D spatial lattice was also studied in [22]. Related results on traveling waves of lattice differential equations (without delay) can be found in Cahn, Chow, and Van Vleck [2], Chen, Fu, and Guo [3], Chen and Guo [4, 5], Chow [10], Mallet-Paret [28], Wu and Zou [42], and references therein. We note that some progress has been made as well for 2-D lattice delay differential equations; see, for example, Cheng, Li, and Wang [8, 9], Shi, Li, and Cheng, [32], and Weng et al. [41]. In addition, Wang, Li, and Ruan [35, 36, 37] studied traveling wave solutions of reaction-diffusion equations with spatial-temporal delay.

The aforementioned studies also suggest that these wave solutions $\phi_c(n + ct)$ are defined for all $t \in \mathbb{R}$. They often determine the long time behavior of the solutions of Cauchy-type problems and constitute an important part of global attractors, which consist of *entire solutions*. However, the global attractors can be quite complicated, and recent studies for reaction-diffusion equations with continuous spatial variables have showed the existence of many new types of entire solutions arising from the

simple traveling wave fronts, and these entire solutions combined provide essential information about the global attractors; see Chen and Guo [6], Chen, Guo, and Nimomiya [7], Fukao, Morita, and Nimomiya [14], Guo and Morita [18], Hamel and Nadirashvili [20, 21], and Yagisita [43]. For the Fisher-KPP nonlinearity and bistable nonlinearity, these entire solutions behave as two (opposite) wave fronts of positive speed(s) approaching each other from both sides of the x -axis and then annihilate in a finite time. Similar results hold true for nonlocal reaction-diffusion equations with delayed monostable and bistable nonlinearities ([24, 38]). Morita and Ninomiya [30] and Guo [19] have constructed other types of entire solutions for reaction-diffusion equations and discrete diffusive equations with bistable nonlinearity, respectively, which are different from those obtained in [6, 7, 14, 18, 20, 21, 24, 38, 43]. In particular, Li, Liu, and Wang [23] established the existence of entire solutions for reaction-advection-diffusion equations in cylinders, where the ignition temperature nonlinearity has been studied. As reported in [30], entire solutions play also very important roles in some other areas, for example, transient dynamics, distinct history of two solutions, etc.

The remainder of this paper is organized as follows: In section 2, we show how systems (1.1) arise from some areas, such as population biology. In section 3, we establish some existence and comparison results, which are needed in what follows. In section 4, we show the existence of heteroclinic orbit $\Gamma(t)$ connecting two equilibria 0 and K . Section 5 is devoted to Theorem 1.1, and then Theorem 1.2 is proved in section 6.

2. Important particular cases. In this section, we derive from a structured population model a particular case of systems along with an explicit formula to calculate $J(i)$.

Consider a single species population with age structures distributed over a patchy environment consisting of all integer nodes of a 1-D lattice. Let $w_n(t)$ be the density of juvenile individuals in the n th patch and at time t , $v_n(t, a)$ be the density of individuals with age a in the n th patch and at time t , and $\tau > 0$ the length of a juvenile period. Then

$$w_n(t) = \int_0^\tau v_n(t, a) da.$$

Let $u_n(t)$ be the density of mature individuals in the n th patch and at time t . Assume that the spatial dispersal of juvenile individuals and mature individuals is isotropic and can be long range (see Murray [31]). Assume that the diffusion rate of juvenile individuals with age a is $\overline{D}(a) \geq 0$ and the diffusion rate of mature individuals is a constant $D > 0$. Let $K(n-i)$ and $I(n-i)$ be the probability distributions of juvenile individuals and mature individuals traveling from the i th patch to the n th patch, respectively. Then we have

$$K(i) \geq 0, \quad I(i) \geq 0, \quad K(i) = K(-i), \quad I(i) = I(-i), \quad \sum_{i \in \mathbb{Z} \setminus \{0\}} K(i) = 1, \quad \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) = 1.$$

Since only the mature population can reproduce, we have

$$(2.1) \quad \begin{cases} \frac{\partial}{\partial t} v_n(t, a) + \frac{\partial}{\partial a} v_n(t, a) = \overline{D}(a) \sum_{i \in \mathbb{Z} \setminus \{0\}} K(n-i) [v_i(t, a) - v_n(t, a)] \\ \quad \quad \quad - \mu(a) v_n(t, a), \quad 0 < a < \tau, \\ v_n(t, 0) = \widehat{b}(u_n(t)), \\ \frac{d}{dt} u_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(n-i) [u_i(t) - u_n(t)] - \widehat{d}(u_n(t)) + v_n(t, \tau), \end{cases}$$

where $\mu(a)$ denotes the death rate of the juvenile individuals with age $a \in (0, \tau)$, $\widehat{b} : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is the birth function, and $\widehat{d} : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is the death function of mature individuals.

For fixed $s \geq -\tau$, let $V_n^s(t) = v_n(t, t - s)$ for $s \leq t \leq s + \tau$. Then $V_n^s(s) = v_n(s, 0) = \widehat{b}(u_n(s))$. From (2.1),

$$(2.2) \quad \begin{aligned} \frac{d}{dt} V_n^s(t) &= \left. \frac{\partial}{\partial t} v_n(t, a) \right|_{a=t-s} + \left. \frac{\partial}{\partial a} v_n(t, a) \right|_{a=t-s} \\ &= \overline{D}(t-s) \sum_{i \in \mathbb{Z} \setminus \{0\}} K(i) [V_{n-i}^s(t) - V_n^s(t)] - \mu(t-s) V_n^s(t). \end{aligned}$$

Note that the grid function $V_n^s(t)$ can be viewed as the discrete spectral of a periodic function $v^s(t, \omega)$ by discrete Fourier transform [15, 34]:

$$(2.3) \quad v^s(t, \omega) = \frac{1}{\sqrt{2\pi}} \sum_{n=-\infty}^{\infty} e^{-i(n\omega)} V_n^s(t),$$

$$(2.4) \quad V_n^s(t) = \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} e^{i(n\omega)} v^s(t, \omega) d\omega,$$

where i is the imaginary unit. Applying (2.2) and (2.3) yields

$$\begin{aligned} \frac{\partial}{\partial t} v^s(t, \omega) &= \left[\overline{D}(t-s) \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) (e^{-ik\omega} - 1) - \mu(t-s) \right] v^s(t, \omega) \\ &= \left[-2\overline{D}(t-s) \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) \sin^2\left(\frac{k\omega}{2}\right) - \mu(t-s) \right] v^s(t, \omega). \end{aligned}$$

Solving the equation, we get

$$\begin{aligned} v^s(t, \omega) &= \exp \left\{ -2 \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) \sin^2\left(\frac{k\omega}{2}\right) \int_s^t \overline{D}(z-s) dz - \int_s^t \mu(z-s) dz \right\} \\ &\quad \times v^s(s, \omega). \end{aligned}$$

By the inverse discrete Fourier transform (2.4), we obtain

$$\begin{aligned} V_n^s(t) &= \frac{1}{\sqrt{2\pi}} e^{-\int_s^t \mu(z-s) dz} \\ &\quad \times \int_{-\pi}^{\pi} e^{i(n\omega)} \exp \left\{ -2\alpha_s \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) \sin^2\left(\frac{k\omega}{2}\right) \right\} v^s(s, \omega) d\omega, \end{aligned}$$

where $\alpha_{t-s} = \int_s^t \overline{D}(z-s) dz = \int_0^{t-s} \overline{D}(z) dz$. Noting that $V_n^s(s) = v_n(s, 0) = \widehat{b}(u_n(s))$, by (2.3) we have

$$v_s(s, \omega) = \frac{1}{\sqrt{2\pi}} \sum_{j \in \mathbb{Z}} e^{-i(j\omega)} V_j^s(t) = \frac{1}{\sqrt{2\pi}} \sum_{j \in \mathbb{Z}} e^{-i(j\omega)} \widehat{b}(u_j(s)).$$

Hence,

$$(2.5) \quad V_n^s(t) = \frac{1}{2\pi} e^{-\int_s^t \mu(z-s) dz} \sum_{j \in \mathbb{Z}} \widehat{b}(u_j(s)) \times \int_{-\pi}^{\pi} e^{i((n-j)\omega)} \exp \left\{ -2\alpha_s \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) \sin^2 \left(\frac{k\omega}{2} \right) \right\} d\omega.$$

Let $t = s + \tau$, $\widehat{\mu} = e^{-\int_s^t \mu(z-s) dz} = e^{-\int_0^\tau \mu(z) dz}$, and $\alpha = \int_0^\tau \overline{D}(z) dz$. Then (2.5) yields

$$v_n(t, \tau) = \frac{\widehat{\mu}}{2\pi} \sum_{j \in \mathbb{Z}} \beta_\alpha(n-j) \widehat{b}(u_j(s)) = \frac{\widehat{\mu}}{2\pi} \sum_{j \in \mathbb{Z}} \beta_\alpha(n-j) \widehat{b}(u_j(t-\tau)),$$

where

$$\beta_\alpha(j) = \int_{-\pi}^{\pi} e^{i(j\omega)} \exp \left\{ -2\alpha \sum_{k \in \mathbb{Z} \setminus \{0\}} K(k) \sin^2 \left(\frac{k\omega}{2} \right) \right\} d\omega.$$

Thus, the last equality of (2.1) becomes

$$(2.6) \quad \frac{d}{dt} u_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - \widehat{d}(u_n(t)) + \frac{\widehat{\mu}}{2\pi} \sum_{j \in \mathbb{Z}} \beta_\alpha(n-j) \widehat{b}(u_j(t-\tau)), \quad t > 0.$$

Let $\widehat{d}(u) = du$, $b(u) = \widehat{\mu} \widehat{b}(u)$, and $J(i) = \frac{1}{2\pi} \beta_\alpha(i)$; then (2.6) reduces to (1.1).

In particular, the case when $I(i) = K(i) = 0$ for $|i| \neq 1$ and $I(1) = K(1) = \frac{1}{2}$ was studied by Weng, Huang, and Wu [40]. In this case, (2.6) reduces to (1.2). When $\overline{D}(a) \equiv 0$, we have $\alpha = 0$, and it follows that (2.6) reduces to

$$\frac{d}{dt} u_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - \widehat{d}(u_n(t)) + \widehat{\mu} \widehat{b}(u_n(t-\tau)), \quad t > 0.$$

When $\tau = 0$, $\alpha = 0$, and $\widehat{\mu} = 1$, we have

$$\frac{d}{dt} u_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - \widehat{d}(u_n(t)) + \widehat{b}(u_j(t)), \quad t > 0,$$

which coincides with (1.3).

When the diffusion rate $\overline{D}(a)$ and death rate $\mu(a)$ of the juvenile individuals are independent of age a , namely, $D_0 \equiv \overline{D}(a)$ and $\gamma \equiv \mu(a)$ for $a \in [0, \tau]$, we have

$$(2.7) \quad \frac{d}{dt} w_n(t) = D_0 \sum_{i \in \mathbb{Z} \setminus \{0\}} K(i) [w_{n-i}(t) - w_n(t)] - \gamma w_n(t) + \widehat{b}(u_n(t)) - \frac{e^{-\gamma\tau}}{2\pi} \sum_{j \in \mathbb{Z}} \beta_\alpha(n-j) \widehat{b}(u_j(t-\tau)), \quad t > 0.$$

We note that it is easy to prove that $\sum_{j \in \mathbb{Z}} \frac{1}{2\pi} \beta_\alpha(j) = 1$. It seems difficult to prove $\beta_\alpha(j) \geq 0$ for general kernel $\sum_{i \in \mathbb{Z} \setminus \{0\}} K(i) = 1$ though it was proved by Weng, Huang, and Wu [40] for the case when $K(i) = 0$ for $|i| \neq 1$ and $K(\pm 1) = \frac{1}{2}$. Nevertheless, in the remainder of this paper, we consider (1.1) for general kernel functions $I(i)$ and $J(i)$ satisfying the assumptions in section 1.

3. Preliminaries. Consider the initial value problem

$$\begin{cases} u'_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - du_n(t) + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau)), \\ u_n(s) = \varphi_n(s), \end{cases} \tag{3.1}$$

where $n \in \mathbb{Z}$, $t > 0$, and $s \in [-\tau, 0]$.

DEFINITION 3.1. A sequence of continuous differentiable functions $\{v_n(t)\}_{n \in \mathbb{Z}}$, $t \in [-\tau, l]$, $l > 0$, is called a supersolution (subsolution) of (3.1) on $[0, l]$ if

$$\begin{aligned} (3.2) \quad v'_n(t) \geq (\leq) D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [v_{n-i}(t) - v_n(t)] - dv_n(t) \\ + \sum_{i \in \mathbb{Z}} J(i) b(v_{n-i}(t - \tau)) \end{aligned}$$

for all $t \in [0, l]$.

LEMMA 3.2. For any $\varphi = \{\varphi_n\}_{n \in \mathbb{Z}}$ with $\varphi_n \in C([-\tau, 0], [0, K])$, (3.1) admits a unique solution $u(t; \varphi) = \{u_n(t; \varphi)\}_{n \in \mathbb{Z}}$ on $[0, +\infty)$ satisfying $u_n(s) = \varphi_n(s)$ and $0 \leq u_n(t) \leq K$ for $s \in [-\tau, 0]$, $t \in [-\tau, +\infty)$, and $n \in \mathbb{Z}$. For any pair of supersolution $w_n^+(t)$ and subsolution $w_n^-(t)$ of (3.1) on $[0, +\infty)$ with $0 \leq w_n^-(t) \leq K$, $0 \leq w_n^+(t) \leq K$ for $t \in [-\tau, +\infty)$, $n \in \mathbb{Z}$, and $w_n^+(s) \geq w_n^-(s)$ for $s \in [-\tau, 0]$, $n \in \mathbb{Z}$, there holds $w_n^+(t) \geq w_n^-(t)$ for $t \geq 0$, $n \in \mathbb{Z}$.

Note that (3.1) is equivalent to

$$\begin{cases} u_n(t) = \varphi_n(0) e^{-(D+d)t} + \int_0^t e^{(D+d)(s-t)} H_n[u](s) ds, & t > 0, \\ u_n(t) = \varphi_n(t), & t \in [-\tau, 0], \end{cases}$$

where $H_n[u](t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) u_{n-i}(t) + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau))$. Lemma 3.2 can be proved using an argument used in [26, Lemma 4.1].

Consider also the following linear initial value problem:

$$\begin{cases} u'_n(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] - du_n(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) u_{n-i}(t - \tau), \\ u_n(s) = \varphi_n(s) \in C([-\tau, 0], \mathbb{R}), \end{cases} \tag{3.3}$$

where $n \in \mathbb{Z}$, $t > 0$, and $s \in [-\tau, 0]$.

Before stating the following theorem, we first define a Banach space l^∞ by

$$l^\infty = \left\{ \xi = \{\xi_i\}_{i \in \mathbb{Z}}, \xi_i \in \mathbb{R} : \sup_{i \in \mathbb{Z}} |\xi_i| < \infty \right\}$$

with the norm $\|\xi\|_{l^\infty} = \sup_{i \in \mathbb{Z}} |\xi_i|$.

THEOREM 3.3. For any $\varphi = \{\varphi_n\}_{n \in \mathbb{Z}}$ with $\varphi \in C([-\tau, 0], l^\infty)$, (3.3) admits a unique solution $u(t) := u(t; \varphi) = \{u_n(t; \varphi)\}_{n \in \mathbb{Z}}$ on $[0, +\infty)$. Furthermore, if $\varphi^1, \varphi^2 \in C([-\tau, 0], l^\infty)$ satisfy $\varphi_n^1(s) \leq \varphi_n^2(s)$ for any $n \in \mathbb{Z}$ and $s \in [-\tau, 0]$, then $u_n(t; \varphi^1) \leq u_n(t; \varphi^2)$ holds for any $n \in \mathbb{Z}$ and $t > 0$.

Proof. Let $X = l^\infty$. Set $X^+ = \{\xi \in l^\infty : \xi_i \geq 0 \text{ for each } i \in \mathbb{Z}\}$. Then it is easy to see that X^+ is a closed cone of X . Let $T(t) = e^{-(D+d)t}$; it is obvious that $\{T(t)\}$ is a strongly continuous semigroup on X . In particular, it is strongly positive. Now let $\mathcal{C} = C([-\tau, 0], X)$ be the Banach space of continuous functions from $[-\tau, 0]$ into X with the supremum norm. Set $\mathcal{C}^+ = \{\Phi \in \mathcal{C} : \Phi(s) \in X^+, s \in [-\tau, 0]\}$. Then \mathcal{C}^+ is a positive cone of \mathcal{C} . For any continuous function $w : [-\tau, +\infty) \rightarrow X$, define $w_t \in \mathcal{C}$,

$t \in [0, +\infty)$, by $w_t(s) = w(t + s)$, $s \in [-\tau, 0]$. Then the map $t \mapsto w_t$ is a continuous function from $[0, +\infty)$ to \mathcal{C} .

Define $f : \mathcal{C} \rightarrow X$ by

$$f(w) = \{f_i(w)\}_{i \in \mathbb{Z}}$$

for $w = \{w_i\}_{i \in \mathbb{Z}} \in \mathcal{C}$, where

$$f_i(w) = D \sum_{j \in \mathbb{Z} \setminus \{0\}} I(j) w_{i-j}(0) + b'(0) \sum_{k \in \mathbb{Z}} J(k) w_{i-k}(-\tau).$$

It is not difficult to verify that $f : \mathcal{C} \rightarrow X$ is globally Lipschitz continuous. Furthermore, since for any $v, w \in \mathcal{C}$ with $v \geq w$ in \mathcal{C} ,

$$\begin{aligned} f_i(v) - f_i(w) &= D \sum_{j \in \mathbb{Z} \setminus \{0\}} I(j) v_{i-j}(0) - D \sum_{j \in \mathbb{Z} \setminus \{0\}} I(j) w_{i-j}(0) \\ &\quad + b'(0) \sum_{k \in \mathbb{Z}} J(k) v_{i-k}(-\tau) - b'(0) \sum_{k \in \mathbb{Z}} J(k) w_{i-k}(-\tau) \\ &\geq 0, \end{aligned}$$

it follows that $f(v) \geq f(w)$ in X for any $v, w \in \mathcal{C}$ with $v \geq w$, which implies that $f : \mathcal{C} \rightarrow X$ is quasi-monotone in the sense that

$$\lim_{h \rightarrow 0} \frac{1}{h} \text{dist}((v(0) - w(0)) + h[f(v) - f(w)], X^+) = 0$$

for any $v, w \in \mathcal{C}$ with $v \geq w$.

Note that (3.3) is equivalent to

$$(3.4) \quad \begin{cases} u(t) = T(t)u(0) + \int_0^t T(t-s)f(u_s)ds, & t > 0, \\ u(t) = \varphi(t), & t \in [-\tau, 0]. \end{cases}$$

Take $M_0 = \max_{t \in [-\tau, 0]} \|\varphi(t)\|_{l^\infty}$. Furthermore, define a vector-valued function $v^+(\cdot) = \{v_n^+(\cdot)\}_{n \in \mathbb{Z}} : [-\tau, +\infty) \rightarrow X$ by

$$(3.5) \quad \begin{cases} v_n^+(t) = M_0, & t \in [-\tau, 0], \\ v_n^+(t) = M_0 e^{(b'(0)-d)t}, & t > 0 \text{ for any } n \in \mathbb{Z}. \end{cases}$$

It is easy to verify that v^+ satisfies

$$(3.6) \quad v^+(t) \geq T(t)v^+(s) + \int_s^t T(t-r)f(v_r^+)dr \text{ for any } t > s \geq 0.$$

Define $v^-(\cdot) = \{v_n^-(\cdot)\}_{i \in \mathbb{Z}} : [-\tau, +\infty) \rightarrow X$ by $v^-(\cdot) = -v^+(\cdot)$. Then v^- satisfies

$$(3.7) \quad v^-(t) \leq T(t)v^-(s) + \int_s^t T(t-r)f(v_r^-)dr \text{ for any } t > s \geq 0.$$

Now we use the conclusions of [29]. By setting $S(t, s) = T(t, s) = T(t - s)$ for any $t \geq s \geq 0$ and $B(t, \Phi) = f(\Phi)$, the existence and uniqueness of the solution $u(t; \varphi)$ follows from [29, Corollary 5].

For any $\varphi^1, \varphi^2 \in \mathcal{C}$ with $\varphi^1 \leq \varphi^2$ in \mathcal{C} , again applying [29, Corollary 5], we have

$$v^-(t) \leq u(t; \varphi^1) \leq u(t; \varphi^2) \leq v^+(t) \text{ in } X \text{ for any } t \geq 0,$$

via letting

$$M_0 = \max \left\{ \max_{s \in [-\tau, 0]} \|\varphi^1(s)\|_{l^\infty}, \max_{s \in [-\tau, 0]} \|\varphi^2(s)\|_{l^\infty} \right\}$$

in (3.5). This implies the solution semiflow is order preserving. The proof is complete. \square

Remark 3.1. Assume that the continuous functions $w^\pm = \{w_n^\pm\}_{n \in \mathbb{Z}} : [-\tau, +\infty) \rightarrow l^\infty$ satisfy (3.6) and (3.7), respectively, and $w_n^+(s) \geq w_n^-(s)$ for any $(n, s) \in \mathbb{Z} \times [-\tau, 0]$; then we have $w_n^+(t) \geq w_n^-(t)$ for any $(n, t) \in \mathbb{Z} \times [0, +\infty)$.

THEOREM 3.4. *Assume that*

$$w_n^-(t) \in C([-\tau, \infty), (-\infty, K]) \quad \text{and} \quad w_n^+(t) \in C([-\tau, \infty), [0, \infty))$$

satisfy $w_n^-(t) \leq w_n^+(t)$ *for any* $t \in [-\tau, 0]$ *and* $n \in \mathbb{Z}$ *and*

$$(3.8) \quad \begin{aligned} \frac{d}{dt} w_n^+(t) &\geq D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n-i}^+(t) - w_n^+(t)] - d w_n^+(t) \\ &\quad + b'(0) \sum_{i \in \mathbb{Z}} J(i) w_{n-i}^+(t - \tau), \end{aligned}$$

$$(3.9) \quad \begin{aligned} \frac{d}{dt} w_n^-(t) &\leq D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n-i}^-(t) - w_n^-(t)] - d w_n^-(t) \\ &\quad + b'(0) \sum_{i \in \mathbb{Z}} J(i) w_{n-i}^-(t - \tau) \end{aligned}$$

for any $t > 0$ *and* $n \in \mathbb{Z}$. *Then there holds* $w_n^+(t) \geq w_n^-(t)$ *for any* $t > 0$ *and* $n \in \mathbb{Z}$.

Proof. Put $w_n(t) := w_n^-(t) - w_n^+(t)$, $n \in \mathbb{Z}$, $t \in [-\tau, +\infty)$. Then $w_n(t)$ is continuous and bounded from above by K , and $\bar{w}(t) := \sup_{n \in \mathbb{Z}} w_n(t)$ is continuous on $[-\tau, \infty)$. We use a contradiction argument to prove the assertion. Suppose that the assertion is not true. Let $M_0 > 0$ be such that $M_0 + d - b'(0)e^{-M_0\tau} > 0$. Then there exists $t_0 > 0$ such that $\bar{w}(t_0) > 0$ and

$$(3.10) \quad \bar{w}(t_0) e^{-M_0 t_0} = \sup_{t \geq -\tau} \{\bar{w}(t) e^{-M_0 t}\} > \bar{w}(s) e^{-M_0 s} \text{ for all } s \in [-\tau, t_0).$$

Let $\{n_j\}_{j \in \mathbb{N}}$ be a sequence so that $w_{n_j}(t_0) > 0$ for all $j \geq 1$ and $\lim_{j \rightarrow \infty} w_{n_j}(t_0) = \bar{w}(t_0)$. Let $\{t_j\}_{j \in \mathbb{N}} \subset (0, t_0)$ so that

$$(3.11) \quad w_{n_j}(t_j) e^{-M_0 t_j} = \max_{t \in [0, t_0]} \{w_{n_j}(t) e^{-M_0 t}\}.$$

Since

$$w_{n_j}(t_0) e^{-M_0 t_0} \leq w_{n_j}(t_j) e^{-M_0 t_j} \leq \bar{w}(t_j) e^{-M_0 t_j} \leq \bar{w}(t_0) e^{-M_0 t_0},$$

we have $\lim_{j \rightarrow +\infty} \bar{w}(t_j) e^{-M_0 t_j} = \bar{w}(t_0) e^{-M_0 t_0}$. Then there must be $\lim_{j \rightarrow +\infty} t_j = t_0$ due to (3.10). In view of $w_{n_j}(t_0) e^{-M_0(t_0 - t_j)} \leq w_{n_j}(t_j) e^{-M_0 t_j} \leq \bar{w}(t_0) e^{-M_0(t_0 - t_j)}$, we obtain $\lim_{j \rightarrow +\infty} w_{n_j}(t_j) = \bar{w}(t_0)$.

Following (3.11), for each $j \geq 1$, we have

$$0 \leq \frac{d}{dt} \{w_{n_j}(t) e^{-M_0 t}\}_{t=t_j} = [w'_{n_j}(t_j) - M_0 w_{n_j}(t_j)] e^{-M_0 t_j},$$

and, hence, $w'_{n_j}(t_j) \geq M_0 w_{n_j}(t_j)$. Then it follows from (3.8) and (3.9) that

$$\begin{aligned} 0 &\geq w'_{n_j}(t_j) - D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n_j-i}(t_j) - w_{n_j}(t_j)] + dw_{n_j}(t_j) \\ &\quad - b'(0) \sum_{i \in \mathbb{Z}} J(i) [w_{n_j-i}^-(t_j - \tau) - w_{n_j-i}^+(t_j - \tau)] \\ &\geq (M_0 + D + d) w_{n_j}(t_j) - D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) w_{n_j-i}(t_j) - b'(0) \bar{w}(t_j - \tau) \\ &\geq (M_0 + D + d) w_{n_j}(t_j) - D \bar{w}(t_j) - b'(0) \bar{w}(t_j - \tau). \end{aligned}$$

Taking $j \rightarrow +\infty$, we have

$$\begin{aligned} 0 &\geq (M_0 + D + d) \bar{w}(t_0) - D \bar{w}(t_0) - b'(0) e^{M_0(t_0-\tau)} [\bar{w}(t_0 - \tau) e^{-M_0(t_0-\tau)}] \\ &\geq (M_0 + d) \bar{w}(t_0) - b'(0) e^{M_0(t_0-\tau)} \bar{w}(t_0) e^{-M_0 t_0} \\ &= [M_0 + d - b'(0) e^{-M_0 \tau}] \bar{w}(t_0). \end{aligned}$$

In view of $M_0 + d - b'(0) e^{-M_0 \tau} > 0$, we obtain that $\bar{w}(t_0) \leq 0$, which contradicts to $\bar{w}(t_0) > 0$. Consequently, we conclude that $w_n^+(t) \geq w_n^-(t)$ for all $n \in \mathbb{Z}$ and $t \in (0, +\infty)$. This completes the proof. \square

4. Existence of heteroclinic orbits. In this section, we show the existence of a heteroclinic orbit connecting the equilibria $u \equiv 0$ and $u \equiv K$ for the following functional differential equation:

$$(4.1) \quad \frac{d}{dt} u(t) = -du(t) + b(u(t - \tau)).$$

There are now various methods developed to establish the existence of such a heteroclinic orbit, for example, Faria, Huang, and Wu [11], Faria and Trofimchuk [12, 13], Li, Wang, and Wu [24], and Smith [33]. However, except for Faria and Trofimchuk [13], these results do not give the exponential decay rate of the heteroclinic orbit connecting the equilibria $u \equiv 0$ and $u \equiv K$ at minus infinity. At the same time, the results in [13] are not directly applicable (see the condition (A1) of [13]) and do not ensure the monotonicity of the heteroclinic orbit.

Define

$$\Lambda(\lambda) = \lambda + d - b'(0) e^{-\lambda \tau};$$

then it is easy to prove that the equation $\Lambda(\lambda) = 0$ has one and only one real root $\lambda_* > 0$ such that $\Lambda(\lambda) < 0$ for any $\lambda < \lambda_*$ and $\Lambda(\lambda) > 0$ for any $\lambda > \lambda_*$.

Define an operator $S : C(\mathbb{R}, [0, K]) \rightarrow C(\mathbb{R}, [0, K])$ by

$$S(u)(t) = e^{-dt} \int_{-\infty}^t e^{ds} b(u(s - \tau)) ds \text{ for any } u \in C(\mathbb{R}, [0, K]).$$

PROPOSITION 4.1.

- (i) If $u \in C(\mathbb{R}, [0, K])$, then $S(u) \in C^1(\mathbb{R}, [0, K])$.
 - (ii) For any $u, v \in C(\mathbb{R}, [0, K])$ with $u \leq v$, $S(u) \leq S(v)$.
 - (iii) For any $u \in C(\mathbb{R}, [0, K])$, if $u(\cdot)$ is increasing in \mathbb{R} , then so is $S(u)(\cdot)$.
- Let $b''_{\max} = \max_{u \in [0, K]} |b''(u)|$. Define

$$\bar{u}(t) = K \min \{e^{\lambda_* t}, 1\} \text{ and } \underline{u}(t) = \max \{Ke^{\lambda_* t} (1 - Me^{\epsilon t}), 0\},$$

where $\epsilon \in (0, \lambda_*)$ and $M > 1$ with

$$1 - \frac{(d + \lambda_*) e^{-\epsilon \tau}}{d + \lambda_* + \epsilon} - \frac{(d + \lambda_*) Ke^{-\lambda_* \tau} b''_{\max}}{Mb'(0)(d + 2\lambda_*)} > 0.$$

LEMMA 4.2. For any $t \in \mathbb{R}$, $S(\bar{u})(t) \leq \bar{u}(t)$ and $\underline{u}(t) \leq S(\underline{u})(t)$.

Proof. First, we prove $S(\bar{u})(t) \leq \bar{u}(t)$. When $t \geq 0$, $\bar{u}(t) = K$. Therefore,

$$\begin{aligned} S(\bar{u})(t) &= e^{-dt} \int_{-\infty}^t e^{ds} b(\bar{u}(s - \tau)) ds \\ &\leq e^{-dt} \int_{-\infty}^t e^{ds} b(K) ds = dKe^{-dt} \int_{-\infty}^t e^{ds} ds = K = \bar{u}(t). \end{aligned}$$

When $t < 0$, $\bar{u}(t) = Ke^{\lambda_* t}$. Noting that $d + \lambda_* = b'(0)e^{-\lambda_* \tau}$, we have

$$\begin{aligned} S(\bar{u})(t) &= e^{-dt} \int_{-\infty}^t e^{ds} b(\bar{u}(s - \tau)) ds \leq e^{-dt} \int_{-\infty}^t e^{ds} b'(0)\bar{u}(s - \tau) ds \\ &\leq b'(0)Ke^{-dt} \int_{-\infty}^t e^{ds} e^{\lambda_*(s-\tau)} ds = \frac{b'(0)e^{-\lambda_* \tau} K}{d + \lambda_*} e^{\lambda_* t} = \bar{u}(t). \end{aligned}$$

Now we prove $\underline{u}(t) \leq S(\underline{u})(t)$. Let $t_0 = \frac{1}{\epsilon} \ln \frac{1}{M} < 0$ such that $1 - Me^{\epsilon t_0} = 0$. When $t \geq t_0$, $\underline{u}(t) = 0$, and, hence, $\underline{u}(t) \leq S(\underline{u})(t)$. When $t < t_0$, $\underline{u}(t) = Ke^{\lambda_* t}(1 - Me^{\epsilon t}) \leq Ke^{\lambda_* t}$. In this case, we have

$$\begin{aligned} S(\underline{u})(t) &= e^{-dt} \int_{-\infty}^t e^{ds} b(\underline{u}(s - \tau)) ds \\ &\geq e^{-dt} \int_{-\infty}^t e^{ds} [b'(0)\underline{u}(s - \tau) - b''_{\max}\underline{u}^2(s - \tau)] ds \\ &\geq e^{-dt} \int_{-\infty}^t e^{ds} \left[b'(0)Ke^{\lambda_*(s-\tau)} (1 - Me^{\epsilon(s-\tau)}) - b''_{\max}K^2e^{2\lambda_*(s-\tau)} \right] ds \\ &= \frac{b'(0)Ke^{-\lambda_* \tau}}{d + \lambda_*} e^{\lambda_* t} - \frac{Mb'(0)Ke^{-(\lambda_* + \epsilon)\tau}}{d + \lambda_* + \epsilon} e^{(\lambda_* + \epsilon)t} - \frac{b''_{\max}K^2e^{-2\lambda_* \tau}}{d + 2\lambda_*} e^{2\lambda_* t} \\ &\geq \underline{u}(t) + \frac{Mb'(0)Ke^{-\lambda_* \tau}}{d + \lambda_*} \left[1 - \frac{(d + \lambda_*)e^{-\epsilon \tau}}{d + \lambda_* + \epsilon} - \frac{(d + \lambda_*)Ke^{-\lambda_* \tau} b''_{\max}}{Mb'(0)(d + 2\lambda_*)} \right] e^{(\lambda_* + \epsilon)t} \\ &\geq \underline{u}(t). \end{aligned}$$

The proof is complete. \square

THEOREM 4.3. *There exists a heteroclinic solution $\Gamma(t)$ of (4.1), which is increasing on \mathbb{R} and satisfies $\lim_{t \rightarrow -\infty} e^{-\lambda_* t} \Gamma(t) = K$, $\Gamma(+\infty) = K$, $\Gamma(t) \leq K e^{\lambda_* t}$, and $\Gamma'(t) > 0$ for every $t \in \mathbb{R}$.*

Proof. By an argument similar to that of [42, Theorem 3.1], we can get a nondecreasing solution $\Gamma(t)$ which meets the theorem except $\Gamma'(t) > 0$ for any $t \in \mathbb{R}$. Since $\Gamma'(t)$ satisfies

$$\Gamma''(t) = -d\Gamma'(t) + b'(\Gamma(t - \tau))\Gamma'(t - \tau) \quad \forall t \in \mathbb{R},$$

we have

$$\Gamma'(t) = e^{-d(t-s)}\Gamma'(s) + \int_s^t e^{-d(t-r)}b'(\Gamma(t - \tau))\Gamma'(r - \tau)dr \quad \text{for any } s < t.$$

Note that $\Gamma'(t) \geq 0$ for any $t \in \mathbb{R}$. Then it is easy to see that if $\Gamma'(t_0) > 0$ for some $t_0 \in \mathbb{R}$, then $\Gamma'(t) > 0$ for all $t > t_0$. In view of $\lim_{t \rightarrow -\infty} e^{-\lambda_* t} \Gamma(t) = K$, we know that there exists a sequence $\{t_i\}$ with $t_i \rightarrow -\infty$ as $i \rightarrow +\infty$ such that $\Gamma(t_i) > 0$ for any $i \in \mathbb{N}$. Hence, we conclude $\Gamma'(t) > 0$ for any $t \in \mathbb{R}$. This completes the proof. \square

5. Proof of Theorem 1.1. In this section, we prove Theorem 1.1.

LEMMA 5.1. *Suppose that $u(t; \varphi) = \{u_n(t; \varphi)\}_{n \in \mathbb{Z}}$ is a solution of (1.1) with initial value $\varphi = \{\varphi_n\}_{n \in \mathbb{Z}}$ with $\varphi_n \in C([-\tau, 0], [0, K])$; then there exists a positive constant $M_* > 0$ such that for any $\varphi = \{\varphi_n\}_{n \in \mathbb{Z}}$ with $\varphi_n \in C([-\tau, 0], [0, K])$ and $t > \tau$, $|u'_n(t; \varphi)| \leq M_*$ and $|u''_n(t; \varphi)| \leq M_*$.*

Proof. Denote $u_n(t; \varphi)$ by $u_n(t)$. Let $M' = 2DK + 2dK$. It is easy to see that $|u'_n(t; \varphi)| \leq M'$ for any $t > 0$. For $t > \tau$, there is

$$\begin{aligned} u''_n(t) &= D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u'_{n-i}(t) - u'_n(t)] \\ &\quad - du'_n(t) + \sum_{i \in \mathbb{Z}} J(i) b'(u_{n-i}(t - \tau)) u'_{n-i}(t - \tau). \end{aligned}$$

Set $M'' = 2DM' + dM' + M'b'(0)$. Then $|u''_n(t; \varphi)| \leq M''$. Note that M' and M'' are independent of φ and $t > \tau$. Take $M_* = \max\{M', M''\}$. This completes the proof. \square

LEMMA 5.2. *Let $u^k(t; \varphi^k) = \{u_n^k(t; \varphi^k)\}_{n \in \mathbb{Z}}$ be a solution of the following initial value problem:*

$$\begin{cases} \frac{d}{dt} u_n^k(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [u_{n-i}(t) - u_n(t)] \\ \quad - du_n(t) + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau)), \quad t > 0, \\ u_n^k(t) = \varphi_n^k(t), \quad t \in [-\tau, 0], \end{cases}$$

where

$$\begin{aligned} \varphi_n^k(t) &= \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(t - k) + \theta_i), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(t - k) + \theta'_j), \chi \Gamma((t - k) + \theta_0) \right\}. \end{aligned}$$

Then $v^k(t; \varphi^k) = \{u_n^k(t+k; \varphi^k)\}_{n \in \mathbb{Z}}$ satisfies

$$(5.1) \quad \limsup_{t > -k, k \rightarrow +\infty} v_n^k(t) \leq \phi_{c_i}(n + c_i t + \theta_i) + \chi K e^{\lambda_*(t+\theta_0)} + \sum_{1 \leq j \leq m, j \neq i} A_{c_j} e^{\lambda_1(c_j)(n+c_j t+\theta_j)} \\ + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j t+\theta'_j)} \quad \text{for } 1 \leq i \leq m,$$

$$(5.2) \quad \limsup_{t > -k, k \rightarrow +\infty} v_n^k(t) \leq \phi_{c'_i}(-n + c'_i t + \theta'_i) + \chi K e^{\lambda_*(t+\theta_0)} + \sum_{1 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j t+\theta_j)} \\ + \sum_{1 \leq j \leq l, j \neq i} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j t+\theta'_j)} \quad \text{for } 1 \leq i \leq l,$$

$$(5.3) \quad \limsup_{t > -k, k \rightarrow +\infty} v_n^k(t) \leq \chi \Gamma(t + \theta_0) + \sum_{1 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j t+\theta_j)} \\ + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j t+\theta'_j)}.$$

Proof. We prove only (5.1), because the proofs of (5.2) and (5.3) are similar to that of (5.1). Assume $m \geq 1$. Consider $i = 1$. Let

$$w_n^k(t) = u_n^k(t) - \phi_{c_1}(n + c_1(t - k) + \theta_1).$$

Then $w_n^k(t)$ satisfies

$$\frac{d}{dt} w_n^k(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n-i}(t) - w_n(t)] - d w_n(t) + \sum_{i \in \mathbb{Z}} J(i) b(u_{n-i}(t - \tau)) \\ - \sum_{i \in \mathbb{Z}} J(i) b(\phi_{c_1}(n - i + c_1(t - \tau) + \theta_1)) \\ \leq D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n-i}(t) - w_n(t)] - d w_n(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) w_{n-i}(t - \tau).$$

Since

$$\chi K e^{\lambda_*(t-k+\theta_0)} + \sum_{2 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j(t-k)+\theta_j)} + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j(t-k)+\theta'_j)}$$

is a solution of (3.3) with

$$\varphi_n(s) = \chi K e^{\lambda_*(s-k+\theta_0)} + \sum_{2 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j(s-k)+\theta_j)} \\ + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j(s-k)+\theta'_j)}$$

for any $s \in [-\tau, 0]$ and $n \in \mathbb{Z}$, then by Theorem 3.4 and

$$\begin{aligned} & \chi K e^{\lambda^*((s-k)+\theta_0)} + \sum_{2 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j(s-k)+\theta_j)} \\ & + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j(s-k)+\theta'_j)} \\ & \geq \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(t - k) + \theta_i), \right. \\ & \quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(t - k) + \theta'_j), \chi \Gamma((t - k) + \theta_0) \right\} \\ & - \phi_{c_1}(n + c_1(s - k) + \theta_1) \end{aligned}$$

for any $s \in [-\tau, 0]$ and $n \in \mathbb{Z}$, we have

$$\begin{aligned} w_n^k(t) & \leq \chi K e^{\lambda^*((t-k)+\theta_0)} + \sum_{2 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j(t-k)+\theta_j)} \\ & + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j(t-k)+\theta'_j)} \quad \text{for any } t > 0. \end{aligned}$$

That is,

$$\begin{aligned} u_n^k(t) & \leq \phi_{c_1}(n + c_1(t - k) + \theta_1) + \chi K e^{\lambda^*((t-k)+\theta_0)} + \sum_{2 \leq j \leq m} A_{c_j} e^{\lambda_1(c_j)(n+c_j(t-k)+\theta_j)} \\ & + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j(t-k)+\theta'_j)} \quad \text{for any } t > 0. \end{aligned}$$

By the arbitrariness of $k \in \mathbb{N}$, we have that (5.1) holds.

When $m = 0$, the inequality (5.1) reduces to the following:

$$\limsup_{t > -k, k \rightarrow +\infty} v_n^k(t) \leq \chi K e^{\lambda^*(t+\theta_0)} + \sum_{1 \leq j \leq l} A_{c'_j} e^{\lambda_1(c'_j)(-n+c'_j t+\theta'_j)},$$

which holds obviously. This completes the proof. \square

Proof of Theorem 1.1. Define $v^k(t) = \{v_n^k(t)\}_{n \in \mathbb{Z}}$ with $v_n^k(t) := u_n(t + k; \psi^k)$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$, where

$$\begin{aligned} \psi^k & = \{\psi_n^k(s)\}_{k \in \mathbb{Z}}, \\ \psi_n^k(s) & = \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(s - k) + \theta_i), \right. \\ & \quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s - k) + \theta'_j), \chi \Gamma((s - k) + \theta_0) \right\} < K \end{aligned}$$

for any $(n, s) \in \mathbb{Z} \times [-\tau, 0]$. Note that

$$\begin{aligned} & \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i t + \theta_i), \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j t + \theta'_j), \chi \Gamma(t + \theta_0) \right\} \\ (5.4) \quad & \leq v_n^k(t) \leq v_n^{k+1}(t) \leq \min \{K, \vartheta_m^+(n, t), \vartheta_l^-(n, t), \vartheta^0(n, t)\} \end{aligned}$$

for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$. From Lemma 5.1 and by a diagonal extraction process, there exists a subsequence $\{v^{k_i}(t) = \{v_n^{k_i}(t)\}_{n \in \mathbb{Z}} : i \in \mathbb{N}\}$ such that $v^{k_i}(t)$

converges to a function $\Phi(t) = \{\Phi_n(t)\}_{n \in \mathbb{Z}}$ in \mathcal{T} ; that is, for any compact set $S \subset \mathbb{Z} \times \mathbb{R}$, $v_n^{k_i}(t)$ and $\frac{d}{dt}v_n^{k_i}(t)$ converge uniformly in $(n, t) \in S$ to $\Phi_n(t)$ and $\frac{d}{dt}\Phi_n(t)$, respectively. In view of $v_n^k(t) \leq v_n^{k+1}(t)$ for any $t > -k$, we have $\lim_{k \rightarrow +\infty} v_n^k(t) = \Phi_n(t)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$. The limit function is unique, whence all of the functions $v^k(t)$ converge to the function $\Phi(t)$ in \mathcal{T} as $k \rightarrow +\infty$. Since $u^{k_i}(t) = \{u_n^{k_i}(t)\}_{n \in \mathbb{Z}}$ satisfies (1.1), the limit function $\Phi(t) = \{\Phi_n(t)\}_{n \in \mathbb{Z}}$ is an entire solution of (1.1). In particular, it follows from (5.4) that (1.8) holds on $(n, t) \in \mathbb{Z} \times \mathbb{R}$.

Now we show (i); that is, $\frac{d}{dt}\Phi_n(t) > 0$ on \mathbb{R} for every $n \in \mathbb{Z}$. Since

$$\begin{aligned} \psi_n^k(s) &= \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(s - k) + \theta_i), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s - k) + \theta'_j), \chi\Gamma((s - k) + \theta_0) \right\} \\ &\leq \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(s + \varepsilon - k) + \theta_i), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s + \varepsilon - k) + \theta'_j), \chi\Gamma((s + \varepsilon - k) + \theta_0) \right\} \\ &= \psi_n^k(s + \varepsilon) \end{aligned}$$

for any $\varepsilon > 0$, $s \in [-\tau, 0]$, and $n \in \mathbb{Z}$, we have $u_n^k(t; \psi^k(\cdot)) \leq u_n^k(t; \psi^k(\cdot + \varepsilon))$ for any $(n, t) \in \mathbb{Z} \times [-\tau, +\infty)$. On the other hand, $\psi_n^k(s + \varepsilon) \leq u_n^k(s + \varepsilon; \psi^k(\cdot))$ for any $\varepsilon > 0$, $s \in [-\tau, 0]$, and $n \in \mathbb{Z}$, and hence,

$$u_n^k(t; \psi^k(\cdot)) \leq u_n^k(t; u_n^k(\cdot + \varepsilon; \psi^k(\cdot))) = u_n^k(t + \varepsilon; \psi^k(\cdot))$$

for any $(n, t) \in \mathbb{Z} \times [-\tau, +\infty)$. Thus, it follows from the arbitrariness of $\varepsilon > 0$; that $u_n^k(t)$ is increasing on t ; that is, $v^k(t)$ is increasing on t . Therefore, $\Phi'_n(t) \geq 0$ on \mathbb{R} for every $n \in \mathbb{Z}$. Since $\Phi'_n(t)$ satisfies

$$\begin{aligned} \Phi''_n(t) &= D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [\Phi'_{n-i}(t) - \Phi'_n(t)] \\ (5.5) \quad &\quad - d\Phi'_n(t) + \sum_{i \in \mathbb{Z}} J(i) b'(\Phi_{n-i}(t - \tau)) \Phi'_{n-i}(t - \tau), \end{aligned}$$

we have that $\Phi'_n(t)$ satisfies

$$\Phi'_n(t) = \Phi'_n(s) e^{-(D+d)(t-s)} + \int_s^t e^{-(D+d)(t-r)} R_n(\Phi)(r) dr \quad \text{for any } s < t,$$

where $R_n(\Phi)(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) \Phi'_{n-i}(t) + \sum_{i \in \mathbb{Z}} J(i) b'(\Phi_{n-i}(t - \tau)) \Phi'_{n-i}(t - \tau) \geq 0$. Obviously, for each $n \in \mathbb{Z}$, if there exists $t_0 \in \mathbb{R}$ such that $\Phi'_n(t_0) > 0$, then $\Phi'_n(t) > 0$ for any $t > t_0$. Therefore, there must be $\Phi'_n(t) > 0$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$. We argue by a contradiction. In fact, assume that, for some $n_1 \in \mathbb{Z}$, there is t_1 such that $\Phi'_{n_1}(t_1) = 0$ and, hence, then $\Phi'_{n_1}(t) = 0$ for any $t \leq t_1$, which implies that $\lim_{t \rightarrow -\infty} \Phi_{n_1}(t) = \Phi_{n_1}(t_1) > 0$. But following from (1.8), we have $\lim_{t \rightarrow -\infty} \Phi_{n_1}(t) = 0$, which yields a contradiction.

Now we prove (vii). For the sake of convenience, we denote

$$\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$$

by $\Phi(t; \theta_0)$ and $\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l}(t)$ by $\Phi(t; -\infty)$. For $\chi \in \{0, 1\}$, let

$$\begin{aligned} \psi^k(t)_\chi &= \left\{ \psi_n^k(s)_\chi \right\}_{k \in \mathbb{Z}}, \\ \psi_n^k(s)_\chi &= \max \left\{ \max_{1 \leq i \leq m} \phi_{c_i}(n + c_i(s - k) + \theta_i), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s - k) + \theta'_j), \chi \Gamma((s - k) + \theta_0) \right\}, \end{aligned}$$

and $v^k(t)_\chi = \{v_n^k(t)_\chi\}_{n \in \mathbb{Z}}$ with $v_n^k(t)_\chi := u_n(t + k; \psi^k(\cdot)_\chi)$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$. Set $\bar{v}^k(t) = v^k(t)_1 - v^k(t)_0 = \{v_n^k(t)_1 - v_n^k(t)_0\}_{n \in \mathbb{Z}}$. Then $\bar{v}^k(t)$ satisfies $0 \leq \bar{v}^k(t) \leq K$ for any $t \in [-\tau - k, +\infty)$ and

$$\frac{d}{dt} \bar{v}_n^k(t) \leq D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [\bar{v}_{n-i}^k(t) - \bar{v}_n^k(t)] - d\bar{v}_n^k(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) \bar{v}_{n-i}^k(t - \tau).$$

Noting that $\bar{v}_n^k(s) \leq Ke^{\lambda_*(s+\theta_0)}$ for any $s \in [-\tau - k, -k]$ and $w_n^k(t) = Ke^{\lambda_*(t+\theta_0)}$ satisfies

$$\frac{d}{dt} w_n^k(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [w_{n-i}^k(t) - w_n^k(t)] - dw_n^k(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) w_{n-i}^k(t - \tau)$$

for any $t \in [-\tau - k, +\infty)$, it follows from Theorem 3.4 that $0 \leq \bar{v}_n^k(t) \leq Ke^{\lambda_*(t+\theta_0)}$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$ and $k \in \mathbb{N}$. Note that $\lim_{k \rightarrow +\infty} v_n^k(t)_0 = \Phi_n(t; -\infty)$ and $\lim_{k \rightarrow +\infty} v_n^k(t)_1 = \Phi_n(t; \theta_0)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$. Therefore, there must be $0 < \Phi_n(t; \theta_0) - \Phi_n(t; -\infty) \leq Ke^{\lambda_*(t+\theta_0)}$ for all $(n, t) \in \mathbb{Z} \times \mathbb{R}$, which implies that $\Phi(t; \theta_0)$ converges uniformly on $(n, t) \in \mathbb{Z} \times (-\infty, a]$ to $\Phi(t; -\infty)$ as $\theta_0 \rightarrow -\infty$ for any $a \in \mathbb{R}$. For any sequence $\theta_0^k \rightarrow -\infty$ ($k \rightarrow +\infty$), the functions $\Phi(t; \theta_0^k)$ converge to a solution of (1.1) in \mathcal{T} , which turns out to be $\Phi(t; -\infty)$. The limit does not depend on the sequence θ_0^k , whence all of the functions $\Phi(t; \theta_0)$ converge to $\Phi(t; -\infty)$ in \mathcal{T} as $\theta_0 \rightarrow -\infty$. The assertion as $\theta_0 \rightarrow +\infty$ is obvious.

We next prove (viii). Assume $\chi = 1$. Similarly to that in (vii), we denote

$$\Phi_{c_1, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$$

by $\Phi(t)_{\theta_i}$ and $\Phi_{c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_m; c'_1, \dots, c'_l; \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_m; \theta'_1, \dots, \theta'_l; \theta_0}(t)$ by $\Phi(t)_\infty$. Set

$$\begin{aligned} \psi^k(t)_{\theta_i} &= \left\{ \psi_n^k(s)_{\theta_i} \right\}_{k \in \mathbb{Z}}, \\ \psi_n^k(s)_{\theta_i} &= \max \left\{ \max_{1 \leq j \leq m} \phi_{c_j}(n + c_j(s - k) + \theta_j), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s - k) + \theta'_j), \Gamma((s - k) + \theta_0) \right\}, \end{aligned}$$

and $v^k(t)_{\theta_i} = \{v_n^k(t)_{\theta_i}\}_{n \in \mathbb{Z}}$ with $v_n^k(t)_{\theta_i} := u_n(t + k; \psi^k(\cdot)_{\theta_i})$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$. Take

$$\begin{aligned} \psi^k(t)_\infty &= \left\{ \psi_n^k(s)_\infty \right\}_{k \in \mathbb{Z}}, \\ \psi_n^k(s)_\infty &= \max \left\{ \max_{j \in \{1, \dots, i-1, i+1, \dots, m\}} \phi_{c_j}(n + c_j(s - k) + \theta_j), \right. \\ &\quad \left. \max_{1 \leq j \leq l} \phi_{c'_j}(-n + c'_j(s - k) + \theta'_j), \Gamma((s - k) + \theta_0) \right\}, \end{aligned}$$

and $v^k(t)_\infty = \{v_n^k(t)_\infty\}_{n \in \mathbb{Z}}$ with $v_n^k(t)_\infty := u_n(t+k; \psi^k(\cdot)_\infty)$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$. Set $\hat{v}^k(t) = v^k(t)_{\theta_i} - v^k(t)_\infty = \{v_n^k(t)_{\theta_i} - v_n^k(t)_\infty\}_{n \in \mathbb{Z}}$. Then $\hat{v}^k(t)$ satisfies $0 \leq \hat{v}_n^k(t) \leq K$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$ and

$$\frac{d}{dt} \hat{v}_n^k(t) \leq D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [\hat{v}_{n-i}^k(t) - \hat{v}_n^k(t)] - d\hat{v}_n^k(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) \hat{v}_{n-i}^k(t - \tau).$$

Noting that $\hat{v}_n^k(s) \leq \phi_{c_i}(n + c_i s + \theta_i) \leq A_{c_i} e^{\lambda_1(c_i)(n + c_i s + \theta_i)}$ for any $s \in [-\tau - k, -k]$ and that $\bar{w}_n^k(t) = A_{c_i} e^{\lambda_1(c_i)(n + c_i t + \theta_i)}$ satisfies

$$\frac{d}{dt} \bar{w}_n^k(t) = D \sum_{i \in \mathbb{Z} \setminus \{0\}} I(i) [\bar{w}_{n-i}^k(t) - \bar{w}_n^k(t)] - d\bar{w}_n^k(t) + b'(0) \sum_{i \in \mathbb{Z}} J(i) \bar{w}_{n-i}^k(t - \tau)$$

for any $t \in [-\tau - k, +\infty)$, it follows from Theorem 3.4 that $0 \leq \hat{v}_n^k(t) \leq A_{c_i} e^{\lambda_1(c_i)(n + c_i t + \theta_i)}$ for any $(n, t) \in \mathbb{Z} \times [-\tau - k, +\infty)$ and $k \in \mathbb{N}$. Since $\lim_{k \rightarrow +\infty} v^k(t)_{\theta_i} = \Phi(t)_{\theta_i}$ and $\lim_{k \rightarrow +\infty} v^k(t)_\infty = \Phi(t)_\infty$, we have $0 < \Phi_n(t)_{\theta_i} - \Phi_n(t)_\infty \leq A_{c_i} e^{\lambda_1(c_i)(n + c_i t + \theta_i)}$ for all $(n, t) \in \mathbb{Z} \times \mathbb{R}$, which implies that $\Phi(t)_{\theta_i}$ converges uniformly on $(n, t) \in \{n : n \leq N_0, n \in \mathbb{Z}\} \times (-\infty, a]$ to $\Phi(t)_\infty$ as $\theta_i \rightarrow -\infty$ for any $N_0 \in \mathbb{Z}$ and $a \in \mathbb{R}$. For any sequence $\theta_i^k \rightarrow -\infty$ ($k \rightarrow +\infty$), the functions $\Phi(t)_{\theta_i^k}$ converge to a solution of (1.1), which must be $\Phi(t)_\infty$. Since the limit is independent of the sequence θ_i^k , all of the functions $\Phi(t)_{\theta_i}$ converge to $\Phi(t)_\infty$ in \mathcal{T} as $\theta_i \rightarrow -\infty$. The assertion as $\theta_j' \rightarrow -\infty$ and the case $\chi = 0$ can be proved similarly.

Using the inequality (1.8), we can prove (ii)–(vi) and (ix) of Theorem 1.1. □

6. Proof of Theorem 1.2. In this section, we prove Theorem 1.2. We prove only the continuous dependence of the entire solution on the parameters c, c', θ, θ' , and θ_0 and the uniqueness of the entire solutions satisfies (1.9). Other conclusions follow immediately from Theorem 1.1.

Consider (1.4) or

$$(6.1) \quad u'_n(t) = \frac{D}{2} [u_{n+1} + u_{n-1} - 2u_n] + f(u_n(t)),$$

where f satisfies the conditions given after (1.4). Let $\phi_c(n + ct)$ be a traveling wave front of (6.1) with wave speed $c > c^*$. As done in section 1, we normalize $\phi_c(n + ct)$ so that $\phi(0) = \frac{1}{2}$. Then the functions $\phi_c(z)$ are continuous with respect to $c \in (c^*, +\infty)$ in the norms $C^1_{loc}(\mathbb{R})$ (see [20, p. 1267] for the definition of these norms). Indeed, if $c_l \rightarrow c \in (c^*, +\infty)$, then by the unique boundedness of $|\phi'_{c_l}(z)|$ and $|\phi''_{c_l}(z)|$ in $z \in \mathbb{R}$ on $l \in \mathbb{N}$ and by a diagonal extraction process, there exists a subsequence c_{l_i} such that $\phi_{c_{l_i}} \rightarrow \phi$ in $C^1_{loc}(\mathbb{R})$, where ϕ is a solution of

$$c\phi'(z) = \frac{D}{2} [\phi(z+1) + \phi(z-1) - 2\phi(z)] + f(\phi(z)) \text{ in } z \in \mathbb{R}.$$

Obviously, ϕ is nondecreasing in \mathbb{R} and is not a constant and $\phi(0) = \frac{1}{2}$. By the assumptions of f , we have $\phi(-\infty) = 0$ and $\phi(+\infty) = 1$. Thus, ϕ is a traveling wave front of (6.1) with wave speed c . Following Chen and Guo [5], we have $\phi \equiv \phi_c$. Finally, the whole sequence $\phi_{c_l} \rightarrow \phi_{c_0}$ in $C^1_{loc}(\mathbb{R})$ as $l \rightarrow +\infty$. In view of [4, 5], we know that α_c and A_c defined by (1.6) and (1.7), respectively, are still valid for (6.1). In addition, there exactly is $\lambda_* = f'(0)$ for (6.1).

LEMMA 6.1. For (6.1), $\alpha_c = \lim_{z \rightarrow -\infty} \phi_c(z) e^{-\lambda_1(c)z}$ is continuous in $c \in (c^*, +\infty)$.

Proof. Fix $c_0 \in (c^*, +\infty)$ and let $c_l \rightarrow c_0$ as $l \rightarrow +\infty$ with $c_l > c^*$ for each $l \in \mathbb{N}$. Then by Chen and Guo [4] (see also Ma, Weng, and Zou [25]), we know that, for each $c \in (c^*, +\infty)$, there exists a unique traveling wave front $\tilde{\phi}_c$ such that $\tilde{\phi}'_c(\cdot) > 0$, $\tilde{\phi}_c(-\infty) = 0$, $\tilde{\phi}_c(+\infty) = 1$, and

$$\lim_{z \rightarrow -\infty} \tilde{\phi}_c(z) e^{-\lambda_1(c)z} = 1.$$

Then we have that $\tilde{\phi}_{c_l} \rightarrow \tilde{\phi}_{c_0}$ in $C^1_{loc}(\mathbb{R})$ as $l \rightarrow +\infty$. In fact, since $c_l \rightarrow c_0$ as $l \rightarrow +\infty$, then there exist a subsequence c_{l_i} and a function $\tilde{\phi}$ such that $\tilde{\phi}_{c_{l_i}} \rightarrow \tilde{\phi}$ in $C^1_{loc}(\mathbb{R})$, where $\tilde{\phi}$ is nondecreasing and satisfies

$$c\tilde{\phi}'(z) = \frac{D}{2} [\tilde{\phi}(z+1) + \tilde{\phi}(z-1) - 2\tilde{\phi}(z)] + f(\tilde{\phi}(z)) \text{ in } z \in \mathbb{R}.$$

On the other hand, by Chen and Guo [4], there exist two constants $q > 1$ and $\beta > 1$, independent of c_l , such that

$$e^{\lambda_1(c_l)z} - qe^{\beta\lambda_1(c_l)z} \leq \tilde{\phi}_{c_l}(z) \leq e^{\lambda_1(c_l)z} + qe^{\beta\lambda_1(c_l)z} \text{ for any } z \in \mathbb{R}$$

and therefore, as $l \rightarrow \infty$,

$$e^{\lambda_1(c_0)z} - qe^{\beta\lambda_1(c_0)z} \leq \tilde{\phi}(z) \leq e^{\lambda_1(c_0)z} + qe^{\beta\lambda_1(c_0)z} \text{ for any } z \in \mathbb{R},$$

which implies that $\tilde{\phi}(z)$ is not a constant and satisfies $\tilde{\phi}(z)e^{-\lambda_1(c_0)z} = 1$. Then it follows from Chen and Guo [5] (see also Ma, Weng, and Zou [25]) that $\tilde{\phi} \equiv \tilde{\phi}_{c_0}$. Consequently, the whole sequence $\tilde{\phi}_{c_l} \rightarrow \tilde{\phi}_{c_0}$ in $C^1_{loc}(\mathbb{R})$ as $l \rightarrow +\infty$.

Now let $\tilde{\phi}_{c_0}(\varsigma_0) = \frac{1}{2}$ and $\tilde{\phi}_{c_l}(\varsigma_l) = \frac{1}{2}$. Then we have that $\varsigma_l \rightarrow \varsigma_0$ as $l \rightarrow +\infty$. Assume that this assertion is not true. Take $\varsigma_l \rightarrow \bar{\varsigma} \neq \varsigma_0$ as $l \rightarrow +\infty$ (up to extraction of some subsequence). If $|\bar{\varsigma}| < \infty$, then, by $\tilde{\phi}_{c_l} \rightarrow \tilde{\phi}_{c_0}$ in $C^1_{loc}(\mathbb{R})$ as $l \rightarrow +\infty$, we have $\tilde{\phi}_{c_0}(\bar{\varsigma}) = \frac{1}{2} = \tilde{\phi}_{c_0}(\varsigma_0)$, which is impossible since $\frac{d}{dz}\tilde{\phi}_{c_0}(z) > 0$ for any $z \in \mathbb{R}$ and $\bar{\varsigma} \neq \varsigma_0$. If $\bar{\varsigma} = +\infty$, then $\tilde{\phi}_{c_l}(\varsigma_0 + 1) < \tilde{\phi}_{c_l}(\varsigma_l) = \frac{1}{2}$ for sufficiently large l implies that $\tilde{\phi}_{c_0}(\varsigma_0 + 1) \leq \frac{1}{2} = \tilde{\phi}_{c_0}(\varsigma_0)$, which is also impossible. Similarly, $\bar{\varsigma} = -\infty$ is impossible. Thus, we conclude that $\varsigma_l \rightarrow \varsigma_0$ as $l \rightarrow +\infty$.

Again applying Chen and Guo [5], we have that $\phi_{c_0}(\cdot) = \tilde{\phi}_{c_0}(\varsigma_0 + \cdot)$ and $\phi_{c_l}(\cdot) = \tilde{\phi}_{c_l}(\varsigma_l + \cdot)$. Since

$$\lim_{z \rightarrow -\infty} \phi_{c_0}(z) e^{-\lambda_1(c_0)(z+\varsigma_0)} = \lim_{z \rightarrow -\infty} \tilde{\phi}_{c_0}(z + \varsigma_0) e^{-\lambda_1(c_0)(z+\varsigma_0)} = 1,$$

we have $\lim_{z \rightarrow -\infty} \phi_{c_0}(z)e^{-\lambda_1(c_0)z} = e^{\lambda_1(c_0)\varsigma_0} = \alpha_{c_0}$. Similarly, we have

$$\lim_{z \rightarrow -\infty} \phi_{c_l}(z) e^{-\lambda_1(c_l)z} = e^{\lambda_1(c_l)\varsigma_l} = \alpha_{c_l}.$$

Finally, there holds $\alpha_{c_l} \rightarrow \alpha_{c_0}$ as $l \rightarrow +\infty$. This completes the proof. \square

Recall that

$$A_c = \inf \left\{ A > 0 : Ae^{\lambda_1(c)z} \geq \phi_c(z) \text{ in } z \in \mathbb{R} \right\}.$$

LEMMA 6.2. For (6.1), A_c is continuous on $c \in (c^*, +\infty)$.

Proof. Fix $c_0 \in (c^*, +\infty)$ and let $c_l \rightarrow c_0$ as $l \rightarrow +\infty$ with $c_l > c^*$ for each $l \in \mathbb{N}$. We prove the theorem by way of contradiction. Assume $A_{c_l} \rightarrow A_0 \in \mathbb{R} \cup \{\infty\}$

as $l \rightarrow \infty$ (up to extraction of some subsequence) and $A_0 \neq A_{c_0}$. Since $A_{c_l} \geq e^{-\lambda_1(c_l)z} \phi_{c_l}(z)$ for any $z \in \mathbb{R}$, $A_0 \geq e^{-\lambda_1(c_0)z} \phi_{c_0}(z)$ and, hence, $A_0 > A_{c_0}$. Fix $b = \min\{\frac{A_0 + A_{c_0}}{2}, A_{c_0} + 1\}$. Then there exists $L \in \mathbb{N}$ such that for any $l > L$, $A_{c_l} > b$. On the other hand, since $\alpha_{c_l} \rightarrow \alpha_{c_0} \leq A_{c_0}$ and $\lambda_1(c_l) \rightarrow \lambda_1(c_0)$, there exists a constant $Z_0 > 0$, independent of c_l , such that $\phi_{c_l}(z)e^{-\lambda_1(c_l)z} \leq b$ for any $|z| > Z_0$. For $z \in [-Z_0, Z_0]$, by $\phi_{c_0}(z)e^{-\lambda_1(c_0)z} \leq A_{c_0}$, $\phi_{c_l}(z) \rightarrow \phi_{c_0}(z)$ in $C^1_{loc}(\mathbb{R})$, and the equicontinuity of $e^{-\lambda_1(c_l)z}$ on l , there exists $L' > L$ such that $\phi_{c_l}(z)e^{-\lambda_1(c_l)z} \leq b$ for any $l > L'$ and $z \in [-Z_0, Z_0]$. Therefore, $\phi_{c_l}(z)e^{-\lambda_1(c_l)z} \leq b$ for any $l > L'$ and $z \in \mathbb{R}$, which contradicts $A_{c_l} > b$ for any $l > L$. \square

Before proving Theorem 1.2, we first consider the following linear Cauchy problem:

$$(6.2) \quad \begin{cases} \frac{d}{dt}u_n(t) = \frac{D}{2} [u_{n+1}(t) + u_{n-1}(t) - 2u_n(t)] + f'(0)u_n(t), & t > 0, \\ u_n(0) = u_n^0, \end{cases}$$

where $u^0 = \{u_n^0\}_{n \in \mathbb{Z}} \in l^\infty$. By Theorem 3.3, we know that (6.2) admits a unique solution $u(t) = u(t; u^0) = \{u_n(t)\}_{n \in \mathbb{Z}}$ on $t \in [0, +\infty)$. By using the discrete Fourier transformation, we can exactly solve the solution $u(t) = \{u_n(t)\}_{n \in \mathbb{Z}}$ of (6.2) as follows:

$$(6.3) \quad u_n(t) = \frac{1}{\pi} e^{f'(0)t} \sum_{i=-\infty}^{+\infty} u_i^0 \int_0^\pi \cos((i-n)\omega) e^{Dt(\cos\omega-1)} d\omega.$$

This formulation is very crucial for the proof of Theorem 1.2.

Proof of Theorem 1.2. We prove only the case $\varrho = \varrho' = \chi = 1$. Consider a sequence

$$(c_k, c'_k, \theta_k, \theta'_k, \theta_{0,k}) \rightarrow (c, c', \theta, \theta', \theta_0) \in (c^*, +\infty)^2 \times \mathbb{R}^3 \quad \text{as } k \rightarrow +\infty.$$

For given $(c_k, c'_k, \theta_k, \theta'_k, \theta_{0,k})$ and $(c, c', \theta, \theta', \theta_0)$, it follows from Theorem 1.1 that there exist entire solutions $\Phi_{c_k; c'_k; \theta_k; \theta'_k; \theta_{0,k}}(t)$ and $\Phi_{c; c'; \theta; \theta'; \theta_0}(t)$ of (1.4) satisfying (1.9). For the sake of convenience, set $\Phi^k(t) = \{\Phi_n^k(t)\}_{n \in \mathbb{Z}} = \Phi_{c_k; c'_k; \theta_k; \theta'_k; \theta_{0,k}}(t)$ and $\Phi(t) = \{\Phi_n(t)\}_{n \in \mathbb{Z}} = \Phi_{c; c'; \theta; \theta'; \theta_0}(t)$.

Using Lemma 5.1, there exists a function $\tilde{\Phi}(t) = \{\tilde{\Phi}_n(t)\}_{n \in \mathbb{Z}}$ such that $\Phi_n^k(t) \rightarrow \tilde{\Phi}_n(t)$ as $k \rightarrow \infty$ (up to extraction of some subsequence) in \mathcal{T} . In particular, the function $\tilde{\Phi}(t) = \{\tilde{\Phi}_n(t)\}_{n \in \mathbb{Z}}$ is also an entire solution of (6.1) (or (1.4)). By passage to the limit $k \rightarrow +\infty$ in (1.9), the function $\tilde{\Phi}(t) = \{\tilde{\Phi}_n(t)\}_{n \in \mathbb{Z}}$ fulfills the estimates

$$(6.4) \quad \begin{aligned} & \max \{ \phi_c(n+ct+\theta), \Gamma(t+\theta_0), \phi_{c'}(-n+c't+\theta') \} \\ & \leq \tilde{\Phi}_n(t) \leq \min \left\{ 1, \phi_c(n+ct+\theta) + e^{\lambda_*(t+\theta_0)} + A_{c'} e^{\lambda_1(c')}(-n+c't+\theta'), \right. \\ & \quad \Gamma(t+\theta_0) + A_c e^{\lambda_1(c)(n+ct+\theta)} + A_{c'} e^{\lambda_1(c')}(-n+c't+\theta'), \\ & \quad \left. \phi_{c'}(-n+c't+\theta') + e^{\lambda_*(t+\theta_0)} + A_c e^{\lambda_1(c)(n+ct+\theta)} \right\} \end{aligned}$$

for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$.

Let us now prove that $\tilde{\Phi}_n(t) = \Phi_n(t)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$. Recall that the functions $v^k(t) = \{v_n^k(t)\}_{n \in \mathbb{N}}$, which are solutions of the Cauchy problems

$$\frac{d}{dt}v_n^k(t) = \frac{D}{2} [v_{n+1}^k(t) + v_{n-1}^k(t) - 2v_n^k(t)] + f(v_n^k(t)), \quad t > -k, \quad n \in \mathbb{N},$$

with the initial conditions

$$v_n^k(-k) = v_{n,0}^k = \max \{ \phi_c(n - ck + \theta), \Gamma(-k + \theta_0), \phi_{c'}(-n - c'k + \theta') \},$$

converge to the function $\Phi(t) = \{ \Phi_n(t) \}_{n \in \mathbb{N}}$ in \mathcal{T} ; see the proof of Theorem 1.1. Let us now compare the functions $\tilde{\Phi}(t)$ to the functions $v^k(t)$ for $t > -k$. Following (6.4), we get that

$$(6.5) \quad \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| \leq \begin{cases} e^{f'(0)(-k+\theta_0)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} & \text{if } n \geq n_k^+ + 2, \\ A_c e^{\lambda_1(c)(n-ck+\theta)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} & \text{if } n_k^- + 1 \leq n \leq n_k^+ - 1, \\ e^{f'(0)(-k+\theta_0)} + A_c e^{\lambda_1(c)(n-ck+\theta)} & \text{if } n \leq n_k^- - 2 \end{cases}$$

for sufficiently large k , where n_k^- and n_k^+ are two integers defined as follows:

$$n_k^- = \text{int} \left[-c'k + \frac{f'(0)}{\lambda_1(c')}k + \theta' - \frac{1}{\lambda_1(c')} \ln \frac{1}{\alpha_{c'}} - \frac{f'(0)}{\lambda_1(c')} \theta_0 \right],$$

$$n_k^+ = \text{int} \left[ck - \frac{f'(0)}{\lambda_1(c)}k - \theta + \frac{1}{\lambda_1(c)} \ln \frac{1}{\alpha_c} + \frac{f'(0)}{\lambda_1(c)} \theta_0 \right].$$

n_k^- is obtained by comparing $\phi_{c'}(-n - c'k + \theta')$ and $\Gamma(-k + \theta_0)$ and using the asymptotic behaviors of $\phi_{c'}$ and Γ . Similarly, we obtain n_k^+ .

For any $x \in \mathbb{R}$, define

$$\text{int}[x] = \max \{ m : m \in \mathbb{Z}, m \leq x \}.$$

Fix $t_0 > -k$. For any $n \in \mathbb{Z}$ and $k \in \mathbb{N}$ with $k > -t_0$, define

$$a_n^k = \frac{1}{2\pi} e^{-D(t_0+k)} \int_0^\pi \cos(nw) e^{D(t_0+k) \cos w} dw.$$

By Weng, Huang, and Wu [40], we know that $a_n^k = a_{-n}^k > 0$ for any $n \in \mathbb{Z}$ and $\sum_{n=-\infty}^{+\infty} a_n^k = 1$. Furthermore, for $n > 1$ there is

$$\begin{aligned} a_n^k &= \frac{1}{2\pi} \int_0^\pi \cos(nw) e^{D(t_0+k)[\cos w - 1]} dw = \frac{1}{2n\pi} \int_0^\pi e^{D(t_0+k)[\cos w - 1]} d \sin(nw) \\ &= \frac{D(t_0+k)}{2n\pi} \int_0^\pi \sin(nw) \sin we^{D(t_0+k)[\cos w - 1]} dw \\ &= \frac{D(t_0+k)}{4n\pi} \int_0^\pi [\cos((n-1)w) - \cos((n+1)w)] e^{D(t_0+k)[\cos w - 1]} dw \\ &= \frac{D(t_0+k)}{2n} (a_{n-1}^k - a_{n+1}^k), \end{aligned}$$

which implies that $a_{n-1}^k > a_{n+1}^k$ for any $n \in \mathbb{N}$. By the symmetry, $a_{-n-1}^k < a_{-n+1}^k$ for any $n \in \mathbb{N}$.

We claim that, for any given $\eta \in (0, +\infty)$, there hold

$$(6.6) \quad a_{\text{int}[\eta k]}^k \rightarrow 0 \quad \text{and} \quad a_{\text{int}[\eta k] - 1}^k \rightarrow 0 \quad \text{as } k \rightarrow +\infty.$$

Consider $a_{\text{int}[\eta k]}^k \rightarrow 0$ as $k \rightarrow +\infty$. If the assertion is false, then we can assume that $a_{\text{int}[\eta k]}^k \rightarrow \delta > 0$ as $k \rightarrow +\infty$ (up to extraction of some subsequence). Taking k sufficiently large such that $\frac{\text{int}[\eta k]}{2} > \frac{4}{\delta} + 1$ and $a_{\text{int}[\eta k]}^k > \frac{\delta}{2}$, by $a_{n-1}^k > a_{n+1}^k$ for $n \in \mathbb{N}$, we have $\sum_{n=0}^{\text{int}[\eta k]} a_n^k > 1$, which contradicts the fact $\sum_{n=-\infty}^{+\infty} a_n^k < 1$. Therefore, $\lim_{k \rightarrow +\infty} a_{\text{int}[\eta k]}^k = 0$. Similarly, we have $\lim_{k \rightarrow +\infty} a_{\text{int}[\eta k]-1}^k = 0$. For any $m > N$,

$$\begin{aligned} \sum_{n=N \geq 1}^m a_n^k &= \frac{D(t_0+k)}{2N} (a_{N-1}^k - a_{N+1}^k) + \frac{D(t_0+k)}{2(N+1)} (a_N^k - a_{N+2}^k) \\ &\quad + \frac{D(t_0+k)}{2(N+2)} (a_{N+1}^k - a_{N+3}^k) \\ &\quad + \dots + \frac{D(t_0+k)}{2(m-1)} (a_{m-2}^k - a_m^k) + \frac{D(t_0+k)}{2m} (a_{m-1}^k - a_{m+1}^k) \\ &= \frac{D(t_0+k)}{2N} a_{N-1}^k + \frac{D(t_0+k)}{2(N+1)} a_N^k - \frac{D(t_0+k)}{N(N+2)} a_{N+1}^k \\ &\quad - \frac{D(t_0+k)}{(N+1)(N+3)} a_{N+2}^k - \dots - \frac{D(t_0+k)}{(m-3)(m-1)} a_{m-2}^k \\ &\quad - \frac{D(t_0+k)}{(m-2)m} a_{m-1}^k - \frac{D(t_0+k)}{2(m-1)} a_m^k - \frac{D(t_0+k)}{2m} a_{m+1}^k \\ &\leq \frac{D(t_0+k)}{2N} a_{N-1}^k + \frac{D(t_0+k)}{2(N+1)} a_N^k. \end{aligned}$$

Therefore,

$$\sum_{n=N \geq 1}^{\infty} a_n^k \leq \frac{D(t_0+k)}{2N} a_{N-1}^k + \frac{D(t_0+k)}{2(N+1)} a_N^k.$$

Consequently, we obtain that for any given $\eta > 0$,

$$(6.7) \quad \sum_{n=\text{int}[\eta k]}^{\infty} a_n^k \rightarrow 0 \quad \text{and} \quad \sum_{n=\text{int}[\eta k]+1}^{\infty} a_n^k \rightarrow 0 \quad \text{as} \quad k \rightarrow +\infty$$

due to $\frac{D(t_0+k)}{[\eta k]} \rightarrow \frac{D}{\eta}$ as $k \rightarrow +\infty$.

Now fix $(n_0, t_0) \in \mathbb{Z} \times \mathbb{R}$; we estimate $\tilde{\Phi}_{n_0}(t_0) - v_{n_0}^k(t_0)$ as $k \rightarrow +\infty$. We compare $\tilde{\Phi}_n(t) - v_n^k(t)$ with a solution of the linear equation

$$\frac{d}{dt} u_n(t) = \frac{D}{2} [u_{n+1}(t) + u_{n-1}(t) - 2u_n(t)] + f'(0) u_n(t), \quad t > -k,$$

with the initial condition $u_n(-k) = |\tilde{\Phi}_n(-k) - v_n^k(-k)|$. Using (6.3), we deduce that

$$\begin{aligned}
 & \left| \tilde{\Phi}_{n_0}(t_0) - v_{n_0}^k(t_0) \right| \\
 \leq & \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=-\infty}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 \leq & \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \int_0^\pi \cos((n_k^+ + 1 - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \int_0^\pi \cos((n_k^+ - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^-+1}^{n_k^+-1} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \int_0^\pi \cos((n_k^- - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \int_0^\pi \cos((n_k^- - 1 - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 & + \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=-\infty}^{n_k^- - 2} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw.
 \end{aligned}$$

Call I, II, III, IV, V, VI, and VII the seven terms on the right-hand side of this last inequality. We have

$$\begin{aligned}
 \text{I} &= \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\
 &\leq \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} e^{f'(0)(-k+\theta_0)} dw \\
 &\quad + \frac{A_{c'}}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} e^{\lambda_1(c')(-n-c'k+\theta')} dw \\
 &= \text{I}_1 + \text{I}_2.
 \end{aligned}$$

Obviously,

$$\begin{aligned}
 \text{I}_1 &= \frac{1}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n - n_0)w) e^{D(t_0+k)[\cos w - 1]} e^{f'(0)(-k+\theta_0)} dw \\
 &= 2e^{f'(0)(t_0+\theta_0)} \sum_{n=n_k^++2}^{+\infty} a_{n-n_0}^k \rightarrow 0 \text{ as } k \rightarrow +\infty
 \end{aligned}$$

due to the fact $\frac{n_k^++2-n_0}{k} > \frac{c\lambda_1(c)-f'(0)}{2\lambda_1(c)} > 0$ for sufficiently large k and (6.7). In

addition, we have

$$\begin{aligned} I_2 &= \frac{A_{c'}}{\pi} e^{f'(0)(t_0+k)} \sum_{n=n_k^++2}^{+\infty} \int_0^\pi \cos((n-n_0)w) e^{D(t_0+k)[\cos w-1]} e^{\lambda_1(c')(-n-c'k+\theta')} dw \\ &= 2A_{c'} e^{(f'(0)-\lambda_1(c')c')k} e^{f'(0)(t_0+\theta')} \sum_{n=n_k^++2}^{+\infty} a_{n-n_0}^k e^{-\lambda_1(c')n} \\ &\leq 2A_{c'} e^{f'(0)(t_0+\theta')} \sum_{n=n_k^++2}^{+\infty} a_{n-n_0}^k \rightarrow 0 \text{ as } k \rightarrow \infty. \end{aligned}$$

Consider II. In this case, let

$$0 \leq \delta = \left(ck - \frac{f'(0)}{\lambda_1(c)}k - \theta + \frac{1}{\lambda_1(c)} \ln \frac{1}{\alpha_c} + \frac{f'(0)}{\lambda_1(c)}\theta_0 \right) - n_k^+ < 1.$$

Then we have

$$\begin{aligned} \text{II} &= 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k \left| \tilde{\Phi}_{n_k^++1}(-k) - v_{n_k^++1}^k(-k) \right| \\ &\leq 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k \left[e^{f'(0)(-k+\theta_0)} + A_{c'} e^{\lambda_1(c')(-n_k^+-1-c'k+\theta')} \right] \\ &\quad + 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k \left[A_{c'} e^{\lambda_1(c')(-n_k^+-1-c'k+\theta')} + A_c e^{\lambda_1(c)(n_k^++1-ck+\theta)} \right] \\ &= 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k \left[e^{f'(0)(-k+\theta_0)} + 2A_{c'} e^{\lambda_1(c')(-n_k^+-1-c'k+\theta')} \right] \\ &\quad + 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k A_c e^{\lambda_1(c) \left\{ \left(ck - \frac{f'(0)}{\lambda_1(c)}k - \theta + \frac{1}{\lambda_1(c)} \ln \frac{1}{\alpha_c} + \frac{f'(0)}{\lambda_1(c)}\theta_0 \right) - \delta + 1 - ck + \theta \right\}} \\ &= 2e^{f'(0)(t_0+k)} a_{n_k^++1-n_0}^k \left[e^{f'(0)(-k+\theta_0)} + 2A_{c'} e^{\lambda_1(c')(-n_k^+-1-c'k+\theta')} \right] \\ &\quad + \frac{2A_c}{\alpha_c} e^{\lambda_1(c)(1-\delta)+f'(0)(t_0+\theta_0)} a_{n_k^++1-n_0}^k \\ &\rightarrow 0 \quad \text{as } k \rightarrow +\infty. \end{aligned}$$

Similarly, we have that III $\rightarrow 0$ as $k \rightarrow +\infty$. For IV, we have

$$\begin{aligned} \text{IV} &= \frac{1}{\pi} e^{f'(0)(t_0+k)} \\ &\quad \times \sum_{n=n_k^-+1}^{n_k^+-1} \int_0^\pi \cos((n-n_0)w) e^{D(t_0+k)[\cos w-1]} \left| \tilde{\Phi}_n(-k) - v_n^k(-k) \right| dw \\ &\leq 2e^{f'(0)(t_0+k)} \sum_{n=n_k^-+1}^{n_k^+-1} a_{n-n_0}^k \left[A_c e^{\lambda_1(c)(n-ck+\theta)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} \right] \\ &= 2e^{f'(0)(t_0+k)} \sum_{n=0}^{n_k^+-1} a_{n-n_0}^k \left[A_c e^{\lambda_1(c)(n-ck+\theta)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} \right] \\ &\quad + 2e^{f'(0)(t_0+k)} \sum_{n=n_k^-+1}^{-1} a_{n-n_0}^k \left[A_c e^{\lambda_1(c)(n-ck+\theta)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} \right] \\ &= \text{IV}_1 + \text{IV}_2. \end{aligned}$$

Consider IV_1 . Then

$$\begin{aligned}
 IV_1 &= 2e^{f'(0)(t_0+k)} \sum_{n=0}^{n_k^+-1} a_{n-n_0}^k \left[A_c e^{\lambda_1(c)(n-ck+\theta)} + A_{c'} e^{\lambda_1(c')(-n-c'k+\theta')} \right] \\
 &\leq 2A_{c'} e^{-(\lambda_1(c')c' - f'(0))k} e^{f'(0)t_0 + \lambda_1(c')\theta'} \sum_{n=0}^{n_k^+} a_{n-n_0}^k \\
 &\quad + 2A_c e^{f'(0)(t_0+k)} \sum_{n=\text{int}[n_k^+/2]}^{n_k^+-1} a_{n-n_0}^k e^{\lambda_1(c)(n-ck+\theta)} \\
 &\quad + 2A_c e^{f'(0)(t_0+k)} \sum_{n=0}^{\text{int}[n_k^+/2]} a_{n-n_0}^k e^{\lambda_1(c)(n-ck+\theta)} \\
 &\leq 2A_{c'} e^{-(\lambda_1(c')c' - f'(0))k} e^{f'(0)t_0 + \lambda_1(c')\theta'} \\
 &\quad + 2A_c e^{f'(0)(t_0+k)} e^{\lambda_1(c)((n_k^+-1)-ck+\theta)} \sum_{n=\text{int}[n_k^+/2]}^{\infty} a_{n-n_0}^k \\
 &\quad + 2A_c e^{\lambda_1(c)(\text{int}[n_k^+/2]-ck+\theta)} e^{f'(0)(t_0+k)} \\
 &= IV_1^1 + IV_1^2 + IV_1^3.
 \end{aligned}$$

It is easy to see that $IV_1^1 \leq 2A_{c'} e^{-(\lambda_1(c')c' - f'(0))k} e^{f'(0)t_0 + \lambda_1(c')\theta'} \rightarrow 0$ as $k \rightarrow +\infty$, because $\lambda_1(c')c' - f'(0) > 0$. By virtue of

$$\begin{aligned}
 &e^{f'(0)(t_0+k)} e^{\lambda_1(c)((n_k^+-1)-ck+\theta)} \\
 &\leq A_c e^{f'(0)(t_0+k)} e^{\lambda_1(c) \left\{ \left(ck - \frac{f'(0)}{\lambda_1(c)}k - \theta + \frac{1}{\lambda_1(c)} \ln \frac{1}{\alpha_c} + \frac{f'(0)}{\lambda_1(c)}\theta_0 \right) - 1 - ck + \theta \right\}} \\
 &= \frac{1}{\alpha_c} e^{f'(0)(t_0+\theta_0) - \lambda_1(c)}
 \end{aligned}$$

and $\frac{n_k^+}{2} - n_0 > \frac{c\lambda_1(c) - f'(0)}{4\lambda_1(c)}k$ for sufficiently large k , we have

$$IV_1^2 = 2A_c e^{f'(0)(t_0+k)} e^{\lambda_1(c)((n_k^+-1)-ck+\theta)} \sum_{n=\text{int}[n_k^+/2]}^{\infty} a_{n-n_0}^k \rightarrow 0 \text{ as } k \rightarrow +\infty.$$

Moreover, we have

$$\begin{aligned}
 IV_1^3 &= 2A_c e^{\lambda_1(c)(\text{int}[n_k^+/2]-ck+\theta)} e^{f'(0)(t_0+k)} \\
 &\leq 2A_c e^{\lambda_1(c) \left(\frac{1}{2} \left[ck - \frac{f'(0)}{\lambda_1(c)}k - \theta + \frac{1}{\lambda_1(c)} \ln \frac{1}{\alpha_c} + \frac{f'(0)}{\lambda_1(c)}\theta_0 \right] - ck + \theta \right)} e^{f'(0)(t_0+k)} \\
 &= 2A_c e^{f'(0)t_0 + \frac{1}{2} \ln \frac{1}{\alpha_c} + \frac{1}{2} f'(0)\theta_0 + \frac{1}{2} \lambda_1(c)\theta} e^{-\frac{1}{2}[\lambda_1(c)c - f'(0)]k} \\
 &\rightarrow 0 \text{ as } k \rightarrow +\infty.
 \end{aligned}$$

We can prove that $IV_2, V, VI,$ and VII converge to zero as $k \rightarrow +\infty$ by arguments similar to those for $IV_1, III, II,$ and I , respectively.

Eventually, $|\tilde{\Phi}_{n_0}(t_0) - v_{n_0}^k(t_0)| \rightarrow 0$ as $k \rightarrow +\infty$. Since $v_{n_0}^k(t_0) \rightarrow \Phi_{n_0}(t_0)$ and $(n_0, t_0) \in \mathbb{Z} \times \mathbb{R}$ is arbitrary, we obtain that $\Phi_n(t) = \tilde{\Phi}_n(t)$ for any $(n, t) \in \mathbb{Z} \times \mathbb{R}$. The limit function being unique, the whole sequence Φ^k converges to Φ as $k \rightarrow +\infty$.

Using the same estimates as above, we can prove that the entire solution of (1.4) satisfying (1.9) is unique. \square

Acknowledgments. The authors thank the referees for their valuable comments and suggestions on the original manuscript.

REFERENCES

- [1] P.W. BATES AND A. CHMAJ, *A discrete convolution model for phase transitions*, Arch. Ration. Mech. Anal., 150 (1999), pp. 281–305.
- [2] J.W. CAHN, S.-N. CHOW, AND E.S. VAN VLECK, *Spatially discrete nonlinear diffusion equations*, Rocky Mountain J. Math., 25 (1995), pp. 87–118.
- [3] X. CHEN, S.-C. FU, AND J.-S. GUO, *Uniqueness and asymptotics of traveling waves of monostable dynamics on lattices*, SIAM J. Math. Anal., 38 (2006), pp. 233–258.
- [4] X. CHEN AND J.-S. GUO, *Existence and asymptotic stability of travelling waves of discrete quasilinear monostable equations*, J. Differential Equations, 184 (2002), pp. 549–569.
- [5] X. CHEN AND J.-S. GUO, *Uniqueness and existence of travelling waves of discrete quasilinear monostable dynamics*, Math. Ann., 326 (2003), pp. 123–146.
- [6] X. CHEN AND J.-S. GUO, *Existence and uniqueness of entire solutions for a reaction-diffusion equation*, J. Differential Equations, 212 (2005), pp. 62–84.
- [7] X. CHEN, J.-S. GUO AND H. NINOMIYA, *Entire solutions of reaction-diffusion equations with balanced bistable nonlinearities*, Proc. Roy. Soc. Edinburgh Sect. A, 136 (2006), pp. 1207–1237.
- [8] C.P. CHENG, W.T. LI, AND Z.C. WANG, *Spreading speeds and traveling waves in a delayed population model with stage structure on a two-dimensional spatial lattice*, IMA J. Appl. Math., 73 (2008), pp. 592–618.
- [9] C.P. CHENG, W.T. LI, AND Z.C. WANG, *Asymptotic Stability of Traveling Wavefronts in a Delayed Population Model with Stage Structure on a Two-dimensional Spatial Lattice*, submitted.
- [10] S.-N. CHOW, *Lattice dynamical systems*, in Dynamical Systems, Lecture Notes in Math. 1822, J.W. Macki and P. Zecca, eds., Springer, Berlin, 2003, pp. 1–102.
- [11] T. FARIA, W. HUANG, AND J. WU, *Travelling waves for delayed reaction-diffusion equations with global response*, Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci., 462A (2006), pp. 229–261.
- [12] T. FARIA AND S. TROFIMCHUK, *Nonmonotone travelling waves in a single species reaction diffusion equation with delay*, J. Differential Equations, 228 (2006), pp. 357–376.
- [13] T. FARIA AND S. TROFIMCHUK, *Positive heteroclinics and traveling waves for scalar population models with a single delay*, Appl. Math. Comput., 185 (2007), pp. 594–603.
- [14] Y. FUKAO, Y. MORITA, AND H. NINOMIYA, *Some entire solutions of the Allen-Cahn equation*, Taiwanese J. Math., 8 (2004), pp. 15–32.
- [15] R.R. GOLDBERG, *Fourier Transform*, Cambridge University Press, New York, 1965.
- [16] S.A. GOURLEY AND J. WU, *Delayed nonlocal diffusive systems in biological invasion and disease spread*, in Nonlinear Dynamics and Evolution Equations, Fields Inst. Commun. 48, American Mathematical Society, Providence, RI, 2006, pp. 137–200.
- [17] S.A. GOURLEY AND J. WU, *Extinction and periodic oscillations in an age-structured population model in a patchy environment*, J. Math. Anal. Appl., 289 (2004), pp. 431–445.
- [18] J.-S. GUO AND Y. MORITA, *Entire solutions of reaction-diffusion equations and an application to discrete diffusive equations*, Discrete Contin. Dyn. Syst., 12 (2005), pp. 193–212.
- [19] Y.-J.L. GUO, *Entire solutions for a discrete diffusive equation*, J. Math. Anal. Appl., 347 (2008), pp. 450–458.
- [20] F. HAMEL AND N. NADIRASHVILI, *Entire solutions of the KPP Equation*, Comm. Pure Appl. Math., 52 (1999), pp. 1255–1276.
- [21] F. HAMEL AND N. NADIRASHVILI, *Travelling fronts and entire solutions of the Fisher-KPP equation in \mathbb{R}^N* , Arch. Ration. Mech. Anal., 157 (2001), pp. 91–163.
- [22] Y. KYRYCHKO, S.A. GOURLEY, AND M.V. BARTUCCCELLI, *Dynamics of a stage-structured population model on an isolated finite lattice*, SIAM J. Math. Anal., 37 (2006), pp. 1688–1708.
- [23] W.T. LI, N.W. LIU, AND Z.C. WANG, *Entire solutions in reaction-advection-diffusion equations in cylinders*, J. Math. Pures Appl., 90 (2008), pp. 492–504.
- [24] W.T. LI, Z.C. WANG, AND J. WU, *Entire solutions of monostable reaction-diffusion equations with delayed nonlinearity*, J. Differential Equations, 245 (2008), pp. 102–129.
- [25] S. MA, P. WENG, AND X. ZOU, *Asymptotic speeds of propagation and traveling wavefronts in a non-local delayed lattice differential equation*, Nonlinear Anal., 65 (2006), pp. 1858–1890.
- [26] S. MA AND X. ZOU, *Propagation and its failure in a lattice delayed differential equation with global interaction*, J. Differential Equations, 212 (2005), pp. 129–190.

- [27] S. MA AND X. ZOU, *Existence, uniqueness and stability of traveling waves in a discrete reaction-diffusion monostable equation with delay*, J. Differential Equations, 217 (2005), pp. 54–87.
- [28] J. MALLET-PARET, *The global structure of traveling waves in spatially discrete dynamical systems*, J. Dynam. Differential Equations, 11 (1999), pp. 49–127.
- [29] R.H. MARTIN AND H.L. SMITH, *Abstract functional equations and reaction-diffusion systems*, Trans. Amer. Math. Soc., 321 (1990), pp. 1–44.
- [30] Y. MORITA AND H. NINOMIYA, *Entire solutions with merging fronts to reaction-diffusion equations*, J. Dynam. Differential Equations, 18 (2006), pp. 841–861.
- [31] J.D. MURRAY, *Mathematical Biology*, Springer, Berlin, 1989.
- [32] Z.X. SHI, W.T. LI, AND C.P. CHENG, *Stability and Uniqueness of Traveling Wavefronts in a Two-dimensional Lattice Differential Equation with Delay*, Appl. Math. Comput., to appear.
- [33] H.L. SMITH, *Monotone Dynamical Systems: An Introduction to the Theory of Competitive and Cooperative Systems*, Math. Surveys and Monogr. 41, American Mathematical Society, Providence, RI, 1995.
- [34] E.C. TITCHMARSH, *Introduction to the Theory of Fourier Integrals*, Clarendon Press, Oxford, 1962.
- [35] Z.C. WANG, W.T. LI, AND S. RUAN, *Travelling wave fronts of reaction-diffusion systems with spatio-temporal delays*, J. Differential Equations, 222 (2006), pp. 185–232.
- [36] Z.C. WANG, W.T. LI, AND S. RUAN, *Existence and stability of traveling wave fronts in reaction advection diffusion equations*, J. Differential Equations, 238 (2007), pp. 153–200.
- [37] Z.C. WANG, W.T. LI, AND S. RUAN, *Traveling fronts in monostable equations with nonlocal delayed effects*, J. Dynam. Differential Equations, 20 (2008), pp. 573–607.
- [38] Z.C. WANG, W.T. LI, AND S. RUAN, *Entire solutions in bistable reaction-diffusion equations with nonlocal delayed nonlinearity*, Trans. Amer. Math. Soc., 361 (2009), pp. 2047–2084.
- [39] Z.C. WANG, W.T. LI, AND S. RUAN, *Entire Solutions in Lattice Delayed Differential Equations with Global Interaction: Bistable Case*, submitted.
- [40] P. WENG, H. HUANG, AND J. WU, *Asymptotic speed of propagation of wave fronts in a lattice delay differential equation with global interaction*, IMA J. Appl. Math., 68 (2003), pp. 409–439.
- [41] P. WENG, J. WU, H. HUANG, AND J. LING, *Asymptotic speed of propagation of wave fronts in a 2D lattice delay differential equation with global interaction*, Can. Appl. Math. Q., 11 (2003), pp. 377–414.
- [42] J. WU AND X. ZOU, *Asymptotic and periodic boundary value problems of mixed FDEs and wave solutions of lattice differential equations*, J. Differential Equations, 135 (1997), pp. 315–357.
- [43] H. YAGISITA, *Backward global solutions characterizing annihilation dynamics of travelling fronts*, Publ. Res. Inst. Math. Sci., 39 (2003), pp. 117–164.
- [44] B. ZINNER, *Existence of travelling wavefront solutions for the discrete Nagumo equation*, J. Differential Equations, 96 (1992), pp. 1–27.

AN ELEMENTARY APPROACH TO A MODEL PROBLEM OF LAGERSTROM*

S. P. HASTINGS[†] AND J. B. MCLEOD[‡]

Abstract. The equation studied is $u'' + \frac{n-1}{r}u' + \varepsilon uu' + ku'^2 = 0$, with boundary conditions $u(1) = 0$, $u(\infty) = 1$. This model equation has been studied by many authors since it was introduced in the 1950s by P. A. Lagerstrom. We use an elementary approach to show that there is an infinite series solution which is uniformly convergent on $1 \leq r < \infty$. The first few terms are easily derived, from which one quickly deduces the inner and outer asymptotic expansions, with no matching procedure or a priori assumptions about the nature of the expansion. We also give a short and elementary existence and uniqueness proof which covers all $\varepsilon > 0$, $k \geq 0$, and $n \geq 1$.

Key words. matched asymptotics, singular perturbation, boundary value problem

AMS subject classifications. Primary, 34E05; Secondary, 34E10, 34E15

DOI. 10.1137/080718759

1. Introduction. The main problem is to investigate the asymptotics as $\varepsilon \rightarrow 0$ of the boundary value problem

$$(1) \quad u'' + \frac{n-1}{r}u' + \varepsilon uu' + ku'^2 = 0,$$

with

$$(2) \quad u(1) = 0, u(\infty) = 1.$$

We consider the cases $k = 0$ and $k = 1$. Our interest in these problems, originally due to Lagerstrom in the 1950s [6], [7], was stimulated by two recent papers by Popovic and Szmolyan [10], [11], who adopt a geometric approach to the problem when $k = 0$, and there are many papers which use methods of matched asymptotics or multiple scales, with varying degrees of rigor. We will review some of this work below. The point of this paper is to give a completely rigorous and relatively short answer to the problem without making any appeal either to geometric methods or to matched asymptotics. We can express the solution as an infinite series, uniformly convergent for all values of the independent variable. From this series we obtain the inner and outer asymptotic expansions with no a priori assumption about the nature of these expansions. An important and, as far as we know, original feature is that there is no “matching.”

Lagerstrom came up with these problems as models of viscous incompressible ($k = 0$) and compressible ($k = 1$) flow, so much of his work centered on $n = 2$ or 3 , but he also discussed general $n \geq 1$ [8]. The infinite series we develop can be obtained for any real number n . What n controls is the rate of convergence of the series.

*Received by the editors March 18, 2008; accepted for publication (in revised form) October 31, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/71875.html>

[†]Department of Mathematics, University of Pittsburgh, Thackeray Hall, Pittsburgh, PA 15360 (sph@pitt.edu).

[‡]Mathematics Institute, University of Oxford, 24-29 St. Giles, Oxford Ox13LB, United Kingdom (mcleod@maths.ox.ac.uk).

For $\varepsilon = k = 0$, there is an obvious distinction between $n > 2$ and $n \leq 2$. If $n > 2$, then the problem (1)–(2) has the unique solution

$$(3) \quad u = 1 - \frac{1}{r^{n-2}},$$

so that the solution with ε small is presumably some sort of perturbation of this. If $n \leq 2$, then there is no such solution. A consequence is that the convergence as $\varepsilon \rightarrow 0$ is more subtle when $n \leq 2$ than when $n > 2$. Our analysis will show that there is little prospect of discussing the behavior for small ε if $n < 2$, but fortunately we can handle all $n \geq 2$. Although it has been thought that finding the asymptotics when $k = 1$ is considerably more difficult than when $k = 0$ [3], we will show that our technique covers each case with comparable effort.

Our methods are not restricted to Lagerstrom's problems (1)–(2). In subsequent work (in preparation), we will show that there is a general method which can yield similar results for a class of singularly perturbed boundary value problems.

We start in section 2 by showing that each of these problems has one and only one solution for any $n \geq 1$ and any $\varepsilon > 0$. This is based on a simple shooting argument plus a comparison principle. These results have been obtained before, but our proof is quite short. In the subsequent sections we develop the integral equation referred to above and show how it leads with relative ease to the inner and outer expansions. These expansions go back to Lagerstrom and Kaplun, with rigorous justification of some of the features to be found in [1] or [11], for example. We find the exposition in Hinch's book [3] particularly clear (though nonrigorous), and make that work our point of comparison in checking that we get the same expansions as were found previously.

2. Existence and uniqueness. As far as we know, the first existence proof was by Hsiao [4], who considered only $n = 1$ and sufficiently small $\varepsilon > 0$. Subsequently Tam gave what seems to be the first proof valid for all $\varepsilon > 0$ and $k \geq 0$ [14]. Subsequent proofs by MacGillivray [9], Cohen, Lagerstrom, and Fokas [1], Hunter, Tajdari, and Boyer [5], each of which covers all $\varepsilon > 0$, and by several other authors, e.g., [12], [10], for restricted ranges of ε , add to the variety of techniques which have been shown to work. Uniqueness is proved in [5] (for $k = 0$) by use of a contraction mapping theorem, and in [1] by essentially a comparison method. The goal of [10] is not to give a short proof, but to illustrate the application of geometric perturbation theory to a much studied problem in matched asymptotic expansions. The proofs we give of existence and uniqueness are considerably shorter than the others we have seen.

THEOREM 1. *There exists a unique solution to the problem (1)–(2) for any $k \geq 0$, $\varepsilon > 0$, and $n \geq 1$.*

Proof. Like some others, starting with [14], we prove existence using a shooting method by considering the initial value problem

$$(4) \quad u'' + \frac{n-1}{r}u' + \varepsilon uu' + ku'^2 = 0,$$

$$(5) \quad u(1) = 0, \quad u'(1) = c,$$

for each $c > 0$. Since $u' = 0$ implies that $u'' = 0$ and u is constant, any solution to this problem is positive and increasing. As was observed in [14],

$$u'' + \varepsilon uu' \leq 0,$$

and so from (5),

$$u' + \frac{1}{2}\varepsilon u^2 \leq c.$$

In particular, since $u' \geq 0$,

$$(6) \quad u \leq \sqrt{\frac{2c}{\varepsilon}},$$

so the solution exists, and satisfies this bound, on $[1, \infty)$. Therefore, $\lim_{r \rightarrow \infty} u(r)$ exists. Writing the equation in the form

$$(7) \quad (r^{n-1}u')' + (\varepsilon u + ku')(r^{n-1}u') = 0$$

and integrating twice gives

$$(8) \quad r^{n-1}u'(r) = ce^{-ku(r) - \varepsilon \int_1^r u(s) ds},$$

$$(9) \quad u(r) = \int_1^r \frac{c}{s^{n-1}} e^{-ku - \int_1^s \varepsilon u dt} ds.$$

If $u(2) < 1$, then since u is increasing, (9) implies that

$$u(2) > p(c) = \int_1^2 \frac{c}{s^{n-1}} e^{-\varepsilon - k} ds.$$

From this and (6), we see that there are c_1 and c_2 , with $0 < c_1 < c_2$, such that if $c = c_1$, then $u(\infty) < 1$, while if $c = c_2$, then $u(\infty) > u(2) \geq 1$. Further, from (8) for any $r > R > 2$,

$$u(r) = u(R) + c \int_R^r \frac{1}{s^{n-1}} e^{-ku - \int_1^s \varepsilon u dt} ds.$$

If $n > 2$, the second term is bounded above by $\frac{c}{(n-2)R^{n-2}}$, while if $1 \leq n \leq 2$ and $R \geq 2$, it is bounded by $c \int_R^\infty e^{-\varepsilon(s-2)p(c)} ds$. Hence, this term tends to zero as $R \rightarrow \infty$, uniformly for $r \geq R$, $c_1 \leq c \leq c_2$. Since $u(R)$ is a continuous function of c , for any R , it follows that $u(\infty)$ is also continuous in c , and so there is a c with $u(\infty) = 1$, giving a solution to (1)–(2).

For uniqueness, suppose that there are two solutions of (1)–(2), say u_1 and u_2 , with $u'_1(1) > u'_2(1) > 0$. Then $u_1 > u_2$ on some maximal interval, say $(1, X)$ where $X \leq \infty$. For the same initial conditions, if $\varepsilon = k = 0$, then direct integration shows that $u_1 > u_2$ on $(1, \infty)$, and, moreover, $u_1(\infty) > u_2(\infty)$. We then raise ε and k , looking for a pair (ε_1, k_1) such that $u_1(X) = u_2(X)$ for some $X \leq \infty$, and if $0 \leq \varepsilon < \varepsilon_1$ or $0 \leq k < k_1$, no such X exists. Hence, at (ε_1, k_1) , $u_1 \geq u_2$ on $[0, \infty)$. If, at (ε_1, k_1) , $X < \infty$, then u_1 and u_2 must be tangent at X , since $u_1 - u_2$ has a minimum there, contradicting the uniqueness of initial value problems for (1). Hence, $X = \infty$, and $u_1 > u_2$ on $(1, \infty)$.

Observe from (7) that if $u'_1(r) = u'_2(r)$ for some r , then $(r^{n-1}(u'_1 - u'_2))' < 0$, since $u_1 > u_2$, so that there cannot be oscillations in $u'_1 - u'_2$. Hence, $u_1(\infty) = u_2(\infty)$ implies that there is an R with $u'_1(R) = u'_2(R)$ and $u'_1 < u'_2$ on (R, ∞) . Integrating

(7), and recalling that $u_1(\infty) = u_2(\infty)$, gives

$$r^{n-1} (u'_1 - u'_2) \Big|_R^\infty = \frac{1}{2} \varepsilon R^{n-1} (u_1^2 - u_2^2) \Big|_R + \frac{1}{2} \varepsilon (n-1) \int_R^\infty s^{n-2} (u_1^2 - u_2^2) ds - k \int_R^\infty s^{n-1} (u_1'^2 - u_2'^2) ds.$$

The left-hand side is zero, and all the terms on the right are positive, giving the necessary contradiction. \square

Remark 1. The existence theorem in [10] has one added part. It is shown there that as $\varepsilon \rightarrow 0$, the solution tends to a so-called singular solution obtained by taking a formal limit as $\varepsilon \rightarrow 0$. See [10] for details. This limit result follows from our rigorous asymptotic expansions given below.

Remark 2. There would seem to be no difficulty in extending the existence proof even to $n < 1$, but the uniqueness proof does use essentially the fact that $n \geq 1$.

3. The infinite series (with $k = 0$, $n \geq 2$). Starting again with (1), and $u(1) = 0$, we first consider the case $k = 0$ and obtain

$$(10) \quad r^{n-1} u' = B e^{-\varepsilon \int_1^r u(t) dt}$$

for some constant B . Since $u(\infty) = 1$, (10) implies that $u'(r)$ is exponentially small as $r \rightarrow \infty$. Hence we can rewrite (10) as

$$r^{n-1} u' = C e^{-\varepsilon r - \varepsilon \int_\infty^r (u-1) dt},$$

so that

$$u - 1 = C \int_\infty^r \frac{1}{t^{n-1}} e^{-\varepsilon t - \varepsilon \int_\infty^t (u-1) ds} dt.$$

Setting $\varepsilon r = \rho$, $\varepsilon t = \tau$, and $\varepsilon s = \sigma$, we obtain

$$(11) \quad u(\rho) - 1 = C \varepsilon^{n-2} \int_\infty^\rho \frac{1}{\tau^{n-1}} e^{-\tau} e^{-\int_\infty^\tau (u(\sigma)-1) d\sigma} d\tau,$$

where we use the arguments ρ and σ to indicate that we mean the rescaled version of u . Here C is a constant satisfying

$$(12) \quad -1 = C \varepsilon^{n-2} \int_\infty^\varepsilon \frac{1}{\tau^{n-1}} e^{-\tau} e^{-\int_\infty^\tau (u-1) d\sigma} d\tau.$$

Since for each ε there is a unique solution, this determines a unique C , dependent on ε .

We now consider the integral

$$(13) \quad \int_\tau^\infty (1 - u(\sigma)) d\sigma,$$

which appears in the exponent in (12). This integral has been seen to converge for each ε , but we need a bit more, namely, that it is bounded uniformly in $\varepsilon \leq \tau < \infty$ as $\varepsilon \rightarrow 0$. To see this, we note that as a function of σ , u satisfies

$$\frac{d^2 u}{d\sigma^2} + \frac{n-1}{\sigma} \frac{du}{d\sigma} + u \frac{du}{d\sigma} = 0$$

$$u = 0 \text{ when } \sigma = \varepsilon, u(\infty) = 1.$$

Denoting the unique solution by $u_\varepsilon(\sigma)$, we claim that if $0 < \varepsilon_1 < \varepsilon_2$, then $u_{\varepsilon_1} > u_{\varepsilon_2}$ for $\varepsilon_2 \leq \sigma < \infty$. If this is false, then ε_1 and ε_2 can be chosen so that $u_{\varepsilon_1}(\sigma_0) = u_{\varepsilon_2}(\sigma_0)$ for some $\sigma_0 \geq \varepsilon_2$. But then the problem

$$\frac{d^2u}{d\sigma^2} + \frac{n-1}{\sigma} \frac{du}{d\sigma} + u \frac{du}{d\sigma} = 0,$$

$$u(\sigma_0) = u_{\varepsilon_1}(\sigma_0), \quad u(\infty) = 1,$$

has two solutions, contradicting our earlier uniqueness proof.

A consequence of this is that $\int_{\varepsilon_2}^\infty (1 - u_{\varepsilon_1}(\sigma)) d\sigma < \int_{\varepsilon_2}^\infty (1 - u_{\varepsilon_2}(\sigma)) d\sigma$, which implies that the integral in the exponent in (12), including the minus sign in front, is bounded below independently of $\tau \geq \varepsilon$ and of ε . We then see that the τ -integral in (12) approaches $-\infty$ as $\varepsilon \rightarrow 0$, and hence that

$$\lim_{\varepsilon \rightarrow 0^+} C\varepsilon^{n-2} = 0.$$

Since $\int_\infty^\tau (u - 1) d\sigma > 0$, it follows from (11) that if

$$E_{n-1}(\rho) = \int_\rho^\infty \frac{1}{\tau^{n-1}} e^{-\tau} d\tau,$$

then

$$(14) \quad |u(\rho) - 1| < C\varepsilon^{n-2} E_{n-1}(\rho).$$

For purposes of future estimates, we make the obvious remark that

$$(15) \quad E_{n-1}(\rho) = \begin{cases} O(\rho^{2-n}) & \text{as } \rho \rightarrow 0 \text{ if } n > 2, \\ O(\log \rho) & \text{as } \rho \rightarrow 0 \text{ if } n = 2, \\ O(\rho^{1-n} e^{-\rho}) & \text{as } \rho \rightarrow \infty. \end{cases}$$

Hence if $n > 2$, there is a constant K such that

$$(16) \quad E_{n-1}(\rho) \leq K \min(\rho^{2-n}, \rho^{1-n} e^{-\rho}).$$

The method now is to work from (11). As observed before, since $u'(r)$ is exponentially small as $r \rightarrow \infty$, the integral term $\int_\rho^\infty (u - 1) d\sigma$ converges. Hence, for given $\varepsilon > 0$ and $\rho_0 > 0$, and any $\rho \geq \rho_0$,

$$(17) \quad u(\rho) - 1 = C\varepsilon^{n-2} \int_\infty^\rho \frac{1}{\tau^{n-1}} e^{-\tau} \left\{ 1 - \int_\infty^\tau (u - 1) d\sigma + \frac{1}{2} \left(\int_\infty^\tau (u - 1) d\sigma \right)^2 - \dots \right\} d\tau,$$

where the series in the integrand converges uniformly for $\rho_0 \leq \tau < \infty$.

In fact, we will need to use this series for all $\rho \geq \varepsilon$. Thus we need to check its convergence in this interval. This follows from (14) and (15), which imply that for any $\rho \geq \varepsilon$, if $n \geq 2$, then

$$(18) \quad \left| \int_\rho^\infty (u(s) - 1) ds \right| < C\varepsilon^{n-2} \int_\varepsilon^\infty E_{n-1}(s) ds$$

and

$$\varepsilon^{n-2} \int_{\varepsilon}^{\infty} E_{n-1}(s) ds = \begin{cases} o(1) & \text{as } \varepsilon \rightarrow 0 \text{ if } n > 2, \\ O(1) & \text{as } \varepsilon \rightarrow 0 \text{ if } n = 2. \end{cases}$$

Hence for $n > 2$ and any C , the series in the integrand of (17) converges uniformly on $[\varepsilon, \infty)$.

Now set

$$\Phi = C\varepsilon^{n-2} \int_{\varepsilon}^{\infty} E_{n-1}(s) ds.$$

We note that if $n > 2$, then $\Phi \rightarrow 0$ as $\varepsilon \rightarrow 0$, while if $n = 2$, then $\Phi \rightarrow 0$ as $C \rightarrow 0$.

We proceed to solve (17) by iteration. Thus, the first approximation is, from (16),

$$u(\rho) - 1 = C\varepsilon^{n-2} \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} d\tau + O(\Phi^2),$$

and we obtain the second approximation by substituting this back into (17). Repeating this, we reach

(19)

$$\begin{aligned} u - 1 = & -C\varepsilon^{n-2} E_{n-1} + (C\varepsilon^{n-2})^2 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left(\int_{\infty}^{\tau} E_{n-1} d\sigma \right) d\tau \\ & + \frac{1}{2} (C\varepsilon^{n-2})^3 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left(\int_{\infty}^{\tau} E_{n-1} d\sigma \right)^2 d\tau \\ & - (C\varepsilon^{n-2})^3 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \int_{\infty}^{\tau} \left\{ \int_{\infty}^{\sigma} \frac{1}{s^{n-1}} e^{-s} \left(\int_{\infty}^s E_{n-1} dt \right) ds \right\} d\sigma d\tau + O(\Phi^4), \end{aligned}$$

as $\Phi \rightarrow 0$.

To obtain C , we need to be able to evaluate each of these terms for small ρ (in particular, for $\rho = \varepsilon$), and this is a matter of integration by parts. Thus, for nonintegral n ,

$$\begin{aligned} E_{n-1}(\rho) &= \int_{\rho}^{\infty} \frac{e^{-\tau}}{\tau^{n-1}} d\tau = -\frac{\rho^{2-n}}{2-n} e^{-\rho} + \frac{1}{2-n} \int_{\rho}^{\infty} \frac{e^{-\tau}}{\tau^{n-2}} d\tau \\ (20) \qquad &= -\frac{\rho^{2-n}}{2-n} e^{-\rho} + \frac{1}{2-n} E_{n-2}, \end{aligned}$$

and this can be repeated to give E_{n-1} as a sum of terms of the form $c_k \rho^k e^{-\rho}$ and E_{n-p} , until $0 < n - p < 1$. Then

$$\begin{aligned} E_{n-p} &= \int_0^{\infty} \frac{e^{-\tau}}{\tau^{n-p}} d\tau - \int_0^{\rho} \frac{e^{-\tau}}{\tau^{n-p}} d\tau \\ &= \Gamma(p + 1 - n) - \int_0^{\rho} \frac{e^{-\tau}}{\tau^{n-p}} d\tau, \end{aligned}$$

and we can then continue to integrate by parts as far as we like. (If n is an integer, we will reach $\int_{\rho}^{\infty} \frac{e^{-\tau}}{\tau} d\tau$, which introduces a logarithm.)

Thus $E_{n-1}(\rho)$ can be expressed as a sum of terms of the form $c_k \rho^k e^{-\rho}$, and so obviously the same is true of E_{n-1}^2 , with $e^{-2\rho}$ in place of $e^{-\rho}$. Also,

$$\begin{aligned} \int_{\infty}^{\rho} E_{n-1}(\tau) d\tau &= \int_{\infty}^{\rho} \left(\int_{\tau}^{\infty} \frac{e^{-\sigma}}{\sigma^{n-1}} d\sigma \right) d\tau \\ &= \left[\tau \left(\int_{\tau}^{\infty} \frac{e^{-\sigma}}{\sigma^{n-1}} d\sigma \right) \right] \Big|_{\infty}^{\rho} + \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau^{n-2}} d\tau \\ (21) \qquad \qquad \qquad &= \rho E_{n-1} - E_{n-2}, \end{aligned}$$

so that $\int_{\infty}^{\rho} E_{n-1} d\tau$ can be expressed as the same type of sum. Hence the second term in (19) gives a sum of terms of the form $E_k(2\rho)$ and the third and fourth terms a sum involving $E_k(3\rho)$.

We now carry the process through in the most interesting cases, $n = 2, 3$.

4. The case $k = 0, n = 2$. When $n = 2$ we are interested in

$$\begin{aligned} E_1(\rho) &= \int_{\rho}^{\infty} \frac{1}{\tau} e^{-\tau} d\tau \\ &= -e^{-\rho} \log \rho + \int_{\rho}^{\infty} e^{-\tau} \log \tau d\tau \\ &= -e^{-\rho} \log \rho + \int_0^{\infty} e^{-\tau} \log \tau d\tau - \int_0^{\rho} e^{-\tau} \log \tau d\tau \\ &= -e^{-\rho} \log \rho - \gamma - \rho(\log \rho - 1) e^{-\rho} + O(\rho^2 \log \rho), \text{ for small } \rho, \\ (22) \qquad \qquad \qquad &= -\log \rho - \gamma + \rho + O(\rho^2 \log \rho). \end{aligned}$$

(See, for example, [2, Chapter 1].) Also, for future purposes, using (20) we obtain

$$(23) \qquad E_2(\rho) = \frac{e^{-\rho}}{\rho} - E_1(\rho)$$

$$(24) \qquad \qquad \qquad = \frac{1}{\rho} + \log \rho + (\gamma - 1) - \frac{1}{2}\rho + O(\rho^2 \log \rho) \text{ as } \rho \rightarrow 0.$$

Looking now at (19), with $\rho = \varepsilon$, we see that as $\varepsilon \rightarrow 0$,

$$C \log \varepsilon \rightarrow -1$$

and

$$C = \frac{1}{\log \frac{1}{\varepsilon}} + O\left(\frac{1}{(\log \frac{1}{\varepsilon})^2}\right).$$

Hence the series in (19) is in powers of $\frac{1}{\log \frac{1}{\varepsilon}}$.

Also, we will work our approximations (in order to compare the results with those of Hinch in [3]) to order $\frac{1}{\log^2(\frac{1}{\varepsilon})}$, so that (for example)

$$u = \frac{a(r)}{\log(\frac{1}{\varepsilon})} + \frac{b(r)}{\log^2(\frac{1}{\varepsilon})} + O\left(\log^{-3} \frac{1}{\varepsilon}\right)$$

for any fixed value of r (ρ of order ε). This, as we shall see, necessitates finding

$$C = \frac{1}{\log\left(\frac{1}{\varepsilon}\right)} \left\{ 1 + \frac{A}{\log\left(\frac{1}{\varepsilon}\right)} + \frac{B}{\log^2\left(\frac{1}{\varepsilon}\right)} + O\left(\log^{-3}\left(\frac{1}{\varepsilon}\right)\right) \right\}$$

and requires use of all the terms in (19).

With this in mind, we look at the second term of (19). Thus from (21),

$$(25) \quad \int_{\rho}^{\infty} E_1 d\tau = -\rho E_1 + e^{-\rho},$$

so that the second term is

$$\begin{aligned} C^2 \int_{\infty}^{\rho} \frac{1}{\tau} e^{-\tau} (\tau E_1 - e^{-\tau}) d\tau &= C^2 \left\{ \int_{\infty}^{\rho} e^{-\tau} E_1 d\tau - \int_{\infty}^{\rho} \frac{e^{-2\tau}}{\tau} d\tau \right\} \\ &= C^2 \left\{ [-e^{-\tau} E_1] \Big|_{\infty}^{\rho} - 2 \int_{\infty}^{\rho} \frac{e^{-2\tau}}{\tau} d\tau \right\} \\ (26) \quad &= C^2 (-e^{-\rho} E_1(\rho) + 2E_1(2\rho)). \end{aligned}$$

From (22), the second term is therefore

$$\begin{aligned} &C^2 (\log \rho + \gamma - 2 \log 2\rho - 2\gamma + O(\rho)) \\ (27) \quad &= C^2 (-\log \rho - \gamma - 2 \log 2 + O(\rho)) \end{aligned}$$

as $\rho \rightarrow 0$.

In the third and fourth terms of (19) we need only the leading terms; i.e., we can ignore the equivalent of $\gamma + 2 \log 2$ in (27). Using (25) the third term becomes

$$(28) \quad \frac{1}{2} C^3 \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau} (e^{-\tau} - \tau E_1)^2 d\tau = \frac{1}{2} C^3 (\log \rho + O(1)) \text{ as } \rho \rightarrow 0.$$

Finally, in the fourth term, the integrand in the τ -integral is just the second term (as a function of σ), so that from (25), the fourth term is

$$(29) \quad M = -C^3 \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau} \left[\int_{\infty}^{\tau} \{-e^{-\sigma} E_1(\sigma) + 2E_1(2\sigma)\} d\sigma \right] d\tau.$$

It is seen from (25) that for any $\tau \leq \infty$, $\int_0^{\tau} E_1(\sigma) d\sigma$ converges. Hence we can write the inner integral above in the form $\int_{\infty}^0 + \int_0^{\tau}$, and it follows that

$$M = -C^3 \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau} \{K + r(\tau)\} d\tau,$$

where K is a constant, r is bounded, and $r(\tau) = O(\tau \log \tau)$ as $\tau \rightarrow 0$. It further follows that

$$M = C^3 (KE_1(\rho) + O(1)) \text{ as } \rho \rightarrow 0.$$

We can evaluate K using (25) and (22):

$$\int_0^\infty E_1(2\sigma) d\sigma = \frac{1}{2} \int_0^\infty E_1(u) du = \frac{1}{2},$$

$$\int_0^\infty e^{-\sigma} E_1(\sigma) d\sigma = [- (e^{-\sigma} - 1) E_1]_0^\infty - \int_0^\infty (e^{-\sigma} - 1) \frac{e^{-\sigma}}{\sigma} d\sigma$$

$$(30) \quad = \lim_{\sigma \rightarrow 0} \{-E_1(2\sigma) + E_1(\sigma)\} = \lim_{\sigma \rightarrow 0} (\log 2\sigma - \log \sigma) = \log 2.$$

Hence, from (29), the fourth term of (19) is

$$(31) \quad C^3 \{E_1(\rho) (\log 2 - 1) + O(1)\} = -C^3 \{(\log 2 - 1) \log \rho + O(1)\} \text{ as } \rho \rightarrow 0.$$

Now setting $\rho = \varepsilon$ and using (27), (28), and (31), we obtain that

$$-1 = -C(-\log \varepsilon - \gamma + O(\varepsilon)) + C^2(-\log \varepsilon - \gamma - 2 \log 2 + O(\varepsilon))$$

$$+ \frac{1}{2} C^3 (\log \varepsilon + O(1)) - C^3 \{(\log 2 - 1) \log \varepsilon + O(1)\}$$

as $\varepsilon \rightarrow 0$. Hence,

$$(32) \quad \frac{1}{\log(\frac{1}{\varepsilon})} = C \left(1 - \frac{\gamma}{\log(\frac{1}{\varepsilon})}\right) - C^2 \left(1 - \frac{\gamma + 2 \log 2}{\log(\frac{1}{\varepsilon})}\right)$$

$$+ C^3 \left(\frac{3}{2} - \log 2\right) + O\left(\log^{-4}\left(\frac{1}{\varepsilon}\right)\right),$$

and

$$C = \frac{1}{\log(\frac{1}{\varepsilon})} + \frac{A}{\log^2(\frac{1}{\varepsilon})} + \frac{B}{\log^3(\frac{1}{\varepsilon})} + O\left(\frac{1}{\log^4(\frac{1}{\varepsilon})}\right),$$

where

$$-\gamma + A - 1 = 0,$$

$$B - \gamma A - 2A + (\gamma + 2 \log 2) + \frac{3}{2} - \log 2 = 0.$$

Hence,

$$A = \gamma + 1,$$

$$B = \gamma^2 + 2\gamma + \frac{1}{2} - \log 2.$$

Thus, for fixed r, ρ of order ε , we have, with $\lambda = \log(\frac{1}{\varepsilon})$,

$$u - 1 = \left(\frac{1}{\lambda} + \frac{\gamma + 1}{\lambda^2} + \frac{(\gamma + 1)^2 - \frac{1}{2} - \log 2}{\lambda^3}\right) (\log r + \log \varepsilon + \gamma)$$

$$+ \frac{1}{\lambda^2} \left(1 + \frac{2(\gamma + 1)}{\lambda}\right) (-\log r - \log \varepsilon - \gamma - 2 \log 2)$$

$$+ \frac{1}{\lambda^3} \left(\frac{3}{2} - \log 2\right) (\log r + \log \varepsilon) + O(\lambda^{-4}),$$

so that, after cancellation,

$$u = \frac{\log r}{\lambda} + \frac{\gamma \log r}{\lambda^2} + O(\lambda^{-3}).$$

This is the “inner expansion.” For the “outer expansion,” i.e., fixed ρ , r of order $\frac{1}{\varepsilon}$, we use (19), truncated to second order, to get

$$u - 1 = -E_1(\rho) \left(\frac{1}{\lambda} + \frac{\gamma + 1}{\lambda^2} \right) + \frac{1}{\lambda^2} (2E_1(2\rho) - e^{-\rho} E_1(\rho)) + O(\lambda^{-3}).$$

These results are in accordance with those of Hinch and of others on this problem.

5. The case $k = 0, n = 3$. Here we are interested in (from (22) and (23))

$$E_2(\rho) = \frac{e^{-\rho}}{\rho} - E_1(\rho) = \frac{1}{\rho} + \log \rho + (\gamma - 1) - \frac{1}{2}\rho + O(\rho^2 \log \rho) \text{ as } \rho \rightarrow 0.$$

Thus, the first term on the right of (19) evaluated at $\rho = \varepsilon$ is

$$-C(1 + \varepsilon \log \varepsilon + (\gamma - 1) + O(\varepsilon^2)) \text{ as } \varepsilon \rightarrow 0.$$

The second term is

$$\begin{aligned} & (C\varepsilon)^2 \int_{\infty}^{\rho} \frac{1}{\tau^2} e^{-\tau} \left(\int_{\infty}^{\tau} E_2 d\sigma \right) d\tau \\ &= (C\varepsilon)^2 \left\{ \left[-E_2(\tau) \int_{\infty}^{\tau} E_2(\sigma) d\sigma \right]_{\infty}^{\rho} + \int_{\infty}^{\rho} E_2^2 d\tau \right\} \\ &= (C\varepsilon)^2 \left\{ -E_2(\rho) \int_{\infty}^{\rho} E_2(\tau) d\tau + \int_{\infty}^{\rho} E_2^2 d\tau \right\}. \end{aligned}$$

From (23) we see that

$$\int_{\infty}^{\rho} E_2^2 d\tau = -\frac{1}{\rho} + \log^2 \rho + O(\log \rho) \text{ as } \rho \rightarrow 0,$$

while from (21),

$$\begin{aligned} \int_{\infty}^{\rho} E_2 d\tau &= \rho E_2 - E_1 = 1 + \log \rho + \gamma + O(\rho \log \rho), \\ E_2 \int_{\infty}^{\rho} E_2 d\tau &= \frac{1}{\rho} \log \rho + \frac{\gamma + 1}{\rho} + O(\log^2 \rho). \end{aligned}$$

In all, the second term is

$$(C\varepsilon)^2 \left\{ -\frac{1}{\rho} \log \rho - \frac{\gamma + 2}{\rho} + O(\log^2 \rho) \right\}.$$

It is readily verified that the third and fourth terms in (19) give $O\{C^3 \varepsilon^3 (\frac{1}{\rho} \log^2 \rho)\}$, which is negligible. Thus, evaluating (19) at $\rho = \varepsilon$, we have

$$-1 = -C\varepsilon \left(\frac{1}{\varepsilon} + \log \varepsilon + \gamma - 1 \right) + (C\varepsilon)^2 \left(-\frac{1}{\varepsilon} \log \varepsilon - \frac{\gamma + 2}{\varepsilon} \right) + O(C^3 \varepsilon^2 \log^2 \varepsilon),$$

so that

$$C = 1 - 2\varepsilon \log \varepsilon - \varepsilon(2\gamma + 1) + O(\varepsilon^2 \log^2 \varepsilon).$$

Then, for fixed r, ρ of order ε , we have

$$\begin{aligned} u - 1 &= -\varepsilon(1 - 2\varepsilon \log \varepsilon - \varepsilon(2\gamma + 1)) \left(\frac{1}{\varepsilon r} + \log \varepsilon + \log r + \gamma - 1 \right) \\ &\quad + \varepsilon^2 \left(-\frac{1}{\varepsilon r} (\log \varepsilon + \log r) - \frac{\gamma + 2}{\varepsilon r} \right) + O(\varepsilon^2 \log^2 \varepsilon), \\ u &= 1 - \frac{1}{r} - \varepsilon \log \varepsilon \left(1 - \frac{1}{r} \right) - \varepsilon \left(\log r + \frac{\log r}{r} \right) + \varepsilon(1 - \gamma) \left(1 - \frac{1}{r} \right) \\ &\quad + O(\varepsilon^2 \log^2 \varepsilon). \end{aligned}$$

For fixed ρ, r of order ε^{-1} , we again use (19), to give

$$(33) \quad \begin{aligned} u - 1 &= -\varepsilon(1 - 2\varepsilon \log \varepsilon - \varepsilon(2\gamma + 1)) E_2(\rho) \\ &\quad + \varepsilon^2 \left\{ E_1(\rho) E_2(\rho) - \rho E_2^2(\rho) - \int_{\rho}^{\infty} E_2^2 d\tau \right\} + O(\varepsilon^3). \end{aligned}$$

Again, these results are in agreement with those of Hinch, and others, although (33) gives one term further.

Remark 3. It is of interest to consider what happens when $n < 2$, since, at least for $n \geq 1$, there still exists a unique solution. Equation (19) is still valid at $\rho = \varepsilon$, but since $E_{n-1}(\rho)$ is no longer singular at $\rho = 0$ for $n < 2$, (19) with $\rho = \varepsilon$ becomes merely an implicit equation for $C\varepsilon^{n-2}$. This tells us that $C \rightarrow 0$, since $\varepsilon^{n-2} \rightarrow \infty$, but we no longer get an asymptotic expansion. In particular, it is no longer obvious that C is unique. Of course, we know this from Theorem 1 if $n \geq 1$.

6. The case $k = 1$. We can in fact treat a generalization, which causes no further difficulties,

$$(34) \quad u'' + \frac{n-1}{r} u' + f(u) u'^2 + \varepsilon u u' = 0,$$

with the same boundary conditions. As before, we will compare our results with those of Hinch in [3].

As remarked in the proof of Theorem 1, the solution will necessarily have $u' > 0$ so that conditions on $f(u)$ are necessary only for $0 \leq u \leq 1$. We require only that f be continuous and positive in this interval.

Then (34) can be written as

$$\frac{(r^{n-1} u')'}{r^{n-1} u'} + f(u) u' + \varepsilon u = 0,$$

so that

$$\log(r^{n-1} u') = -F(u) - \varepsilon \int_1^r u dt + A$$

for some constant A , where

$$F(u) = \int_0^u f(s) ds.$$

This becomes

$$e^{F(u)} u' = \frac{C}{r^{n-1}} e^{-\varepsilon r - \varepsilon \int_\infty^r (u-1) ds},$$

or, on integration,

$$G(u) - G(1) = C \int_\infty^r \frac{1}{t^{n-1}} e^{-\varepsilon t - \varepsilon \int_\infty^t (u-1) ds} dt,$$

where

$$G(u) = \int_0^u e^{F(v)} dv.$$

In order to keep the manipulations simple and effect comparisons, we will consider from here the Lagerstrom model, where $f(u) = 1$, $F(u) = u$, $G(u) = e^u - 1$. Then, with $\varepsilon r = \rho$, $\varepsilon t = \tau$, we have

$$(35) \quad e^u - e = C\varepsilon^{n-2} \int_\infty^\rho \frac{1}{\tau^{n-1}} e^{-\tau} e^{-\int_\infty^\tau (u-1) d\sigma} d\tau,$$

and writing

$$u - 1 = \frac{u - 1}{e^u - e} (e^u - e),$$

we get

$$(36) \quad e^u - e = C\varepsilon^{n-2} \int_\infty^\rho \frac{1}{\tau^{n-1}} e^{-\tau} e^{-\int_\infty^\tau \frac{u-1}{e^u - e} (e^u - e) d\sigma} d\tau.$$

As in section 3, we can integrate by parts, and since $0 \leq \frac{u-1}{e^u - e} \leq 1$ in $0 \leq u < 1$, we will develop a convergent series as before. To get the first three terms (necessary to give Hinch's accuracy when $n = 2$), we have from (36) that

$$(37) \quad e^u - e = C\varepsilon^{n-2} \int_\infty^\rho \frac{1}{\tau^{n-1}} e^{-\tau} \left\{ 1 - \int_\infty^\tau \frac{u-1}{e^u - e} (e^u - e) d\sigma \right. \\ \left. + \frac{1}{2} \left(\int_\infty^\tau \frac{u-1}{e^u - e} (e^u - e) d\sigma \right)^2 + \dots \right\} d\tau.$$

As before, since $e^u - e \rightarrow 0$ exponentially fast as $\rho \rightarrow \infty$, the series in the integrand converges uniformly for large τ , so that (37) is valid for large ρ . But again we need to extend it down to $\rho = \varepsilon$. From (35) we have

$$e^u - e \leq C\varepsilon^{n-2} E_{n-1}(\rho),$$

and so the convergence proof is the same as that preceding (19).

Before proceeding further with $n = 2$, we make a couple of remarks about the simpler case $n > 2$. Then, as we saw in section 5, only two terms are necessary to give the required accuracy, and then (37) gives

$$e^u - u = C\varepsilon^{n-2} \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau^{n-1}} \left\{ 1 - \int_{\infty}^{\tau} \frac{u-1}{e^u - e} (e^u - e) d\sigma + \dots \right\},$$

and since $\frac{u-1}{e^u - e}$ appears in what is already the highest order term, we can replace it by its limit as $u \rightarrow 1$, i.e., $\frac{1}{e}$. Thus we get, to the required order,

$$e^u - e = -C\varepsilon^{n-2} E_{n-1} - \left(\frac{C\varepsilon^{n-2}}{e} \right) \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau^{n-1}} \int_{\infty}^{\tau} (e^u - e) d\sigma d\tau.$$

(We will proceed more carefully for $n = 2$.) This, apart from the factor $\frac{1}{e}$, is the same equation dealt with in section 5 (with $e^u - e$ in place of $u - 1$), and the solution can be written down from there. If we had a general function f in place of 1, we would get

$$e^{F(u)} - e^{F(1)} = -C\varepsilon^{n-2} E_{n-1} - \frac{C\varepsilon^{n-2}}{e^{F(1)} f(1)} \int_{\infty}^{\rho} \frac{e^{-\tau}}{\tau^{n-1}} \left(\int_{\infty}^{\tau} (e^{F(u)} - e^{F(1)}) d\sigma \right) d\tau.$$

Turning now to the case $n = 2$ and $F(u) = u$, we need three terms on the right of (37). Thus,

$$(38) \quad \frac{u-1}{e^u - e} = \frac{1}{e} - \frac{1}{2e^2} (e^u - e) + O(e^u - e)^2 \text{ as } u \rightarrow 1.$$

We follow the method used just before (19) and obtain from (37) that

$$\begin{aligned} e^u - e &= -C\varepsilon^{n-2} E_{n-1} + \frac{1}{e} (C\varepsilon^{n-2})^2 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left(\int_{\infty}^{\tau} E_{n-1} d\sigma \right) d\tau \\ &+ \frac{1}{2e^2} (C\varepsilon^{n-2})^3 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left(\int_{\infty}^{\tau} E_{n-1}^2 d\sigma \right) d\tau \\ &+ \frac{1}{2e^2} (C\varepsilon^{n-2})^3 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left(\int_{\infty}^{\tau} E_{n-1} d\sigma \right)^2 d\tau \\ &- \frac{1}{e^2} (C\varepsilon^{n-2})^3 \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \int_{\infty}^{\tau} \left\{ \int_{\infty}^{\sigma} \frac{1}{s^{n-1}} e^{-s} \left(\int_{\infty}^s E_{n-1} dt \right) ds \right\} d\sigma d\tau + O(\Phi^4) \\ &= -C\varepsilon^{n-2} E_{n-1} + F_1 + F_2 + F_3 + F_4 + O(\Phi^4), \end{aligned}$$

say. As before, if $n > 2$, then this is valid for any C as $\varepsilon \rightarrow 0$, uniformly in $\rho \geq \varepsilon$, while if $n = 2$, it is valid as $C \rightarrow 0$.

For $n = 2$ we can continue to follow the argument in section 4. Thus, as $\rho \rightarrow 0$,

$$\begin{aligned} F_1 &= \frac{1}{e} C^2 (-\log \rho - \gamma - 2 \log 2 + O(\rho)), \\ F_3 &= \frac{1}{2e^2} C^3 (\log \rho + O(1)), \\ F_4 &= -\frac{1}{e^2} C^3 [(\log 2 - 1) \log \rho + O(1)]. \end{aligned}$$

The term F_2 did not appear before. Only the highest order term is needed for our expansion and this is

$$-\frac{1}{2e^2}C^3 \left(\int_{\infty}^{\rho} \frac{1}{\tau} e^{-\tau} d\tau \right) \int_0^{\infty} E_1^2 d\sigma.$$

Now

$$\begin{aligned} \int_0^{\infty} E_1^2 d\sigma &= [\tau E_1^2]_0^{\infty} + 2 \int_0^{\infty} \tau \frac{e^{-\tau}}{\tau} E_1 d\tau \\ &= 2 \int_0^{\infty} e^{-\tau} E_1 d\tau = 2 \log 2 \text{ from (30)}. \end{aligned}$$

Thus,

$$F_2 = \frac{1}{e^2} C^3 E_1 (\log 2 + O(\rho \log^2 \rho)) = -\frac{1}{e^2} C^3 ((\log 2) \log \rho + O(1)),$$

and, evaluating at $\rho = \varepsilon$, we have

$$\begin{aligned} 1 - e &= C (\log \varepsilon + \gamma + O(\varepsilon)) \\ &\quad - \frac{1}{e} C^2 (\log \varepsilon + \gamma + 2 \log 2 + O(\varepsilon)) \\ &\quad + \frac{1}{2e^2} C^3 \log \varepsilon (-2 \log 2 + 1 - 2 \log 2 + 2) + O(C^3), \\ \frac{e-1}{\log \frac{1}{\varepsilon}} &= C \left(1 - \frac{\gamma}{\log(\frac{1}{\varepsilon})} \right) - \frac{1}{e} C^2 \left(1 - \frac{\gamma + 2 \log 2}{\log(\frac{1}{\varepsilon})} \right) \\ &\quad + \frac{1}{2e^2} C^3 \left(3 - 4 \log 2 + O\left(\frac{C\varepsilon}{\log \varepsilon}\right) + O\left(\frac{C^2\varepsilon}{\log \varepsilon}\right) + O\left(\frac{C^3}{\log \varepsilon}\right) \right). \end{aligned}$$

Hence if

$$C = \frac{e-1}{\log(\frac{1}{\varepsilon})} + \frac{A}{\log^2(\frac{1}{\varepsilon})} + \frac{B}{\log^3(\frac{1}{\varepsilon})} + O\left(\log^{-4}\left(\frac{1}{\varepsilon}\right)\right),$$

then

$$-\gamma(e-1) + A - \frac{(e-1)^2}{e} = 0,$$

$$A = \frac{e-1}{e} (\gamma e + e - 1),$$

$$B - A\gamma + \frac{(e-1)^2}{e} (\gamma + 2 \log 2) - \frac{2A(e-1)}{e} + \frac{1}{2e^2} (e-1)^3 (3 - 4 \log 2) = 0.$$

We can of course calculate B , but in fact its value will be irrelevant to the level of approximation that we take.

Then, for fixed r (ρ of order ε), we have, with $l = \log(\frac{1}{\varepsilon})$,

$$\begin{aligned}
 e^u - e &= (e - 1) \left\{ \frac{1}{l} + \frac{\gamma + 1 - \frac{1}{e}}{l^2} + \frac{B/(e - 1)}{l^3} \right\} (\log \varepsilon + \log r + \gamma) \\
 &\quad + \frac{1}{e} (e - 1)^2 \left\{ \frac{1}{l^2} + \frac{2(\gamma + 1 - \frac{1}{e})}{l^3} \right\} (-\log \varepsilon - \log r - \gamma - 2 \log 2) \\
 &\quad + \frac{1}{2e^2} \frac{(e - 1)^3}{l^3} (3 - 4 \log 2) (\log \varepsilon + \log r) + O(l^{-3}) \\
 (39) \quad &= 1 - e + \frac{(e - 1) \log r}{l} + \frac{\gamma(e - 1) \log r}{l^2} + O(l^{-3}).
 \end{aligned}$$

(Note that the definitions of A and B were such that $u = 0$ at $r = 1$ up to and including order l^{-2} , so that to that order there can be only terms in $\log r$, not constant terms. We do not need the explicit value of B .) To obtain u , we have to invert, so that

$$u = \log \left\{ 1 + \frac{e - 1}{l} \log r + \frac{\gamma(e - 1)}{l^2} \log r + O(l^{-3}) \right\}.$$

For fixed ρ, r of order ε^{-1} , we have

$$e^u - e = -\frac{e - 1}{l} \left(1 + \frac{\gamma + 1 - \frac{1}{e}}{l} \right) E_1(\rho) + \frac{(e - 1)^2}{el^2} (2E_1(2\rho) - e^{-\rho} E_1(\rho)) + O(l^{-3}).$$

Thus

$$\begin{aligned}
 u - 1 &= \frac{1}{e} (e^u - e) - \frac{1}{2e^2} (e^u - e)^2 + \dots \\
 &= -\frac{e - 1}{e} \left(1 + \frac{\gamma + 1 - \frac{1}{e}}{l} \right) \frac{E_1(\rho)}{l} + \frac{(e - 1)^2}{e^2} \frac{(2E_1(2\rho) - e^{-\rho} E_1(\rho))}{l^2} \\
 (40) \quad &- \frac{(e - 1)^2}{2e^2} \frac{E_1^2(\rho)}{l^2} + O(l^{-3}).
 \end{aligned}$$

Again, these results are consistent with those of Hinch and others, except that Hinch has an algebraic mistake which in (40) replaces $\gamma + 1 - \frac{1}{e}$ by $\gamma - 1 + \frac{1}{e}$.

7. Final remarks. Starting with Lagerstrom, the terms involving $\log \varepsilon$ in the inner expansions have been considered difficult to explain. They are often called “switchback” terms, because there is nothing obvious in the equation which indicates the need for such terms, and because, starting with an expansion in powers of ε , one finds inconsistent results which are resolved only by adding terms of lower order, that is, powers of $\varepsilon \log \varepsilon$. The recent approach to the problem by geometric perturbation theory explains this by reference to a “resonance phenomenon,” which is too complicated for us to describe here [10], [11].

In our work, the necessity for such terms is seen already from (11) and the resulting expansion (17):

$$u(\rho) - 1 = C\varepsilon^{n-2} \int_{\infty}^{\rho} \frac{1}{\tau^{n-1}} e^{-\tau} \left\{ 1 - \int_{\infty}^{\tau} (u - 1) d\sigma + \frac{1}{2} \left(\int_{\infty}^{\tau} (u - 1) d\sigma \right)^2 - \dots \right\} d\tau.$$

In the existence proof it was seen in (9) that $C = O(1)$ as $\varepsilon \rightarrow 0$. On the right of (17) the first term is simply $-C\varepsilon^{n-2}E_{n-1}(\rho)$, and the simple expansions given for E_1 and E_2 show immediately the need for the logarithmic terms. There is no “switchback,” because the procedure does not start with any assumption about the nature of the expansion, and there is no need for a “matching.”

A number of authors have noted that the outer expansion is a uniformly valid asymptotic expansion on $[1, \infty)$, and therefore it “contains” the inner expansion [5], though this is more subtle when $k = 1$ [13]. Our twist on this is that both expansions are contained in the uniformly convergent series defined implicitly by (17). The simple derivation of this series via the integral equation (11) is new, as far as we know.

Acknowledgment. We thank the referees for some very helpful comments. In particular, they called our attention to earlier proofs of the existence and uniqueness results, in some cases by techniques similar to ours.

REFERENCES

- [1] D. S. COHEN, A. FOKAS, AND P. A. LAGERSTROM, *Proof of some asymptotic results for a model equation for low Reynolds number flow*, SIAM J. Appl. Math., 35 (1978), pp. 187–207.
- [2] A. ERDÉLYI, W. MAGNUS, F. OBERHETTINGER, AND F. G. TRICOMI, *Higher Transcendental Functions*, Vol. 1, McGraw-Hill, New York, 1953.
- [3] E. J. HINCH, *Perturbation Methods*, Cambridge University Press, Cambridge, UK, 1991.
- [4] G. C. HSIAO, *Singular perturbations for a nonlinear differential equation with a small parameter*, SIAM J. Math. Anal., 4 (1973), pp. 283–301.
- [5] C. HUNTER, M. TAJDARI, AND S. D. BOYER, *On Lagerstrom’s model of slow incompressible viscous flow*, SIAM J. Appl. Math., 50 (1990), pp. 48–63.
- [6] S. KAPLUN AND P. A. LAGERSTROM, *Asymptotic expansions of Navier-Stokes solutions for small Reynolds number*, J. Math. Mech., 6 (1957), pp. 585–593.
- [7] P. A. LAGERSTROM AND R. G. CASTEN, *Basic concepts underlying singular perturbation techniques*, SIAM Rev., 14 (1972), pp. 63–120.
- [8] P. A. LAGERSTROM AND C. A. REINELT, *Note on logarithmic switchback terms in regular and singular perturbation expansions*, SIAM J. Appl. Math. 44 (1984), pp. 451–462.
- [9] A. D. MACGILLIVRAY, *On a model equation of Lagerstrom*, SIAM J. Appl. Math., 34 (1978), pp. 804–812.
- [10] N. POPOVIC AND P. SZMOLYAN, *A geometric analysis of the Lagerstrom model problem*, J. Differential Equations, 199 (2004), pp. 290–325.
- [11] N. POPOVIC AND P. SZMOLYAN, *Rigorous asymptotic expansions for Lagerstrom’s model equations, a geometric approach*, Nonlinear Anal., 59 (2004), pp. 531–565.
- [12] S. ROSENBLAT AND J. SHEPHERD, *On the asymptotic solution of the Lagerstrom model equation*, SIAM J. Appl. Math., 29 (1975), pp. 110–120.
- [13] L. A. SKINNER, *Note on the Lagerstrom singular perturbation models*, SIAM J. Appl. Math., 41 (1981), pp. 362–364.
- [14] K. TAM, *On the Lagerstrom model for flow at low Reynolds numbers*, J. Math. Anal. Appl., 49 (1975), pp. 286–294.

VECTOR FIELDS FOR MEAN VALUE COORDINATES*

S. L. LEE†

Abstract. We find vector fields \mathbf{F} that provide the representation $\mathbf{v} = \int_{\partial D} \mathbf{r} \mathbf{F} \cdot d\mathbf{r} / \int_{\partial D} \mathbf{F} \cdot d\mathbf{r}$ for any compact 2-D manifold $D \subset \mathbb{R}^2$ with piecewise smooth boundary ∂D and $\mathbf{v} \in \mathbb{R}^2$, and the representation $\mathbf{v} = \iint_{\partial M} \mathbf{r} \mathbf{F} \cdot d\mathbf{S} / \iint_{\partial M} \mathbf{F} \cdot d\mathbf{S}$ for any compact 3-D manifold $M \subset \mathbb{R}^3$ and $\mathbf{v} \in \mathbb{R}^3$. Our method exploits properties of conservative fields in \mathbb{R}^2 and divergence free vector fields in \mathbb{R}^3 . Discrete versions, which are more general than Floater’s mean value coordinates, are derived from the above representations with a special choice of \mathbf{F} , either by taking points on the boundaries of $D \subset \mathbb{R}^2$ and $M \subset \mathbb{R}^3$ or by considering representations on boundaries of polygons in \mathbb{R}^2 or polyhedra in \mathbb{R}^3 .

Key words. star-shaped region, mean value coordinates, conservative and divergence free vector fields, Stoke’s theorem, divergence theorem, spherical harmonics

AMS subject classifications. 41A05, 52A30, 52B10, 52B70

DOI. 10.1137/070694144

1. Introduction. Representation of points and functionals on a set by its extreme points or boundary is an important problem in mathematics and its applications. Barycentric coordinates, the Krein–Millman theorem, and Choquet’s theorem are examples of such a representation. Recently, in conjunction with the construction of one-one transformations and parametrizations of meshes in \mathbb{R}^3 , Floater [2] has found an explicit formula for the representation of points that lie in the kernel of star-shaped polygons in \mathbb{R}^2 in terms of extreme points of the polygon. Because of its usefulness in computer graphics and its potential applications in functional approximation and interpolation, the idea has been quickly extended to star-shaped polyhedrons in \mathbb{R}^3 (see [3], [4], [5], [6], [7]). Floater’s original idea was motivated by the mean value property of harmonic functions; i.e., if ϕ is harmonic in a region $\Omega \subset \mathbb{R}^2$, then for any disc $D(\mathbf{v}, r) \subset \Omega$ with center at \mathbf{v} and radius r ,

$$(1.1) \quad \phi(\mathbf{v}) = \frac{1}{2\pi r} \int_{\partial D} \phi(\mathbf{r}) ds,$$

where $ds = \left| \frac{d\mathbf{r}}{dt} \right| dt$ for a parametrization $\mathbf{r} = \mathbf{r}(t)$ of ∂D . He observed that a class of real-valued piecewise linear functions defined on a triangular mesh in \mathbb{R}^2 , which he calls *convex combination functions*, shares discretely some properties of harmonic functions. Forcing the mean value property (1.1) on the convex combination functions produces the new coordinates, which he calls the *mean value coordinates*. The problem of computing these coordinates then reduces to integrating elementary trigonometric functions over the unit circle or the unit sphere (see [2], [3], [7]).

This paper explores the connection between mean value coordinates with conservative vector fields on compact 2-D manifolds in \mathbb{R}^2 and divergence free vector fields on compact 3-D manifolds in \mathbb{R}^3 that provide a mathematical foundation, and puts

*Received by the editors June 9, 2007; accepted for publication (in revised form) November 8, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/69414.html>

†Department of Mathematics, National University of Singapore, 2 Science Drive 2, Singapore 117543 (matleesl@nus.edu.sg).

mean value coordinates in a more general mathematical setting. In particular we find vector fields \mathbf{F} that provide the representation

$$\mathbf{v} = \frac{\int_{\partial D} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}{\int_{\partial D} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}$$

for any compact 2-D manifold $D \subset \mathbb{R}^2$ with piecewise smooth boundary ∂D and for any vector $\mathbf{v} \in \mathbb{R}^2$, and

$$\mathbf{v} = \frac{\iint_{\partial M} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}{\iint_{\partial M} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}$$

for any compact 3-D manifold $M \subset \mathbb{R}^3$ and $\mathbf{v} \in \mathbb{R}^3$. We shall refer to these as *mean value representations*. Our method exploits properties of conservative and divergence free vector fields. With a special choice of \mathbf{F} and restricting to points on the boundaries of the manifolds or by considering the representations on the boundaries of polygons in \mathbb{R}^2 and polyhedra in \mathbb{R}^3 , discrete versions of mean value representation are obtained, which are more general than Floater's mean value coordinates.

In section 2 a search for vector fields that provide the mean value representation in \mathbb{R}^2 leads to a class of vector fields

$$\mathbf{F} = \left[\frac{1}{xy} \Phi \left(\frac{x}{r}, \frac{y}{r} \right), \frac{-1}{y^2} \Phi \left(\frac{x}{r}, \frac{y}{r} \right) \right]^T,$$

where Φ is an arbitrary real-valued differentiable function on the unit circle \mathbb{S}^1 . For these \mathbf{F} , the vector fields $x\mathbf{F}$, $y\mathbf{F}$ are conservative on $\mathbb{R}^2 \setminus \{\mathbf{0}\}$. Further, if

$$\Phi(\cos \theta, \sin \theta) = -\cos \theta \sin^2 \theta \left(1 + \sum_{k=2}^{\infty} a_k \cos k\theta + b_k \sin k\theta \right),$$

with $|\sum_{k=2}^{\infty} a_k \cos k\theta + b_k \sin k\theta| \leq 1$, then the corresponding \mathbf{F} provides a mean value representation. A more general form of Floater's mean value coordinates is derived from the particular choice of $\Phi(\cos \theta, \sin \theta) = -\cos \theta \sin^2 \theta$, i.e., $\Phi\left(\frac{x}{r}, \frac{y}{r}\right) = -\frac{xy^2}{r^3}$, and hence $\mathbf{F}(x, y) := [-y/r^3, x/r^3]^T$. The derivation exploits properties of path independent integrals of conservative fields. A similar search in section 3 for vector fields that provide the mean value representations in \mathbb{R}^3 leads to a class

$$\mathbf{F}(\mathbf{r}) = \left[\frac{1}{xyz} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right), \frac{1}{x^2z} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right), \frac{1}{x^2y} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right) \right]^T,$$

where $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}$ is an arbitrary differentiable function and $\rho := (x^2 + y^2 + z^2)^{1/2}$. In this case $x\mathbf{F}$, $y\mathbf{F}$, and $z\mathbf{F}$ are divergence free on $\mathbb{R}^3 \setminus \{\mathbf{0}\}$. Further, if

$$\begin{aligned} \Phi \left(\frac{x}{\rho}, \frac{x}{\rho}, \frac{x}{\rho} \right) &= \frac{x^2yz}{\rho^4} \left(1 + \Psi \left(\frac{x}{\rho}, \frac{x}{\rho}, \frac{x}{\rho} \right) \right) \\ &= \sin^3 \phi \cos^2 \theta \sin \theta \cos \phi (1 + \Psi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)), \end{aligned}$$

where Ψ is orthogonal to the first order spherical harmonics and $|\Psi| \leq 1$, then \mathbf{F} provides a mean value representation in \mathbb{R}^3 . The derivation of Floater's type mean value coordinates in \mathbb{R}^3 is carried out in section 4. Here we take $\Psi = 0$, and hence

$\Phi(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho}) = \frac{x^2yz}{\rho^4}$ and $\mathbf{F} = \frac{\mathbf{r}}{\rho^4}$, and consider compact 3-D manifolds with polyhedral faces of any topology. The integral of $\iint_{\partial M} \mathbf{rF} \cdot d\mathbf{S}$ over the boundary of such a manifold M is the sum of the integrals over its faces. The integral $\iint_S \mathbf{rF} \cdot d\mathbf{S}$ over each face S is evaluated using Stoke's theorem. This is done by finding vector fields $\mathbf{A}, \mathbf{B}, \mathbf{C}$ such that $x\mathbf{F} = \text{curl}(\mathbf{A})$, $y\mathbf{F} = \text{curl}(\mathbf{B})$, $z\mathbf{F} = \text{curl}(\mathbf{C})$, which is possible since $\text{div}(x\mathbf{F}) = \text{div}(y\mathbf{F}) = \text{div}(z\mathbf{F}) = 0$. This leads to a simple formula for the integrals over the faces of the boundary of M , from which discrete versions of the mean value representation or Floater's type mean value coordinates are obtained. If all the faces are triangular, the procedure leads to a unique solution. In all other cases the solutions are not unique.

2. Representation on boundaries of compact 2-D manifolds in \mathbb{R}^2 . Let $D \subset \mathbb{R}^2$ be a compact 2-D manifold with piecewise smooth boundary ∂D , which is assumed to be oriented in the positive direction; i.e., if one travels along this direction, the region D lies on his left. We want to find a vector field $\mathbf{F} = [F_1(x, y), F_2(x, y)]^T$ such that for any $D \subset \mathbb{R}^2$ and $\mathbf{v} \in \mathbb{R}^2$,

$$(2.1) \quad \int_{\partial D} (\mathbf{r} - \mathbf{v})\mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} = 0$$

and

$$(2.2) \quad \int_{\partial D} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} \neq 0,$$

so that

$$(2.3) \quad \mathbf{v} = \frac{\int_{\partial D} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}{\int_{\partial D} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}.$$

Without loss of generality, we may assume that $\mathbf{v} = \mathbf{0}$ so that the problem reduces to finding \mathbf{F} such that for any D

$$(2.4) \quad \int_{\partial D} \mathbf{r} \mathbf{F} \cdot d\mathbf{r} = \left[\int_{\partial D} xF_1(x, y)dx + xF_2(x, y)dy, \int_{\partial D} yF_1(x, y)dx + yF_2(x, y)dy \right]^T = \mathbf{0}$$

and (2.2) holds, i.e., $\int_{\partial D} \mathbf{F}(\mathbf{r}) \cdot d\mathbf{r} \neq 0$.

First, suppose that \mathbf{F} is continuously differentiable on a domain that contains D . By Green's theorem, (2.4) holds if

$$\frac{\partial(xF_1)}{\partial y} = \frac{\partial(xF_2)}{\partial x} \quad \text{and} \quad \frac{\partial(yF_1)}{\partial y} = \frac{\partial(yF_2)}{\partial x},$$

or, equivalently,

$$(2.5) \quad x\frac{\partial F_2}{\partial x} - x\frac{\partial F_1}{\partial y} = -F_2 \quad \text{and} \quad y\frac{\partial F_2}{\partial x} - y\frac{\partial F_1}{\partial y} = F_1.$$

The equations in (2.5) give $F_2 = \frac{-x}{y}F_1$, which together with the first equation of (2.5) lead to

$$(2.6) \quad x\frac{\partial F_1}{\partial x} + y\frac{\partial F_1}{\partial y} = -2F_1.$$

The general solution of (2.6) that vanishes at infinity is $F_1(x, y) = \frac{1}{xy}\Phi\left(\frac{x}{r}, \frac{y}{r}\right)$, where Φ is an arbitrary real-valued differentiable function defined on the unit circle \mathbb{S}^1 and $r := \|\mathbf{r}\| = (x^2 + y^2)^{1/2}$. Further, $\Phi\left(\frac{x}{r}, \frac{y}{r}\right)$ is the general solution of the corresponding homogeneous equation.

PROPOSITION 2.1. *Let*

$$(2.7) \quad \mathbf{F} = [F_1, F_2]^T := \left[\frac{1}{xy}\Phi\left(\frac{x}{r}, \frac{y}{r}\right), \frac{-1}{y^2}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) \right]^T,$$

where Φ is an arbitrary real-valued differentiable function on the unit circle \mathbb{S}^1 . Then

$$(2.8) \quad x\mathbf{F} = \left[\frac{1}{y}\Phi\left(\frac{x}{r}, \frac{y}{r}\right), \frac{-x}{y^2}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) \right]^T \quad \text{and} \quad y\mathbf{F} = \left[\frac{1}{x}\Phi\left(\frac{x}{r}, \frac{y}{r}\right), \frac{-1}{y}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) \right]^T$$

are conservative vector fields on any region not containing the origin with potential functions ϕ_1 and ϕ_2 , respectively, where

$$(2.9) \quad \phi_1(x, y) = \int \frac{1}{y}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) dx = - \int \frac{x}{y^2}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) dy,$$

$$(2.10) \quad \phi_2(x, y) = \int \frac{1}{x}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) dx = - \int \frac{1}{y}\Phi\left(\frac{x}{r}, \frac{y}{r}\right) dy.$$

Proof. Using the fact that $x\frac{\partial\Phi}{\partial x} + y\frac{\partial\Phi}{\partial y} = 0$, it is straightforward to verify that $x\mathbf{F} = \nabla\phi_1$ and $y\mathbf{F} = \nabla\phi_2$. \square

It follows from Proposition 2.1 that the integrals of $x\mathbf{F}$ and $y\mathbf{F}$ over the boundary of any 2-D manifold in \mathbb{R}^2 that does not contain the origin are zero. We also require their integrals over any closed curve that encloses the origin to be zero, and it suffices to integrate over the unit circle \mathbb{S}^1 . We have

$$(2.11) \quad \int_{\mathbb{S}^1} x\mathbf{F} \cdot d\mathbf{r} = - \int_0^{2\pi} \frac{1}{\sin^2\theta}\Phi(\cos\theta, \sin\theta)d\theta,$$

$$(2.12) \quad \int_{\mathbb{S}^1} y\mathbf{F} \cdot d\mathbf{r} = - \int_0^{2\pi} \frac{1}{\cos\theta\sin\theta}\Phi(\cos\theta, \sin\theta)d\theta.$$

If

$$(2.13) \quad \Phi(\cos\theta, \sin\theta) = -\cos\theta\sin^2\theta(1 + \psi(\cos\theta, \sin\theta)),$$

where

$$\psi(\cos\theta, \sin\theta) := \sum_{k=2}^{\infty} a_k \cos k\theta + b_k \sin k\theta$$

is a differentiable function, then (2.11) and (2.12) give

$$(2.14) \quad \int_{\mathbb{S}^1} x\mathbf{F} \cdot d\mathbf{r} = \int_0^{2\pi} \cos\theta(1 + \psi(\cos\theta, \sin\theta))d\theta = 0,$$

$$(2.15) \quad \int_{\mathbb{S}^1} y\mathbf{F} \cdot d\mathbf{r} = \int_0^{2\pi} \sin\theta(1 + \psi(\cos\theta, \sin\theta))d\theta = 0,$$

since $1 + \psi$ is orthogonal to $\cos \theta$ and $\sin \theta$.

We now state the results and complete the proof.

THEOREM 2.2. *Let \mathbf{F} be as in Proposition 2.1 with Φ in (2.13). Then for any compact 2-D manifold $D \subset \mathbb{R}^2$ with piecewise smooth boundary ∂D and for $\mathbf{v} \in \mathbb{R}^2 \setminus \partial D$,*

$$(2.16) \quad \int_{\partial D} (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} = \mathbf{0}.$$

Further, if $|\psi| \leq 1$,

$$(2.17) \quad \mathbf{v} = \frac{\int_{\partial D} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}{\int_{\partial D} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r}}$$

for any $\mathbf{v} \in \mathbb{R}^2$.

We first prove a lemma, from which one can deduce that (2.2) holds if $|\psi| \leq 1$.

LEMMA 2.3. *For any compact 2-D manifold $D \subset \mathbb{R}^2$ with piecewise smooth boundary,*

$$(2.18) \quad \int_{\partial D} \mathbf{F} \cdot d\mathbf{r} = \begin{cases} -\iint_D \frac{1}{r^3} (1 + \psi(\frac{x}{r}, \frac{y}{r})) \, dx dy & \text{if } \mathbf{0} \notin D, \\ \iint_{\mathbb{R}^2 \setminus D} \frac{1}{r^3} (1 + \psi(\frac{x}{r}, \frac{y}{r})) \, dx dy & \text{if } \mathbf{0} \in D. \end{cases}$$

In particular, $\int_{\partial D} \mathbf{F} \cdot d\mathbf{r} \neq 0$ if $|\psi| \leq 1$.

Proof. If $\mathbf{0} \notin D$, by Green's theorem

$$(2.19) \quad \begin{aligned} \int_{\partial D} \mathbf{F} \cdot d\mathbf{r} &= \iint_D \frac{\partial}{\partial x} \left(\frac{-1}{y^2} \Phi \right) - \frac{\partial}{\partial y} \left(\frac{1}{xy} \Phi \right) \, dx dy \\ &= \iint_D \frac{1}{xy^2} \Phi \left(\frac{x}{r}, \frac{y}{r} \right) - \frac{1}{xy^2} \left\{ x \frac{\partial \Phi}{\partial x} + y \frac{\partial \Phi}{\partial y} \right\} \, dx dy \\ &= \iint_D \frac{1}{xy^2} \Phi \left(\frac{x}{r}, \frac{y}{r} \right) \, dx dy \\ &= - \iint_D \frac{1}{r^3} \left(1 + \psi \left(\frac{x}{r}, \frac{y}{r} \right) \right) \, dx dy. \end{aligned}$$

If $\mathbf{0} \in D$, take a disc D_R , with center at the origin and radius R , that contains D . By (2.19),

$$\int_{\partial D_R} \mathbf{F} \cdot d\mathbf{r} - \int_{\partial D} \mathbf{F} \cdot d\mathbf{r} = - \iint_{D_R \setminus D} \frac{1}{r^3} \left(1 + \psi \left(\frac{x}{r}, \frac{y}{r} \right) \right) \, dx dy$$

so that

$$\begin{aligned} \int_{\partial D} \mathbf{F} \cdot d\mathbf{r} &= \int_{\partial D_R} \mathbf{F} \cdot d\mathbf{r} + \iint_{D_R \setminus D} \frac{1}{r^3} \left(1 + \psi \left(\frac{x}{r}, \frac{y}{r} \right) \right) \, dx dy \\ &= -\frac{1}{R} \int_0^{2\pi} 1 + \psi(\cos \theta, \sin \theta) d\theta + \iint_{D_R \setminus D} \frac{1}{r^3} \left(1 + \psi \left(\frac{x}{r}, \frac{y}{r} \right) \right) \, dx dy. \end{aligned}$$

Taking the limit as $R \rightarrow \infty$ gives the second formula in (2.18). □

Proof of Theorem 2.2. If $\mathbf{v} \notin \partial D$, (2.16) follows from Proposition 2.1 and (2.14) and (2.15). Equation (2.17) follows from (2.16) and Lemma 2.3.

If $\mathbf{v} \in \partial D$, for any $\epsilon > 0$, we replace an appropriate segment of ∂D that contains \mathbf{v} by a circular arc with center at \mathbf{v} and radius ϵ , so that \mathbf{v} lies outside the resulting piecewise smooth curve, which we call ∂D_ϵ . Then (2.17) holds for ∂D_ϵ for all $\epsilon > 0$, and therefore holds in the limit as $\epsilon \rightarrow 0$. \square

If we choose $\psi = 0$ in (2.13) so that

$$(2.20) \quad \Phi(\cos \theta, \sin \theta) = -\cos \theta \sin^2 \theta,$$

a more general form of Floater’s mean value coordinates [2] can be deduced as a corollary. The proof here, which is different from the existing methods (see [2], [3], [5] [6]), provides a better understanding why they work.

COROLLARY 2.4. *Let $D \subset \mathbb{R}^2$ be a compact 2-D manifold with piecewise smooth boundary ∂D , and let $\mathbf{v}_j \in \partial D$, $j = 1, 2, \dots, k > 2$, be arranged in the positive orientation. Then for any $\mathbf{v} \in \mathbb{R}^2$,*

$$(2.21) \quad \mathbf{v} = \frac{\sum_{j=1}^k w_j \mathbf{v}_j}{\sum_{j=1}^k w_j},$$

where

$$(2.22) \quad w_j := \frac{\tan(\alpha_{j-1}/2) + \tan(\alpha_j/2)}{\|\mathbf{v}_j - \mathbf{v}\|}, \quad j = 1, 2, \dots, k - 1,$$

and α_j is the angle at \mathbf{v} of the oriented triangle $[\mathbf{v}, \mathbf{v}_j, \mathbf{v}_{j+1}]$ with $|\alpha_j| < \pi$ and takes a positive value if the vector $\frac{\mathbf{v}_j}{\|\mathbf{v}_j\|} - \frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|}$ is in the counterclockwise orientation and a negative value otherwise.

Proof. Take $\psi = 0$ in (2.13) so that $\Phi\left(\frac{x}{r}, \frac{y}{r}\right) = -\frac{xy^2}{r^3}$ and $\mathbf{F}(\mathbf{r}) \equiv \mathbf{F}(x, y) := [-y/r^3, x/r^3]^T$. By (2.9) and (2.10), $x\mathbf{F} = \nabla\phi_1$ and $y\mathbf{F} = \nabla\phi_2$, where

$$(2.23) \quad \phi_1(x, y) = -\frac{1}{y} \int \frac{xy^2}{r^3} dx = \frac{y}{r},$$

$$(2.24) \quad \phi_2(x, y) = \int \frac{xy^2}{yr^3} dy = \frac{-x}{r}.$$

Hence the vector $[\phi_1(x, y), \phi_2(x, y)]^T = [y, -x]^T/r$ is the rotation of the $[x, y]^T/r$ along the unit circle in the clockwise direction through an angle of $\pi/2$. Let $\mathbf{v}_j = [x_j, y_j]$, $j = 1, 2, \dots, k$, $\mathbf{v}_{k+1} := \mathbf{v}_1$. Then for $j = 1, 2, \dots, k$,

$$\begin{aligned} \int_{\mathbf{v}_j}^{\mathbf{v}_{j+1}} \mathbf{r} \mathbf{F} \cdot d\mathbf{r} &= [\phi_1(\mathbf{v}_{j+1}), \phi_2(\mathbf{v}_{j+1})]^T - [\phi_1(\mathbf{v}_j), \phi_2(\mathbf{v}_j)]^T \\ &= [y_{j+1}, -x_{j+1}]^T/\|\mathbf{v}_{j+1}\| - [y_j, -x_j]^T/\|\mathbf{v}_j\|. \end{aligned}$$

The vector $[y_{j+1}, -x_{j+1}]^T/\|\mathbf{v}_{j+1}\| - [y_j, -x_j]^T/\|\mathbf{v}_j\|$ is parallel to $\mathbf{v}_{j+1}/\|\mathbf{v}_{j+1}\| + \mathbf{v}_j/\|\mathbf{v}_j\|$ and

$$\begin{aligned} \left\| \frac{[y_{j+1}, -x_{j+1}]^T}{\|\mathbf{v}_{j+1}\|} - \frac{[y_j, -x_j]^T}{\|\mathbf{v}_j\|} \right\| &= \left\| \frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|} - \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|} \right\| \\ &= \tan(\alpha_j/2) \left\| \frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|} + \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|} \right\|, \end{aligned}$$

where α_j is the angle at the origin, $\mathbf{0}$, of the triangle $[\mathbf{0}, \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|}, \frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|}]$. Hence

$$\int_{\mathbf{v}_j}^{\mathbf{v}_{j+1}} \mathbf{r} \mathbf{F} \cdot d\mathbf{r} = \tan(\alpha_j/2) \left\{ \frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|} + \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|} \right\}$$

with the convention that α_j takes the positive sign if $\frac{\mathbf{v}_{j+1}}{\|\mathbf{v}_{j+1}\|} - \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|}$ is counterclockwise and the negative sign if it is clockwise. Hence for $\mathbf{v} \in \mathbb{R}^2$,

$$\int_{\mathbf{v}_j}^{\mathbf{v}_{j+1}} (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} = \tan(\alpha_j/2) \left\{ \frac{\mathbf{v}_j - \mathbf{v}}{\|\mathbf{v}_j - \mathbf{v}\|} + \frac{\mathbf{v}_{j+1} - \mathbf{v}}{\|\mathbf{v}_{j+1} - \mathbf{v}\|} \right\},$$

where α_j is the angle at \mathbf{v} of triangle $[\mathbf{v}, \mathbf{v}_j, \mathbf{v}_{j+1}]$. Substituting this into

$$\sum_{j=1}^k \int_{\mathbf{v}_j}^{\mathbf{v}_{j+1}} (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} = \int_{\partial D} (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{r} = \mathbf{0}$$

gives (2.21) with w_j given by (2.22). □

3. Representation on boundaries of compact 3-D manifolds in \mathbb{R}^3 . As

in the previous section, we want to find a differentiable vector field $\mathbf{F}(\mathbf{r}) = [F_1(x, y, z), F_2(x, y, z), F_3(x, y, z)]^T$, $\mathbf{r} = [x, y, z]^T \in \mathbb{R}^3 \setminus \{\mathbf{0}\}$, that satisfies

$$(3.1) \quad \iint_{\partial M} \mathbf{r} \mathbf{F}(\mathbf{r}) \cdot d\mathbf{S} = \mathbf{0}$$

and

$$(3.2) \quad \iint_{\partial M} \mathbf{F}(\mathbf{r}) \cdot d\mathbf{S} \neq 0$$

for any compact 3-D manifold $M \in \mathbb{R}^3$ with piecewise smooth boundary ∂M . We assume that the positive orientation of ∂M is one in which its normal points away from M . Then if (3.1) and (3.2) hold, then, for any $\mathbf{v} \in \mathbb{R}^3$,

$$(3.3) \quad \mathbf{v} = \frac{\iint_{\partial M} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}{\iint_{\partial M} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}.$$

A similar search as in the previous section requiring that $x\mathbf{F}$, $y\mathbf{F}$, and $z\mathbf{F}$ be divergence free, i.e., $\text{div}(x\mathbf{F}) = \text{div}(y\mathbf{F}) = \text{div}(z\mathbf{F}) = 0$, shows that F_i , $i = 1, 2, 3$, are solutions of the partial differential equation

$$(3.4) \quad x \frac{\partial U}{\partial x} + y \frac{\partial U}{\partial y} + z \frac{\partial U}{\partial z} = -3U.$$

The general solution of (3.4) that vanishes at infinity is

$$(3.5) \quad U(x, y, z) = \frac{1}{xyz} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right),$$

where $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}$ is an arbitrary differentiable function on the unit sphere $\mathbb{S}^2 \subset \mathbb{R}^3$, $\rho := (x^2 + y^2 + z^2)^{1/2}$, and $\Phi(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho})$ is the general solution of the corresponding homogeneous equation. To state the results, we define spherical harmonics

$$\begin{aligned} C_\ell^m(\phi, \theta) &:= N_\ell^m P_\ell^m(\cos \phi) \cos m\theta, \quad m = 0, 1, \dots, \ell, \\ S_\ell^m(\phi, \theta) &:= N_\ell^m P_\ell^m(\cos \phi) \sin m\theta, \quad m = 1, 2, \dots, \ell, \end{aligned}$$

where P_ℓ^m are the associated Legendre polynomials and N_ℓ^m are normalization constants (see [9]).

THEOREM 3.1. *Let*

(3.6)

$$\mathbf{F}(\mathbf{r}) \equiv \mathbf{F}(x, y, z) := \left[\frac{1}{xyz} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right), \frac{1}{x^2z} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right), \frac{1}{x^2y} \Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right) \right]^T,$$

where $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}$ is an arbitrary differentiable function. Then $x\mathbf{F}$, $y\mathbf{F}$, and $z\mathbf{F}$ are divergence free on $\mathbb{R}^3 \setminus \{\mathbf{0}\}$.

Further, if

$$\begin{aligned} \Phi \left(\frac{x}{\rho}, \frac{x}{\rho}, \frac{x}{\rho} \right) &:= \frac{x^2yz}{\rho^4} \left(1 + \Psi \left(\frac{x}{\rho}, \frac{x}{\rho}, \frac{x}{\rho} \right) \right) \\ (3.7) \quad &= \sin^3 \phi \cos^2 \theta \sin \theta \cos \phi \left(1 + \Psi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi) \right) \end{aligned}$$

and

$$(3.8) \quad \Psi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi) := \sum_{\ell=2}^{\infty} \left\{ \sum_{m=0}^{\ell} a_\ell^m C_\ell^m(\phi, \theta) + \sum_{m=1}^{\ell} b_\ell^m S_\ell^m(\phi, \theta) \right\},$$

then for any compact 3-D manifold M in \mathbb{R}^3 with piecewise smooth boundary and for $\mathbf{v} \in \mathbb{R}^3 \setminus \partial M$,

$$(3.9) \quad \iint_{\partial M} (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S} = \mathbf{0}.$$

In addition, if $|\Psi| \leq 1$, then for any $\mathbf{v} \in \mathbb{R}^3$,

$$(3.10) \quad \mathbf{v} = \frac{\iint_{\partial M} \mathbf{r} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}{\iint_{\partial M} \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S}}.$$

Proof. Suppose \mathbf{F} is as in (3.6), where $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}$ is an arbitrary differentiable function on the unit sphere \mathbb{S}^2 . Then

$$\operatorname{div}(x\mathbf{F}) = \frac{1}{xyz} \left(x \frac{\partial \Phi}{\partial x} + y \frac{\partial \Phi}{\partial y} + z \frac{\partial \Phi}{\partial z} \right) = 0, \quad (x, y, z) \in \mathbb{R}^3 \setminus \{\mathbf{0}\},$$

since $\Phi(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho})$ satisfies

$$x \frac{\partial \Phi}{\partial x} + y \frac{\partial \Phi}{\partial y} + z \frac{\partial \Phi}{\partial z} = 0.$$

Similarly, $\operatorname{div}(y\mathbf{F}) = \operatorname{div}(z\mathbf{F}) = 0$.

Suppose Φ is given by (3.7) and (3.8). To prove (3.9) we may assume $\mathbf{v} = \mathbf{0}$, so it reduces to proving (3.1). Since $x\mathbf{F}$, $y\mathbf{F}$, and $z\mathbf{F}$ are divergence free, (3.1) holds for any compact 3-D manifold M that does not contain the origin. To prove that (3.1) holds for any M that contains the origin, it suffices to prove that it holds for the unit

ball B with boundary $\partial B = \mathbb{S}^2$. In spherical polar coordinates,

$$(3.11) \quad \iint_{\mathbb{S}^2} x \mathbf{F} \cdot d\mathbf{S} = \int_0^{2\pi} \int_0^\pi \frac{\Phi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)}{\sin \phi \cos \phi \sin \theta \cos \theta} d\phi d\theta,$$

$$(3.12) \quad \iint_{\mathbb{S}^2} y \mathbf{F} \cdot d\mathbf{S} = \int_0^{2\pi} \int_0^\pi \frac{\Phi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)}{\sin \phi \cos \phi \cos^2 \theta} d\phi d\theta,$$

$$(3.13) \quad \iint_{\mathbb{S}^2} z \mathbf{F} \cdot d\mathbf{S} = \int_0^{2\pi} \int_0^\pi \frac{\Phi(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)}{\sin^2 \phi \cos^2 \theta \sin \theta} d\phi d\theta.$$

By (3.7) and (3.11)–(3.13),

$$\begin{aligned} \iint_{\mathbb{S}^2} \mathbf{r} \mathbf{F} \cdot d\mathbf{S} &\equiv \iint_{\mathbb{S}^2} [x, y, z]^T \mathbf{F} \cdot d\mathbf{S} \\ &= \int_0^{2\pi} \int_0^\pi (1 + \Psi) [C_1^1(\phi, \theta), S_1^1(\phi, \theta), C_1^0(\phi, \theta)]^T \sin \phi d\phi d\theta, \end{aligned}$$

where

$$C_1^0(\phi, \theta) = \cos \phi, \quad C_1^1(\phi, \theta) = \sin \phi \cos \theta, \quad S_1^1(\phi, \theta) = \sin \phi \sin \theta$$

are the first order spherical harmonics. Since $1 + \Psi$ is orthogonal to C_1^0 , C_1^1 , and S_1^1 , it follows that $\iint_{\mathbb{S}^2} \mathbf{r} \mathbf{F} \cdot d\mathbf{S} = 0$.

To prove (3.10) we first show that (3.2) holds. The same argument as in the proof of Lemma 2.3 using the divergence theorem gives

$$(3.14) \quad \iint_{\partial M} \mathbf{F} \cdot d\mathbf{S} = \begin{cases} \iiint_M \frac{1}{\rho^4} \left(1 + \Psi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right) \right) dx dy dz & \text{if } \mathbf{0} \notin M, \\ - \iiint_{\mathbb{R}^3 \setminus M} \frac{1}{\rho^4} \left(1 + \Psi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right) \right) dx dy dz & \text{if } \mathbf{0} \in M. \end{cases}$$

In particular, $\iint_{\partial M} \mathbf{F} \cdot d\mathbf{S} \neq 0$ if $|\Psi| \leq 1$. Hence (3.10) follows from (3.9) for $\mathbf{v} \notin \partial M$. A limiting argument as in the proof of Theorem 2.2 shows that it also holds for $\mathbf{v} \in \partial M$. \square

4. Mean value coordinates in \mathbb{R}^3 . We now apply the results in Theorem 3.1 to compute mean value coordinates in \mathbb{R}^3 . An explicit expression for mean value coordinates for points that lie in the kernel of a star-shaped polyhedron in \mathbb{R}^3 is given in [3]. An algorithm for computing these coordinates for all points in \mathbb{R}^3 can be found in [6], but no explicit expression is given there. In this section we derive a new formula for mean value coordinates for points in \mathbb{R}^3 with respect to the boundary of any compact 3-D manifold with polyhedral faces of any topology. In the case where all the faces of ∂M are triangular, the formula is equivalent to that in [3].

Consider $\mathbf{F} = [F_1, F_2, F_3]^T$ in (3.6), where Ψ is given by (3.7) as in Theorem 3.1. Throughout this section we shall assume that $\Psi \equiv 0$ so that Φ takes the simplest form:

$$\Phi \left(\frac{x}{\rho}, \frac{y}{\rho}, \frac{z}{\rho} \right) = \frac{x^2 y z}{\rho^4}$$

and

$$(4.1) \quad \mathbf{F} \equiv \mathbf{F}(\mathbf{r}) = \frac{\mathbf{r}}{\rho^4}, \quad \rho \neq 0.$$

Now

$$\begin{aligned} \operatorname{curl}(\mathbf{F}) &= \nabla \left(\frac{1}{\rho^4} \right) \times \mathbf{r} + \frac{1}{\rho^4} \operatorname{curl}(\mathbf{r}) \\ &= \frac{-4}{\rho^6} \mathbf{r} \times \mathbf{r} = \mathbf{0} \end{aligned}$$

and

$$\operatorname{div}(\mathbf{F}) = \nabla \left(\frac{1}{\rho^4} \right) \cdot \mathbf{r} + \frac{1}{\rho^4} \operatorname{div}(\mathbf{r}) = -\frac{1}{\rho^4}, \quad \rho \neq 0.$$

Since

$$\operatorname{div}(x\mathbf{F}) = \operatorname{div}(y\mathbf{F}) = \operatorname{div}(z\mathbf{F}) = 0,$$

we want to find vector fields $\mathbf{A}, \mathbf{B}, \mathbf{C}$ such that

$$(4.2) \quad x\mathbf{F} = \operatorname{curl}(\mathbf{A}), \quad y\mathbf{F} = \operatorname{curl}(\mathbf{B}), \quad z\mathbf{F} = \operatorname{curl}(\mathbf{C}), \quad \rho \neq 0,$$

so that Stoke's theorem can be applied to evaluate the integrals of $x\mathbf{F}$, $y\mathbf{F}$, $z\mathbf{F}$ over the polyhedral faces of ∂M .

By (4.2)

$$\operatorname{curl}(x\mathbf{F}) = \nabla (\operatorname{div}(\mathbf{A})) - \nabla^2 \mathbf{A},$$

so that if $\operatorname{div}(\mathbf{A}) = 0$, the vector components of \mathbf{A} satisfy Poisson's equations

$$(4.3) \quad \nabla^2 \mathbf{A} = -\operatorname{curl}(x\mathbf{F}), \quad \rho \neq 0.$$

Similarly, if $\operatorname{div}(\mathbf{B}) = \operatorname{div}(\mathbf{C}) = 0$,

$$(4.4) \quad \nabla^2 \mathbf{B} = -\operatorname{curl}(y\mathbf{F}),$$

$$(4.5) \quad \nabla^2 \mathbf{C} = -\operatorname{curl}(z\mathbf{F}), \quad \rho \neq 0.$$

Now

$$\operatorname{curl}(x\mathbf{F}) = \nabla x \times \mathbf{F} + x \operatorname{curl}(\mathbf{F}) = [0, -F_3, F_2]^T = \frac{1}{\rho^4} [0, -z, y]^T,$$

$$\operatorname{curl}(y\mathbf{F}) = \nabla y \times \mathbf{F} + y \operatorname{curl}(\mathbf{F}) = [F_3, 0, -F_1]^T = \frac{1}{\rho^4} [z, 0, -x]^T,$$

$$\operatorname{curl}(z\mathbf{F}) = \nabla z \times \mathbf{F} + z \operatorname{curl}(\mathbf{F}) = [-F_2, F_1, 0]^T = \frac{1}{\rho^4} [-y, x, 0]^T,$$

so that by (4.3)–(4.5),

$$(4.6) \quad \nabla^2 \mathbf{A} = \frac{1}{\rho^4} [0, z, -y]^T,$$

$$(4.7) \quad \nabla^2 \mathbf{B} = \frac{1}{\rho^4} [-z, 0, x]^T,$$

$$(4.8) \quad \nabla^2 \mathbf{C} = \frac{1}{\rho^4} [y, -x, 0]^T, \quad \rho \neq 0.$$

We look for divergence free solutions, $\mathbf{A}, \mathbf{B}, \mathbf{C}$, of Poisson’s equations (4.6), (4.7), (4.8), which satisfy (4.2) and vanish at infinity.

PROPOSITION 4.1. *The vector fields*

$$(4.9) \quad \mathbf{A} = \frac{1}{2\rho^2}[0, z, -y]^T, \quad \mathbf{B} = \frac{1}{2\rho^2}[-z, 0, x]^T, \quad \mathbf{C} = \frac{1}{2\rho^2}[y, -x, 0]^T$$

are divergence free and satisfy (4.2). Further, they are unique divergence free solutions of (4.6), (4.7), (4.8) that satisfy (4.2) and vanish at infinity.

Proof. It is straightforward to verify that $\mathbf{A}, \mathbf{B}, \mathbf{C}$ satisfy (4.2) and $\text{div}(\mathbf{A}) = \text{div}(\mathbf{B}) = \text{div}(\mathbf{C}) = 0$. Hence they are particular solutions of (4.6), (4.7), (4.8), respectively. To prove uniqueness, we consider (4.6), in which the general solution is

$$\mathbf{A}_g = \left[U, U + \frac{z}{2\rho^2}, U - \frac{y}{2\rho^2} \right]^T,$$

where U is the general solution of the corresponding homogeneous equation, which is the Laplace equation. The conditions that $\text{div}(\mathbf{A}_g) = 0$, $\text{curl}(\mathbf{A}_g) = x\mathbf{F}$, and \mathbf{A}_g vanishes at infinity imply $U = 0$. Hence $\mathbf{A}_g = \mathbf{A}$. The proof of uniqueness of \mathbf{B} and \mathbf{C} is the same. \square

Consider a 3-D manifold, $M \subset \mathbb{R}^3$, with polyhedral faces (not necessarily planar) of arbitrary topology, and each face has piecewise linear boundaries. Let \mathcal{F} be the set of all its faces, which are oriented with their normals pointing away from M . The orientation of each face induces an orientation on the edges that form its boundary. Let \mathcal{V} be the set of all vertices of ∂M .

PROPOSITION 4.2. *Let S be a face of ∂M whose boundary is the union of oriented line segments $[\mathbf{v}_j, \mathbf{v}_{j+1}]$, $j = 1, 2, \dots, k = k(S)$, $\mathbf{v}_{k+1} = \mathbf{v}_1$. Then for any $\mathbf{v} \notin \partial S$,*

$$(4.10) \quad \iint_S (\mathbf{r} - \mathbf{v})\mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S} = \frac{1}{2} \sum_{j=1}^k \alpha_j \mathbf{n}_j,$$

where $\alpha_j \in (0, \pi)$ is the angle at \mathbf{v} of the triangle $[\mathbf{v}, \mathbf{v}_j, \mathbf{v}_{j+1}]$ and

$$(4.11) \quad \mathbf{n}_j = \frac{(\mathbf{v}_{j+1} - \mathbf{v}) \times (\mathbf{v}_j - \mathbf{v})}{\|(\mathbf{v}_{j+1} - \mathbf{v}) \times (\mathbf{v}_j - \mathbf{v})\|}.$$

Proof. We assume without loss of generality that $\mathbf{v} = \mathbf{0}$. By Stoke’s theorem

$$(4.12) \quad \begin{aligned} \iint_S \mathbf{r} \mathbf{F} \cdot d\mathbf{S} &= \left[\iint_S x \mathbf{F} \cdot d\mathbf{S}, \iint_S y \mathbf{F} \cdot d\mathbf{S}, \iint_S z \mathbf{F} \cdot d\mathbf{S} \right]^T \\ &= \left[\int_{\partial S} \mathbf{A} \cdot d\mathbf{r}, \int_{\partial S} y\mathbf{B} \cdot d\mathbf{r}, \int_{\partial S} \mathbf{C} \cdot d\mathbf{r} \right]^T \end{aligned}$$

$$(4.13) \quad = \sum_{j=1}^k \left[\int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{A} \cdot d\mathbf{r}, \int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{B} \cdot d\mathbf{r}, \int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{C} \cdot d\mathbf{r} \right]^T.$$

Let $\mathbf{v}_j = [x_j, y_j, z_j]^T$. A straightforward computation gives

$$\int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{A} \cdot d\mathbf{r} = - \begin{vmatrix} y_j & z_j \\ y_{j+1} & z_{j+1} \end{vmatrix} \left(\frac{\alpha_j}{2\|\mathbf{v}_{j+1} \times \mathbf{v}_j\|} \right),$$

$$\int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{B} \cdot d\mathbf{r} = \begin{vmatrix} x_j & z_j \\ x_{j+1} & z_{j+1} \end{vmatrix} \left(\frac{\alpha_j}{2\|\mathbf{v}_{j+1} \times \mathbf{v}_j\|} \right),$$

$$\int_{[\mathbf{v}_j, \mathbf{v}_{j+1}]} \mathbf{C} \cdot d\mathbf{r} = - \begin{vmatrix} x_j & y_j \\ x_{j+1} & y_{j+1} \end{vmatrix} \left(\frac{\alpha_j}{2\|\mathbf{v}_{j+1} \times \mathbf{v}_j\|} \right),$$

which, by (4.13), leads to (4.10) and (4.11). \square

The next step in the computation of the mean value coordinates is to express the sum on the right of (4.10) in terms of $\mathbf{v}_j - \mathbf{v}$, $j = 1, 2, \dots, k$:

$$(4.14) \quad \frac{1}{2} \sum_{j=1}^k \alpha_j \mathbf{n}_j = \sum_{j=1}^k a_j(S)(\mathbf{v}_j - \mathbf{v}),$$

where $a_j(S) \equiv a_{\mathbf{v}_j}(S)$ so that

$$\mathbf{0} = \sum_{S \in \mathcal{F}} \iint_S (\mathbf{r} - \mathbf{v}) \mathbf{F}(\mathbf{r} - \mathbf{v}) \cdot d\mathbf{S} = \sum_{S \in \mathcal{F}} \sum_{j=1}^{k(S)} a_j(S)(\mathbf{v}_j - \mathbf{v}),$$

which would lead to

$$(4.15) \quad \mathbf{v} = \frac{\sum_{S \in \mathcal{F}} \sum_{j=1}^{k(S)} a_j(S) \mathbf{v}_j}{\sum_{S \in \mathcal{F}} \sum_{j=1}^{k(S)} a_j(S)} = \frac{\sum_{\mathbf{p} \in \mathcal{V}} \lambda_{\mathbf{p}} \mathbf{p}}{\sum_{\mathbf{p} \in \mathcal{V}} \lambda_{\mathbf{p}}},$$

where

$$(4.16) \quad \lambda_{\mathbf{p}} := \sum_{\substack{S \in \mathcal{F} \\ \mathbf{p} \in S}} a_{\mathbf{p}}(S).$$

If $k = 3$, $a_j(S)$ can be obtained uniquely by taking the inner products of the expressions in (4.14) with \mathbf{n}_ℓ or $\mathbf{v}_\ell - \mathbf{v}$, $\ell = 1, 2, 3$. Taking the inner product with \mathbf{n}_ℓ gives

$$a_{\ell+2}(S) = \frac{1}{2 \mathbf{n}_\ell \cdot (\mathbf{v}_{\ell+2} - \mathbf{v})} \sum_{j=1}^3 \mathbf{n}_\ell \cdot \mathbf{n}_j \alpha_j, \quad \ell = 1, 2, 3,$$

where $a_j = a_{j+3}$ and $\mathbf{n}_j = \mathbf{n}_{j+3}$ for all j . These coordinates are the same as those obtained earlier in [3] for the representation of points in the kernel of 3-D star-shaped manifolds. We have shown that they are applicable to all points in \mathbb{R}^3 with respect to the boundaries of more general manifolds. If $k > 3$, the representation (4.14) exists, but is not unique. We shall derive a formula for $a_j(S)$ in (4.14) for the general case, which reduces to the above formula when $k = 3$. The idea is to express

$$(4.17) \quad \frac{1}{2} \sum_{j=1}^k \alpha_j \mathbf{n}_j = b_{\ell,0}(\mathbf{v} - \mathbf{v}_\ell) + b_{\ell,1}(\mathbf{v} - \mathbf{v}_{\ell+1}) + b_{\ell,2}(\mathbf{v} - \mathbf{v}_{\ell+2})$$

for $\ell = 1, 2, \dots, k$, where $b_{k+1,i} = b_{1,i}$, $i = 0, 1, 2$, and define

$$a_j(S) := \frac{1}{k}(b_{j,0} + b_{j-1,1} + b_{j-2,2}), \quad j = 1, 2, \dots, k.$$

Then

$$\begin{aligned} \frac{1}{2} \sum_{j=1}^k \alpha_j \mathbf{n}_j &= \frac{1}{k} \sum_{\ell=1}^k \sum_{i=0}^2 b_{\ell,i}(\mathbf{v} - \mathbf{v}_{\ell+i}) \\ &= \frac{1}{k} \sum_{j=1}^k \sum_{i=0}^2 b_{j-i,i}(\mathbf{v} - \mathbf{v}_j) \\ &= \sum_{j=1}^k a_j(S)(\mathbf{v} - \mathbf{v}_j), \end{aligned}$$

as required by (4.14) so that (4.15) holds.

To compute the coefficients $b_{\ell,i}$, $i = 0, 1, 2$, taking the inner product of (4.17) with \mathbf{n}_ℓ , $\mathbf{n}_{\ell+1}$, and $\mathbf{n}_{\ell+2}$ leads to

$$\begin{aligned} b_{\ell,2} &= \frac{1}{2 \mathbf{v}_{\ell+2} \cdot \mathbf{n}_\ell} \sum_{j=1}^k \mathbf{n}_\ell \cdot \mathbf{n}_j \alpha_j, \\ b_{\ell,0} &= \frac{1}{2 \mathbf{v}_\ell \cdot \mathbf{n}_{\ell+1}} \sum_{j=1}^k \mathbf{n}_{\ell+1} \cdot \mathbf{n}_j \alpha_j, \end{aligned}$$

and

$$b_{\ell,1} = \frac{1}{2 \mathbf{v}_{\ell+1} \cdot \mathbf{n}_{\ell+2}} \sum_{j=1}^k \mathbf{n}_{\ell+2} \cdot \mathbf{n}_j \alpha_j - \frac{\mathbf{v}_\ell \cdot \mathbf{n}_{\ell+2}}{2 \mathbf{v}_\ell \cdot \mathbf{n}_{\ell+1} \mathbf{v}_{\ell+1} \cdot \mathbf{n}_{\ell+2}} \sum_{j=1}^k \mathbf{n}_{\ell+1} \cdot \mathbf{n}_j \alpha_j$$

for $\ell = 1, 2, \dots, k$.

Remark. In general, the functions Φ and Ψ in (3.7) and (3.8), respectively, define the vector field

$$\mathbf{F}(\mathbf{r}) = \frac{g(\mathbf{r})}{\rho^4} \mathbf{r}, \quad \mathbf{r} = (x, y, z) \in \mathbb{R}^3 \setminus \{\mathbf{0}\},$$

where

$$g(\mathbf{r}) := 1 + \Psi \left(\frac{\mathbf{r}}{\rho} \right),$$

and \mathbf{F} provides mean value representations. Since

$$(4.18) \quad \mathbf{F}(\mathbf{r}) = -g(\mathbf{r}) \nabla \left(\frac{1}{2\rho^2} \right)$$

and

$$(4.19) \quad \operatorname{div}(x\mathbf{F}(\mathbf{r})) = 0, \quad \operatorname{div}(y\mathbf{F}(\mathbf{r})) = 0, \quad \operatorname{div}(z\mathbf{F}(\mathbf{r})) = 0,$$

the system (4.18), (4.19) can be discretized by the mimetic finite difference schemes (see [1], [8], and the references therein) to give a large class of mean value type coordinates.

Acknowledgments. Thanks to an anonymous referee for pointing out the connection with mimetic finite difference schemes. Thanks also to Rick Beatson for useful discussions that generated my interest in the subject and to the Abdus Salam International Centre for Theoretical Physics, Trieste, for providing an excellent environment and financial support during my visit in March/April 2007.

REFERENCES

- [1] F. BREZZI, K. LIPNIKOV, M. SHASHKOV, AND V. SIMONCINI, *A new discretization methodology for diffusion problems on generalized polyhedral meshes*, *Comput. Methods Appl. Mech. Engrg.*, 196 (2007), pp. 3682–3692.
- [2] M. S. FLOATER, *Mean value coordinates*, *Comput. Aided Geom. Design*, 20 (2003), pp. 19–27.
- [3] M. S. FLOATER, G. KOS, AND M. REIMERS, *Mean value coordinates in 3D*, *Comput. Aided Geom. Design*, 22 (2005), pp. 623–631.
- [4] K. HORMANN AND M. S. FLOATER, *Mean value coordinates for arbitrary planar polygons*, *ACM Trans. Graph.*, 25 (2006), pp. 1424–1441.
- [5] T. JU, S. SCHAEFER, AND J. WARREN, *Mean value coordinates for closed triangular meshes*, *ACM Trans. Graph.*, 24 (2005), pp. 561–566.
- [6] T. LANGER, A. BELYAEV, AND H.-P. SEIDEL, *Mean value coordinates for arbitrary spherical polygons and polyhedra in \mathbb{R}^3* , in *Curve and Surface Design: Avignon, 2006*, P. Chenin, T. Lyche, and L. L. Schumaker, eds., Nashboro Press, Brentwood, TN, 2007, pp. 193–202.
- [7] T. LANGER, A. BELYAEV, AND H.-P. SEIDEL, *Spherical barycentric coordinates*, in *Fourth Eurographics Symposium on Geometry Processing*, A. Sheffer and K. Polthier, eds., Eurographics Association, 2006, pp. 81–88.
- [8] K. LIPNIKOV, M. SHASHKOV, AND D. SVYATSKIY, *The mimetic finite difference discretization of diffusion problem on unstructured polyhedral meshes*, *J. Comput. Phys.*, 211 (2006), pp. 473–491.
- [9] E. T. WHITTAKER AND G. N. WATSON, *Solution of Laplace’s equation involving Legendre functions*, in *A Course in Modern Analysis*, 4th ed., Cambridge University Press, 1990, pp. 391–395.

COLLISIONS IN THREE-DIMENSIONAL FLUID STRUCTURE INTERACTION PROBLEMS*

MATTHIEU HILLAIRET[†] AND TAKÉO TAKAHASHI[‡]

Abstract. This paper deals with a system composed of a rigid ball moving into a viscous incompressible fluid over a fixed horizontal plane. The equations of motion for the fluid are the Navier–Stokes equations, and the equations for the motion of the rigid ball are obtained by applying Newton’s laws. We show that for any weak solution of the corresponding system satisfying the energy inequality, the rigid ball never touches the plane. This result is the extension of that obtained in [M. Hillairet, *Comm. Partial Differential Equations*, 32 (2007), pp. 1345–1371] in the two-dimensional setting.

Key words. fluid structure interactions, Cauchy theory, qualitative properties, collisions

AMS subject classifications. 35R35, 76D03, 76D05

DOI. 10.1137/080716074

1. Introduction. In the last decade, several studies showed that collisions between rigid bodies in a fluid would lead to great difficulties in the mathematical treatment of fluid-structure interaction models. For example, in [4, p. 287], Feireisl constructs a solution in which a sphere remains stuck to the ceiling of the cavity regardless of the intensity of the gravity. This example emphasizes that collisions would lead to unphysical solutions to standard mathematical systems. Indeed, Starovoitov proves also that several solutions exist when contact occurs [11, 12]. Therefore, at least one of these solutions does not represent a physical configuration.

Before handling the description of these collisions, several studies proved they do not occur in fluid structure systems. In [15], Vázquez and Zuazua prove no collision can occur between particles for a one-dimensional toy model. Then Starovoitov obtains a criterion for the velocity field of solutions [12] in the multidimensional setting. Namely, he proves no collision can occur if the gradient of the velocity field is sufficiently integrable. Finally, two parallel studies [8, 9] prove a no-collision result when there is only one body in a bounded (or partially bounded) two-dimensional cavity. In the first case, the author considers a rigid disk inside a bigger disk. In the second case, the author considers a rigid disk above a ramp. The aim of the present study is to extend these two-dimensional results to three-dimensional comparable configurations, i.e., for a rigid sphere above a ramp in \mathbb{R}^3 .

1.1. Mathematical model. We consider a homogeneous rigid sphere \mathcal{B} with radius 1 and density $\rho_{\mathcal{B}}$. We denote by \mathbf{G} its center (of mass), by \mathbf{V} (resp., $\boldsymbol{\omega}$) its translation (resp., angular) velocity, and by m (resp., \mathbb{J}) its mass (resp., inertia). Notice that $\boldsymbol{\omega}$ is a vector in \mathbb{R}^3 and $\mathbb{J} = J\mathbf{I}_3$, where $J \in (0, \infty)$, and \mathbf{I}_3 is the 3×3

*Received by the editors February 19, 2008; accepted for publication (in revised form) November 8, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/71607.html>

[†]Université de Toulouse, UPS, IMT, Laboratoire MIP, Université Paul Sabatier Toulouse 3, 31062 Toulouse Cedex 9, France (matthieu.hillairet@math.univ-toulouse.fr). This author was supported by ANR grant Jeunes Chercheurs “RUGO-Analyse et Calcul des Effets de Rugosité sur les Écoulements.”

[‡]Institut Elie Cartan UMR 7502, INRIA–Nancy–Université–CNRS, POB 239, Vandœuvre-lès-Nancy 54506, France, and Team-Project CORIDA, INRIA Nancy Grand-Est, 615 rue du Jardin Botanique, POB 54600, Villers-lès-Nancy, France (takeo8@gmail.com). This author was supported in part by ANR grants JCJC06_137283 and BLAN07-2_202879.

identity matrix. The velocity field of \mathcal{B} reads $\mathbf{V} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G})$ for all $\mathbf{x} \in \mathcal{B}$. The sphere evolves over a ramp \mathcal{P} . The remainder of the cavity \mathbb{R}_+^3 is denoted by \mathcal{F} . It contains an incompressible viscous and Newtonian fluid which does not slip on boundaries and has constant density $\rho_{\mathcal{F}} = 1$ and viscosity μ . The whole system evolves only through the interactions between solid and fluid without any external force field.

The evolution of the fluid is described by (\mathbf{u}, p) , a velocity/pressure field satisfying the incompressible Navier–Stokes equations:

$$(1.1) \quad \begin{cases} \partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u} = \operatorname{div} \mathbb{T}(\mathbf{u}, p) \\ \operatorname{div} \mathbf{u} = 0 \end{cases} \quad \text{in } \mathcal{F},$$

where $\mathbb{T}(\mathbf{u}, p)$ is the Newtonian stress tensor:

$$\mathbb{T}(\mathbf{u}, p) = 2\mu D(\mathbf{u}) - p\mathbb{I}_3.$$

Here $D(\mathbf{u})$ stands for the symmetric part of the gradient of \mathbf{u} .

To describe the evolution of \mathcal{B} , we apply Newton's laws assuming continuity of the stress-tensor of the fluid on $\partial\mathcal{B}$. It yields

$$(1.2) \quad \begin{cases} - \int_{\partial\mathcal{B}} \mathbb{T}(\mathbf{u}, p) \mathbf{n} \, d\sigma = m \dot{\mathbf{V}}, \\ - \int_{\partial\mathcal{B}} \mathbb{T}(\mathbf{u}, p) \mathbf{n} \times (\mathbf{x} - \mathbf{G}) \, d\sigma = J \dot{\boldsymbol{\omega}}. \end{cases}$$

Here \mathbf{n} stands for the normal to $\partial\mathcal{B}$ directed towards \mathcal{B} . As \mathbb{J} is a scalar matrix, the inertial term $\mathbb{J}\boldsymbol{\omega} \times \boldsymbol{\omega}$ vanishes in the conservation of momentum.

This system is complemented with the boundary conditions

$$(1.3) \quad \mathbf{u}|_{\partial\mathcal{B}} = \mathbf{V} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G}), \quad \mathbf{u}|_{\mathcal{P}} = 0, \quad \mathbf{u}|_{\infty} = 0$$

and initial conditions

$$(1.4) \quad \mathbf{u}(0, \cdot) = \mathbf{u}^0, \quad \mathbf{V}(0) = \mathbf{V}^0, \quad \boldsymbol{\omega}(0) = \boldsymbol{\omega}^0, \quad \mathbf{G}(0) = \mathbf{G}^0.$$

For short, we shall refer to the whole system (1.1)–(1.3) as (FSIS) for fluid-solid interaction system. We emphasize that this system is strongly coupled. On the one hand, the position of the sphere \mathcal{B} fixes the domain \mathcal{F} where the incompressible Navier–Stokes equations (1.1) have to be solved and the movement of \mathcal{B} fixes the boundary conditions for (1.1) on $\partial\mathcal{B}$. In particular, we lay stress upon these time-dependences, denoting by $\mathcal{F}(t)$ (resp., $\mathcal{B}(t)$) the domain occupied by the fluid (resp., the solid body) at time t in the following. We shall reserve the notation \mathcal{B} (resp., \mathcal{F}) for the sphere (resp., the fluid) as “actors” in the scenarios provided by our solutions to (FSIS). On the other hand, the solution (\mathbf{u}, p) prescribes the displacement of \mathcal{B} via the computation of the forces and torques applied to \mathcal{B} :

$$- \int_{\partial\mathcal{B}} \mathbb{T}(\mathbf{u}, p) \mathbf{n} \, d\sigma, \quad - \int_{\partial\mathcal{B}} \mathbb{T}(\mathbf{u}, p) \mathbf{n} \times (\mathbf{x} - \mathbf{G}) \, d\sigma.$$

Our main result is that no collision can occur between \mathcal{B} and \mathcal{P} in finite time in solutions to (FSIS). This reads as follows.

THEOREM 1.1. *Given $T > 0$, let (\mathbf{u}, \mathbf{G}) be a weak solution to (FSIS) over $(0, T)$ with initial data $(\mathbf{u}^0, \mathbf{G}^0)$. Then there exists a decreasing function $h_{\min} \in$*

$\mathcal{C}([0, T]; (0, \infty))$ depending only on initial data $(\mathbf{u}^0, \mathbf{G}^0)$ such that $h(t) := \text{dist}(\mathcal{B}(t), \mathcal{P})$ satisfies

$$h(t) \geq h_{\min}(t) \quad \forall t \in (0, T).$$

This result has been expected ever since the computations of Cooley and O'Neill [1] in the slow motion regime. However, until now no rigorous mathematical result has been available in the full nonlinear case. We emphasize that this result still holds true when one adds a reasonable external force field \mathbf{f} , for example, the gravity or $\mathbf{f} \in L^2((0, T) \times \mathbb{R}_+^3)$. In section 3, we provide an interpretation for the weak formulation of (FSIS) explaining how the distance can be estimated from below with a suitable test function. Then we construct a test function explicitly. In section 4, the interpretation of the weak formulation is applied to the constructed test function. Technical details are postponed to the appendix.

As mentioned in Theorem 1.1, our result applies to any weak solution to (FSIS). However, the Cauchy theory of weak solutions has been developed only in bounded and in exterior domains (to our knowledge) so that we do not know whether one solution exists to (FSIS). For the sake of completeness, we extend classical results for the Cauchy theory of (FSIS) to a half-space in the next section. Eventually, we obtain that, for a sufficiently large class of initial data $(\mathbf{u}^0, \mathbf{G}^0)$, the system (FSIS) predicts that no collision can occur between \mathcal{B} and \mathcal{P} .

1.2. Notation. Throughout, bold symbols stand for vectors. Given $\mathbf{a} \in \mathbb{R}^3$ we denote by $\mathbf{a} \otimes \mathbf{a}$ the symmetric matrix with entries $a_i a_j$. Coordinates $\mathbf{x} = (x_1, x_2, x_3)$ are centered on \mathcal{P} . For example, we have $\mathcal{P} := \{(x_1, x_2, 0), (x_1, x_2) \in \mathbb{R}^2\}$. The half-space above \mathcal{P} is \mathbb{R}_+^3 and \mathbb{R}_{++}^3 stands for $\{(x_1, x_2, x_3) \text{ with } x_3 > 1\}$. This is the domain where the center of mass \mathbf{G} can evolve as long as no collision between \mathcal{B} and \mathcal{P} occurs.

Given $\mathbf{G} \in \mathbb{R}^3$ and $\delta > 0$, we denote by $\mathcal{B}(\mathbf{G}, \delta)$ the sphere with center \mathbf{G} and radius δ . For short, we also set $\mathcal{B}_{\mathbf{G}} = \mathcal{B}(\mathbf{G}, 1)$. This is the domain occupied by \mathcal{B} when its center of mass meets \mathbf{G} . In this case, the fluid domain $\mathcal{F}_{\mathbf{G}}$ is the complementary of $\mathcal{B}_{\mathbf{G}}$ in \mathbb{R}_+^3 . If the orthogonal projection of \mathbf{G} on \mathcal{P} is the center of coordinates, we have $\mathbf{G} = \mathbf{G}_h = (0, 0, 1 + h)$ with $h = \text{dist}(\mathcal{B}_{\mathbf{G}}, \mathcal{P})$. In this case the suitable parameter is h and not \mathbf{G} . Thus, when using notation with h as subscript instead of \mathbf{G} , we implicitly mean that the subscript should be \mathbf{G}_h . For example, $\mathcal{B}_h := \mathcal{B}_{\mathbf{G}_h}$.

In the whole paper, we denote by $\eta : [0, \infty) \rightarrow [0, 1]$ a smooth function such that

$$\eta(s) = \begin{cases} 1 & \text{if } s < \frac{1}{2}, \\ 0 & \text{if } s > 1, \end{cases}$$

and we set $\eta_\alpha = \eta(\cdot/\alpha)$ for all parameters $\alpha > 0$.

We use the classical Lebesgue and Sobolev spaces $L^\alpha(A)$, $W^{\beta, \alpha}(A)$, $H^\beta(A)$ with A an open set, $\alpha \geq 1$, and $\beta \geq 0$. We define

$$\mathcal{H} = \{\phi \in L^2(\mathbb{R}_+^3) ; \text{div } \phi = 0, \phi \cdot \mathbf{n} = 0 \text{ on } \mathcal{P}\},$$

$$\mathcal{V} = \{\phi \in H^1(\mathbb{R}_+^3) ; \text{div } \phi = 0, \phi = 0 \text{ on } \mathcal{P}\}.$$

We recall that \mathcal{H} and \mathcal{V} are closed subspaces of $L^2(\mathbb{R}_+^3)$ and $H_0^1(\mathbb{R}_+^3)$, respectively. Thus, they form Hilbert spaces with respect to the induced inner products. For an open subset $A \subset \mathbb{R}_+^3$, we also consider the following Hilbert spaces:

$$\mathbb{H}(A) = \{\phi \in \mathcal{H} ; D(\phi) = 0 \text{ in } A\},$$

$$\mathbb{V}(A) = \{\phi \in \mathcal{V} ; D(\phi) = 0 \text{ in } A\}.$$

To simplify, if $\mathbf{G} \in \mathbb{R}_{++}^3$, we set

$$\mathbb{H}(\mathbf{G}) = \mathbb{H}(\mathcal{B}_{\mathbf{G}}), \quad \mathbb{V}(\mathbf{G}) = \mathbb{V}(\mathcal{B}_{\mathbf{G}}).$$

For all $\mathbf{G} \in \mathbb{R}_{++}^3$, we also denote by $\rho_{\mathbf{G}}$ the function

$$\rho_{\mathbf{G}}(\mathbf{x}) = \begin{cases} \rho_{\mathcal{B}} & \text{if } \mathbf{x} \in \mathcal{B}_{\mathbf{G}}, \\ 1 & \text{if } \mathbf{x} \in \mathcal{F}_{\mathbf{G}}. \end{cases}$$

If $\mathbf{v} \in \mathbb{H}(\mathbf{G})$, from [14, p. 18], there exists a unique pair $(\mathbf{V}[\mathbf{v}], \boldsymbol{\omega}[\mathbf{v}]) \in \mathbb{R}^3 \times \mathbb{R}^3$ such that

$$\mathbf{v}|_{\mathcal{B}_{\mathbf{G}}} = \mathbf{V}[\mathbf{v}] + \boldsymbol{\omega}[\mathbf{v}] \times (\mathbf{x} - \mathbf{G}).$$

In particular, if $(\mathbf{u}, \mathbf{v}) \in \mathbb{H}(\mathbf{G})^2$,

$$\int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{v} \, d\mathbf{x} = \int_{\mathbb{R}_+^3 \setminus \mathcal{B}_{\mathbf{G}}} \mathbf{u} \cdot \mathbf{v} \, d\mathbf{x} + m \mathbf{V}[\mathbf{u}] \cdot \mathbf{V}[\mathbf{v}] + J \boldsymbol{\omega}[\mathbf{u}] \cdot \boldsymbol{\omega}[\mathbf{v}].$$

2. Cauchy theory. First, we give the definition of weak solution to (FSIS).

DEFINITION 2.1. *Given $\mathbf{G}^0 \in \mathbb{R}_{++}^3$ and $\mathbf{u}^0 \in \mathbb{H}(\mathbf{G}^0)$, a pair (\mathbf{u}, \mathbf{G}) is called a weak solution to (FSIS) on $(0, T)$ with initial data $(\mathbf{u}^0, \mathbf{G}^0)$ if*

$$(2.1) \quad \mathbf{G} \in W^{1,\infty}(0, T; \mathbb{R}_{++}^3), \quad \text{with } \mathbf{G}(0) = \mathbf{G}^0,$$

$$(2.2) \quad \mathbf{u} \in L^\infty(0, T; \mathcal{H}) \cap L^2(0, T; \mathcal{V}),$$

$$(2.3) \quad \mathbf{u} = \mathbf{V} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G}) \quad \text{in } \mathcal{B}_{\mathbf{G}}, \quad \text{with } \mathbf{V} = \dot{\mathbf{G}};$$

if for all $\mathbf{v} \in \mathcal{C}([0, T]; H_0^1(\mathbb{R}_+^3)) \cap H^1(0, T; L^2(\mathbb{R}_+^3))$ with compact support in $(0, T) \times \mathbb{R}_+^3$ and such that $\mathbf{v} \in \mathbb{V}(\mathbf{G}(t))$ for all $t \in [0, T]$,

$$(2.4) \quad - \int_0^T \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \partial_t \mathbf{v} \, d\mathbf{y} \, dt + 2\mu \int_0^T \int_{\mathbb{R}_+^3} D(\mathbf{u}) : D(\mathbf{v}) \, d\mathbf{y} \, dt \\ - \int_0^T \int_{\mathbb{R}_+^3} \mathbf{u} \otimes \mathbf{u} : D(\mathbf{v}) \, d\mathbf{y} \, dt = 0;$$

if for all $\mathbf{v} \in \mathcal{C}([0, T]; L^2(\mathbb{R}_+^3))$ with compact support in $[0, T) \times \mathbb{R}_+^3$ and such that $\mathbf{v} \in \mathbb{H}(\mathbf{G}(t))$ for all $t \in [0, T]$ we have

$$(2.5) \quad W : t \mapsto \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{v} \, d\mathbf{x} \in \mathcal{C}([0, T]) \quad \text{with } W(0) = \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^0} \mathbf{u}^0 \cdot \mathbf{v} \, d\mathbf{x};$$

if the energy estimate holds true:

$$\frac{1}{2} \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} |\mathbf{u}|^2 \, d\mathbf{x} + 2\mu \int_0^t \int_{\mathbb{R}_+^3} |D(\mathbf{u})|^2 \, d\mathbf{x} \, ds \leq \frac{1}{2} \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^0} |\mathbf{u}^0|^2 \, d\mathbf{x} \quad \text{for a.a. } t \in (0, T).$$

Before going to our existence result for such weak solutions, let us recall some of their straightforward properties. First, combining (2.2) and (2.3) yields that $\mathbf{u}(t, \cdot) \in$

$\mathbb{V}(\mathbf{G}(t))$ for almost all $t \in (0, T)$. Moreover, it follows from standard arguments that the pair $(\mathbf{V}, \boldsymbol{\omega})$ such that (2.3) holds satisfies

$$(2.6) \quad |\mathbf{V}|^2 + |\boldsymbol{\omega}|^2 \leq C \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} |\mathbf{u}|^2 \, d\mathbf{x},$$

where C depends only on ρ_B . In particular, the pair $(\mathbf{V}, \boldsymbol{\omega})$ associated to a weak solution (\mathbf{u}, \mathbf{G}) belongs to $L^\infty(0, T)$.

The result of well-posedness we obtain for (FSIS) can be stated as follows.

THEOREM 2.2. *Assuming $\mathbf{G}^0 \in \mathbb{R}_{++}^3$ and $\mathbf{u}^0 \in \mathbb{H}(\mathbf{G}^0)$, there exists at least one maximal weak solution $(T_0, (\mathbf{U}, \mathbf{G}))$ to (FSIS) with initial data $(\mathbf{U}^0, \mathbf{G}^0)$. Moreover, we have the alternative:*

- $T_0 = \infty$,
- $T_0 < \infty$ and $G_3(t) \rightarrow 1$ as $t \rightarrow T_0$.

The proof of Theorem 2.2 given in what follows is inspired by methods developed in other papers, and since we use many similar arguments, we choose to present here only the main ideas and refer to the appropriate references to avoid repeating technical calculations which are not the main interest of this paper. From now on $(\mathbf{G}^0, \mathbf{u}^0) \in \mathbb{R}_{++}^3 \times \mathbb{H}(\mathbf{G}^0)$ is a fixed initial condition. We denote $(\mathbf{V}^0, \boldsymbol{\omega}^0) = (\mathbf{V}[\mathbf{u}^0], \boldsymbol{\omega}[\mathbf{u}^0])$ and $d^0 := \text{dist}(\mathcal{B}_{\mathbf{G}^0}, \mathcal{P})$. Due to our assumption, there holds $d^0 > 0$. The remainder of the section is devoted to the proof of Theorem 2.2.

2.1. Strong solutions for an approximate system. As in [7], we prove the existence of weak solutions by first obtaining strong solutions for an approximate problem of (FSIS). More precisely, we consider an even nonnegative function $\kappa \in C_0^\infty(\mathbb{R})$ such that $\kappa(s) = 0$ if $|s| \geq 1$. We define for all $\varepsilon > 0$

$$(2.7) \quad K_\varepsilon(\mathbf{x}) = \frac{c}{\varepsilon^3} \kappa\left(\left|\frac{\mathbf{x}}{\varepsilon}\right|^2\right) \quad (\mathbf{x} \in \mathbb{R}^3).$$

The constant c is chosen so that

$$\int_{\mathbb{R}^3} K_\varepsilon(\mathbf{x}) \, d\mathbf{x} = 1 \quad \forall \varepsilon > 0.$$

Then, for all $\mathbf{u} \in L^2((0, T) \times \mathbb{R}_+^3)$ and for all $\varepsilon > 0$, we set

$$\mathbf{u}_\varepsilon(t, \mathbf{x}) = \int_{\mathbb{R}_+^3} K_\varepsilon(\mathbf{x} - \mathbf{y}) \mathbf{u}(t, \mathbf{y}) \, d\mathbf{y}.$$

Let us consider the following problem, which approximates (FSIS):

$$(2.8) \quad \partial_t \mathbf{u} - \mu \Delta \mathbf{u} + (\mathbf{u}_\varepsilon \cdot \nabla) \mathbf{u} + \nabla p = 0 \quad \text{in } \mathcal{F}_{\mathbf{G}(t)}, \quad t \in (0, T),$$

$$(2.9) \quad \text{div } \mathbf{u} = 0 \quad \text{in } \mathcal{F}_{\mathbf{G}(t)}, \quad t \in (0, T),$$

$$(2.10) \quad \mathbf{u} = 0 \quad \text{on } \mathcal{P}, \quad t \in (0, T),$$

$$(2.11) \quad \mathbf{u}|_\infty = 0, \quad t \in (0, T),$$

$$(2.12) \quad \mathbf{u} = \dot{\mathbf{G}} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G}) \quad \text{on } \partial \mathcal{B}_{\mathbf{G}}, \quad t \in (0, T),$$

$$(2.13) \quad m\ddot{\mathbf{G}} = - \int_{\partial\mathcal{B}_{\mathbf{G}}} \mathbb{T}(\mathbf{u}, p)\mathbf{n} \, d\sigma + \frac{1}{2} \int_{\partial\mathcal{B}_{\mathbf{G}}} ((\mathbf{u}_\varepsilon - \mathbf{u}) \cdot \mathbf{n})\mathbf{u} \, d\sigma \quad \text{in } (0, T),$$

$$(2.14) \quad \begin{aligned} J\dot{\boldsymbol{\omega}} &= - \int_{\partial\mathcal{B}_{\mathbf{G}}} (\mathbf{x} - \mathbf{G}) \times \mathbb{T}(\mathbf{u}, p)\mathbf{n} \, d\sigma \\ &+ \frac{1}{2} \int_{\partial\mathcal{B}_{\mathbf{G}}} ((\mathbf{u}_\varepsilon - \mathbf{u}) \cdot \mathbf{n})(\mathbf{x} - \mathbf{G}) \times \mathbf{u} \, d\sigma \quad \text{in } (0, T). \end{aligned}$$

We complete the system with the initial conditions

$$(2.15) \quad \mathbf{u}(0, \cdot) = \mathbf{u}^0, \quad \dot{\mathbf{G}}(0) = \mathbf{V}^0, \quad \boldsymbol{\omega}(0) = \boldsymbol{\omega}^0, \quad \mathbf{G}(0) = \mathbf{G}^0.$$

We define the space $\widehat{H}^1(A)$ by

$$\widehat{H}^1(A) = \{q \in L^2_{loc}(\overline{A}) ; \nabla q \in L^2(A)\}.$$

We denote

$$\mathcal{F}_T = \{(t, \mathbf{x}) \in [0, T] \times \mathbb{R}^3 ; \mathbf{x} \in \mathcal{F}_{\mathbf{G}(t)}\}.$$

Consider a smooth mapping $\mathbf{X} : \mathbb{R}^3_{++} \times \mathcal{F}_{\mathbf{G}^0} \rightarrow \mathbb{R}^3$ such that for all $\mathbf{G} \in \mathbb{R}^3_{++}$, $\mathbf{X}(\mathbf{G}, \cdot)$ is a C^∞ -diffeomorphism from $\mathcal{F}_{\mathbf{G}^0}$ onto $\mathcal{F}_{\mathbf{G}}$. Moreover, suppose that the mappings

$$(\mathbf{G}, \mathbf{y}) \mapsto D_{\mathbf{G}}D_{\mathbf{y}}^\alpha \mathbf{X}(\mathbf{G}, \mathbf{y}), \quad \alpha \in \mathbb{N}^3,$$

exist and are continuous and compactly supported in $\mathcal{F}_{\mathbf{G}^0}$. For any $\mathbf{g} : \mathcal{F}_T \rightarrow \mathbb{R}^3$, we denote by $\mathbf{g}_{\mathbf{X}} : [0, T] \times \mathcal{F}_{\mathbf{G}^0} \rightarrow \mathbb{R}^3$ the mapping $\mathbf{g}_{\mathbf{X}}(t, \mathbf{y}) = \mathbf{g}(t, \mathbf{X}(\mathbf{G}(t), \mathbf{y}))$ for all $t \geq 0$ for all $\mathbf{y} \in \mathcal{F}_{\mathbf{G}^0}$. We use similar notation for $g : \mathcal{F}_T \rightarrow \mathbb{R}$.

We introduce the following function spaces in a variable domain:

$$\begin{aligned} L^2(0, T; H^2(\mathcal{F}(t))) &= \{\mathbf{u} ; \mathbf{u}_{\mathbf{X}} \in L^2(0, T; H^2(\mathcal{F}_{\mathbf{G}^0}))\}, \\ H^1(0, T; L^2(\mathcal{F}(t))) &= \{\mathbf{u} ; \mathbf{u}_{\mathbf{X}} \in H^1(0, T; L^2(\mathcal{F}_{\mathbf{G}^0}))\}, \\ C([0, T], H^1(\mathcal{F}(t))) &= \{\mathbf{u} ; \mathbf{u}_{\mathbf{X}} \in C([0, T], H^1(\mathcal{F}_{\mathbf{G}^0}))\}, \\ L^2(0, T; \widehat{H}^1(\mathcal{F}(t))) &= \{p ; p_{\mathbf{X}} \in L^2(0, T; \widehat{H}^1(\mathcal{F}_{\mathbf{G}^0}))\}. \end{aligned}$$

THEOREM 2.3. *Assume the initial conditions satisfy*

$$\text{dist}(\mathcal{B}_{\mathbf{G}^0}, \mathcal{P}) = d^0 > 0, \quad \mathbf{u}^0 \in \mathbb{V}(\mathbf{G}^0).$$

Then, given $\varepsilon > 0$, there exist a time $T > 0$ depending only on $\|\mathbf{u}^0\|_{L^2(\mathbb{R}^3_+)}$ and a 4-uplet $(\mathbf{u}, \mathbf{G}, \boldsymbol{\omega}, p)$ satisfying

$$(2.16) \quad \mathbf{G} \in H^2(0, T), \quad \text{dist}(\mathcal{B}_{\mathbf{G}(t)}, \mathcal{P}) \geq \frac{d^0}{2} > 0 \quad \forall t \in [0, T],$$

$$(2.17) \quad \mathbf{u} \in L^2(0, T; H^2(\mathcal{F}(t))) \cap C([0, T]; H^1(\mathcal{F}(t))) \cap H^1(0, T; L^2(\mathcal{F}(t))),$$

$$(2.18) \quad p \in L^2(0, T; \widehat{H}^1(\mathcal{F}(t))), \quad \boldsymbol{\omega} \in H^1(0, T),$$

and satisfying (2.8)–(2.14) almost everywhere on $[0, T]$ or in the trace sense.

Proof. One can obtain local existence of strong solutions to (2.8)–(2.14) via arguments similar to those in the proof of Theorem 1.1 in [2] (see also [3, 13, 10] for

results obtained applying similar techniques). For the sake of brevity, we compute only energy estimates here in order to prove that these local strong solutions can be continued up to collision between \mathcal{B} and \mathcal{P} . We refer the reader to the mentioned articles for technical details.

First, we multiply (2.8) by \mathbf{u} , (2.13) by $\dot{\mathbf{G}}$, and (2.14) by $\boldsymbol{\omega}$. We deduce the energy estimate

$$(2.19) \quad \frac{1}{2} \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}(t)} |\mathbf{u}|^2(t) \, d\mathbf{x} + 2\mu \int_0^t \int_{\mathbb{R}_+^3} |D(\mathbf{u})|^2 \, d\mathbf{x} \, ds = \frac{1}{2} \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^0} |\mathbf{u}^0|^2 \, d\mathbf{x}.$$

From the above estimate and (2.6) we obtain a time $T > 0$ depending on $\int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^0} |\mathbf{u}^0|^2 \, d\mathbf{x}$ such that (2.16) holds for all ε .

Then we introduce a smooth function Υ with compact support in $\mathcal{B}(\mathbf{G}^0, 1 + d_0/4)$ and such that $\Upsilon = 1$ on $\mathcal{B}_{\mathbf{G}^0}$. We also set

$$\begin{aligned} \mathbf{w}_R &= \frac{1}{2} \left(\dot{\mathbf{G}} \times (\mathbf{x} - \mathbf{G}) + |\mathbf{x} - \mathbf{G}|^2 \boldsymbol{\omega} \right), \\ \mathbf{u}_R(t, \mathbf{x}) &= \operatorname{curl} \left[\Upsilon (\mathbf{x} - \mathbf{G}(t) + \mathbf{G}^0) \mathbf{w}_R(t, \mathbf{x}) \right]. \end{aligned}$$

The function \mathbf{u}_R has a compact support in $\mathcal{B}(\mathbf{G}^0, 1 + d_0/4)$ and satisfies

$$\operatorname{div} \mathbf{u}_R = 0, \quad \mathbf{u}_R = \dot{\mathbf{G}} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G}) \text{ on } \mathcal{B}(t).$$

We multiply (2.8) by

$$\boldsymbol{\varphi} = \partial_t \mathbf{u} + (\mathbf{u}_R \cdot \nabla) \mathbf{u} - (\mathbf{u} \cdot \nabla) \mathbf{u}_R,$$

which yields

$$(2.20) \quad \int_{\mathcal{F}(t)} (\partial_t \mathbf{u} + \mathbf{u}_\varepsilon \cdot \nabla \mathbf{u}) \cdot \boldsymbol{\varphi} \, d\mathbf{x} = \int_{\mathcal{F}(t)} \operatorname{div} \mathbb{T}(\mathbf{u}, p) \cdot \boldsymbol{\varphi} \, d\mathbf{x}.$$

Using the regularization of the nonlinear term, we obtain the existence of a constant C_0 depending only on $\|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}$ such that the left-hand side of (2.20) can be written as

$$LHS = \int_{\mathcal{F}(t)} |\partial_t \mathbf{u}|^2 \, d\mathbf{x} + R_l \quad \text{with} \quad |R_l| \leq C_0 \int_{\mathcal{F}(t)} |\nabla \mathbf{u}|^2 \, d\mathbf{x} + \frac{1}{2} \int_{\mathcal{F}(t)} |\partial_t \mathbf{u}|^2 \, d\mathbf{x}.$$

Concerning the right-hand side of (2.20), we apply the same arguments as in the proof of Lemma 4.3 in [3] and obtain (see [3, (4.24)])

$$RHS = -\mu \frac{d}{dt} \int_{\mathcal{F}(t)} |D(\mathbf{u})|^2 \, d\mathbf{x} + I_r + R_r,$$

where

$$I_r = \int_{\partial \mathcal{B}(t)} \mathbb{T}(\mathbf{u}, p) \mathbf{n} \cdot \left[\ddot{\mathbf{G}} + \dot{\boldsymbol{\omega}} \times (\mathbf{x} - \mathbf{G}) - \boldsymbol{\omega} \times \dot{\mathbf{G}} \right] \, d\sigma,$$

and, with the same convention as previously for C_0 , we obtain

$$|R_r| \leq C_0 \left[1 + \int_{\mathcal{F}(t)} |\nabla \mathbf{u}|^2 \, d\mathbf{x} \right].$$

Finally, using (2.13)–(2.14), we deduce

$$I_r = - \left[m|\ddot{\mathbf{G}}|^2 + J|\dot{\boldsymbol{\omega}}|^2 - m\ddot{\mathbf{G}} \cdot (\boldsymbol{\omega} \times \dot{\mathbf{G}}) \right] + \tilde{R}_r,$$

with

$$\tilde{R}_r = \frac{1}{2} \int_{\partial B(t)} (\mathbf{u} - \mathbf{u}_\varepsilon) \cdot \mathbf{n} \left(\dot{\mathbf{G}} + \boldsymbol{\omega} \times (\mathbf{x} - \mathbf{G}) \right) \cdot \left(\ddot{\mathbf{G}} - \boldsymbol{\omega} \times \dot{\mathbf{G}} + \dot{\boldsymbol{\omega}} \times (\mathbf{x} - \mathbf{G}) \right) \, d\sigma.$$

In this last integral, extending the rigid velocity fields to $\mathcal{B}(\mathbf{G}, 1 + d_0/4)$ in a similar fashion to that in \mathbf{u}_R , we obtain, after integration by parts,

$$|\tilde{R}_r| \leq C_0 + \frac{m}{2} |\ddot{\mathbf{G}}|^2 + \frac{J}{2} |\dot{\boldsymbol{\omega}}|^2.$$

Combining all these estimates, we finally deduce

$$\begin{aligned} & \mu \frac{d}{dt} \left[\int_{\mathcal{F}(t)} |\nabla \mathbf{u}|^2 \, d\mathbf{x} \right] + \frac{1}{2} \left[m|\ddot{\mathbf{G}}|^2 + J|\dot{\boldsymbol{\omega}}|^2 + \int_{\mathcal{F}(t)} |\partial_t \mathbf{u}|^2 \, d\mathbf{x} \right] \\ & \leq C_0 \left[1 + \int_{\mathbb{R}_+^3} |\nabla \mathbf{u}|^2 \, d\mathbf{x} \right]. \end{aligned}$$

As a consequence, we have obtained that the mapping $t \mapsto \|\mathbf{u}\|_{H^1(\mathcal{F}(t))}$ is bounded on $[0, T]$. \square

2.2. Convergences. As in the assumptions of our theorem, we assume that the initial condition \mathbf{u}^0 belongs to $\mathbb{H}(\mathbf{G}^0)$; we introduce a sequence $\mathbf{u}^{0k} \in \mathbb{V}(\mathbf{G}^0)$ such that

$$\mathbf{u}^{0k} \rightarrow \mathbf{u}^0 \quad \text{in } \mathbb{H}(\mathbf{G}^0).$$

We also take a sequence $\varepsilon^k \rightarrow 0$. Then, applying Theorem 2.3, we can consider a (uniform) time T such that, for all k , the corresponding solutions $(\mathbf{u}^k, \mathbf{G}^k, \boldsymbol{\omega}^k, p^k)$ exist on $[0, T]$ and satisfy

$$\text{dist}(\mathcal{B}_{\mathbf{G}^k(t)}, \mathcal{P}) > \frac{d^0}{2} \quad \forall t \in [0, T], \forall k.$$

In the following, we extend \mathbf{u}^k with the value of the rigid velocity field $\mathbf{V}^k + \boldsymbol{\omega}^k \times (\mathbf{x} - \mathbf{G}^k)$ on the solid domain. From (2.19), the following holds:

$$\mathbf{u}^k \text{ is bounded in } L^\infty(0, T; L^2(\mathbb{R}_+^3)) \cap L^2(0, T; H_0^1(\mathbb{R}_+^3)),$$

so that, up to a subsequence (which we do not relabel), we can assume

$$(2.21) \quad \mathbf{u}^k \rightharpoonup \mathbf{u} \quad \text{in } L^2(0, T; H^1(\mathbb{R}_+^3))\text{-weak and } L^\infty(0, T; L^2(\mathbb{R}_+^3))\text{-weak}^*,$$

$$(2.22) \quad \mathbf{G}^k \rightarrow \mathbf{G} \quad \text{in } C([0, T]; \mathbb{R}_{++}^3).$$

Taking any $\mathbf{U} \in \mathcal{D}((0, T) \times \mathbb{R}_+^3)$ such that $D(\mathbf{U}) = 0$ in a neighborhood of $\mathcal{B}_{\mathbf{G}^k(t)}$ for all $t \in (0, T)$, we can multiply (2.8) by \mathbf{U} for k sufficiently large. Integrating by

parts and using Reynolds' transport theorem yields

$$\begin{aligned}
 (2.23) \quad & - \int_0^T \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \partial_t \mathbf{U} \, d\mathbf{x} \, ds \\
 & + \int_0^T \int_{\mathbb{R}_+^3} (\mathbf{u}^k \otimes \mathbf{u}^k) : D(\mathbf{U}) \, d\mathbf{x} \, ds + 2\mu \int_0^T \int_{\mathbb{R}_+^3} D(\mathbf{u}^k) : D(\mathbf{U}) \, d\mathbf{x} \, ds \\
 = & \int_0^T \int_{\mathcal{F}_{\mathbf{G}^k(t)}} [((\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k) \cdot \nabla) \mathbf{U}] \cdot \mathbf{u}^k \, d\mathbf{x} - \frac{1}{2} \int_0^T \int_{\partial \mathcal{B}_{\mathbf{G}^k(t)}} (\mathbf{U} \cdot \mathbf{u}^k) ((\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k) \cdot \mathbf{n}) \, d\sigma.
 \end{aligned}$$

In order to pass to the limit in this weak formulation, we need to prove L^2 -compactness of the \mathbf{u}^k . As usual, this is the main difficulty of the proof. The procedure to prove this compactness property follows closely the method developed in [7]. The main difference here is that our cavity is unbounded. Therefore, we do not look for a compactness property of the sequence \mathbf{u}^k on the whole domain but only locally in space. Indeed, with the help of Friedrichs' lemma (see [6, Lemma II.4.2]) we have the following: for all relatively compact $\mathcal{O} \subset \mathbb{R}_+^3$, for any $\gamma > 0$ there exist $I = I(\gamma, \mathcal{O}) \in \mathbb{N}$ and functions $\psi_j \in L^\infty(\mathcal{O})$, $j = 1, \dots, I$, such that

$$\begin{aligned}
 (2.24) \quad & \|\mathbf{u}^k - \mathbf{u}\|_{L^2(0,T;L^2(\mathcal{O}))}^2 \\
 & \leq \sum_{j=1}^I \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}} (\mathbf{u}^k(t) - \mathbf{u}(t)) \cdot \psi_j \, d\mathbf{y} \right)^2 dt + \gamma \|\nabla \mathbf{u}^k - \nabla \mathbf{u}\|_{L^2(0,T;L^2(\mathcal{O}))}^2.
 \end{aligned}$$

Due to the uniform bound on \mathbf{u}^k in $L^2(0, T; H_0^1(\mathbb{R}_+^3))$, our remaining task is to prove that, for any $\psi \in L^\infty(\mathcal{O})$, there exists a subsequence (which we do not relabel) for which there holds

$$(2.25) \quad \lim_{k \rightarrow \infty} \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}} (\mathbf{u}^k(t) - \mathbf{u}(t)) \cdot \psi \, d\mathbf{y} \right)^2 dt = 0.$$

As a first step, we obtain results similar to those of (2.25) for another family of test functions (not included in $L^\infty(\mathcal{O})$). To this end, we divide the segment $[0, T]$ into N segments $[t_{i-1}, t_i]$, with $\Delta t = t_i - t_{i-1} = T/N$, $i = 1, \dots, N$. For all i and for $\delta < \frac{d_0}{2}$, we consider an orthonormal basis $(\mathbf{e}_j^{i,\delta})$ of $\mathbb{V}(\mathcal{B}(\mathbf{G}(t_i), 1 + \delta))$. Without further restrictions, we assume all the $\mathbf{e}_j^{i,\delta}$ with compact support. We also consider the set of piecewise linear functions in t :

$$(2.26) \quad \mathbf{U}^\delta(t, \mathbf{x}) = \mathbf{e}_j^{i-1,\delta}(\mathbf{x}) + \frac{t - t_i}{\Delta t} \left(\mathbf{e}_l^{i,\delta}(\mathbf{x}) - \mathbf{e}_j^{i-1,\delta}(\mathbf{x}) \right)$$

for $t \in [t_{i-1}, t_i]$, $j, l \in \mathbb{N}$, $i \in \{1, \dots, N\}$. There exists a countable set of functions satisfying (2.26).

It is worth noting that, since \mathbf{G} is uniformly continuous, for N big enough, the functions \mathbf{U}^δ of the above form satisfy $D(\mathbf{U}^\delta(t)) = 0$ in $\mathcal{B}(\mathbf{G}(t), 1 + \delta/2)$ for all t . In particular, these functions are in $C([0, T]; \mathbb{V}(\mathbf{G}(t)))$. These functions are important to approximate continuous functions in time with value in $\mathbb{H}(\mathbf{G}(t))$ and thus functions in $L^2(0, T; \mathbb{H}(\mathbf{G}(t)))$ (see Lemmas 4.1 and 4.2 in [7]). From (2.22), there exists $k_0 = k_0(\delta)$ such that for $k \geq k_0$,

$$D(\mathbf{U}^\delta) = 0 \quad \text{in } \mathcal{B}_{\mathbf{G}^k(t)} \quad \forall t \in [0, T]$$

for all functions satisfying (2.26).

We multiply (2.8) by \mathbf{U}^δ satisfying (2.26). After similar computations to those leading to (2.23), we obtain

$$(2.27) \quad \frac{d}{dt} \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \mathbf{U}^\delta \, d\mathbf{x} = -2\mu \int_{\mathcal{F}^k(t)} D(\mathbf{u}^k) : D(\mathbf{U}^\delta) \, d\mathbf{x} \\ + \int_{\mathcal{F}^k(t)} (\partial_t \mathbf{U}^\delta + (\mathbf{u}_\varepsilon^k \cdot \nabla) \mathbf{U}^\delta) \cdot \mathbf{u}^k \, d\mathbf{x} - \frac{1}{2} \int_{\partial \mathcal{B}_{\mathbf{G}^k(t)}} (\mathbf{U}^\delta \cdot \mathbf{u}^k) ((\mathbf{u}_\varepsilon^k - \mathbf{u}^k) \cdot \mathbf{n}) \, d\sigma.$$

Following the estimates of [7], using Arzelà and Ascoli and the diagonal Cantor procedure, we obtain that for all \mathbf{U}^δ satisfying (2.26), we have

$$(2.28) \quad \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \mathbf{U}^\delta \, d\mathbf{x} \rightarrow \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{U}^\delta \, d\mathbf{x} \quad \text{in } C([0, T]).$$

However, as $\mathbf{G}^k \rightarrow \mathbf{G}$ in $\mathcal{C}([0, T])$, and \mathbf{u}^k is bounded in $L^\infty(0, T; L^2(\mathbb{R}_+^3)) \cap L^2(0, T; H_0^1(\mathbb{R}_+^3))$, this also leads to

$$(2.29) \quad \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \mathbf{U}^\delta \, d\mathbf{x} \rightarrow \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{U}^\delta \, d\mathbf{x} \quad \text{in } C([0, T]).$$

Via a density argument (see Lemmas 4.1 and 4.2 in [7] for details), we might extend this convergence result to all $\mathbf{U} \in C([0, T]; L^2(\mathbb{R}_+^3))$, $\operatorname{div} \mathbf{U} = 0$, $D(\mathbf{U}) = 0$ in $\mathcal{B}(\mathbf{G}(t), 1)$,

$$(2.30) \quad \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \mathbf{U} \, d\mathbf{x} \rightarrow \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{U} \, d\mathbf{x} \quad \text{in } C([0, T]),$$

and to all $\mathbf{U} \in L^2(0, T; L^2(\mathbb{R}_+^3))$, $\operatorname{div} \mathbf{U} = 0$, $D(\mathbf{U}) = 0$ in $\mathcal{B}(\mathbf{G}(t), 1)$,

$$(2.31) \quad \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}^k} \mathbf{u}^k \cdot \mathbf{U} \, d\mathbf{x} \rightarrow \int_{\mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \mathbf{U} \, d\mathbf{x} \quad \text{in } L^2(0, T).$$

Relation (2.30) implies in particular that \mathbf{u} satisfies the initial conditions.

However, in (2.25), no assumption is made on the velocity field ψ over $\mathcal{B}_{\mathbf{G}}$. In order to reduce (2.25) to the above convergence result, we need to project such $\psi \in L^\infty(\mathcal{O})$ on velocity fields which are rigid over $\mathcal{B}_{\mathbf{G}}$. Similarly, \mathbf{u}^k is rigid on $\mathcal{B}_{\mathbf{G}^k}$, so we also modify \mathbf{u}^k slightly to obtain a function with the same properties but which is rigid in $\mathcal{B}_{\mathbf{G}}$. To this end, we extend \mathbf{u}^k by 0 in $\mathbb{R}^3 \setminus \mathbb{R}_+^3$, and since $\mathbf{u}^k = 0$ on \mathcal{P} , we have $\mathbf{u}^k \in L^2(0, T; H^1(\mathbb{R}^3))$. We set

$$\widehat{\mathbf{u}}^k(t, \mathbf{y}) = \mathbf{u}^k(t, \mathbf{y} + \mathbf{G}^k(t) - \mathbf{G}(t)).$$

We have

$$\operatorname{div} \widehat{\mathbf{u}}^k = 0, \quad \widehat{\mathbf{u}}^k = 0 \quad \text{on } \{\mathbf{y} \in \mathbb{R}^3; y_3 = G_3 - G_3^k\}, \quad D(\mathbf{u}^k) = 0 \quad \text{in } \mathcal{B}_{\mathbf{G}}$$

and

$$(2.32) \quad \|\widehat{\mathbf{u}}^k - \mathbf{u}^k\|_{L^2((0, T) \times \mathbb{R}^3)} \leq \|\mathbf{G}^k - \mathbf{G}\|_{L^\infty(0, T)} \|\nabla \mathbf{u}^k\|_{L^2((0, T) \times \mathbb{R}^3)} \rightarrow 0.$$

We define

$$\mathcal{L}^2(\mathbb{R}^3) = \{\mathbf{v} \in L^2(\mathbb{R}^3); \operatorname{div} \mathbf{v} = 0 \quad \text{in } \mathbb{R}^3\}.$$

Following [7] (see (4.37) in [7]), there exists a function $\Lambda : \mathcal{L}^2(\mathbb{R}^3) \rightarrow \mathcal{L}^2(\mathbb{R}^3)$ such that for all $\mathbf{v} \in \mathcal{L}^2(\mathbb{R}^3)$, $\mathbf{u} = \Lambda(\mathbf{v})$ satisfies

$$\mathbf{u} = \mathbf{v} \quad \text{in } \mathcal{B}_{\mathbf{G}^0}, \quad \mathbf{u} = 0 \quad \text{in } \mathbb{R}^3 \setminus \overline{\mathcal{B}\left(\mathbf{G}^0, 1 + \frac{d^0}{4}\right)},$$

$$\|\mathbf{u}\|_{L^2(\mathbb{R}^3)} \leq c\|\mathbf{v}\|_{L^2(\mathbb{R}^3)}.$$

Moreover,

$$\text{if } \mathbf{v} \in C([0, T]; L^2(\mathbb{R}^3)), \quad \text{then } \Lambda\mathbf{v} \in C([0, T]; L^2(\mathbb{R}^3)),$$

and

$$\text{if } \mathbf{v} \in L^2(0, T; L^2(\mathbb{R}^3)), \quad \text{then } \Lambda\mathbf{v} \in L^2(0, T; L^2(\mathbb{R}^3)).$$

We also define $\check{\Lambda}^t : \mathcal{L}^2(\mathbb{R}^3) \rightarrow \mathcal{L}^2(\mathbb{R}^3)$ as follows:

$$\check{\Lambda}^t \check{\mathbf{v}}(\mathbf{y}) = [\Lambda \check{\mathbf{v}}](\mathbf{y} + \mathbf{G}^0 - \mathbf{G}(t)),$$

where

$$\check{\mathbf{v}}(\mathbf{x}) = \mathbf{v}(\mathbf{x} + \mathbf{G}(t) - \mathbf{G}^0).$$

Then, as in [7], we consider

$$\bar{\mathbf{u}}^k(\mathbf{x}, t) = \mathbf{u}^k + \check{\Lambda}^t(\hat{\mathbf{u}}^k - \mathbf{u}^k).$$

This function is rigid in $\mathcal{B}_{\mathbf{G}(t)}$ and

$$(2.33) \quad \|\bar{\mathbf{u}}^k - \mathbf{u}^k\|_{L^2((0, T) \times \mathbb{R}^3)} \rightarrow 0.$$

We are now in a position to prove (2.25) by using (2.31). Assume that $\psi \in L^\infty(\mathcal{O})$. We set

$$I_k := \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}}(\mathbf{u}^k(t) - \mathbf{u}(t)) \cdot \psi \, d\mathbf{y} \right)^2 dt$$

and

$$\bar{I}_k = \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}}(\bar{\mathbf{u}}^k(t) - \mathbf{u}(t)) \cdot \psi \, d\mathbf{y} \right)^2 dt.$$

From (2.33), to prove that I_k goes to 0, it is sufficient to prove that \bar{I}_k goes to 0. Then, given $\mathbf{G} \in \mathbb{R}_+^3$ such that $\text{dist}(\mathcal{B}_{\mathbf{G}}, \mathcal{P}) > 0$, we denote by $P(\mathbf{G})$ the orthogonal projection from $L^2(\mathbb{R}_+^3, \rho_{\mathbf{G}} d\mathbf{x})$ onto $\mathbb{H}(\mathbf{G})$ and introduce

$$\tilde{\psi}(t, \mathbf{x}) = P(\mathbf{G}(t))\psi(\mathbf{x}).$$

We notice that $\tilde{\psi} \in L^\infty(0, T; L^2(\mathbb{R}_+^3))$, $\text{div } \tilde{\psi} = 0$, $D(\tilde{\psi}) = 0$ in $\mathcal{B}(\mathbf{G}(t), 1)$ and

$$\bar{I}_k := \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}}(\bar{\mathbf{u}}^k(t) - \mathbf{u}(t)) \cdot \tilde{\psi} \, d\mathbf{y} \right)^2 dt.$$

We set

$$\tilde{I}_k := \int_0^T \left(\int_{\mathcal{O}} \rho_{\mathbf{G}} (\mathbf{u}^k(t) - \mathbf{u}(t)) \cdot \tilde{\boldsymbol{\psi}} \, dy \right)^2 dt$$

and notice that, as before, if \tilde{I}_k goes to 0, then \bar{I}_k goes to 0. However, since $\tilde{\boldsymbol{\psi}} \in L^\infty(0, T; L^2(\mathbb{R}_+^3))$ we can take $\mathbf{U} = \tilde{\boldsymbol{\psi}}$ in (2.31). Consequently, we obtain that there exists a subsequence we do not relabel such that $\bar{I}_k \rightarrow 0$. This ends the proof of L^2 -compactness of the sequence \mathbf{u}^k . In particular we conclude from (2.24) that

$$\|\mathbf{u}^k - \mathbf{u}\|_{L^2(0, T; L^2(\mathcal{O}))} \rightarrow 0.$$

Combining (2.33) and the above relation, we obtain that

$$\|\bar{\mathbf{u}}^k - \mathbf{u}\|_{L^2(0, T; L^2(\mathcal{B}_{\mathbf{G}}))} \rightarrow 0.$$

This implies that

$$\dot{\mathbf{G}}^k \rightarrow \dot{\mathbf{G}} \quad \text{and} \quad \boldsymbol{\omega}^k \rightarrow \boldsymbol{\omega} \quad \text{in } L^2(0, T).$$

Using the above compactness, we can pass to the limit in (2.23) in the first three terms. For the two terms of the last line of (2.23), we proceed as follows. First, we can notice that

$$\int_{\mathcal{F}_{\mathbf{G}^k(t)}} |\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k|^2 \, d\mathbf{x} \leq C\varepsilon^k \|\mathbf{u}^k\|_{H^1(\mathcal{F}_{\mathbf{G}^k(t)})}^2$$

and thus

$$\int_0^T \int_{\mathcal{F}_{\mathbf{G}^k(t)}} [((\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k) \cdot \nabla) \mathbf{U}] \cdot \mathbf{u}^k \, d\mathbf{x} \rightarrow 0,$$

as $k \rightarrow \infty$.

Second, from the choice of κ and K_ε (see (2.7)), we can easily check that

$$\mathbf{u}_{\varepsilon^k}^k = \dot{\mathbf{G}}^k + \boldsymbol{\omega}^k \times (\mathbf{x} - \mathbf{G}^k) \quad \text{in } \mathcal{B}(\mathbf{G}^k(t), 1 - \varepsilon^k).$$

As a consequence, using Lemma 4.10 in [5], we conclude that

$$\int_{\partial \mathcal{B}_{\mathbf{G}^k(t)}} |\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k|^2 \, d\sigma \leq C\varepsilon^k \int_{\mathcal{B}_{\mathbf{G}^k(t)}} |\nabla \mathbf{u}_{\varepsilon^k}^k - \nabla \mathbf{u}^k|^2 \, d\mathbf{x}.$$

The above inequality implies that

$$-\frac{1}{2} \int_0^T \int_{\partial \mathcal{B}_{\mathbf{G}^k(t)}} (\mathbf{U} \cdot \mathbf{u}^k) ((\mathbf{u}_{\varepsilon^k}^k - \mathbf{u}^k) \cdot \mathbf{n}) \, d\sigma \rightarrow 0,$$

as $k \rightarrow \infty$.

Finally, we can pass to the limit in (2.23) and obtain the weak formulation (2.4) for smooth test functions. We can pass from smooth test functions to the required regularity for \mathbf{v} by applying the same approximation technique as when we obtained (2.29).

3. Constructing test functions. Let us begin with some notation. We introduce (r, θ, z) , the cylindrical coordinates associated to (x_1, x_2, x_3) :

$$x_1 = r \cos(\theta), \quad x_2 = r \sin(\theta), \quad x_3 = z.$$

Given $h > 0$ and $l > 0$, we denote by $\Omega_{h,l}$ the cylindric domain under \mathcal{B}_h with radius l :

$$(3.1) \quad \Omega_{h,l} := \{(r, \theta, z) \in \mathcal{F}_h \text{ such that } r \in [0, l], z \in (0, 1 + h)\}.$$

We notice that whenever $l < 1$, the upper boundary of $\Omega_{h,\delta}$ is parametrized by (r, θ) :

$$(r, \theta, z) \in \partial\Omega_{h,l} \cap \partial\mathcal{B}_h \Leftrightarrow \{r \in [0, \delta] \text{ and } z = \delta_h(r)\},$$

where, for arbitrary nonnegative h ,

$$(3.2) \quad \delta_h(s) := 1 + h - \sqrt{1 - s^2} \quad \forall s \in [0, 1].$$

As in [9], we estimate the distance between \mathcal{B} and \mathcal{P} from below with a suitable choice of test function in the weak formulation. To this end, we introduce an approximation of the Stokes solution for a given position of \mathcal{B} in \mathbb{R}_+^3 (namely, \mathcal{B}_h). We call these approximations “static functions” and denote them by $(\mathbf{w}_h)_{h>0}$. Given a weak solution (\mathbf{u}, \mathbf{G}) to (FSIS) in $(0, T)$, we construct admissible test functions by setting

$$(3.3) \quad \begin{aligned} \tilde{\mathbf{w}} : (0, T) \times \mathbb{R}_+^3 &\longrightarrow \mathbb{R}^3, \\ (t, \mathbf{x}) &\longmapsto \zeta(t) \mathbf{w}_{h(t)}(x_1 - G_1(t), x_2 - G_2(t), x_3) \end{aligned}$$

for arbitrary $\zeta \in \mathcal{D}(0, T)$. In this definition $h(t)$ stands for the distance between the sphere and the ramp at time t .

Applying the weak formulation, we obtain

$$\int_0^T \int_{\mathbb{R}_+^3} [\rho \mathbf{G} \mathbf{u} \cdot \tilde{\mathbf{w}}_t + (\mathbf{u} \otimes \mathbf{u} - 2\mu \mathbf{D}(\mathbf{u})) : \mathbf{D}(\tilde{\mathbf{w}})] \, \mathbf{d}\mathbf{x} \, dt = 0.$$

In this equation, the key ingredient is

$$\int_{\mathbb{R}_+^3} \mathbf{D}(\mathbf{u}) : \mathbf{D}(\mathbf{w}_h) \, \mathbf{d}\mathbf{x}.$$

It shall behave like \dot{h}/h^α with an exponent α to be made precise. The other terms appear as remainders. We shall bound them by an integrable (in time) function. This relies on the following lemma.

LEMMA 3.1. *Given $h > 0$, $r_0 > 0$, and $(\mathbf{u}, \mathbf{w}) \in H_0^1(\mathbb{R}_+^3) \times (\mathbb{H}(\mathbf{G}_h) \cap \mathcal{C}^\infty(\mathcal{F}_h))$, we assume \mathbf{w} is with compact support. Then there exists C depending only on the size of the support of \mathbf{w} such that*

$$(3.4) \quad \left| \int_{\mathbb{R}_+^3} \mathbf{u} \cdot \mathbf{w} \, \mathbf{d}\mathbf{x} \right| \leq C \|\nabla \mathbf{u}\|_{L^2(\mathbb{R}_+^3)} \left[\|\mathbf{w}\|_{2,2} + \|\mathbf{w}\|_{L^2(\mathbb{R}_+^3 \setminus \Omega_{h,r_0})} \right],$$

where

$$\|\mathbf{w}\|_{2,2}^2 = \int_0^{r_0} \left(\delta_h(r)^2 \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0, 2\pi)} \{|\mathbf{w}(r, \theta, z)|^2\} \, dz \right] \right) r \, dr.$$

If, moreover, $\mathbf{w} \in \mathbb{V}(\mathbf{G}_h)$, we have

$$(3.5) \quad \left| \int_{\mathbb{R}_+^3} \mathbf{u} \otimes \mathbf{u} : D(\mathbf{w}) \, d\mathbf{x} \right| \leq C \|\nabla \mathbf{u}\|_{L^2(\mathbb{R}_+^3)}^2 \left[\|\mathbf{w}\|_{\infty,2} + \|D(\mathbf{w})\|_{L^\infty(\mathcal{F}_h \setminus \Omega_{h,r_0})} \right],$$

where

$$\|\mathbf{w}\|_{\infty,2} = \sup_{r \in (0,r_0)} \left\{ \delta_h(r)^{\frac{3}{2}} \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0,2\pi)} \{ |\nabla \mathbf{w}(r, \theta, z)|^2 \} \, dz \right]^{\frac{1}{2}} \right\}.$$

Proof. We denote by I_1 and I_2 the two integrals we want to estimate in (3.5) and (3.4).

We first deal with I_1 . As $D(\mathbf{w}) = 0$ in \mathcal{B}_h , we might restrict the integration domain to \mathcal{F}_h . We split the integral into an integral in $\mathcal{F}_h \setminus \Omega_{h,r_0}$ and an integral in Ω_{h,r_0} : $I_1 = I_1^{in} + I_1^{out}$ with

$$|I_1^{out}| = \left| \int_{\mathcal{F}_h \setminus \Omega_{h,r_0}} \mathbf{u} \otimes \mathbf{u} : D(\mathbf{w}) \, d\mathbf{x} \right| \leq \|\mathbf{u}\|_{L^2(\text{Supp}(\mathbf{w}))}^2 \|D(\mathbf{w})\|_{L^\infty(\mathcal{F}_h \setminus \Omega_{h,r_0})}.$$

Because $\text{Supp}(\mathbf{w})$ is bounded and $\mathbf{u} \in H_0^1(\mathbb{R}_+^3)$, we can use the Poincaré inequality. Concerning the integral in Ω_{h,r_0} , we have

$$I_1^{in} = \int_0^{2\pi} \int_0^{r_0} \int_0^{\delta_h(r)} [\mathbf{u}(r, \theta, z) \otimes \mathbf{u}(r, \theta, z) : D(\mathbf{w})(r, \theta, z)] r \, dz \, dr \, d\theta.$$

Using a Hölder inequality with respect to the z -variable, we deduce

$$|I_1^{in}| \leq C \int_0^{2\pi} \int_0^{r_0} \left[\int_0^{\delta_h(r)} |\mathbf{u}(r, \theta, z)|^4 \, dz \right]^{\frac{1}{2}} \left[\int_0^{\delta_h(r)} |D(\mathbf{w})|^2 \, dz \right]^{\frac{1}{2}} r \, dr \, d\theta.$$

Then a direct generalization of the Poincaré inequality (see Lemma 12 in [9]) implies

$$\left[\int_0^{\delta_h(r)} |\mathbf{u}(r, \theta, z)|^4 \, dz \right]^{\frac{1}{2}} \leq C \delta_h(r)^{\frac{3}{2}} \left[\int_0^{\delta_h(r)} |\nabla \mathbf{u}(r, \theta, z)|^2 \, dz \right].$$

Substituting this in I_1^{in} and using again a Hölder inequality, we then obtain (3.5).

To estimate I_2 , we decompose it in the same manner as I_1 , and with the same proof, we deduce that there exists $C = C(\text{Supp}(\mathbf{w}))$ such that

$$|I_2^{out}| \leq C \|\nabla \mathbf{u}\|_{L^2(\mathbb{R}_+^3)} \|\mathbf{w}\|_{L^2(\mathbb{R}_+^3 \setminus \Omega_{h,r_0})}.$$

It remains to estimate the integral in Ω_{h,r_0} :

$$I_2^{in} = \int_0^{2\pi} \int_0^{r_0} \int_0^{\delta_h(r)} [\mathbf{u}(r, \theta, z) \cdot \mathbf{w}(r, \theta, z)] r \, dz \, dr \, d\theta.$$

As above, a Hölder inequality in the z -variable associated to the Poincaré inequality implies

$$|I_2^{in}| \leq C \int_0^{2\pi} \int_0^{r_0} \left[\int_0^{\delta_h(r)} |\nabla \mathbf{u}(r, \theta, z)|^2 \, dz \right]^{\frac{1}{2}} \delta_h(r) \left[\int_0^{\delta_h(r)} |\mathbf{w}|^2 \, dz \right]^{\frac{1}{2}} r \, dr \, d\theta.$$

We conclude by using a Cauchy–Schwarz inequality. \square

3.1. Explicit formula. From now on h is a fixed positive parameter. As in [9], we introduce a velocity field which is a good approximation (in a sense to be made precise) to the solution to the Stokes problem

$$(3.6) \quad \begin{cases} \operatorname{div} \mathbb{T}(\mathbf{w}, q) = 0 \\ \operatorname{div} \mathbf{w} = 0 \\ \mathbf{w}|_{\mathcal{P}} = 0, \\ \mathbf{w}|_{\partial \mathcal{B}_h} = \mathbf{e}_3. \end{cases} \quad \text{in } \mathcal{F}_h,$$

At first, we focus on the divergence-free and boundary conditions. So we introduce a potential vector field \mathbf{a}_h and set $\mathbf{w}_h = \operatorname{curl} \mathbf{a}_h$. One choice for \mathbf{a}_h could be

$$\mathbf{a}_h^s(\mathbf{x}) := \frac{\eta_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1)}{2} (\mathbf{e}_3 \times (\mathbf{x} - \mathbf{G}_h)) \quad \forall \mathbf{x} \in \mathcal{F}_h, \quad \text{with } h_0 > 0.$$

The field $\mathbf{w}_h^s := \operatorname{curl} \mathbf{a}_h^s$ satisfies the divergence-free and boundary conditions regardless of the value of $h_0 < h$. However, when h goes to 0, this particular velocity field does not take advantage of the particular shape of the aperture between \mathcal{B} and \mathcal{P} . Thus, we need to find another velocity field, especially in this aperture, in $\Omega_{h,1/2}$.

As we want to obtain an approximation of the Stokes problem, we construct a velocity field in which the fluid escapes from under the sphere with the most efficiency. Consequently, we want the velocity field to be planar and radial in each plane. Thus, our potential vector field reads, in cylindrical coordinates,

$$\mathbf{a}_h^d(r, \theta, z) = (-\phi_h^d(r, z) \sin(\theta), \phi_h^d(r, z) \cos(\theta), 0)^\top \quad \forall (r, \theta, z) \in \Omega_{h,1/2},$$

so that, for all $(r, \theta, z) \in \Omega_{h,1/2}$,

$$\mathbf{w}_h^d(r, \theta, z) = \left(-\partial_z \phi_h^d(r, z) \cos(\theta), -\partial_z \phi_h^d(r, z) \sin(\theta), \partial_r \phi_h^d(r, z) + \frac{\phi_h^d(r, z)}{r} \right)^\top.$$

We set, in order to fit boundary conditions (this shall be critical in Lemma 3.2),

$$\phi_h^d(r, z) = r \chi_o \left(\frac{z}{\delta_h(r)} \right), \quad \text{with } \chi_o(s) = \frac{s^2(3 - 2s)}{2} \quad \forall s \in (0, 1).$$

From now on, we set $h_0 = (\sqrt{17/16} - 1)/2$. It remains to interpolate \mathbf{w}_h^s and \mathbf{w}_h^d so that we obtain

$$\mathbf{a}_h(\mathbf{x}) = \begin{cases} \eta_{1/2}(r) \mathbf{a}_h^d(r, \theta, z) + (1 - \eta_{1/2}(r)) \mathbf{a}_h^s(\mathbf{x}) & \text{in } \Omega_{h,1/2}, \\ \mathbf{a}_h^s(\mathbf{x}) & \text{in } \mathbb{R}_+^3 \setminus \Omega_{h,1/2} \end{cases}$$

and $\mathbf{w}_h = \operatorname{curl} \mathbf{a}_h$. Explicitly, in $\Omega_{h,1/2}$, we have

$$(3.7) \quad \mathbf{w}_h(r, \theta, z) = \eta_{1/2}(r) \mathbf{w}_h^d(r, \theta, z) + (1 - \eta_{1/2}(r)) \mathbf{w}_h^s(\mathbf{x}) + \mathbf{rem}_0(\mathbf{x}),$$

where, denoting by $\mathbf{n}\boldsymbol{\pi}_3(\mathbf{x}) = (x_1, x_2, 0)^\top / \sqrt{x_1^2 + x_2^2}$, we have

$$\mathbf{rem}_0(\mathbf{x}) = \eta'_{1/2}(r) \mathbf{n}\boldsymbol{\pi}_3(\mathbf{x}) \times (\mathbf{a}_h^d(r, \theta, z) - \mathbf{a}_h^s(\mathbf{x})) \quad \text{in } \Omega_{h,1/2}.$$

3.2. From static to moving test function. The main point in this subsection is to prove that, given a weak solution to (FSIS) (\mathbf{u}, \mathbf{G}) and $\zeta \in \mathcal{D}(0, T)$, the function $\tilde{\mathbf{w}}$ constructed in (3.3) is a suitable test function. To this end, we need to extend \mathbf{w}_h on \mathbb{R}_+^3 first. This is possible thanks to the following technical result.

LEMMA 3.2. *Given $h > 0$, we have*

$$\begin{aligned} \mathbf{w}_h(\mathbf{x}) &= \mathbf{e}_3, & \mathbf{a}_h(\mathbf{x}) &= (\mathbf{e}_3 \times \mathbf{x})/2 & \forall \mathbf{x} \in \partial\mathcal{B}_h, \\ \mathbf{w}_h(\mathbf{x}) &= 0, & \mathbf{a}_h(\mathbf{x}) &= 0 & \forall \mathbf{x} \in \mathcal{P}. \end{aligned}$$

Proof. We set $\lambda = z/\delta_h(r)$ and differentiations of λ by subscripts. We have

$$(3.8) \quad \partial_z \phi_h^d(r, z) = r \lambda_z \chi_o'(\lambda), \quad \partial_r \phi_h^d(r, z) = \chi_o(\lambda) + r \lambda_r \chi_o'(\lambda).$$

Computing with the value of λ yields

$$\lambda_z = \frac{1}{\delta_h(r)}, \quad \lambda_r = -\frac{z \delta_h'(r)}{(\delta_h(r))^2}.$$

Our choice for χ_o implies that

$$\chi_o(0) = \chi_o'(0) = 0, \quad \chi_o(1) = \frac{1}{2}, \quad \chi_o'(1) = 0.$$

Replacing λ by 0 in (3.8) yields

$$\phi_h^d(r, z) = \partial_z \phi_h^d(r, z) = \partial_r \phi_h^d(r, z) = 0 \quad \text{on } \mathcal{P}.$$

Consequently, $\mathbf{a}_h^d = \mathbf{w}_h^d = 0$ on \mathcal{P} . Replacing λ by 1,

$$\partial_z \phi_h^d(r, z) = 0, \quad \phi_h^d(r, z) = \frac{r}{2}, \quad \partial_r \phi_h^d(r, z) = \frac{1}{2} \quad \text{on } \partial\mathcal{B}_h.$$

Consequently, $\mathbf{a}_h^d(\mathbf{x}) = (\mathbf{e}_3 \times \mathbf{x})/2$ and $\mathbf{w}_h^d = \mathbf{e}_3$ on \mathcal{B}_h .

Concerning the smooth part, a straightforward computation leads to

$$\mathbf{w}_h^s(\mathbf{x}) = \eta_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) \mathbf{e}_3 + \eta'_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) \frac{\mathbf{x} - \mathbf{G}_h}{|\mathbf{x} - \mathbf{G}_h|} \times \frac{(\mathbf{e}_3 \times (\mathbf{x} - \mathbf{G}_h))}{2}.$$

Due to our choice, we have

$$\eta_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) = 1, \quad \eta'_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) = 0 \quad \text{on } \partial\mathcal{B}_h.$$

Consequently $\mathbf{a}_h^s(\mathbf{x}) = (\mathbf{e}_3 \times \mathbf{x})/2$ and $\mathbf{w}_h^s = \mathbf{e}_3$ on $\partial\mathcal{B}_h$. Then

$$\eta_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) = 0, \quad \eta'_{h_0}(|\mathbf{x} - \mathbf{G}_h| - 1) = 0 \quad \text{if } |\mathbf{x} - \mathbf{G}_h| \geq 1 + 2h_0 = \sqrt{17/16}.$$

Moreover, if $\mathbf{x} \in \mathcal{P} \setminus \overline{\Omega_{h,1/4}}$, we have $r > 1/4$ and, as $\mathbf{G}_h = (0, 0, 1 + h)$,

$$|\mathbf{x} - \mathbf{G}_h|^2 > (1 + h)^2 + (1/4)^2 > 17/16.$$

Consequently, $\mathbf{a}_h^s(\mathbf{x}) = \mathbf{w}_h^s(\mathbf{x}) = 0$ for $\mathbf{x} \in \mathcal{P} \setminus \overline{\Omega_{h,1/4}}$.

It remains to check that boundary conditions are satisfied in the transition region, i.e., when $\mathbf{x} \in \overline{\Omega_{h,1/2}} \setminus \Omega_{h,1/4}$. On \mathcal{P} , we remark that $\mathbf{w}_h^d(\mathbf{x}) = \mathbf{w}_h^s(\mathbf{x}) = 0 = \mathbf{a}_h^s(\mathbf{x}) =$

$\mathbf{a}_h^d(\mathbf{x}) = 0$. Interpolating the potential vector fields, we obtain $\mathbf{w}_h = 0$ on \mathcal{P} . Finally, on $\overline{\mathcal{B}_h} \cap \overline{\Omega_{h,1/4}}$ we have already computed

$$\mathbf{w}_h^d(\mathbf{x}) = \mathbf{w}_h^s(\mathbf{x}) = \mathbf{e}_3 \quad \text{and} \quad \mathbf{a}_h^s(\mathbf{x}) = \mathbf{a}_h^d(\mathbf{x}) = (\mathbf{e}_3 \times \mathbf{x})/2.$$

Interpolating the potential vector fields, we deduce $\mathbf{w}_h(\mathbf{x}) = \mathbf{e}_3$. This concludes the proof. \square

Remark 3.1. According to this lemma, we extend \mathbf{a}_h (resp., \mathbf{w}_h) to \mathbb{R}_+^3 with the value $(\mathbf{e}_3 \times \mathbf{x})/2$ (resp., \mathbf{e}_3) in \mathcal{B}_h . In what follows, we consider the functions $\mathbf{a} : (h, \mathbf{x}) \rightarrow \mathbf{a}_h(\mathbf{x})$ and $\mathbf{w} : (h, \mathbf{x}) \rightarrow \mathbf{w}_h(\mathbf{x})$. Denoting by $\mathcal{Q}_c = \{(h, \mathbf{x}) \in (0, 1) \times \mathbb{R}^3 ; \mathbf{x} \in \mathcal{B}_h\}$, standard analytic arguments imply $\mathbf{a} \in \mathcal{C}^\infty(\mathcal{Q}_c) \cap \mathcal{C}^\infty(((0, 1) \times \mathbb{R}_+^3) \setminus \overline{\mathcal{Q}_c})$. We note that \mathbf{w}_h vanishes as soon as $|\mathbf{x} - \mathbf{G}_h| > (\sqrt{17/16} - 1)/2$ and $|\mathbf{x}| > 1/2$. Consequently, the above lemma implies $\mathbf{w} \in H^1(\overline{(h, 1)} \times \mathbb{R}_+^3)$ for any $\bar{h} > 0$ and, after standard composition arguments, this yields

$$\tilde{\mathbf{w}} \in \mathcal{C}([0, T]; H^1(\mathbb{R}_+^3)) \cap H^1(0, T; L^2(\mathbb{R}_+^3))$$

as long as $h(t) \in (\bar{h}, 1]$ for all $t \in (0, T)$. So, $\tilde{\mathbf{w}}$ is a suitable test function for the weak formulation as long as $h(0, T) \subset (\bar{h}, 1)$.

3.3. Estimate of remainder terms. In order to exploit the weak formulation with our test function, we need to dominate remainder terms according to Lemma 3.1. We begin with estimates on Sobolev norms of \mathbf{w}_h .

By construction, our test functions behave differently under the sphere (in $\Omega_{h,\delta}$) and above the sphere (in $\mathcal{F}_h \setminus \Omega_{h,\delta}$ for arbitrary fixed $\delta > 0$). Above the sphere we have the following.

LEMMA 3.3. *Given $\alpha \geq 0$ and $\delta > 0$ there exists $C(\alpha, \delta) < \infty$ such that*

$$\|\mathbf{a}_h\|_{H^\alpha(\mathcal{F}_h \setminus \Omega_{h,\delta})} \leq C(\alpha, \delta) \quad \forall h \in (0, 1).$$

Proof. By construction the restriction $\mathbf{a} : \mathcal{Q}_{c,\delta} \rightarrow \mathbb{R}^3$, with

$$\mathcal{Q}_{c,\delta} := \{(h, \mathbf{x}) \in [0, 1] \times \overline{\mathbb{R}_+^3} \quad \text{with } \mathbf{x} \notin \Omega_{h,\delta}\},$$

is smooth and with compact support. \square

Inside $\Omega_{h,1/4}$, estimates rely essentially on dominations of integrals:

$$\int_0^{\frac{1}{4}} \frac{r^\alpha \, dr}{[\delta_h(r)]^\beta}.$$

We refer the reader to the appendix for such computations.

LEMMA 3.4. *The family $(\mathbf{w}_h)_{0 < h < 1}$ is uniformly bounded in $L^2(\mathbb{R}_+^3)$.*

Proof. Because of the previous lemma, we focus on \mathbf{w}_h^d inside $\Omega_{h,1/4}$.

In this region, we have

$$(\mathbf{w}_h^d(r, \theta, z))_1 = -\partial_z \phi_h^d(r, z) \cos(\theta), \quad (\mathbf{w}_h^d(r, \theta, z))_2 = -\partial_z \phi_h^d(r, z) \sin(\theta),$$

and

$$(\mathbf{w}_h^d(r, \theta, z))_3 = \partial_r \phi_h^d(r, z) + \frac{\phi_h^d(r, z)}{r}.$$

Thus

$$|\mathbf{w}_h^d(r, \theta, z)| \leq |\partial_z \phi_h^d(r, z)| + |\partial_r \phi_h^d(r, z)| + \frac{|\phi_h^d(r, z)|}{r}.$$

Applying Lemma A.3, this leads to

$$|\mathbf{w}_h^d(r, \theta, z)| \leq C \left(1 + \frac{r}{\delta_h(r)} \right) \quad \forall (r, \theta, z) \in \Omega_{h,1/4}, \quad \forall h \in (0, 1).$$

The result then follows from Lemma A.1 with $(\alpha, \beta) = (3, 1)$. \square

As a technical device for applying Lemma 3.1, we have the following.

LEMMA 3.5. *Let us define*

$$w_h(r, \theta, z) = |\partial_r \mathbf{w}_h^d(r, \theta, z)| + \frac{|\partial_\theta \mathbf{w}_h^d(r, \theta, z)|}{r} + |\partial_h \mathbf{w}_h^d(r, \theta, z)| \quad \forall (r, \theta, z) \in \Omega_{h,1/4}.$$

Then

$$(3.9) \quad \int_0^{\frac{1}{4}} \left(\delta_h(r)^2 \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0, 2\pi)} |w_h(r, \theta, z)|^2 dz \right] \right) r dr$$

is uniformly bounded for $h \in (0, 1)$.

Proof. A straightforward computation yields, for all $(r, \theta, z) \in \Omega_{h,1/4}$,

$$|\partial_r \mathbf{w}_h^d(r, \theta, z)| \leq C \left(|\partial_{rz} \phi_h^d(r, z)| + |\partial_{rr} \phi_h^d(r, z)| + \left| \frac{\partial_r \phi_h^d(r, z)}{r} - \frac{\phi_h^d(r, z)}{r^2} \right| \right)$$

and

$$\begin{aligned} \frac{|\partial_\theta \mathbf{w}_h^d(r, \theta, z)|}{r} &\leq \frac{|\partial_z \phi_h^d(r, z)|}{r}, \\ |\partial_h \mathbf{w}_h^d(r, \theta, z)| &\leq \frac{|\partial_h \phi_h^d(r, z)|}{r} + |\partial_{hr} \phi_h^d(r, z)| + |\partial_{hz} \phi_h^d(r, z)|. \end{aligned}$$

Combining the above inequalities with Lemma A.3, we deduce there exists a constant C independent of h such that

$$|w_h(r, \theta, z)| \leq C \left(\frac{1}{\delta_h(r)} + \frac{r}{\delta_h(r)^2} \right)$$

for all $(r, \theta, z) \in \Omega_{h,1/4}$. Consequently,

$$\int_0^{1/4} \left(\delta_h(r)^2 \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0, 2\pi)} |w_h(r, \theta, z)|^2 dz \right] \right) r dr \leq C \int_0^{1/4} (\delta_h(r) + 1)r dr.$$

As the last integral remains bounded when h goes to 0, the same holds for the integral of w_h . \square

Then, to dominate the trilinear form, we need the following result.

LEMMA 3.6. *We set*

$$dw_h(r, \theta, z) = |\partial_r \mathbf{w}_h^d(r, \theta, z)| + \frac{|\partial_\theta \mathbf{w}_h^d(r, \theta, z)|}{r} + |\partial_z \mathbf{w}_h^d(r, \theta, z)| \quad \forall (r, \theta, z) \in \Omega_{h,1/4}.$$

Then the quantity

$$(3.10) \quad \sup_{r \in (0, 1/4)} \left\{ \delta_h(r)^{\frac{3}{2}} \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0, 2\pi)} \{|dw_h(r, \theta, z)|^2\} dz \right]^{\frac{1}{2}} \right\}$$

is uniformly bounded for $h \in (0, 1)$.

Proof. As in the previous proof, there exists a constant C independent of h such that

$$|dw_h(r, \theta, z)| \leq C \left(\frac{1}{\delta_h(r)} + \frac{r}{\delta_h(r)^2} \right)$$

for all $(r, \theta, z) \in \Omega_{h,1/4}$. Therefore

$$\int_0^{\delta_h(r)} |dw_h(r, z)|^2 dz \leq C \left(\frac{1}{\delta_h(r)} + \frac{r^2}{\delta_h(r)^3} \right)$$

and

$$\sup_{r \in (0,1/4)} \left\{ \delta_h(r)^{\frac{3}{2}} \left[\int_0^{\delta_h(r)} \sup_{\theta \in (0,2\pi)} \{|dw_h(r, \theta, z)|^2\} dz \right]^{\frac{1}{2}} \right\} \leq C \sup_{r \in (0,1/4)} (\delta_h(r) + r),$$

which is uniformly bounded when $h \in (0, 1)$. \square

Finally, there holds the following lemma, which is reminiscent of works by Starovoitov.

LEMMA 3.7. *There exists a constant $C > 0$ such that*

$$|\nabla \mathbf{w}_h|_2^2 \geq \frac{C}{h} \quad \forall h \in (0, 1).$$

Proof. Given $h > 0$ we already noticed that

$$\mathbf{w}_h(\mathbf{x}) = \mathbf{w}_h^d(\mathbf{x}) \quad \forall \mathbf{x} \in \Omega_{h,1/4}.$$

Consequently,

$$|\nabla \mathbf{w}_h(\mathbf{x})| \geq |\partial_z \mathbf{w}_h^d(\mathbf{x})| \quad \forall \mathbf{x} \in \Omega_{h,1/4}.$$

With the explicit formula for \mathbf{w}_h^d , we have $|\partial_z \mathbf{w}_h^d| \geq |\partial_{zz} \phi_h^d|$, where $|\partial_{zz} \phi_h^d(r, z)| = \frac{r}{\delta_h^2(r)} \chi''(\lambda)$. Consequently,

$$|\nabla \mathbf{w}_h|_2^2 \geq 2\pi \int_0^{\frac{1}{4}} \frac{r^3 dr}{\delta_h(r)^3} \int_0^1 |\chi''(s)|^2 ds.$$

As χ is a polynomial with degree 3, its second derivative does not vanish, and neither does the s -integral. Then we obtain the result by applying Lemma A.1 with $\alpha = \beta = 3$. \square

3.4. Our test function and the Stokes problem. First, we prove that our choice is a good one because it is a good approximation of the solution to the Stokes problem.

LEMMA 3.8. *There exist $q_h \in C^\infty(\overline{\mathcal{F}_h})$ and $\mathbf{f}_h \in C_c^\infty(\overline{\mathcal{F}_h})$ such that*

$$(3.11) \quad \begin{cases} \mu \Delta \mathbf{w}_h - \nabla q_h = \mathbf{f}_h \\ \operatorname{div} \mathbf{w}_h = 0 \end{cases} \quad \text{in } \mathcal{F}_h,$$

where there exists an absolute constant C for which

$$\int_0^{2\pi} \int_0^{\frac{1}{4}} \left(\delta_h(r)^2 \left[\int_0^{\delta_h(r)} |\mathbf{f}_h|^2 dz \right] \right) r dr d\theta + \|\mathbf{f}_h\|_{L^2(\mathcal{F}_h \setminus \Omega_{h,1/4})}^2 \leq C.$$

Proof. By construction, we have $\mathbf{w}_h = \text{curl} \tilde{\mathbf{a}}_h^d + \text{curl} \tilde{\mathbf{a}}_h^s$, where

$$\tilde{\mathbf{a}}_h^d(\mathbf{x}) = \begin{cases} \eta_{1/2}(r)\mathbf{a}_h^d(\mathbf{x}), & \mathbf{x} \in \Omega_{h,1/2}, \\ 0 & \text{else,} \end{cases} \quad \tilde{\mathbf{a}}_h^s(\mathbf{x}) = \begin{cases} (1 - \eta_{1/2}(r))\mathbf{a}_h^s(\mathbf{x}), & \mathbf{x} \in \Omega_{h,1/2}, \\ \mathbf{a}_h^s(\mathbf{x}) & \text{else.} \end{cases}$$

Then, according to Lemma 3.3, the smooth part $\tilde{\mathbf{a}}_h^s$ is bounded in any Sobolev space uniformly in h . Consequently, $\tilde{\mathbf{f}}_h = \mu\Delta \text{curl} \tilde{\mathbf{a}}_h^s$ is bounded in all Sobolev spaces. We have

$$\mu\Delta\mathbf{w}_h = \mu\Delta \text{curl} \tilde{\mathbf{a}}_h^d + \tilde{\mathbf{f}}_h.$$

In the following we write $\tilde{\phi}_h^d(r, z) = \eta_{1/2}(r)\phi_h^d(r, z)$ for all $(r, \theta, z) \in \Omega_{h,1/2}$. Let us recall that in cylindrical coordinates we have

$$\Delta = \frac{\partial_r[r\partial_r]}{r} + \frac{\partial_{\theta\theta}}{r^2} + \partial_{zz}.$$

Consequently,

$$[\Delta \text{curl} \tilde{\mathbf{a}}_h^d]_1 = - \left[\partial_{rrz}\tilde{\phi}_h^d + \frac{\partial_{rz}\tilde{\phi}_h^d}{r} - \frac{\partial_z\tilde{\phi}_h^d}{r^2} + \partial_{zzz}\tilde{\phi}_h^d \right] \cos(\theta)$$

and

$$[\Delta \text{curl} \tilde{\mathbf{a}}_h^d]_2 = - \left[\partial_{rrz}\tilde{\phi}_h^d + \frac{\partial_{rz}\tilde{\phi}_h^d}{r} - \frac{\partial_z\tilde{\phi}_h^d}{r^2} + \partial_{zzz}\tilde{\phi}_h^d \right] \sin(\theta),$$

$$[\Delta \text{curl} \tilde{\mathbf{a}}_h^d]_3 = \partial_{rrr}\tilde{\phi}_h^d + 2\frac{\partial_{rr}\tilde{\phi}_h^d}{r} - \frac{\partial_r\tilde{\phi}_h^d}{r^2} + \frac{\tilde{\phi}_h^d}{r^3} + \partial_{zz} \left[\partial_r\tilde{\phi}_h^d + \frac{\tilde{\phi}_h^d}{r} \right].$$

We remark here that

$$\partial_{zzz}\tilde{\phi}_h^d(r, z) = -6\frac{r\eta_{1/2}(r)}{\delta_h^3(r)}.$$

Consequently, denoting by Φ a primitive of $s \mapsto -6s\eta_{1/2}(s)/(\delta_h(s))^3$, we have

$$\nabla\Phi(r) = (\partial_{zzz}\phi_h^d \cos(\theta), \partial_{zzz}\phi_h^d \sin(\theta), 0)^\top.$$

We set

$$q_h(\mathbf{x}) = \mu\Phi(r) + \mu\partial_z \left[\partial_r\tilde{\phi}_h^d + \frac{\tilde{\phi}_h^d}{r} \right], \quad \check{\mathbf{f}}_h = \mu\Delta \text{curl} \tilde{\mathbf{a}}_h^d - \nabla q_h.$$

In particular $\mu\Delta\mathbf{w}_h - \nabla q_h = \tilde{\mathbf{f}}_h + \check{\mathbf{f}}_h$, so that our result follows from the same result for $\check{\mathbf{f}}_h$. Denoting by $\check{f}_1, \check{f}_2, \check{f}_3$ the Cartesian components of $\check{\mathbf{f}}_h$, straightforward computations show that

$$|\check{f}_1|^2 + |\check{f}_2|^2 \leq 4 \left[\partial_{rrz}\tilde{\phi}_h^d + \frac{\partial_{rz}\tilde{\phi}_h^d}{r} - \frac{\partial_z\tilde{\phi}_h^d}{r^2} \right]^2.$$

As $\eta'_{1/2}$ vanishes uniformly in $\Omega_{h,1/4}$, Lemma 3.3 implies there exists a universal constant C such that

$$|\check{f}_1|^2 + |\check{f}_2|^2 \leq C \left[1 + \left| \partial_{rrz} \phi_h^d + \frac{\partial_{rz} \phi_h^d}{r} - \frac{\partial_z \phi_h^d}{r^2} \right| \right]^2.$$

Then, for the same reasons, we have

$$|\check{f}_3|^2 \leq C \left[1 + |\partial_{rrr} \phi_h^d| + \left| \frac{2\partial_{rr} \phi_h^d}{r} \right| + \left| \frac{\phi_h^d}{r^3} - \frac{\partial_r \phi_h^d}{r^2} \right| \right]^2.$$

Replacing with the size computed in Lemma A.3, we obtain

$$|\check{\mathbf{f}}_h(\mathbf{x})|^2 \leq C \left(\frac{r}{\delta_h^2(r)} + \frac{1}{\delta_h(r)} \right)^2.$$

Consequently, for arbitrary $r \in (0, 1/2)$

$$\int_0^{\delta_h(r)} |\check{\mathbf{f}}_h|^2 \, dz \leq C \left(\frac{r^2}{\delta_h(r)^3} + \frac{1}{\delta_h(r)} \right)$$

and

$$\int_0^{2\pi} \int_0^{1/2} \delta_h(r)^2 \int_0^{\delta_h(r)} |\check{\mathbf{f}}_h|^2 \, dz \, dr \, d\theta \leq C \int_0^{1/2} \left(\frac{r^3}{\delta_h(r)} + r\delta_h(r) \right) \, dr,$$

which is uniformly bounded for $h \in (0, 1)$. This concludes the proof. \square

As a direct corollary, we get the following lemma.

LEMMA 3.9. *There exist $K_m < \infty$ and a function $\tilde{n}_3 : [0, 1] \rightarrow \mathbb{R}_+$ such that, for any $h < 1$ and $\mathbf{w} \in \mathbb{V}(\mathbf{G}_h)$ such that $\mathbf{w} = \mathbf{V}_\mathbf{w} + \mathbf{R}_\mathbf{w} \times (\mathbf{x} - \mathbf{G}_h)$ in \mathcal{B}_h , we have*

$$(3.12) \quad \left| 2\mu \int_{\mathbb{R}_+^3} D(\mathbf{w}_h) : D(\mathbf{w}) \, d\mathbf{x} - \tilde{n}_3(h) \mathbf{V}_\mathbf{w} \cdot \mathbf{e}_3 \right| \leq K_m \|\nabla \mathbf{w}\|_{L^2(\mathbb{R}_+^3)}.$$

Moreover, there exist $h_m > 0$ and a constant $c > 0$ such that $\tilde{n}_3(h) \geq c/h$ for all $h < h_m$.

Proof. Given $h > 0$ and $\mathbf{w} \in \mathbb{V}(\mathbf{G}_h)$, we apply the Stokes identity with (3.11) and obtain

$$(3.13) \quad 2\mu \int_{\mathbb{R}_+^3} D(\mathbf{w}_h) : D(\mathbf{w}) \, d\mathbf{x} = \int_{\partial \mathcal{B}_h} \mathbb{T}(\mathbf{w}_h, p_h) \mathbf{n} \cdot \mathbf{w} \, d\sigma - \int_{\mathcal{F}_h} \mathbf{f}_h \cdot \mathbf{w} \, d\mathbf{x}.$$

For symmetry reasons, there exists $\tilde{n}_3 : (0, 1) \rightarrow \mathbb{R}$ such that

$$\int_{\partial \mathcal{B}_h} \mathbb{T}(\mathbf{w}_h, p_h) \mathbf{n} \, d\sigma = \tilde{n}_3(h) \mathbf{e}_3$$

and

$$\int_{\partial \mathcal{B}_h} (\mathbf{x} - \mathbf{G}_h) \times \mathbb{T}(\mathbf{w}_h, p_h) \mathbf{n} \, d\sigma = 0.$$

On the other hand, applying (3.4) and Lemma 3.8, we also deduce

$$\left| \int_{\mathcal{F}_h} \mathbf{f}_h \cdot \mathbf{w} \, d\mathbf{x} \right| \leq C \|\nabla \mathbf{w}\|_{L^2(\mathbb{R}_+^3)} \quad \forall h \in (0, 1),$$

where C is a positive constant. Finally, we have obtained the existence of a constant K such that, for arbitrary $h \in (0, 1)$ and $\mathbf{w} \in \mathbb{V}(\mathbf{G})$,

$$\left| 2\mu \int_{\mathbb{R}_+^3} D(\mathbf{w}_h) : D(\mathbf{w}) \, d\mathbf{x} - \tilde{n}_3(h) \mathbf{V}_\mathbf{w} \cdot \mathbf{e}_3 \right| \leq K \|\nabla \mathbf{w}\|_{L^2(\mathbb{R}_+^3)}.$$

In order to estimate \tilde{n}_3 , we take $\mathbf{w} = \mathbf{w}_h$ in (3.13) and obtain

$$(3.14) \quad \tilde{n}_3 = \int_{\partial \mathcal{B}_h} \mathbb{T}(\mathbf{w}_h, p_h) \mathbf{n} \cdot \mathbf{e}_3 \, d\sigma = 2\mu \int_{\mathbb{R}_+^3} |D(\mathbf{w}_h)|^2 \, d\mathbf{x} + \int_{\mathcal{F}_h} \mathbf{f}_h \cdot \mathbf{w}_h \, d\mathbf{x}.$$

Dealing as previously with the last integral, we deduce that

$$\left| \int_{\mathcal{F}_h} \mathbf{f}_h \cdot \mathbf{w}_h \, d\mathbf{x} \right| \leq K \|\nabla \mathbf{w}_h\|_{L^2(\mathbb{R}_+^3)} \quad \forall h \in (0, 1).$$

But, applying Lemma 3.7, we have that

$$2\mu \int_{\mathbb{R}_+^3} |D(\mathbf{w}_h)|^2 \, d\mathbf{x} = \mu \int_{\mathbb{R}_+^3} |\nabla \mathbf{w}_h|^2 \, d\mathbf{x} \geq \frac{C}{h} \quad \forall h \in (0, 1).$$

Consequently, the asymptotic behavior of the right-hand side in (3.14) when h goes to 0 is prescribed by the first integral. Hence, there exist $h_m > 0$ and constants $\tilde{c}, c > 0$ such that

$$\tilde{n}_3(h) \geq \tilde{c} \int_{\mathbb{R}_+^3} |D(\mathbf{w}_h)|^2 \, d\mathbf{x} \geq \frac{c}{h} \quad \forall h \in (0, h_m). \quad \square$$

4. Proof of Theorem 1.1. We let the reader convince himself that Theorem 1.1 is a direct consequence of the following theorem.

THEOREM 4.1. *Given (\mathbf{u}, \mathbf{G}) a weak solution to (FSIS) on $(0, T)$ with initial data $(\mathbf{u}^0, \mathbf{G}^0)$, we assume there exists $0 \leq \tau_0 < \tau_1 \leq T$ for which*

$$h(t) := \text{dist}(\mathcal{B}(t), \mathcal{P}) \leq 1 \quad \forall t \in [\tau_0, \tau_1].$$

Then there exists $C(\|\mathbf{u}^0\|_{L^2(\mathbb{R}_+^3)}) < \infty$ depending only on the L^2 -norm of initial data such that

$$h(t) \geq h(\tau_0) \exp \left[-C(\|\mathbf{u}^0\|_{L^2(\mathbb{R}_+^3)})(1 + \sqrt{T}) \right] \quad \forall t \in (\tau_0, \tau_1).$$

The remainder of this paper is devoted to the proof of this result. From now on (\mathbf{u}, \mathbf{G}) is a given weak solution to (FSIS) with initial data $(\mathbf{u}^0, \mathbf{G}^0)$. For simplicity, we assume that $h(t) \leq 1$ for all $t \in (0, T)$. This means that $\tau_0 = 0$ and $\tau_1 = T$ in the assumptions of our theorem.

As mentioned before, we estimate the distance h from below with our approximation of the Stokes problem. So, from now on, $(\mathbf{w}_h)_{h \in (0,1)}$ are the approximations constructed in section 3.1. Given $0 < t_0 < t_1 < 1$, we set

$$\zeta_\varepsilon(t) = \eta_\varepsilon(\text{dist}(t, [t_0, t_1])).$$

Then $\zeta_\varepsilon \in \mathcal{D}(0, T)$ whenever ε is sufficiently small. Consequently, according to Remark 3.1, for ε sufficiently small

$$\begin{aligned} \tilde{\mathbf{w}}_\varepsilon : (0, T) \times \mathbb{R}_+^3 &\longrightarrow \mathbb{R}^3, \\ (t, \mathbf{x}) &\longmapsto \zeta_\varepsilon(t) \mathbf{w}_{h(t)}(x_1 - G_1(t), x_2 - G_2(t), x_3) \end{aligned}$$

can be taken as a test function in (2.4):

$$(4.1) \quad \int_{(0,T) \times \mathbb{R}_+^3} [\rho_{\mathbf{G}} \mathbf{u} \cdot \partial_t \tilde{\mathbf{w}}_\varepsilon + (\mathbf{u} \otimes \mathbf{u} - 2\mu D(\mathbf{u})) : D(\tilde{\mathbf{w}}_\varepsilon)] \, d\mathbf{x} \, dt = 0.$$

In the following, we set

$$\begin{aligned} I_1 &:= \int_{(0,T) \times \mathbb{R}_+^3} \rho_{\mathbf{G}} \mathbf{u} \cdot \partial_t \tilde{\mathbf{w}}_\varepsilon \, dt \, d\mathbf{x}, \\ I_2 &:= \int_{(0,T) \times \mathbb{R}_+^3} \mathbf{u} \otimes \mathbf{u} : D(\tilde{\mathbf{w}}_\varepsilon) \, dt \, d\mathbf{x}, \\ I_3 &:= \int_{(0,T) \times \mathbb{R}_+^3} D(\mathbf{u}) : D(\tilde{\mathbf{w}}_\varepsilon) \, dt \, d\mathbf{x}. \end{aligned}$$

After a change of variables, we have for almost all $t \in (0, T)$

$$\int_{\mathbb{R}_+^3} D(\mathbf{u})(t, \cdot) : D(\tilde{\mathbf{w}}_\varepsilon)(t, \cdot) \, d\mathbf{x} = \zeta_\varepsilon(t) \int_{\mathbb{R}_+^3} D(\mathbf{u})(t, x_1 + G_1, x_2 + G_2, x_3) : D(\mathbf{w}_{h(t)}) \, d\mathbf{x}.$$

Thus, applying Lemma 3.9,

$$\int_{\mathbb{R}_+^3} D(\mathbf{u})(t, x_1 + G_1, x_2 + G_2, x_3) : D(\mathbf{w}_{h(t)}) \, d\mathbf{x} = \dot{h} \tilde{n}_3(h) + E(t),$$

where $|E(t)| = K_M |\nabla \mathbf{u}(t, \cdot)|_2$. Consequently,

$$(4.2) \quad I_3 = \int_0^T \zeta_\varepsilon(t) \dot{h}(t) \tilde{n}_3(h(t)) \, dt + \tilde{E},$$

where

$$(4.3) \quad |\tilde{E}| \leq K_M \sqrt{T} \|\mathbf{u}_0\|_2.$$

Similarly, for almost all $t \in (0, T)$,

$$\begin{aligned} &\int_{\mathbb{R}_+^3} [\mathbf{u} \otimes \mathbf{u}](t, \cdot) : D(\tilde{\mathbf{w}}_\varepsilon)(t, \cdot) \, d\mathbf{x} \\ &= \zeta_\varepsilon(t) \int_{\mathbb{R}_+^3} [\rho_{\mathbf{G}_h} \mathbf{u} \otimes \mathbf{u}](t, x_1 + G_1, x_2 + G_2, x_3) : D(\mathbf{w}_{h(t)})(\mathbf{x}) \, d\mathbf{x}. \end{aligned}$$

Consequently, applying Lemma 3.1 together with Lemmas 3.6 and 3.3, we obtain

$$\left| \int_{\mathbb{R}_+^3} [\mathbf{u} \otimes \mathbf{u}](t, \cdot) : D(\tilde{\mathbf{w}}_\varepsilon)(t, \cdot) \, d\mathbf{x} \right| \leq K_m \|\nabla \mathbf{u}\|_{L^2(\mathbb{R}_+^3)}^2.$$

Thus,

$$(4.4) \quad |I_2| \leq K_m \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2.$$

Finally, computing $\partial_t \tilde{\mathbf{w}}_\varepsilon$ we have, after our change of variable,

$$I_1 = I_1^X + I_1^w,$$

where

$$I_1^\chi := \int_0^T \int_{\mathbb{R}_+^3} [\rho_{\mathbf{G}_h} \mathbf{u}](t, x_1 + G_1, x_2 + G_2, x_3) \cdot \zeta'_\varepsilon(t) \mathbf{w}_{h(t)}(\mathbf{x}) \, d\mathbf{x} \, dt$$

and

$$I_1^w := \int_0^T \zeta_\varepsilon(t) \int_{\mathbb{R}_+^3} [\rho_{\mathbf{G}_h} \mathbf{u}](x_1 + G_1, x_2 + G_2, x_3) \cdot \left[\dot{h} \partial_h \mathbf{w}_h - V_1 \partial_{x_1} \mathbf{w}_h - V_2 \partial_{x_2} \mathbf{w}_h \right] (\mathbf{x}) \, d\mathbf{x} \, dt.$$

Applying the Cauchy–Schwarz inequality and Lemma 3.4 on \mathbf{w}_h , we deduce that

$$|I_1^\chi| \leq C \left[\int_{t_0-\varepsilon}^{t_0} |\zeta'_\varepsilon(t)| \|\mathbf{u}(t, \cdot)\|_{L^2(\mathbb{R}_+^3)} \, dt + \int_{t_1}^{t_1+\varepsilon} |\zeta'_\varepsilon(t)| \|\mathbf{u}(t, \cdot)\|_{L^2(\mathbb{R}_+^3)} \, dt \right],$$

and therefore, using the uniform L^2 -bound on \mathbf{u} , we obtain

$$(4.5) \quad |I_1^\chi| \leq K \|\mathbf{u}^0\|_{L^2(\mathbb{R}_+^3)}.$$

Finally, applying (3.4) in Lemma 3.1 together with (3.9) in Lemma 3.5, we conclude that

$$|I_1^w| \leq K_m \int_0^T [|\dot{h}|^2 + |\mathbf{V}_1|^2 + |\mathbf{V}_2|^2]^{\frac{1}{2}} \|\nabla \mathbf{u}\|_{L^2(\mathbb{R}_+^3)} \, dt,$$

so that, with standard energy estimate,

$$(4.6) \quad |I_1^w| \leq K_m \sqrt{T} \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2.$$

Gathering (4.2), (4.4), (4.5), (4.6) with (4.1) yields

$$\left| \int_0^T \zeta_\varepsilon(t) \dot{h}(t) \tilde{n}_3(h(t)) \, dt \right| \leq K_m (1 + \sqrt{T}) \left\{ \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)} + \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2 \right\},$$

where we emphasize that K_m depend only on our choice for the approximation of the solution to the Stokes problem, but not on ε . Thus, letting ε go to 0, as h and \tilde{n}_3 are continuous functions, we obtain

$$|N_3(h(t_1)) - N_3(h(t_0))| \leq K_m (1 + \sqrt{T}) \left\{ \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)} + \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2 \right\},$$

where N_3 is a primitive of \tilde{n}_3 which vanishes in $h = 1$, for example. Applying Lemma 3.9, we have $\tilde{n}_3(h) \geq c/h$ when $0 < h < h_m$ for some $c > 0$ and $h_m > 0$, and we finally deduce

$$|\ln(h(t)/h(t_0))| \leq K_m (1 + \sqrt{T}) \left\{ \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)} + \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2 \right\}.$$

Because h is continuous, letting t_0 tend to 0, we finally obtain

$$h(t) \geq h(0) \exp \left[-K_m (1 + \sqrt{T}) \left\{ \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)} + \|\mathbf{u}_0\|_{L^2(\mathbb{R}_+^3)}^2 \right\} \right].$$

This is the expected result.

Appendix. Detailed description of ϕ_h^d . In this section we estimate the size of ϕ_h^d and its derivatives. We recall that

$$\phi_h^d(r, \theta, z) = r\chi_o(z/\delta_h(r)), \quad \text{with} \quad \chi_o(s) = \frac{s^2(3-2s)}{2}.$$

In order to compare functions in the following, we introduce the following conventions. Given families $(f_h : \Omega_{h,1/4} \rightarrow \mathbb{R})_{h \in (0,1)}$ and $(g_h : \Omega_{h,1/4} \rightarrow \mathbb{R})_{h \in (0,1)}$ we denote $f_h \prec g_h$ if there exists an absolute constant such that

$$|f_h(\mathbf{x})| \leq Cg_h(\mathbf{x}) \quad \forall \mathbf{x} \in \Omega_{h,1/4} \text{ and } h < 1.$$

Given nonnegative functions $f : (0, 1) \rightarrow \mathbb{R}^+$ and $g : (0, 1) \rightarrow \mathbb{R}^+$, we also denote

$$f(s) \sim g(s) \quad \forall s \in (0, 1)$$

if there exist two positive constants c and C such that

$$cf(s) \leq g(s) \leq Cf(s) \quad \forall s \in (0, 1).$$

First, we compute typical $L^1(0, 1/4)$ -sizes of functions $r \mapsto r^\alpha / (\delta_h(r))^\beta$.

LEMMA A.1. *Given $(\alpha, \beta) \in (\mathbb{R}_+)^2$, we have the following estimations for all $h \in (0, 1)$:*

$$\int_0^{1/4} \frac{r^\alpha}{\delta_h(r)^\beta} dr \sim \begin{cases} 1 & \text{if } \alpha > 2\beta - 1, \\ h^{\frac{(\alpha+1)-2\beta}{2}} & \text{if } \alpha < 2\beta - 1. \end{cases}$$

Proof. As in [9], we remark that, for all $h \in (0, 1)$, we have

$$h + \frac{s^2}{2} \leq \delta_h(s) \leq h + s^2 \quad \forall s \in (0, 1/4).$$

Consequently, we can replace $\delta_h(r)$ by $h + \gamma r^2$ with some generic parameter $\gamma > 0$, and we are bound to calculate

$$I_{\alpha,\beta} := \int_0^{1/4} \frac{r^\alpha}{(h + \gamma r^2)^\beta} dr,$$

in which we set $r = \sqrt{hs}$. This yields

$$I_{\alpha,\beta} := h^{\frac{(\alpha+1)-2\beta}{2}} \int_0^{\frac{1}{4\sqrt{h}}} \frac{s^\alpha}{(1 + \gamma s^2)^\beta} ds.$$

Consequently, if $\alpha > 2\beta - 1$, the integral behaves like $Ch^{-\frac{(\alpha+1)-2\beta}{2}}$, and we obtain the first case, while if $\alpha < 2\beta - 1$, the integral goes to a finite positive value as $h \rightarrow \infty$, and we obtain the second case. \square

We now compare $\lambda(r, z, h) = z/\delta_h(r)$ to member functions $(r, \theta, z) \mapsto r^\alpha / (\delta_h(r))^\beta$ in $\Omega_{h,1/4}$.

LEMMA A.2. *We have the following sizes:*

$$\begin{aligned} \lambda &\prec 1, & \lambda_r &\prec r/\delta_h, & \lambda_z &\prec 1/\delta_h, & \lambda_h &\prec 1/\delta_h, \\ \lambda_{rh} &\prec r/\delta_h^2, & \lambda_{zh} &\prec 1/\delta_h^2, & \lambda_{rr} &\prec 1/\delta_h, & \lambda_{rz} &\prec r/\delta_h^2, \\ \lambda_{rrz} &\prec 1/\delta_h^2, & & & \lambda_{rrr} &\prec r/\delta_h^2. & & \end{aligned}$$

Proof. The reader may rapidly check that all the derivatives of δ_h are independent of h and that all the odd ones are bounded by r over $(0, 1/4)$. Then in $\Omega_{h,1/4}$ we have $z \in (0, \delta_h(r))$, and consequently $\lambda < 1$. Then

$$\lambda_r = -\frac{\lambda\delta'_h}{\delta_h}, \quad \lambda_z = \frac{1}{\delta_h}, \quad \lambda_h = -\frac{\lambda}{\delta_h}.$$

As δ' is bounded by r necessarily independent of h , we get $\lambda_r < r/\delta_h$ and $\lambda_z < 1/\delta_h$, $\lambda_h < 1/\delta_h$. To the next order, we get

$$\lambda_{rz} = -\frac{\delta'_h}{\delta_h^2}, \quad \lambda_{rr} = \lambda \left(2\frac{(\delta'_h)^2}{\delta_h^2} - \frac{\delta''_h}{\delta_h} \right), \quad \lambda_{rh} = -2\frac{\lambda\delta'_h}{\delta_h^2}, \quad \lambda_{zh} = -\frac{1}{\delta_h^2}.$$

As δ'' is bounded independently of h and $r^2 \leq h + r^2$, we obtain

$$\lambda_{rz} < \frac{r}{\delta_h^2}, \quad \lambda_{rr} < \frac{1}{\delta_h}, \quad \lambda_{rh} < \frac{r}{\delta_h}, \quad \lambda_{zh} < \frac{1}{\delta_h^2}.$$

Finally,

$$\lambda_{rrz} = \frac{1}{\delta_h} \left(2\frac{(\delta'_h)^2}{\delta_h^2} - \frac{\delta''_h}{\delta_h} \right), \quad \lambda_{rrr} = \lambda \left(6\frac{\delta''_h\delta'_h}{\delta_h^2} - 6\frac{(\delta'_h)^3}{\delta_h^3} - \frac{\delta_h^{(3)}}{\delta_h} \right).$$

And, as $\delta_h^{(3)}$ is bounded by r and $r^2 \leq \delta_h$,

$$\lambda_{rrz} < \frac{1}{\delta_h^2}, \quad \lambda_{rrr} < \frac{r}{\delta_h^2}. \quad \square$$

Then we obtain the following lemma.

LEMMA A.3. *We have the following sizes:*

$$\begin{aligned} \phi_h^d &< r, & \partial_r \phi_h^d &< 1, & \partial_z \phi_h^d &< r/\delta_h, & \partial_r \phi_h^d / r - \phi_h^d / r^2 &< r/\delta_h, \\ \partial_h \phi_h^d &< r/\delta_h, & \partial_{rh} \phi_h^d &< 1/\delta_h, & \partial_{zh} \phi_h^d &< r/\delta_h^2, & \partial_{rz} \phi_h^d / r - \partial_z \phi_h^d / r^2 &< r/\delta_h^2, \\ \partial_{rr} \phi_h^d &< r/\delta_h, & \partial_{rz} \phi_h^d &< 1/\delta_h, & \partial_{zz} \phi_h^d &< r/\delta_h^2, \\ \partial_{rrr} \phi_h^d &< 1/\delta_h, & \partial_{rzz} \phi_h^d &< 1/\delta_h^2, & \partial_{rrz} \phi_h^d &< r/\delta_h^2, & \partial_{zzz} \phi_h^d &< r/\delta_h^3. \end{aligned}$$

Proof. By definition, we have $\phi_h^d(r, z) = r\chi_o(\lambda)$, where χ_o is a fixed polynomial and, according to the previous lemma, λ is bounded. Consequently, we obtain $\phi_h^d < r$.

In the following, we shall drop all dependencies of χ_o in λ . Due to the same argument as for χ_o , all those quantities depending only on χ_o are bounded independently of (h, r, z) in $\Omega_{h,1/4}$.

So, we compute

$$\partial_r \phi_h^d = \chi_o + r\lambda_r \chi'_o, \quad \partial_z \phi_h^d = r\lambda_z \chi'_o, \quad \partial_h \phi_h^d = r\lambda_h \chi'_o.$$

Applying the previous lemma and $r^2 \leq \delta_h(r)$, we get

$$\partial_r \phi_h^d < 1, \quad \partial_z \phi_h^d < r/\delta_h, \quad \partial_h \phi_h^d < r/\delta_h, \quad \partial_r \phi_h^d / r - \phi_h^d / r^2 = \lambda_r \chi'_o < r/\delta_h.$$

To the next order, we obtain, as λ_z is independent of z ,

$$\begin{aligned} \partial_{rr} \phi_h^d &= (2\lambda_r + r\lambda_{rr})\chi'_o + r(\lambda_r)^2 \chi''_o, \\ \partial_{rz} \phi_h^d &= (\lambda_z + r\lambda_{rz})\chi'_o + r\lambda_r \lambda_z \chi''_o, \quad \partial_{zz} \phi_h^d = r(\lambda_z)^2 \chi''_o, \\ \partial_{zh} \phi_h^d &= r\lambda_z \lambda_h \chi''_o + r\lambda_{hz} \chi'_o, \quad \text{and} \quad \partial_{rh} \phi_h^d = \lambda_h \chi'_o + r\lambda_{rh} \chi'_o + r\lambda_r \lambda_h \chi''_o. \end{aligned}$$

As above,

$$\partial_{rr}\phi_h^d \prec r/\delta_h, \quad \partial_{rz}\phi_h^d \prec 1/\delta_h, \quad \partial_{zz}\phi_h^d \prec r/\delta_h^2, \quad \partial_{rh}\phi_h^d \prec 1/\delta_h, \quad \partial_{zh}\phi_h^d \prec r/\delta_h^2$$

and

$$\frac{\partial_{rz}\phi_h^d}{r} - \frac{\partial_z\phi_h^d}{r^2} = \lambda_{rz}\chi_o' + \lambda_r\lambda_z\chi_o'' \prec r/\delta_h^2.$$

To the next order, we obtain

$$\partial_{rrr}\phi_h^d = (3\lambda_{rr} + r\lambda_{rrr})\chi_o' + (3\lambda_r^2 + 3r\lambda_{rr}\lambda_r)\chi_o'' + r(\lambda_r)^3\chi_o^{(3)},$$

and thus $\partial_{rrr}\phi_h^d \prec 1/\delta_h$; and

$$\partial_{rzz}\phi_h^d = (\lambda_z^2 + 2r\lambda_{rz}\lambda_z)\chi_o'' + r(\lambda_z)^2\lambda_r\chi_o^{(3)},$$

so $\partial_{rzz}\phi_h^d \prec 1/\delta_h^2$; and

$$\partial_{rrz}\phi_h^d = (2\lambda_{rz} + r\lambda_{rrz})\chi_o' + (2\lambda_r\lambda_z + r(\lambda_{rr}\lambda_z + 2\lambda_{rz}\lambda_r))\chi_o'' + r(\lambda_r)^2\lambda_z\chi_o^{(3)},$$

so $\partial_{rrz}\phi_h^d \prec r/\delta_h^2$. Finally, $\partial_{zzz}\phi_h^d = r(\lambda_z)^3\chi_o^{(3)}$, so that $\partial_{zzz}\phi_h^d \prec r/\delta_h^3$. \square

REFERENCES

- [1] M. D. A. COOLEY AND M. E. O'NEILL, *On the slow motion generated in a viscous fluid by the approach of a sphere to a plane wall or stationary sphere*, *Mathematika*, 16 (1969), pp. 37–49.
- [2] P. CUMSILLE AND T. TAKAHASHI, *Wellposedness for the system modelling the motion of a rigid body of arbitrary form in an incompressible viscous fluid*, *Czechoslovak Math. J.*, 58 (2008), pp. 961–992.
- [3] P. CUMSILLE AND M. TUCSNAK, *Wellposedness for the Navier-Stokes flow in the exterior of a rotating obstacle*, *Math. Methods Appl. Sci.*, 29 (2006), pp. 595–623.
- [4] E. FEIREISL, *On the motion of rigid bodies in a viscous incompressible fluid*, *J. Evol. Equ.*, 3 (2003), pp. 419–441.
- [5] H. FUJITA AND N. SAUER, *On existence of weak solutions of the Navier-Stokes equations in regions with moving boundaries*, *J. Fac. Sci. Univ. Tokyo Sect. I*, 17 (1970), pp. 403–420.
- [6] G. P. GALDI, *An introduction to the mathematical theory of the Navier-Stokes equations, Vol. I, Linearized Steady Problems*, Springer Tracts Nat. Philos. 38, Springer-Verlag, New York, 1994.
- [7] M. D. GUNZBURGER, H.-C. LEE, AND G. A. SEREGIN, *Global existence of weak solutions for viscous incompressible flows around a moving rigid body in three dimensions*, *J. Math. Fluid Mech.*, 2 (2000), pp. 219–266.
- [8] T. I. HESLA, *Collisions of Smooth Bodies in Viscous Fluids: A Mathematical Investigation*, Ph.D. thesis, University of Minnesota, revised version, 2005.
- [9] M. HILLAIRET, *Lack of collision between solid bodies in a 2D incompressible viscous flow*, *Comm. Partial Differential Equations*, 32 (2007), pp. 1345–1371.
- [10] A. INOUE AND M. WAKIMOTO, *On existence of solutions of the Navier-Stokes equation in a time dependent domain*, *J. Fac. Sci. Univ. Tokyo Sect. IA Math.*, 24 (1977), pp. 303–319.
- [11] V. N. STAROVOÏTOV, *Nonuniqueness of a solution to the problem on motion of a rigid body in a viscous incompressible fluid*, *J. Math. Sci.*, 130 (2005), pp. 4893–4898.
- [12] V. N. STAROVOITOV, *Behavior of a rigid body in an incompressible viscous fluid near a boundary*, in *Free Boundary Problems (Trento, 2002)*, *Internat. Ser. Numer. Math.* 147, Birkhäuser, Basel, 2004, pp. 313–327.
- [13] T. TAKAHASHI, *Analysis of strong solutions for the equations modeling the motion of a rigid-fluid system in a bounded domain*, *Adv. Differential Equations*, 8 (2003), pp. 1499–1532.
- [14] R. TEMAM, *Problèmes mathématiques en plasticité*, Gauthier-Villars, Montrouge, 1983.
- [15] J. L. VÁZQUEZ AND E. ZUAZUA, *Lack of collision in a simplified 1D model for fluid-solid interaction*, *Math. Models Methods Appl. Sci.*, 16 (2006), pp. 637–678.

ASYMPTOTIC STABILITY OF PERIODIC SOLUTIONS FOR NONSMOOTH DIFFERENTIAL EQUATIONS WITH APPLICATION TO THE NONSMOOTH VAN DER POL OSCILLATOR*

ADRIANA BUICĂ†, JAUME LLIBRE‡, AND OLEG MAKARENKOV§

Abstract. In this paper we study the existence, uniqueness, and asymptotic stability of the periodic solutions of the Lipschitz system $\dot{x} = \varepsilon g(t, x, \varepsilon)$, where $\varepsilon > 0$ is small. Our results extend the classical second Bogoliubov theorem for the existence of stable periodic solutions to nonsmooth differential systems. As an application we prove the existence of asymptotically stable 2π -periodic solutions of the nonsmooth van der Pol oscillator $\ddot{u} + \varepsilon(|u| - 1)\dot{u} + (1 + a\varepsilon)u = \varepsilon\lambda \sin t$. Moreover, we construct the so-called resonance curves that describe the dependence of the amplitude of these solutions as a function of the parameters a and λ . Finally we compare such curves with the resonance curves of the classical van der Pol oscillator $\ddot{u} + \varepsilon(u^2 - 1)\dot{u} + (1 + a\varepsilon)u = \varepsilon\lambda \sin t$.

Key words. periodic solution, asymptotic stability, averaging theory, nonsmooth differential system, nonsmooth van der Pol oscillator

AMS subject classifications. 34C29, 34C25, 47H11

DOI. 10.1137/070701091

1. Introduction. In this paper we study the existence, uniqueness, and asymptotic stability of the T -periodic solutions of the system

$$(1.1) \quad \dot{x} = \varepsilon g(t, x, \varepsilon),$$

where $\varepsilon > 0$ is a small parameter, and the function $g \in C^0(\mathbb{R} \times \mathbb{R}^k \times [0, 1], \mathbb{R}^k)$ is T -periodic in the first variable and locally Lipschitz with respect to the second. For this class of differential systems, the study of the T -periodic solutions can be made using the averaging function

$$(1.2) \quad g_0(v) = \int_0^T g(\tau, v, 0) d\tau$$

and looking for the periodic solutions that starts near some $v_0 \in g_0^{-1}(0)$.

In the case that g is of class C^1 , we recall the stable periodic case of the second Bogoliubov's theorem [6, Chap. 1, section 5, Theorem II] which states *If $\det(g_0)'(v_0) \neq 0$ and $\varepsilon > 0$ is sufficiently small, then system (1.1) has a unique T -periodic solution in a neighborhood of v_0 . Moreover, if all the eigenvalues of the Jacobian matrix*

*Received by the editors August 24, 2007; accepted for publication (in revised form) November 8, 2008; published electronically February 25, 2009.

<http://www.siam.org/journals/sima/40-6/70109.html>

†Department of Applied Mathematics, Babeş-Bolyai University, Cluj-Napoca, Romania (abuica@math.ubbcluj.ro).

‡Departament de Matemàtiques, Universitat Autònoma de Barcelona, 08193 Bellaterra, Barcelona, Spain (llibre@mat.uab.cat). This author was partially supported by MEC/FEDER grant MTM2008-03437 and by CICYT grant 2005SGR 00550.

§Research Institute of Mathematics, Voronezh State University, Voronezh, Russia (omakarenkov@math.vsu.ru). This author was partially supported by grant BF6M10 of the Russian Federation Ministry of Education and U.S. CRDF (BRHE), by RFBR grant 06-01-72552, by the President of the Russian Federation Young Researcher grant MK-1620.2008.1, and by Marie Curie IIF GA-2008-221331.

$(g_0)'(v_0)$ have negative real part, then this periodic solution is asymptotically stable. This theorem has a long history and includes results by Fatou [15], Mandelstam and Papaleksi [31], and Krylov and Bogoliubov [25, section 2].

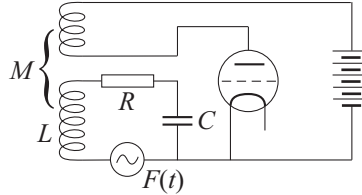


FIG. 1.1. Circuit scheme for the classical triode oscillator (see [3, Chap. VIII, section 2, Figure 348], [30, Chap. I, section 5, Figure 1], and [36, section 3.1.7, Figures 3–5]).

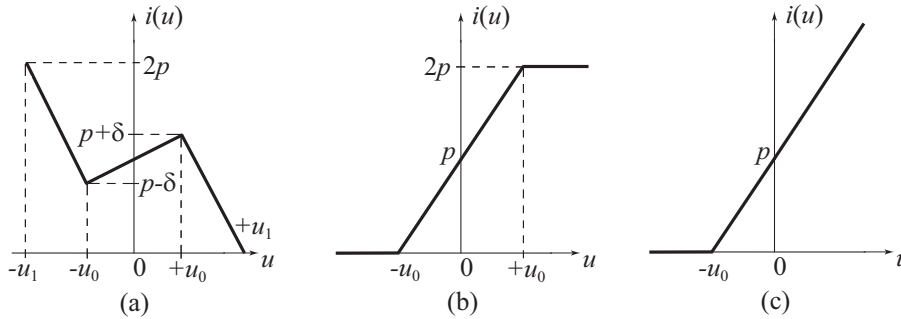


FIG. 1.2. Characteristics of the triode of the circuit of Figure 1.1. (a) Triode in a harsh regime (see [3, Chap. IV, section 7, Figure 212b], [30, Chap. I, section 5, comments for eqs. 5.3–5.4]); (b) triode with saturation (see [3, Chap. VIII, section 3, Figure 364]); (c) triode without saturation (see [3, Chap. IX, section 7, Figure 482]).

The Bogoliubov result provided a theoretical justification for some resonance phenomena which appear in many real physical systems. One of the most significant examples is the classical triode oscillator whose scheme is drawn in Figure 1.1 and whose current u is described by the second order differential equation

$$(1.3) \quad \ddot{u} + \frac{1}{LC} (RC - Mi'(u)) \dot{u} + \omega^2 u = \frac{1}{LC} F(t),$$

where $R = \varepsilon R_0$, $M = \varepsilon M_0$, $\omega^2 = 1 + \varepsilon b$, $F(t) = \varepsilon \lambda \sin t$, $\varepsilon > 0$, is assumed to be small and the triode characteristic $i(u)$ is drawn in Figure 1.2(a). The analysis of the diagram of bifurcation of the periodic solutions in this system is performed in almost every book on nonlinear oscillations (see Andronov, Witt, and Khaikin [3, Chap. VIII, section 2], Malkin [30, Chap. I, section 5], and Nayfeh and Mook [36, section 3.1.7]) but with the smooth approximation $i(u) = i_{(a)}(u) = S_0 + S_1 u - \frac{1}{3} S_3 u^3$ (leading to the classical van der Pol equation). Therefore it is natural to look for a technique that permits one to avoid this smooth approximation and allows one to work with the original shape of the triode characteristic drawn in Figure 1.2(a).

Though the unforced equation (1.3) (i.e., for $F = 0$) with i described by Figure 1.2(b) and Figure 1.2(c) is well studied (see [3, Chap. VIII, section 3 and Chap. IX, section 7]), the question about resonances in these equations when $F \neq 0$ (e.g., $F(t) = \varepsilon \lambda \sin t$) is still partially open. In this direction Levinson [29] uncovered a

family of solutions of (1.3) of remarkable singular structure and Levi [28] completed the study of the limit behavior of all solutions. The present paper complements these results by describing the location of asymptotically stable periodic solutions of (1.3).

Studying (1.3) with the triode characteristic given by Figure 1.2(a), 1.2(b), or 1.2(c) we finally note that there exists a change of variables (see, for example, how Levinson changed equation 2.0 in [29]) that allows one to rewrite (1.3) into the form (1.1) with some function g that is not C^1 but is Lipschitz with respect to the second variable. Therefore the goal of this work is to generalize the results on the existence of a stable periodic solution of the second Bogoliubov theorem to the case that the function g of (1.1) is only Lipschitz.

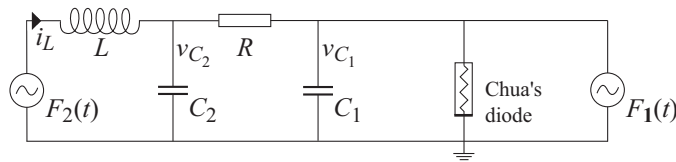


FIG. 1.3. Forced Chua's circuit (see [5, 11, 20, 35, 40]).

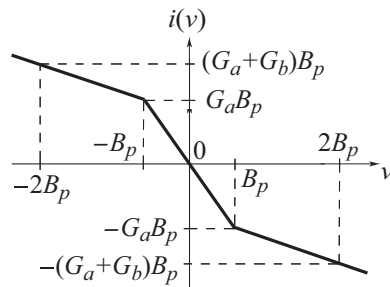


FIG. 1.4. Nonlinear characteristic of the Chua's diode of the circuit drawn at Figure 1.3 given by $i(v) = G_b v + (1/2)(G_a - G_b)(|v + B_p| - |v - B_p|)$, where $G_a, G_b, B_p \in \mathbb{R}$ are some constants depending on the properties of the Chua's diode (see [10]).

Another motivation of this paper comes from the forced Chua's circuit (see Figure 1.3) studied in a large number of papers in the modern electrical engineering. This circuit is described by the three-dimensional system

$$(1.4) \quad \begin{aligned} C_1 \frac{dv_{C_1}}{dt} &= \frac{v_{C_2} - v_{C_1}}{R} - i(v_{C_1}) + F_1(t), \\ C_2 \frac{dv_{C_2}}{dt} &= \frac{v_{C_1} - v_{C_2}}{R} + i_L, \\ L \frac{di_L}{dt} &= -v_{C_2} + F_2(t, v_{C_2}), \end{aligned}$$

where $i(v)$ (the characteristic of the Chua's diode) is a piecewise linear function, as it is represented in Figure 1.4. The recent literature provides insight into the numerical simulations of (1.4) (see [40, 20], where $F_1 \neq 0$ and $F_2 \neq 0$, [5, 35], where $F_1 = 0$ and F_2 is periodic, or [11] where both F_1 and F_2 are periodic). Generalization of the Bogoliubov result for (1.1) with Lipschitz right-hand part will allow for the first

time the theoretical detection of asymptotically stable periodic solutions of (1.4) in the case that C_1 is large enough. Eventually this theoretical analysis may provide new interesting parameters of the forced Chua’s circuit for doing additional numerical experiments.

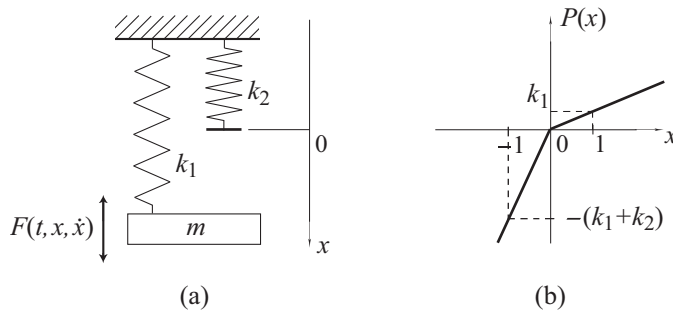


FIG. 1.5. A prototypic device presented in (a) where a driven mass is attached to a immovable beam via a spring with piecewise linear stiffness like in (b); see, e.g., [7], [24, Chap. I, p. 16 and Chap. IV, p. 100], and [38].

On the other hand, part of the interest in generalizing the Bogoliubov result comes from mechanics, where differential systems with piecewise linear stiffness describe various oscillating processes. One of these systems is exhibited by the device drawn in Figure 1.5(a), where a forced mass is attached to a spring whose stiffness changes from k_1 to $k_1 + k_2$ when the mass coordinate crosses 0 in the negative direction. This device is governed by the second order differential equation

$$(1.5) \quad m\ddot{x} + P(x) = F(t, x, \dot{x}),$$

where the piecewise linear stiffness P is drawn in Figure 1.5(b). Depending on the particular configuration of the device of Figure 1.5(a), different expressions for F in (1.5) must be considered. Thus we have that $F(t, x, \dot{x}) = -f(x)\dot{x} + M \cos \omega t$ with piecewise constant f for a shock-absorber and jiggging conveyor (see [24, Chap. I, p. 16 and Chap. IV, p. 100]), where the original Bogoliubov result is employed without justification). The function F takes the simpler form $F(t, x, \dot{x}) = -c\dot{x} + MQ(t)$ for an impact resonator, and $F(t, x, \dot{x}) = -c\dot{x} + M \sin \omega t$ for a cracked-body model (see [38, 7], where only numerical experiments are performed). In each of these situations (1.5) can be rewritten in the form (1.1) with g Lipschitz, provided that the constant k_2 and the amplitude of the force F are sufficiently small. Therefore the extension of the Bogoliubov result to the nonsmooth case that we shall do will allow one to justify the resonances that appeared in all these results. We note that the recent report by Los Alamos National Laboratory [13] describes the increasing interest in a specific form of the model of Figure 1.5(a) called the cracked-body model and, particularly, in the suspension bridge models. Consequently the results of this paper can be applied to such models.

A first model of a one-dimensional suspended bridge is drawn in Figure 1.6(a). It is represented (see [16, 27]) by the beam bending under its own weight and being supported by cables whose restoring force due to elasticity is proportional to u^+ (see Figure 1.6(b)), where $u = u(t, x)$ is the displacement of a point at a distance x from one end of the bridge at time t and u is measured in the downward direction. Looking for u of the form $u(t, x) = z(t) \sin(\pi x/L)$ and considering $F(x, t) = h(t) \sin(\pi x/L)$,

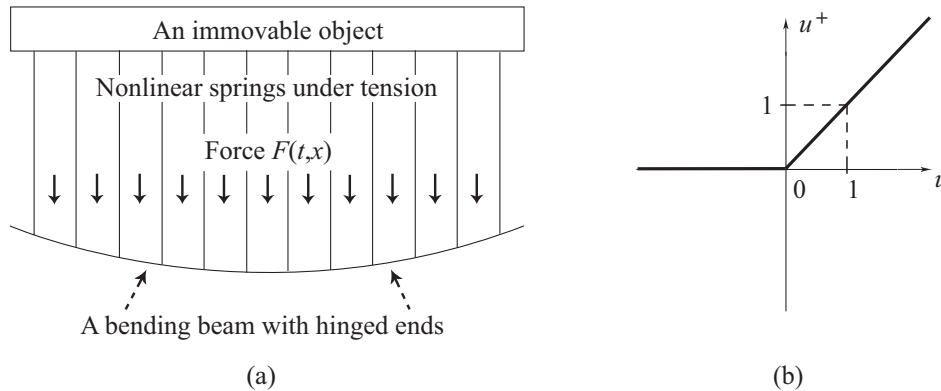


FIG. 1.6. (a) *The first idealization of the suspension bridge: the beam bending under its own weight is supported by the nonlinear cables (see [27, Figure 2]);* (b) *characteristic of stiffness of nonlinear springs.*

we arrive (see [16]) at the following particular case of differential equation (1.5):

$$(1.6) \quad m\ddot{z} + \delta\dot{z} + c(\pi/L)^4 z + dz^+ = mg + h(t),$$

where the constant $m > 0$ is the mass per unit of length, $\delta > 0$ is a small viscous damping coefficient, $c > 0$ measures the flexibility or stiffness of the bridge, $L > 0$ is the length of the bridge, $d > 0$ represents the stiffness of nonlinear springs, and h is a continuous T -periodic force modelling wind, marching troops, or cattle (see [19] for details). Considering $c > 0$ and $d > 0$ fixed and assuming that either $c > 0$ and $h(t)$ are sufficiently small, or that $c > 0$ is fixed and $h(t)$ is sufficiently large, or that $c > 0$ is sufficiently small and $h(t)$ fixed, Glover, Lazer, and McKenna [16], Lazer and McKenna [27], and Fabry [14] proved various theorems on the location of asymptotically stable T -periodic solutions in (1.6). The question *What happens with these solutions when $d > 0$, $\delta > 0$, and $h(t)$ are all sufficiently small?* was open and can be solved using the generalization of the Bogoliubov result that we provide. Lazer and McKenna proved in [26] that the Poincaré map for (1.6) is differentiable, but we note that this is not sufficient for applying the original Bogoliubov result.

We end the list of possible applications noting that system (1.4) describing the Chua's circuit (Figure 1.3) appeared recently for studying the so-called negative slope mechanical systems (see Awrejcewicz [4, section 8.2.2]). So our results can also be applied to these mechanical systems.

These applications require generalizations of the second Bogoliubov theorem for Lipschitz right-hand parts. To the best of our knowledge Mitropol'skii was the first to consider such a kind of generalization. Assuming that g is Lipschitz, $g_0 \in C^3(\mathbb{R}^k, \mathbb{R}^k)$, and all the eigenvalues of the matrix $(g_0)'(v_0)$ have negative real part, Mitropol'skii [34] developed the Bogoliubov result proving the existence and uniqueness of a T -periodic solution of system (1.1) in a neighborhood of v_0 . There was great progress weakening the assumptions of the existence result (see Samoilenko [39] and Mawhin [32]), but this progress did not take place in the case of the uniqueness. Moreover, the asymptotic stability of the T -periodic solution remained unstudied in the case of Lipschitz systems for a long time. It has been done recently by Buică and Daniilidis in [8] for Lipschitz systems (1.1) assuming that the function $v \mapsto g(t, v, 0)$ is differentiable at v_0 for almost any $t \in [0, T]$ and that the eigenvectors of the matrix $(g_0)'(v_0)$ are orthogonal.

In section 2, assuming that g is piecewise differentiable in the second variable, we prove in Theorem 2.5 that the stable periodic solution of the Bogoliubov theorem persists when g is not necessary C^1 . Theorem 2.5 follows from this more general Theorem 2.1 whose hypotheses do not use any differentiability—neither of g , nor of g_0 . Assuming only continuity for g , we show in Theorem 2.9 the existence of a nonasymptotically stable T -periodic solution of system (1.1) if the Brouwer topological degree of $-g_0$ is negative. In section 3 we illustrate our results constructing the resonance curves of the nonsmooth van der Pol oscillator, also studied in [18], and compare these curves with the resonance curves of the classical van der Pol oscillator, which were constructed by Andronov and Witt [1, 2].¹

2. Main results. Throughout the paper $\Omega \subset \mathbb{R}^k$ will be an open set. For any $\delta > 0$ we denote $B_\delta(v_0) = \{v \in \mathbb{R}^k : \|v - v_0\| \leq \delta\}$. We have the following main result on the existence, uniqueness, and asymptotic stability of T -periodic solutions for system (1.1).

THEOREM 2.1. *Let $g \in C^0(\mathbb{R} \times \Omega \times [0, 1], \mathbb{R}^k)$ and $v_0 \in \Omega$. Assume the following four conditions.*

- (i) *For some $L > 0$ we have that $\|g(t, v_1, \varepsilon) - g(t, v_2, \varepsilon)\| \leq L \|v_1 - v_2\|$ for any $t \in [0, T]$, $v_1, v_2 \in \Omega$, $\varepsilon \in [0, 1]$.*
- (ii) *For any $\gamma > 0$ there exists $\delta > 0$ such that*

$$\left\| \int_0^T g(\tau, v_1 + u(\tau), \varepsilon) d\tau - \int_0^T g(\tau, v_2 + u(\tau), \varepsilon) d\tau - \int_0^T g(\tau, v_1, 0) d\tau + \int_0^T g(\tau, v_2, 0) d\tau \right\| \leq \gamma \|v_1 - v_2\|$$

for any $u \in C^0([0, T], \mathbb{R}^k)$, $\|u\| \leq \delta$, $v_1, v_2 \in B_\delta(v_0)$, and $\varepsilon \in [0, \delta]$.

- (iii) *Let g_0 be the averaged function given by (1.2) and consider that $g_0(v_0) = 0$.*
- (iv) *There exist $q \in [0, 1]$, $\alpha, \delta_0 > 0$, and a norm $\|\cdot\|_0$ on \mathbb{R}^k such that $\|v_1 + \alpha g_0(v_1) - v_2 - \alpha g_0(v_2)\|_0 \leq q \|v_1 - v_2\|_0$ for any $v_1, v_2 \in B_{\delta_0}(v_0)$.*

Then there exists $\delta_1 > 0$ such that for every $\varepsilon \in (0, \delta_1]$, system (1.1) has exactly one T -periodic solution x_ε with $x_\varepsilon(0) \in B_{\delta_1}(v_0)$. Moreover, the solution x_ε is asymptotically stable and $x_\varepsilon(0) \rightarrow v_0$ as $\varepsilon \rightarrow 0$.

When the solution $x(\cdot, v, \varepsilon)$ of system (1.1) with the initial condition $x(0, v, \varepsilon) = v$ is well defined on $[0, T]$ for any $v \in B_{\delta_0}(v_0)$, the map $v \mapsto x(T, v, \varepsilon)$ is also well defined and is called the *Poincaré map* at time T of system (1.1). In order to prove the existence, uniqueness, and stability of the T -periodic solutions of system (1.1) stated in Theorem 2.1, it is sufficient to study the same properties for the fixed points of this Poincaré map.

Before proving Theorem 2.1 we state and prove two lemmas. In order to state the first lemma, we need to introduce the function

$$g_\varepsilon(v) = \int_0^T g(\tau, x(\tau, v, \varepsilon), \varepsilon) d\tau$$

¹At the final stage of publishing of this paper, we have been informed by Prof. Michael Guevara that the resonance curves mentioned appeared for the first time in [Balth. van der Pol, *Tijdschr. Ned Rad Gen.*, (1924) (in Dutch)], an English translation appeared in [*Phil. Mag.*, vol. 3, 1927, p. 65]. We thank Prof. Guevara for calling our attention to those papers and some other historical background on resonance curves.

and to note that by writing the equivalent integral equation of system (1.1) we have

$$x(T, v, \varepsilon) = v + \varepsilon g_\varepsilon(v).$$

LEMMA 2.2. *Let $g \in C^0(\mathbb{R} \times \Omega \times [0, 1], \mathbb{R}^k)$ and $\delta_0 > 0$ be such that $B_{\delta_0}(v_0) \subset \Omega$. If (i) is satisfied, then there exist $\delta \in [0, \delta_0]$ and $L_1 > 0$ such that the map $(v, \varepsilon) \mapsto g_\varepsilon(v)$ is well defined and continuous on $B_{\delta_0}(v_0) \times [0, \delta]$ and*

$$\|g_\varepsilon(v_1) - g_\varepsilon(v_2)\| \leq L_1 \|v_1 - v_2\| \quad \text{for any } \varepsilon \in [0, \delta], \quad v_1, v_2 \in B_{\delta_0}(v_0).$$

If both (i) and (ii) are satisfied, then for any $\gamma > 0$ there exists $\delta \in [0, \delta_0]$ such that

$$\|g_\varepsilon(v_1) - g_0(v_1) - g_\varepsilon(v_2) + g_0(v_2)\| \leq \gamma \|v_1 - v_2\|$$

for any $v_1, v_2 \in B_\delta(v_0)$ and $\varepsilon \in [0, \delta]$.

Proof. Using the continuity of the solution of a differential system with respect to the initial data and the parameter (see [37, Chap. 4, section 23, statements G and D]), we obtain the existence of $\varepsilon_0 > 0$ such that $x(t, v, \varepsilon) \in \Omega$ for any $t \in [0, T]$, $v \in B_{\delta_0}(v_0)$, and $\varepsilon \in [0, \varepsilon_0]$. Using the Grönwall–Bellman lemma (see [17, Lemma 6.2] or [12, Chap. II, section 11]) from the representation $x(t, v, \varepsilon) = v + \varepsilon \int_0^t g(\tau, x(\tau, v, \varepsilon), \varepsilon) d\tau$ and the property (i), we obtain $\|x(t, v_1, \varepsilon) - x(t, v_2, \varepsilon)\| \leq e^{\varepsilon LT} \|v_1 - v_2\|$ for all $t \in [0, T]$, $v_1, v_2 \in B_{\delta_0}(v_0)$, and $\varepsilon \in [0, \varepsilon_0]$. Therefore $y(t, v, \varepsilon) = \int_0^t g(\tau, x(\tau, v, \varepsilon), \varepsilon) d\tau$ satisfies the property

$$(2.1) \quad \|y(t, v_1, \varepsilon) - y(t, v_2, \varepsilon)\| \leq L_1 \|v_1 - v_2\|$$

for all $t \in [0, T]$, $v_1, v_2 \in B_{\delta_0}(v_0)$, $\varepsilon \in [0, \varepsilon_0]$, and $L_1 = LT e^{\varepsilon_0 LT}$. Since $g_\varepsilon(v) = y(T, v, \varepsilon)$ the first part of the lemma has been proven.

Taking into account that $x(t, v, \varepsilon) = v + \varepsilon y(t, v, \varepsilon)$, we have

$$(2.2) \quad y(T, v_1, \varepsilon) - y(T, v_1, 0) - y(T, v_2, \varepsilon) + y(T, v_2, 0) = I_1(v_1, v_2, \varepsilon) + I_2(v_1, v_2, \varepsilon),$$

where

$$\begin{aligned} I_1(v_1, v_2, \varepsilon) &= \int_0^T [g(\tau, v_2 + \varepsilon y(\tau, v_1, \varepsilon), \varepsilon) - g(\tau, v_2 + \varepsilon y(\tau, v_2, \varepsilon), \varepsilon)] d\tau, \\ I_2(v_1, v_2, \varepsilon) &= \int_0^T [(g(\tau, v_1 + \varepsilon y(\tau, v_1, \varepsilon), \varepsilon) - g(\tau, v_2 + \varepsilon y(\tau, v_1, \varepsilon), \varepsilon))] d\tau \\ &\quad - \int_0^T (g(\tau, v_1, 0) - g(\tau, v_2, 0)) d\tau. \end{aligned}$$

Since $(t, v, \varepsilon) \mapsto y(t, v, \varepsilon)$ is bounded on $[0, T] \times B_{\delta_0}(v_0) \times [0, \varepsilon_0]$, we have that $\varepsilon y(t, v, \varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$ uniformly with respect to $t \in [0, T]$ and $v \in B_{\delta_0}(v_0)$. Decreasing $\varepsilon_0 > 0$, if necessary, we get that $v_2 + \varepsilon y(t, v_1, \varepsilon) \in \Omega$ for any $t \in [0, T]$, $v_1, v_2 \in B_{\delta_0}(v_0)$, $\varepsilon \in [0, \varepsilon_0]$. By assumption (i) and relation (2.1) we obtain that $\|I_1(v_1, v_2, \varepsilon)\| \leq T \cdot \varepsilon L L_1 \|v_1 - v_2\|$ for all $\varepsilon \in [0, \varepsilon_0]$, $v_1, v_2 \in B_{\delta_0}(v_0)$.

We fix $\gamma > 0$ and take $\delta > 0$ given by (ii). Without loss of generality we can consider that $\delta \leq \min\{\delta_0, \varepsilon_0, \gamma/(2TLL_1)\}$. Therefore assumption (ii) implies that $\|I_2(v_1, v_2, \varepsilon)\| \leq (\gamma/2) \|v_1 - v_2\|$ for any $\varepsilon \in [0, \delta]$, $v_1, v_2 \in B_\delta(v_0)$. Substituting the obtained estimations for I_1 and I_2 into (2.2) we have $\|y(T, v_1, \varepsilon) - y(T, v_1, 0) - y(T, v_2, \varepsilon) + y(T, v_2, 0)\| \leq \gamma \|v_1 - v_2\|$.

$y(T, v_2, 0) \leq (\varepsilon TLL_1 + \gamma/2)\|v_1 - v_2\| \leq \gamma\|v_1 - v_2\|$ for any $\varepsilon \in [0, \delta]$, $v_1, v_2 \in B_\delta(v_0)$. Hence the proof is complete. \square

LEMMA 2.3. *Let $g_0 : \Omega \rightarrow \mathbb{R}^k$, satisfying assumption (iv) for some $q \in (0, 1)$, $\alpha, \delta_0 > 0$, and a norm $\|\cdot\|_0$ on \mathbb{R}^k . Then $\|v_1 + \varepsilon g_0(v_1) - v_2 - \varepsilon g_0(v_2)\|_0 \leq (1 - \varepsilon(1 - q)/\alpha)\|v_1 - v_2\|_0$ for any $v_1, v_2 \in B_{\delta_0}(v_0)$ and any $\varepsilon \in [0, \alpha]$.*

Proof. Indeed, the equality $v + \varepsilon g_0(v) = (1 - \varepsilon/\alpha)v + \varepsilon/\alpha(v + \alpha g_0(v))$ implies that the Lipschitz constant of the function $I + \varepsilon g_0$ with respect to the norm $\|\cdot\|_0$ is $(1 - \varepsilon/\alpha) + \varepsilon/\alpha q = 1 - \varepsilon(1 - q)/\alpha$. \square

Proof of Theorem 2.1. By Lemma 2.2 we have that there exists $\delta_1 \in [0, \delta_0]$ such that

$$(2.3) \quad \|g_\varepsilon(v_1) - g_0(v_1) - g_\varepsilon(v_2) + g_0(v_2)\|_0 \leq ((1 - q)/(2\alpha))\|v_1 - v_2\|_0$$

for any $\varepsilon \in [0, \delta_1]$, $v_1, v_2 \in B_{\delta_1}(v_0)$. First we prove that there exists $\varepsilon_1 \in [0, \delta_1]$ such that for every $\varepsilon \in [0, \varepsilon_1]$ there exists $v_\varepsilon \in B_{\delta_1}(v_0)$ such that $x(\cdot, v_\varepsilon, \varepsilon)$ is a T -periodic solution of (1.1) by showing that there exists v_ε such that $x(T, v_\varepsilon, \varepsilon) = v_\varepsilon$. Using (iii) and (iv) we have

$$\|v + \alpha g_0(v) - v_0\|_0 \leq q\|v - v_0\|_0 \quad \text{for any } v \in B_{\delta_1}(v_0).$$

Therefore we have that the map $I + \alpha g_0$ maps $B_{\delta_1}(v_0)$ into itself. From Lemma 2.2 we have that there exists $\varepsilon_0 > 0$ such that the map $(v, \varepsilon) \mapsto g_\varepsilon(v)$ is well defined and continuous on $B_{\delta_1}(v_0) \times [0, \varepsilon_0]$. We deduce that there exists $\varepsilon_1 > 0$ sufficiently small such that, for every $\varepsilon \in [0, \varepsilon_1]$, the map $I + \alpha g_\varepsilon$ maps $B_{\delta_1}(v_0)$ into itself as well. Therefore, by the Brouwer theorem (see, for example, [23, Theorem 3.1]) we have that $B_{\delta_1}(v_0)$ contains at least one fixed point of the map $I + \alpha g_\varepsilon$ for any $\varepsilon \in [0, \varepsilon_1]$. Denote this fixed point by v_ε . Then we have $g_\varepsilon(v_\varepsilon) = 0$ and $x(T, v_\varepsilon, \varepsilon) = v_\varepsilon$ for any $\varepsilon \in [0, \varepsilon_1]$.

Now we prove that $x(\cdot, v_\varepsilon, \varepsilon)$ is the only T -periodic solution of (1.1) starting near v_0 and that, moreover, it is asymptotically stable. Knowing that $x(T, v, \varepsilon) = v + \varepsilon g_\varepsilon(v)$ we write the following identity:

$$(2.4) \quad x(T, v, \varepsilon) = v + \varepsilon g_0(v) + \varepsilon(g_\varepsilon(v) - g_0(v)).$$

Using Lemma 2.3 we have from (2.3) and (2.4) that

$$\begin{aligned} \|x(T, v_1, \varepsilon) - x(T, v_2, \varepsilon)\|_0 &\leq (1 - \varepsilon(1 - q)/\alpha + \varepsilon(1 - q)/(2\alpha))\|v_1 - v_2\|_0 \\ &= (1 - \varepsilon(1 - q)/(2\alpha))\|v_1 - v_2\|_0 \end{aligned}$$

for all $v_1, v_2 \in B_{\delta_1}(v_0)$ and $\varepsilon \in [0, \delta_1]$. We proved before that there exists $\varepsilon_1 > 0$ such that for every $\varepsilon \in [0, \varepsilon_1]$ there exists $v_\varepsilon \in B_{\delta_1}(v_0)$ with $x(\cdot, v_\varepsilon, \varepsilon)$ a T -periodic solution of (1.1). Since $\varepsilon(1 - q)/(2\alpha) > 0$ and $\varepsilon_1 \leq \delta_1$, the last inequality implies that for each $\varepsilon \in [0, \delta_1]$, the T -periodic solution $x(\cdot, v_\varepsilon, \varepsilon)$ is the only T -periodic solution of (1.1) in $B_{\delta_1}(v_0)$ and, moreover (see [23, Lemma 9.2]), it is asymptotically stable. \square

Remark 2.4. We note that a similar result close to Theorem 2.1 is contained in [8, Theorem 3.5]. But instead of the assumption (iv) with a fixed $\alpha > 0$, it is assumed in [8] with any $\alpha > 0$ sufficiently small. Anyway, notice that Lemma 2.3 implies that it is the same to assume (iv) for only one $\alpha > 0$ or for all $\alpha > 0$ sufficiently small. The advantage of our Theorem 2.1 is that it does not require differentiability of $g(t, \cdot, \varepsilon)$ at any point, while [8] needs it at v_0 . See also Remark 2.8.

In general it is not easy to check assumptions (ii) and (iv) in the applications of Theorem 2.1. Thus we also give the following theorem based on Theorem 2.1 which assumes certain type of piecewise differentiability instead of (ii) and deals with properties of the matrix $(g_0)'(v_0)$ instead of the Lipschitz constant of g_0 .

For any set $M \subset [0, T]$ measurable in the sense of Lebesgue we denote by $\text{mes}(M)$ the Lebesgue measure of M (see [21, Chap. V, section 3]).

THEOREM 2.5. *Let $g \in C^0(\mathbb{R} \times \Omega \times [0, 1], \mathbb{R}^k)$ satisfying (i). Let g_0 be the averaged function given by (1.2) and consider $v_0 \in \Omega$ such that $g_0(v_0) = 0$. Assume that*

- (v) *given any $\tilde{\gamma} > 0$ there exist $\tilde{\delta} > 0$ and $M \subset [0, T]$ with $\text{mes}(M) < \tilde{\gamma}$ such that for every $v \in B_{\tilde{\delta}}(v_0)$, $t \in [0, T] \setminus M$, and $\varepsilon \in [0, \tilde{\delta}]$ we have that $g(t, \cdot, \varepsilon)$ is differentiable at v and $\|g'_v(t, v, \varepsilon) - g'_v(t, v_0, 0)\| \leq \tilde{\gamma}$.*

Finally assume that

- (vi) *g_0 is continuously differentiable in a neighborhood of v_0 and the real parts of all the eigenvalues of $(g_0)'(v_0)$ are negative.*

Then there exists $\delta_1 > 0$ such that for every $\varepsilon \in (0, \delta_1]$, system (1.1) has exactly one T -periodic solution x_ε with $x_\varepsilon(0) \in B_{\delta_1}(v_0)$. Moreover, the solution x_ε is asymptotically stable and $x_\varepsilon(0) \rightarrow v_0$ as $\varepsilon \rightarrow 0$.

For proving Theorem 2.5 we need two preliminary lemmas.

LEMMA 2.6. *Let $g \in C^0(\mathbb{R} \times \Omega \times [0, 1], \mathbb{R}^k)$ satisfying (i). If (v) holds, then (ii) is satisfied.*

Proof. Let $\gamma > 0$ be an arbitrary number. We show that (ii) holds with $\delta = \tilde{\delta}/2$, where $\tilde{\delta}$ is given by (v) applied with $\tilde{\gamma} = \min\{\gamma/(4L), \gamma/(4T)\}$. We consider also $M \subset [0, T]$ given by (v) applied with the same value of $\tilde{\gamma}$.

Let $u \in C^0([0, T], \mathbb{R}^k)$, $\|u\| \leq \delta$, and $F(v) = \int_0^T g(\tau, v+u(\tau), \varepsilon)d\tau - \int_0^T g(\tau, v, 0)d\tau$. Let $v_1, v_2 \in B_\delta(v_0)$ and $\varepsilon \in [0, \delta]$. We have $F(v) = F_1(v) + F_2(v)$, where $F_1(v) = \int_M (g(\tau, v+u(\tau), \varepsilon) - g(\tau, v, 0))d\tau$ and $F_2(v) = \int_{[0, T] \setminus M} (g(\tau, v+u(\tau), \varepsilon) - g(\tau, v, 0))d\tau$. By (i) we have that $\|F_1(v_1) - F_1(v_2)\| \leq 2L \cdot \text{mes}(M)\|v_1 - v_2\| < 2L\tilde{\gamma}\|v_1 - v_2\| \leq (\gamma/2)\|v_1 - v_2\|$. On the other hand, using (v), we will prove that a similar relation holds for F_2 . In order to do this, we denote $h(\tau, v) = g(\tau, v+u(\tau), \varepsilon) - g(\tau, v, 0)$. Notice that for each $\tau \in [0, T] \setminus M$ we can write $h'_v(\tau, v) = (g'_v(\tau, v+u(\tau), \varepsilon) - g'_v(\tau, v_0, 0)) - (g'_v(\tau, v, 0) - g'_v(\tau, v_0, 0))$. As a direct consequence of (v) we deduce that $\|h'_v(\tau, v)\| \leq 2\tilde{\gamma}$ for all $v \in B_\delta(v_0)$ and $\tau \in [0, T] \setminus M$. Now applying the mean value theorem for the function $h(\tau, \cdot)$, we have $\|h(\tau, v_1) - h(\tau, v_2)\| \leq 2\tilde{\gamma}\|v_1 - v_2\|$ for all $\tau \in [0, T] \setminus M$ and all $v_1, v_2 \in B_\delta(v_0)$. Then $\|F_2(v_1) - F_2(v_2)\| \leq \int_{[0, T] \setminus M} \|h(\tau, v_1) - h(\tau, v_2)\|d\tau \leq 2T\tilde{\gamma}\|v_1 - v_2\| \leq (\gamma/2)\|v_1 - v_2\|$. Therefore we have proved that $\|F(v_1) - F(v_2)\| \leq \gamma\|v_1 - v_2\|$, which coincides with (ii). \square

LEMMA 2.7. *Let $g_0 : \Omega \rightarrow \mathbb{R}^k$ satisfying assumption (vi) for some $v_0 \in \Omega$. Then there exist $q \in [0, 1)$, $\alpha, \delta_0 > 0$ and a norm $\|\cdot\|_0$ on \mathbb{R}^k such that (iv) is satisfied.*

Proof. If λ is an eigenvalue of $\alpha(g_0)'(v_0)$, then $\lambda + 1$ is an eigenvalue of $I + (\alpha g_0)'(v_0)$. Since the eigenvalues of $\alpha(g_0)'(v_0)$ tend to 0 as $\alpha \rightarrow 0$ and have negative real parts, then there exists $\alpha \in [0, 1)$ such that the absolute values of all the eigenvalues of $I + \alpha(g_0)'(v_0)$ are less than one. Therefore (see [22, p. 90, Lemma 2.2]) there exist $\tilde{q} \in [0, 1)$ and a norm $\|\cdot\|_0$ on \mathbb{R}^k such that $\sup_{\|\xi\|_0 \leq 1} \|\xi + \alpha(g_0)'(v_0)\xi\|_0 \leq \tilde{q}$.

By continuous differentiability of g_0 in a neighborhood of v_0 we have that $\|g_0(v_1) - g_0(v_2) - (g_0)'(v_0)(v_1 - v_2)\| / \|v_1 - v_2\| \leq \|g_0(v_1) - g_0(v_2) - (g_0)'(v_2)(v_1 - v_2)\| + \|(g_0)'(v_2)(v_1 - v_2) - (g_0)'(v_0)(v_1 - v_2)\| / \|v_1 - v_2\| \rightarrow 0$ as $\max\{\|v_1 - v_0\|, \|v_2 - v_0\|\} \rightarrow 0$. Therefore taking into account that all norms on \mathbb{R}^k are equivalent, there exists $\delta_0 > 0$ such that $\|g_0(v_1) - g_0(v_2) - (g_0)'(v_0)(v_1 - v_2)\|_0 \leq (1 - \tilde{q})/(2\alpha)\|v_1 - v_2\|_0$ for all

$v_1, v_2 \in B_{\delta_0}(v_0)$. Then

$$\begin{aligned} & \|v_1 + \alpha g_0(v_1) - v_2 - \alpha g_0(v_2)\|_0 \\ & \leq \alpha \|g_0(v_1) - g_0(v_2) - (g_0)'(v_0)(v_1 - v_2)\|_0 + \|v_1 - v_2 + \alpha (g_0)'(v_0)(v_1 - v_2)\|_0 \\ & \leq (1 + \tilde{q})/2 \|v_1 - v_2\|_0 \end{aligned}$$

for all $v_1, v_2 \in B_{\delta_0}(v_0)$. \square

Proof of Theorem 2.5. Lemmas 2.6 and 2.7 imply that assumptions (ii) and (iv) of Theorem 2.1 are satisfied. Therefore the conclusion of the theorem follows by applying Theorem 2.1. \square

It was observed by Mitropol'skii in [34] that in spite of the fact that $g(t, \cdot, \varepsilon)$ in (1.1) is only Lipschitz, sometimes the function g_0 turns out to be differentiable in applications. In particular we will see in section 3 that this is the case for the nonsmooth van der Pol oscillator.

Clearly if $g \in C^1(\mathbb{R} \times \mathbb{R}^k \times [0, 1], \mathbb{R}^k)$, then (i) and (v) hold in any open bounded set $\Omega \subset \mathbb{R}^k$. Therefore Theorem 2.5 is a generalization of the stable periodic case of the second Bogoliubov theorem formulated in the introduction.

Remark 2.8. Theorem 2.5 does not require the eigenvectors of $(g_0)'(v_0)$ to be orthogonal as in [8, Theorem 3.6]. Moreover, assumption (H_2) of [8] is more restrictive than (v).

For completeness we also give the following theorem on the existence of nonasymptotically stable T -periodic solutions for (1.1). In the theorem below, $d(F, V)$ denotes the Brouwer topological degree of the vector field $F \in C^0(\mathbb{R}^k, \mathbb{R}^k)$ on the open and bounded set $V \subset \mathbb{R}^k$ (see [23, Chap. 2, section 5.2]).

THEOREM 2.9. *Let $g \in C^0(\mathbb{R} \times \mathbb{R}^k \times [0, 1], \mathbb{R}^k)$. Assume that there exists an open bounded set $V \subset \mathbb{R}^k$ such that $g_0(v) \neq 0$ for any $v \in \partial V$ and*

(vii) $d(-g_0, V) < 0$.

Then there exists $\varepsilon_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0]$ system (1.1) has at least one nonasymptotically stable T -periodic solutions x_ε with $x_\varepsilon(0) \in V$.

Proof. Since $g_0(v) \neq 0$ for any $v \in \partial V$, then from Mawhin's theorem [32] (or [33, section 5]) we have that there exists $\varepsilon_0 > 0$ such that

(2.5) $d(-g_0, V) = d(I - x(T, \cdot, \varepsilon), V)$ for any $\varepsilon \in (0, \varepsilon_0]$.

By [23, Theorem 9.6] for any asymptotically stable T -periodic solution x_ε of (1.1) we have that $d(I - x(T, \cdot, \varepsilon), B_\delta(x_\varepsilon(0))) = 1$ for $\delta > 0$ sufficiently small. Therefore if all the possible T -periodic solutions of (1.1) with $\varepsilon \in (0, \varepsilon_0]$ had been asymptotically stable, then the degree $d(I - x(T, \cdot, \varepsilon), V)$ would have been nonnegative, contradicting (vii) and (2.5). \square

Remark 2.10. Assumptions (iii) and (iv) imply that $d(-g_0, V) = 1$ (see [23, Theorem 5.16]).

Finally thinking in terms of the application to the nonsmooth van der Pol oscillator, we formulate the following theorem which combines Mawhin's theorem (see [32] or [33, Theorem 3]) and Theorems 2.5 and 2.9. In this theorem $([g_0]_i)'_{(j)}$ stays for the derivative of the i th component of the function g_0 with respect to the j th variable.

THEOREM 2.11. *Let $g \in C^0(\mathbb{R} \times \Omega \times [0, 1], \mathbb{R}^2)$. Let $v_0 \in \Omega$ be such a point that $g_0(v_0) = 0$ and g_0 is continuously differentiable in a neighborhood of v_0 .*

- (a) *If $\det(g_0)'(v_0) \neq 0$, then there exists $\varepsilon_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0]$ system (1.1) has at least one T -periodic solution x_ε satisfying $x_\varepsilon(0) \rightarrow v_0$ as $\varepsilon \rightarrow 0$.*

(b) If (i) and (v) hold and

$$(2.6) \quad \det (g_0)'(v_0) > 0 \quad \text{and} \quad ([g_0]_1)'_{(1)}(v_0) + ([g_0]_2)'_{(2)}(v_0) < 0,$$

then there exists $\varepsilon_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0]$ system (1.1) has exactly one T -periodic solution x_ε such that $x_\varepsilon(0) \rightarrow v_0$ as $\varepsilon \rightarrow 0$. Moreover, the solution x_ε is asymptotically stable.

(c) If $\det (g_0)'(v_0) < 0$, then there exists $\varepsilon_0 > 0$ such that for any $\varepsilon \in (0, \varepsilon_0]$ system (1.1) has at least one nonasymptotically stable T -periodic solution x_ε such that $x_\varepsilon(0) \rightarrow v_0$ as $\varepsilon \rightarrow 0$.

Proof. Statement (a) is added for the completeness of the formulation of Theorem 2.11 and it follows from Mawhin’s theorem (see [32] or [33, Theorem 3]).

On the other hand, it is a simple calculation to show that (2.6) implies that all the eigenvalues of $(g_0)'(v_0)$ have negative real part. Therefore assumption (vi) of Theorem 2.5 is also satisfied and statement (b) follows from this theorem.

Statement (c) follows from Theorem 2.9. Indeed, since $\det (g_0)'(v_0) < 0$, it implies (see [23, Theorem 5.9]) that $d(g_0, B_\rho(v_0))$ is defined for any $\rho > 0$ sufficiently small and that $d(g_0, B_\rho(v_0)) = \det(g_0)'(v_0) < 0$. \square

3. Application to the nonsmooth van der Pol oscillator. In [18] Hogan first demonstrated the existence of a limit cycle for the nonsmooth van der Pol equation $\ddot{u} + \varepsilon(|u| - 1)\dot{u} + u = 0$. This equation governs the circuit drawn at Figure 1.1 with the triode characteristic $i(u) = S_0 + S_1u - S_2u|u|$ whose derivative $i'(u) = S_1 - 2S_2|u|$ is nondifferentiable (see Nayfeh and Mook [36, section 3.3.4], where the same stiffness characteristic appears in mechanics). In this paper we extend this study by considering the van der Pol problem on the location of stable and unstable periodic solutions of the perturbed equation

$$(3.1) \quad \ddot{u} + \varepsilon(|u| - 1)\dot{u} + (1 + a\varepsilon)u = \varepsilon\lambda \sin t,$$

where a is a detuning parameter and $\varepsilon\lambda \sin t$ is an external force. We assume that $\varepsilon > 0$ is sufficiently small, and we consider that the parameters a and λ vary in \mathbb{R} .

We finally note that standard change of variables (see example in section 3) brings (1.3) into the form (1.1), with g sufficiently smooth to satisfy the hypotheses of the second Bogoliubov theorem. But we remind the reader that our aim is to apply directly Theorem 2.5, in the same way that Andronov and Witt applied Bogoliubov theorem to the classical van der Pol oscillator

$$(3.2) \quad \ddot{u} + \varepsilon(u^2 - 1)\dot{u} + (1 + a\varepsilon)u = \varepsilon\lambda \sin t,$$

which can be found in [1, Figure 4] or in [30, Chap. I, section 16, Figure 15].

A function u is a solution of (3.1) if and only if $(z_1, z_2) = (u, \dot{u})$ is a solution of the system

$$(3.3) \quad \begin{aligned} \dot{z}_1 &= z_2, \\ \dot{z}_2 &= -z_1 + \varepsilon[-az_1 - (|z_1| - 1)z_2 + \lambda \sin t]. \end{aligned}$$

After the change of variables

$$\begin{pmatrix} z_1(t) \\ z_2(t) \end{pmatrix} = \begin{pmatrix} \cos t & \sin t \\ -\sin t & \cos t \end{pmatrix} \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix},$$

system (3.3) takes the form

$$\begin{aligned}
 \dot{x}_1 &= \varepsilon \sin(-t) [-a(x_1 \cos t + x_2 \sin t) \\
 &\quad - (|x_1 \cos t + x_2 \sin t| - 1)(-x_1 \sin t + x_2 \cos t) + \lambda \sin t], \\
 \dot{x}_2 &= \varepsilon \cos(-t) [-a(x_1 \cos t + x_2 \sin t) \\
 &\quad - (|x_1 \cos t + x_2 \sin t| - 1)(-x_1 \sin t + x_2 \cos t) + \lambda \sin t].
 \end{aligned}
 \tag{3.4}$$

The corresponding averaged function g_0 , calculated using (1.2), is given by

$$\begin{aligned}
 [g_0]_1(M, N) &= \pi a N - \pi \lambda + \pi M - \frac{4}{3} M \sqrt{M^2 + N^2}, \\
 [g_0]_2(M, N) &= -\pi a M + \pi N - \frac{4}{3} N \sqrt{M^2 + N^2},
 \end{aligned}
 \tag{3.5}$$

and it is continuously differentiable in $\mathbb{R}^2 \setminus \{0\}$.

In short, by statement (a) of Theorem 2.11, the zeros $(M, N) \in \mathbb{R}^2$ of this function with the property that $\det(g_0)'(M, N) \neq 0$ determine the 2π -periodic solutions of (3.3) emanating from the solution of the unperturbed system

$$\begin{aligned}
 u_1(t) &= M \cos t + N \sin t, \\
 u_2(t) &= -M \sin t + N \cos t.
 \end{aligned}
 \tag{3.6}$$

We have the following expression for the determinant:

$$\det(g_0)'(M, N) = \pi^2(1 + a^2) + \frac{32}{9}(M^2 + N^2) - 4\pi\sqrt{M^2 + N^2}.
 \tag{3.7}$$

Following Andronov and Witt [1] we are concerned with the dependence of the amplitude of the solution (3.6) with respect to a and λ . Thus we decompose this solution as follows:

$$u_1(t) = A \sin(t + \phi), \quad u_2(t) = A \cos(t + \phi),
 \tag{3.8}$$

where (M, N) is related to (A, ϕ) by

$$M = A \sin \phi, \quad N = A \cos \phi.
 \tag{3.9}$$

Substituting (3.9) into (3.5) and (3.7) we obtain

$$\begin{aligned}
 [g_0((A \sin \phi, A \cos \phi))]_1 &= -(4/3) \cdot A|A| \sin \phi + \pi a A \cos \phi + \pi A \sin \phi - \pi \lambda, \\
 [g_0((A \sin \phi, A \cos \phi))]_2 &= -(4/3) \cdot A|A| \cos \phi - \pi a A \sin \phi + \pi A \cos \phi,
 \end{aligned}
 \tag{3.10}$$

and, respectively,

$$\det(g_0)'((A \sin \phi, A \cos \phi)) = \pi^2(1 + a^2) + \frac{32}{9}A^2 - 2\pi|A|.
 \tag{3.11}$$

Looking for the zeros (A, ϕ) of (3.10), we find the implicit formula

$$A^2 \left(a^2 + \left(1 - \frac{4}{3\pi}|A| \right)^2 \right) = \lambda^2
 \tag{3.12}$$

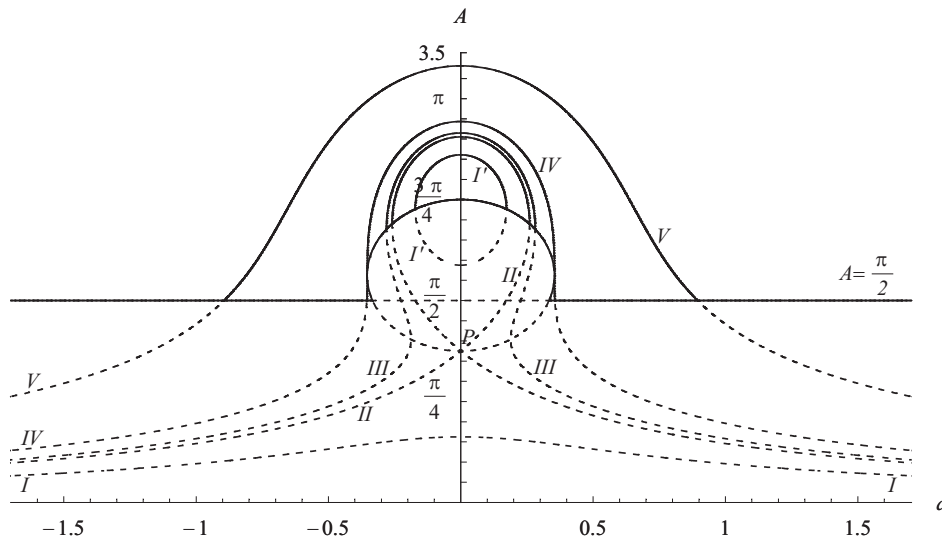


FIG. 3.1. Dependence of the amplitude of stable (solid curves) and unstable (dash curves) 2π -periodic solutions of the nonsmooth periodically perturbed van der Pol equation (3.1) on the detuning parameter a obtained over formulas (3.12), (3.16), and (3.17) for different values of λ . The curve I is plotted with $\lambda = 0.4$, II with $\lambda = 3\pi/16$, III with some $\lambda = \sqrt{0.4} \in (3\pi/16, 9\sqrt{3}\pi/64)$, IV with $\lambda = 9\sqrt{3}\pi/64$, and V with $\lambda = 1.5$. Point P is $2/\sqrt{3}$.

for determining A . Observe that the number of positive zeros of (3.12) coincide with the number of zeros of the equation $A^2(a^2 + (1 - \frac{4}{3\pi}A)^2) = \lambda^2$. To estimate this number we define

$$f(A) = A^2 \left(a^2 + \left(1 - \frac{4}{3\pi}A \right)^2 \right) - \lambda^2,$$

and we have

$$f'(A) = 2A \left(a^2 + \left(1 - \frac{4}{3\pi}A \right)^2 \right) - \frac{8}{3\pi}A^2 \left(1 - \frac{4}{3\pi}A \right).$$

Since f' has one or two zeros, then (3.12) has one, two, or three positive solutions A for any fixed a and λ . In order to understand the different situations that can appear, we follow Andronov and Witt, who suggested in [1] (see also [2]) to construct the so-called *resonance curves*, namely, the curves A in function of a , for λ fixed. The equation of this curve is given by formula (3.12). Some curves (3.12) corresponding to different values of λ are drawn in Figure 3.1. The way for describing these resonance curves (3.12) is borrowed from [30, Ch. 1, section 5], where the classical van der Pol equation is considered.

When $\lambda = 0$ the curve (3.12) is formed by the axis $A = 0$ and the isolated point $(0, 3\pi/4)$. When $\lambda > 0$ but sufficiently small, the resonance curve consists of two branches: instead of $A = 0$ we have the curve of the type $I - I$, and instead of the point $(0, 3\pi/4)$ we obtain an oval $I' - I'$ surrounding this point. When $\lambda > 0$ increases, the oval $I' - I'$ and the branch $I - I$ tend to each other, and, for a certain λ , there exists only one branch $II - II$ with a double point P . The value of this λ can be obtained assuming that (3.12) has for $a = 0$ a double root and, therefore, (3.11)

should be zero. Solving jointly (3.12) and (3.11) with $a = 0$ we obtain $\lambda = 3\pi/16$ and $P = 2\pi/8$. If $\lambda > 3\pi/16$, then we have curves of the type *III* which take form *V* when $\lambda > 0$ crosses the value $\lambda = 9\sqrt{3}\pi/64$. From here, if $\lambda < 3\pi/16$, then (3.12) has three real roots when $|a|$ is sufficiently small, and only one root when $|a|$ is greater than a certain number which depends on λ . When $3\pi/16 < \lambda < 9\sqrt{3}\pi/64$ (3.12) has one, three, and again one solution according to whether $a < a_1$, $a_1 < a < a_2$, and $a > a_2$, respectively, where a_1, a_2 depend on λ . The amplitude curves of type *V* provide exactly one solution of (3.12) for any value of a . The value $\lambda = 9\sqrt{3}\pi/64$, which separates the curves where (3.12) has three solutions from the curves where (3.12) has one solution, is obtained from the property that (3.12) with this λ has a double root for some a and thus this value of a vanishes (3.11). Therefore $\lambda = 9\sqrt{3}\pi/64$ is the point separating the interval $(0, \lambda)$ where the system formed by (3.12) and

$$(3.13) \quad \pi^2(1 + a^2) + \frac{32}{9}A^2 - 2\pi|A| = 0$$

has at least one solution from the interval (λ, ∞) where (3.12)–(3.13) has no solutions.

In short, we have studied the amplitudes of the 2π -periodic solutions of system (3.3) depending on a and λ , when a physical system described by (3.3) possesses 2π -periodic oscillations and when some of them are asymptotically stable. To find the answer we have used statement (b) of Theorem 2.11. Assumption (i) is obviously satisfied with $\Omega = \mathbb{R}^2$. The next statement shows that the right-hand side of system (3.4) satisfies (v).

PROPOSITION 3.1. *Let $v_0 \in \mathbb{R}^2, v_0 \neq 0$. Then the right-hand side of (3.4) satisfies (v) for any $a, \lambda \in \mathbb{R}$.*

The proof of Proposition 3.1 is given in section 4.

Thus we have to study the signs of (3.11) and $([g_0]_1)'_M(A \sin \phi, A \cos \phi) + ([g_0]_2)'_N(A \sin \phi, A \cos \phi)$. We have

$$([g_0]_1)'_M(M, N) + ([g_0]_2)'_N(M, N) = 2 \left(\pi - 2\sqrt{M^2 + N^2} \right),$$

and therefore the conditions for the asymptotic stability of the 2π -periodic solutions of (3.3) near (3.6) are

$$(3.14) \quad \pi^2(1 + a^2) + \frac{32}{9}(M^2 + N^2) - 4\pi\sqrt{M^2 + N^2} > 0$$

and

$$(3.15) \quad 2 \left(\pi - 2\sqrt{M^2 + N^2} \right) < 0.$$

Substituting (3.9) into the inequalities (3.14) and (3.15), we obtain the following equivalent inequalities in terms of the amplitude A :

$$(3.16) \quad \pi^2(1 + a^2) + \frac{32}{9}A^2 - 2\pi|A| > 0$$

and

$$(3.17) \quad 2\pi - 4|A| < 0.$$

Conditions (3.16) and (3.17) mean that the asymptotically stable 2π -periodic solutions of (3.3) correspond to those parts of resonance curves under consideration which are

outside the ellipse (3.13) and above the line $A = \pi/2$. All the results are collected in Figure 3.1, where it is easy to see that for any detuning parameter a and any amplitude $\lambda > 0$, (3.1) possesses at least one asymptotically stable 2π -periodic solution with amplitude close to A obtained from (3.12). Among all the asymptotically stable 2π -periodic solutions of (3.1), there exists exactly one whose fixed neighborhood does not contain any nonasymptotically stable 2π -periodic solution of (3.1) for sufficiently small $\varepsilon > 0$. The amplitude of this asymptotically stable 2π -periodic solution is obtained from (3.16)–(3.17).

To compare the changes due to nonsmoothness in the behavior of the resonance curves, we give in Figure 3.2 the resonance curves of the classical van der Pol oscillator (3.2).

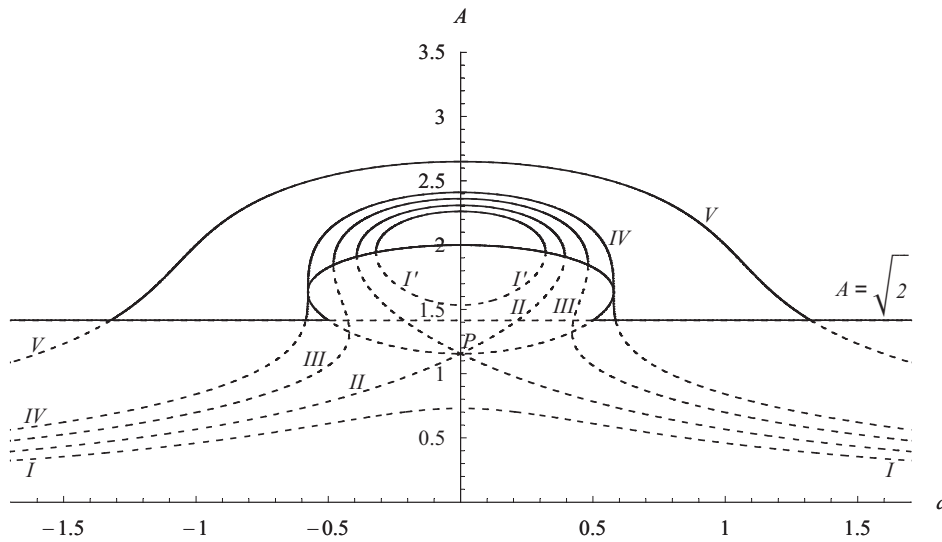


FIG. 3.2. Dependence of the amplitude of stable (solid curves) and unstable (dash curves) 2π -periodic solutions of the classical periodically perturbed van der Pol equation (3.2) on the detuning parameter a for different values of λ . Following Andronov and Witt (see [1, Figure 4]), curve I is plotted with $\lambda = \sqrt{0.4}$, II with $\lambda = 4\sqrt{3}/9$, III with some $4\sqrt{3}/9 < \lambda < \sqrt{32/27}$, IV with $\lambda = \sqrt{32/27}$, and V with $\lambda = 2$. Point P is $2/\sqrt{3}$.

The formulas of Figure 3.1 can be compared with the formulas for Figure 3.2. In fact, the corresponding expressions (3.12)–(3.13) and (3.14)–(3.15) are (see the formulas (5.21)–(5.22) and (16.6)–(16.7) of [30])

$$A^2 \left(a^2 + \left(1 - \frac{A^2}{4} \right)^2 \right) = \lambda^2,$$

$$1 - a^2 - A^2 + \frac{3}{16}A^4 = 0,$$

and

$$1 + a^2 - (M^2 + N^2) + \frac{3}{16}(M^2 + N^2)^2 > 0,$$

$$2 - (M^2 + N^2) < 0,$$

respectively, when we consider the classical van der Pol equation (3.2).

It can be checked that the eigenvectors of the matrix $(g_0)'((A \sin \phi, A \cos \phi))$ are orthogonal only for $A = 0$, so Theorem 3.6 from [8] cannot be applied. At the same time assumption (H2) from [8] is not satisfied for our problem (see Remark 2.8).

4. Appendix.

Proof of Proposition 3.1. As before, $[v]_i$ is the i th component of the vector $v \in \mathbb{R}^2$. Let $g(t, v) = |[v]_1 \cos t + [v]_2 \sin t|$, and notice that it is enough to prove that $g : [0, 2\pi] \times \mathbb{R}^2 \rightarrow \mathbb{R}$ satisfies (v). In the case that $[v_0]_2 \neq 0$, denote $\theta(v) = \arctan(-[v]_1/[v]_2)$ if $[v_0]_2 = 0$, denote $\theta(v) = \arctan(-[v]_1/[v]_2)$ if $[v_0]_1[v]_2 < 0$, $\theta(v_0) = \pi/2$, and, respectively, $\theta(v) = \arctan(-[v]_1/[v]_2) + \pi$ if $[v_0]_1[v]_2 > 0$. In any case notice that the function $v \mapsto \theta(v)$ is continuous in every sufficiently small neighborhood of v_0 . Fix $\tilde{\gamma} > 0$. Let M be the union of two intervals centered in $\theta(v_0)$ (respectively, $\theta(v_0) + 2\pi$ if $\theta(v_0) < 0$) and in $\theta(v_0) + \pi$, each of length $\tilde{\gamma}/2$. Denote them M_1 and M_2 . Take $\tilde{\delta} > 0$ such that $\theta(v) \in M_1$ for all $v \in B_{\tilde{\delta}}(v_0)$. Of course, also $\theta(v) + \pi \in M_2$ for all $\|v - v_0\| \leq \tilde{\delta}$. This implies that for fixed $t \in [0, 2\pi] \setminus M$, $[v]_1 \cos t + [v]_2 \sin t$ has constant sign for all $v \in B_{\tilde{\delta}}(v_0)$, which further gives that $g(t, \cdot)$ is differentiable and $g'_v(t, v) = g'_v(t, v_0)$ for all $v \in B_{\tilde{\delta}}(v_0)$. Hence (v) is fulfilled. \square

Acknowledgments. The authors are grateful to the anonymous referee who motivated us to include a list of several applications of Theorem 2.5 in the introduction—that definitely improved the paper. The authors also thank Aris Daniilidis for helpful discussions and Rafael Ortega who called our attention to the change of variables used in the Levinson paper [29]. Part of this work was done during a visit of the first and the third author to the Centre de Recerca Matemàtica, Barcelona (CRM). They express their gratitude to the CRM for providing very nice working conditions.

Finally we thank Martin Golubitsky and Andre Vanderbauwhede, who invited us to present the paper at their minisimposia “Recent developments in bifurcation theory” of Equadiff 2007 (see [9]), giving a significant impact to its recognition.

REFERENCES

- [1] A. ANDRONOV AND A. WITT, *On mathematical theory of entrainment*, Z. Prikl. Phys., 6 (1930), pp. 3–17 (in Russian).
- [2] A. ANDRONOV AND A. WITT, *Zur Theorie des Mitnehmens von van der Pol*, Arch. Elektrotechnik, 24 (1930), pp. 99–110 (in German).
- [3] A. A. ANDRONOV, A. A. WITT, AND S. E. KHAIKIN, *Theory of Oscillators*, (translated from the Russian by F. Immirzi; translation edited and abridged by W. Fishwick) Pergamon Press, Oxford, 1966.
- [4] J. AWREJCEWICZ AND C. H. LAMARQUE, *Bifurcation and Chaos in Nonsmooth Mechanical Systems*, World Scientific, River Edge, NJ, 2003.
- [5] M. S. BAPTISTA, T. P. SILVA, J. C. SARTORELLI, AND I. L. CALDAS, *Phase synchronization in the perturbed Chua circuit*, Phys. Rev. E (3), 74 (2006), no. 056707.
- [6] N. N. BOGOLIUBOV, *On Some Statistical Methods in Mathematical Physics*, Akademiya Nauk Ukrainskoi SSR, Kiev, 1945 (in Russian).
- [7] A. P. BOVSUNOVSKII, *Comparative analysis of nonlinear resonances of a mechanical system with unsymmetrical piecewise characteristic of restoring force*, Strength Materials, 39 (2007), pp. 159–169.
- [8] A. BUICĂ AND A. DANILIDIS, *Stability of periodic solutions for Lipschitz systems obtained via the averaging method*, Proc. Amer. Math. Soc., 135 (2007), pp. 3317–3327.
- [9] A. BUICĂ, J. LLIBRE, AND O. MAKARENKOV, *Lipschitz constant of integral with respect to a parameter versus derivative of integral with respect to a parameter in the theory of nonsmooth bifurcations*, Book of abstracts, Equadiff 2007, Vienna University of Technology, 2007, pp. 88–89.

- [10] L. O. CHUA, *Global unfolding of Chua's circuit*, IEICE Trans. Fundamentals, E76-A, 5 (1993), pp. 704–734.
- [11] D. COFAGNA AND G. GRASSI, *Chaotic beats in a modified Chua's circuit: Dynamic behaviour and circuit design*, Internat. J. Bifur. Chaos Appl. Sci. Engrg., 17 (2007), pp. 209–226.
- [12] B. P. DEMIDOVICH, *Lectures on the Mathematical Theory of Stability*, Izdat. Nauka, Moscow, 1967 (in Russian).
- [13] S. W. S. W. DOEBLING, C. R. FARRAR, M. B. PRIME, AND D. W. SHEVITZ, *Damage Identification and Health Monitoring of Structural and Mechanical Systems from Changes in Their Vibration Characteristics: A Literature Review*, LA-13070-MS, UC-900, Los Alamos National Laboratory, 1996.
- [14] C. FABRY, *Large-amplitude oscillations of a nonlinear asymmetric oscillator with damping*, Nonlinear Anal. Ser. A: Theory Methods, 44 (2001), pp. 613–626.
- [15] P. FATOU, *Sur le mouvement d'un système soumis à des forces à courte période*, Bull. Soc. Math. France, 56 (1928), pp. 98–139 (in French).
- [16] J. GLOVER, A. C. LAZER, AND P. J. MCKENNA, *Existence and stability of large scale nonlinear oscillations in suspension bridges*, Z. Angew. Math. Phys., 40 (1989), pp. 172–200.
- [17] J. K. HALE, *Ordinary Differential Equations*, Robert E. Krieger Publishing Co., Inc., Huntington, NY, 1980.
- [18] S. J. HOGAN, *Relaxation oscillations in a system with a piecewise smooth drag coefficient*, J. Sound Vibration, 263 (2003), pp. 467–471.
- [19] A. A. JAKKULA, *A history of suspension bridges in bibliographical form*, Bull. Agricultural and Mechanical College of Texas 4th ser., vol. 12 Federal Works Agency, Washington, D.C., 1941.
- [20] G. A. JOHNSTON AND E. R. HUNT, *Derivative control of the steady state in Chua's circuit driven in the chaotic region*, IEEE Trans. Circuits Systems I Fund. Theory Appl., 40 (2000), pp. 833–835.
- [21] A. N. KOLMOGOROV AND S. V. FOMIN, *Introductory Real Analysis* (translated and edited by R. A. SILVERMAN, Selected Russian Publications in the Mathematical Sciences, Prentice-Hall, Englewood Cliffs, NJ, 1970).
- [22] M. A. KRASNOSEL'SKII, *Positive Solutions of Operator Equations* (translated from the Russian by R. E. Flaherty; edited by Leo F. Boron), P. Noordhoff Ltd., Groningen, The Netherlands, 1964.
- [23] M. A. KRASNOSEL'SKII, *The Operator of Translation Along the Trajectories of Differential Equations*, Translations of Mathematical Monographs 19 (translated from the Russian by Scripta Technica), American Mathematical Society, Providence, RI, 1968.
- [24] B. I. KRYUKOV, *Dynamics of Resonance-Type Vibration Machines*, Naukova Dnka, Kiev, 1967 (in Russian).
- [25] N. M. KRYLOV AND N. N. BOGOLIUBOV, *Introduction to Non-linear Mechanics*, Akademiya Nauk Ukrainskoi SSR, 1937 (in Russian); English translation: Annals of Mathematics Studies 11, Princeton University Press, Princeton, NJ, 1943.
- [26] A. C. LAZER AND P. J. MCKENNA, *Existence, uniqueness, and stability of oscillations in differential equations with asymmetric nonlinearities*, Trans. Amer. Math. Soc., 315 (1989), pp. 721–739.
- [27] A. C. LAZER AND P. J. MCKENNA, *Large-amplitude periodic oscillations in suspension bridges: Some new connections with nonlinear analysis*, SIAM Rev., 32 (1990), pp. 537–578.
- [28] M. LEVI, *Qualitative analysis of the periodically forced relaxation oscillations*, Mem. Amer. Math. Soc., 32 (1981), no. 244.
- [29] N. LEVINSON, *A second order differential equation with singular solutions*, Ann. Math. (2), 50 (1949), pp. 127–153.
- [30] I. G. MALKIN, *Some problems of the theory of nonlinear oscillations*, Gosudarstv. Izdat. Tehn.-Teor. Lit., Moscow, 1956 (in Russian); English translation: AEC tr 3766 (book 1), US Atomic Energy Commission, 1959.
- [31] L. I. MANDELSTAM AND N. D. PAPALEKSI, *On justification of one approximation method for solving differential equations*, Zh. Eksper. Teoret. Fiz., IV (1934), no. 117.
- [32] J. MAWHIN, *Le Problème des Solutions Périodiques en Mécanique non Linéaire*, Thèse de doctorat en sciences, Université de Liège, 1969; published in *Degré topologique et solutions périodiques des systèmes différentiels non linéaires*, Bull. Soc. Roy. Sci. Liège, 38 (1969), pp. 308–398.
- [33] J. MAWHIN, *Periodic solutions in the golden sixties: The birth of a continuation theorem*, in Ten Mathematical Essays on Approximation in Analysis and Topology, Elsevier B. V., Amsterdam, 2005, pp. 199–214.

- [34] YU. A. MITROPOL'SKII, *On periodic solutions of systems of nonlinear differential equations with non-differentiable right-hand sides*, Ukrain. Mat. Ž., 11 (1959), pp. 366–379 (in Russian).
- [35] K. MURALI AND M. LAKSHMANAN, *Chaotic dynamics of the driven Chua's circuit*, IEEE Trans. Circuits Systems I Fund. Theory Appl., 50 (2003), pp. 1503–1508.
- [36] A. H. NAYFEH AND D. T. MOOK, *Nonlinear Oscillations*, Wiley-Interscience, New York, 1979.
- [37] L. S. PONTRJAGIN, *Ordinary differential equations* (translated from the Russian by L. Kacinskas and W. B. Counts), Adiwes International Series in Mathematics, Addison-Wesley Publishing, Reading, MA, 1962.
- [38] R. PUERS, J. BIENSTMAN, AND J. VANDEWALLE, *The autonomous impact resonator: A new operation principle for a silicon resonant strain gauge*, in Int. Conf. on Solid-State Sensors and Actuators, Digest of Technical Papers, Chicago, 1997, pp. 1105–1108.
- [39] A. M. SAMOILENKO, *On periodic solutions of differential equations with nondifferentiable right-hand sides*, Ukrain. Mat. Ž., 15 (1963), pp. 328–332 (in Russian).
- [40] E. SÁNCHEZ, M. A. MATÍAS, AND V. PÉREZ-MUÑUZURI, *Chaotic synchronization in small assemblies of driven Chua's circuits*, IEEE Trans. Circuits Systems I Fund. Theory Appl., 47 (2000), pp. 644–654.

SPATIAL DYNAMICS OF A NONLOCAL PERIODIC REACTION-DIFFUSION MODEL WITH STAGE STRUCTURE*

YU JIN[†] AND XIAO-QIANG ZHAO[†]

Abstract. In this paper, we investigate a nonlocal periodic reaction-diffusion population model with stage structure. In the case of unbounded spatial domain, we establish the existence of the asymptotic speed of spread and show that it coincides with the minimal wave speed for monotone periodic traveling waves. In the case of bounded spatial domain, we obtain a threshold result on the global attractivity of either zero or a positive periodic solution.

Key words. nonlocal periodic model, spreading speed, traveling waves, positive periodic solution, global attractivity

AMS subject classifications. 34K05, 35K57, 37N25, 92B05

DOI. 10.1137/070709761

1. Introduction. Age structure has been an interesting topic in population dynamics (see, e.g., [1, 2, 3, 5, 6, 7, 9, 14, 17, 18, 20] and the references therein), since we can investigate the separate quantities of immature and mature populations in an age-structured population model. To derive a model for a single species of population with age structure and diffusion, we usually assume that individuals move around not only after maturity, but also while immature. For a standard argument, Metz and Diekmann [14] give

$$\frac{\partial u}{\partial t} + \frac{\partial u}{\partial a} = D(a) \frac{\partial^2 u}{\partial x^2} - \mu(a)u,$$

where $u(t, a, x)$ is the density of the population of the species at time $t \geq 0$, age $a \geq 0$, and location x in a spatial domain Ω ; $D(a) \geq 0$ and $\mu(a) \geq 0$ are the diffusion rate and the death rate of the population at age a , respectively.

To study the behaviors of immature individuals and mature individuals, we can also divide the population of a species into two groups: immature population and mature population. For simplicity, we assume that the maturation time (or the length of the juvenile period) is the same for all juvenile individuals, denoted by $\tau \geq 0$. For distributed maturation delay, see, e.g., [2, 3] and the references therein. Assume that the diffusion rate and death rate are age-dependent for immature individuals, but age-independent for mature individuals. As a result, we have the following system for a single species of population with age structure and diffusion (see also [5, 17, 18, 20]):

(1.1)

$$\left\{ \begin{array}{l} \partial_t u(t, a, x) + \partial_a u(t, a, x) = d_j(a) \Delta u - \mu_j(a)u(t, a, x), \quad t > 0, 0 < a < \tau, x \in \Omega, \\ u(t, 0, x) = f(u_m(t, x)), \quad t \geq -\tau, x \in \Omega, \\ \partial_t u_m(t, x) = d_m \Delta u_m - g(u_m(t, x)) + u(t, \tau, x), \quad t > 0, x \in \Omega, \end{array} \right.$$

*Received by the editors November 30, 2007; accepted for publication (in revised form) November 17, 2008; published electronically February 25, 2009. This research was supported in part by the NSERC of Canada and MITACS of Canada.

<http://www.siam.org/journals/sima/40-6/70976.html>

[†]Department of Mathematics and Statistics, Memorial University of Newfoundland, St. John's, NL, A1C 5S7, Canada (yuj@math.mun.ca, xzhao@math.mun.ca).

where $u(t, a, x)$ is the density of the population at time $t \geq -\tau$, age $a \geq 0$, location $x \in \Omega$, $u_m(t, x)$ is the density of the mature population, $f(u_m)$ and $g(u_m)$ are the birth rate and the mortality rate of mature individuals, respectively, $d_j(a) \geq 0$ is the diffusion rate of the immature individuals at age $a \in (0, \tau)$, $d_m \geq 0$ is the diffusion rate of the mature individuals, $\mu_j(a) > 0$ denotes the per capita mortality rate of juveniles at age a , $u(t, \tau, x)$ is the adults recruitment term for those of maturation age τ , Δ is the Laplacian operator.

In fact, the dynamics of many populations is influenced greatly by the time varying environments (e.g., due to seasonal variation). For example, in a one year period, the birth rate may be high in spring and summer and low in winter, while in winter more individuals might be at risk of death because of low temperature, lack of food, or some other reasons. Moreover, populations usually like to move in warm weather during the spring and summer time. Therefore, it is more realistic to consider a nonautonomous version of (1.1) for population dynamics. In particular, a periodic model, in which the birth rate, mortality rates, and diffusion rates are assumed to be periodic in time, is probably the simplest but nonetheless interesting and realistic case. In this paper, we consider the following model:

(1.2)

$$\begin{cases} \partial_t u(t, a, x) + \partial_a u(t, a, x) = d_j(t, a)\Delta u - \mu_j(t, a)u(t, a, x), & t > 0, a \in (0, \tau), x \in \Omega, \\ u(t, 0, x) = f(t, u_m(t, x)), & t \geq -\tau, x \in \Omega, \\ \partial_t u_m(t, x) = d_m(t)\Delta u_m - g(t, u_m(t, x)) + u(t, \tau, x), & t > 0, x \in \Omega, \end{cases}$$

where $d_j(t, a) \geq 0$ and $\mu_j(t, a) \geq 0$ denote the diffusion rate and the per capita mortality rate of juveniles at age a at time t , respectively; $d_m(t) \geq 0$ denotes the diffusion rate of mature individuals at time t ; $f(t, u_m)$ and $g(t, u_m)$ are the birth and mortality rates of mature individuals at time t , respectively.

Similarly as in [18] (see also [16]), we integrate along characteristics to reduce the system (1.2) to one equation with nonlocal terms. Let $v(r, a, x) = u(a + r, a, x)$, where r is regarded as a parameter. It follows that

$$\begin{cases} \partial_a v(r, a, x) = [\partial_t u(t, a, x) + \partial_a u(t, a, x)]_{t=r+a} \\ \qquad \qquad \qquad = d_j(a + r, a)\Delta v(r, a, x) - \mu_j(a + r, a)v(r, a, x), \\ v(r, 0, x) = f(r, u_m(r, x)). \end{cases}$$

Integrating the last equation, we obtain

$$v(r, a, x) = \int_{\Omega} \Gamma(\zeta(r, a), x - y)F(r, a)f(r, u_m(r, y))dy,$$

where Γ is the fundamental solution associated with the partial differential operator $\partial_t - \Delta$ and

$$(1.3) \quad \zeta(r, a) = \int_0^a d_j(s + r, s)ds, \quad F(r, a) = \exp\left(-\int_0^a \mu_j(s + r, s)ds\right).$$

Since $u(t, a, x) = v(t - a, a, x)$, it follows that

$$(1.4) \quad u(t, a, x) = \int_{\Omega} \Gamma(\zeta(t - a, a), x - y)F(t - a, a)f(t - a, u_m(t - a, y))dy.$$

Set

$$a(t) := \zeta(t - \tau, \tau), \quad b(t) := F(t - \tau, \tau), \quad f_{-\tau}(t, u) := f(t - \tau, u).$$

Substituting (1.4) into the equation for u_m in (1.2), we finally reduce the age-structured population model (1.2) to the following time-delayed reaction-diffusion equation for mature individuals:

(1.5)

$$\begin{cases} \partial_t u_m(t, x) \\ = d_m(t) \Delta u_m - g(t, u_m(t, x)) + b(t) \int_{\Omega} \Gamma(a(t), x - y) f_{-\tau}(t, u_m(t - \tau, y)) dy, \\ u_m(s, x) = \phi(s, x), \quad s \in [-\tau, 0], \quad x \in \Omega, \end{cases}$$

where $\phi(t, x)$ is an initial function to be specified later. For simplicity, dropping all m 's and writing $u_m(t, x)$ as $u(t, x)$, we investigate the following system:

(1.6)

$$\begin{cases} \partial_t u(t, x) = d(t) \Delta u - g(t, u(t, x)) + b(t) \int_{\Omega} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y)) dy, \\ u(s, x) = \phi(s, x), \quad s \in [-\tau, 0], \quad x \in \Omega. \end{cases}$$

Basically we assume that $d_j(t, a)$ and $\mu_j(t, a)$ are periodic in $t \geq 0$ with the period $\omega > 0$ for $a \in (0, \tau)$, and that $d(t)$, $g(t, u)$, and $f(t, u)$ are periodic in t with the period $\omega > 0$ for $u \in \mathbb{R}_+$. This implies that $a(t) = a(t + \omega)$, $d(t) = d(t + \omega)$, $b(t) = b(t + \omega)$, $g(t, u) = g(t + \omega, u)$, and $f(t, u) = f(t + \omega, u)$ for all $t \geq 0$, $u \in \mathbb{R}_+$. Moreover, we assume $d(t) \geq d > 0$ for all $t \geq 0$, and

- (H1) $f \in C^1([-\tau, +\infty) \times \mathbb{R}_+, \mathbb{R}_+)$, $g \in C^1(\mathbb{R}_+^2, \mathbb{R}_+)$, $f(t, 0) = 0$ for $t \geq -\tau$, $f_u(t, u) > 0$ for all $t \geq -\tau$ and $u \geq 0$, $g(t, 0) = 0$ for $t \geq 0$, and there exists $l_1 > 0$ such that $|g(t, u) - g(t, v)| \leq l_1 |u - v|$ for all $t \geq 0$ and $u, v \in \mathbb{R}_+$;
- (H2) $G(t, u, v) := -g(t, u) + b(t) f_{-\tau}(t, v)$ is strictly subhomogeneous in (u, v) in the sense that for any $\alpha \in (0, 1)$, $G(t, \alpha u, \alpha v) > \alpha G(t, u, v)$ for all $u, v \geq 0$;
- (H3) there exists positive number $L > 0$ such that $G(t, \bar{L}, \bar{L}) \leq 0$ for all $t \geq 0$, $\bar{L} \geq L$.

The purpose of this paper is to study the asymptotic speed of spread and periodic traveling waves of (1.6) in the infinite spatial domain, and the global attractivity of zero or a positive periodic solution of (1.6) in a bounded spatial domain. The asymptotic speed of spread (in short, spreading speed) was first introduced by Aronson and Weinberger [4] for reaction-diffusion equations and has been an important ecological metric in a wide range of ecological applications; see, e.g., [10, 11, 18] and the references therein. Intuitively, the spreading speed c^* in a spatial epidemic model can be interpreted as: if one runs at a speed $c > c^*$, then one will leave the epidemic behind; whereas if one runs at a speed $c < c^*$, then one will eventually be surrounded by the epidemic. Traveling wave solutions have also been investigated extensively for a variety of evolution systems; see, e.g., [5, 10, 11, 17, 18] and the references therein. For the autonomous case of (1.6), the dynamics, including spreading speed and traveling waves, have been studied extensively. So, Wu and Zou [17] investigated traveling wave fronts in the case where $\Omega = \mathbb{R}$, $g(u) = \beta u$. Gourley and Kuang [5] established the linear stabilities of two spatially homogeneous equilibrium solutions, studied traveling

wave fronts in the case where $\Omega = \mathbb{R}$, $f(u) = \alpha u$, and $g(u) = \beta u^2$, and obtained a global convergence theorem in the case of bounded intervals. Thieme and Zhao [18] studied the traveling wave solutions, minimal wave speed, and asymptotic speed of spread in the case of $\Omega = \mathbb{R}^n$. Xu and Zhao [20] established a threshold dynamic and global attractivity of the positive steady state when Ω is a bounded domain in \mathbb{R}^n .

This paper is organized as follows. In section 2, we first establish the well-posedness and the comparison principle for (1.6) with $\Omega = \mathbb{R}$, then prove the existence of the spreading speed c^* for solutions of (1.6) with initial data having compact supports, and show that it coincides with the minimal wave speed for monotone periodic traveling waves, by appealing to the theory of the spreading speed and traveling waves for monotone periodic semiflows developed in [10, 11]. In section 3, we use the theory of monotone and subhomogeneous dynamical systems to investigate the global dynamics of (1.6) in a bounded domain $\Omega \subseteq \mathbb{R}^n$, and obtain a threshold result for global attractivity of either zero or a positive periodic solution.

2. Spreading speed and traveling waves. In this section, we consider that the population diffuses in an unbounded spatial domain and study (1.6) with $\Omega = \mathbb{R}$:

(2.1)

$$\begin{cases} \partial_t u(t, x) = d(t)\Delta u - g(t, u(t, x)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y)) dy, \\ u(t, x) = \phi(t, x), \quad t \in [-\tau, 0], x \in \mathbb{R}. \end{cases}$$

In the following, we first apply the threshold dynamics in a scalar periodic and time-delayed equation, developed by Xu and Zhao [21], to the spatially homogeneous system associated with (2.1) to find a periodic solution of (2.1). Then we use the theory of abstract functional differential equations and reaction-diffusion systems to establish the existence of solutions to (2.1) and a comparison principle. Finally, we prove that the solution periodic semiflow of (2.1) satisfies all the assumptions on monotone periodic semiflows in [10], and hence, we obtain the existence of the spreading speed and traveling wave solutions for (2.1).

Let \mathbb{Y} be the space of all continuous functions from $[-\tau, 0]$ to \mathbb{R} with the usual supreme norm $\|\cdot\|_{\mathbb{Y}}$ (i.e., $\mathbb{Y} = C([-\tau, 0], \mathbb{R})$), and let $\mathbb{Y}_+ = C([-\tau, 0], \mathbb{R}_+)$. Then $(\mathbb{Y}, \mathbb{Y}_+)$ is an ordered Banach space. For $\varphi, \psi \in \mathbb{Y}$, we write $\varphi \leq \psi$ if $\psi - \varphi \in \mathbb{Y}_+$, $\varphi < \psi$ if $\psi - \varphi \in \mathbb{Y}_+ \setminus \{0\}$, $\varphi \ll \psi$ if $\psi - \varphi \in \text{int}(\mathbb{Y}_+)$. Moreover, we define $\mathbb{Y}_r = \{\varphi \in \mathbb{Y} : 0 \leq \varphi \leq r\}$ for any $r \in \mathbb{Y}$ with $r \gg 0$.

Let \mathbb{X} be the set of all bounded and continuous functions from \mathbb{R} into \mathbb{R} and $\mathbb{X}_+ = \{\varphi \in \mathbb{X}; \varphi(x) \geq 0 \text{ for all } x \in \mathbb{R}\}$. For $\varphi, \psi \in \mathbb{X}$, we write $\varphi \leq \psi$ ($\varphi \ll \psi$) if $\varphi(x) \leq \psi(x)$ ($\varphi(x) < \psi(x)$) for all $x \in \mathbb{R}$, $\varphi < \psi$ if $\varphi \leq \psi$ but $\varphi \neq \psi$. It is easy to see that \mathbb{X}_+ is a positive cone of \mathbb{X} . Define $\mathbb{X}_r = \{\varphi \in \mathbb{X} : 0 \leq \varphi \leq r\}$ for any $r \in \mathbb{X}$ with $r \gg 0$. We equip \mathbb{X} with the compact open topology, i.e., $u^m \rightarrow u$ in \mathbb{X} means that the sequence of $u^m(x)$ converges to $u(x)$ as $m \rightarrow \infty$ uniformly for x in any compact set on \mathbb{R} . Define

$$\|u\|_{\mathbb{X}} = \sum_{k=1}^{\infty} \frac{\max_{|x| \leq k} |u(x)|}{2^k} \quad \forall u \in \mathbb{X},$$

where $|\cdot|$ denotes the usual norm in \mathbb{R} . Then $(\mathbb{X}, \|\cdot\|_{\mathbb{X}})$ is a normed space. Let $d_{\mathbb{X}}(\cdot, \cdot)$ be the distance induced by the norm $\|\cdot\|_{\mathbb{X}}$. It follows that the topology in the metric space $(\mathbb{X}, d_{\mathbb{X}})$ is the same as the compact open topology in \mathbb{X} . Moreover, $(\mathbb{X}_r, d_{\mathbb{X}})$ is a complete metric space.

Let C be the set of continuous functions from $[-\tau, 0]$ into \mathbb{X} , $C_+ = \{\varphi \in C, \varphi(s) \in \mathbb{X}_+, s \in [-\tau, 0]\}$ and $C_r = \{\varphi \in C : 0 \leq \varphi \leq r\}$ for any $r \in \mathbb{Y}$ with $r \gg 0$. Then C_+ is a positive cone of C . For convenience, we also identify an element $\varphi \in C$ as a function from $[-\tau, 0] \times \mathbb{R}$ into \mathbb{R} defined by $\varphi(s, x) = \varphi(s)(x)$ for any $s \in [-\tau, 0]$ and $x \in \mathbb{R}$. For $\varphi, \psi \in C$, we write $\varphi \leq \psi$ ($\varphi \ll \psi$) if $\varphi(s, x) \leq \psi(s, x)$ ($\varphi(s, x) < \psi(s, x)$) for all $s \in [-\tau, 0], x \in \mathbb{R}$, $\varphi < \psi$ if $\varphi \leq \psi$ but $\varphi \neq \psi$. For any continuous function $w(\cdot) : [-\tau, b) \rightarrow \mathbb{X}, b > 0$, we define $w_t \in C$ by $w_t(s) = w(t + s)$ for all $t \in [0, b), s \in [-\tau, 0]$. It is then easy to see that $t \rightarrow w_t$ is a continuous function from $[0, b)$ to C . Moreover, we also equip C with the compact open topology and define the norm on C :

$$\|u\|_C = \sum_{k=1}^{\infty} \frac{\max_{|x| \leq k, s \in [-\tau, 0]} |u(s, x)|}{2^k} \quad \forall u \in C,$$

where $|\cdot|$ denotes the usual norm in \mathbb{R} .

For any constant $N > 0$, \widehat{N} denotes the constant function with value N in \mathbb{Y}, \mathbb{X} , or C .

Now we consider the spatially homogeneous system associated with (2.1). Letting $u(t, x) = w(t)$, we have

$$(2.2) \quad \begin{cases} \frac{dw(t)}{dt} = -g(t, w(t)) + b(t)f_{-\tau}(t, w(t - \tau)), \\ w(t) = \varphi(t), \quad t \in [-\tau, 0], \quad \varphi \in \mathbb{Y}_+. \end{cases}$$

The linearized equation associated with (2.2) at $w = 0$ is

$$(2.3) \quad \begin{cases} \frac{dw(t)}{dt} = -g_u(t, 0)w(t) + b(t)\partial_u f_{-\tau}(t, 0)w(t - \tau), \\ w(t) = \varphi(t), \quad t \in [-\tau, 0], \quad \varphi \in \mathbb{Y}_+. \end{cases}$$

Since g, b , and $f_{-\tau}$ are periodic functions in $t \geq 0$, we can easily see that for any $\varphi \in \mathbb{Y}_+$, (2.3) admits a unique solution $w(t, \varphi)$ existing for all $t \geq -\tau$ with $w(s, \varphi) = \varphi(s)$ for $s \in [-\tau, 0]$, and $w_t(\varphi) \in \mathbb{Y}_+$ for all $t \geq 0$, where $\{w_t\}_{t \geq 0}$ is the solution semiflow for (2.3) defined by $w_t(\varphi)(s) = w(t + s, \varphi)$ for all $s \in [-\tau, 0], t > 0$.

Define the Poincaré map of (2.3) $P : \mathbb{Y}_+ \rightarrow \mathbb{Y}_+$ by $P(\varphi) = w_\omega(\varphi)$ for all $\varphi \in \mathbb{Y}_+$, and let $r = r(P)$ be the spectral radius of P . The following two results come from [21].

PROPOSITION 2.1 (see [21, Proposition 2.1]). *$r = r(P)$ is positive and is an eigenvalue of P with a positive eigenfunction φ^* . Moreover, if $\tau = k\omega$ for some integer $k \geq 0$, then $r - 1$ has the same sign as $\int_0^\omega [-g_u(t, 0) + b(t)\partial_u f_{-\tau}(t, 0)]dt$.*

THEOREM 2.2 (see [21, Theorem 2.1]). *Let (H1)–(H3) hold. The following statements are valid.*

- (i) *If $r \leq 1$, then zero solution is globally asymptotically stable for (2.2) with respect to \mathbb{Y}_+ .*
- (ii) *If $r > 1$, then (2.2) has a unique positive ω -periodic solution $\beta^*(t)$, and $\beta^*(t)$ is globally asymptotically stable with respect to $\mathbb{Y}_+ \setminus \{0\}$.*

In the remainder of this section, we further assume that

(H4) $r = r(P) > 1$.

By the proof of [21, Theorem 2.1] and (H3), it is easy to see that $\beta^*(t) \in [0, L]$ for all $t \geq -\tau$ and $[\widehat{0}, \widehat{L}]$ is positively invariant for (2.1). Define $\beta_0^* \in \mathbb{Y}_{\widehat{L}}$ as $\beta_0^*(s) = \beta^*(s)$ for all $s \in [-\tau, 0]$.

Consider

$$(2.4) \quad \begin{cases} \partial_t u(t, x) = d(t)\Delta u, & t > 0, \\ u(0, x) = \phi(x), & x \in \mathbb{R}, \phi \in \mathbb{X}. \end{cases}$$

The solution of (2.4) can be expressed as

$$(2.5) \quad u(t, x, \phi) = \int_{\mathbb{R}} \Gamma(\eta(t), x - y)\phi(y)dy, \quad t \geq 0,$$

where $\eta(t) = \int_0^t d(s)ds$. According to [8, Chapter II], (2.4) admits an evolution operator $U(t, s) : \mathbb{X} \rightarrow \mathbb{X}$, $0 \leq s \leq t$, which satisfies $U(t, t) = I$, $U(t, s)U(s, \rho) = U(t, \rho)$ for all $0 \leq \rho \leq s \leq t$, and $U(t, 0)(\phi)(x) = u(t, x, \phi)$ for $t \geq 0$, $x \in \mathbb{R}$, and $\phi \in \mathbb{X}$, where $u(t, x, \phi)$ is the solution of (2.4). Moreover, for any $0 \leq s < t$, $U(t, s)$ is a compact and positive operator on \mathbb{X} , and $U(t, s)(\phi)(x) > 0$ for all $0 \leq s < t$, $x \in \mathbb{R}$, and $\phi \in \mathbb{X}$ provided that $\phi(x) \geq 0$ and $\phi \not\equiv 0$.

Define $B : [0, \infty) \times C \rightarrow \mathbb{X}$ by $B(t, \phi) := -g(t, \phi(0, \cdot)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), \cdot - y)f_{-\tau}(t, \phi(-\tau, y))dy$ for any $t \in [0, \infty)$, $\phi \in C$. Then (2.1) becomes

$$(2.6) \quad \begin{cases} \partial_t u(t, x) = d(t)\Delta u + B(t, u_t), & t > 0, \\ u(t, x) = \phi(t, x), & t \in [-\tau, 0], x \in \mathbb{R}, \end{cases}$$

which can be written as an integral equation

$$(2.7) \quad u(t, \cdot, \phi) = U(t, 0)\phi(0, \cdot) + \int_0^t U(t, s)B(s, u_s)ds, \quad t \geq 0, \quad \phi \in C,$$

whose solutions are called mild solutions to (2.6).

THEOREM 2.3. *Let (H1)–(H4) hold. For any $\phi \in C_{\hat{L}}$, system (2.1) has a unique mild solution $u(t, x, \phi)$ with $u_0(\cdot, \cdot, \phi) = \phi$ and $u_t(\cdot, \cdot, \phi) \in C_{\hat{L}}$ for all $t \geq 0$, and $u(t, x, \phi)$ is a classic solution when $t > \tau$. Moreover, if $\hat{u}(t, x)$ and $\bar{u}(t, x)$ are a pair of lower and upper solutions of (2.1), respectively, with $\hat{u}_0(\cdot, \cdot) \leq \bar{u}_0(\cdot, \cdot)$, then $\hat{u}_t(\cdot, \cdot) \leq \bar{u}_t(\cdot, \cdot)$ for all $t \geq 0$.*

Proof. We first show that B is quasi-monotone on $[0, \infty) \times C_{\hat{L}}$ in the sense that

$$(2.8) \quad \lim_{h \rightarrow 0^+} d(\psi(0, \cdot) - \phi(0, \cdot) + h[B(t, \psi) - B(t, \phi)], \mathbb{X}_+) = 0$$

for all $\phi, \psi \in C_{\hat{L}}$ with $\phi(s, x) \leq \psi(s, x)$ for all $s \in [-\tau, 0]$, $x \in \mathbb{R}$. In fact, for any $\phi, \psi \in C_{\hat{L}}$ with $\phi(s, x) \leq \psi(s, x)$ for all $(s, x) \in [-\tau, 0] \times \mathbb{R}$, we have

$$\begin{aligned} & \psi(0, \cdot) - \phi(0, \cdot) + h[B(t, \psi) - B(t, \phi)] \\ &= \psi(0, \cdot) - \phi(0, \cdot) + h[-(g(t, \psi(0, \cdot)) - g(t, \phi(0, \cdot)))] \\ &+ h \left[\int_{\mathbb{R}} \Gamma(a(t), \cdot - y)b(t)(f_{-\tau}(t, \psi(-\tau, y)) - f_{-\tau}(t, \phi(-\tau, y)))dy \right] \\ &\geq \psi(0, \cdot) - \phi(0, \cdot) - h(g(t, \psi(0, \cdot)) - g(t, \phi(0, \cdot))) \\ &\geq (1 - hl_1)(\psi(0, \cdot) - \phi(0, \cdot)). \end{aligned}$$

Thus, for $1 - hl_1 > 0$, $\psi(0, \cdot) - \phi(0, \cdot) + h[B(t, \psi) - B(t, \phi)] \in \mathbb{X}_+$, and hence, (2.8) holds. Then by [13, Corollary 5] (for $v^- = 0$, $v^+ = \hat{L}$, $S^+ = S^- = S = T \equiv U$,

$B^+ = B^- = B$), (2.1) admits a unique mild solution $u(t, \cdot, \phi)$ on $[-\tau, \infty)$ for any $\phi \in C_{\widehat{L}}$ and $u_t(\cdot, \cdot, \phi) \in C_{\widehat{L}}$ for all $t \geq 0$. Moreover, the comparison principle holds for lower and upper solutions. \square

In order to study the spreading speed and traveling waves, we introduce the assumptions in [10, 11]. Let $u \in C$. Define the reflection operator \mathcal{R} by $\mathcal{R}[u](\theta, x) := u(\theta, -x)$ for all $\theta \in [-\tau, 0], x \in \mathbb{R}$. Given $y \in \mathbb{R}$, define the translation operator T_y by $T_y[u](\theta, x) := u(\theta, x - y)$ for all $\theta \in [-\tau, 0], x \in \mathbb{R}$. Let $Q : C_{b^*} \rightarrow C_{b^*}$ be a map, where $b^* \in \mathbb{Y}$ with $b^* \gg 0$. Assume the following:

- (A1) $Q[\mathcal{R}[u]] = \mathcal{R}[Q[u]], T_y[Q[u]] = Q[T_y[u]] \forall y \in \mathbb{R}$.
- (A2) $Q : C_{b^*} \rightarrow C_{b^*}$ is continuous with respect to the compact open topology.
- (A4) $Q : C_{b^*} \rightarrow C_{b^*}$ is monotone in the sense that $Q[u] \geq Q[v]$ whenever $u \geq v$ in C_{b^*} .
- (A5) $Q : \mathbb{Y}_{b^*} \rightarrow \mathbb{Y}_{b^*}$ admits exactly two fixed points 0 and b^* , and for any positive number ε , there is an $\alpha \in \mathbb{Y}_{b^*}$ with $\|\alpha\|_{\mathbb{Y}} < \varepsilon$ such that $Q[\alpha] \gg \alpha$.
- (A6) One of the following two conditions holds:
 - (a) $Q[C_{b^*}]$ is precompact in C_{b^*} .
 - (b') The set $Q[C_{b^*}](0, \cdot)$ is precompact in \mathbb{X} , and there is a positive number $\varsigma \leq \tau$ such that $Q[u](\theta, x) = u(\theta + \varsigma, x)$ for $-\tau \leq \theta \leq -\varsigma$, the operator

$$S[u](\theta, x) := \begin{cases} u(0, x), & -\tau \leq \theta < -\varsigma, \\ Q[u](\theta, x), & -\varsigma \leq \theta \leq 0 \end{cases}$$

is continuous on C_{b^*} , and $S[D](\cdot, 0)$ is precompact in \mathbb{Y} for any T -invariant set $D \subseteq C_{b^*}$ with $D(0, \cdot)$ being precompact in \mathbb{X} . A set $W \subseteq C_{b^*}$ is said to be T -invariant if $T_y W = W$ for all $y \in \mathbb{R}$.

Recall that a family of operators $\{\Phi_t\}_{t \geq 0}$ is an ω -periodic semiflow on a metric space (X, ρ) with the metric ρ , provided that $\{\Phi_t\}$ satisfies

- (i) $\Phi_0(v) = v \forall v \in X$;
- (ii) $\Phi_t(\Phi_\omega(v)) = \Phi_{t+\omega}(v) \forall t \geq 0, v \in X$;
- (iii) $\Phi(t, v) = \Phi_t(v)$ is continuous in (t, v) on $[0, +\infty) \times X$.

Define a family of operators $\{Q_t\}_{t \geq 0}$ on $C_{\widehat{L}}$ by

$$Q_t(\phi)(s, x) = u(t + s, x, \phi) \quad \forall t \geq 0, \quad s \in [-\tau, 0], \quad x \in \mathbb{R}, \quad \phi \in C_{\widehat{L}},$$

where $u(t, x, \phi)$ is the mild solution of (2.1) with $u(s, x) = \phi(s, x)$ for $s \in [-\tau, 0], x \in \mathbb{R}$. Note that for any $(t_0, \phi_0) \in \mathbb{R}_+ \times C_{\widehat{L}}$, we have

$$\|Q_t(\phi) - Q_{t_0}(\phi_0)\|_C \leq \|Q_t(\phi) - Q_t(\phi_0)\|_C + \|Q_t(\phi_0) - Q_{t_0}(\phi_0)\|_C.$$

Note that $U(t, 0)\varphi$ is continuous in $(t, \varphi) \in [0, \infty) \times \mathbb{X}$ with respect to the compact open topology. By a similar argument as in [12, Theorem 8.5.2], it follows that $Q_t(\phi)$ is continuous at (t_0, ϕ_0) with respect to the compact open topology. Thus, $\{Q_t\}_{t \geq 0}$ is an ω -periodic semiflow on $C_{\widehat{L}}$.

LEMMA 2.4. *For each $t > 0$, Q_t is strictly subhomogeneous.*

Proof. For any $\phi \in C_{\widehat{L}}$ with $\phi \not\equiv 0$, let $u(t, x, \phi)$ be the solution of (2.1) with $u(s, x) = \phi(s, x)$ for $s \in [-\tau, 0], x \in \mathbb{R}$. Fix $k \in (0, 1)$. Since $G(t, u, v)$ is strictly subhomogeneous in (u, v) , we have

$$\begin{aligned} & \partial_t(ku(t, x)) \\ &= d(t)\Delta(ku) - kg(t, u(t, x)) + kb(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y)) dy \\ &\leq d(t)\Delta(ku) - g(t, ku(t, x)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, ku(t - \tau, y)) dy. \end{aligned}$$

Thus, $ku(t, x, \phi)$ is a lower solution of (2.1) with $ku(s, x, \phi) = k\phi(s, x)$ for $s \in [-\tau, 0], x \in \mathbb{R}$. Then, $ku(t, x, \phi) \leq u(t, x, k\phi)$ for $t \geq 0$, where $u(t, x, k\phi)$ is the solution of (2.1) with $u(s, x, k\phi) = k\phi(s, x)$ for $(s, x) \in [-\tau, 0] \times \mathbb{R}$.

Let $w(t, x) = u(t, x, k\phi) - ku(t, x, \phi)$. Then $w(s, x) = 0$ for $(s, x) \in [-\tau, 0] \times \mathbb{R}$ and $w(s, x) \geq 0$ for $(s, x) \in [-\tau, \infty) \times \mathbb{R}$. We further show that $w(t, x) > 0$ for all $t > 0, x \in \mathbb{R}$. For simplicity, we write $\tilde{F}(t, u(t, x), v(t, x)) = -g(t, u(t, x)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, v(t, y)) dy$. It follows that

(2.9)

$$\begin{aligned} & \frac{\partial w(t, x)}{\partial t} \\ &= \frac{\partial u(t, x, k\phi)}{\partial t} - k \frac{\partial u(t, x, \phi)}{\partial t} \\ &= d(t)\Delta u(t, x, k\phi) + \tilde{F}(t, u(t, x, k\phi), u(t - \tau, x, k\phi)) \\ &\quad - k[d(t)\Delta u(t, x, \phi) + \tilde{F}(t, u(t, x, \phi), u(t - \tau, x, \phi))] \\ &= d(t)\Delta w(t, x) + [\tilde{F}(t, u(t, x, k\phi), u(t - \tau, x, k\phi)) - \tilde{F}(t, ku(t, x, \phi), ku(t - \tau, x, \phi))] \\ &\quad + [\tilde{F}(t, ku(t, x, \phi), ku(t - \tau, x, \phi)) - k\tilde{F}(t, u(t, x, \phi), u(t - \tau, x, \phi))] \\ &= d(t)\Delta w(t, x) - g(t, u(t, x, k\phi)) + g(t, ku(t, x, \phi)) + h(t, x) \\ &\quad + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) [f_{-\tau}(t, u(t - \tau, y, k\phi)) - f_{-\tau}(t, ku(t - \tau, y, \phi))] dy \\ &\geq d(t)\Delta w(t, x) - g(t, u(t, x, k\phi)) + g(t, ku(t, x, \phi)) + h(t, x) \\ &\geq d(t)\Delta w(t, x) - l_1 w(t, x) + h(t, x), \end{aligned}$$

where $h(t, x) = \tilde{F}(t, ku(t, x, \phi), ku(t - \tau, x, \phi)) - k\tilde{F}(t, u(t, x, \phi), u(t - \tau, x, \phi))$. Let $\tilde{U}(t, s) : \mathbb{X} \rightarrow \mathbb{X}, 0 \leq s \leq t$, be the evolution operator of

$$\begin{cases} \partial_t u(t, x) = d(t)\Delta u - l_1 u(t, x), & t > 0, \\ u(0, x) = \psi(x), & x \in \mathbb{R}, \psi \in \mathbb{X}. \end{cases}$$

Then $\tilde{U}(t, s)(\psi)(x) = e^{-l_1(t-s)}U(t, s)(\psi)(x)$ for all $t \geq s \geq 0, x \in \Omega, \psi \in C$, where $U(t, s)$ is the evolution operator of (2.4). Thus, the equation

$$(2.10) \quad \begin{cases} \partial_t u(t, x) = d(t)\Delta u - l_1 u(t, x) + h(t, x), \\ u(0, x) = \psi(x), & x \in \mathbb{R}, \psi \in \mathbb{X} \end{cases}$$

can be written as

$$(2.11) \quad u(t, x, \psi) = \tilde{U}(t, 0)(\psi)(x) + \int_0^t \tilde{U}(t, s)h(s, x)ds, \quad t \geq 0, x \in \mathbb{R}, \psi \in C.$$

By (H2), we have $h(t, x) > 0$ for all $t > 0, x \in \mathbb{R}$. It then follows from (2.11) and the property of $\tilde{U}(t, s)$ that for any $\psi \geq 0$ with $\psi \not\equiv 0$, the solution of (2.10) satisfies $u(t, x, \psi) > 0$ for all $t > 0, x \in \mathbb{R}$. Then by (2.9) and the comparison principle, we have $w(t, x) > 0$ for all $t > 0, x \in \mathbb{R}$. Therefore, $u(t, x, k\phi) > ku(t, x, \phi)$ for all $t > 0, x \in \mathbb{R}$, and hence, $Q_t(k\phi) > kQ_t(\phi)$ for all $t > 0$, which indicates that for each $t > 0, Q_t$ is strictly subhomogeneous. \square

LEMMA 2.5. For any $\varphi \in C_{\widehat{L}}$ with $\varphi \not\equiv 0$, $u(t, x, \varphi) > 0$ for all $t \geq \tau$, $x \in \mathbb{R}$.

Proof. Let $\varphi \in C_{\widehat{L}}$ with $\varphi \not\equiv 0$. By Theorem 2.3, $u(t, x, \varphi) \geq 0$ for all $t \geq 0$ and $x \in \mathbb{R}$. It follows from (H1) that for any $t > 0$, $u(t, x, \varphi)$ satisfies

$$\begin{aligned} \partial_t u(t, x) &= d(t)\Delta u - g(t, u(t, x, \varphi)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y, \varphi)) \\ &\geq d(t)\Delta u - g(t, u(t, x, \varphi)) \\ &\geq d(t)\Delta u - l_1 u(t, x, \varphi). \end{aligned}$$

By [19, Theorem 5.5.4], $u(t, x, \varphi) > 0$ for all $t > 0$, $x \in \mathbb{R}$, provided that $\varphi(0, \cdot) > 0$.

Now we show that for any $\varphi \in C_{\widehat{L}}$ with $\varphi \not\equiv 0$ and $\varphi(0, \cdot) = 0$, there exists $t_0 = t_0(\varphi) \in [0, \tau]$ such that $u(t_0, \cdot, \varphi) > 0$. Assume, by contradiction, that for some $\varphi \in C_{\widehat{L}}$ with $\varphi \not\equiv 0$ and $\varphi(0, \cdot) = 0$ we have $u(t, \cdot, \varphi) \equiv 0$ for all $t \in [0, \tau]$. It follows from (2.7) that

$$0 = \int_0^t U(t, s) b(s) \int_{\mathbb{R}} \Gamma(a(s), x - y) f_{-\tau}(s, u_s(-\tau, y)) dy ds, \quad t \in [0, \tau],$$

which implies that $\int_{\mathbb{R}} \Gamma(a(s), x - y) f_{-\tau}(s, u_s(-\tau, y)) dy = 0$ for any $s \in [0, \tau]$, and hence, $f_{-\tau}(s, u_s(-\tau, y)) = 0$ for any $s \in [0, \tau]$, $y \in \mathbb{R}$. Then by (H1), $u_s(-\tau, y) = 0$ for any $s \in [0, \tau]$, $y \in \mathbb{R}$. That is, $\varphi \equiv 0$, a contradiction. Thus, we have $u(t_0, \cdot, \varphi) > 0$ for some $t_0 = t_0(\varphi) \in [0, \tau]$. Then for any $t > t_0$, $\widetilde{U}(t, t_0)[u(t_0, \cdot, \varphi)](x) = e^{-l_1(t-t_0)} U(t, t_0)[u(t_0, \cdot, \varphi)](x) > 0$, and hence, by the comparison principle, we have $u(t, x, \varphi) > 0$ for all $t > t_0$, $x \in \mathbb{R}$.

Therefore, for any $\varphi \in C_{\widehat{L}}$ with $\varphi \not\equiv 0$, $u(t, x, \varphi) > 0$ for all $t > \tau$, $x \in \mathbb{R}$. □

LEMMA 2.6. For any $t > 0$, Q_t satisfies (A1), (A2), (A4), and (A6) with $b^* = \widehat{L}$, and Q_ω satisfies (A5) with $b^* = \beta_0^*$, where $\beta_0^* \in \mathbb{Y}_{\widehat{L}}$ with $\beta_0^*(s) = \beta^*(s)$ for all $s \in [-\tau, 0]$.

Proof. It is easy to see that Q_t satisfies (A1), (A2), and (A4) with $b^* = \widehat{L}$ for any $t > 0$.

Let $\widehat{Q}_t = Q_t|_{\mathbb{Y}_{\widehat{L}}}$. Then $\widehat{Q}_t : \mathbb{Y}_{\widehat{L}} \rightarrow \mathbb{Y}_{\widehat{L}}$ is the ω -periodic semiflow generated by (2.2). Moreover, it is not difficult to see that \widehat{Q}_t is strictly monotone for any $t \geq \tau$ and strongly monotone for any $t \geq 2\tau$ on $\mathbb{Y}_{\widehat{L}}$. Note that (2.2) has a positive ω -periodic solution $\beta^*(t)$ which is globally asymptotically stable in $\mathbb{Y}_{\widehat{L}} \setminus \{0\}$. We see that \widehat{Q}_ω has only two fixed points 0 and β_0^* in $\mathbb{Y}_{\widehat{L}}$, where $\beta_0^*(s) = \beta^*(s)$ for all $s \in [-\tau, 0]$. Thus, by the Dancer–Hess connecting orbit lemma (see, e.g., [23]), the map \widehat{Q}_ω admits a strictly monotone full orbit $\{\varphi_n\}_{-\infty}^\infty \subseteq \mathbb{Y}_{\beta_0^*}$ connecting 0 to β_0^* and $\varphi_n < \varphi_{n+1}$ for any $n = 0, \pm 1, \pm 2, \dots$. For any $\bar{n} \in \mathbb{N}$ such that $\bar{n}\omega \geq 2\tau$, since $\widehat{Q}_{\bar{n}\omega}$ is strongly monotone, we have $\widehat{Q}_{\bar{n}\omega}(\varphi_n) = \widehat{Q}_\omega^{\bar{n}}(\varphi_n) \ll \widehat{Q}_\omega^{\bar{n}}(\varphi_{n+1}) = \widehat{Q}_{\bar{n}\omega}(\varphi_{n+1})$ for any $n = 0, \pm 1, \pm 2, \dots$. That is, $\varphi_{n+\bar{n}\omega} \ll \varphi_{n+1+\bar{n}\omega}$ for any $n = 0, \pm 1, \pm 2, \dots$. Therefore, $\varphi_n \ll \varphi_{n+1}$ for any $n = 0, \pm 1, \pm 2, \dots$, and hence, Q_ω satisfies (A5) with $b^* = \beta_0^*$.

Now we show that Q_t satisfies (A6)(a) with $b^* = \widehat{L}$ for $t > \tau$. Fix $t_0 > \tau$ and set $a = t_0 - \tau$, $b = t_0$. Let $u(t, \varphi)$ be the solution of (2.1) with $u_0(\varphi) = \varphi \in C_{\widehat{L}}$ and define the Kuratowski measure of noncompactness of a subset A of \mathbb{X} as

$$\alpha(A) = \inf\{r > 0 : A \text{ has a finite cover of diameter } \leq r\}.$$

First we prove that $\overline{\{u(t, \varphi) : a \leq t \leq b, \varphi \in C_{\widehat{L}}\}}$ is compact in \mathbb{X} . By (2.7), for any $\epsilon \in (0, a), t \in [a, b]$, and $\varphi \in C_{\widehat{L}}$, we have

$$\begin{aligned} &u(t, \varphi) \\ &= U(t, 0)\varphi(0, \cdot) + \int_0^{t-\epsilon} U(t, s)B(s, u_s)ds + \int_{t-\epsilon}^t U(t, s)B(s, u_s)ds \\ &= U(t, t-\epsilon) \left[U(t-\epsilon, 0)\varphi(0, \cdot) + \int_0^{t-\epsilon} U(t-\epsilon, s)B(s, u_s)ds \right] + \int_{t-\epsilon}^t U(t, s)B(s, u_s)ds \\ &= U(t, t-\epsilon)u(t-\epsilon, \varphi) + \int_{t-\epsilon}^t U(t, s)B(s, u_s)ds. \end{aligned}$$

Since $\{u(t-\epsilon, \varphi), t \in [a, b], \varphi \in C_{\widehat{L}}\}$ is bounded in \mathbb{X}_+ and $U(t, t-\epsilon)$ is compact, we have

$$\alpha(\{U(t, t-\epsilon)u(t-\epsilon, \varphi), t \in [a, b], \varphi \in C_{\widehat{L}}\}) = 0.$$

It is easy to see $\{U(t, s)B(s, u_s) : t \in [a, b], s \in [0, t], \varphi \in C_{\widehat{L}}\}$ is bounded in \mathbb{X}_+ . Let $N > 0$ such that $\|U(t, s)B(s, u_s)\|_{\mathbb{X}} \leq N$ for all $t \in [a, b], s \in [0, t], \varphi \in C_{\widehat{L}}$. By the fact of $\alpha(A) \leq \delta(A)$, where $\delta(A)$ is the diameter of $A \subseteq \mathbb{X}$, we have

$$\alpha\left(\left\{\int_{t-\epsilon}^t U(t, s)B(s, u_s)ds : t \in [a, b], s \in [t-\epsilon, t], \varphi \in C_{\widehat{L}}\right\}\right) \leq 2\epsilon N.$$

Thus,

$$\begin{aligned} &\alpha(\{u(t, \varphi) : t \in [a, b], \varphi \in C_{\widehat{L}}\}) \\ &\leq \alpha(\{U(t, t-\epsilon)u(t-\epsilon, \varphi), t \in [a, b], \varphi \in C_{\widehat{L}}\}) \\ &\quad + \alpha\left(\left\{\int_{t-\epsilon}^t U(t, s)B(s, u_s)ds : t \in [a, b], s \in [t-\epsilon, t], \varphi \in C_{\widehat{L}}\right\}\right) \\ &\leq 2\epsilon N. \end{aligned}$$

Letting $\epsilon \rightarrow 0$, we have $\alpha(\{u(t, \varphi) : t \in [a, b], \varphi \in C_{\widehat{L}}\}) = 0$, and hence, $\{u(t, \varphi) : t \in [a, b], \varphi \in C_{\widehat{L}}\}$ is precompact in \mathbb{X} .

Given a compact interval $I \subseteq \mathbb{R}$, let $K = \min\{K_1 > 0 : I \subseteq [-K_1, K_1]\}$. Since $\{u(t, \varphi) : t \in [a, b], \varphi \in C_{\widehat{L}}\}$ is precompact in \mathbb{X} , $\{u(t, \varphi)|_I : t \in [a, b], \varphi \in C_{\widehat{L}}\}$ is equicontinuous in \mathbb{X} , and hence, for any $\epsilon > 0$, there exists $\delta > 0$ such that

$$(2.12) \quad |u(t, x_1, \varphi) - u(t, x_2, \varphi)| < \epsilon$$

for all $t \in [a, b]$ and $\varphi \in C_{\widehat{L}}$, provided that $x_1, x_2 \in I$ and $|x_1 - x_2| < \delta$.

Let $[a_1, b_1]$ be any bounded interval on \mathbb{R} with $a_1 > 0$ and let $U_0(t)$ be the semigroup generated by $u_t = \Delta u$. Then $U_0(t)\varphi(x) = \int_{-\infty}^{+\infty} \Gamma(t, x-y)\varphi(y)dy$ for all $t > 0, x \in \mathbb{R}, \varphi \in \mathbb{X}$.

By the properties of Γ , we can find an $N_0 > 0$ such that $\int_{|y| \geq N_0} \Gamma(b_1, y)dy \leq \epsilon$. Since $\frac{\partial \Gamma(t, y)}{\partial t} > 0$ for all $t > 0$ and $y^2 > 2t$, we have $\int_{|y| \geq N_1} \Gamma(t, y)dy \leq \epsilon$ for all $t \in [a_1, b_1]$, where $N_1 = \max\{N_0, \sqrt{2b_1}\}$. Moreover, since $\int_{-N_1}^{N_1} \Gamma(t, y)dy$ is continuous in $t \in [a_1, b_1]$, there is a $\delta_1 > 0$ such that $|\int_{-N_1}^{N_1} (\Gamma(t_1, y) - \Gamma(t_2, y))dy| < \epsilon$ provided that

$t_1, t_2 \in [a_1, b_1]$ and $|t_1 - t_2| < \delta_1$. Therefore, for any $t_1, t_2 \in [a_1, b_1]$ and $|t_1 - t_2| < \delta_1$, $\psi \in \mathbb{X}_{\widehat{L}}, x \in I$,

$$\begin{aligned} & |(U_0(t_1)\psi)(x) - (U_0(t_2)\psi)(x)| \\ &= \left| \int_{\mathbb{R}} \Gamma(t_1, x - y)\psi(y)dy - \int_{\mathbb{R}} \Gamma(t_2, x - y)\psi(y)dy \right| \\ &= \left| \int_{\mathbb{R}} (\Gamma(t_1, y) - \Gamma(t_2, y))\psi(x - y)dy \right| \\ &\leq \left| \int_{|y| \leq N_1} (\Gamma(t_1, y) - \Gamma(t_2, y))\psi(x - y)dy \right| + \left| \int_{|y| \geq N_1} (\Gamma(t_1, y) - \Gamma(t_2, y))\psi(x - y)dy \right| \\ &< 2\epsilon L. \end{aligned}$$

It follows from the continuity of $\eta(t)$ in $t \in \mathbb{R}_+$ and definitions of $U_0(t)$ and $U(t, s)$ that there exists $\delta_2 > 0$ such that

$$|(U(t_1, 0)\varphi(0, \cdot))(x) - (U(t_2, 0)\varphi(0, \cdot))(x)| < 2\epsilon L$$

for all $x \in I$, $\varphi \in C_{\widehat{L}}$, provided that $t_1, t_2 \in [a, b]$ and $|t_1 - t_2| < \delta_2$. Let $\bar{\delta} \in (0, \min\{\epsilon, \delta_2\})$. Then for $x \in I$, $\varphi \in C_{\widehat{L}}$, $t_1, t_2 \in [a, b]$, and $|t_1 - t_2| < \bar{\delta}$, we have

$$\begin{aligned} & |u(t_1, x, \varphi) - u(t_2, x, \varphi)| \\ &\leq |(U(t_1, 0)\varphi(0, \cdot))(x) - (U(t_2, 0)\varphi(0, \cdot))(x)| \\ (2.13) \quad & + \left| \int_0^{t_1} (U(t_1, s)B(s, u_s))(x)ds - \int_0^{t_2} (U(t_2, s)B(s, u_s))(x)ds \right| \\ &\leq 2L\epsilon + 2N \cdot 2^K \bar{\delta} \\ &\leq 2(L + 2^K N)\epsilon, \end{aligned}$$

where N was defined in the former paragraph of this proof. This implies that $u(t, x, \varphi)$ is equicontinuous in $t \in [a, b]$ for $x \in I$ and $\varphi \in C_{\widehat{L}}$.

Consequently, by (2.12) and (2.13), for any $\varphi \in C_{\widehat{L}}$, $\theta_1, \theta_2 \in [-\tau, 0]$, $x_1, x_2 \in I$ with $|\theta_1 - \theta_2| < \bar{\delta}$ and $|x_1 - x_2| < \delta$, we have

$$\begin{aligned} & |u_{t_0}(\varphi)(\theta_1, x_1) - u_{t_0}(\varphi)(\theta_2, x_2)| \\ &= |u(t_0 + \theta_1, x_1, \varphi) - u(t_0 + \theta_2, x_2, \varphi)| \\ &\leq |u(t_0 + \theta_1, x_1, \varphi) - u(t_0 + \theta_1, x_2, \varphi)| + |u(t_0 + \theta_1, x_2, \varphi) - u(t_0 + \theta_2, x_2, \varphi)| \\ &\leq (2L + 2^{K+1}N + 1)\epsilon, \end{aligned}$$

which indicates that $\{u_{t_0}(\varphi) : \varphi \in C_{\widehat{L}}\}$ is equicontinuous for $(\theta, x) \in [-\tau, 0] \times I$. Therefore, $\{u_{t_0}(\varphi) : \varphi \in C_{\widehat{L}}\}$ is precompact in $C_{\widehat{L}}$ and (A6)(a) follows from $Q_{t_0}(C_{\widehat{L}}) = \{u_{t_0}(\varphi) : \varphi \in C_{\widehat{L}}\}$ for $t_0 > \tau$.

Finally, we show that Q_t satisfies (A6)(b') with $b^* = \widehat{L}$ for $0 < t \leq \tau$. Fix $t_1 \in (0, \tau]$ and define

$$S[\varphi](\theta, x) := \begin{cases} \varphi(0, x), & -\tau \leq \theta < -t_1, \\ Q_{t_1}(\varphi)(\theta, x), & -t_1 \leq \theta \leq 0 \end{cases}$$

for any $\varphi \in C_{\widehat{L}}$. By the above analysis, we know that $\{u(t, \varphi) : a \leq t \leq b, \varphi \in C_{\widehat{L}}\}$ is precompact in \mathbb{X} for any $0 < a \leq b$. In particular, fixing $a = b = t_1$, we can easily see that $\{u_{t_1}(\varphi)(0, \cdot), \varphi \in C_{\widehat{L}}\} = \{u(t_1, \cdot, \varphi), \varphi \in C_{\widehat{L}}\}$ is precompact in \mathbb{X} , that is, $Q_{t_1}[C_{\widehat{L}}](0, \cdot)$ is precompact in \mathbb{X} .

Since Q_t is an ω -periodic semiflow, it is easy to see that $S[\varphi]$ is continuous on $C_{\widehat{L}}$. Let D be a T -invariant subset of $C_{\widehat{L}}$ (i.e., $T_y D = D$ for all $y \in \mathbb{R}$) with $D(0, \cdot)$ being precompact in \mathbb{X} . Now we show that for any given compact interval $I \subseteq \mathbb{R}$, $S[D]$ is equicontinuous on $[-\tau, 0] \times I$, that is, for any $\epsilon > 0$, there exist $\delta_1, \delta_2 > 0$ such that $|S[\varphi](\theta_1, x_1) - S[\varphi](\theta_2, x_2)| < \epsilon$ for any $\varphi \in D$ if $\theta_1, \theta_2 \in [-\tau, 0], x_1, x_2 \in I$, and $|\theta_1 - \theta_2| < \delta_1, |x_1 - x_2| < \delta_2$.

Since $S[\varphi](\theta, x) = \varphi(0, x)$ for all $\varphi \in D, \theta \in [-\tau, -t_1], x \in I$, and $D(0, \cdot)$ is precompact in \mathbb{X} , it is obvious that $S[D]$ is equicontinuous on $[-\tau, -t_1] \times I$.

Note that there exists $N > 0$ such that $\|U(t, s)B(s, u_s)\|_{\mathbb{X}} \leq N$ for all $t \in [0, t_1], s \in [0, t], \varphi \in C_{\widehat{L}}$. Let $\delta_0 = \min\{\epsilon/(2^K N), t_1\}$. Then for any $t < \delta_0, x \in I$, and $\varphi \in D$, we have

$$(2.14) \quad \left| \int_0^t U(t, s)B(s, u_s)(x)ds \right| < 2^K N \delta_0 = \epsilon.$$

Let $\mathcal{F}(t, \psi) := U(t, 0)\psi$ for $(t, \psi) \in [0, \delta_0] \times D(0, \cdot)$. Then \mathcal{F} is continuous on $[0, \delta_0] \times D(0, \cdot)$ and $\mathcal{F}([0, \delta_0] \times D(0, \cdot))$ is precompact in \mathbb{X} . Thus, for the above I , there exists $\delta_2 > 0$ such that for $x_1, x_2 \in I$ and $|x_1 - x_2| < \delta_2$, we have

$$(2.15) \quad |U(t, 0)\psi(x_1) - U(t, 0)\psi(x_2)| < \epsilon \quad \forall t \in [0, \delta_0], \psi \in D(0, \cdot).$$

Moreover, since \mathcal{F} is uniformly continuous on $[0, \delta_0] \times D(0, \cdot)$, there exists $\delta_1 > 0, \delta_3 > 0$ such that $\|\mathcal{F}(\bar{t}_1, \psi_1) - \mathcal{F}(\bar{t}_2, \psi_2)\|_{\mathbb{X}} < \epsilon/2^K$ if $\bar{t}_1, \bar{t}_2 \in [0, \delta_0], \psi_1, \psi_2 \in D(0, \cdot)$, and $|\bar{t}_1 - \bar{t}_2| < \delta_1, \|\psi_1 - \psi_2\|_{\mathbb{X}} < \delta_3$. In particular, we have $\|U(\bar{t}_1, 0)\psi - U(\bar{t}_2, 0)\psi\|_{\mathbb{X}} < \epsilon/2^K$ if $\bar{t}_1, \bar{t}_2 \in [0, \delta_0], \psi \in D(0, \cdot)$, and $|\bar{t}_1 - \bar{t}_2| < \delta_1$. Then

$$(2.16) \quad |U(\bar{t}_1, 0)\psi(x) - U(\bar{t}_2, 0)\psi(x)| < \epsilon \quad \forall \psi \in D(0, \cdot), x \in I, \bar{t}_1, \bar{t}_2 \in [0, \delta_0], \text{ and } |\bar{t}_1 - \bar{t}_2| < \delta_1.$$

By (2.14)–(2.16), we can easily obtain that if $\theta_1, \theta_2 \in [-t_1, \delta_0 - t_1], x_1, x_2 \in I$ and $|\theta_1 - \theta_2| < \delta_1, |x_1 - x_2| < \delta_2$, then for any $\varphi \in D$,

$$\begin{aligned} & |S[\varphi](\theta_1, x_1) - S[\varphi](\theta_2, x_2)| \\ &= |Q_{t_1}[\varphi](\theta_1, x_1) - Q_{t_1}[\varphi](\theta_2, x_2)| \\ &= |u(t_1 + \theta_1, x_1, \varphi) - u(t_1 + \theta_2, x_2, \varphi)| \\ &\leq |(U(t_1 + \theta_1, 0)\varphi(0, \cdot))(x_1) - (U(t_1 + \theta_2, 0)\varphi(0, \cdot))(x_2)| \\ &\quad + \left| \int_0^{t_1 + \theta_1} U(t_1 + \theta_1, s)B(s, u_s)(x_1)ds - \int_0^{t_1 + \theta_2} U(t_1 + \theta_2, s)B(s, u_s)(x_2)ds \right| \\ &\leq |(U(t_1 + \theta_1, 0)\varphi(0, \cdot))(x_1) - (U(t_1 + \theta_1, 0)\varphi(0, \cdot))(x_2)| \\ &\quad + |(U(t_1 + \theta_1, 0)\varphi(0, \cdot))(x_2) - (U(t_1 + \theta_2, 0)\varphi(0, \cdot))(x_2)| + 2\epsilon \\ &< 4\epsilon, \end{aligned}$$

which implies that $S[D]$ is equicontinuous on $[-t_1, \delta_0 - t_1] \times I$.

By a similar argument as for (A6)(a), it is easy to see that $S[D]$ is equicontinuous on $[\delta_0 - t_1, 0] \times I$.

Therefore, $S[D]$ is equicontinuous on $[-\tau, 0] \times I$, and hence, $S[D]$ is precompact in $C_{\widehat{L}}$. Thus, (A6)(b') is valid for $Q_t, t \in (0, \tau]$. \square

It then follows from Lemma 2.6 and [11, Theorems 2.11 and 2.15] that Q_ω has an asymptotic speed of spread $c_\omega^* > 0$.

Consider the linearized system of (2.1) at the zero solution:

$$(2.17) \quad \begin{cases} \partial_t u(t, x) = d(t)\Delta u - g_u(t, 0)u(t, x) + b(t)\partial_u f_{-\tau}(t, 0) \int_{\mathbb{R}} \Gamma(a(t), x - y)u(t - \tau, y)dy, \\ u(t, x) = \phi(t, x), \quad t \in [-\tau, 0], \quad x \in \mathbb{R}. \end{cases}$$

For $\alpha > 0$, let $u(t, x) = e^{-\alpha x}v(t)$. Substituting $u(t, x)$ into (2.17) yields

$$e^{-\alpha x}v'(t) = d(t)\alpha^2 e^{-\alpha x}v(t) - g_u(t, 0)v(t)e^{-\alpha x} + b(t)\partial_u f_{-\tau}(t, 0)v(t - \tau) \int_{\mathbb{R}} \Gamma(a(t), y)e^{-\alpha(x-y)} dy.$$

Since $\Gamma(t, x)$ is even in x and by [18, Proposition 4.2], we obtain

$$(2.18) \quad \begin{aligned} v'(t) &= d(t)\alpha^2 v(t) - g_u(t, 0)v(t) + b(t)\partial_u f_{-\tau}(t, 0)v(t - \tau) \int_{\mathbb{R}} \Gamma(a(t), y)e^{\alpha y} dy, \\ &= d(t)\alpha^2 v(t) - g_u(t, 0)v(t) + b(t)\partial_u f_{-\tau}(t, 0)v(t - \tau) \int_{\mathbb{R}} \Gamma(a(t), y)e^{-\alpha y} dy, \\ &= d(t)\alpha^2 v(t) - g_u(t, 0)v(t) + b(t)\partial_u f_{-\tau}(t, 0)v(t - \tau)e^{\alpha^2 a(t)}. \end{aligned}$$

Then $u(t, x) = e^{-\alpha x}v(t)$ satisfies (2.17) with $\phi(s, x) = e^{-\alpha x}v(s)$ for $s \in [-\tau, 0]$ and $x \in \mathbb{R}$ if $v(t)$ satisfies (2.18) for $t \geq 0$.

Let M_t be the linear solution map defined by (2.17) and let $v(t, v_0)$ be the solution of (2.18) with $v(s, v_0) = v_0(s)$ for $s \in [-\tau, 0], v_0 \in \mathbb{Y}$. Define $B_\alpha^t(v_0) := M_t(v_0 e^{-\alpha x})(0)$. It is not difficult to see that $B_\alpha^t(v_0) = v(t, v_0)$, and hence, B_α^t is the solution map associated with (2.18) on \mathbb{Y} .

Let $\gamma(\alpha)$ be the spectral radius of the Poincaré map associated with (2.18), and [21, Proposition 2.1] implies that $\gamma(\alpha) > 0$. It follows from the proof of [21, Proposition 2.1] that there exists a positive ω -periodic function $w(t)$ such that $v(t) = e^{\lambda(\alpha)t}w(t)$ is a solution of (2.18), where $\lambda(\alpha) = \frac{\ln \gamma(\alpha)}{\omega}$. Define $\psi \in \mathbb{Y}$ by $\psi(\theta) = e^{\lambda(\alpha)\theta}w(\theta)$ for all $\theta \in [-\tau, 0]$. Clearly, $v(t, \psi) = e^{\lambda(\alpha)t}w(t)$ for all $t \geq 0$. Then we have

$$B_\alpha^t(\psi)(\theta) = v(t + \theta, \psi) = e^{\lambda(\alpha)t}e^{\lambda(\alpha)\theta}w(t + \theta) \quad \forall \theta \in [-\tau, 0], t \geq 0.$$

By the ω -periodicity of $w(t)$, it follows that

$$B_\alpha^\omega(\psi)(\theta) = e^{\lambda(\alpha)\omega}e^{\lambda(\alpha)\theta}w(\theta) = e^{\lambda(\alpha)\omega}\psi(\theta) \quad \forall \theta \in [-\tau, 0],$$

that is, $B_\alpha^\omega(\psi) = e^{\lambda(\alpha)\omega}\psi$. This implies that $e^{\lambda(\alpha)\omega}$ is the principle eigenvalue of B_α^ω with positive eigenfunction ψ .

Let $\Phi(\alpha) := \frac{1}{\alpha} \ln e^{\lambda(\alpha)\omega} = \frac{\lambda(\alpha)\omega}{\alpha} = \frac{\ln \gamma(\alpha)}{\alpha}$. Then we have the following result.

PROPOSITION 2.7. Assume that (H1)–(H4) hold. Let c_ω^* be the asymptotic speed of spread of Q_ω . Then $c_\omega^* = \inf_{\alpha>0} \Phi(\alpha) = \inf_{\alpha>0} \frac{\ln \gamma(\alpha)}{\alpha}$.

Proof. When $\alpha = 0$, (2.18) becomes (2.3). It follows from (H4) that $\gamma(0) > 1$, and hence (C7) in [11] is satisfied. Now we prove that $\Phi(\infty) = \infty$. By (2.18), we have

$$v'(t) \geq [\alpha^2 d(t) - g_u(t, 0)]v(t) \quad \forall t \geq 0,$$

and hence,

$$\frac{w'(t)}{w(t)} \geq \alpha^2 d(t) - g_u(t, 0) - \lambda(\alpha).$$

Then

$$0 = \int_0^\omega \frac{w'(t)}{w(t)} dt \geq \int_0^\omega (\alpha^2 d(t) - g_u(t, 0)) dt - \lambda(\alpha)\omega,$$

which implies that

$$\lambda(\alpha)\omega \geq \alpha^2 \int_0^\omega d(t) dt - \int_0^\omega g_u(t, 0) dt.$$

Therefore,

$$\Phi(\alpha) = \frac{\lambda(\alpha)\omega}{\alpha} \geq \alpha \int_0^\omega d(t) dt - \frac{\int_0^\omega g_u(t, 0) dt}{\alpha}.$$

Letting $\alpha \rightarrow \infty$, we can easily obtain $\Phi(\infty) = \infty$.

Since $G(t, \cdot, \cdot)$ is subhomogeneous in (u, v) , it follows from [23, Lemma 2.3.2] that $G(t, u, v) \leq G_u(t, 0, 0)u + G_v(t, 0, 0)v$, that is,

$$-g(t, u) + b(t)f_{-\tau}(t, v) \leq -g_u(t, 0)u + b(t)\partial_u f_{-\tau}(t, 0)v,$$

and hence, we have

$$\begin{aligned} & -g(t, u(t, x)) + b(t) \int_{\mathbb{R}} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y)) dy \\ & \leq -g_u(t, 0)u(t, x) + b(t)\partial_u f_{-\tau}(t, 0) \int_{\mathbb{R}} \Gamma(a(t), x - y) u(t - \tau, y) dy. \end{aligned}$$

By the comparison principle, we have $Q_\omega(\varphi) \leq M_\omega(\varphi)$ for any $\varphi \in C_{\beta_0^*}$. Thus, [11, Theorem 3.10] implies that $c_\omega^* \leq \inf_{\alpha>0} \Phi(\alpha)$.

Let $K > 0$ such that $K - g_u(t, 0) > 0$ for all $t \in [0, \omega]$. Set $\bar{G}(t, u, v) = Ku + G(t, u, v)$. Then $\bar{G}_u(t, 0, 0) > 0$, $\bar{G}_v(t, 0, 0) > 0$ for all $t \in [0, \omega]$. It is easy to see that for any $\epsilon \in (0, 1)$, there exists $\delta = \delta(\epsilon) \in (0, L)$ such that

$$\bar{G}(t, u, v) \geq (1 - \epsilon)\bar{G}_u(t, 0, 0)u + (1 - \epsilon)\bar{G}_v(t, 0, 0)v \quad \forall (u, v) \in [0, \delta]^2,$$

and hence, for any $(u, v) \in [0, \delta]^2$,

$$G(t, u, v) = -Ku + \bar{G}(t, u, v) \geq [(1 - \epsilon)G_u(t, 0, 0) - \epsilon K]u + (1 - \epsilon)G_v(t, 0, 0)v.$$

Moreover, there exists $\xi = \xi(\delta) > 0$ such that for any $\varphi \in C_{\hat{\xi}}$, we have

$$0 \leq u(t, x, \varphi) \leq u(t, x, \hat{\xi}) < \delta \quad \forall x \in \mathbb{R}, t \in [0, \omega].$$

Thus, for any $\varphi \in C_{\hat{\xi}}$, $u(t, x, \varphi)$ satisfies

$$\begin{aligned} \partial_t u(t, x) &\geq d(t)\Delta u(t, x) + [(1 - \epsilon)g_u(t, 0) - \epsilon K]u(t, x) \\ &\quad + (1 - \epsilon)b(t)\partial_u f_{-\tau}(t, 0) \int_{\mathbb{R}} \Gamma(a(t), x - y)u(t - \tau, y)dy \quad \forall t \in [0, \omega]. \end{aligned}$$

Let $M_t^\epsilon, t \geq 0$, be the solution maps associated with the linear system

$$\begin{aligned} \partial_t u(t, x) &= d(t)\Delta u(t, x) + [(1 - \epsilon)g_u(t, 0) - \epsilon K]u(t, x) \\ &\quad + (1 - \epsilon)b(t)\partial_u f_{-\tau}(t, 0) \int_{\mathbb{R}} \Gamma(a(t), x - y)u(t - \tau, y)dy \quad \forall t \in [0, \omega]. \end{aligned}$$

The comparison principle implies that $M_t^\epsilon(\varphi) \leq Q_t(\varphi)$ for all $\varphi \in C_{\hat{\xi}}, t \in [0, \omega]$. In particular, $M_\omega^\epsilon(\varphi) \leq Q_\omega(\varphi)$ for all $\varphi \in C_{\hat{\xi}}$. By a similar analysis for M_t^ϵ as for M_t , it follows from [11, Theorem 3.10] that $\inf_{\alpha > 0} \Phi_\epsilon(\alpha) \leq c_\omega^*$.

Therefore, $\inf_{\alpha > 0} \Phi_\epsilon(\alpha) \leq c_\omega^* \leq \inf_{\alpha > 0} \Phi(\alpha)$ for all $\epsilon \in (0, 1)$. Letting $\epsilon \rightarrow 0$, we have $c_\omega^* = \inf_{\alpha > 0} \Phi(\alpha)$. \square

Let $c^* = \frac{c_\omega^*}{\omega} = \frac{1}{\omega} \inf_{\alpha > 0} \Phi(\alpha) = \frac{1}{\omega} \inf_{\alpha > 0} \frac{\ln \gamma(\alpha)}{\alpha}$. The following result shows that c^* is the spreading speed of solutions of (2.1) with initial functions having compact support.

THEOREM 2.8. *Assume that (H1)–(H4) hold and let $c^* = c_\omega^*/\omega$. Then the following statements are valid.*

- (1) *For any $c > c^*$, if $\varphi \in C_{\beta_0^*}$ with $0 \leq \varphi \ll \beta_0^*$ and $\varphi(\cdot, x) = 0$ for x outside a bounded interval, then*

$$\lim_{t \rightarrow \infty, |x| \geq ct} u(t, x, \varphi) = 0.$$

- (2) *For any $c < c^*$, if $\varphi \in C_{\beta_0^*}$ with $\varphi \not\equiv 0$, then*

$$\lim_{t \rightarrow \infty, |x| \leq ct} (u(t, x, \varphi) - \beta^*(t)) = 0.$$

Proof. Conclusion (1) follows from [10, Theorem 2.1]. By Lemma 2.4 and [10, Theorem 2.1], for any $c < c^*$, there is a positive number σ such that, if $\varphi \in C_{\beta_0^*}$ with $\varphi(\cdot, x) > 0$ for x on an interval of length 2σ , then $\lim_{t \rightarrow \infty, |x| \leq ct} (u(t, x, \varphi) - \beta^*(t)) = 0$. It follows from Lemma 2.5 that for any $\varphi \in C_{\beta_0^*}$ with $\varphi \not\equiv 0$, $Q_t(\varphi) \gg 0$ for all $t > 2\tau$. We can fix $t_0 > 2\tau$ and take $Q_{t_0}(\varphi)$ as a new initial value for $u(t, x, \varphi)$. Then by the above analysis, conclusion (2) is valid. \square

We say $u(t, x) = \mathcal{U}(t, x - ct)$ is an ω -periodic traveling wave of (2.1) connecting $\beta^*(t)$ to 0 if it is a solution of (2.1), $\mathcal{U}(t, \xi)$ is ω -periodic in t , and $\mathcal{U}(t, -\infty) = \beta^*(t)$ and $\mathcal{U}(t, \infty) = 0$ uniformly for $t \in [0, \omega]$. By [10, Theorems 2.2 and 2.3], we have the following result about traveling waves of (2.1).

THEOREM 2.9. *Assume that (H1)–(H4) hold. Let c^* be defined as $c^* = c_\omega^*/\omega$. Then for any $c \geq c^*$, (2.1) has an ω -periodic traveling wave solution $\mathcal{U}(t, x - ct)$ connecting*

$\beta^*(t)$ to 0 such that $U(t, s)$ is continuous and nonincreasing in s . Moreover, for any $c < c^*$, (2.1) has no ω -periodic traveling wave $U(t, x - ct)$ connecting $\beta^*(t)$ to 0.

3. Dynamics in a bounded domain. In this section, we consider (1.6) in a bounded spatial domain

$$(3.1) \begin{cases} \partial_t u(t, x) = d(t)\Delta u - g(t, u) + b(t) \int_{\Omega} \Gamma(a(t), x - y) f_{-\tau}(t, u(t - \tau, y)) dy, \\ (t, x) \in (0, \infty) \times \Omega, \\ Bu(t, x) = 0 \text{ on } (0, \infty) \times \partial\Omega, \\ u(t, x) = \phi(t, x), \quad t \in [-\tau, 0], \quad x \in \Omega, \end{cases}$$

where $\Omega \subseteq \mathbb{R}^n$ is a bounded domain with boundary $\partial\Omega$ of class $C^{1+\theta}$ ($0 < \theta \leq 1$), the boundary condition is either $Bu = u$ (Dirichlet boundary condition) or $Bu = (\partial u / \partial \nu) + \alpha(x)u$ (Robin type boundary condition) for some nonnegative function $\alpha \in C^{1+\theta}(\partial\Omega, \mathbb{R})$, and $\partial u / \partial \nu$ denotes the differentiation in the direction of outward normal ν to $\partial\Omega$.

Let $p \in (1, \infty)$ be fixed. For each $\beta \in (\frac{1}{2} + \frac{1}{2p}, 1)$, let \mathbb{X}_β be the fractional power space of $L^p(\Omega)$ with respect to $-\Delta$ and the boundary condition $Bu = 0$ (see, e.g., [8]). Then \mathbb{X}_β is an ordered Banach space with the positive cone \mathbb{X}_β^+ consisting of all nonnegative functions in \mathbb{X}_β , and \mathbb{X}_β^+ has nonempty interior $\text{int}(\mathbb{X}_\beta^+)$. Moreover, $\mathbb{X}_\beta \subseteq C^{1+\nu}(\bar{\Omega})$ with continuous inclusion for $\nu \in [0, 2\beta - 1 - \frac{1}{p}]$. Denote the norm on \mathbb{X}_β by $\|\cdot\|_\beta$. Then there exists a constant $k_\beta > 0$ such that $\|\phi\|_\infty := \max_{x \in \bar{\Omega}} |\phi(x)| \leq k_\beta \|\phi\|_\beta$ for all $\phi \in \mathbb{X}_\beta$.

Let $\bar{C} = C([-\tau, 0], \mathbb{X}_\beta)$ and $\bar{C}^+ = C([-\tau, 0], \mathbb{X}_\beta^+)$. For convenience, we identify an element $\phi \in \bar{C}$ as a function from $[-\tau, 0] \times \bar{\Omega}$ to \mathbb{R} defined by $\phi(s, x) = \phi(s)(x)$. For any $N \geq L$, let $\bar{C}_N = \{\phi \in \bar{C} : 0 \leq \phi(s, x) \leq N, (s, x) \in [-\tau, 0] \times \bar{\Omega}\}$. For any function $y(\cdot) : [-\tau, b] \rightarrow \mathbb{X}_\beta$, where $b > 0$, define $y_t \in \bar{C}$, by $y_t(s) = y(t + s)$ for all $s \in [-\tau, 0]$, $t \in [0, b]$.

Note that the differential operator Δ generates an analytic semigroup $\bar{U}_0(t)$ on $L^p(\Omega)$ and that the standard parabolic maximum principle (see, e.g., [15, Corollary 7.2.3]) implies that the semigroup $\bar{U}_0(t) : \mathbb{X}_\beta \rightarrow \mathbb{X}_\beta$ is strongly positive in the sense that $\bar{U}_0(t)(\mathbb{X}_\beta^+ \setminus \{0\}) \subseteq \text{int}(\mathbb{X}_\beta^+)$ for all $t > 0$. By a similar analysis as in section 2, we can write (3.1) as an integral equation (2.7) with $u_0 = \phi \in \bar{C}^+$. It then follows from [13, Corollary 5] that, for any $\phi \in \bar{C}_L^+$, (3.1) has a unique mild solution $u(t, x, \phi)$ with $u_0(\cdot, \cdot, \phi) = \phi$ and $u_t(\cdot, \cdot, \phi) \in \bar{C}_L^+$ for all $t \geq 0$. Moreover, $u(t, x, \phi)$ is a classic solution when $t > \tau$ and the comparison theorem holds for (3.1).

Define a family of operators $\{Q_t\}_{t \geq 0}$ on \bar{C}^+ by

$$Q_t(\phi)(s, x) = u(t + s, x, \phi) \quad \forall \phi \in \bar{C}^+, x \in \bar{\Omega}, \quad t \geq 0, \quad s \in [-\tau, 0].$$

Similarly as in section 2, we can show that $\{Q_t\}_{t \geq 0}$ is a monotone ω -periodic semiflow on \bar{C}^+ ; $u(t, x, \phi) > 0$ for $t > \tau$, $x \in \bar{\Omega}$, $\phi \in \bar{C}^+$ with $\phi \not\equiv 0$, and hence, Q_t is strongly positive for $t > 2\tau$; moreover, Q_t is compact on \bar{C}^+ for all $t > \tau$. Let $n_1 = \min\{n \in \mathbb{N}, n\omega > 2\tau\}$. Then $Q_{n_1\omega}$ is compact and strongly positive on \bar{C}^+ . We

can further show that the periodic semiflow $\{Q_t\}_{t \geq 0}$ is point dissipative on \bar{C}^+ . By [23, Theorem 1.1.3], we have the following result.

LEMMA 3.1. *Let (H1)–(H3) hold. Then $Q_{n_1\omega}$ admits a global attractor on \bar{C}^+ . Consider the linearized system of (3.1) at the zero solution*

$$(3.2) \quad \begin{cases} \partial_t \tilde{u}(t, x) = d(t)\Delta \tilde{u} - g_u(t, 0)\tilde{u}(t, x) + b(t)\partial_u f_{-\tau}(t, 0) \int_{\Omega} \Gamma(a(t), x - y)\tilde{u}(t - \tau, y)dy, \\ t > 0, x \in \Omega, \\ B\tilde{u}(t, x) = 0, \quad t > 0, x \in \partial\Omega, \\ \tilde{u}(s, x) = \phi(s, x), \quad s \in [-\tau, 0], x \in \Omega, \phi \in \bar{C}. \end{cases}$$

Similarly as in Theorem 2.3, we can show that the comparison principle holds for (3.2), and hence, the solution map \tilde{u}_t of (3.2) is monotone increasing for all $t \geq 0$.

Now we consider (3.1) and (3.2) as $n_1\omega$ -periodic systems. Define the Poincaré map of (3.2) $P_1 : \bar{C} \rightarrow \bar{C}$ by $P_1(\phi) = \tilde{u}_{n_1\omega}(\phi)$ for all $\phi \in \bar{C}$, where $\tilde{u}_{n_1\omega}(\phi)(s, x) = \tilde{u}(n_1\omega + s, x, \phi)$ for all $(s, x) \in [-\tau, 0] \times \bar{\Omega}$, and $\tilde{u}(t, x, \phi)$ is the solution of (3.2) with $\tilde{u}(s, x) = \phi(s, x)$ for all $(s, x) \in [-\tau, 0] \times \bar{\Omega}$. Similarly as in section 2, we can obtain that P_1 is also compact and strongly positive. Let $r_1 = r(P_1)$ be the spectral radius of P_1 . By the Krein–Rutman theorem (see, e.g., [8, Theorem 7.2]), $r_1 > 0$ and P_1 has a positive eigenfunction $\bar{\phi} \in \text{int}(\bar{C}^+)$ corresponding to r_1 .

LEMMA 3.2. *Let $\mu = -\frac{1}{n_1\omega} \ln r_1$. Then there exists a positive $n_1\omega$ -periodic function $v(t, x)$ such that $e^{-\mu t}v(t, x)$ is a solution of (3.2).*

Proof. By the definitions of r_1 and $\bar{\phi}$, we have $P_1\bar{\phi} = r_1\bar{\phi}$. Let $\tilde{u}(t, x, \bar{\phi})$ be the solution of (3.2) with $\tilde{u}(s, x) = \bar{\phi}(s, x)$ for all $s \in [-\tau, 0], x \in \Omega$. Since $\bar{\phi} \gg 0$, it is not difficult to see that $\tilde{u}(\cdot, \cdot, \bar{\phi}) \gg 0$. Let $\mu = -\frac{1}{n_1\omega} \ln r_1$ and $v(t, x) = e^{\mu t}\tilde{u}(t, x, \bar{\phi})$ for all $t \geq -\tau, x \in \Omega$. Then $r_1 = e^{-n_1\omega\mu}$ and $v(t, x) > 0$ for all $t \in [-\tau, \infty), x \in \Omega$. Moreover,

$$(3.3) \quad \begin{aligned} &v_t(t, x) \\ &= e^{\mu t}\tilde{u}_t(t, x, \bar{\phi}) + \mu e^{\mu t}\tilde{u}(t, x, \bar{\phi}) \\ &= e^{\mu t}[d(t)\Delta \tilde{u} - g_u(t, 0)\tilde{u}(t, x, \bar{\phi}) \\ &\quad + b(t)\partial_u f_{-\tau}(t, 0) \int_{\Omega} \Gamma(a(t), x - y)\tilde{u}(t - \tau, y, \bar{\phi})dy] + \mu v \\ &= d(t)\Delta v - g_u(t, 0)v(t, x) + e^{\mu\tau}b(t)\partial_u f_{-\tau}(t, 0) \int_{\Omega} \Gamma(a(t), x - y)v(t - \tau, y)dy + \mu v \end{aligned}$$

for all $(t, x) \in (0, \infty) \times \Omega$. Thus, $v(t, x)$ is a solution of $n_1\omega$ -periodic equation (3.3) with $Bv = 0$ on $(0, \infty) \times \partial\Omega$ and $v(s, x) = e^{\mu s}\bar{\phi}(s, x)$ for all $s \in [-\tau, 0], x \in \Omega$.

For any $\theta \in [-\tau, 0], x \in \Omega$, we have

$$v(n_1\omega + \theta, x) = e^{\mu(n_1\omega + \theta)} \cdot P_1(\bar{\phi})(\theta, x) = e^{\mu(n_1\omega + \theta)} \cdot r_1\bar{\phi}(\theta, x) = e^{\mu\theta} \cdot \tilde{u}(\theta, x, \bar{\phi}) = v(\theta, x).$$

Therefore, $v_0(\theta, \cdot) = v_{n_1\omega}(\theta, \cdot)$ for all $\theta \in [-\tau, 0]$, and hence, the existence and uniqueness of solutions of (3.3) imply that

$$v(t, x) = v(t + n_1\omega, x) \quad \forall t \geq -\tau, x \in \Omega,$$

that is, $v(t, x)$ is an $n_1\omega$ -periodic solution of (3.3). Clearly, $e^{-\mu t}v(t, x)$ is a solution of (3.2). \square

Define $P_0 : \bar{C} \rightarrow \bar{C}$ by $P_0(\phi) = \tilde{u}_\omega(\phi)$ for all $\phi \in \bar{C}$, where $\tilde{u}(t, x, \phi)$ is the solution of (3.2) with $\tilde{u}(s, x) = \phi(s, x)$ for all $s \in [-\tau, 0]$, $x \in \Omega$. Let $r_0 = r(P_0)$ be the spectral radius of P_0 .

THEOREM 3.3. *Let (H1)–(H3) hold. For any $\phi \in \bar{C}^+$, denote by $u(t, x, \phi)$ the solution of (3.1) with $u(s, x) = \phi(s, x)$ for all $(t, x) \in [-\tau, 0] \times \Omega$. Then the following two statements are valid.*

- (i) *If $r_0 < 1$, then $\lim_{t \rightarrow \infty} \|u(t, \cdot, \phi)\|_\beta = 0$ for every $\phi \in \bar{C}^+$.*
- (ii) *If $r_0 > 1$, then (3.1) admits a unique positive ω -periodic solution $u^*(t, x)$ and $\lim_{t \rightarrow \infty} \|u(t, \cdot, \phi) - u^*(t, \cdot)\|_\beta = 0$ for all $\phi \in \bar{C}^+ \setminus \{0\}$.*

Proof. Since $P_1 = \tilde{u}_{n_1\omega}$, $P_0 = \tilde{u}_\omega$, and $\tilde{u}_{n_1\omega} = \tilde{u}_\omega^{n_1}$, where \tilde{u}_t is the solution map of (3.2), by the properties of spectral radius of linear operators, we know that $r(P_1) = (r(P_0))^{n_1}$, i.e., $r_1 = (r_0)^{n_1}$. Note that the qualitative solutions of (3.1) and (3.2) do not change whether we consider them as $n_1\omega$ -periodic systems or ω -periodic systems. The conditions in Theorem 3.3 can be replaced by $r_1 < 1$ and $r_1 > 1$, respectively. In the following, we will consider (3.1) and (3.2) as $n_1\omega$ -periodic systems and prove the theorem under the conditions of $r_1 < 1$ and $r_1 > 1$.

In the case where $r_1 < 1$, we have $\mu = -\frac{1}{n_1\omega} \ln r_1 > 0$. By Lemma 3.2, (3.2) has a solution $\tilde{u}(t, x) := \tilde{u}(t, x, \bar{\phi}) = e^{-\mu t}v(t, x)$ with $\tilde{u}(s, x) = \bar{\phi}(s, x)$ for all $(s, x) \in [-\tau, 0] \times \Omega$, where $\bar{\phi} \in \text{int}(\bar{C}^+)$ is the positive eigenfunction of P_1 corresponding to r_1 and $v(t, x)$ is $n_1\omega$ -periodic in $t \geq -\tau$. Then v is bounded on $[-\tau, \infty) \times \bar{\Omega}$, and hence, there exists $\rho > 0$ such that $\|v(t, \cdot)\|_\infty \leq \rho$ for all $t \geq -\tau$. Thus, $\lim_{t \rightarrow \infty} \|\tilde{u}(t, \cdot)\|_\infty = 0$. By the basic analysis of solutions of (3.2), it follows that $\lim_{t \rightarrow \infty} \|\tilde{u}(t, \cdot)\|_\beta = 0$.

Given $\phi \in \bar{C}^+$, since $\lim_{\delta \rightarrow 0^+} (\bar{\phi} - \delta\phi) = \bar{\phi} \in \text{int}(\bar{C}^+)$ for any $\epsilon > 0$, there exists $\delta_\phi > 0$, such that $\bar{\phi} - \delta\phi \in B_\epsilon(\bar{\phi}) \subseteq \bar{C}^+$ for $0 < \delta \leq \delta_\phi$, where $B_\epsilon(\bar{\phi})$ is an open ball in \bar{C}^+ centered at $\bar{\phi}$ with radius ϵ . Therefore, $\bar{\phi} \geq \delta_\phi\phi$ in \bar{C}^+ . It then follows from the comparison principle that $\tilde{u}(t, x) \geq \delta_\phi\tilde{u}(t, x, \phi)$ for all $t \geq -\tau, x \in \bar{\Omega}$, where $\tilde{u}(t, \cdot, \phi)$ is the solution of (3.2) with $\tilde{u}(s, x) = \phi(s, x)$ for all $(s, x) \in [-\tau, 0] \times \Omega$. Thus, $\lim_{t \rightarrow \infty} \|\tilde{u}(t, \cdot, \phi)\|_\infty = 0$, and hence, $\lim_{t \rightarrow \infty} \|\tilde{u}(t, \cdot, \phi)\|_\beta = 0$ for any $\phi \in \bar{C}^+$.

Note that a solution of (3.1) satisfies

$$\partial_t u(t, x) \leq d(t)\Delta u - g_u(t, 0)u(t, x) + b(t)\partial_u f_{-\tau}(t, 0) \int_\Omega \Gamma(a(t), x - y)u(t - \tau, y)dy$$

for any $t > 0, x \in \Omega$. Similarly to the proof of Theorem 2.3, we can show that the comparison theorem for abstract functional differential equations [13, Proposition 3] can be applied to (3.1) and (3.2). Therefore, for any $\phi \in \bar{C}^+, u(t, \cdot, \phi) \leq \tilde{u}(t, \cdot, \phi)$ for all $t \geq -\tau$, where $u(t, \cdot, \phi)$ and $\tilde{u}(t, \cdot, \phi)$ are solutions of (3.1) and (3.2), respectively. It then follows that solutions of (3.1) satisfy $\lim_{t \rightarrow \infty} \|u(t, \cdot, \phi)\|_\beta = 0$ for all $\phi \in \bar{C}^+$.

In the case where $r_1 > 1$, we have $\mu < 0$. Let $\bar{C}_0 = \{\phi \in \bar{C}^+ : \phi \not\equiv 0\}$, $\partial\bar{C}_0 = \bar{C}^+ \setminus \bar{C}_0 = \{0\}$. Similarly to the proof of Lemma 2.5, we can show that for any $\phi \in \bar{C}_0$, the solution $u(t, x, \phi)$ of (3.1) satisfies $u(t, x, \phi) > 0$ for all $t > \tau, x \in \Omega$. It follows that $Q_t(\bar{C}_0) \subseteq \text{int}(\bar{C}^+)$ for all $t > 2\tau$. Clearly, $Q_t(0) = 0$ for all $t \geq 0$. We now have the following claim.

Claim. Zero is a uniform weak repeller for \bar{C}_0 in the sense that there exists $\delta_0 > 0$ such that $\lim_{t \rightarrow \infty} \sup \|Q_t(\phi)\|_\beta \geq \delta_0$ for all $\phi \in \bar{C}_0$.

Indeed, we consider the following system:

$$(3.4) \quad \begin{cases} \partial_t u^\varepsilon(t, x) = d(t)\Delta u^\varepsilon - (g_u(t, 0) + \varepsilon)u^\varepsilon(t, x) \\ \quad + b(t)(\partial_u f_{-\tau}(t, 0) - \varepsilon) \int_\Omega \Gamma(a(t), x - y)u^\varepsilon(t - \tau, y)dy, \\ Bu^\varepsilon(t, x) = 0, \quad t > 0, x \in \partial\Omega, \\ u^\varepsilon(s, x) = \phi(s, x), \quad s \in [-\tau, 0], x \in \Omega, \phi \in \bar{C}. \end{cases}$$

Define the Poincaré map of (3.4) $P_\varepsilon : \bar{C} \rightarrow \bar{C}$ by

$$P_\varepsilon(\phi) = u_{n_1\omega}^\varepsilon(\phi) \quad \forall \phi \in \bar{C},$$

where

$$u_{n_1\omega}^\varepsilon(\phi)(s, x) = u^\varepsilon(n_1\omega + s, x, \phi) \quad \forall (s, x) \in [-\tau, 0] \times \bar{\Omega}$$

and $u^\varepsilon(t, x, \phi)$ is the solution of (3.4) with $u^\varepsilon(s, x) = \phi(s, x)$ for all $s \in [-\tau, 0], x \in \Omega$. Let $r_\varepsilon = r(P_\varepsilon)$ be the spectral radius of P_ε . Since $r_1 = r(P_1) > 1$, there exists a sufficiently small positive number ε_1 such that $r_\varepsilon > 1$ for all $\varepsilon \in [0, \varepsilon_1]$. We fix an $\varepsilon \in (0, \varepsilon_1)$. Since $\lim_{u \rightarrow 0^+} \frac{g(t, u)}{u} = g_u(t, 0)$ and $\lim_{u \rightarrow 0^+} \frac{f_{-\tau}(t, u)}{u} = \partial_u f_{-\tau}(t, 0)$ uniformly for $t \in [0, n_1\omega]$, there exists $\delta_\varepsilon > 0$ such that $g(t, u) < (g_u(t, 0) + \varepsilon)u$ and $f_{-\tau}(t, u) > (\partial_u f_{-\tau}(t, 0) - \varepsilon)u$ for $u \in (0, \delta_\varepsilon), t \in [0, n_1\omega]$. Let $\delta_0 = \delta_\varepsilon/k_\beta$. Suppose, by contradiction, that there exists $\phi_0 \in \bar{C}_0$ such that $\lim_{t \rightarrow \infty} \sup \|Q_t(\phi)\|_\beta < \delta_0$. Then there exists $t_0 > \tau$ such that $\|u(t, \cdot, \phi_0)\|_\infty \leq k_\beta \|u(t, \cdot, \phi_0)\|_\beta < \delta_\varepsilon$ for all $t \geq t_0$. Therefore, $u(t, x, \phi_0)$ satisfies

$$(3.5) \quad \begin{aligned} &\partial_t u(t, x) \\ &> d(t)\Delta u - (g_u(t, 0) + \varepsilon)u(t, x) + b(t)(\partial_u f_{-\tau}(t, 0) - \varepsilon) \int_\Omega \Gamma(a(t), x - y)u(t - \tau, y)dy \end{aligned}$$

for $t \geq t_0, x \in \Omega$. Let $\bar{\phi}_\varepsilon$ be the positive eigenfunction of P_ε associated with r_ε and $\mu_\varepsilon = -\frac{1}{n_1\omega} \ln r_\varepsilon$. Then by Lemma 3.2, the solution $u^\varepsilon(t, x, \bar{\phi}_\varepsilon)$ of (3.4) with $u^\varepsilon(s, x) = \bar{\phi}_\varepsilon(s, x)$ for all $s \in [-\tau, 0], x \in \Omega$, satisfies $u^\varepsilon(t, x, \bar{\phi}_\varepsilon) = e^{-\mu_\varepsilon t} v_\varepsilon(t, x)$, where $v_\varepsilon(t, x)$ is a positive $n_1\omega$ -periodic function in $t \geq -\tau$. Since $u(t, x, \phi_0) > 0$ for all $t \geq \tau, x \in \Omega$, there exists $\zeta > 0$ such that

$$u(t_0 + s, x, \phi_0) \geq \zeta u^\varepsilon(s, x, \bar{\phi}_\varepsilon) = \zeta \bar{\phi}_\varepsilon(s, x) \quad \forall s \in [-\tau, 0], x \in \bar{\Omega}.$$

By (3.5) and the comparison theorem, we have

$$u(t, x, \phi_0) \geq \zeta u^\varepsilon(t - t_0, x, \bar{\phi}_\varepsilon) = \zeta e^{-\mu_\varepsilon(t-t_0)} v_\varepsilon(t, x) \quad \forall t \geq t_0, x \in \bar{\Omega}.$$

Since $\mu_\varepsilon < 0$, it follows that $u(t, x, \phi_0)$ is unbounded, a contradiction. Thus, the claim is true.

By the claim above, $Q_{n_1\omega}$ is weakly uniformly persistent with respect to $(\bar{C}_0, \partial\bar{C}_0)$. Since $Q_{n_1\omega}$ admits a global attractor on \bar{C}^+ , it follows from [23, Theorem 1.3.3] that $Q_{n_1\omega}$ is uniformly persistent with respect to $(\bar{C}_0, \partial\bar{C}_0)$ in the sense that there exists $\delta_1 > 0$ such that $\lim_{n \rightarrow \infty} \inf \|Q_{n_1\omega}^n(\phi)\|_\beta \geq \delta_1$ for all $\phi \in \bar{C}_0$.

Note that $Q_{n_1\omega}$ is compact, point dissipative, and uniformly persistent. It follows from [23, Theorem 1.3.6] that $Q_{n_1\omega} : \bar{C}_0 \rightarrow \bar{C}_0$ admits a global attractor A_0 and

has a fixed point $\hat{\phi}$ in A_0 . Similarly as in the proof of Lemma 2.4, we can show that $Q_{n_1\omega}$ is strictly subhomogeneous. Then [22, Lemma 1] implies that $Q_{n_1\omega}$ has at most one fixed point. Thus, $Q_{n_1\omega}$ has a unique equilibrium $\hat{\phi}$ in \bar{C}_0 . Clearly, by the strong monotonicity of $Q_{n_1\omega}$, we have $\hat{\phi} \in \text{int}(\bar{C}^+)$. Moreover, it follows from [23, Theorem 2.3.2] that $A_0 = \{\hat{\phi}\}$ since $Q_{n_1\omega}$ is strongly monotone and strictly subhomogeneous. Thus, $\hat{\phi}$ is globally attractive in \bar{C}_0 for $Q_{n_1\omega}$.

Let $u(t, x, \hat{\phi})$ be the solution of (3.1) with $u(s, x) = \hat{\phi}(s, x)$ for all $(s, x) \in [-\tau, 0] \times \Omega$. Since $\hat{\phi}$ is a fixed point of $Q_{n_1\omega}$ and is globally attractive in \bar{C}_0 , $u(t, x, \hat{\phi})$ is an $n_1\omega$ -periodic solution of (3.1) which attracts all solutions of (3.1) in $\bar{C}^+ \setminus \{0\}$. That is,

$$\lim_{t \rightarrow \infty} \|u(t, \cdot, \phi) - u(t, \cdot, \hat{\phi})\|_{\beta} = 0 \quad \forall \phi \in \bar{C}_0.$$

Now we show that $u(t, x, \hat{\phi})$ is also ω -periodic. Since $Q_{n_1\omega}(\hat{\phi}) = \hat{\phi}$, we have $Q_{\omega}(Q_{n_1\omega}(\hat{\phi})) = Q_{\omega}(\hat{\phi})$, i.e., $Q_{n_1\omega}(Q_{\omega}(\hat{\phi})) = Q_{\omega}(\hat{\phi})$, which implies that $Q_{\omega}(\hat{\phi})$ is also a fixed point of $Q_{n_1\omega}$. By the fact that $\hat{\phi} \gg 0$ and the fact that Q_{ω} is monotone, it follows that $Q_{\omega}(\hat{\phi}) \gg 0$. Note that $Q_{n_1\omega}$ has a unique fixed point in $\text{int}(\bar{C}^+)$. Then $Q_{\omega}(\hat{\phi}) = \hat{\phi}$, that is, $\hat{\phi}$ is a fixed point of Q_{ω} , and hence, $u(t, x, \hat{\phi})$ is an ω -periodic solution of (3.1). Thus, $u^*(t, x) := u(t, x, \hat{\phi})$ for all $(t, x) \in [-\tau, \infty) \times \bar{\Omega}$, is the desired ω -periodic solution. \square

Acknowledgments. We are grateful to two anonymous referees for their careful reading and helpful suggestions, which led to an improvement of our original manuscript.

REFERENCES

- [1] W. G. AIELLO AND H. I. FREEDMAN, *A time-delay model of single-species growth with stage structure*, Math. Biosci., 101 (1990), pp. 139–153.
- [2] J. F. M. AL-OMARI AND S. A. GOURLEY, *A nonlocal reaction-diffusion model for a single species with stage structure and distributed maturation delay*, European J. Appl. Math., 16 (2005), pp. 37–51.
- [3] J. F. M. AL-OMARI AND S. A. GOURLEY, *Monotone wave-fronts in a structured population model with distributed maturation delay*, IMA J. Appl. Math., 70 (2005), pp. 858–879.
- [4] D. G. ARONSON AND H. F. WEINBERGER, *Nonlinear diffusion in population genetics, combustion, and nerve pulse propagation*, in Partial Differential Equations and Related Topics, J. A. Goldstein, ed., Lecture Notes in Math. 446, Springer-Verlag, Berlin, 1975, pp. 5–49.
- [5] S. A. GOURLEY AND Y. KUANG, *Wavefronts and global stability in a time-delayed population model with stage structure*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 459 (2003), pp. 1563–1579.
- [6] S. A. GOURLEY AND J. W. H. SO, *Dynamics of a food-limited population model incorporating nonlocal delays on a finite domain*, J. Math. Biol., 44 (2002), pp. 49–78.
- [7] S. A. GOURLEY AND J. WU, *Delayed non-local diffusive systems in biological invasion and disease spread*, Fields Inst. Commun., 48 (2006), pp. 137–200.
- [8] P. HESS, *Periodic-Parabolic Boundary Value Problems and Positivity*, Longman Scientific and Technical, Harlow, UK, 1991.
- [9] D. LIANG, J. WU, AND F. ZHANG, *Modelling population growth with delayed nonlocal reaction in 2-dimensions*, Math. Biosci. Eng., 2 (2005), pp. 111–132.
- [10] X. LIANG, Y. YI, AND X.-Q. ZHAO, *Spreading speeds and traveling waves for periodic evolution systems*, J. Differential Equations, 231 (2006), pp. 57–77.
- [11] X. LIANG AND X.-Q. ZHAO, *Asymptotic speeds of spread and traveling waves for monotone semiflows with applications*, Comm. Pure Appl. Math., 60 (2007), pp. 1–40; Erratum: 61 (2008), pp. 137–138.
- [12] R. H. MARTIN, *Nonlinear Operators and Differential Equations in Banach Spaces*, John Wiley and Sons, New York, 1976.

- [13] R. H. MARTIN AND H. L. SMITH, *Abstract functional-differential equations and reaction-diffusion systems*, Trans. Amer. Math. Soc., 321 (1990), pp. 1–44.
- [14] J. A. J. METZ AND O. DIEKMANN, *The dynamics of physiologically structured populations*, J. A. J. Metz and O. Diekmann, eds., Springer-Verlag, Berlin, 1986.
- [15] H. L. SMITH, *Monotone Dynamical Systems: An Introduction to the Theory of Competitive and Cooperative Systems*, American Mathematical Society, 1995.
- [16] H. L. SMITH AND H. R. THIEME, *Strongly order preserving semiflows generated by functional-differential equations*, J. Differential Equations, 93 (1991), pp. 332–363.
- [17] J. W.-H. SO, J. WU, AND X. ZOU, *A reaction-diffusion model for a single species with age structure. I. Travelling wavefronts on unbounded domains*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 457 (2001), pp. 1841–1853.
- [18] H. R. THIEME AND X.-Q. ZHAO, *Asymptotic speeds of spread and traveling waves for integral equations and delayed reaction-diffusion models*, J. Differential Equations, 195 (2003), pp. 430–470.
- [19] A. I. VOLPERT, V. A. VOLPERT, AND V. A. VOLPERT, *Traveling Wave Solutions of Parabolic Systems*, AMS, Providence, RI, 1994.
- [20] D. XU AND X.-Q. ZHAO, *A nonlocal reaction-diffusion population model with stage structure*, Can. Appl. Math. Q., 11 (2003), pp. 303–319.
- [21] D. XU AND X.-Q. ZHAO, *Dynamics in a periodic competitive model with stage structure*, J. Math. Anal. Appl., 311 (2005), pp. 417–438.
- [22] X.-Q. ZHAO, *Global attractivity and stability in some monotone discrete dynamical system*, Bull. Austral. Math. Soc., 53 (1996), pp. 305–324.
- [23] X.-Q. ZHAO, *Dynamical Systems in Population Biology*, Springer-Verlag, New York, 2003.

EXISTENCE OF SOLUTIONS FOR A MODEL DESCRIBING THE DYNAMICS OF JUNCTIONS BETWEEN DISLOCATIONS*

NICOLAS FORCADEL[†] AND RÉGIS MONNEAU[‡]

Abstract. We study a dynamical version of a multiphase field model of Koslowski and Ortiz for planar dislocation networks. We consider a two-dimensional vector field which describes phase transitions between constant phases. Each phase transition corresponds to a dislocation line, and the vectorial field description allows the formation of junctions between dislocations. This vector field is assumed to satisfy a nonlocal vectorial Hamilton–Jacobi equation with nonzero viscosity. For this model, we prove the existence for all time of a weak solution.

Key words. dislocation dynamics, nonlocal equations, junctions, parabolic system of equations

AMS subject classifications. 35K15, 74K30

DOI. 10.1137/070710925

1. Introduction.

1.1. Physical motivation. Dislocations are line defects in crystal, and their motion is at the origin of plastic properties of metals (see, for instance, the book of Hirth and Lothe [13]). The typical length of these dislocation lines is of the order of the micrometer, and their typical thickness is of the order of the nanometer. Dislocation lines exist in almost all metals. At low temperatures, these lines are contained in their crystallographic plane, called the slip plane. When an elastic stress field is applied, these lines can move in their slip plane. The normal velocity of these lines is then proportional to the effective stress in the material.

Another important property of dislocations is that each dislocation line is characterized by a vector quantity, called the Burgers vector (see [13] for more details). To explain this briefly, let us say that this Burgers vector reflects the microscopic nature of the dislocation defect in a crystal lattice. To fix the ideas, let us consider three unit vectors b^1 , b^2 , and b^3 in the plane such that

$$(1.1) \quad b^1 + b^2 + b^3 = 0.$$

Then each dislocation line can be of three different natures: with Burgers vector b^1 , with Burgers vector b^2 , or with Burgers vector b^3 . The Burgers vector is an invariant of the dislocation line and then does not change during the evolution of the dislocation.

A consequence of the existence of this Burgers vector is the possibility for the dislocation lines to create some triple junctions, with three dislocation lines, respectively, of Burgers vector b^1 , b^2 , and b^3 , satisfying (1.1). Indeed, in real crystals, we can observe such junctions. We can even observe networks of dislocations related by junctions. Those self-organized structures are called Frank networks (see Figure 1.1). See, also, for instance, p. 190 in Hull and Bacon [14] for such networks in body-

*Received by the editors December 14, 2007; accepted for publication (in revised form) October 21, 2008; published electronically March 20, 2009. This work was supported by the ACI JC 1025.

<http://www.siam.org/journals/sima/40-6/71092.html>

[†]Projet Commands, CMAP-INRIA Futurs, Ecole Polytechnique, 91128 Palaiseau and ENSTA, UMA, 32 Bd Victor, 75739 Paris Cedex 15, France (forcadel@casemade.dauphine.fr).

[‡]CERMICS, Paris Est-ENPC, 6 & 8 avenue Blaise Pascal, Cit Descartes, Champs sur Marne, 77455 Marne la Valle Cedex 2, France (monneau@cermics.enpc.fr).

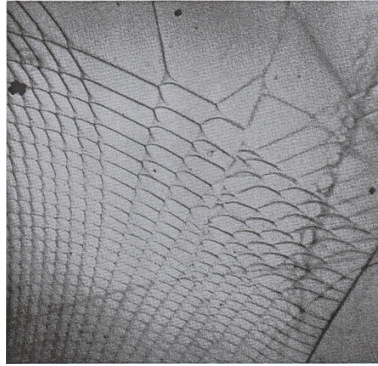


FIG. 1.1. Frank networks observed with electron microscopy.

centered cubic (BCC) iron or p. 188 for hexagonal networks in face-centered cubic (FCC) crystals.

In the present paper, we consider a special case of a set of dislocations contained in a single slip plane, where the dislocations can move. We are interested, in particular, in the motion of the junctions between dislocations, which has not been studied a lot, both from the modeling point of view and from the mathematical analysis point of view (see, for instance, the work of Rodney, Le Bouar, and Finel [20]). Let us mention, for the stationary case, the work of Cacace and Garroni [9]. The goal of the present paper is to propose and to study a model for the dynamics of junctions of dislocations.

Nevertheless, the question of junctions has several other physical applications, and there is a lot of literature on this subject. Let us mention, for instance, the problem of crystal growth or grain growth (see Taylor [23, 24] and Bronsard and Reitich [8]). We also refer to Bonnet [7] for problems concerning the minimization of the Mumford-Shah functional.

1.2. A phase field model for the dynamics of junctions. In a phase field model, the dislocation can be represented as the phase transition of a phase parameter $\rho(x) = \rho_1(x)e^1 + \rho_2(x)e^2 \in \mathbb{R}^2$ defined for $x = x_1e^1 + x_2e^2$ in the plane \mathbb{R}^2 with (e^1, e^2) an orthonormal basis. As in the work of Koslowski and Ortiz [15], we consider only the case of a periodic distribution of dislocations, reducing the problem on the torus, for $x \in \mathbb{T}^2 = \mathbb{R}^2/\mathbb{Z}^2$. Then the energy of the dislocations, in the presence of a constant exterior applied resolved stress $\sigma^0 \in \mathbb{R}^2$, is then given (see [15]) by

$$(1.2) \quad \mathcal{E}(\rho) = \int_{\mathbb{T}^2} -\frac{1}{2} (C^0 \star \rho) \cdot \rho - \sigma^0 \cdot \rho + W(\rho).$$

In (1.2) and throughout the paper, we denote by $A \cdot B$ the scalar product between two vectors $A, B \in \mathbb{R}^2$. The precise meaning of the whole expression (1.2) will be explained later.

For any phase transition between two states $\rho = A$ and $\rho = B$, the difference $B - A$ needs physically to be the Burgers vector of the dislocation, i.e., a vector of the lattice $\Lambda = \mathbb{Z}a^1 + \mathbb{Z}a^2$ of the crystal we are considering, with a general given basis (a^1, a^2) . This information is encoded in the potential $W : \mathbb{R}^2 \rightarrow \mathbb{R}_+$, which is assumed to be minimal on Λ and to have the periodicity of the lattice Λ :

$$(1.3) \quad W(\rho + a) = W(\rho) \quad \text{for any } a \in \Lambda.$$

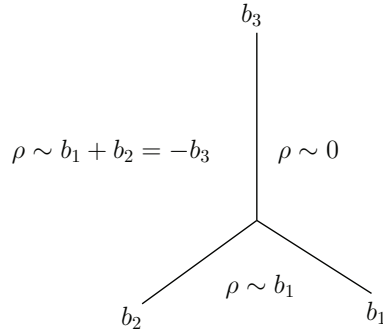


FIG. 1.2. The junction of three dislocations as phase transitions of ρ .

In other words, because ρ can be defined up to addition of any vector of the lattice Λ , only the transitions of ρ between two phases, i.e., two constant vectors of the lattice Λ , are physically meaningful. In particular, in this model, junctions of three dislocations of Burgers vectors $b^1, b^2, b^3 \in \Lambda$ with $b^1 + b^2 + b^3 = 0$ are expected, like, for instance, as the phase transitions between the states $0, b^1, -b^3$ (see Figure 1.2).

In the expression giving the energy (1.2), the kernel $C^0(x)$ is a 2×2 symmetric matrix which takes into account the long range elastic interactions between dislocations and

$$(C^0 \star \rho)_i = \sum_{j=1,2} C_{ij}^0 \star \rho_j \quad \text{for } i = 1, 2,$$

where \star denotes the usual convolution on the torus \mathbb{T}^2 . For instance, in the particular case of isotropic elasticity, for $k = (k_1, k_2) \in \mathbb{Z}^2$, the k -Fourier coefficient of the matrix C^0 is given (see [15] and also a limit case of the Peierls–Nabarro model in Alvarez et al. [3]) by

$$(1.4) \quad \widehat{C}^0(k) = -\frac{\mu}{2|k|} \left(\frac{1}{(1-\nu)} k \otimes k + k^\perp \otimes k^\perp \right),$$

where $k^\perp = (-k_2, k_1)$ is the vector obtained by a rotation of k of angle $\pi/2$. Here, $\mu > 0$ is a Lamé coefficient and $\nu \in (-1, 1/2)$ is the Poisson ratio of the material.

The fact that the matrix $-\widehat{C}^0(k)$ is nonnegative is related to the fact that the elastic energy is nonnegative. This insures that the elastic part of the energy $\int_{\mathbb{T}^2} -\frac{1}{2}(C^0 \star \rho) \cdot \rho$ created by the dislocations is nonnegative. The fact that the matrix $\widehat{C}^0(k)$ is not proportional to the identity reflects the fact that the existence of a Burgers vector associated to the dislocation line (that can be seen in the phase transition of ρ) creates some anisotropic elastic stress, even if we work in the framework of isotropic elasticity. The stress created by a dislocation line has somehow a preferred direction which is given by its Burgers vector.

When the material is submitted to an exterior shear stress, it makes the dislocations move. The dynamics of a given dislocation line is physically given by its normal velocity, which is called the resolved Peach–Koehler force. This force is the sum of the resolved exterior shear stress and the stress created by all of the dislocation lines, including the line itself. The total resolved stress $\sigma[\rho]$ is then formally given by the opposite of the gradient of the energy, i.e., $-\mathcal{E}'(\rho)$, and can be expressed as the

following nonlocal quantity:

$$(1.5) \quad \sigma[\rho] = \sigma^0 + C^0 \star \rho - W'_\rho(\rho).$$

Let us remark that one mathematical difficulty in the computation of this stress comes from the term $C^0 \star \rho$, which, from the point of view of the regularity, behaves like $\nabla \rho$, because of the linear growth of $\widehat{C}^0(k)$ in k . We have now to write an evolution equation for ρ , keeping in mind that the normal velocity has to be proportional to the stress, which means roughly speaking that

$$\rho = \bar{\rho} b^1 \quad \text{with} \quad \bar{\rho}(t, x) \in \mathbb{R}$$

has to satisfy

$$(1.6) \quad \frac{\bar{\rho}_t}{|\nabla \bar{\rho}|} = \bar{\sigma}[\bar{\rho}] \quad \text{with} \quad \bar{\sigma}[\bar{\rho}] = b^1 \cdot \sigma[\bar{\rho} b^1].$$

Indeed, (1.6) is an equation very difficult to study mathematically, because $\bar{\sigma}[\bar{\rho}]$ behaves like $\nabla \bar{\rho}$. Up to our knowledge no results exist for (1.6), even in a framework of viscosity solutions for nonlocal equations. Indeed, one of the mathematical difficulties is due to the fact that, even putting the matrix C^0 to zero, (1.6) would become an equation like the Burgers equation, because of the presence of the derivative of the potential W in $\bar{\sigma}[\bar{\rho}]$, and could then create shocks in finite time. The presence of C^0 in the stress, even if this corresponds to a kind of degenerate diffusion (when $\nabla \bar{\rho} = 0$) in the equation, does not help sufficiently to get a good framework for a suitable notion of the solution to (1.6).

Because of these high mathematical difficulties and because our goal is really to deal with the dynamics of junctions of dislocations, we simplify mathematically the equation adding an artificial small viscosity $\varepsilon \in (0, 1)$ and considering, instead of (1.6), the following equation:

$$(1.7) \quad \bar{\rho}_t = \bar{\sigma}[\bar{\rho}] |\nabla \bar{\rho}| + \varepsilon \Delta \bar{\rho}.$$

If the solution for the limit case $\varepsilon = 0$ is smooth enough, it seems reasonable to think that the solution to (1.7) is a good approximation. Moreover, any reasonable solution of the limit equation for $\varepsilon = 0$ should probably be seen as a limit of the solution to the ε -equation when ε goes to zero. Even numerically, if we would like to compute the solution in the limit case $\varepsilon = 0$, every classical numerical method introduces a numerical diffusion, which can be more or less interpreted as the additional viscosity $\varepsilon > 0$ in (1.7). For all of these reasons, (1.7) seems a good candidate for the evolution equation. As we will see in this paper, it is then mathematically possible to deal with this equation and to prove the existence of a global solution to (1.7). Nevertheless, even the study of this equation with ε -viscosity is not so simple, because, for instance, the uniqueness of the solution to (1.7) remains an open (and probably very difficult) problem in general.

As explained above, our goal is really to generalize (1.7) to take into account the evolution of junctions. To this end, we assume that the phase parameter $\rho(t, x) \in \mathbb{R}^2$ satisfies the following evolution equation:

$$(1.8) \quad \begin{cases} (\rho_k)_t = |\nabla \rho|^{-1} \sum_{i=1,2} \sum_{j=1,2} (\sigma[\rho])_i \nabla_j \rho_i \nabla_j \rho_k + \varepsilon \Delta \rho_k & \text{for } k = 1, 2 \\ \rho(0, x) = \rho^0(x) & \text{in } (0, T) \times \mathbb{T}^2, \\ & \text{on } \mathbb{T}^2, \end{cases}$$

where σ is given in (1.5), $\rho_t = \frac{\partial \rho}{\partial t}$ and $\nabla_j \rho_i = \frac{\partial \rho_i}{\partial x_j}$ for $i, j = 1, 2$, and

$$|\nabla \rho|^2 = \sum_{i=1,2} \sum_{j=1,2} |\nabla_j \rho_i|^2.$$

Here, we can check easily that (1.8) reduces to (1.7) when $\rho = \bar{\rho} b^1$ for any Burgers vector b^1 satisfying $|b^1| = 1$. For this reason, (1.8) seems to be a nice model to describe the general dynamics of junctions of dislocations.

Finally, let us mention that our model (1.8) has some similarities with the model of Allen and Cahn [2] on the motion of curved boundaries in which they consider gradient flow associated with a free-energy functional. This led to the study of scalar Ginzburg–Landau-type diffusion equations like

$$u_t = \Delta u - W'(u).$$

1.3. Main result. We need to introduce some assumptions that we will comment on below.

We make the following assumption on the kernel $C^0 : \mathbb{T}^2 \rightarrow \mathbb{R}_{sym}^{2 \times 2}$.

(A) We assume that there exists a constant $m > 0$ such that, for any $k \in \mathbb{Z}^2$, the Fourier coefficients of the kernel $\widehat{C}^0(k) = \int_{\mathbb{T}^2} dx e^{-2i\pi k \cdot x} C^0(x)$ satisfy $\widehat{C}^0(k) = M(k)$, where for any $\xi = (\xi_1, \xi_2) \in \mathbb{R}^2$ and any $p = (p_1, p_2) \in \mathbb{R}^2$

$$(1.9) \quad \begin{cases} M \in C^\infty(\mathbb{R}^2 \setminus \{0\}; \mathbb{R}_{sym}^{2 \times 2}), & M(-\xi) = M(\xi), & M(\xi) = |\xi| M\left(\frac{\xi}{|\xi|}\right), \\ \frac{|\xi| |p|^2}{m} \geq - \sum_{i=1,2} \sum_{j=1,2} p_i \cdot M_{ij}(\xi) \cdot p_j \geq m |\xi| |p|^2 & \text{with } |p|^2 = \sum_{i=1,2} (p_i)^2. \end{cases}$$

We also make the following assumption on the potential $W : \mathbb{T}^2 \rightarrow \mathbb{R}_+$.

(B) We assume that $W \in C^2$ and W satisfies (1.3).

Condition (1.3) of assumption (B) is natural for the potential W , as explained in subsection 1.2. In assumption (B), we assume moreover that W is smooth enough (here C^2).

Assumption (A) requires more comments. First of all, remark that assumption (A) is satisfied by the matrix C^0 given by Koslowski and Ortiz [15] for isotropic elasticity (see (1.4)). Moreover, for general elasticity, the matrix $-\widehat{C}^0(k)$ has to be symmetric and nonnegative, because it corresponds physically to the nonnegative elastic energy of the dislocations. The last point concerns the 1-homogeneity in k of the matrix $\widehat{C}^0(k)$. Indeed, this is a general fact, when we compute $\widehat{C}^0(k)$ for general linear elasticity (see, for instance, the computation in the case of cubic elasticity done in Alvarez et al. [3]). Indeed, we start with a three-dimensional (3D) elastic energy, which behaves like the square of the H^1 -norm of the displacement. When we consider dislocations in a single plane, as in the present model, this naturally reduces the problem from a 3D to a 2D problem and creates an energy as the square of the $H^{\frac{1}{2}}$ -norm of the phase parameter, i.e., an energy like $\sum_{k \in \mathbb{Z}^2} |k| |\widehat{\rho}|^2$, which is directly reflected into the 1-homogeneity in k of the matrix $\widehat{C}^0(k)$.

Then we have the following result for the model of dynamics of junctions between dislocations.

THEOREM 1.1 (existence of a solution). *Under assumptions (A) and (B), if $\rho^0 \in (H^1(\mathbb{T}^2))^2$, then, for any constant applied stress $\sigma^0 \in \mathbb{R}^2$ and for any time $T > 0$, there exists a solution ρ of (1.8) with $\rho \in C^0([0, T]; (L^{\frac{4}{3}}(\mathbb{T}^2))^2)$.*

As mentioned above, the uniqueness of the solution is not known. Let us also mention that (1.8) is a nonlocal system of scalar equations and can be sketched as the following equation:

$$(1.10) \quad v_t = |\nabla v|^2 + \Delta v.$$

Indeed, this comes from our assumption (A) that the convolution with the kernel behaves like a first order operator. A lot of work has been done on equations (or systems) like (1.10). Let us mention, for instance, the works of Boccardo, Murat, and Puel [4, 5, 6] in which they study general equations including (1.10) and prove the existence result.

Equation (1.10) is also similar to the Navier–Stokes equations written for the potential A such that the velocity of the fluid is given by $u = \text{curl } A$ (see, for instance, Leray [17]).

1.4. Organization of the paper. In section 2, we study an approximate problem of (1.8) where the right-hand side is approached by some term at most linear in the solution. The main result is proved in section 3. In a first subsection, we give some a priori estimates for the solution of the approximate problem, and then, in a second subsection, we pass to the limit in the approximate problem.

1.5. Notation. In what follows, we will denote by C a generic constant, which will then satisfy $C + C = C$, $C \cdot C = C$, and so on. We also use the following set:

$$W^{2,1;p}(Q_T) = \left\{ u \in L^p(Q_T); u_t \in L^p(Q_T) \text{ and } \frac{\partial^2 u}{\partial x_i \partial x_j} \in L^p(Q_T) \text{ for } i, j = 1, 2 \right\},$$

where $Q_T = (0, T) \times \mathbb{T}^2$.

2. An approximate problem. We first start to approximate the right-hand side of (1.8) by some term at most linear in the solution. To this end, we introduce a function h^n defined by

$$h^n(r) = h^0(r - n)$$

with

$$h^0(r) = \begin{cases} 1 & \text{if } r \leq 0, \\ 1 - r & \text{if } 0 \leq r \leq 1, \\ 0 & \text{if } r \geq 1. \end{cases}$$

We then look at the following approximate problem:

$$(2.1) \quad \begin{cases} \rho_t - \epsilon \Delta \rho = f^n[\rho] & \text{on } Q_T := (0, T) \times \mathbb{T}^2, \\ \rho(0, \cdot) = \rho^0 & \text{on } \mathbb{T}^2, \end{cases}$$

where

$$f^n[\rho] = h^n(|\nabla \rho|) |\nabla \rho|^{-1} (\nabla \rho)^T \cdot \nabla \rho \cdot \sigma[\rho]$$

and $\sigma[\rho]$ is given in (1.5) and is at most linear in ρ .

The natural idea to find a solution to (2.1) is to define the map Φ which associates to any function u , the solution $\rho = \Phi(u)$ of

$$(2.2) \quad \begin{cases} \rho_t - \epsilon \Delta \rho = f^n[u] & \text{on } Q_T := (0, T) \times \mathbb{T}^2, \\ \rho(0, \cdot) = \rho^0 & \text{on } \mathbb{T}^2, \end{cases}$$

and to prove that Φ has a fixed point in a suitable space. This way, we will prove the following result.

THEOREM 2.1 (existence of a solution for the approximate problem). *If $\rho^0 \in (H^1(\mathbb{T}^2))^2$, then, for any $n \geq 1$ and any $T > 0$, there exists a solution ρ^n of (2.1) with $\rho^n \in L^2((0, T); (H^2(\mathbb{T}^2))^2) \cap C^0([0, T]; (L^2(\mathbb{T}^2))^2)$.*

In this section, we will give the proof of this theorem. In a first subsection, we will collect some preliminary results, and, in a second subsection, we will prove that Φ has a fixed point.

2.1. Preliminary results. The following lemma will be important.

LEMMA 2.2 (estimate on $C^0 \star \rho$). *For any $p \in (1, +\infty)$, there exists a constant C (depending on p and on the constant m defined in assumption (A)) such that, for any $\rho \in (W^{1,p}(\mathbb{T}^2))^2$, we have*

$$(2.3) \quad \|C^0 \star \rho\|_{(L^p(\mathbb{T}^2))^2} \leq C \|\nabla \rho\|_{(L^p(\mathbb{T}^2))^{2 \times 2}}.$$

Partial proof of Lemma 2.2. Let us make the proof for $p = 2$. We have, with $\sigma = \sigma[\rho]$,

$$\begin{aligned} |C^0 \star \rho|_{(L^2(\mathbb{T}^2))^2}^2 &= \sum_{k \in \mathbb{Z}^2} \left| (\widehat{C^0 \star \rho})(k) \right|^2 \\ &= \sum_{k \in \mathbb{Z}^2} \left| \widehat{C}^0(k) \cdot \widehat{\rho}(k) \right|^2 \\ &\leq \frac{1}{m^2} \sum_{k \in \mathbb{Z}^2} |k|^2 |\widehat{\rho}(k)|^2 \\ &\leq \frac{1}{(2\pi m)^2} \sum_{k \in \mathbb{Z}^2} \left| \widehat{\nabla \rho}(k) \right|^2 \\ &= \frac{1}{(2\pi m)^2} |\nabla \rho|_{(L^2(\mathbb{T}^2))^{2 \times 2}}^2 \end{aligned}$$

which provides the result in the case $p = 2$.

The proof for the general case $p \in (1, +\infty)$ is given in Appendix A. □

An immediate corollary is the following estimate on the stress.

COROLLARY 2.3 (estimate on $\sigma[\rho]$). *For any $p \in (1, +\infty)$, there exists a constant C (depending on p , on the constant σ^0 , on the potential W , and on the constant m defined in assumption (A)) such that, for any $\rho \in (W^{1,p}(\mathbb{T}^2))^2$, we have*

$$(2.4) \quad \|\sigma[\rho]\|_{(L^p(\mathbb{T}^2))^2} \leq C (1 + \|\nabla \rho\|_{(L^p(\mathbb{T}^2))^{2 \times 2}}).$$

We will also need the following result.

LEMMA 2.4 (estimate on $f^n[u]$). *If $u \in (H^1(\mathbb{T}^2))^2$, then $f^n[u] \in (L^2(\mathbb{T}^2))^2$ with the following estimate:*

$$\|f^n[u]\|_{(L^2(\mathbb{T}^2))^2} \leq C(n + 1) (1 + \|\nabla u\|_{L^2(\mathbb{T}^2)^{2 \times 2}}),$$

where the constant C depends on σ^0 , on the potential W , and on the constant m defined in assumption (A).

Proof of Lemma 2.4. Since $\text{supp}(h^n) \subset [0, n + 1]$, the following holds:

$$(2.5) \quad |f^n[u]| \leq (n + 1)|\sigma[u]|,$$

where we have used the fact that $|B^T \cdot B \cdot p| \leq |B|^2|p|$ for $B \in \mathbb{R}^{2 \times 2}$ and $p \in \mathbb{R}^2$. Then

$$\begin{aligned} \|f^n[u]\|_{(L^2(\mathbb{T}^2))^2} &\leq (n + 1)\|\sigma[u]\|_{(L^2(\mathbb{T}^2))^2} \\ &\leq C(n + 1) (1 + \|\nabla u\|_{L^2(\mathbb{T}^2)^{2 \times 2}}), \end{aligned}$$

where we have used Corollary 2.3. \square

We now recall some classical results. We start with the following parabolic estimates for the following equation:

$$(2.6) \quad \begin{cases} g_t - \epsilon \Delta g = f & \text{on } Q_T := (0, T) \times \mathbb{T}^2, \\ g(0, \cdot) = g^0 & \text{on } \mathbb{T}^2. \end{cases}$$

PROPOSITION 2.5 (parabolic estimates for the heat equation). *Let $g^0 \in H^1(\mathbb{T}^2)$ and $f \in L^2(Q_T)$. Then there exists a unique solution g to (2.6) with*

$$(2.7) \quad g \in L^2((0, T); H^2(\mathbb{T}^2)) \cap L^\infty((0, T); H^1(\mathbb{T}^2)), \quad g_t \in L^2(Q_T).$$

We have the following estimate:

$$(2.8) \quad \begin{aligned} &\sup_{0 \leq t \leq T} \|g(t)\|_{H^1(\mathbb{T}^2)} + \|g\|_{L^2((0, T); H^2(\mathbb{T}^2))} + \|g_t\|_{L^2((0, T); L^2(\mathbb{T}^2))} \\ &\leq C_T (\|f\|_{L^2(Q_T)} + \|g^0\|_{H^1(\mathbb{T}^2)}), \end{aligned}$$

where the constant C_T depends only on T and ϵ .

Moreover, we have

$$(2.9) \quad \sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} g^2(t) + 4\epsilon \int_0^T \int_{\mathbb{T}^2} |\nabla g|^2 \leq 4 \int_{\mathbb{T}^2} (g^0)^2 + 16T \int_0^T \int_{\mathbb{T}^2} f^2.$$

Proof of Proposition 2.5. For the proof of (2.7)–(2.8), we refer to Evans [11, Theorem 5, p. 360].

To prove (2.9), we simply multiply (2.6) by g and integrate over \mathbb{T}^2 and $(0, t)$, taking the supremum for $0 \leq t \leq T$. We get

$$(2.10) \quad \sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} \frac{g^2(t)}{2} \leq \sup_{0 \leq t \leq T} \left(\int_{\mathbb{T}^2} \frac{g^2(t)}{2} + \epsilon \int_0^t \int_{\mathbb{T}^2} |\nabla g|^2 \right) \leq \int_{\mathbb{T}^2} \frac{(g^0)^2}{2} + \int_0^T \int_{\mathbb{T}^2} |g f|$$

and

$$(2.11) \quad \epsilon \int_0^T \int_{\mathbb{T}^2} |\nabla g|^2 \leq \int_{\mathbb{T}^2} \frac{g^2(T)}{2} + \epsilon \int_0^T \int_{\mathbb{T}^2} |\nabla g|^2 \leq \int_{\mathbb{T}^2} \frac{(g^0)^2}{2} + \int_0^T \int_{\mathbb{T}^2} |g f|.$$

Summing (2.10) and (2.11), we finally get

$$\sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} \frac{g^2(t)}{2} + \epsilon \int_0^T \int_{\mathbb{T}^2} |\nabla g|^2 \leq \int_{\mathbb{T}^2} (g^0)^2 + 2 \int_0^T \int_{\mathbb{T}^2} |g f|.$$

We now use the fact that

$$\begin{aligned} 2 \int_0^T \int_{\mathbb{T}^2} |g f| &\leq 2 \left(\int_0^T \int_{\mathbb{T}^2} g^2 \right)^{\frac{1}{2}} \cdot \left(\int_0^T \int_{\mathbb{T}^2} f^2 \right)^{\frac{1}{2}} \\ &\leq 2 \left(T \sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} g^2(t) \right)^{\frac{1}{2}} \cdot \left(\int_0^T \int_{\mathbb{T}^2} f^2 \right)^{\frac{1}{2}} \\ &\leq \frac{1}{4} \sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} g^2(t) + 4T \int_0^T \int_{\mathbb{T}^2} f^2, \end{aligned}$$

which implies the result. \square

We also recall the following theorem.

THEOREM 2.6 (Schaefer’s fixed point theorem). *Let X be a real Banach space. Suppose that*

$$\Phi : X \rightarrow X$$

is a continuous and compact mapping. Assume further that the set

$$\{u \in X, \quad u = \lambda \Phi(u) \quad \text{for some } \lambda \in [0, 1]\}$$

is bounded. Then Φ has a fixed point.

For the proof of this theorem, we refer to Evans [11, Theorem 4, p. 504].

Finally, we will need some compactness argument and a weak continuity property contained in the following two classical results.

PROPOSITION 2.7 (compactness). *We recall that*

$$W^{2,1;2}(Q_T) = \{g \in L^2((0, T); H^2(\mathbb{T}^2)), \quad g_t \in L^2(Q_T)\}.$$

Then the injection

$$W^{2,1;2}(Q_T) \longrightarrow L^2((0, T); H^1(\mathbb{T}^2)) \quad \text{is compact.}$$

For the proof of this result, we refer to Lions [18, Theorem 5.1, p. 58].

PROPOSITION 2.8 (continuity). *With the notation of Proposition 2.7, let us consider a sequence $(g^m)_m$ such that*

$$g^m \rightharpoonup g \quad \text{weakly in } W^{2,1;2}(Q_T).$$

We assume, also, that $g^m_{|t=0} = \rho^0$. Then

$$g_{|t=0} = \rho^0.$$

This result is classical, but for the reader’s convenience we give the proof in Appendix A.

2.2. Proof of Theorem 2.1. We are now ready to make the proof of Theorem 2.1. To this end, for any $T > 0$, we set

$$X_T = L^2((0, T); H^1(\mathbb{T}^2)).$$

In all that follows, the index n is assumed fixed. We first remark that if $u \in X_T^2$, then $f^n[u] \in (L^2(Q_T))^2$, and then we can consider the solution ρ of

$$(2.12) \quad \begin{cases} \rho_t - \epsilon \Delta \rho = f^n[u] & \text{on } Q_T := (0, T) \times \mathbb{T}^2, \\ \rho(0, \cdot) = \rho^0 & \text{on } \mathbb{T}^2, \end{cases}$$

which satisfies $\rho \in X_T^2$ because of the parabolic estimates of Proposition 2.5. Then we set $\Phi(u) = \rho$ and see that Φ maps X_T^2 into X_T^2 . We will prove that Φ admits a fixed point using Schaefer’s fixed point theorem. We do the proof in four steps.

Step 1 (weak continuity of Φ). Let us consider sequences $(u^m)_m, (\rho^m)_m$ such that

$$\begin{cases} u^m \in X_T^2, & \rho^m = \Phi(u^m), \\ u^m \longrightarrow u & \text{in } X_T^2. \end{cases}$$

From Lemma 2.4, we deduce that

$$(2.13) \quad \|f^n[u^m]\|_{(L^2(Q_T))^2} \leq C(n+1) \left(1 + \|u^m\|_{X_T^2}\right).$$

From the parabolic estimates (Proposition 2.5), we deduce that ρ^m is bounded in $(W^{2,1;2}(Q_T))^2$; i.e., there exists a constant $C > 0$ such that

$$(2.14) \quad \|\rho^m\|_{(W^{2,1;2}(Q_T))^2} \leq C.$$

Therefore, up to a subsequence, we have

$$\rho^m \rightharpoonup \rho \quad \text{in } (W^{2,1;2}(Q_T))^2,$$

and, from Proposition 2.8, we deduce that

$$\rho|_{t=0} = \rho^0 \quad \text{on } \mathbb{T}^2.$$

We now claim that

$$(2.15) \quad f^n[u^m] \longrightarrow f^n[u] \quad \text{in } L^1(Q_T).$$

Indeed, we can write

$$f^n[u] = g^n(\nabla u) \cdot \sigma[u] \quad \text{with } g^n(\nabla u) := h^n(|\nabla u|) |\nabla u|^{-1} (\nabla u)^T \cdot \nabla u.$$

From Lemma 2.2, for $p = 2$, and the continuity of W' , we already deduce that

$$(2.16) \quad \sigma[u^m] \longrightarrow \sigma[u] \quad \text{in } L^2(Q_T).$$

From the convergence of u^m to u in X_T^2 , we deduce that, up to a subsequence, we have $\nabla u^m \longrightarrow \nabla u$ a.e. in Q_T . Now, from the fact that g^n is continuous and bounded, we deduce, in particular, that

$$(2.17) \quad g^n(\nabla u^m) \longrightarrow g^n(\nabla u) \quad \text{in } L^2(Q_T).$$

Then the convergence (2.15) follows from (2.16) and (2.17).

Therefore, we conclude that ρ solves (2.12). Finally, by uniqueness of the solutions of (2.12), we deduce that the limit ρ does not depend on the choice of the subsequence and then that the full sequence converges:

$$\rho^m \rightharpoonup \rho \quad \text{weakly in } (W^{2,1;2}(Q_T))^2, \quad \text{with } \rho = \Phi(u).$$

Step 2 (compactness of Φ). The compactness (and the usual strong continuity) of Φ follows from the compactness of the injection $(W^{2,1;2}(Q_T))^2 \rightarrow X_T^2$ (see Proposition 2.7).

Step 3 (a priori bounds on the solutions of $u = \lambda\Phi(u)$ for T small). Let us consider a solution u of

$$(2.18) \quad u = \lambda\Phi(u) \quad \text{for some } \lambda \in [0, 1].$$

Then, from the parabolic estimates (2.9), we have

$$\begin{aligned} & \sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} |u(t)|^2 + 4\varepsilon \int_0^T \int_{\mathbb{T}^2} |\nabla u|^2 \\ & \leq 4 \int_{\mathbb{T}^2} |\rho^0|^2 + 16T \int_0^T \int_{\mathbb{T}^2} |\lambda f^n[u]|^2 \\ & \leq 4 \int_{\mathbb{T}^2} |\rho^0|^2 + 32TC^2(n+1)^2 \left(T + \int_0^T \int_{\mathbb{T}^2} |\nabla u|^2 \right), \end{aligned}$$

where, in the third line, we have used Lemma 2.4 and the fact that $|\lambda| \leq 1$. Therefore, for

$$(2.19) \quad T \leq T^* := (16C^2(n+1)^2)^{-1} \varepsilon$$

we have

$$\sup_{0 \leq t \leq T} \int_{\mathbb{T}^2} |u(t)|^2 + 2\varepsilon \int_0^T \int_{\mathbb{T}^2} |\nabla u|^2 \leq 4 \int_{\mathbb{T}^2} |\rho^0|^2 + 2\varepsilon T,$$

which proves that there exists a constant $C > 0$ such that any solution of (2.18) satisfies

$$\|u\|_{X_T^2} \leq C.$$

We can then apply Schaefer’s fixed point theorem (Theorem 2.6) to deduce that Φ has a fixed point on X_T^2 , and, therefore, there is a solution ρ of (2.1) on the time interval $(0, T)$ if T satisfies (2.19), i.e., if T is small enough independently on the initial data ρ^0 .

Step 4 (solution for any time). Let us call $\rho(\rho^0, t)$ the function $\rho(t, \cdot)$ obtained at Step 3 as a solution of (2.1) on the time interval $[0, T^*)$ with initial data ρ^0 . From the parabolic estimates (Proposition 2.5), we also know that $\rho(t, \cdot) \in (H^1(\mathbb{T}^2))^2$ for any $t \in [0, T^*)$. Then we can define with $\tau = T^*/2$

$$u(0) = \rho^0 \quad \text{and} \quad u(t) = \rho(u(k\tau), t) \quad \text{if } k\tau \leq t < (k+1)\tau \quad \text{with } k \in \mathbb{N}.$$

Using the fact that $u_t \in L^2_{loc}((0, +\infty); (L^2(\mathbb{T}^2))^2)$ and the fact that the problem is invariant by translation in time, we can easily check that u solves (2.1) for any $T > 0$ and provides the desired solution $\rho^n = u$ of Theorem 2.1.

This ends the proof of Theorem 2.1. \square

3. A priori estimates and proof of Theorem 1.1.

3.1. A priori estimates. We have the following a priori estimates.

LEMMA 3.1 (a priori estimates). *There exists a constant $C > 0$ such that, for all $T > 0, n \geq 1$, and $0 < \varepsilon < 1$, any solution ρ^n of (2.1) given by Theorem 2.1 satisfies*

$$(3.1) \quad \|\rho^n\|_{L^\infty((0,T);(H^{\frac{1}{2}}(\mathbb{T}^2))^2)}^2 \leq C e^{\frac{CT}{\varepsilon}},$$

$$(3.2) \quad \|\rho^n\|_{L^2((0,T);(H^{\frac{3}{2}}(\mathbb{T}))^2)}^2 \leq \frac{C}{\varepsilon} e^{\frac{CT}{\varepsilon}},$$

and

$$(3.3) \quad \left\| h^n (|\nabla \rho^n|) |\nabla \rho^n|^{-\frac{1}{2}} \nabla \rho^n \cdot \sigma[\rho^n] \right\|_{(L^2(Q_T))^2}^2 \leq C e^{\frac{CT}{\varepsilon}}.$$

Proof of Lemma 3.1.

Step 1 (preliminaries on the energy). We first recall the expression of the energy for a general \mathbb{Z}^2 -periodic smooth function $\rho(x) = (\rho_1(x), \rho_2(x))$:

$$\mathcal{E}(\rho) = \int_{\mathbb{T}^2} -\frac{1}{2} (C^0 \star \rho) \cdot \rho - \sigma^0 \cdot \rho + W(\rho).$$

For future use, we start to evaluate from below the first term in the energy, using Fourier series

$$\begin{aligned} \int_{\mathbb{T}^2} - (C^0 \star \rho) \cdot \rho &= \sum_{k \in \mathbb{Z}^2} -(\widehat{C^0 \star \rho})(k) \cdot \widehat{\rho}^*(k) \\ &= \sum_{k \in \mathbb{Z}^2} -(\widehat{C^0}(k) \cdot \widehat{\rho}(k)) \cdot \widehat{\rho}^*(k) \\ &\geq m \sum_{k \in \mathbb{Z}^2} |k| |\widehat{\rho}(k)|^2, \end{aligned}$$

where we have used, in the first line, the fact that ρ and C^0 are real, and, in the last line, we have used assumption (A). Then we define

$$\|\rho\|_{(H^{\frac{1}{2}}(\mathbb{T}^2))^2}^2 := \sum_{k \in \mathbb{Z}^2} |k| |\widehat{\rho}(k)|^2.$$

Similarly, we compute

$$\begin{aligned} &\int_{\mathbb{T}^2} - (C^0 \star (\nabla \rho)^T) : \nabla \rho \\ &= (2\pi)^2 \sum_{k \in \mathbb{Z}^2} -(\widehat{C^0}(k) \cdot \widehat{\rho}(k) \otimes (ik)) : (ik)^* \otimes \widehat{\rho}^*(k) \\ &= (2\pi)^2 \sum_{k \in \mathbb{Z}^2} -|k|^2 (\widehat{C^0}(k) \cdot \widehat{\rho}(k)) : \widehat{\rho}^*(k) \\ &\geq (2\pi)^2 m \sum_{k \in \mathbb{Z}^2} |k|^3 |\widehat{\rho}(k)|^2, \end{aligned}$$

where we have used assumption (A) in the last line. Then we define

$$\|\rho\|_{\left(\dot{H}^{\frac{3}{2}}(\mathbb{T}^2)\right)^2}^2 := \sum_{k \in \mathbb{Z}^2} |k|^3 |\widehat{\rho}(k)|^2.$$

Step 2 (estimate on the time derivative of the energy). Let us fix $T > 0$. We know that any solution ρ^n given by Theorem 2.1 belongs to the space $W^{2,1;2}(Q_T)$. In particular, using the following general fact (because of assumption (A))

$$\int_{\mathbb{T}^2} -(C^0 \star \rho) \cdot \rho = \operatorname{Re} \left(\sum_{k \in \mathbb{Z}^2} -|k| \left(\widehat{C}^0 \left(\frac{k}{|k|} \right) \cdot \widehat{\rho}(k) \right) \cdot \widehat{\rho}^*(k) \right)$$

we deduce that the energy $\mathcal{E}(\rho^n(t))$ is well-defined for almost every $t \in [0, T)$ and that, for almost every time $t \in [0, T)$, we have

$$\begin{aligned} \frac{d}{dt} \mathcal{E}(\rho^n(t)) &= \int_{\mathbb{T}^2} -\sigma[\rho^n] \cdot \rho_t^n \\ &= \int_{\mathbb{T}^2} -h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 - \varepsilon \sigma[\rho^n] \cdot \Delta \rho^n \\ &= \int_{\mathbb{T}^2} -h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 \\ &\quad - \int_{\mathbb{T}^2} \varepsilon \{W''(\rho^n) : ((\nabla \rho^n)^T \cdot \nabla \rho^n) - (C^0 \star (\nabla \rho^n)^T) : \nabla \rho^n\}. \end{aligned}$$

Therefore,

$$\begin{aligned} &\frac{d}{dt} \mathcal{E}(\rho^n(t)) + \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 \\ (3.4) \quad &\leq C\varepsilon \left\{ \int_{\mathbb{T}^2} |\nabla \rho^n|^2 + (C^0 \star (\nabla \rho^n)^T) : \nabla \rho^n \right\}. \end{aligned}$$

But now (with a generic constant C)

$$\begin{aligned} \|\nabla \rho^n\|_{(L^2(\mathbb{T}^2))^{2 \times 2}}^2 &\leq C \sum_{k \in \mathbb{Z}^2} |k|^2 |\widehat{\rho}^n(k)|^2 \\ &\leq C \sum_{k \in \mathbb{Z}^2} |k|^{\frac{3}{2}} |\widehat{\rho}^n(k)| \cdot |k|^{\frac{1}{2}} |\widehat{\rho}^n(k)| \\ &\leq C \left(\sum_{k \in \mathbb{Z}^2} \frac{1}{2\alpha} |k|^3 |\widehat{\rho}^n(k)|^2 + \frac{\alpha}{2} |k| |\widehat{\rho}^n(k)|^2 \right) \\ &\leq C \left(\int_{\mathbb{T}^2} -\frac{1}{\alpha} (C^0 \star (\nabla \rho^n)^T) : \nabla \rho^n + \int_{\mathbb{T}^2} -\alpha (C^0 \star \rho^n) \cdot \rho^n \right), \end{aligned}$$

where $\widehat{\rho}^n(k)$ are the Fourier coefficients of ρ^n and α is a constant which will be precised later. We then deduce finally that

$$\begin{aligned} &\frac{d}{dt} \mathcal{E}(\rho^n(t)) + \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 \\ (3.5) \quad &\leq -C\varepsilon \left(1 - \frac{1}{\alpha} \right) \int_{\mathbb{T}^2} -(C^0 \star (\nabla \rho^n)^T) : \nabla \rho^n + C\varepsilon \alpha \int_{\mathbb{T}^2} -(C^0 \star \rho^n) \cdot \rho^n \\ &\leq -C\varepsilon \|\rho^n(t)\|_{\left(\dot{H}^{\frac{3}{2}}(\mathbb{T}^2)\right)^2}^2 + C\varepsilon \left(1 + \mathcal{E}(\rho^n(t)) + |\sigma^0| \left| \int_{\mathbb{T}^2} \rho^n(t) \right| \right) \end{aligned}$$

for α chosen large enough, with C a suitable positive constant.

Step 3 (estimate on the time derivative of the mean value of the solution). Integrating in space equation (2.1), we get

$$\frac{d}{dt} \int_{\mathbb{T}^2} \rho^n(t) = \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} (\nabla \rho^n)^T \cdot \nabla \rho^n \cdot \sigma[\rho^n]$$

and then

$$\begin{aligned} (3.6) \quad \frac{d}{dt} \left| \int_{\mathbb{T}^2} \rho^n(t) \right| &\leq \int_{\mathbb{T}^2} (h^n(|\nabla \rho^n|) |\nabla \rho^n|)^{\frac{1}{2}} \cdot \left((h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1})^{\frac{1}{2}} |\nabla \rho^n \cdot \sigma[\rho^n]| \right) \\ &\leq \int_{\mathbb{T}^2} (1 + |\sigma^0|) h^n(|\nabla \rho^n|) |\nabla \rho^n| \\ &\quad + \frac{1}{4(1 + |\sigma^0|)} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2. \end{aligned}$$

Step 4 (estimate on the energy). Setting

$$(3.7) \quad F^n(t) = 1 + \mathcal{E}(\rho^n(t)) + (1 + |\sigma^0|) \left| \int_{\mathbb{T}^2} \rho^n(t) \right| + (1 + |\sigma^0|)^4,$$

we deduce from (3.5) and (3.6) that

$$\begin{aligned} &\frac{d}{dt} F^n(t) + \frac{3}{4} \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 \\ &\leq -C\varepsilon \|\rho^n(t)\|_{\dot{H}^{\frac{3}{2}}(\mathbb{T}^2)}^2 + C\varepsilon \left(1 + \mathcal{E}(\rho^n(t)) + |\sigma^0| \left| \int_{\mathbb{T}^2} \rho^n(t) \right| \right) \\ &\quad + (1 + |\sigma^0|)^2 \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|. \end{aligned}$$

Now we remark that

$$\begin{aligned} (1 + |\sigma^0|)^2 \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n| &\leq (1 + |\sigma^0|)^2 \int_{\mathbb{T}^2} |\nabla \rho^n| \\ &\leq \frac{C\varepsilon}{2} \int_{\mathbb{T}^2} |\nabla \rho^n|^2 + \frac{(1 + |\sigma^0|)^4}{2C\varepsilon}. \end{aligned}$$

Using the fact that (since the domain is bounded)

$$\int_{\mathbb{T}^2} |\nabla \rho^n|^2 \leq \|\rho^n(t)\|_{\dot{H}^{\frac{3}{2}}(\mathbb{T}^2)}^2,$$

we get

$$\begin{aligned} (3.8) \quad \frac{d}{dt} F^n(t) + \frac{3}{4} \int_{\mathbb{T}^2} h^n(|\nabla \rho^n|) |\nabla \rho^n|^{-1} |\nabla \rho^n \cdot \sigma[\rho^n]|^2 + \frac{C\varepsilon}{2} \|\rho^n(t)\|_{\dot{H}^{\frac{3}{2}}(\mathbb{T}^2)}^2 \\ \leq C\varepsilon \left(1 + \mathcal{E}(\rho^n(t)) + |\sigma^0| \left| \int_{\mathbb{T}^2} \rho^n(t) \right| \right) + \frac{(1 + |\sigma^0|)^4}{2C\varepsilon} \\ \leq \frac{C}{\varepsilon} F^n(t). \end{aligned}$$

This implies, using the Gronwall lemma,

$$(3.9) \quad F^n(t) \leq F^n(0)e^{\frac{C}{\varepsilon}t}.$$

Step 5 (estimate on ρ^n). Let us first remark that

$$(3.10) \quad \mathcal{E}(\rho^n(t)) \geq \int -\frac{1}{2} (C^0 \star \rho^n) \cdot \rho^n - |\sigma^0| \left| \int_{\mathbb{T}^2} \rho^n(t) \right|.$$

Using (3.9), (3.10), and the definition of $F^n(t)$ yields

$$\int_{\mathbb{T}^2} -\frac{1}{2} (C^0 \star \rho^n) \cdot \rho^n + \left| \int_{\mathbb{T}^2} \rho^n(t) \right| \leq Ce^{\frac{C}{\varepsilon}t}.$$

Using Step 1, we then get

$$\|\rho^n\|_{L^\infty((0,T);(\dot{H}^{\frac{1}{2}}(\mathbb{T}^2))^2)}^2 \leq Ce^{\frac{C}{\varepsilon}T} \quad \text{and} \quad \left| \int_{\mathbb{T}^2} \rho^n(t) \right| \leq Ce^{\frac{C}{\varepsilon}T}.$$

This implies (3.1). Taking the integral \int_0^T in (3.8) and using the fact that $\forall t \leq T$, $F^n(t) \geq 0$, we get

$$\|h_n(|\nabla \rho^n|)|\nabla \rho^n|^{-\frac{1}{2}}|\nabla \rho^n \cdot \sigma[\rho^n]|\|_{(L^2(Q_T))^2}^2 + \varepsilon C \|\rho^n\|_{L^2((0,T);(\dot{H}^{\frac{3}{2}}(\mathbb{T}^2))}^2 \leq Ce^{\frac{C}{\varepsilon}T},$$

which implies (3.2) and (3.3). \square

3.2. Proof of Theorem 1.1. We are now able to prove Theorem 1.1. In this section, we denote by C a generic constant which can depend on ρ^0, ε , and T but which does not depend on n .

Proof of Theorem 1.1. Let $T > 0$. The idea of the proof is to pass to the limit in (2.1). The only difficulty is to prove that the nonlinear term $f^n[\rho^n]$ converges in a certain sense to $|\nabla \rho|^{-1}(\nabla \rho)^T \cdot \nabla \rho \cdot \sigma[\rho]$, where ρ is the limit of ρ^n in an appropriate norm. The proof is decomposed into five steps.

Step 1 (a priori bound on $f^n[\rho^n]$). We have the following estimate on $f^n[\rho^n]$:

$$(3.11) \quad \|f^n[\rho^n]\|_{(L^{\frac{4}{3}}(Q_T))^2} \leq C.$$

To prove this, let us write

$$f^n[\rho^n] = (|\nabla \rho^n|^{-1}(\nabla \rho^n)^T) \cdot \left(|\nabla \rho^n|^{\frac{1}{2}} \right) \left(h_n(|\nabla \rho^n|)|\nabla \rho^n|^{-\frac{1}{2}}\nabla \rho^n \cdot \sigma[\rho^n] \right).$$

Using (3.3), we have that the last term is bounded in $(L^2(Q_T))^2$ by C . Moreover, the first term is bounded by 1 in $(L^\infty(Q_T))^{2 \times 2}$, and then we just have to bound the term $|\nabla \rho^n|^{\frac{1}{2}}$ in $L^4(Q_T)$. Using (3.2), we have

$$(3.12) \quad \| |\nabla \rho^n|^{\frac{1}{2}} \|_{L^4(Q_T)} = \left(\int_{Q_T} |\nabla \rho^n|^2 \right)^{\frac{1}{4}} \leq \|\rho^n\|_{L^2((0,T);(\dot{H}^{\frac{3}{2}}(\mathbb{T}))^2)}^{\frac{1}{2}} \leq C.$$

This implies (3.11).

Step 2 (strong convergence of $\nabla\rho^n$ in $L^2((0, T); (L^{\frac{4}{3}}(\mathbb{T}^2))^{2 \times 2})$). Using the parabolic estimates for the heat equation (see [16, Chapter 4.3, p. 80 and Chapter 4.9, p. 341]) and Step 1, we get

$$(3.13) \quad \|\nabla\rho^n\|_{W^{\frac{1}{2}; \frac{4}{3}}((0, T); (L^{\frac{4}{3}}(\mathbb{T}^2))^{2 \times 2})} \leq C$$

where, for a Banach space B ,

$$W^{\frac{1}{2}; p}((0, T); B) = \left\{ g \in L^p((0, T); B), \int_0^T \int_0^T \frac{\|g(t) - g(s)\|_B^p}{|t - s|^{\frac{1}{2}p+1}} dt ds < \infty \right\}$$

is equipped with the following norm:

$$\|g\|_{W^{\frac{1}{2}; p}((0, T); B)} := \left(\int_0^T \int_0^T \frac{\|g(t) - g(s)\|_B^p}{|t - s|^{\frac{1}{2}p+1}} dt ds \right)^{\frac{1}{p}}.$$

Moreover, using (3.2) we get

$$(3.14) \quad \|\nabla\rho^n\|_{L^2((0, T); (H^{\frac{1}{2}}(\mathbb{T}^2))^{2 \times 2})} \leq C.$$

We then use the following lemma.

LEMMA 3.2 (compactness result). *Let $(g_n)_n$ be a sequence uniformly bounded in*

$$L^2\left((0, T); H^{\frac{1}{2}}(\mathbb{T}^2)\right) \cap W^{\frac{1}{2}; \frac{4}{3}}\left((0, T); L^{\frac{4}{3}}(\mathbb{T}^2)\right);$$

then, for a subsequence,

$$g_n \rightarrow g \text{ strongly in } L^2\left((0, T); L^{\frac{4}{3}}(\mathbb{T}^2)\right).$$

Formally, the proof uses the fact that $H^{\frac{1}{2}} \subset L^{\frac{4}{3}}$ with compact injection in space, while the compactness in time comes from (3.13). We refer to Simon [21, Corollary 5, p. 86] for a more general result and for the proof of this lemma.

Using (3.13), (3.14), and Lemma 3.2, we then deduce that, for a subsequence, $\nabla\rho^n \rightarrow \nabla\rho$ strongly in $L^2((0, T); (L^{\frac{4}{3}}(\mathbb{T}^2))^2)$ and almost everywhere.

Step 3. (weak convergence of $\sigma[\rho^n]$ in $L^2((0, T); (L^4(\mathbb{T}^2))^2)$). We have $H^{\frac{1}{2}}(\mathbb{T}^2) \subset L^4(\mathbb{T}^2)$ with continuous injection (see Adams [1, Theorem 7.57, p. 217]). So $L^2((0, T); H^{\frac{1}{2}}(\mathbb{T}^2)) \subset L^2((0, T); L^4(\mathbb{T}^2))$ with continuous injection. We then deduce from (3.2) that

$$(3.15) \quad \|\nabla\rho^n\|_{L^2((0, T); (L^4(\mathbb{T}^2))^{2 \times 2})} \leq C.$$

Using Lemma 2.2, we then get

$$(3.16) \quad \|C^0 \star \rho^n\|_{L^2((0, T); (L^4(\mathbb{T}^2))^2)} \leq C.$$

Using the fact that the application $W_{x,t}^{2,1, \frac{4}{3}}(Q_T) \mapsto L^{\frac{4}{3}}(Q_T)$ is compact and the converse of the Lebesgue theorem, we deduce that $W'(\rho^n) \rightarrow W'(\rho)$ almost everywhere. This implies that $\sigma[\rho^n] \rightharpoonup \sigma[\rho]$ in $L^2((0, T); (L^4(\mathbb{T}^2))^2)$.

Step 4. (passing to the limit). Using Steps 2 and 3 and the fact that $|\nabla\rho^n|^{-1}\nabla\rho^n$ is bounded by 1, we deduce that

$$f_n[\rho^n] \rightarrow |\nabla\rho|^{-1}(\nabla\rho)^T \cdot \nabla\rho \cdot \sigma[\rho] \quad \text{in the distributions sense.}$$

By passing to the limit in (2.1), we obtain

$$(3.17) \quad \rho_t - \epsilon\Delta\rho = |\nabla\rho|^{-1}(\nabla\rho)^T \cdot \nabla\rho \cdot \sigma[\rho] \text{ in } \mathcal{D}'((0, T) \times \mathbb{T}^2).$$

Step 5 (initial condition). Using the fact that ρ_t^n are bounded uniformly in $L^{\frac{4}{3}}(Q_T)$ (by parabolic estimates for the heat equation and Step 1), we deduce that (uniformly in n)

$$\|\rho^n(t+h) - \rho^n(t)\|_{(L^{\frac{4}{3}}(\mathbb{T}^2))^2} \leq Ch^{\frac{1}{4}}\|\rho_t^n\|_{L^{\frac{4}{3}}((0,T);(L^{\frac{4}{3}}(\mathbb{T}^2))^2)}$$

and then $\rho \in C^0((0, T); (L^{\frac{4}{3}}(\mathbb{T}^2))^2)$ and $\rho|_{t=0} = \rho^0$.

This achieves the proof of Theorem 1.1. \square

Appendix A.

Full proof of Lemma 2.2. Here we do the proof for any $p \in (1, +\infty)$. Under assumption (A), there exists a constant $C > 0$ depending only on p such that the following result holds for any $\tilde{\rho} \in W^{1,p}(\mathbb{R}^2)$:

$$|\tilde{C}^0 \star_{\mathbb{R}^2} \tilde{\rho}|_{(L^p(\mathbb{R}^2))^2} \leq \frac{C}{m} |\nabla\tilde{\rho}|_{(L^p(\mathbb{R}^2))^{2 \times 2}},$$

where the Fourier transform of \tilde{C}^0 satisfies $\widehat{\tilde{C}^0} = M$ with M as in (1.9).

This result can be found in the scalar case on \mathbb{R}^n in Stein [22, Proposition 5, p. 251] or Coifman, Meyer [10, Theorem 9, p. 39 and Proposition 2, p. 41]. See, also, Calderon–Zygmund inequalities Theorem 2.7.2 in Morrey [19]. Here the convolution by \tilde{C}^0 is a multiplier operator in the class S^1 of pseudodifferential operators. We then get the result in the vectorial case, summing the scalar components. See, also, the book of Garroni and Menaldi [12] for complements on integro-differential operators.

The fact that the result holds on the torus \mathbb{T}^2 is then classical. We prove it for the convenience of the reader. To this end, we consider a smooth function φ such that

$$\varphi(x) = 1 \quad \text{on } [-1/3, 1/3]^2, \quad \text{supp } \varphi \subset [-2/3, 2/3]^2, \quad \text{and } 0 \leq \varphi \leq 1$$

such that

$$\sum_{k \in \mathbb{Z}^2} \varphi(x - k) = 1.$$

For any smooth function $\rho : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, which is \mathbb{Z}^2 -periodic, we then set for $K > 0$

$$(S_{2K}\rho)(x) = \sum_{|k| \leq 2K, k \in \mathbb{Z}^2} \varphi(x - k)\rho(x).$$

Therefore, we get for $K > 0$ large enough

$$\begin{aligned} & |B_K| \left\{ |\tilde{C}^0 \star_{\mathbb{R}^2} \rho|_{L^p((-1/2, 1/2)^2)} + 0(1/K) \right\} \\ & \leq |\tilde{C}^0 \star_{\mathbb{R}^2} \rho|_{L^p(B_K)} \\ & \leq |\tilde{C}^0 \star_{\mathbb{R}^2} (S_{2K}\rho)|_{L^p(\mathbb{R}^2)} + |\tilde{C}^0 \star_{\mathbb{R}^2} (\rho - (S_{2K}\rho))|_{L^p(B_K)} \\ & \leq \frac{C}{m} |\nabla(S_{2K}\rho)|_{L^p(\mathbb{R}^2)} + |\tilde{C}^0 \star_{\mathbb{R}^2} (\rho - (S_{2K}\rho))|_{L^p(B_K)} \\ & \leq \frac{C}{m} |B_{2K}| \left\{ |\nabla\rho|_{(L^p((-1/2, 1/2)^{2 \times 2}))} + 0(1/K) \right\} + |\rho|_{(L^\infty(\mathbb{R}^2))^2} |B_K| \int_{|z| \geq K-1} |\tilde{C}^0(z)|. \end{aligned}$$

Using the fact that $\int_{|z| \geq K-1} |\tilde{C}^0(z)| = O(1/K)$, dividing by $|B_K|$, and taking the limit as $K \rightarrow +\infty$, we get

$$|\tilde{C}^0 \star_{\mathbb{R}^2} \rho|_{(L^p(\mathbb{T}^2))^2} \leq \frac{C}{m} \frac{|B_2|}{|B_1|} |\nabla\rho|_{(L^p(\mathbb{T}^2))^{2 \times 2}},$$

i.e.,

$$|C^0 \star_{\mathbb{T}^2} \rho|_{(L^p(\mathbb{T}^2))^2} \leq \frac{4C}{m} |\nabla\rho|_{(L^p(\mathbb{T}^2))^{2 \times 2}}$$

with

$$C^0(x) = \sum_{k \in \mathbb{Z}^2} \tilde{C}^0(x - k).$$

We then get the final result by density of smooth functions in $(W^{1,p}(\mathbb{T}^2))^2$. \square

Proof of Proposition 2.8. For simplicity of notation, we denote by $g(t)$ the function $x \mapsto g(t, x)$. We have

$$\begin{aligned} \|g^m(t) - \rho^0\|_{(L^2(\mathbb{T}^2))^2} & \leq \int_0^t ds \|g_t^m(s)\|_{(L^2(\mathbb{T}^2))^2} \\ & \leq \sqrt{t} \|g_t^m\|_{(L^2(Q_T))^2}. \end{aligned}$$

Using the fact that g^m is bounded uniformly in $W^{2,1;2}(Q_T)$ (this is a consequence of the fact that $g^m \rightharpoonup g$ in $W^{2,1;2}(Q_T)$), we get

$$(A.1) \quad \|g^m(t) - \rho^0\|_{(L^2(\mathbb{T}^2))^2} \leq C\sqrt{t}.$$

Now let $\varphi \in C_c^\infty([0, +\infty), \mathbb{R})$ be such that $\varphi \geq 0$. Using (A.1), we get that

$$\int_0^t ds \|g^m(s) - \rho^0\|_{(L^2(\mathbb{T}^2))^2}^2 \varphi(s) \leq C \int_0^t ds s \varphi(s).$$

Using Fatou's lemma, we deduce that

$$\int_0^t \left(\|g(s) - \rho^0\|_{(L^2(\mathbb{T}^2))^2}^2 - Cs \right) \varphi(s) \leq 0.$$

Using that $\varphi \geq 0$ is arbitrary, we deduce that, for almost every t , we have

$$\|g(t) - \rho^0\|_{(L^2(\mathbb{T}^2))^2}^2 \leq \sqrt{Ct}.$$

This implies the result. \square

Acknowledgments. The authors thank A. El Hajj for fruitful discussion in the preparation of this paper. They also thank the referees for their comments and questions that helped to improve the manuscript.

REFERENCES

- [1] R. A. ADAMS, *Sobolev Spaces*, Pure Appl. Math. 65, Academic Press, New York, London, 1975.
- [2] S. ALLEN AND J. CAHN, *A microscopic theory for the antiphase boundary motion and its application to antiphase domain coarsening*, Acta Metallurgica, 27 (1979), pp. 1085–1095.
- [3] O. ALVAREZ, P. HOCH, Y. LE BOUAR, AND R. MONNEAU, *Dislocation dynamics: Short time existence and uniqueness of the solution*, Arch. Ration. Mech. Anal., 85 (2006), pp. 371–414.
- [4] L. BOCCARDO, F. MURAT, AND J.-P. PUEL, *Existence de solutions faibles pour des équations elliptiques quasi-linéaires à croissance quadratique*, in Nonlinear Partial Differential Equations and Their Applications. Collège de France Seminar Vol. IV (Paris, 1981/1982), Res. Notes in Math. 84, Pitman, Boston, MA, 1983, pp. 19–73.
- [5] L. BOCCARDO, F. MURAT, AND J.-P. PUEL, *Résultats d'existence pour certains problèmes elliptiques quasilineaires*, Ann. Sc. Norm. Super. Pisa Cl. Sci., 11 (1984), pp. 213–235.
- [6] L. BOCCARDO, F. MURAT, AND J.-P. PUEL, *Existence results for some quasilinear parabolic equations*, Nonlinear Anal., 13 (1989), pp. 373–392.
- [7] A. BONNET, *On the regularity of the edge set of Mumford-Shah minimizers*, in Variational Methods for Discontinuous Structures (Como, 1994), Progr. Nonlinear Differential Equations Appl. 25, Birkhäuser, Basel, 1996, pp. 93–103.
- [8] L. BRONSARD AND F. REITICH, *On three-phase boundary motion and the singular limit of a vector-valued Ginzburg-Landau equation*, Arch. Ration. Mech. Anal., 124 (1993), pp. 355–379.
- [9] S. CACACE AND A. GARRONI, *A Multi-phase Transition Model for Dislocations with Interfacial Microstructure*, Interfaces and Free Boundaries, to appear.
- [10] R. R. COIFMAN AND Y. MEYER, *Au Delà des Opérateurs Pseudo-différentiels*, Astérisque 57, Société Mathématique de France, Paris, 1978.
- [11] L. C. EVANS, *Partial Differential Equations*, Grad. Stud. Math. 19, American Mathematical Society, Providence, RI, 1998.
- [12] M. G. GARRONI AND J.-L. MENALDI, *Green Functions for Second Order Parabolic Integro-differential Problems*, Pitman Res. Notes Math. 275, Longman Scientific & Technical, Harlow, UK, 1992.
- [13] J. P. HIRTH AND J. LOTHE, *Theory of Dislocations*, 2nd ed., Krieger, Malabar, FL, 1992.
- [14] D. HULL AND D. BACON, *Introduction to Dislocations*, 4th ed., Butterworth-Heinemann, Oxford, 2001.
- [15] M. KOSLOWSKI AND M. ORTIZ, *A multi-phase field model of planar dislocation networks*, Model. Simul. Material Sci. Eng., 12 (2004), pp. 1087–1097.
- [16] O. A. LADYŽENSKAJA, V. A. SOLONNIKOV, AND N. N. URAL'CEVA, *Linear and Quasilinear Equations of Parabolic Type*, Transl. Math. Monogr. 23, American Mathematical Society, Providence, RI, 1967.
- [17] J. LERAY, *Sur le mouvement d'un liquide visqueux emplissant l'espace*, Acta Math., 63 (1934), pp. 193–248.
- [18] J.-L. LIONS, *Quelques Méthodes de Résolution des Problèmes aux Limites Non Linéaires*, Dunod, Paris, 1969.
- [19] J. C. B. MORREY, *Multiple Integrals in the Calculus of Variations*, Grundlehren Math. Wiss. 130, Springer-Verlag New York, 1966.
- [20] D. RODNEY, Y. LE BOUAR, AND A. FINEL, *Phase field methods and dislocations*, Acta Materialia, 51 (2003), pp. 17–30.
- [21] J. SIMON, *Compact sets in the space $L^p(0, T; B)$* , Ann. Mat. Pura Appl., 146 (1987), pp. 65–96.
- [22] E. M. STEIN, *Harmonic Analysis: Real-variable Methods, Orthogonality, and Oscillatory Integrals*, Princeton Math. Ser. 43, Princeton University Press, Princeton, NJ, 1993.
- [23] J. E. TAYLOR, *The motion of multiple-phase junctions under prescribed phase-boundary velocities*, J. Differential Equations, 119 (1995), pp. 109–136.
- [24] J. E. TAYLOR, *A variational approach to crystalline triple-junction motion*, J. Stat. Phys., 95 (1999), pp. 1221–1244.

MINIMIZATION OF ELECTROSTATIC FREE ENERGY AND THE POISSON–BOLTZMANN EQUATION FOR MOLECULAR SOLVATION WITH IMPLICIT SOLVENT*

BO LI†

Abstract. In an implicit-solvent description of the solvation of charged molecules (solutes), the electrostatic free energy is a functional of concentrations of ions in the solvent. The charge density is determined by such concentrations together with the point charges of the solute atoms, and the electrostatic potential is determined by the Poisson equation with a variable dielectric coefficient. Such a free-energy functional is considered in this work for both the case of point ions and that of ions with a uniform finite size. It is proved for each case that there exists a unique set of equilibrium concentrations that minimize the free energy and that are given by the corresponding Boltzmann distributions through the equilibrium electrostatic potential. Such distributions are found to depend on the boundary data for the Poisson equation. Pointwise upper and lower bounds are obtained for the free-energy minimizing concentrations. Proofs are also given for the existence and uniqueness of the boundary-value problem of the resulting Poisson–Boltzmann equation that determines the equilibrium electrostatic potential. Finally, the equivalence of two different forms of such a boundary-value problem is proved.

Key words. implicit solvent, electrostatic free energy, ionic concentrations, electrostatic potentials, the Poisson–Boltzmann equation, variational methods, nonlinear elliptic interface problems

AMS subject classifications. 35J, 35Q, 49S, 82D, 92C

DOI. 10.1137/080712350

1. Introduction. It has long been realized that the electrostatic potential of a charged molecular system extremizes an electrostatic free-energy functional [3, 6, 12, 13, 15, 18, 20, 26, 28, 29]. In a simple setting, this functional is given by

$$F[c_1, \dots, c_M; \psi] = \int \left\{ -\frac{\varepsilon}{8\pi} |\nabla\psi|^2 + \rho\psi + \beta^{-1} \sum_{j=1}^M c_j [\ln(\Lambda^3 c_j) - 1] - \sum_{j=1}^M \mu_j c_j \right\} dx,$$

where c_1, \dots, c_M are ionic concentrations, ψ is an electrostatic potential, ε is the dielectric constant, ρ is the charge density defined to be a linear combination of the ionic concentrations, β is the inverse thermal energy, Λ is the thermal de Broglie wavelength, and μ_j is the chemical potential of the j th ionic species. Throughout, we use the electrostatics CGS units. We also use $\log x$ to denote the natural logarithm of $x > 0$. Extremizing this functional with respect to the concentrations and the potential lead to the Boltzmann distribution of concentrations and the Poisson equation for the equilibrium potential, respectively [3, 6, 12, 13, 15, 26, 28]. Notice, however, that this free-energy functional is concave with respect to the electrostatic potential. Therefore, the extremizing concentrations and potential do not minimize this free-energy functional, rather they form an unstable saddle point of the system [1, 6, 12, 13]. This flaw of theory is removed in the free-energy minimization approach that was proposed

*Received by the editors January 3, 2008; accepted for publication (in revised form) December 9, 2008; published electronically March 20, 2009. This work was supported by the U.S. National Science Foundation (NSF) through the grant DMS-0451466 and DMS-0811259, by the NSF Center for Theoretical Biological Physics with the NSF grant PHY-0822283, and by the U.S. Department of Energy through the grant DE-FG02-05ER25707.

<http://www.siam.org/journals/sima/40-6/71235.html>

†Department of Mathematics and the NSF Center for Theoretical Biological Physics, University of California, San Diego, 9500 Gilman Drive, Mail code: 0112, La Jolla, CA 92093-0112 (bli@math.ucsd.edu).

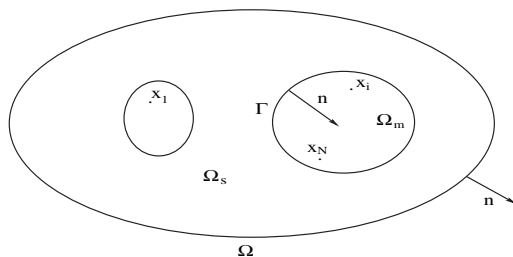


FIG. 1. *The geometry of a solvation system with an implicit solvent.*

in [12, 20]. The key point in this new approach is that the electrostatic free-energy functional depends solely on the ionic concentrations and the electrostatic potential is determined by such concentrations through the Poisson equation. In the recent article [5], this free-energy minimization approach was revisited and applied to the implicit-solvent (or continuum-solvent) description of solvation.

The present work is a mathematical study of the free-energy minimization approach to the electrostatics applied to the solvation of molecules with an implicit-solvent. Such application introduces additional mathematical complications due to the presence of point charges in solutes and the dielectric boundaries. We consider both the case of point ions—ions modeled as points without volumes—and that of ions with a uniform finite size. The finite-size effect of ions is known to be important in continuum modeling of electrostatics in molecular systems. Our analysis shows particularly that the free-energy minimizing ionic concentrations are uniformly bounded from above and away from zero at each spatial point. This uniform boundedness, which is proved by somewhat tedious constructions, is a consequence of the property that the free-energy minimizing concentrations have a large entropy. We do not consider the more general case of ions with different sizes for which there seems no explicit Boltzmann distributions.

Consider now the solvation of charged molecules with an implicit solvent [27]. We divide the entire region Ω of the solvation system into the region of solute molecules $\Omega_m \subset \mathbb{R}^3$ that is possibly multiply connected, the region of solvent (such as salted water) $\Omega_s \subset \mathbb{R}^3$, and the solute-solvent interface $\Gamma = \partial\Omega_m \cap \partial\Omega_s$; cf. Figure 1. This interface Γ serves as the dielectric boundary. Assume the solutes consist of N atoms with the i th one located at x_i and carrying a charge Q_i . Assume also there are M ionic species in the solvent with $q_j = ez_j$ the buck charge of the j th ionic species, where e is the elementary charge and z_j the valence of j th ionic species. Denote by $c_j = c_j(x)$ the local concentration at $x \in \Omega_s$ of the j th ionic species. Following the common assumption that the mobile ions in the solvent cannot penetrate the dielectric boundary Γ , we define $c_j(x) = 0$ for all $x \in \Omega_m$ and $1 \leq j \leq M$.

We consider two mean-field approximations of the electrostatic free energy of the solvation system as functionals of the local ionic concentrations $c = (c_1, \dots, c_M)$ in the solvent region. In the first one, point ions are assumed, and the related electrostatic free-energy functional is given by [5, 12, 19, 20, 26]

$$\begin{aligned}
 F_0[c] = & \frac{1}{2} \sum_{i=1}^N Q_i (\psi - \psi_{vac})(x_i) + \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) \psi dx \\
 (1.1) \quad & + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} c_j [\log(a^3 c_j) - 1] dx - \sum_{j=1}^M \int_{\Omega_s} \mu_j c_j dx.
 \end{aligned}$$

In the second approximation, all ions are assumed to have a uniform linear size, and the related free-energy functional is given by [3, 20]

$$(1.2) \quad F_a[c] = \frac{1}{2} \sum_{i=1}^N Q_i (\psi - \psi_{vac})(x_i) + \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) \psi dx \\ + \beta^{-1} \sum_{j=0}^M \int_{\Omega_s} c_j [\log(a^3 c_j) - 1] dx - \sum_{j=1}^M \int_{\Omega_s} \mu_j c_j dx,$$

where the summation in the β^{-1} term starts from $j = 0$ and

$$(1.3) \quad c_0(x) = a^{-3} \left[1 - \sum_{j=1}^M a^3 c_j(x) \right] \quad \forall x \in \Omega_s.$$

In (1.1) and (1.2), ψ is the electrostatic potential of the solvation system,

$$(1.4) \quad \psi_{vac}(x) = \sum_{i=1}^N \frac{Q_i}{\varepsilon_m |x - x_i|}$$

defines the electrostatic potential generated by all the point charges Q_i at x_i in a medium with the dielectric constant ε_m (usually taken as that in the vacuum), $a > 0$ is a constant, and μ_j is the constant chemical potential of the j th ionic species. The constant $a > 0$ represents in (1.1) a nonphysical cut-off which is often chosen to be the thermal de Broglie wavelength and in (1.2) the uniform linear size of ions.

The electrostatic potential ψ is determined by the Poisson equation

$$(1.5) \quad \nabla \cdot \varepsilon_\Gamma \nabla \psi = -4\pi\rho \quad \text{in } \Omega,$$

where ε_Γ is the dielectric coefficient and ρ is the charge density, together with a boundary condition which is usually taken to be

$$(1.6) \quad \psi = \psi_0 \quad \text{on } \partial\Omega,$$

where ψ_0 is a given function. The dielectric coefficient is defined to be

$$(1.7) \quad \varepsilon_\Gamma(x) = \begin{cases} \varepsilon_m & \text{if } x \in \Omega_m, \\ \varepsilon_s & \text{if } x \in \Omega_s, \end{cases}$$

where ε_m and ε_s are the dielectric constants of the solutes and the solvent, respectively. The charge density is given by

$$(1.8) \quad \rho = \sum_{i=1}^N Q_i \delta_{x_i} + \sum_{j=1}^M q_j c_j \quad \text{in } \Omega,$$

where δ_{x_i} denotes the Dirac delta function centered at x_i .

The first two terms in (1.1) or (1.2) represent the internal electrostatic energy, which are often written formally as the integral of $\rho\psi/2$ over the entire region Ω . Based on Born's definition [2], the contribution to the electrostatic free energy due to the solute point charges is given as the first term in (1.1) or (1.2) though the reaction

field $\psi - \psi_{vac}$. The β^{-1} term represents the ideal gas entropy. The term $1 - \sum_{j=1}^M a^3 c_j$ in (1.2) is the concentration of solvent molecules. It describes the effect of finite size of ions. The last term in (1.1) or (1.2) accounts for a constant chemical potential in the system. The osmotic pressure from the mobile ions is dropped, since it is only an additive constant to the free-energy functional in the present setting. We remark that the use of notations F_0 and F_a does not indicate that we can obtain the functional F_0 by simply setting $a = 0$ in F_a .

In this work, we prove the following results:

- (1) For each of the free-energy functionals (1.1) and (1.2), there admits a unique minimizer c_1, \dots, c_M , which is also the unique equilibrium, in an admissible set of concentrations. Moreover, such concentrations and the corresponding equilibrium electrostatic potential ψ are related by the boundary-data dependent Boltzmann distributions

$$(1.9) \quad c_j(x) = \begin{cases} c_j^\infty e^{-\beta q_j [\psi(x) - \hat{\psi}_0(x)/2]} & \text{for point ions,} \\ \frac{c_j^\infty e^{-\beta q_j [\psi(x) - \hat{\psi}_0(x)/2]}}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i [\psi(x) - \hat{\psi}_0(x)/2]}} & \text{for finite-size ions,} \end{cases}$$

for a.e. $x \in \Omega_s$ and $1 \leq j \leq M$, where $c_j^\infty = a^{-3} e^{\beta \mu_j}$ and $\hat{\psi}_0 \in H^1(\Omega)$ is determined by

$$(1.10) \quad \begin{cases} \int_{\Omega} \varepsilon_{\Gamma} \nabla \hat{\psi}_0 \cdot \nabla \eta \, dx = 0 & \forall \eta \in H_0^1(\Omega), \\ \hat{\psi}_0 = \psi_0 & \text{on } \partial\Omega. \end{cases}$$

The free-energy minimizing concentrations are shown to be uniformly bounded above and below away from zero. These results are summarized in Theorems 2.3–2.5 and Lemmas 3.4 and 3.5.

- (2) The equilibrium electrostatic potential ψ is the unique solution to the boundary-data dependent Poisson–Boltzmann equation (PBE) [3, 4, 10, 11, 16, 17, 20, 31], together with the boundary condition (1.6),

$$(1.11) \quad \nabla \cdot \varepsilon_{\Gamma} \nabla \psi + 4\pi \chi_{\Omega_s} \sum_{j=1}^M q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} \quad \text{in } \Omega$$

for the case of point ions, and

$$(1.12) \quad \nabla \cdot \varepsilon_{\Gamma} \nabla \psi + 4\pi \chi_{\Omega_s} \sum_{j=1}^M \frac{q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)}}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)}} = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} \quad \text{in } \Omega$$

for the case of finite-size ions, where χ_{Ω_s} is the characteristic function of Ω_s . These equations can be written together as

$$(1.13) \quad \nabla \cdot \varepsilon_{\Gamma} \nabla \psi - 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} \quad \text{in } \Omega,$$

where B' is the derivative of the function $B : \mathbb{R} \rightarrow \mathbb{R}$ defined by (1.14)

$$B(\psi) = \begin{cases} \sum_{j=1}^M \beta^{-1} c_j^\infty e^{-\beta q_j \psi} & \text{for point ions,} \\ \beta^{-1} a^{-3} \log \left(1 + a^3 \sum_{j=1}^M c_j^\infty e^{-\beta q_j \psi} \right) & \text{for finite-size ions.} \end{cases}$$

See Theorem 2.1.

- (3) The boundary-value problem of the PBE (1.13) and (1.6) is equivalent to the elliptic interface problem

$$(1.15) \quad \begin{cases} \nabla \cdot \varepsilon_m \nabla \psi = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} & \text{in } \Omega_m, \\ \nabla \cdot \varepsilon_s \nabla \psi - 4\pi B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) = 0 & \text{in } \Omega_s, \\ \llbracket \psi \rrbracket = \llbracket \varepsilon_\Gamma \nabla \psi \cdot n \rrbracket = 0 & \text{on } \Gamma, \\ \psi = \psi_0 & \text{on } \Omega. \end{cases}$$

Here and below, we denote for any function u on Ω , $u_m = u|_{\Omega_m}$, $u_s = u|_{\Omega_s}$, and $\llbracket u \rrbracket = u_s - u_m$ on Γ . See Theorem 2.2.

Two variations of the PBE (1.11) with $\psi_0 = 0$ are commonly used [8, 15, 29]. First, we have by the Taylor expansion and the electrostatic neutrality $\sum_{j=1}^M c_j^\infty = 0$ that

$$\sum_{j=1}^M q_j c_j^\infty e^{-\beta q_j \psi} \approx - \left(\sum_{j=1}^M \beta q_j^2 c_j^\infty \right) \psi,$$

if $|\psi|$ is small, leading to the linearized PBE [9]

$$\nabla \cdot \varepsilon_\Gamma \nabla \psi - \varepsilon_s \kappa^2 \chi_{\Omega_s} \psi = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} \quad \text{in } \Omega,$$

where $\kappa = \sqrt{4\pi\beta \sum_{i=1}^M q_i^2 c_i^\infty / \varepsilon_s^2}$ is the ionic strength or the inverse Debye–Hückel screening length. Clearly, all of our results for the nonlinear PBE (1.11) hold true for the linearized PBE. Second, for the common $z : -z$ type of salt such as NaCl in the solution, we have $M = 2$, $c_1^\infty = c_2^\infty$, and $q_1 = -q_2 = ze$. The PBE (1.11) reduces to the following sinh PBE:

$$\nabla \cdot \varepsilon_\Gamma \nabla \psi - 8\pi z e c_1^\infty \chi_{\Omega_s} \sinh(\beta z e \psi) = -4\pi \sum_{i=1}^N Q_i \delta_{x_i} \quad \text{in } \Omega.$$

In proving the existence of minimizers of the functionals F_0 and F_a , we use de la Vallée Poussin’s criterion [25] of the sequential compactness in $L^1(\Omega)$. The uniqueness of such minimizers follows basically from the convexity of these functionals. A crucial step in defining and deriving equilibriums of F_0 and F_a is the construction of

L^∞ -concentrations that are bounded below in Ω_s by a positive constant and that have low free energies. Such constructions are made by increasing the entropy of ionic concentrations through their small perturbations. The effect of inhomogeneous Dirichlet boundary data to the Boltzmann distributions and, hence, to the PBE can be useful to guide practical numerical computations. The equivalence of the two formulations is a common property for many physical problems. The interface formulation of the boundary-value problem of the PBE has been used for numerical computations using boundary integral method [22–24]. The finite-size effect is important in modeling electrostatics [3, 20].

The rest of the paper is organized as follows: In section 2, we state our main results; in section 3, we provide some lemmas; in section 4, we prove our theorems on the boundary-value problem of PBE; in section 5, we prove our theorems on the free-energy minimization. Finally, in Appendix, we give the proof of two lemmas.

2. Main results. Throughout the rest of the paper, we make the following assumptions:

- A1. The set $\Omega \subset \mathbb{R}^3$ is nonempty, bounded, open, and connected. The sets $\Omega_m \subset \mathbb{R}^3$ and $\Omega_s \subset \mathbb{R}^3$ are nonempty, bounded, and open, and satisfy that $\overline{\Omega_m} \subset \Omega$ and $\Omega_s = \Omega \setminus \overline{\Omega_m}$. The N points x_1, \dots, x_N for some integer $N \geq 1$ belong to Ω_m . Both $\partial\Omega$ and Γ are of C^2 . The unit exterior normal at the boundary of Ω_s is denoted by n ; cf. Figure 1.
- A2. $M \geq 2$ is an integer. All $a > 0$, $\beta > 0$, $Q_i \in \mathbb{R}$ ($1 \leq i \leq N$), $q_j \in \mathbb{R}$ and $\mu_j \in \mathbb{R}$ ($1 \leq j \leq M$), $\varepsilon_m > 0$, and $\varepsilon_s > 0$ are constants;
- A3. The functions ψ_{vac} and ε_Γ are defined in (1.4) and (1.7), respectively. The boundary data ψ_0 is the trace of a given function, also denoted by ψ_0 , in $W^{2,\infty}(\Omega)$.

Boundary values are understood as traces. When no confusion arises, the capital letter C , with or without a subscript, denotes a positive constant that can depend on all $\Omega_m, \Omega_s, \Omega, \Gamma, \varepsilon_m, \varepsilon_s, a, \beta, N, M, x_i$, and Q_i ($1 \leq i \leq N$), q_j and μ_j ($1 \leq j \leq M$), and ψ_0 .

For any open set $U \subseteq \mathbb{R}^3$ that contains all x_1, \dots, x_N , we denote

$$H_*^1(U) = \{u \in W^{1,1}(U) : u|_{U_\alpha} \in H^1(U_\alpha) \forall \alpha > 0\},$$

where $U_\alpha = U \setminus (\cup_{i=1}^N \overline{B(x_i, \alpha)})$ and $B(x_i, \alpha)$ denotes the ball centered at x_i with radius α .

DEFINITION 2.1. A function $\psi \in H_*^1(\Omega)$ is a weak solution to the boundary-value problem of the PBE (1.13) and (1.6), if $\psi = \psi_0$ on $\partial\Omega$, $\chi_{\Omega_s} B(\psi) \in L^2(\Omega_s)$, and

$$\int_{\Omega} \left[\varepsilon_\Gamma \nabla \psi \cdot \nabla \eta + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \eta \right] dx = 4\pi \sum_{i=1}^N Q_i \eta(x_i) \quad \forall \eta \in C_c^\infty(\Omega).$$

We remark that if $\phi \in H^1(U)$ for some bounded and smooth domain $U \subset \mathbb{R}^3$, then e^ϕ and, hence, $B(\phi)$ may not be in $L^1(U)$. For example, let $U = B(0, 1)$ be the unit ball of \mathbb{R}^3 and $\alpha \in (0, 1/2)$. Define $\phi(x) = |x|^{-\alpha}$ for any $x \in U$. Then $\phi \in H^1(U)$ and that $e^\phi \notin L^1(U)$. Notice by (1.14) that $\chi_{\Omega_s} B(\psi) \in L^2(\Omega_s)$ is equivalent to $\chi_{\Omega_s} e^{-\beta q_j \psi} \in L^2(\Omega_s)$ or $\chi_{\Omega_s} e^{-\beta q_j (\psi - \hat{\psi}_0/2)} \in L^2(\Omega_s)$ ($j = 1, \dots, M$), which in turn are equivalent to $\chi_{\Omega_s} B(\psi - \hat{\psi}_0/2) \in L^2(\Omega_s)$.

THEOREM 2.1. There exists a unique weak solution $\psi \in H_*^1(\Omega)$ to the boundary-value problem of the PBE (1.13) and (1.6). Moreover, $\psi \in C(\overline{\Omega} \setminus (\cup_{i=1}^N \overline{B(x_i, \alpha)}))$

for any $\alpha > 0$ such that the closure of $\cup_{i=1}^N B(x_i, \alpha)$ is contained in Ω_m , and $\psi \in C^\infty((\Omega_m \setminus \{x_1, \dots, x_N\}) \cup \Omega_s)$.

DEFINITION 2.2. A function $\psi : \Omega \rightarrow \mathbb{R}$ is a weak solution of the interface problem (1.15), if the following are satisfied: $\psi_m \in H_*^1(\Omega_m)$ and

$$(2.2) \quad \int_{\Omega_m} \varepsilon_m \nabla \psi \cdot \nabla \eta dx = 4\pi \sum_{i=1}^N Q_i \eta(x_i) \quad \forall \eta \in C_c^\infty(\Omega_m);$$

$\psi_s \in H^1(\Omega_s)$, $\chi_{\Omega_s} B(\psi) \in L^2(\Omega_s)$, and

$$(2.3) \quad \int_{\Omega_s} \left[\varepsilon_s \nabla \psi \cdot \nabla \eta + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \eta \right] dx = 0 \quad \forall \eta \in C_c^\infty(\Omega_s);$$

and the third and fourth equations in (1.15) hold true.

THEOREM 2.2. A function $\psi : \Omega \rightarrow \mathbb{R}$ is a weak solution to the boundary-value problems (1.13) and (1.6), if and only if it is a weak solution to the boundary-value problem (1.15).

Let $U \subset \mathbb{R}^3$ be a nonempty, bounded, and open set. Let $f \in L^1(U)$. Assume

$$(2.4) \quad \sup_{0 \neq \xi \in L^\infty(U) \cap H_0^1(U)} \frac{\int_U f \xi dx}{\|\xi\|_{H^1(U)}} < \infty.$$

Since $L^\infty(U) \cap H_0^1(U)$ is dense in $H_0^1(U)$, we can identify f as an element in $H^{-1}(U)$, the dual of $H_0^1(U)$, with

$$\langle f, \xi \rangle = \int_U f \xi dx \quad \forall \xi \in L^\infty(U) \cap H_0^1(U),$$

and we write $f \in L^1(U) \cap H^{-1}(U)$. The $H^{-1}(U)$ norm of f is given by (2.4). We define

$$X = \left\{ c = (c_1, \dots, c_M) \in L^1(\Omega, \mathbb{R}^M) : c = 0 \text{ a.e. } \Omega_m \text{ and } \sum_{j=1}^M q_j c_j \in H^{-1}(\Omega) \right\},$$

$$\|c\|_X = \sum_{j=1}^M \|c_j\|_{L^1(\Omega_s)} + \left\| \sum_{j=1}^M q_j c_j \right\|_{H^{-1}(\Omega)} \quad \forall c = (c_1, \dots, c_M) \in X.$$

Clearly, $(X, \|\cdot\|_X)$ is a Banach space.

Let $\alpha \in \mathbb{R}$ and define $S_\alpha : [0, \infty) \rightarrow \mathbb{R}$ by $S_\alpha(0) = 0$ and $S_\alpha(u) = u(\alpha + \log u)$ if $u > 0$. It is easy to see that S_α is bounded below on $[0, \infty)$ and strictly convex on $(0, \infty)$. Define

$$V_0 = \left\{ (c_1, \dots, c_M) \in X : c_j \geq 0 \text{ a.e. } \Omega_s \text{ and } \int_\Omega S_0(c_j) dx < \infty, j = 1, \dots, M \right\},$$

$$W_0 = \left\{ (c_1, \dots, c_M) \in V_0 : \text{there exists } p > \frac{3}{2} \text{ such that } c_j \in L^p(\Omega), j = 1, \dots, M \right\},$$

$$V_a = \left\{ (c_1, \dots, c_M) \in V_0 : c_0 = a^{-3} \left(1 - \sum_{j=1}^M a^3 c_j \right) \geq 0 \text{ a.e. } \Omega_s \right\}.$$

Clearly, all $V_0, W_0,$ and V_a are nonempty and convex. For any $c = (c_1, \dots, c_M) \in V_0,$ there exists a unique weak solution $\psi = \psi(c)$ of the boundary-value problem (1.5) and (1.6) with the charge density ρ given by (1.8); in particular, $\psi - \psi_{vac}$ is harmonic in $\Omega_m,$ cf. Lemma 3.2. We shall call $\psi = \psi(c)$ the electrostatic potential corresponding to $c.$ Therefore, $F_0 : V_0 \rightarrow \mathbb{R}$ and $F_a : V_a \rightarrow \mathbb{R}$ are well defined. We use V, F to denote V_0, F_0 or W_0, F_0 or $V_a, F_a.$

DEFINITION 2.3. *An element $c = (c_1, \dots, c_M) \in V$ is an equilibrium of $F : V \rightarrow \mathbb{R},$ if*

$$(2.5) \quad \text{there exist } \gamma_1 > 0 \text{ and } \gamma_2 > 0 \text{ such that } \gamma_1 \leq c_j(x) \leq \gamma_2 \text{ a.e. } x \in \Omega_s, \quad j = 1, \dots, M,$$

for the case of point ions, or

$$(2.6) \quad \text{there exists } \theta_0 \in (0, 1) \text{ such that } a^3 c_j(x) \geq \theta_0 \text{ a.e. } x \in \Omega_s, \quad j = 0, 1, \dots, M,$$

for the case of finite-size ions; and

$$\delta F[c]e := \lim_{t \rightarrow 0} \frac{F[c + te] - F[c]}{t} = 0 \quad \forall e \in X \cap L^\infty(\Omega, \mathbb{R}^M).$$

DEFINITION 2.4. *An element $c \in V$ is a local minimizer of $F : V \rightarrow \mathbb{R},$ if there exists $\varepsilon > 0$ such that $F[d] \geq F[c]$ for any $d \in V$ with $\|d - c\|_X < \varepsilon.$*

THEOREM 2.3. *There exists a unique minimizer of $F_0 : V_0 \rightarrow \mathbb{R}.$ It is also the unique local minimizer of $F_0 : V_0 \rightarrow \mathbb{R}.$*

It is an open question if the unique minimizer of $F_0 : V_0 \rightarrow \mathbb{R}$ is an equilibrium of $F_0 : V_0 \rightarrow \mathbb{R}$ as defined in Definition 2.3. The answer to this question would be yes if this minimizer were in W_0 or if $\min_{d \in V_0} F_0[d] = \min_{d \in W_0} F_0[d],$ neither of which is clearly true. This is the reason we introduce the class of concentrations $W_0.$ See the proof of Lemma 3.4 in Appendix.

THEOREM 2.4.

- (1) *There exists a unique equilibrium $c = (c_1, \dots, c_M)$ of $F_0 : W_0 \rightarrow \mathbb{R}.$ It is also the unique global minimizer and the unique local minimizer of $F_0 : W_0 \rightarrow \mathbb{R}.$*
- (2) *If $\psi = \psi(c)$ is the corresponding electrostatic potential, then the Boltzmann distributions (1.9) for point ions holds true and ψ is the unique weak solution to the corresponding boundary-value problem of PBE (1.11) and (1.6). Moreover,*

$$(2.7) \quad \begin{aligned} \min_{d \in W_0} F_0[d] &= \frac{1}{2} \sum_{i=1}^N Q_i (\psi - \psi_{vac})(x_i) + \int_{\Gamma} \frac{1}{8\pi} (\psi - \hat{\phi}_0) \varepsilon_{\Gamma} \partial_n (\psi - \hat{\psi}_0) dS \\ &\quad - \int_{\Omega_s} \frac{\varepsilon_s}{8\pi} |\nabla (\psi - \hat{\psi}_0)|^2 dx - \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} c_j^{\infty} e^{-\beta q_j (\psi - \hat{\psi}_0/2)} dx. \end{aligned}$$

THEOREM 2.5.

- (1) *There exists a unique equilibrium $c = (c_1, \dots, c_M)$ of $F_a : V_a \rightarrow \mathbb{R}.$ It is also the unique global minimizer and the unique local minimizer of $F_a : V_a \rightarrow \mathbb{R}.$*
- (2) *If $\psi = \psi(c)$ is the corresponding electrostatic potential, then the Boltzmann distributions (1.9) for finite-size ions holds true and ψ is the unique weak solution to the corresponding boundary-value problem of PBE (1.12) and (1.6).*

Moreover,

$$\begin{aligned}
 \min_{d \in V_a} F_a[d] &= \frac{1}{2} \sum_{i=1}^N Q_i(\psi - \psi_{vac})(x_i) + \int_{\Gamma} \frac{1}{8\pi} (\psi - \hat{\phi}_0) \varepsilon_{\Gamma} \partial_n (\psi - \hat{\psi}_0) dS \\
 (2.8) \quad &- \int_{\Omega_s} \frac{\varepsilon_{\Gamma}}{8\pi} |\nabla (\psi - \hat{\psi}_0)|^2 dx - \beta^{-1} a^{-3} \int_{\Omega_s} \\
 &\left[1 + \log \left(1 + a^3 \sum_{j=1}^M c_j^{\infty} e^{-\beta q_j (\psi - \hat{\psi}_0/2)} \right) \right] dx.
 \end{aligned}$$

3. Some lemmas. The key point of our first lemma below is the existence and continuity across the interface Γ of the normal flux for a solution of an elliptic interface problem. In terms of electrostatics, this means that the electrostatic potential and the normal component of electrostatic displacement are continuous across dielectric boundaries. These seem to be known results. For completeness, we give a proof here.

LEMMA 3.1. *Let $U \subset \mathbb{R}^3$ be an open set such that $\Gamma \subset U \subseteq \Omega$. Let $g \in L^1(U) \cap H^{-1}(U)$. Suppose $u \in H^1(U)$ satisfies*

$$(3.1) \quad \int_U \varepsilon_{\Gamma} \nabla u \cdot \nabla \eta dx = \int_U g \eta dx \quad \forall \eta \in C_c^{\infty}(U).$$

Then $[[u]] = 0$ on Γ . If in addition $g \in L^2(U)$, then $[[\varepsilon_{\Gamma} \partial_n u]] = 0$ on Γ .

Proof. Fix an open ball $B \subset U$ such that $\Gamma \cap B \neq \emptyset$. Let $\eta \in C_c^{\infty}(U)$ with $\text{supp } \eta \subset B$. Let n_j with $1 \leq j \leq 3$ be the j th component of n , the unit normal at the Γ , pointing from Ω_s to Ω_m . It follows from the fact that $u \in H^1(\Omega)$ and integration by parts that

$$\begin{aligned}
 - \int_B u \partial_j \eta dx &= \int_B (\partial_j u) \eta dx \\
 &= \int_{B \cap \Omega_m} (\partial_j u) \eta dx + \int_{B \cap \Omega_s} (\partial_j u) \eta dx \\
 &= - \int_{B \cap \Omega_m} u \partial_j \eta dx - \int_{\Gamma \cap B} u_m \eta n_j dS - \int_{B \cap \Omega_s} u \partial_j \eta dx + \int_{\Gamma \cap B} u_s \eta n_j dS \\
 &= - \int_B u \partial_j \eta dx + \int_{\Gamma \cap B} (u_s - u_m) \eta n_j dS, \quad j = 1, 2, 3.
 \end{aligned}$$

This and the arbitrariness of η imply $[[u]] = 0$ on Γ .

To show the continuity of $\varepsilon_{\Gamma} \nabla u \cdot n$ across Γ , we fix an open set $U_0 \subset \mathbb{R}^3$ such that $\Gamma \subset U_0 \subset \overline{U_0} \subset U$ and that the boundary ∂U_0 is C^2 . By the fact that $u \in H^1(U)$ and $g \in L^2(U)$, and by (3.1), we have

$$\begin{aligned}
 (\varepsilon_t \nabla u_t)|_{U_0 \cap \Omega_t} &\in L^2(U_0 \cap \Omega_t, \mathbb{R}^3) \quad \text{and} \\
 (\nabla \cdot \varepsilon_t \nabla u_t)|_{U_0 \cap \Omega_t} &= -g \in L^2(U_0 \cap \Omega_t), \quad t = m, s.
 \end{aligned}$$

Therefore, by Theorem 1.2 in [30], the trace of $(\varepsilon_t \nabla u_t)|_{U_0 \cap \Omega_t} \cdot \nu \in H^{-1/2}(\partial(U_0 \cap \Omega_t))$ exists, and also by (3.1),

$$(3.2) \quad \int_{U_0 \cap \Omega_t} \varepsilon_t \nabla u_t \cdot \nabla \eta dx = \int_{U_0 \cap \Omega_t} g \eta dx + \int_{\partial(U_0 \cap \Omega_t)} (\varepsilon_t \nabla u_t \cdot \nu) \eta dS \quad \forall \eta \in C_c^{\infty}(U_0 \cap \Omega_t),$$

where ν denotes the unit exterior normal of the boundary $\partial(U_0 \cap \Omega_t)$ which contains Γ and $t = m, s$. Notice that the normals ν at Γ from both sides $U_0 \cap \Omega_m$ and $U_0 \cap \Omega_s$ are in opposite directions.

These traces are determined independent of the choice of U_0 . In fact, if $Q_0 \subset \mathbb{R}^3$ is another open set such that $\Gamma \subset Q_0 \subset \overline{Q_0} \subset U$ and the boundary ∂Q_0 is C^2 , then the traces $(\varepsilon_m \nabla u_m)|_{Q_0 \cap \Omega_m} \cdot \nu \in H^{-1/2}(\partial(Q_0 \cap \Omega_m))$ and $(\varepsilon_s \nabla u_s)|_{Q_0 \cap \Omega_s} \cdot \nu \in H^{-1/2}(\partial(Q_0 \cap \Omega_s))$ exist, and (3.2) holds true for $t = m, s$ when U_0 is replaced by Q_0 . Consider now (3.2) with $t = m$. Choose any $\eta \in C_c^1(U_0 \cap Q_0)$ such that $\eta = 0$ on $\partial(U_0 \cap \Omega_m) \setminus \Gamma$ and on $\partial(Q_0 \cap \Omega_m) \setminus \Gamma$. Extend η by $\eta = 0$ to outside $U_0 \cap Q_0$. By (3.2) with $t = m$ and the corresponding equation with U_0 replaced by Q_0 , we obtain that

$$\int_{\partial(U_0 \cap \Omega_m)} (\varepsilon_m \nabla u_m \cdot \nu) \eta dS = \int_{\partial(Q_0 \cap \Omega_m)} (\varepsilon_m \nabla u_m \cdot \nu) \eta dS.$$

The arbitrariness of η then implies that the trace of $(\varepsilon_m \nabla u_m \cdot \nu)|_{U_0 \cap \Omega_m}$ on Γ determined by U_0 is the same as that determined by Q_0 . By the same argument, we see that the trace of $(\varepsilon_s \nabla u_s \cdot n)|_{U_0 \cap \Omega_s}$ on Γ determined by U_0 is the same as that determined by Q_0 .

Now, by the fact that $U_0 = (U_0 \cap \Omega_m) \cup (U_0 \cap \Omega_s)$ and $(U_0 \cap \Omega_m) \cap (U_0 \cap \Omega_s) = \emptyset$, and by our convention for the direction of the unit normal n along Γ , we obtain from (3.1) and (3.2) that for any $\eta \in C_c^\infty(U)$ with $\text{supp } \eta \subset U_0$

$$\int_{\Gamma} \left(\varepsilon_m \frac{\partial u_m}{\partial n} - \varepsilon_s \frac{\partial u_s}{\partial n} \right) \eta dS = 0.$$

The arbitrariness of η implies $[\varepsilon_\Gamma \partial_n u] = 0$ on Γ . □

Let $L : H^{-1}(\Omega) \rightarrow H^1(\Omega)$ be the linear operator defined as follows: for any $\xi \in H^{-1}(\Omega)$, $L\xi \in H_0^1(\Omega)$ is the unique function in $H_0^1(\Omega)$ that satisfies

$$(3.3) \quad \frac{1}{4\pi} \int_{\Omega} \varepsilon_\Gamma \nabla(L\xi) \cdot \nabla v \, dx = \xi(v) \quad \forall v \in H_0^1(\Omega).$$

It is easy to see that $\langle \xi, \eta \rangle = \xi(L\eta)$ defines an inner product of $H^{-1}(\Omega)$. Denote by $\|\cdot\|$ the corresponding norm of $H^{-1}(\Omega)$, i.e., $\|\xi\| = \sqrt{\langle \xi, \xi \rangle} = \sqrt{\xi(L\xi)}$ for any $\xi \in H^{-1}(\Omega)$. One can verify that there exist $C_1 = C_1(\Omega, \varepsilon_m, \varepsilon_s) > 0$ and $C_2 = C_2(\varepsilon_m, \varepsilon_s) > 0$ such that

$$(3.4) \quad C_1 \|\xi\| \leq \|\xi\|_{H^{-1}(\Omega)} \leq C_2 \|\xi\| \quad \forall \xi \in H^{-1}(\Omega).$$

It follows from [21] (with minor modifications) that there exists a unique $G \in H_*^1(\Omega)$ such that $G = 0$ on $\partial\Omega$ and

$$(3.5) \quad \int_{\Omega} \varepsilon_\Gamma \nabla G \cdot \nabla \eta \, dx = 4\pi \sum_{i=1}^N Q_i \eta(x_i) \quad \forall \eta \in C_c^\infty(\Omega).$$

Clearly, $G - \psi_{vac}$ is harmonic in Ω_m and $G \in W^{1,p}(\Omega)$ for any $p \in [1, 3/2)$. Notice that the function $\hat{\psi}_0 \in H^1(\Omega)$ defined in (1.10) is harmonic in $\Omega_m \cup \Omega_s$.

The next lemma gives a solution decomposition of the Poisson equation (1.5) with its right-hand side consisting of Dirac masses and a function in $H^{-1}(\Omega)$ that represents the density of ionic charges. This decomposition is a mathematical formulation of the Born cycle [2].

LEMMA 3.2. *Let $f \in L^1(\Omega) \cap H^{-1}(\Omega)$ be such that $f = 0$ in Ω_m . Then $\psi := G + \hat{\psi}_0 + Lf$ is the unique function in $H_*^1(\Omega)$ that satisfies $\psi = \psi_0$ on $\partial\Omega$ and*

$$(3.6) \quad \int_{\Omega} \varepsilon_{\Gamma} \nabla \psi \cdot \nabla \eta dx = 4\pi \sum_{i=1}^N Q_i \eta(x_i) + 4\pi \int_{\Omega_s} f \eta dx \quad \forall \eta \in C_c^{\infty}(\Omega).$$

Moreover, Lf and $\psi - \psi_{vac}$ are harmonic in Ω_m , and

$$(3.7) \quad \sum_{i=1}^N Q_i (Lf)(x_i) = \int_{\Omega_s} G f dx.$$

Proof. From the definition of $G, \hat{\psi}_0$, and L (cf. (3.5), (1.10), and (3.3)), we easily verify that the function ψ is in $H_*^1(\Omega)$, $\psi = \psi_0$ on $\partial\Omega$, and (3.6) holds true. If $\bar{\psi} \in H_*^1(\Omega)$ satisfies $\bar{\psi} = 0$ on $\partial\Omega$ and $\int_{\Omega} \varepsilon_{\Gamma} \nabla \bar{\psi} \cdot \nabla \eta dx = 0$ for all $\eta \in C_c^{\infty}(\Omega)$, then clearly $\bar{\psi} \in H_0^1(\Omega)$ and in fact $\bar{\psi} = 0$ a.e. Ω . This proves the needed uniqueness.

By the fact that $f = 0$ in Ω_m and the definition of L (cf. (3.3)), Lf is harmonic in Ω_m . The fact that $\psi - \psi_{vac}$ is harmonic in Ω_m follows from (3.6) with $\eta \in C_c^{\infty}(\Omega)$ so chosen that $\text{supp } \eta \subset \Omega_m$ and

$$\int_{\Omega_m} \varepsilon_m \nabla \psi_{vac} \cdot \nabla \eta dx = 4\pi \sum_{i=1}^N Q_i \eta(x_i) \quad \forall \eta \in C_c^{\infty}(\Omega_m).$$

It remains to prove (3.7). Denote $\psi_c = Lf \in H_0^1(\Omega)$. Let $\alpha > 0$ be sufficiently small and let $B_{\alpha} = \cup_{i=1}^N B(x_i, \alpha)$. By the fact that G is harmonic in $\Omega_m \setminus B_{\alpha}$ and $G - \psi_{vac}$ is harmonic in Ω_m , we obtain by a series of routine calculations that

$$\begin{aligned} \int_{\Omega_m} \varepsilon_m \nabla G \cdot \nabla \psi_c dx &= \int_{\Omega_m \setminus B_{\alpha}} \varepsilon_m \nabla G \cdot \nabla \psi_c dx + \int_{B_{\alpha}} \varepsilon_m \nabla G \cdot \nabla \psi_c dx \\ &= - \int_{\Omega_m \setminus B_{\alpha}} \varepsilon_m (\Delta G) \psi_c dx + \int_{\partial(\Omega_m \setminus B_{\alpha})} \varepsilon_m \psi_c \frac{\partial G}{\partial \nu} dS + O(\alpha) \\ &= - \int_{\Gamma} \varepsilon_m \psi_c |m| \frac{\partial G|_m}{\partial n} dS + \int_{\partial B_{\alpha}} \varepsilon_m \psi_c \frac{\partial G}{\partial \nu} dS + O(\alpha) \\ &= - \int_{\Gamma} \varepsilon_m \psi_c |m| \frac{\partial G|_m}{\partial n} dS + \int_{\partial B_{\alpha}} \varepsilon_m \psi_c \frac{\partial (G - \psi_{vac})}{\partial \nu} dS \\ &\quad + \sum_{i=1}^N \int_{\partial B(x_i, \alpha)} \varepsilon_m \psi_c \frac{\partial \psi_{vac}}{\partial \nu} dS + O(\alpha) \\ &\rightarrow - \int_{\Gamma} \varepsilon_m \psi_c |m| \frac{\partial G|_m}{\partial n} dS + 4\pi \sum_{i=1}^N Q_i \psi_c(x_i) \quad \text{as } \alpha \rightarrow 0, \end{aligned}$$

where ν is the exterior unit normal of $\partial(\Omega_m \setminus B_{\alpha})$ and $\nu = -n$ on Γ by our convention for the direction of n . Consequently, by the continuity of $\varepsilon_{\Gamma} \nabla G \cdot n$ across Γ (cf.

Lemma 3.1), the fact that G is harmonic in Ω_s , and $G = \psi_c = 0$ on $\partial\Omega$, we obtain

$$\begin{aligned}
 4\pi \sum_{i=1}^N Q_i \psi_c(x_i) &= \int_{\Gamma} \varepsilon_s \psi_c|_s \frac{\partial G|_s}{\partial n} dS + \int_{\Omega_m} \varepsilon_m \nabla G \cdot \nabla \psi_c dx \\
 &= \int_{\Omega_s} \varepsilon_s (\Delta G) \psi_c dx + \int_{\Omega_s} \varepsilon_s \nabla G \cdot \nabla \psi_c dx + \int_{\Omega_m} \varepsilon_m \nabla G \cdot \nabla \psi_c dx \\
 (3.8) \quad &= \int_{\Omega} \varepsilon_{\Gamma} \nabla G \cdot \nabla \psi_c dx.
 \end{aligned}$$

Since $\psi_c = Lf \in H_0^1(\Omega)$ is harmonic in Ω_m , we also have by the properties of G (cf. [21]) and integration by parts that

$$\begin{aligned}
 \int_{\Omega_m} \varepsilon_m \nabla G \cdot \nabla \psi_c dx &= \int_{\Omega_m \setminus B_\alpha} \varepsilon_m \nabla G \cdot \nabla \psi_c dx + \int_{B_\alpha} \varepsilon_m \nabla G \cdot \nabla \psi_c dx \\
 &= - \int_{\Omega_m \setminus B_\alpha} \varepsilon_m G \Delta \psi_c dx + \int_{\partial(\Omega_m \setminus B_\alpha)} \varepsilon_m G \frac{\partial \psi_c}{\partial \nu} dS + O(\alpha) \\
 &= - \int_{\Gamma} \varepsilon_m G|_m \frac{\partial \psi_c|_m}{\partial n} dS + \int_{\partial B_\alpha} \varepsilon_m G \frac{\partial \psi_c}{\partial \nu} dS + O(\alpha) \\
 (3.9) \quad &\rightarrow - \int_{\Gamma} \varepsilon_m G|_m \frac{\partial \psi_c|_m}{\partial n} dS \quad \text{as } \alpha \rightarrow 0.
 \end{aligned}$$

Let $\hat{G} \in H^1(\Omega_m)$ be such that $\hat{G} = G$ on $\Gamma = \partial\Omega_m$. Replacing G in (3.9) by \hat{G} and repeating the same calculations, we obtain

$$(3.10) \quad \int_{\Omega_m} \varepsilon_m \nabla G \cdot \nabla \psi_c dx = \int_{\Omega_m} \varepsilon_m \nabla \hat{G} \cdot \nabla \psi_c dx.$$

Define $\overline{G} : \Omega \rightarrow \mathbb{R}$ by $\overline{G}(x) = \hat{G}(x)$ if $x \in \overline{\Omega_m}$ and by $\overline{G}(x) = G(x)$ if $x \in \Omega_s$. Clearly, $\overline{G} \in H_0^1(\Omega)$. Since $\psi_c = Lf$ and $f = 0$ in Ω_m , we, thus, have by (3.8) and (3.10) that

$$\begin{aligned}
 4\pi \sum_{i=1}^N Q_i \psi_c(x_i) &= \int_{\Omega_m} \varepsilon_m \nabla \hat{G} \cdot \nabla \psi_c dx + \int_{\Omega_s} \varepsilon_s \nabla G \cdot \nabla \psi_c dx \\
 &= \int_{\Omega} \varepsilon_{\Gamma} \nabla \overline{G} \cdot \nabla \psi_c dx = 4\pi \int_{\Omega} \overline{G} f dx = 4\pi \int_{\Omega_s} G f dx.
 \end{aligned}$$

This implies (3.7). \square

By Lemma 3.2, the potential $\psi = \psi(c_1, \dots, c_M)$ corresponding to a set of concentrations (c_1, \dots, c_M) is well defined with $f = \sum_{j=1}^M q_j c_j$ and is given by

$$(3.11) \quad \psi(c_1, \dots, c_M) = G + \hat{\psi}_0 + L \left(\sum_{j=1}^M q_j c_j \right).$$

Moreover, the functional $F_0 : V_0 \rightarrow \mathbb{R}$ and $F_a : V_a \rightarrow \mathbb{R}$ can be rewritten as

$$\begin{aligned}
 F_0[c] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx + \sum_{j=1}^M \int_{\Omega_s} \mu_{0j} c_j dx \\
 (3.12) \quad &+ \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} S_{-1}(c_j) dx + E_0, \quad \forall c = (c_1, \dots, c_M) \in V_0,
 \end{aligned}$$

$$\begin{aligned}
 F_a[c] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx + \sum_{j=1}^M \int_{\Omega_s} \mu_{aj} c_j dx \\
 (3.13) \quad &+ \beta^{-1} \sum_{j=0}^M \int_{\Omega_s} S_{-1}(c_j) dx + E_a \quad \forall c = (c_1, \dots, c_M) \in V_a,
 \end{aligned}$$

respectively, where

$$(3.14) \quad \mu_{0j}(x) = q_j G(x) + \frac{1}{2} q_j \hat{\psi}_0(x) + 3\beta^{-1} \log a - \mu_j \quad \forall x \in \Omega_s, \quad j = 1, \dots, M,$$

$$(3.15) \quad E_0 = \frac{1}{2} \sum_{i=1}^N Q_i \left(G + \hat{\psi}_0 - \psi_{vac} \right) (x_i),$$

$$(3.16) \quad \mu_{aj}(x) = q_j G(x) + \frac{1}{2} q_j \hat{\psi}_0(x) - \mu_j \quad \forall x \in \Omega_s, \quad j = 1, \dots, M,$$

$$(3.17) \quad E_a = \frac{1}{2} \sum_{i=1}^N Q_i \left(G + \hat{\psi}_0 - \psi_{vac} \right) (x_i) + 3\beta^{-1} a^{-3} (\log a) |\Omega_s|,$$

where $|E|$ denotes the Lebesgue measure of a Lebesgue measurable set $E \subset \mathbb{R}^3$.

LEMMA 3.3. *Let $D \subset \mathbb{R}^3$ be a bounded and open set. Let $\alpha \in \mathbb{R}$. Let $\{u^{(k)}\}$ be a sequence of functions in $L^1(D)$ such that $u^{(k)} \geq 0$ a.e. D for each $k \geq 1$ and that*

$$\sup_{k \geq 1} \int_D S_\alpha \left(u^{(k)} \right) dx < \infty.$$

Then there exists a subsequence $\{u^{(k_j)}\}$ of $\{u^{(k)}\}$ such that $\{u^{(k_j)}\}$ converges weakly in $L^1(D)$ to some $u \in L^1(D)$ with $u \geq 0$ a.e. D and

$$\int_D S_\alpha(u) dx \leq \liminf_{k \rightarrow \infty} \int_D S_\alpha \left(u^{(k)} \right) dx.$$

Proof. Since $S_\alpha : [0, \infty) \rightarrow \mathbb{R}$ is bounded below, by passing to a subsequence if necessary, we may assume that the limit

$$(3.18) \quad A := \lim_{k \rightarrow \infty} \int_D S_\alpha \left(u^{(k)} \right) dx = \liminf_{k \rightarrow \infty} \int_D S_\alpha \left(u^{(k)} \right) dx$$

exists and is finite. Since $S_\alpha(\lambda)/\lambda \rightarrow +\infty$ as $\lambda \rightarrow +\infty$, $\{u^{(k)}\}$ is weakly sequentially compact in $L^1(D)$ by de la Vallée Poussin’s criterion [25]. Therefore, this sequence has a subsequence, not relabeled, that converges weakly in $L^1(D)$ to some $u \in L^1(D)$. Clearly, $u \geq 0$ a.e. D .

Let $\varepsilon > 0$. By (3.18), there exists an integer $K > 0$ such that

$$(3.19) \quad \int_D S_\alpha(u^{(k)}) dx \leq A + \varepsilon \quad \forall k > K.$$

By Mazur’s theorem [7, 32], there exist convex combinations $v^{(k)}$ of $u^{(K+1)}, \dots, u^{(K+k)}$ for all $k \geq 1$ such that $v^{(k)} \rightarrow u$ in $L^1(D)$. Let $v^{(k)} = \sum_{j=1}^k \lambda_{k,j} u^{(K+j)}$ with $\lambda_{k,j} \geq 0$ for all j and k , and $\sum_{j=1}^k \lambda_{k,j} = 1$ for all k . Since $S_\alpha : [0, \infty) \rightarrow \mathbb{R}$ is convex, we have by Jensen’s inequality and (3.19) that

$$(3.20) \quad S_\alpha(v^{(k)}) \leq \sum_{j=1}^k \lambda_{k,j} S_\alpha(u^{(K+j)}) \leq \sum_{j=1}^k \lambda_{k,j} (A + \varepsilon) = A + \varepsilon \quad \forall k \geq 1.$$

Since $v^{(k)} \rightarrow u$ in $L^1(D)$, there exists a subsequence $\{v^{(k_j)}\}$ of $\{v^{(k)}\}$ such that $v^{(k_j)}(x) \rightarrow u(x)$ a.e. $x \in D$. Consequently, since $S_\alpha : [0, \infty) \rightarrow \mathbb{R}$ is continuous and bounded below, we have by Fatou’s lemma and (3.20) that

$$\int_D S_\alpha(u(x)) dx = \int_D \lim_{j \rightarrow \infty} S_\alpha(v^{(k_j)}(x)) dx \leq \liminf_{j \rightarrow \infty} \int_D S_\alpha(v^{(k_j)}(x)) dx \leq A + \varepsilon,$$

concluding the proof by the arbitrariness of $\varepsilon > 0$. \square

The next two lemmas state some boundedness of concentrations that have low free energies. Their proofs are somewhat tedious, and are given in Appendix A.

LEMMA 3.4. *Let $c = (c_1, \dots, c_M) \in W_0$ satisfy that $c \notin L^\infty(\Omega, \mathbb{R}^M)$ or there exists $j \in \{1, \dots, M\}$ with $|\{x \in \Omega_s : c_j(x) < \alpha\}| > 0$ for all $\alpha > 0$. Then for any $\varepsilon > 0$ there exist $\hat{c} = (\hat{c}_1, \dots, \hat{c}_M) \in W_0$ that satisfies (2.5) with c replaced by \hat{c} , $\|\hat{c} - c\|_X < \varepsilon$, and $F_0[\hat{c}] < F_0[c]$.*

LEMMA 3.5. *Let $c = (c_1, \dots, c_M) \in V_a$ and c_0 be defined by (1.3). Assume there exists $j \in \{0, 1, \dots, M\}$ such that $|\{x \in \Omega_s : a^3 c_j(x) < \alpha\}| > 0$ for all $\alpha > 0$. Let $\varepsilon > 0$. Then there exists $\hat{c} = (\hat{c}_1, \dots, \hat{c}_M) \in V_a$ that satisfies (2.6) with c replaced by \hat{c} , $\|\hat{c} - c\|_X < \varepsilon$, and $F_a[\hat{c}] < F_a[c]$.*

4. The Poisson–Boltzmann equation: Proof of Theorems 2.1 and 2.2.

Proof of Theorem 2.1. It is easy to verify that the function $B : \mathbb{R} \rightarrow \mathbb{R}$ defined in (1.14) is convex for both the case of point ions and that of finite-size ions. Let

$$\mathcal{K} := \{u \in H^1(\Omega) : u = \psi_0 \text{ on } \partial\Omega \text{ and } \chi_{\Omega_s} B(u) \in L^2(\Omega)\}.$$

Clearly, $\mathcal{K} \neq \emptyset$ since $\psi_0 \in \mathcal{K}$ and \mathcal{K} is convex since $B : \mathbb{R} \rightarrow \mathbb{R}$ is convex. We show now that \mathcal{K} is closed in $H^1(\Omega)$. Let $u_k \in \mathcal{K}$ ($k = 1, 2, \dots$) and $u_k \rightarrow u$ in $H^1(\Omega)$

for some $u \in H^1(\Omega)$. Clearly, $u = \psi_0$ on $\partial\Omega$. Up to a subsequence, not relabeled, $u_k(x) \rightarrow u(x)$ a.e. $x \in \Omega$. Since $B : \mathbb{R} \rightarrow \mathbb{R}$ is convex and positive, we have

$$\frac{d^2}{dv^2} ([B(v)]^2) = 2[B'(v)]^2 + 2B(v)B''(v) > 0 \quad \forall v \in \mathbb{R}.$$

Thus, $v \mapsto [B(v)]^2$ is convex. It then follows from Fatou's lemma, Jensen's inequality, and the $H^1(\Omega)$ -boundedness of $\{u_k\}$ that

$$\begin{aligned} \frac{1}{|\Omega_s|} \int_{\Omega_s} [B(u)]^2 dx &\leq \liminf_{k \rightarrow \infty} \frac{1}{|\Omega_s|} \int_{\Omega_s} [B(u_k)]^2 dx \\ &\leq \liminf_{k \rightarrow \infty} \left[B \left(\frac{1}{|\Omega_s|} \int_{\Omega_s} u_k dx \right) \right]^2 < \infty. \end{aligned}$$

This implies that $u \in \mathcal{K}$. Therefore, \mathcal{K} is closed in $H^1(\Omega)$. Since \mathcal{K} is convex, it is also weakly closed in $H^1(\Omega)$.

Define now $J : \mathcal{K} \rightarrow \mathbb{R}$ by

$$J[u] = \int_{\Omega} \left[\frac{\varepsilon_{\Gamma}}{2} |\nabla u|^2 + 4\pi \chi_{\Omega_s} B \left(u + G - \frac{\hat{\psi}_0}{2} \right) \right] dx \quad \forall u \in \mathcal{K},$$

where G and $\hat{\psi}_0$ are defined in (3.5) and (1.10), respectively. Note that $\psi_0 \in \mathcal{K}$ and that $J[\psi_0] < \infty$. By the Poincaré inequality, there exist constants $C_3 > 0$ and $C_4 \geq 0$ such that $J[u] \geq C_3 \|u\|_{H^1(\Omega)}^2 - C_4$ for all $u \in \mathcal{K}$. Thus, $\alpha := \inf_{u \in \mathcal{K}} J[u]$ is finite. Let $v_k \in \mathcal{K}$ ($k = 1, 2, \dots$) be such that $\lim_{k \rightarrow \infty} J[v_k] = \alpha$. Then, $\{v_k\}$ is bounded in $H^1(\Omega)$ and, hence, it has a subsequence, not relabeled, that weakly converges to some $v \in H^1(\Omega)$. Since \mathcal{K} is weakly closed, $v \in \mathcal{K}$. Since the embedding $H^1(\Omega) \hookrightarrow L^2(\Omega)$ is compact, up to a further subsequence, again not relabeled, $v_k \rightarrow v$ a.e. in Ω . Therefore, since $B : \mathbb{R} \rightarrow \mathbb{R}$ is continuous and nonnegative, Fatou's lemma implies

$$\liminf_{k \rightarrow \infty} \int_{\Omega} \chi_{\Omega_s} B \left(v_k + G - \frac{\hat{\psi}_0}{2} \right) dx \geq \int_{\Omega} \chi_{\Omega_s} B \left(v + G - \frac{\hat{\psi}_0}{2} \right) dx.$$

Since $u \mapsto \int_{\Omega} \varepsilon_{\Gamma} |\nabla u|^2 dx$ is convex and $H^1(\Omega)$ continuous, it is sequentially weakly lower semicontinuous. Consequently, $\liminf_{k \rightarrow \infty} J[v_k] \geq J[v]$. Thus, v is a minimizer of $J : \mathcal{K} \rightarrow \mathbb{R}$.

Notice that $\chi_{\Omega_s} B'(v + G - \hat{\psi}_0/2) \in L^2(\Omega_s)$. Simple calculations of the first variation of $J : \mathcal{K} \rightarrow \mathbb{R}$ at any $\eta \in C_c^\infty(\Omega)$ leads to

$$\int_{\Omega} \left[\varepsilon_{\Gamma} \nabla v \cdot \nabla \eta + 4\pi \chi_{\Omega_s} B' \left(v + G - \frac{\hat{\psi}_0}{2} \right) \eta \right] dx = 0 \quad \forall \eta \in C_c^\infty(\Omega).$$

The function $\psi = v + G$ is, thus, a needed solution.

We now prove the uniqueness. Let ϕ be another weak solution. Let $\xi = \psi - \phi$. Then, $\xi \in H_*^1(\Omega)$, $\xi = 0$ on $\partial\Omega$, and

$$\int_{\Omega} \left\{ \varepsilon_{\Gamma} \nabla \xi \cdot \nabla \eta + 4\pi \chi_{\Omega_s} \left[B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) - B' \left(\phi - \frac{\hat{\psi}_0}{2} \right) \right] \eta \right\} dx = 0 \quad \forall \eta \in C_c^\infty(\Omega).$$

Choosing the test functions $\eta \in C_c^\infty(\Omega)$ so that $\text{supp } \eta \subset \Omega_m$, we find that ξ is harmonic in Ω_m . This and the fact that $\xi \in H_*^1(\Omega)$ imply that $\xi \in H_0^1(\Omega)$. Thus, the above test functions η can be chosen from $H_0^1(\Omega)$. In particular, setting $\eta = \xi$ and using the convexity of $B : \mathbb{R} \rightarrow \mathbb{R}$, we obtain that $\xi = 0$ and, hence, $\psi = \phi$ in $H^1(\Omega)$.

Let $\sigma > 0$ be such that the closure of $B_\sigma := \cup_{i=1}^N B(x_i, \sigma)$ is contained in Ω_m . Clearly, the unique weak solution $\psi \in H_*^1(\Omega)$ satisfies

$$(4.1) \quad \int_{\Omega \setminus \overline{B_\sigma}} \varepsilon_\Gamma \nabla \psi \cdot \nabla \eta \, dx = \int_{\Omega \setminus \overline{B_\sigma}} g \eta \, dx \quad \forall \eta \in C_c^\infty(\Omega \setminus \overline{B_\sigma}),$$

where $g = -4\pi \chi_{\Omega_s} B'(\psi - \hat{\psi}_0/2) \in L^2(\Omega \setminus \overline{B_\sigma})$. Therefore, $\psi \in C(\overline{\Omega} \setminus B_\sigma)$ by the standard regularity theory [14]. Since $\varepsilon_\Gamma = \varepsilon_m$ in Ω_m and $\varepsilon_\Gamma = \varepsilon_s$ in Ω_s , ψ is harmonic in $\Omega_m \setminus \overline{B_\sigma}$. Hence, $\psi \in C^\infty(\Omega_m \setminus \overline{B_\sigma})$. Notice that $B \in C^\infty(\mathbb{R})$; thus, we have $\psi \in C^\infty(\Omega_s)$ by a standard bootstrapping argument. \square

Proof of Theorem 2.2. Let $\psi \in H_*^1(\Omega)$ be a weak solution to the boundary-value problem (1.13) and (1.6). Clearly, $\psi_m \in H_*^1(\Omega_m)$. For any $\eta \in C_c^\infty(\Omega_m)$, we extend η to the entire Ω by defining $\eta = 0$ outside Ω_m . Then, we obtain (2.2) from (2.1). Since all $x_i \in \Omega_m$ ($i = 1, \dots, N$), we have $\psi_s \in H^1(\Omega_s)$. Since $\chi_{\Omega_s} B(\psi) \in L^2(\Omega_s)$, it follows from (1.14) that $\chi_{\Omega_s} B'(\psi) \in L^2(\Omega_s)$. For any $\eta \in C_c^\infty(\Omega_s)$, we, again, extend η to Ω by defining $\eta = 0$ outside Ω_s . Then, we obtain (2.3) from (2.1). Finally, by Lemma 3.1, (4.1), and (1.6), the last two equations in (1.15) hold true. Hence, ψ is a weak solution to (1.15).

Now let $\psi : \Omega \rightarrow \mathbb{R}$ be a solution to the boundary-value problem (1.15). We first show that $\psi \in H_*^1(\Omega)$. Let $\sigma > 0$ be small enough so that the closure of $B_\sigma := \cup_{i=1}^N B(x_i, \sigma)$ is contained in Ω_m . Since $\psi_m \in H_*^1(\Omega_m)$, $\psi_m \in H^1(\Omega_m \setminus \overline{B_\sigma})$. Thus, the trace $\psi_m|_\Gamma \in L^2(\Gamma)$, and is independent on the choice of σ . Similarly, $\psi_s \in H^1(\Omega_s)$, and, hence, $\psi_s|_\Gamma \in L^2(\Gamma)$. Fix $j \in \{1, 2, 3\}$. Define $\xi_j : \Omega \setminus \overline{B_\sigma} \rightarrow \mathbb{R}$ by $\xi_j = \partial_j \psi_m$ in $\Omega_m \setminus \overline{B_\sigma}$ and $\xi_j = \partial_j \psi_s$ in Ω_s . Clearly, $\xi_j \in L^2(\Omega \setminus \overline{B_\sigma})$. Let n_j be the j th component of the unit exterior normal n at Γ , pointing from Ω_s to Ω_m . Then, for any $\eta \in C_c^\infty(\Omega \setminus \overline{B_\sigma})$, we have

$$\begin{aligned} \int_{\Omega \setminus \overline{B_\sigma}} \xi_j \eta \, dx &= \int_{\Omega_m \setminus \overline{B_\sigma}} (\partial_j \psi_m) \eta \, dx + \int_{\Omega_s} (\partial_j \psi_s) \eta \, dx \\ &= - \int_{\Omega_m \setminus \overline{B_\sigma}} \psi \partial_j \eta \, dx - \int_\Gamma \psi_m \eta n_j \, dS - \int_{\Omega_s} \psi \partial_j \eta \, dx + \int_\Gamma \psi_s \eta n_j \, dS \\ &= - \int_{\Omega \setminus \overline{B_\sigma}} \psi \partial_j \eta \, dx + \int_\Gamma [\![\psi]\!] n_j \eta \, dS \\ &= - \int_{\Omega \setminus \overline{B_\sigma}} \psi \partial_j \eta \, dx, \end{aligned}$$

where in the last step we used the fact that $[\![\psi]\!] = 0$ on Γ . Thus, $\xi_j = \partial_j \psi \in L^2(\Omega \setminus \overline{B_\sigma})$, and, hence, $\psi \in H_*^1(\Omega)$ by the arbitrariness of $\sigma > 0$.

Clearly, $\chi_{\Omega_s} B'(\psi) \in L^2(\Omega_s)$ and $\psi = \psi_0$ on $\partial\Omega$. It remains to show that (2.1) holds true. Let $\eta \in C_c^\infty(\Omega)$. Let V_1 and V_2 be two open sets in \mathbb{R}^3 such that ∂V_1 and ∂V_2 are of C^2 , $x_i \notin V_2$ for $i = 1, \dots, N$, and $\Gamma \subset V_1 \subset \overline{V_1} \subset V_2 \subset \overline{V_2} \subset \Omega$. Let

$\zeta \in C_c^\infty(\Omega)$ be such that $\text{supp } \zeta \subset V_2$ and $\zeta = 1$ on V_1 . Then, $(1 - \zeta)\eta|_{\Omega_m} \in C_c^\infty(\Omega_m)$, $(1 - \zeta)\eta|_{\Omega_s} \in C_c^\infty(\Omega_s)$, and $(1 - \zeta(x_i))\eta(x_i) = \eta(x_i)$, $i = 1, \dots, N$. We, thus, have by (2.2) and (2.3) that

$$\begin{aligned}
 & \int_{\Omega} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla \eta + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \eta \right] dx \\
 &= \int_{\Omega} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla ((1 - \zeta)\eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) (1 - \zeta)\eta \right] dx \\
 &\quad + \int_{\Omega} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla (\zeta\eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta\eta \right] dx \\
 &= \int_{\Omega_m} \varepsilon_{\Gamma} \nabla \psi \cdot \nabla ((1 - \zeta)\eta) dx \\
 &\quad + \int_{\Omega_s} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla ((1 - \zeta)\eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) (1 - \zeta)\eta \right] dx \\
 &\quad + \int_{V_2} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla (\zeta\eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta\eta \right] dx \\
 (4.2) \quad &= 4\pi \sum_{i=1}^N Q_i \eta(x_i) + \int_{V_2} \left[\varepsilon_{\Gamma} \nabla \psi \cdot \nabla (\zeta\eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta\eta \right] dx.
 \end{aligned}$$

We now show that the second term in (4.2) is zero. Notice that $\psi|_{V_2} \in H^1(V_2)$ and $x_i \notin V_2$ ($1 \leq i \leq N$). Denoting $V_m = V_2 \cap \Omega_m$ and $V_s = V_2 \cap \Omega_s$, we have by (2.2) and (2.3) that

$$\begin{aligned}
 & \int_{V_m} \varepsilon_m \nabla \psi \cdot \nabla \xi \, dx = 0 \quad \forall \xi \in C_c^\infty(V_m), \\
 & \int_{V_s} \varepsilon_s \nabla \psi \cdot \nabla \xi \, dx = -4\pi \int_{V_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \xi \, dx \quad \forall \xi \in C_c^\infty(V_s).
 \end{aligned}$$

Consequently, since $\chi_{\Omega_s} B'(\psi - \hat{\psi}_0/2) \in L^2(\Omega_s)$, we infer from the regularity theory of elliptic boundary-value problems [14] that $\psi|_{V_m} \in H^2(V_m)$ and $\psi|_{V_s} \in H^2(V_s)$, and that

$$\begin{aligned}
 & \nabla \cdot \varepsilon_m \nabla \psi = 0 \quad \text{a.e. } V_m, \\
 & \nabla \cdot \varepsilon_s \nabla \psi - 4\pi B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) = 0 \quad \text{a.e. } V_s.
 \end{aligned}$$

Therefore, the trace of $\varepsilon_m \nabla \psi \cdot n$ and that of $\varepsilon_s \nabla \psi \cdot n$ on Γ both exist. Moreover,

$$\begin{aligned}
 & \int_{V_2} \left[\varepsilon_\Gamma \nabla \psi \cdot \nabla(\zeta \eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta \eta \right] dx \\
 &= \int_{V_m} \varepsilon_m \nabla \psi \cdot \nabla(\zeta \eta) dx + \int_{V_s} \left[\varepsilon_s \nabla \psi \cdot \nabla(\zeta \eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta \eta \right] dx \\
 &= - \int_{V_m} (\nabla \cdot \varepsilon_m \nabla \psi) \zeta \eta dx - \int_\Gamma (\varepsilon_m \nabla \psi \cdot n) \zeta \eta dS \\
 &\quad + \int_{V_s} \left[-(\nabla \cdot \varepsilon_s \nabla \psi)(\zeta \eta) + 4\pi \chi_{\Omega_s} B' \left(\psi - \frac{\hat{\psi}_0}{2} \right) \zeta \eta \right] dx + \int_\Gamma (\varepsilon_s \nabla \psi \cdot n) \zeta \eta dS \\
 &= \int_\Gamma \llbracket \varepsilon \nabla \psi \cdot n \rrbracket \zeta \eta dS \\
 &= 0,
 \end{aligned}
 \tag{4.3}$$

where in the last step, we used the third equation of (1.15). Now, since $\eta \in C_c^\infty(\Omega)$ is arbitrary, we obtain (2.1) from (4.2) and (4.3). \square

5. Minimization of the electrostatic free energy: Proof of Theorems 2.3, 2.4, and 2.5.

Proof of Theorem 2.3. Let $t = 1 + \beta \max_{1 \leq j \leq M} \|\mu_{0j}\|_{L^\infty(\Omega_s)}$, where μ_{0j} ($j = 1, \dots, M$) are defined in (3.14). It follows from (3.12) and (3.4) that there exists $C_5 > 0$ such that

$$F_0[c] \geq C_5 \left\| \sum_{j=1}^M q_j c_j \right\|_{H^{-1}(\Omega)}^2 + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} S_{-t}(c_j) dx + E_0 \quad \forall c = (c_1, \dots, c_M) \in V_0,
 \tag{5.1}$$

where E_0 is defined in (3.15). Let $z = \inf_{c \in V_0} F_0[c]$. Since $S_{-t} : [0, \infty) \rightarrow \mathbb{R}$ is bounded below, z is finite.

Let $c^{(k)} = (c_1^{(k)}, \dots, c_M^{(k)}) \in V_0$ ($k = 1, 2, \dots$) be such that $\lim_{k \rightarrow \infty} F_0[c^{(k)}] = z$. It follows from (5.1) that $\{\int_{\Omega_s} S_{-t}(c_j^{(k)}) dx\}$ is bounded for each $j = 1, \dots, M$. Therefore, by Lemma 3.3, up to a subsequence that is not relabeled, $\{c_j^{(k)}\}$ converges weakly in $L^1(\Omega_s)$ to some $c_j \in L^1(\Omega_s)$, and

$$\int_{\Omega_s} S_{-t}(c_j) dx \leq \liminf_{k \rightarrow \infty} \int_{\Omega_s} S_{-t}(c_j^{(k)}) dx < \infty \quad j = 1, \dots, M.
 \tag{5.2}$$

Define $c_j = 0$ on Ω_m for all $j = 1, \dots, M$. By (5.1), $\{\sum_{j=1}^M q_j c_j^{(k)}\}$ is bounded in $H^{-1}(\Omega)$. Since $H^{-1}(\Omega)$ is a Hilbert space, $\{\sum_{j=1}^M q_j c_j^{(k)}\}$ has a subsequence, again not relabeled, that weakly converges to some $F \in H^{-1}(\Omega)$. Let $\xi \in L^\infty(\Omega) \cap H_0^1(\Omega)$. We have

$$F(\xi) = \lim_{k \rightarrow \infty} \int_\Omega \left(\sum_{j=1}^M q_j c_j^{(k)} \right) \xi dx = \int_\Omega \left(\sum_{j=1}^M q_j c_j \right) \xi dx.$$

Therefore, $\sum_{j=1}^M q_j c_j \in H^{-1}(\Omega)$ and, hence, $c = (c_1, \dots, c_M) \in V_0$.

By (5.2) and the fact that the norm of a Banach space is sequentially weakly lower semicontinuous, we have $z = \liminf_{k \rightarrow \infty} F_0[c^{(k)}] \geq F_0[c] \geq z$. This implies that $c \in V_0$ is a global minimizer of $F_0 : V_0 \rightarrow \mathbb{R}$.

Let $d = (d_1, \dots, d_M) \in V_0$ be a local minimizer of $F_0 : V_0 \rightarrow \mathbb{R}$. Then for $\lambda \in (0, 1)$ close to 0, we have by the convexity of $F_0 : V_0 \rightarrow \mathbb{R}$ that

$$F_0[d] \leq F_0[\lambda c + (1 - \lambda)d] \leq \lambda F_0[c] + (1 - \lambda)F_0[d],$$

leading to $F_0[d] \leq F_0[c]$. Thus, d is also a global minimizer of $F_0 : V_0 \rightarrow \mathbb{R}$. Clearly, $(c + d)/2 \in V_0$. Consequently, it follows from the definition of the norm $\|\cdot\|$ and the Cauchy-Schwarz inequality with respect to the inner product $\langle \xi, \eta \rangle = \xi(L\eta)$ ($\xi, \eta \in H^{-1}(\Omega)$) that

$$\begin{aligned} 0 &\leq F_0\left[\frac{c+d}{2}\right] - \min_{e \in V_0} F_0[e] \\ &= F_0\left[\frac{c+d}{2}\right] - \frac{1}{2}F_0[c] - \frac{1}{2}F_0[d] \\ &= \frac{1}{8} \left\| \sum_{j=1}^M q_j(c_j + d_j) \right\|^2 - \frac{1}{4} \left\| \sum_{j=1}^M q_j c_j \right\|^2 - \frac{1}{4} \left\| \sum_{j=1}^M q_j d_j \right\|^2 \\ &\quad + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} \left[S_{-1}\left(\frac{c_j + d_j}{2}\right) - \frac{1}{2}S_{-1}(c_j) - \frac{1}{2}S_{-1}(d_j) \right] dx \\ &\leq \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} \left[S_{-1}\left(\frac{c_j + d_j}{2}\right) - \frac{1}{2}S_{-1}(c_j) - \frac{1}{2}S_{-1}(d_j) \right] dx. \end{aligned}$$

This, together with the convexity of S_{-1} on $[0, \infty)$, implies that

$$S_{-1}\left(\frac{c_j(x) + d_j(x)}{2}\right) = \frac{1}{2}S_{-1}(c_j(x)) + \frac{1}{2}S_{-1}(d_j(x)) \quad \forall j = 1, \dots, M, \forall x \in \Omega_s \setminus \omega_s,$$

for some $\omega_s \subset \Omega_s$ with $|\omega_s| = 0$. Let $x \in \Omega_s \setminus \omega_s$. Then it follows from the definition of $S_{-1} : [0, \infty) \rightarrow \mathbb{R}$ that $c_j(x) = 0$ if and only if $d_j(x) = 0$ for all $j = 1, \dots, M$. The strict convexity of S_0 on $(0, \infty)$ then implies that $c = d$ a.e. Ω_s . Hence, $c = d$ in V_0 . \square

Proof of Theorem 2.4. (1) Let $c = (c_1, \dots, c_M) \in W_0$. We show that the following four statements are equivalent:

- (i) c is an equilibrium of $F_0 : W_0 \rightarrow \mathbb{R}$;
- (ii) The property (2.5) holds true, and

$$(5.3) \quad q_j L \left(\sum_{j=1}^M q_j c_j \right) + \mu_{0j} + \beta^{-1} \log c_j = 0 \quad \text{a.e. } \Omega_s, \quad j = 1, \dots, M;$$

- (iii) c is a global minimizer of $F_0 : W_0 \rightarrow \mathbb{R}$;
- (iv) c is a local minimizer of $F_0 : W_0 \rightarrow \mathbb{R}$.

Assume (i) is true. Then (2.5) holds true by Definition 2.3. Let $e = (e_1, \dots, e_M) \in X \cap L^\infty(\Omega, \mathbb{R}^M)$. Notice that $S'_{-1}(u) = \log u$ for any $u > 0$. Thus, for each $j \in \{1, \dots, M\}$ and each $x \in \Omega_s$, the mean-value theorem implies the existence of $\theta_j(x) \in [0, 1]$ such that

$$S_{-1}(c_j(x) + te_j(x)) - S_{-1}(c_j(x)) = te_j(x) \log(c_j(x) + t\theta_j(x)e_j(x)).$$

Hence, by the Lebesgue dominated convergence theorem,

$$\lim_{t \rightarrow 0} \int_{\Omega_s} \frac{S_{-1}(c_j + te_j) - S_{-1}(c_j)}{t} dx = \int_{\Omega_s} e_j \log c_j dx, \quad j = 1, \dots, M.$$

Therefore, it follows from Definition 2.3, the definition of the norm $\|\cdot\|$, (3.12), and (3.4) that

$$\begin{aligned} 0 &= \delta F_0[c]e \\ &= \lim_{t \rightarrow 0} \frac{F_0[c + te] - F_0[c]}{t} \\ &= \lim_{t \rightarrow 0} \left\{ \frac{1}{2} t \left\| \sum_{j=1}^M q_j e_j \right\|^2 + \int_{\Omega_s} \left(\sum_{j=1}^M q_j e_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx \right. \\ &\quad \left. + \sum_{j=1}^M \int_{\Omega_s} \mu_{0j} e_j dx + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} \frac{1}{t} [S_{-1}(c_j + te_j) - S_{-1}(c_j)] dx \right\} \\ &= \sum_{j=1}^M \int_{\Omega_s} \left[q_j L \left(\sum_{j=1}^M q_j c_j \right) + \mu_{0j} + \beta^{-1} \log c_j \right] e_j dx \quad \forall e \in X \cap L^\infty(\Omega, \mathbb{R}^M). \end{aligned} \tag{5.4}$$

This implies (5.3). Hence, (ii) is true.

Assume (ii) is true. We show that (iii) is true. By Lemma 3.4, we need only to show that $F_0[c] \leq F_0[d]$ for any fixed $d = (d_1, \dots, d_M) \in W_0$ that satisfies (2.5) with c replaced by d . In fact, setting $e = (e_1, \dots, e_M) = d - c \in X \cap L^\infty(\Omega, \mathbb{R}^M)$, we have by the convexity of $S_{-1} : [0, \infty) \rightarrow \mathbb{R}$ that

$$S_{-1}(d_j) - S_{-1}(c_j) \geq (d_j - c_j) S'_{-1}(c_j) = e_j \log c_j \quad \text{a.e. } \Omega_s.$$

Therefore, it follows from (3.12) and (5.3) that

$$\begin{aligned} F_0[d] - F_0[c] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j e_j \right) L \left(\sum_{j=1}^M q_j e_j \right) dx + \int_{\Omega_s} \left(\sum_{j=1}^M q_j e_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx \\ &\quad + \sum_{j=1}^M \int_{\Omega_s} \mu_{0j} e_j dx + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} [S_{-1}(d_j) - S_{-1}(c_j)] dx \\ &\geq \sum_{j=1}^M \int_{\Omega_s} \left[q_j L \left(\sum_{i=1}^M q_i c_i \right) + \mu_{0j} + \beta^{-1} \log c_j \right] e_j dx \\ &= 0. \end{aligned}$$

Hence, $F_0[c] \leq F_0[d]$, and (iii) is true.

Clearly, (iii) implies (iv).

Finally, assume (iv) is true. By Lemma 3.4, (2.5) holds true. For any $e \in X \cap L^\infty(\Omega, \mathbb{R}^M)$, it is easy to see that $\delta F_0[c]e$ exists, cf. (5.4). Since $F_0[c + te] \geq F_0[c]$ for $|t|$ small enough, we have $\delta F_0[c]e = 0$. Therefore, c is an equilibrium of $F_0 : W_0 \rightarrow \mathbb{R}$, and (i) is true.

Let now $\psi \in H_*^1(\Omega)$ be the unique weak solution to the boundary-value problem of the PBE (1.11) and (1.6), cf. Theorem 2.1. Define $c = (c_1, \dots, c_M) : \Omega \rightarrow \mathbb{R}$ by (1.9) for point ions and $c_j(x) = 0$ for all $x \in \Omega_m$ and all $j = 1, \dots, M$. Clearly, $c \in W_0$. Moreover, by Theorem 2.1, $\psi|_{\Omega_s} \in C(\overline{\Omega_s})$. This implies (2.5). It follows from (1.11), (1.6), (1.9) for point ions, and Lemma 3.2 with $f = \sum_{j=1}^M q_j c_j$ that ψ is the electrostatic potential corresponding to c ; i.e., $\psi = G + \hat{\psi}_0 + L(\sum_{j=1}^M q_j c_j)$. This, together with the Boltzmann relations (1.9) for point ions and (3.14), implies (5.3). Hence, c is an equilibrium, and, thus, a local and global minimizer, of $F_0 : W_0 \rightarrow \mathbb{R}$. The uniqueness of equilibria or local minimizers is equivalent to that of global minimizers, and can be proved by the same argument used in the proof of Theorem 2.3.

(2) It is clear that from our definition of c and ψ that we need only to prove (2.7). Since c is the unique minimizer of $F_0 : W_0 \rightarrow \mathbb{R}$ and ψ is the corresponding electrostatic potential determined by (3.11), we have by (1.1) and (1.9) for point ions that

$$\begin{aligned}
 \min_{d \in W_0} F_0[d] &= F_0[c] \\
 &= \frac{1}{2} \sum_{i=1}^N Q_i(\psi - \psi_{vac})(x_i) + \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) \psi dx \\
 &\quad + \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} c_j [\log(a^3 c_j) - 1] dx - \sum_{j=1}^M \int_{\Omega_s} \mu_j c_j dx \\
 &= \frac{1}{2} \sum_{i=1}^N Q_i(\psi - \psi_{vac})(x_i) - \frac{1}{2} \sum_{j=1}^M \int_{\Omega_s} q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} (\psi - \hat{\psi}_0) dx \\
 (5.5) \quad &\quad - \beta^{-1} \sum_{j=1}^M \int_{\Omega_s} c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} dx.
 \end{aligned}$$

Since ψ is the unique solution to the boundary-value problem of PBE (1.11) and (1.6), and since $\hat{\psi}_0$ is harmonic in Ω_s by (1.10), we have

$$\varepsilon_s \Delta (\psi - \hat{\psi}_0) + 4\pi \sum_{j=1}^M q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} = 0 \quad \text{a.e. } \Omega_s.$$

Multiplying both sides of this equation by $\psi - \hat{\psi}_0$ and integrate the resulting terms over Ω_s , we obtain by integration by parts and the fact that by Lemma 3.1 both $\psi - \hat{\psi}_0$ and $\varepsilon_\Gamma \partial_n (\psi - \hat{\psi}_0)$ are continuous across Γ ,

$$\begin{aligned}
 & - \int_{\Omega_s} \varepsilon_s |\nabla (\psi - \hat{\psi}_0)|^2 dx + \int_\Gamma \varepsilon_\Gamma (\psi - \hat{\psi}_0) \partial_n (\psi - \hat{\psi}_0) dS \\
 & \quad + 4\pi \sum_{j=1}^M \int_{\Omega_s} q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} (\psi - \hat{\psi}_0) dx = 0.
 \end{aligned}$$

This and (5.5) imply (2.7). \square

Proof of Theorem 2.5. (1) We first show that $F_a : V_a \rightarrow \mathbb{R}$ is a convex functional. Define $T_M = \{(u_1, \dots, u_M) \in \mathbb{R}^M : u_j > 0 \text{ for } j = 1, \dots, M, \text{ and } \sum_{j=1}^M u_j < 1\}$ and

$$h(u) = \left(1 - \sum_{j=1}^M u_j\right) \left[\log \left(1 - \sum_{j=1}^M u_j\right) - 1\right] \quad \forall u = (u_1, \dots, u_M) \in T_M.$$

Clearly, T_M is convex. We have $\partial_{u_i u_j} h(u) = (1 - \sum_{k=1}^M u_k)^{-1}$ for all $1 \leq i, j \leq M$. Let $H(u) = (\partial_{u_i u_j} h)$ be the Hessian of $h : T_M \rightarrow \mathbb{R}$. Then, for any $y = (y_1, \dots, y_M) \in \mathbb{R}^M$, we have $y \cdot H(u)y = (\sum_{k=1}^M y_k)^2 / (1 - \sum_{k=1}^M u_k) \geq 0$. Therefore, $H(u)$ is symmetric, semidefinite for any $u \in T_M$. Hence, $h : T_M \rightarrow \mathbb{R}$ is convex. Consequently, since V_a is a convex subset of X and $S_{-1} : [0, \infty) \rightarrow \mathbb{R}$ is convex, we conclude by (3.13) that $F_a : V_a \rightarrow \mathbb{R}$ is convex.

Let now $c = (c_1, \dots, c_M) \in V_a$. By the same argument used in the proof of Theorem 2.4, we obtain the equivalence of the following four statements:

- (i) c is an equilibrium of $F_a : V_a \rightarrow \mathbb{R}$;
- (ii) The property (2.6) holds true, and

$$(5.6) \quad q_j L \left(\sum_{i=1}^M q_i c_i\right) + \mu_{aj} + \beta^{-1} \log \left(\frac{a^3 c_j}{1 - a^3 \sum_{i=1}^M c_i}\right) = 0$$

a.e. $\Omega_s, \quad j = 1, \dots, M;$

- (iii) c is a global minimizer of $F_a : V_a \rightarrow \mathbb{R}$;
- (iv) c is a local minimizer of $F_a : V_a \rightarrow \mathbb{R}$.

Let $\psi \in H_*^1(\Omega)$ be the unique weak solution to the boundary-value problem of the PBE (1.12) and (1.6), cf. Theorem 2.1. Define $c = (c_1, \dots, c_M) : \Omega \rightarrow \mathbb{R}$ by (1.9) for finite-size ions and $c_j(x) = 0$ for all $x \in \Omega_m$ and all $j = 1, \dots, M$. Clearly, $c \in V_a$. Moreover, by Theorem 2.1, $\psi|_{\Omega_s} \in C(\overline{\Omega_s})$. This implies (2.6). By (1.9) for finite-size ions, we have

$$a^3 \sum_{j=1}^M c_j(x) = 1 - \frac{1}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)}}.$$

This together with (1.9) for finite-size ions imply that

$$(5.7) \quad \frac{c_j}{1 - a^3 \sum_{i=1}^M c_i} = c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)}, \quad j = 1, \dots, M.$$

It follows from (1.12), (1.6), (1.9) for finite-size ions, and Lemma 3.2 with $f = \sum_{j=1}^M q_j c_j$ that ψ is the electrostatic potential corresponding to c , i.e., $\psi = G + \hat{\psi}_0 + L(\sum_{j=1}^M q_j c_j)$. This, together with (5.7) and (3.16), implies (5.6). Hence, c is an equilibrium, and, thus, a local and global minimizer, of $F_a : V_a \rightarrow \mathbb{R}$. The uniqueness of equilibria or local minimizers is equivalent to that of global minimizers, and can be proved by the same argument used in the proof of Theorem 2.3.

(2) We need only to prove (2.8). Since c is the unique minimizer of $F_a : V_a \rightarrow \mathbb{R}$ and ψ is the corresponding electrostatic potential determined by (3.11), we have by

(1.2), (1.3), and (1.9) for finite-size ions that

$$\begin{aligned}
 \min_{d \in V_a} F_a[d] &= F_a[c] \\
 &= \frac{1}{2} \sum_{i=1}^N Q_i(\psi - \psi_{vac})(x_i) + \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) \psi dx \\
 &\quad + \beta^{-1} \sum_{j=0}^M \int_{\Omega_s} c_j [\log(a^3 c_j) - 1] dx - \sum_{j=1}^M \int_{\Omega_s} \mu_j c_j dx \\
 &= \frac{1}{2} \sum_{i=1}^N Q_i(\psi - \psi_{vac})(x_i) - \frac{1}{2} \sum_{j=1}^M \int_{\Omega_s} \frac{q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} (\psi - \hat{\psi}_0)}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)}} dx \\
 (5.8) \quad &\quad - \beta^{-1} a^{-3} \int_{\Omega_s} \left[1 + \log \left(1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)} \right) \right] dx.
 \end{aligned}$$

Since ψ is the unique solution to the boundary-value problem of PBE (1.12) and (1.6), and since $\hat{\psi}_0$ is harmonic in Ω_s by (1.10), we have

$$\varepsilon_s \Delta (\psi - \hat{\psi}_0) + 4\pi \sum_{j=1}^M \frac{q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)}}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)}} = 0 \quad \text{a.e. } \Omega_s.$$

Multiplying both sides of this equation by $\psi - \hat{\psi}_0$ and integrating the resulting terms over Ω_s , we obtain by integration by parts and the fact that by Lemma 3.1 both $\psi - \hat{\psi}_0$ and $\varepsilon_\Gamma \partial_n (\psi - \hat{\psi}_0)$ are continuous across Γ ,

$$\begin{aligned}
 & - \int_{\Omega_s} \varepsilon_s |\nabla (\psi - \hat{\psi}_0)|^2 dx + \int_\Gamma \varepsilon_\Gamma (\psi - \hat{\psi}_0) \partial_n (\psi - \hat{\psi}_0) dS \\
 & + 4\pi \sum_{j=1}^M \int_{\Omega_s} \frac{q_j c_j^\infty e^{-\beta q_j (\psi - \hat{\psi}_0/2)} (\psi - \hat{\psi}_0)}{1 + a^3 \sum_{i=1}^M c_i^\infty e^{-\beta q_i (\psi - \hat{\psi}_0/2)}} dx = 0.
 \end{aligned}$$

This and (5.8) imply (2.8). \square

Appendix A.

We now prove Lemma 3.4 and Lemma 3.5 by constructing ionic concentrations that satisfy required conditions and that have lower free energies. The key idea here is based on the following observation: the function $S_\alpha : [0, \infty) \rightarrow \mathbb{R}$, defined for any $\alpha \in \mathbb{R}$ by $S_\alpha(0) = 0$ and $S_\alpha(u) = u(\alpha + \log u)$ if $u > 0$, has a unique minimizer which is a positive number. Moreover, the magnitude $|S'_\alpha(u)|$ is very large if u is close to 0 or ∞ . Notice that $-S_\alpha$ represents the entropy of the system. Therefore, small changes of concentrations near zero or infinity can largely increase the corresponding entropy and, hence, decrease the free energy.

Proof of Lemma 3.4. We first construct $\bar{c} \in W_0$ that satisfies

$$(A.1) \quad \bar{c}_j(x) \leq \gamma'_2 \quad \text{a.e. } x \in \Omega_s, \quad j = 1, \dots, M,$$

for some constant $\gamma'_2 > 0$, $\|\bar{c} - c\|_X < \varepsilon/2$, and $F_0[\bar{c}] \leq F_0[c]$ with a strict inequality if $c \notin L^\infty(\Omega, \mathbb{R}^M)$. Let $A > 0$. Define $\bar{c} = (\bar{c}_1, \dots, \bar{c}_M) : \Omega \rightarrow \mathbb{R}$ by

$$(A.2) \quad \bar{c}_j(x) = \begin{cases} c_j(x) & \text{if } c_j(x) \leq A \\ 0 & \text{if } c_j(x) > A \end{cases} \quad \forall x \in \Omega, \quad j = 1, \dots, M.$$

Clearly, $\bar{c} \in W_0$ and (A.1) holds true with $\gamma'_2 = A$. Moreover, $\sum_{j=1}^M \|\bar{c}_j - c_j\|_{L^1(\Omega)} < \varepsilon/4$ for $A > 0$ large enough.

Denote

$$\tau_j(A) = \{x \in \Omega_s : c_j(x) > A\}, \quad j = 1, \dots, M.$$

Since $c \in W_0$, there exists $p > 3/2$ such that each $c_j \in L^p(\Omega)$ ($1 \leq j \leq M$). Thus,

$$\sum_{j=1}^M q_j \bar{c}_j - \sum_{j=1}^M q_j c_j = - \sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \rightarrow 0 \quad \text{in } L^p(\Omega) \quad \text{as } A \rightarrow \infty.$$

By the definition of $L : H^{-1}(\Omega) \rightarrow H_0^1(\Omega)$ and the regularity theory for elliptic problems [14], we have $L(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j)|_{\Omega_s} \in W^{2,p}(\Omega_s)$ and

$$\left\| L \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) \right\|_{W^{2,p}(\Omega_s)} \leq C \left\| \sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right\|_{L^p(\Omega_s)} \rightarrow 0 \quad \text{as } A \rightarrow \infty.$$

Hence, by (3.4) and the embedding $W^{2,p}(\Omega_s) \hookrightarrow L^\infty(\Omega_s)$ that

$$\begin{aligned} \left\| \sum_{j=1}^M q_j \bar{c}_j - \sum_{j=1}^M q_j c_j \right\|_{H^{-1}(\Omega)}^2 &\leq C \int_{\Omega_s} \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) L \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) dx \\ &\leq C \left\| \sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right\|_{L^1(\Omega_s)} \left\| L \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) \right\|_{L^\infty(\Omega_s)} \\ &\leq C \left\| \sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right\|_{L^p(\Omega_s)} \left\| L \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) \right\|_{W^{2,p}(\Omega_s)} \\ &\leq C \left\| \sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right\|_{L^p(\Omega_s)}^2 \\ &\rightarrow 0 \quad \text{as } A \rightarrow \infty. \end{aligned}$$

Therefore, $\|\bar{c} - c\|_X < \varepsilon$ if $A > 0$ is large enough.

Notice that $\bar{c}_j = c_j - \chi_{\tau_j(A)} c_j$ for all $j = 1, \dots, M$. Thus,

$$\begin{aligned} &\frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \bar{c}_j \right) L \left(\sum_{j=1}^M q_j \bar{c}_j \right) dx - \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx \\ &= -\frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \chi_{\tau_j(A)} c_j \right) L \left(\sum_{j=1}^M q_j c_j + \sum_{j=1}^M q_j \bar{c}_j \right) dx \\ \text{(A.3)} \quad &\leq \frac{1}{2} \sum_{j=1}^M |q_j| d_j(A) \int_{\tau_j(A)} c_j dx, \end{aligned}$$

where

$$d_j(A) = \left\| L \left(\sum_{j=1}^M q_j c_j \right) \right\|_{L^\infty(\Omega_s)} + \left\| L \left(\sum_{j=1}^M q_j \bar{c}_j \right) \right\|_{L^\infty(\Omega_s)}.$$

Since

$$\begin{aligned} \left\| L \left(\sum_{j=1}^M q_j \bar{c}_j \right) \right\|_{L^\infty(\Omega_s)} &\leq C \left\| L \left(\sum_{j=1}^M q_j \bar{c}_j \right) \right\|_{W^{2,p}(\Omega_s)} \\ &\leq C \left\| \sum_{j=1}^M q_j \bar{c}_j \right\|_{L^p(\Omega_s)} \rightarrow \left\| \sum_{j=1}^M q_j c_j \right\|_{L^p(\Omega_s)} \end{aligned}$$

as $A \rightarrow \infty$, we have $\max_{1 \leq j \leq M} d_j(A) \leq C$ as $A > 0$ large enough. For each fixed $j \in \{1, \dots, M\}$ and $x \in \tau_j(A)$, we also have

$$(A.4) \quad S_{-1}(\bar{c}_j(x)) - S_{-1}(c_j(x)) = -S_{-1}(c_j(x)) = -c_j(x) \log c_j(x) \leq -c_j(x) \log A.$$

Therefore, it follows from (3.12), (A.3), and (A.4) that

$$F_0[\bar{c}] - F_0[c] \leq \sum_{j=1}^M \left(\frac{1}{2} |q_j| d_j(A) + \|\mu_{0j}\|_{L^\infty(\Omega_s)} - \beta^{-1} \log A \right) \int_{\tau_j(A)} c_j dx.$$

If $A > 0$ is large enough, this is nonpositive. If $c \notin L^\infty(\Omega, \mathbb{R}^M)$, then there exists $j \in \{1, \dots, M\}$ such that $|\tau_j(A)| > 0$ for all $A > 0$. In this case, we have the strict inequality $F_0[\bar{c}] < F_0[c]$.

We now construct $\hat{c} \in W_0$ that satisfies (2.5) with c replaced by \hat{c} , $\|\hat{c} - \bar{c}\|_X < \varepsilon/2$, and $F_0[\hat{c}] \leq F_0[\bar{c}]$ with a strict inequality if there exists $j \in \{1, \dots, M\}$ such that $|\{x \in \Omega_s : c_j(x) < \alpha\}| > 0$ for all $\alpha > 0$, all these implying that \hat{c} satisfies all the desired properties. If there exists $\gamma'_1 > 0$ such that $c_j(x) \geq \gamma'_1$ for a.e. $x \in \Omega_s$ and $j = 1, \dots, M$, then $\hat{c} = \bar{c}$ with $A \geq \gamma'_1$ (cf. (A.2)) satisfies all the desired properties with $\gamma_1 = \gamma'_1$ and $\gamma_2 = \gamma'_2$. Assume otherwise there exists $j_0 \in \{1, \dots, M\}$ such that $|\{x \in \Omega_s : c_{j_0}(x) < \alpha\}| > 0$ for all $\alpha > 0$. This means that $|\{x \in \Omega_s : \bar{c}_{j_0}(x) < \alpha\}| > 0$ for all $\alpha > 0$.

Define

$$\begin{aligned} \rho_j(\alpha) &= \{x \in \Omega_s : \bar{c}_j(x) < \alpha\} \quad \forall \alpha > 0, \quad j = 1, \dots, M, \\ I_0 &= \{j \in \{1, \dots, M\} : |\rho_j(\alpha)| > 0 \forall \alpha > 0\}, \\ I_1 &= \{1, \dots, M\} \setminus I_0. \end{aligned}$$

Clearly, $I_0 \neq \emptyset$. If $I_1 \neq \emptyset$, then there exists $\alpha_1 > 0$ such that

$$\bar{c}_j(x) \geq \alpha_1 \quad \text{a.e. } x \in \Omega_s, \quad \forall j \in I_1.$$

Define for $0 < \alpha < \alpha_1$ and $1 \leq j \leq M$

$$\hat{c}_j(x) = \begin{cases} \bar{c}_j(x) + \alpha \chi_{\rho_j(\alpha)}(x) & \text{if } j \in I_0 \\ \bar{c}_j(x) & \text{if } j \in I_1 \end{cases} \quad \forall x \in \Omega.$$

Clearly, $\hat{c} = (\hat{c}_1, \dots, \hat{c}_M) \in W_0$ and (2.5) holds true with c replaced by \hat{c} , $\gamma_1 = \alpha$, and $\gamma_2 = \gamma'_2 + \alpha$. Moreover, $\sum_{j=1}^M \|\hat{c}_j - \bar{c}_j\|_{L^1(\Omega)} < \varepsilon/4$ if $\alpha > 0$ is small enough.

Furthermore,

$$\begin{aligned}
 \left\| \sum_{j=1}^M q_j \hat{c}_j - \sum_{j=1}^M q_j \bar{c}_j \right\|_{H^{-1}(\Omega)} &= \alpha \left\| \sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right\|_{H^{-1}(\Omega)} \leq \alpha \left\| \sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right\|_{L^2(\Omega)} \\
 \text{(A.5)} \qquad \qquad \qquad &\leq \alpha \sum_{j \in I_0} |q_j| \sqrt{|\rho_j(\alpha)|} \rightarrow 0 \quad \text{as } \alpha \rightarrow 0.
 \end{aligned}$$

Hence, $\|\hat{c} - \bar{c}\|_X < \varepsilon/2$ if $\alpha > 0$ is small enough.

By the mean-value theorem and the fact that $S'_{-1}(u) = \log u$ for any $u > 0$,

$$\begin{aligned}
 \sum_{j=1}^M \int_{\Omega_s} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx &= \sum_{j \in I_0} \int_{\rho_j(\alpha)} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx \\
 &\leq \alpha \log(2\alpha) \sum_{j \in I_0} |\rho_j(\alpha)|.
 \end{aligned}$$

Consequently, it follows from (3.13), (3.4), (A.5) that

$$\begin{aligned}
 F_0[\hat{c}] - F_0[\bar{c}] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \bar{c}_j + \alpha \sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right) L \left(\sum_{j=1}^M q_j \bar{c}_j + \alpha \sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right) dx \\
 &\quad - \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \bar{c}_j \right) L \left(\sum_{j=1}^M q_j \bar{c}_j \right) dx + \alpha \sum_{j \in I_0} \int_{\rho_j(\alpha)} \mu_{0j} dx \\
 &\quad + \beta^{-1} \sum_{j \in I_0} \int_{\rho_j(\alpha)} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx \\
 &\leq \frac{\alpha^2}{2} \int_{\Omega_s} \left(\sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right) L \left(\sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right) dx \\
 &\quad + \alpha \int_{\Omega_s} \left(\sum_{j \in I_0} q_j \chi_{\rho_j(\alpha)} \right) L \left(\sum_{j=1}^M q_j \bar{c}_j \right) dx \\
 &\quad + \alpha \sum_{j \in I_0} \|\mu_{0j}\|_{L^\infty(\Omega_s)} |\rho_j(\alpha)| + \alpha \sum_{j \in I_0} \beta^{-1} \log(2\alpha) |\rho_j(\alpha)| \\
 &\leq \frac{\alpha^2}{2} \left(\sum_{i \in I_0} q_i^2 \right) \sum_{j \in I_0} |\rho_j(\alpha)| + \alpha \left\| L \left(\sum_{j=1}^M q_j \bar{c}_j \right) \right\|_{L^\infty(\Omega_s)} \sum_{j \in I_0} |q_j| |\rho_j(\alpha)| \\
 &\quad + \alpha \sum_{j \in I_0} \|\mu_{0j}\|_{L^\infty(\Omega_s)} |\rho_j(\alpha)| + \alpha \sum_{j \in I_0} \beta^{-1} \log(2\alpha) |\rho_j(\alpha)| \\
 &= \alpha \sum_{j \in J_0} \left[\frac{\alpha}{2} \left(\sum_{j \in I_0} q_j^2 \right) + |q_j| \left\| L \left(\sum_{j=1}^M q_j \bar{c}_j \right) \right\|_{L^\infty(\Omega_s)} \right. \\
 &\quad \left. + \|\mu_{0j}\|_{L^\infty(\Omega_s)} + \beta^{-1} \log(2\alpha) \right] |\rho_j(\alpha)|.
 \end{aligned}$$

Since $I_0 \neq \emptyset$, this is strictly negative if $\alpha > 0$ is small enough.

Proof of Lemma 3.5. We first construct $\bar{c} = (\bar{c}_1, \dots, \bar{c}_M) \in V_a$ such that

$$(A.6) \quad a^3 \bar{c}_0(x) = 1 - a^3 \sum_{j=1}^M \bar{c}_j(x) \geq \theta_1 \quad \text{a.e. } x \in \Omega_s$$

for some constant $\theta_1 \in (0, 1)$, $\|\bar{c} - c\|_X < \varepsilon/2$, and $F_a[\bar{c}] \leq F_a[c]$ with a strict inequality if $|\{x \in \Omega_s : a^3 c_0(x) < \alpha\}| > 0$ for all $\alpha > 0$.

Denote for any $\alpha > 0$

$$\omega_0(\alpha) = \{x \in \Omega_s : a^3 c_0(x) < \alpha\}.$$

If there exists a constant $\alpha_1 > 0$ such that $|\omega_0(\alpha_1)| = 0$, i.e., $a^3 c_0(x) \geq \alpha_1$ a.e. Ω_s , then $(\bar{c}_1, \dots, \bar{c}_M) = (c_1, \dots, c_M) \in V_a$ satisfies all the desired properties with $\theta_1 = \alpha_1/(1 + \alpha_1) \in (0, 1)$. Suppose $|\omega_0(\alpha)| > 0$ for any $\alpha > 0$. Let $0 < \alpha < 1/(4M)$. Let $x \in \omega_0(\alpha)$. Then there exists some $j = j(x) \in \{1, \dots, M\}$ such that $a^3 c_j(x) \geq 1/(2M)$. In fact, if this were not true, then $a^3 c_i(x) < 1/(2M)$ for all $i = 1, \dots, M$. Hence, $a^3 c_0(x) = 1 - a^3 \sum_{i=1}^M c_i(x) > 1/2 > \alpha$. This would mean that $x \notin \omega_0(\alpha)$, a contradiction. Denoting

$$H_j(\alpha) = \left\{ x \in \omega_0(\alpha) : a^3 c_j(x) \geq \frac{1}{2M} \right\}, \quad j = 1, \dots, M,$$

we, thus, have $\omega_0(\alpha) = \cup_{j=1}^M H_j(\alpha)$. Since $|\omega_0(\alpha)| > 0$, we have $|H_{j_1}(\alpha)| > 0$ for some j_1 ($1 \leq j_1 \leq M$). If $|H_j(\alpha) \setminus H_{j_1}(\alpha)| = 0$ for all $j \neq j_1$, then we have $\omega_0(\alpha) = \tilde{K}_1(\alpha) \cup H_{j_1}(\alpha)$ for some $\tilde{K}_1(\alpha) \subset \omega_0(\alpha)$ with $|\tilde{K}_1(\alpha)| = 0$. Otherwise, $|H_{j_2}(\alpha) \setminus H_{j_1}(\alpha)| > 0$ for some $j_2 \neq j_1$. In case $|\omega_0(\alpha) \setminus [H_{j_1}(\alpha) \cup H_{j_2}(\alpha)]| = 0$, we have $\omega_0(\alpha) = \tilde{K}_2(\alpha) \cup H_{j_1}(\alpha) \cup [H_{j_2}(\alpha) \setminus H_{j_1}(\alpha)]$ for some $\tilde{K}_2(\alpha) \subset \omega_0(\alpha)$ with $|\tilde{K}_2(\alpha)| = 0$. By induction, we see that there exist $m \in \{1, \dots, M\}$, $\tilde{K}_m(\alpha) \subset \omega_0(\alpha)$ with $|\tilde{K}_m(\alpha)| = 0$, and mutually disjoint sets $K_{j_1}(\alpha), \dots, K_{j_m}(\alpha) \subseteq \omega_0(\alpha)$ such that $K_{j_i}(\alpha) \subseteq H_{j_i}(\alpha)$ and $|K_{j_i}(\alpha)| > 0$ for $i = 1, \dots, m$, and $\omega_0(\alpha) = \tilde{K}_m(\alpha) \cup [\cup_{i=1}^m K_{j_i}(\alpha)]$. By relabeling, we may assume that $j_i = i$ for $i = 1, \dots, m$.

Define now

$$(A.7) \quad \bar{c}_j(x) = \begin{cases} c_j(x) - \alpha a^{-3} \chi_{K_j(\alpha)}(x) & \forall x \in \Omega, \quad j = 1, \dots, m, \\ c_j(x) & \forall x \in \Omega, \quad j = m + 1, \dots, M, \end{cases}$$

$$\bar{c}_0(x) = a^{-3} \left[1 - a^3 \sum_{j=1}^M \bar{c}_j(x) \right] \quad \forall x \in \Omega_s.$$

It is easy to see that $(\bar{c}_1, \dots, \bar{c}_M) \in V_a$. Moreover,

$$(A.8) \quad a^3 \bar{c}_0(x) = a^3 c_0(x) + \alpha \chi_{\omega_0(\alpha)}(x) \geq \alpha \quad \text{a.e. } x \in \Omega_s,$$

implying (A.6) with $\theta_1 = \alpha$. Clearly, $\sum_{j=1}^M \|\bar{c}_j - c_j\|_{L^1(\Omega)} \leq \alpha a^{-3} \sum_{j=1}^m |K_j(\alpha)|$. Moreover,

$$\left\| \sum_{j=1}^M q_j \bar{c}_j - \sum_{j=1}^M q_j c_j \right\|_{H^{-1}(\Omega)} \leq \alpha a^{-3} \left\| \sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right\|_{L^2(\Omega)} \leq \alpha a^{-3} \sqrt{\sum_{j=1}^m q_j^2 |K_j(\alpha)|}.$$

Therefore, $\|\bar{c} - c\|_X < \varepsilon/2$, provided that $\alpha > 0$ is small enough.

If $x \in K_j(\alpha)$ for some j with $1 \leq j \leq m$, then $c_j(x) \geq 1/(2Ma^3)$, and $\bar{c}_j(x) \geq 1/(4Ma^3)$ since $0 < \alpha < 1/(4M)$. By the mean-value theorem and the fact that $S'_{-1}(u) = \log u$ for any $u > 0$, there exists $\eta_j(x)$ with $\bar{c}_j(x) \leq \eta_j(x) \leq c_j(x)$ such that

$$\begin{aligned} S_{-1}[\bar{c}_j(x)] - S_{-1}[c_j(x)] &= [\bar{c}_j(x) - c_j(x)] \log \eta_j(x) \\ &\leq -\alpha a^{-3} \log \bar{c}_j(x) \leq \alpha a^{-3} \log(4Ma^3). \end{aligned}$$

By the same argument using (A.8) and the definition of $\omega_0(\alpha)$, we obtain

$$S_{-1}(\bar{c}_0(x)) - S_{-1}(c_0(x)) \leq \alpha a^{-3} \log(a^{-3}\alpha) \quad \text{a.e. } x \in \omega_0(\alpha).$$

Consequently, we have by (3.13), (3.4), and the embedding $L^2(\Omega_s) \hookrightarrow H^{-1}(\Omega_s)$ that

$$\begin{aligned} F_a[\bar{c}] - F_a[c] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j - \alpha a^{-3} \sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right) \\ &\quad L \left(\sum_{j=1}^M q_j c_j - \alpha a^{-3} \sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right) dx \\ &\quad - \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j c_j \right) L \left(\sum_{j=1}^M q_j c_j \right) dx - \alpha a^{-3} \sum_{j=1}^m \int_{K_j(\alpha)} \mu_{aj} dx \\ &\quad + \beta^{-1} \sum_{j=0}^m \int_{\omega_0(\alpha)} [S_{-1}(\bar{c}_j) - S_{-1}(c_j)] dx \\ &\leq \frac{1}{2} \alpha^2 a^{-6} \int_{\Omega_s} \left(\sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right) L \left(\sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right) dx \\ &\quad - \alpha a^{-3} \int_{\Omega_s} \left(\sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right) L \left(\sum_{j=1}^m q_j c_j \right) dx \\ &\quad + \alpha a^{-3} \sum_{j=1}^m \|\mu_{aj}\|_{L^\infty(\Omega_s)} |K_j(\alpha)| \\ &\quad + \beta^{-1} \alpha a^{-3} \log(a^{-3}\alpha) |\omega_0(\alpha)| + \beta^{-1} \alpha a^{-3} \log(4Ma^3) |\omega_0(\alpha)| \\ &\leq C \alpha^2 a^{-6} \left\| \sum_{j=1}^m q_j \chi_{K_j(\alpha)} \right\|_{L^2(\Omega_s)}^2 \\ &\quad + \alpha a^{-3} \left\| L \left(\sum_{j=1}^M q_j c_j \right) \right\|_{L^\infty(\Omega)} \sum_{j=1}^m |q_j| |K_j(\alpha)| \\ &\quad + \alpha a^{-3} \sum_{j=1}^m \|\mu_{aj}\|_{L^\infty(\Omega_s)} |K_j(\alpha)| + \beta^{-1} \alpha a^{-3} \log(8Ma^3\alpha) \sum_{j=1}^m |K_j(\alpha)| \\ &= \alpha \sum_{j=1}^M \left[C \alpha a^{-6} q_j^2 + a^{-3} |q_j| \left\| L \left(\sum_{j=1}^M q_j c_j \right) \right\|_{L^\infty(\Omega)} + a^{-3} \|\mu_{aj}\|_{L^\infty(\Omega_s)} \right. \\ &\quad \left. + \beta^{-1} a^{-3} \log(8Ma^3\alpha) \right] |K_j(\alpha)|, \end{aligned}$$

where $C > 0$ is a constant independent of α . Thus, $F_a[\bar{c}] - F_a[c]$ is nonpositive for $\alpha > 0$ sufficiently small. It is strictly negative, if $|\omega_0(\alpha)| = \sum_{j=1}^m |K_j(\alpha)| > 0$ for all $\alpha > 0$, i.e., if $|\{x \in \Omega_s : a^3 c_0(x) < \alpha\}| > 0$ for all $\alpha > 0$.

We now construct $\hat{c} \in V_a$ that satisfies (2.6) with c replaced by \hat{c} , $\|\hat{c} - \bar{c}\|_X < \varepsilon/2$, and $F_a[\hat{c}] \leq F_a[\bar{c}]$ with a strict inequality if there exists $j \in \{1, \dots, M\}$ such that $|\{x \in \Omega_s : a^3 c_j(x) < \alpha\}| > 0$ for all $\alpha > 0$, all these implying that \hat{c} satisfies all the desired properties. If there exists $\theta_2 \in (0, 1)$ such that $c_j(x) \geq \theta_2$ for a.e. $x \in \Omega_s$ and all $j = 1, \dots, M$, then $\hat{c} = \bar{c}$ with $0 < \alpha < \theta_2/2$ (cf. (A.7)) satisfies all the desired properties with $\theta_0 = \min(\theta_1, \theta_2/2)$. Assume otherwise there exists $j_0 \in \{1, \dots, M\}$ such that $|\{x \in \Omega_s : c_{j_0}(x) < \alpha\}| > 0$ for all $\alpha > 0$. This means that $|\{x \in \Omega_s : \bar{c}_{j_0}(x) < \alpha\}| > 0$ for all $\alpha > 0$.

Define

$$\begin{aligned} \sigma_j(\alpha) &= \{x \in \Omega_s : a^3 \bar{c}_j(x) < \alpha\} \quad \forall \alpha > 0, \quad j = 1, \dots, M, \\ J_0 &= \{j \in \{1, \dots, M\} : |\sigma_j(\alpha)| > 0 \quad \forall \alpha > 0\}, \\ J_1 &= \{1, \dots, M\} \setminus J_0. \end{aligned}$$

Clearly, $J_0 \neq \emptyset$. If $J_1 \neq \emptyset$, then there exists $\alpha_2 > 0$ such that

$$a^3 \bar{c}_j(x) \geq \alpha_2 \quad \text{a.e. } x \in \Omega_s, \quad \forall j \in J_1.$$

Define for $0 < \alpha < \min\{\alpha_2, \theta_1/M\}$ and $1 \leq j \leq M$

$$\begin{aligned} \hat{c}_j(x) &= \begin{cases} \bar{c}_j(x) + \alpha a^{-3} \chi_{\sigma_j(\alpha)}(x) & \text{if } j \in J_0 \\ \bar{c}_j(x) & \text{if } j \in J_1 \end{cases} \quad \forall x \in \Omega. \\ \hat{c}_0(x) &= a^{-3} \left[1 - \sum_{j=1}^M a^3 \hat{c}_j(x) \right] \quad \forall x \in \Omega_s. \end{aligned}$$

Notice by (A.6) that

$$\begin{aligned} a^3 \hat{c}_0(x) &= 1 - \sum_{j=1}^M a^3 \hat{c}_j(x) = 1 - \sum_{j=1}^M a^3 \bar{c}_j(x) - \alpha \sum_{j \in J_0} \chi_{\sigma_j(\alpha)} \\ &\geq \theta_1 - \alpha M > 0 \quad \text{a.e. } x \in \Omega_s. \end{aligned}$$

Thus, $\hat{c} = (\hat{c}_1, \dots, \hat{c}_M) \in V_a$. Clearly, (2.6) holds true for $\theta_0 = \min\{\alpha, \alpha_2, \theta_1 - \alpha M\}$. Applying the same argument used above, we obtain that $\|\hat{c} - \bar{c}\|_X < \varepsilon/2$ for $\alpha > 0$ small enough.

We have now by the mean-value theorem that

$$\begin{aligned} \sum_{j=1}^M \int_{\Omega_s} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx &= \sum_{j \in J_0} \int_{\sigma_j(\alpha)} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx \\ &\leq \alpha a^{-3} \log(2\alpha a^{-3}) \sum_{j \in J_0} |\sigma_j(\alpha)|. \end{aligned}$$

Similarly, we have by (2.6) that

$$\int_{\Omega_s} [S_{-1}(\hat{c}_0) - S_{-1}(\bar{c}_0)] dx \leq -\alpha a^{-3} \log(a^{-3} \theta_0) \sum_{j \in J_0} |\sigma_j(\alpha)|.$$

Consequently, we have by (3.13) and a similar argument that

$$\begin{aligned}
F_a[\hat{c}] - F_a[\bar{c}] &= \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \bar{c}_j + \alpha a^{-3} \sum_{j \in J_0} q_j \chi_{\sigma_j(\alpha)} \right) \\
&\quad L \left(\sum_{j=1}^M q_j \bar{c}_j + \alpha a^{-3} \sum_{j \in J_0} q_j \chi_{\sigma_j(\alpha)} \right) dx \\
&\quad - \frac{1}{2} \int_{\Omega_s} \left(\sum_{j=1}^M q_j \bar{c}_j \right) L \left(\sum_{j=1}^M q_j \bar{c}_j \right) dx + \alpha a^{-3} \sum_{j \in J_0} \int_{\sigma_j(\alpha)} \mu_{aj} dx \\
&\quad + \beta^{-1} \sum_{j \in J_0} \int_{\sigma_j(\alpha)} [S_{-1}(\hat{c}_j) - S_{-1}(\bar{c}_j)] dx \\
&\quad + \beta^{-1} \int_{\Omega_s} [S_{-1}(\hat{c}_0) - S_{-1}(\bar{c}_0)] dx \\
&\leq C \alpha^2 a^{-6} \left\| \sum_{j \in J_0} q_j \chi_{\sigma_j(\alpha)} \right\|_{L^2(\Omega_s)}^2 + \alpha a^{-3} \left\| \sum_{j=1}^M q_j \bar{c}_j \right\|_{L^\infty(\Omega_s)} \sum_{j \in J_0} |q_j| |\sigma_j(\alpha)| \\
&\quad + \alpha a^{-3} \sum_{j \in J_0} \|\mu_{aj}\|_{L^\infty(\Omega_s)} |\sigma_j(\alpha)| + \beta^{-1} \alpha a^{-3} \log(2\alpha/\theta_0) \sum_{j \in J_0} |\sigma_j(\alpha)| \\
&\leq \alpha \sum_{j \in J_0} \left[C \alpha a^{-6} q_j^2 + a^{-3} |q_j| \left\| \sum_{j=1}^M q_j \bar{c}_j \right\|_{L^\infty(\Omega_s)} + a^{-3} \|\mu_{aj}\|_{L^\infty(\Omega_s)} \right. \\
&\quad \left. + \beta^{-1} a^{-3} \log(2\alpha/\theta_0) \right] |\sigma_j(\alpha)|.
\end{aligned}$$

Since $J_0 \neq \emptyset$, this is strictly negative if $\alpha > 0$ is sufficiently small. The case that $J_1 = \emptyset$ can be treated similarly. \square

Acknowledgments. The author thanks Dr. Jianwei Che, Dr. Joachim Dzubiella, and Dr. Benzhuo Lu for helpful discussions on the background of the underlying problem.

REFERENCES

- [1] R. ALLEN, J.-P. HANSEN, AND S. MELCHIONNA, *Electrostatic potential inside ionic solutions confined by dielectrics: A variational approach*, Phys. Chem. Chem. Phys., 3 (2001), pp. 4177–4186.
- [2] M. Z. BORN, *Volumen und Hydratationswärme der Ionen*, Phys., 1 (1920), pp. 45–48.
- [3] I. BORUKHOV, D. ANDELMAN, AND H. ORLAND, *Steric effects in electrolytes: A modified Poisson-Boltzmann equation*, Phys. Rev. Lett., 79 (1997), pp. 435–438.
- [4] D. L. CHAPMAN, *A contribution to the theory of electrocapillarity*, Phil. Mag., 25 (1913), pp. 475–481.
- [5] J. CHE, J. DZUBIELLA, B. LI, AND J. A. MCCAMMON, *Electrostatic free energy and its variations in implicit solvent models*, J. Phys. Chem. B, 112 (2008), pp. 3058–3069.
- [6] T. CHOU, *Physics of cellular materials: Basic electrostatics*, IPAM Lecture Notes (2002).
- [7] J. B. CONWAY, *A Course in Functional Analysis*, Springer-Verlag, Berlin, 1985.

- [8] M. E. DAVIS AND J. A. MCCAMMON, *Electrostatics in biomolecular structure and dynamics*, Chem. Rev., 90 (1990), pp. 509–521.
- [9] P. DEBYE AND E. HÜCKEL, *Zur theorie der elektrolyte*, Physik. Zeitschr., 24 (1923), pp. 185–206.
- [10] B. DERJAGUIN AND L. LANDAU, *A theory of the stability of strongly charged lyophobic sols and the coalescence of strongly charged particles in electrolytic solution*, Acta Phys.-Chim. USSR, 14 (1941), pp. 633–662.
- [11] M. FIXMAN, *The Poisson-Boltzmann equation and its applications to polyelectrolytes*, J. Chem. Phys., 70 (1979), pp. 4995–5005.
- [12] F. FOGOLARI AND J. M. BRIGGS, *On the variational approach to Poisson-Boltzmann free energies*, Chem. Phys. Lett., 281 (1997), pp. 135–139.
- [13] F. FOGOLARI, P. ZUCCATO, G. ESPOSITO, AND P. VIGLINO, *Biomolecular electrostatics with the linearized Poisson-Boltzmann equation*, Biophys. J., 76 (1999), pp. 1–16.
- [14] D. GILBARG AND N. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, 2nd ed., Springer-Verlag, Berlin, 1998.
- [15] M. K. GILSON, M. E. DAVIS, B. A. LUTY, AND J. A. MCCAMMON, *Computation of electrostatic forces on solvated molecules using the Poisson-Boltzmann equation*, J. Phys. Chem., 97 (1993), pp. 3591–3600.
- [16] M. GOUY, *Sur la constitution de la charge électrique a la surface d'un électrolyte*, J. de Phys., 9 (1910), pp. 457–468.
- [17] D. C. GRAHAME, *The electrical double layer and the theory of electrocapillarity*, Chem. Rev., 32 (1947), pp. 441–501.
- [18] P. GROCHOWSKI AND J. TRYLSKA, *Continuum molecular electrostatics, salt effects and counterion binding—A review of the Poisson-Boltzmann model and its modifications*, Biopolymers, 89 (2008), pp. 93–113.
- [19] J.-P. HANSEN AND H. LÖWEN, *Effective interactions between electric double layers*, Annu. Rev. Phys. Chem., 51 (2000), pp. 209–242.
- [20] V. KRALJ-IGLIC AND A. IGLIC, *A simple statistical mechanical approach to the free energy of the electric double layer including the excluded volume effect*, J. Phys. II (France), 6 (1996), pp. 477–491.
- [21] W. LITTMAN, G. STAMPACCHIA, AND H. F. WEINBERGER, *Regular points for elliptic equations with discontinuous coefficients*, Ann. Scuola Norm. Sup. Pisa, 3 (1963), pp. 43–77.
- [22] B. LU, X. CHENG, T. HOU, AND J. A. MCCAMMON, *Calculation of the Maxwell stress tensor and the Poisson-Boltzmann force on a solvated molecular surface using hypersingular boundary integrals*, J. Chem. Phys., 123 (2005), p. 084904.
- [23] B. LU, X. CHENG, J. HUANG, AND J. A. MCCAMMON, *Order N algorithm for computation of electrostatic interaction in biomolecular systems*, PNAS, 103 (2006), pp. 19314–19319.
- [24] B. LU, D. ZHANG, AND J. A. MCCAMMON, *Computation of electrostatic forces between solvated molecules determined by Poisson-boltzmann equations using a boundary element method*, J. Chem. Phys., 122 (2005), p. 214102.
- [25] E. J. MACSHANE, *Integration*, Princeton University Press, Princeton, NJ, 1947.
- [26] E. S. REINER AND C. J. RADKE, *Variational approach to the electrostatic free energy in charged colloidal suspensions: General theory for open systems*, J. Chem. Soc. Faraday Trans., 86 (1990), pp. 3901–3912.
- [27] B. ROUX AND T. SIMONSON, *Implicit solvent models*, Biophys. Chem., 78 (1999), pp. 1–20.
- [28] K. A. SHARP AND B. HONIG, *Calculating total electrostatic energies with the nonlinear Poisson-Boltzmann equation*, J. Phys. Chem., 94 (1990), pp. 7684–7692.
- [29] K. A. SHARP AND B. HONIG, *Electrostatic interactions in macromolecules: Theory and applications*, Annu. Rev. Biophys. Biophys. Chem., 19 (1990), pp. 301–332.
- [30] R. TEMAM, *Navier-Stokes Equations: Theory and Numerical Analysis*, 3rd ed., North-Holland, Amsterdam, 1984.
- [31] E. J. W. VERWEY AND J. T. G. OVERBEEK, *Theory of the Stability of Lyophobic Colloids*, Elsevier, Amsterdam, 1948.
- [32] K. YOSIDA, *Functional Analysis*, 6th ed., Springer-Verlag, Berlin, 1994.

WAVE BREAKING AND PERSISTENCE PROPERTIES FOR THE DISPERSIVE ROD EQUATION*

ZHENG GUANG GUO[†] AND YONG ZHOU[‡]

Abstract. This paper is concerned with some aspects of blow up of solutions and persistence properties. Firstly, we will try to give sufficient conditions on the initial data, which guarantee finite time singularity formation for the corresponding solutions. Then a particular class of initial data for the periodic case is also considered in this paper. Finally, we investigate the persistence properties of the strong solutions.

Key words. rod equation, convolution problem, blow up, persistence property

AMS subject classifications. 30C70, 37L05, 35Q58, 58E35

DOI. 10.1137/080734704

1. Introduction. Although a rod is always three-dimensional, if its diameter is much less than the axial length scale, one-dimensional equations can give a good description of the motion of the rod. Recently Dai [18] derived a new (one-dimensional) nonlinear dispersive equation including extra nonlinear terms involving second-order and third-order derivatives for a compressible hyperelastic material. The equation reads

$$v_\tau + \sigma_1 v v_\xi + \sigma_2 v \xi \xi_\tau + \sigma_3 (2v_\xi v \xi_\xi + v v \xi \xi \xi) = 0,$$

where $v(\xi, \tau)$ represents the radial stretch relative to a prestressed state, $\sigma_1 \neq 0$, $\sigma_2 < 0$, and $\sigma_3 \leq 0$ are constants determined by the prestress and the material parameters. If one introduces the following transformations

$$\tau = \frac{3\sqrt{-\sigma_2}}{\sigma_1} t, \quad \xi = \sqrt{-\sigma_2} x,$$

then the above equation turns into

$$(1.1) \quad u_t - u_{xxt} + 3uu_x = \gamma(2u_x u_{xx} + uu_{xxx}),$$

where $\gamma = 3\sigma_3/(\sigma_1\sigma_2)$. In [19], the authors derived that the value range of γ is from -29.4760 to 3.4174 for some special compressible materials. In this paper, from the mathematical viewpoint, we regard γ as a real number.

When $\gamma = 1$ in (1.1), we recover the shallow water (Camassa–Holm) equation derived physically by Camassa and Holm in [4] by approximating directly the Hamiltonian for Euler’s equations in the shallow water regime, where $u(x, t)$ represents the free surface above a flat bottom. Recently, the alternative derivations of the

*Received by the editors September 7, 2008; accepted for publication (in revised form) December 9, 2008; published electronically March 20, 2009. This work was partially supported by the Shanghai Rising-Star (08QH14006), Fok Ying Tung Education Foundation (111002), and Shuguang Project (07SG29).

<http://www.siam.org/journals/sima/40-6/73470.html>

[†]Department of Mathematics, Zhejiang Normal University, Jinhua, Zhejiang 321004, China and Department of Mathematics, East China Normal University, Shanghai 200026, China (gzg19801213@yahoo.com.cn).

[‡]Corresponding author. Department of Mathematics, Zhejiang Normal University, Jinhua, Zhejiang 321004, China (yzhoumath@zjnu.edu.cn).

Camassa–Holm equation were presented in [5, 25, 13]. Some satisfactory results have been obtained for this shallow water equation recently. Local well-posedness for the initial datum $u_0(x) \in H^s$, with $s > 3/2$ was proved by several authors [10, 26, 32, 34]. For the initial data with lower regularity, we refer to Molinet’s paper [30] and also the recent paper [3]. Moreover, wave breaking for a large class of initial data has been established in [8, 9, 10, 11, 26, 29, 36, 40, 41]. Later, in [20], the first author considered the Camassa–Holm equation with weakly dissipative term and established blow-up criteria and sufficient conditions on global existence of the solutions. Recently, in [24], among others, Himonas et al. showed the infinite propagation speed for the Camassa–Holm equation in the sense that a strong solution of the Cauchy problem with compact initial profile cannot be compactly supported at any later time unless it is the zero solution, which is an improvement of previous results in this direction obtained in [6, 21, 22]. These results were recently extended in [23] for a large range of physically important equations. However, in [38], global existence of weak solutions is proved, but uniqueness is obtained only under an a priori assumption that is known to hold only for initial data $u_0(x) \in H^1$ such that $u_0 - u_{0xx}$ is a sign-definite Radon measure (under this condition, global existence and uniqueness was also shown in [14], in which case a stronger concept of weak solutions can be introduced). Actually, there are approaches towards unique global weak solutions for initial data in H^1 , the important distinction being between conservative solutions [3] and dissipative solutions [2]. The solitary waves of the Camassa–Holm equation are peaked solitons [4] with $u(x, t) = e^{-|x-ct|}$, where $c \in \mathbb{R}$ is the wave speed. The orbital stability of the peakons was shown by Constantin and Strauss [16]. An alternative approach is proposed in the paper [15]. Here it is worth it to point out that the peakons replicate a feature that is characteristic for the waves of great height-waves of largest amplitude that are exact solutions of the governing equations for water waves; cf. [7, 35, 12].

If $\gamma = 0$, (1.1) is the Benjamin–Bona–Mahony (BBM) equation, a well-known model for surface waves in a canal [1], and its solutions are global.

For general $\gamma \in \mathbb{R}$, much of the early results were proven by Constantin and Strauss in [17] first. Local well-posedness of strong solutions to (1.1) was established by applying Kato’s theory, and some sufficient conditions on the initial data were found to guarantee the finite blow up of the corresponding solutions for the spatially nonperiodic case. Later, Zhou [43] proved the well-posedness result in detail, and various refined sufficient conditions on the initial data were found to guarantee the finite time blow up of the corresponding solutions for both the spatially periodic and nonperiodic case; he also applied the best constant of a convolution problem to establish the corresponding blow-up criteria [44]. Recently, blow-up criteria were presented in [37, 39, 42]. It should be mentioned that Liu and Zhou improved their previous results by applying the optimal constant on a kind of Sobolev inequality in [27] (see also a recent paper [28]). Mustafa [31] established a solution with constant H^1 energy. For $\gamma < 1$, the stability of solitary waves was established in [17, 45].

We now finish this introduction by outlining the rest of the paper. In section 2, we recall the local well-posedness for (1.1) with initial datum $u_0 \in H^s$, $s > 3/2$, and the lifespan of the corresponding solution is finite if and only if its first-order derivative blows up. In section 3, we investigate various sufficient conditions of the initial datum to guarantee the finite time blow up. Persistence properties are established for (1.1) in section 4.

2. Preliminaries. In this section, we would like to list some useful theorems for later use.

THEOREM 2.1 (see [39]). *Let the initial datum $u_0(x) \in H^s(\Omega)$, $s > 3/2$. Then there exist $T = T(\|u_0\|) > 0$ and a unique solution $u(x, t)$, which depends continuously on the initial datum u_0 to (1.1) such that*

$$u \in C([0, T]; H^s(\Omega)) \cap C^1([0, T]; H^{s-1}(\Omega)).$$

Moreover, the following two quantities $E_1(u)$ and $E_2(u)$ are invariants with respect to time t for (1.1):

$$E_1(u) = \int_{\Omega} (u^2 + u_x^2) dx, \quad E_2(u) = \int_{\Omega} (u^3 + \gamma u u_x^2) dx.$$

Actually, the local well-posedness was proved for both the periodic and nonperiodic case in the above paper; these two invariants play an important role in considering blow-up phenomenon.

The maximum value of T in Theorem 2.1 is called the lifespan of the solution, in general. If $T < \infty$, that is, $\limsup_{t \uparrow T} \|u(\cdot, t)\|_{H^s} = \infty$, we say that the solution blows up in finite time. The following theorem tells us that the solution blows up if and only if the first-order derivative blows up. This phenomenon coincides physically with the rod breaking.

THEOREM 2.2 (see [17]). *Let $u_0(x) \in H^s(\Omega)$, $s > 3/2$, and $u(x, t)$ be the corresponding solution to problem (1.1), with lifespan T . Then*

$$(2.1) \quad \sup_{x \in \Omega, 0 \leq t < T} |u(x, t)| \leq C(\|u_0\|_{H^1}).$$

T is bounded if and only if

$$(2.2) \quad \liminf_{t \uparrow T} \inf_{x \in \Omega} \{\gamma u_x(x, t)\} = -\infty.$$

For $\gamma \neq 0$, we set

$$(2.3) \quad m(t) := \inf_{x \in \Omega} (u_x(x, t) \text{sign}\{\gamma\}), \quad t > 0,$$

where $\text{sign}\{a\}$ is the sign function of $a \in \mathbb{R}$ and we set $m_0 := m(0)$. Then for every $t \in [0, T)$, there exists at least one point $\xi(t) \in \Omega$, with $m(t) = u_x(\xi(t), t)$. Just as the proof given in [11], one can show the following property of $m(t)$.

THEOREM 2.3 (see [11]). *Let $u(t)$ be the solution to (1.1) on $[0, T)$, with initial data $u_0 \in H^s(\Omega)$, $s > 3/2$, as given by Theorem 2.1. Then the function $m(t)$ is almost everywhere differentiable on $[0, T)$, with*

$$\frac{dm(t)}{dt} = u_{xt}(\xi(t), t), \quad \text{a.e. on } (0, T).$$

To consider the quantity $m(t)$ for wave breaking comes from an idea of Seliger [33] originally, the rigorous regularity proof is given in [11] for the Camassa–Holm equation. As we know, the operator $(I - \partial_x^2)^{-1}$ can be expressed by

$$(I - \partial_x^2)^{-1} f = G * f = \int_{\Omega} G(x - y) f(y) dy$$

for all $f \in L^2(\Omega)$, where $G(x)$ is the associated Green’s function. For the periodic case and the whole line case $G(x) = \frac{\cosh(x - [x] - \frac{1}{2})}{2 \sinh(\frac{1}{2})}$ and $G(x) = \frac{1}{2} e^{-|x|}$, respectively, where $[x]$ denotes the integer part of x .

In this paper, $\mathbb{S} := \mathbb{R}/\mathbb{Z}$ stands for the unit circle and $\Omega = \mathbb{R}$ or \mathbb{S} . Then (1.1) can be rewritten as

$$(2.4) \quad u_t + \gamma uu_x + \partial_x G * \left(\frac{3-\gamma}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) = 0.$$

In what follows, just for simplicity, we assume that $0 < \gamma < 3$ and introduce the following notation:

$$F(u) = \frac{3-\gamma}{2} u^2 + \frac{\gamma}{2} u_x^2.$$

3. Blow-up phenomena. In this section, we investigate sufficient conditions on the initial data which guarantee the finite time blow up of the corresponding solutions to (1.1). Firstly, we consider the following convolution problem. We can compute

$$\begin{aligned} G * \left(\frac{\gamma\alpha^2}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) (x) &= \frac{1}{2} e^{-x} \int_{-\infty}^x e^y \left(\frac{\gamma\alpha^2}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) dy \\ &\quad + \frac{1}{2} e^x \int_x^{\infty} e^{-y} \left(\frac{\gamma\alpha^2}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) dy \\ &\geq \frac{1}{2} e^{-x} \int_{-\infty}^x \gamma\alpha e^y uu_x dy - \frac{1}{2} e^x \int_x^{\infty} \gamma\alpha e^{-y} uu_x dy \\ &= \frac{\gamma\alpha}{2} u^2 - \frac{\gamma\alpha}{4} \int_{-\infty}^x e^y u^2 dy - \frac{\gamma\alpha}{4} \int_x^{\infty} e^{-y} u^2 dy, \end{aligned}$$

where α is a positive number. So we obtain

$$G * \left(\frac{\gamma\alpha^2 + \gamma\alpha}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) (x) \geq \frac{\gamma\alpha}{2} u^2(x).$$

Now we let

$$(3.1) \quad \alpha^2 + \alpha = \frac{3-\gamma}{\gamma}.$$

Therefore,

$$(3.2) \quad G * \left(\frac{3-\gamma}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) (x) \geq \frac{\gamma\alpha}{2} u^2(x).$$

Moreover, $\frac{\gamma\alpha}{2}$ is the best constant if and only if

$$\alpha u = u_x \quad \text{in } (-\infty, x) \quad \text{and} \quad -\alpha u = u_x \quad \text{in } (x, \infty),$$

which can be solved as $u = \lambda e^{-\alpha|x-y|}$ for some $\lambda, y \in \mathbb{R}$, α is a positive root determined by (3.1).

The following theorem refines the result first obtained in [17] by admitting a larger class of initial data.

THEOREM 3.1. *Assume $u_0(x) \in H^3(\mathbb{R})$ is odd, and $u'_0(0) < 0$. Then the corresponding solution of (1.1) blows up in finite time.*

Proof. Let $T > 0$ is the maximal time of existence of the solution $u(x, t)$ to (1.1), with initial data u_0 . As one can check, the function $-u(-x, t)$ is also a solution to

(1.1), with initial data $-u_0(-x)$. By this fact, $-u_0(-x) = u_0(x)$ and the uniqueness of solution to (1.1), we get

$$u(x, t) = -u(-x, t) \quad \text{for all } t \in [0, T), \quad x \in \mathbb{R}.$$

Differentiating both sides of (2.4) with respect to variable x , we obtain

$$(3.3) \quad u_{xt} = -\frac{\gamma}{2}u_x^2 + \frac{3-\gamma}{2}u^2 - \gamma uu_{xx} - \left[G * \left(\frac{3-\gamma}{2}u^2 + \frac{\gamma}{2}u_x^2 \right) \right] (x, t).$$

In view of the above analysis, we have

$$u(0, t) = u_{xx}(0, t) = 0, \quad t \in [0, T).$$

Therefore,

$$\frac{du_x(0, t)}{dt} \leq -\frac{\gamma}{2}u_x^2(0, t).$$

From the hypothesis, we have $u'_0(0) < 0$. Therefore, $u_x(0, t) < 0$ for all $t \in [0, T)$. Solving the above inequality, we get

$$-\frac{1}{u_x(0, t)} + \frac{1}{u'_0(0)} \leq -\frac{\gamma}{2}t.$$

We conclude that there exists T and

$$T \leq -\frac{2}{\gamma u'_0(0)}$$

such that $\lim_{t \uparrow T} u_x(0, t) = -\infty$. Thus, we finish the proof by Theorem 2.2. \square

The following theorem gives the blow-up phenomenon via initial potential data y_0 , which reads as follows.

COROLLARY 3.1. *Suppose $y_0 \in H^1(\mathbb{R})$ is an odd function, satisfies $\int_0^\infty e^{-\xi} y_0(\xi) d\xi < 0$. Then the maximum time of existence of the corresponding solution $u(x, t)$ to (1.1) is finite.*

In fact, y_0 is an odd function; we can easily check that u_0 is also odd. At this time, $u'_0(0) = \int_0^\infty e^{-\xi} y_0(\xi) d\xi$. By the hypothesis and theorem above, we get the desired result. We omit the detailed proof for conciseness.

THEOREM 3.2. *Assume $u_0 \in H^3(\mathbb{S})$, $u_0 \not\equiv 0$, $\int_{\mathbb{S}} (u_0^3 + \gamma u_0 u_{0x}^2) dx = 0$, and one of the following conditions is satisfied:*

- (i) $\gamma \in \left(\frac{3 \sinh(\frac{1}{2})}{2 + \sinh(\frac{1}{2})}, 3 \right)$,
 - (ii) $\gamma \in \left(0, \frac{3 \sinh(\frac{1}{2})}{2 + \sinh(\frac{1}{2})} \right]$,
- $$m(0) < -\frac{\sqrt{2}}{2} \alpha \|u_0\|_{H^1(\mathbb{S})}.$$

Then the corresponding solution to (1.1) blows up in finite time.

Proof. By assumption and the invariance property of $E_2(u)$, we have

$$\int_{\mathbb{S}} u^3 + \gamma uu_x^2 dx = \int_{\mathbb{S}} u_0^3 + \gamma u_0 u_{0x}^2 dx = 0.$$

Therefore, $u(x, t)$ must change sign, so there exists at least one zero point on \mathbb{S} . Then for each $t \in [0, T)$, suppose there is a $\xi_t \in [0, 1]$ such that $u(\xi_t, t) = 0$, for $x \in \mathbb{S}$ we have

$$(3.4) \quad u^2(x, t) = \left(\int_{\xi_t}^x u_x dx \right)^2 \leq (x - \xi_t) \int_{\xi_t}^x u_x^2 dx, \quad x \in \left[\xi_t, \xi_t + \frac{1}{2} \right].$$

Thus, the relation above and an integration by parts yield

$$(3.5) \quad \begin{aligned} \int_{\xi_t}^{\xi_t + \frac{1}{2}} u^2 u_x^2 dx &\leq \int_{\xi_t}^{\xi_t + \frac{1}{2}} (x - \xi_t) u_x^2 \left(\int_{\xi_t}^x u_x^2 \right) dx \\ &= \int_{\xi_t}^{\xi_t + \frac{1}{2}} (x - \xi_t) \left(\int_{\xi_t}^x u_x^2 \right) d \left(\int_{\xi_t}^x u_x^2 \right) \\ &\leq \frac{1}{4} \left(\int_{\xi_t}^{\xi_t + \frac{1}{2}} u_x^2 dx \right)^2. \end{aligned}$$

Doing a similar estimate on $[\xi_t + \frac{1}{2}, \xi_t + 1]$, we obtain

$$(3.6) \quad \int_{\mathbb{S}} u^2 u_x^2 dx \leq \frac{1}{4} \left(\int_{\mathbb{S}} u_x^2 dx \right)^2.$$

In view of (3.4), we also have

$$(3.7) \quad \| u(x, t) \|_{L^\infty}^2 \leq \frac{1}{2} \int_{\mathbb{S}} u_x^2 dx.$$

Let

$$q(t) = \int_{\mathbb{S}} u_x^3 dx, \quad t \geq 0.$$

Multiplying both sides of (3.3) with u_x^2 , then integrating by parts, we obtain the equation for $q(t)$ as

$$(3.8) \quad \begin{aligned} \frac{dq(t)}{dt} &= -\frac{\gamma}{2} \int_{\mathbb{S}} u_x^4 dx + \frac{3(3-\gamma)}{2} \int_{\mathbb{S}} u^2 u_x^2 dx - 3 \int_{\mathbb{S}} u_x^2 G * F(u) dx \\ &\leq -\frac{\gamma}{2} \int_{\mathbb{S}} u_x^4 dx - \left(\frac{3\gamma}{4 \sinh(\frac{1}{2})} - \frac{3(3-\gamma)}{8} \right) \left(\int_{\mathbb{S}} u_x^2 \right)^2, \end{aligned}$$

where we use the fact $\frac{1}{2 \sinh(\frac{1}{2})} \leq G(x) \leq \frac{\cosh(\frac{1}{2})}{2 \sinh(\frac{1}{2})}$.

If $\gamma \in (\frac{3 \sinh(\frac{1}{2})}{2 + \sinh(\frac{1}{2})}, 3)$, then $\frac{3\gamma}{4 \sinh(\frac{1}{2})} - \frac{3(3-\gamma)}{8} \geq 0$. On the other hand, we have $\int_{\mathbb{S}} u_x^2 dx \geq \frac{2}{3} \| u_0 \|_{H^1}^2$. Therefore, the above inequality yields

$$(3.9) \quad \begin{aligned} \frac{dq(t)}{dt} dx &= -\frac{\gamma}{2} \int_{\mathbb{S}} u_x^4 dx + \frac{3(3-\gamma)}{2} \int_{\mathbb{S}} u^2 u_x^2 dx - 3 \int_{\mathbb{S}} u_x^2 G * F(u) dx \\ &\leq -\frac{\gamma}{2} \int_{\mathbb{S}} u_x^4 dx - \frac{4}{9} \left(\frac{3\gamma}{4 \sinh(\frac{1}{2})} - \frac{3(3-\gamma)}{8} \right) \| u_0 \|_{H^1}^4. \end{aligned}$$

For simplicity, we set $\mu = \frac{4}{9}(\frac{3\gamma}{4\sinh(\frac{1}{2})} - \frac{3(3-\gamma)}{8})$ and in view of Holder’s inequality, we get

$$(3.10) \quad \begin{aligned} \frac{dq(t)}{dt} &= -\frac{\gamma}{2} \int_{\mathbb{S}} u_x^4 dx + \frac{3(3-\gamma)}{2} \int_{\mathbb{S}} u^2 u_x^2 dx - 3 \int_{\mathbb{S}} u_x^2 G * F(u) dx \\ &\leq -\frac{\gamma}{2} q^{\frac{4}{3}}(t) - \mu \|u_0\|_{H^1}^4. \end{aligned}$$

First, since $q(t) \leq q_0 - \mu \|u_0\|_{H^1}^4 t$, it is easy to find that there exist a time $t_0 \geq 0$ such that $q(t_0) < 0$. Then for $t > t_0$, we have

$$\frac{dq(t)}{dt} \leq -\frac{\gamma}{2} q^{\frac{4}{3}}(t), \quad \text{with } q(t_0) < 0.$$

So the solution to (3.10) satisfies

$$q(t) \leq \left(q^{-\frac{1}{3}}(t_0) + \frac{\gamma}{6}(t - t_0) \right)^{-3},$$

which reaches to $-\infty$ before t arrives at $T_0 = -6q^{-\frac{1}{3}}(t_0)/\gamma + t_0$.

On the other hand, since

$$\left| \int_{\mathbb{S}} u_x^3 dx \right| \leq C(s) \|u_0\|_{H^1}^2 \|u\|_{H^s} \quad \text{for } s \in \left(\frac{3}{2}, 3 \right],$$

then

$$\int_{\mathbb{S}} u_x^3 dx \geq \inf u_x(x, t) \|u_0\|_{H^1}^2$$

shows that $\lim_{t \rightarrow T_0} \inf_{x \in \mathbb{S}} u_x(x, t) = -\infty$.

If $\gamma \in (0, \frac{3\sinh(\frac{1}{2})}{2+\sinh(\frac{1}{2})}]$, we can derive an equation for $m(t)$ as

$$\frac{dm}{dt} = -\frac{\gamma}{2} m^2 + \frac{3-\gamma}{2} u^2(\xi(t), t) - [G * F(u)](\xi(t), t) = 0.$$

Now combining the above equation and (3.2) together, we have

$$(3.11) \quad \frac{dm(t)}{dt} \leq -\frac{\gamma}{2} m^2(t) + \frac{3-\gamma-\gamma\alpha}{4} \|u_0\|_{H^1}^2.$$

If

$$m(0) < - \left(\frac{3-\gamma-\gamma\alpha}{2\gamma} \right)^{\frac{1}{2}} \|u_0\|_{H^1} = -\frac{\sqrt{2}}{2} \alpha \|u_0\|_{H^1},$$

then we can conclude the solution to (3.11) goes to $-\infty$ in finite time. This completes the proof. \square

Remark 3.1. $\int_{\mathbb{S}} u_0 dx = 0$ or $\int_{\mathbb{S}} y_0 dx = 0$ can be a substitute of the condition $\int_{\mathbb{S}} (u_0^3 + \gamma u_0 u_{0x}^2) dx = 0$; the theorem still holds. If $\gamma = 1$, we recover the case for the Camassa–Holm equation [10]. Moreover, one may find Theorem 3.2 is different from Theorem 3.4 in [42]. However, in [39], Yin didn’t consider the case (ii) of this theorem for (1.1).

4. Persistence properties and unique continuation. In this section, we shall investigate persistence properties of the solution to (1.1) in L^∞ -space. The main idea comes from a recent work of Himonas et al. [24].

THEOREM 4.1. *Assume that $u_0(x) \in H^s(\mathbb{R})$ and $s > \frac{3}{2}$, if for some $\theta \in (0, 1)$,*

$$|u_0(x)|, |u_{0x}(x)| \sim O(e^{-\theta x}) \quad \text{as } x \uparrow \infty.$$

Then the corresponding strong solution $u(x, t) \in C([0, T]; H^s(\mathbb{R}))$ to (1.1) satisfies that

$$|u(x, t)|, |u_x(x, t)| \sim O(e^{-\theta x}) \quad \text{as } x \uparrow \infty$$

uniformly in the time interval $[0, T]$ before blow up.

Notation.

$$|u(x)| \sim O(e^{-\theta x}) \quad \text{as } x \uparrow \infty \quad \text{if} \quad \lim_{x \rightarrow \infty} \frac{|u(x)|}{e^{-\theta x}} = L,$$

and

$$|u(x)| \sim o(e^{-\beta x}) \quad \text{as } x \uparrow \infty \quad \text{if} \quad \lim_{x \rightarrow \infty} \frac{|u(x)|}{e^{-\beta x}} = 0.$$

Proof. The proof is organized as follows. Firstly, we will give estimates on $\|u(x, t)\|_\infty$ and $\|u_x(x, t)\|_\infty$. Here $\|\cdot\|_p$ means the L^p norm. Secondly, we use the weight function to obtain the desired result.

Multiplying (1.1) by u^{2n-1} , with $n \in \mathbb{Z}^+$, then integrating both sides with respect to x variable, we can get

$$(4.1) \quad \int_{\mathbb{R}} u^{2n-1} u_t dx + \gamma \int_{\mathbb{R}} u^{2n-1} u u_x dx + \int_{\mathbb{R}} u^{2n-1} \partial_x G * F(u) dx = 0.$$

The first term of the above identity is

$$(4.2) \quad \int_{\mathbb{R}} u^{2n-1} u_t dx = \frac{1}{2n} \frac{d}{dt} \|u(t)\|_{2n}^{2n} = \|u(t)\|_{2n}^{2n-1} \frac{d}{dt} \|u(t)\|_{2n}$$

and

$$(4.3) \quad \left| \int_{\mathbb{R}} u^{2n-1} u u_x dx \right| \leq \|u_x(t)\|_\infty \|u(t)\|_{2n}^{2n}.$$

In view of Holder’s inequality

$$(4.4) \quad \int_{\mathbb{R}} u^{2n-1} \partial_x G * F(u) dx \leq \|u(t)\|_{2n}^{2n-1} \|\partial_x G * F(u)\|_{2n},$$

so

$$(4.5) \quad \frac{d}{dt} \|u(t)\|_{2n} \leq \gamma \|u_x(t)\|_\infty \|u(t)\|_{2n} + \|\partial_x G * F(u)\|_{2n}.$$

In view of the Soblev embedding theorem, then there exists a constant $M > 0$ such that applying Gronwall’s gives us

$$(4.6) \quad \|u(t)\|_{2n} \leq e^{Mt} \left(\|u(0)\|_{2n} + \int_0^t \|\partial_x G * F(u)\|_{2n} d\tau \right).$$

Note that

$$(4.7) \quad \lim_{p \rightarrow \infty} \|f\|_p = \|f\|_\infty \quad \text{when} \quad f \in L^1(\mathbb{R}) \cap L^\infty(\mathbb{R}).$$

Taking limits in (4.6) we obtain

$$(4.8) \quad \|u(t)\|_\infty \leq e^{Mt} \left(\|u(0)\|_\infty + \int_0^t \|\partial_x G * F(u)\|_\infty d\tau \right).$$

Next, we will establish an estimate on $\|u_x(t)\|_\infty$ using the same method as above. Differentiating (1.1) with respect to x variable produces the following equation:

$$(4.9) \quad u_{xt} + \gamma uu_{xx} + \gamma u_x^2 + \partial_x^2 G * F(u) = 0.$$

Multiplying the above identity by u_x^{2n-1} , considering the second term with integration by parts,

$$(4.10) \quad \int_{\mathbb{R}} uu_{xx}u_x^{2n-1} dx = -\frac{1}{2n} \int_{\mathbb{R}} u_x^{2n} u_x dx,$$

so we get

$$(4.11) \quad \int_{\mathbb{R}} u_x^{2n-1} u_{xt} dx + \gamma \int_{\mathbb{R}} u_x^{2n+1} dx - \frac{\gamma}{2n} \int_{\mathbb{R}} u_x^{2n} u_x dx + \int_{\mathbb{R}} u_x^{2n-1} \partial_x^2 G * F(u) dx = 0.$$

Similarly, one can get the inequality

$$(4.12) \quad \frac{d}{dt} \|u_x(t)\|_{2n} \leq 2\gamma \|u_x(t)\|_\infty \|u_x(t)\|_{2n} + \|\partial_x^2 G * F(u)\|_{2n},$$

and therefore, as before, we obtain

$$(4.13) \quad \|u_x(t)\|_{2n} \leq e^{2Mt} \left(\|u_x(0)\|_{2n} + \int_0^t \|\partial_x^2 G * F(u)\|_{2n} d\tau \right).$$

Taking limits in (4.13) to obtain

$$(4.14) \quad \|u_x(t)\|_\infty \leq e^{2Mt} \left(\|u_x(0)\|_\infty + \int_0^t \|\partial_x^2 G * F(u)\|_\infty d\tau \right).$$

In order to get the desired result, we introduce the function $\psi_N(x)$, which is independent on t as follows:

$$(4.15) \quad \psi_N(x) = \begin{cases} 1, & x \leq 0, \\ e^{\theta x}, & x \in (0, N), \\ e^{\theta N}, & x \geq N, \end{cases}$$

where $N \in \mathbb{Z}^+$. From (1.1) we obtain

$$(4.16) \quad \psi_N u_t + \psi_N \gamma uu_x + \psi_N \partial_x G * F(u) = 0,$$

while for (4.9), we get

$$(4.17) \quad u_{xt} \psi_N + \psi_N \gamma uu_{xx} + \psi_N \gamma u_x^2 + \psi_N \partial_x^2 G * F(u) = 0.$$

In order to get the estimate on $u_x\psi_N$, we need to remove the second derivatives, by using integration by parts, we obtain

$$\begin{aligned}
 \left| \int_{\mathbb{R}} \psi_N u u_{xx} (u_x \psi_N)^{2n-1} dx \right| &= \left| \int_{\mathbb{R}} u (\psi_N u_x)^{2n-1} ((u_x \psi_N)_x - u_x \psi'_N) dx \right| \\
 &= \left| \int_{\mathbb{R}} u \left(\frac{(u_x \psi_N)^{2n}}{2n} \right)_x dx - \int_{\mathbb{R}} u u_x \psi'_N (u_x \psi_N)^{2n-1} dx \right| \\
 (4.18) \qquad \qquad \qquad &\leq 2(\|u\|_{\infty} + \|u_x(t)\|_{\infty}) \|u_x \psi_N\|_{2n}^{2n},
 \end{aligned}$$

where we use the fact $0 \leq \psi'_N(x) \leq \psi_N(x)$ a.e. $x \in \mathbb{R}$. Next, we will do estimates on $\|u\psi_N\|_{\infty}$ and $\|u_x\psi_N\|_{\infty}$ step-by-step as before what were done on $\|u\|_{\infty}$ and $\|u_x\|_{\infty}$, we may get

$$\begin{aligned}
 \|u(t)\psi_N\|_{\infty} + \|u_x\psi_N\|_{\infty} &\leq e^{2Mt}(\|u(0)\psi_N\|_{\infty} + \|u_x(0)\psi_N\|_{\infty}) \\
 (4.19) \qquad \qquad \qquad &+ e^{2Mt} \int_0^t (\|\psi_N \partial_x G * F(u)\|_{\infty} \\
 &+ \|\psi_N \partial_x^2 G * F(u)\|_{\infty}) d\tau.
 \end{aligned}$$

On the other hand, computing the integral we see there exists a $c_1 > 0$, depending only on $\theta \in (0, 1)$ such that

$$(4.20) \qquad \qquad \qquad \psi_N(x) \int_{\mathbb{R}} e^{-|x-y|} \frac{1}{\psi_N(y)} dy \leq c_1.$$

Therefore, one gets

$$\begin{aligned}
 |\psi_N \partial_x G * g^2(x)| &\leq \frac{1}{2} \psi_N(x) \int_{\mathbb{R}} e^{-|x-y|} \frac{1}{\psi_N(y)} \psi_N(y) g(y) g(y) dy \\
 &\leq \frac{1}{2} \|\psi_N g\|_{\infty} \|g\|_{\infty} \left(\psi_N(x) \int_{\mathbb{R}} e^{-|x-y|} \frac{1}{\psi_N(y)} dy \right) \\
 (4.21) \qquad \qquad \qquad &\leq c_1 \|\psi_N g\|_{\infty} \|g\|_{\infty}.
 \end{aligned}$$

Similarly,

$$(4.22) \qquad \qquad \qquad |\psi_N \partial_x^2 G * g^2(x)| \leq c_1 \|\psi_N g\|_{\infty} \|g\|_{\infty}.$$

Thus, combining (4.21) and (4.22) with (4.19), it follows that there exists a constant $C_1 = C_1(M, T) > 0$ such that

$$\begin{aligned}
 \|u(t)\psi_N\|_{\infty} + \|u_x\psi_N\|_{\infty} &\leq C_1(\|u(0)\psi_N\|_{\infty} + \|u_x(0)\psi_N\|_{\infty}) \\
 &+ C_1 \int_0^t (\|u(\tau)\|_{\infty} + \|u_x(\tau)\|_{\infty}) \\
 &(\|u(\tau)\psi_N\|_{\infty} + \|u_x(\tau)\psi_N\|_{\infty}) d\tau \\
 (4.23) \qquad \qquad \qquad &\leq C_1 \left(\|u(0)\psi_N\|_{\infty} + \|u_x(0)\psi_N\|_{\infty} \right. \\
 &\left. + \int_0^t (\|u(\tau)\psi_N\|_{\infty} + \|u_x(\tau)\psi_N\|_{\infty}) d\tau \right).
 \end{aligned}$$

Then for any $N \in \mathbb{Z}^+$ and any $t \in [0, T]$, $x > 0$, we have

$$\begin{aligned}
 \|u(t)\psi_N\|_{\infty} + \|u_x\psi_N\|_{\infty} &\leq C_1(\|u(0)\psi_N\|_{\infty} + \|u_x(0)\psi_N\|_{\infty}) \\
 (4.24) \qquad \qquad \qquad &\leq C_1(\|u(0)e^{\theta x}\|_{\infty} + \|u_x(0)e^{\theta x}\|_{\infty}).
 \end{aligned}$$

Taking the limit as N goes to infinity in (4.24), we obtain

$$(4.25) \quad (|u(x, t)e^{\theta x}| + |u_x(x, t)e^{\theta x}|) \leq C_1 (\|u(0)e^{\theta x}\|_\infty + \|u_x(0)e^{\theta x}\|_\infty).$$

This completes the proof. \square

The following result ensures that the only solution which can decay at a determined rate, at any two distinct times, is the trivial solution $u \equiv 0$.

THEOREM 4.2. *Assume that $u_0(x) \in H^s(\mathbb{R})$ and $s > \frac{3}{2}$, satisfies that for some $\delta \in (\frac{1}{2}, 1)$,*

$$|u_0(x)| \sim o(e^{-x}) \quad \text{and} \quad |u_{0x}(x)| \sim O(e^{-\delta x}) \quad \text{as } x \uparrow \infty,$$

$u(x, t) \in C([0, T]; H^s(\mathbb{R}))$ is the corresponding strong solution to (1.1), and there exists $t_1 \in (0, T]$ for some $T > 0$ such that

$$(4.26) \quad |u(x, t_1)| \sim o(e^{-x}) \quad \text{as } x \uparrow \infty,$$

then $u \equiv 0$.

Proof. Integrating (1.1) from 0 to t_1 , we get

$$(4.27) \quad u(x, t_1) - u(x, 0) + \gamma \int_0^{t_1} uu_x d\tau + \int_0^{t_1} \partial_x G * F(u) d\tau = 0.$$

According to the hypothesis, we easily have

$$(4.28) \quad u(x, t_1) - u(x, 0) \sim o(e^{-x}) \quad \text{as } x \uparrow \infty.$$

At the same time, in view of Theorem 4.1 it follows that

$$(4.29) \quad \int_0^{t_1} uu_x d\tau \sim O(e^{-2\delta x}) \quad \text{as } x \uparrow \infty,$$

and so

$$(4.30) \quad \int_0^{t_1} uu_x d\tau \sim o(e^{-x}) \quad \text{as } x \uparrow \infty.$$

If $u(x, t) \neq 0$, the following deduction tells us the last term of (4.27) is infinitesimal with the same order not higher order of e^{-x} . Thus, a contradiction occurs.

$$(4.31) \quad \begin{aligned} \int_0^{t_1} \partial_x G * F(u) d\tau &= \partial_x G * \int_0^{t_1} F(u) d\tau \\ &= \partial_x G * f(x). \end{aligned}$$

However,

$$(4.32) \quad 0 \leq f(x) \sim O(e^{-2\delta x}) \quad \text{so that} \quad f(x) \sim o(e^{-x}) \quad \text{as } x \uparrow \infty.$$

Therefore,

$$(4.33) \quad \partial_x G * f(x) = -\frac{1}{2}e^{-x} \int_{-\infty}^x e^y f(y) dy + \frac{1}{2}e^x \int_x^\infty e^{-y} f(y) dy.$$

From (4.32) it follows that

$$e^x \int_x^\infty e^{-y} f(y) dy = e^x \int_x^\infty e^{-y} o(e^{-y}) dy = o(1) e^x \int_x^\infty e^{-2y} dy \sim o(1) e^{-x} \sim o(e^{-x}).$$

If $f(x) \neq 0$, one has that

$$(4.34) \quad \int_{-\infty}^x e^y f(y) dy \geq C_0 \quad \text{for } x \text{ large enough.}$$

Hence,

$$(4.35) \quad \begin{aligned} -\partial_x G * f(x) &= \frac{1}{2} e^{-x} \int_{-\infty}^x e^y f(y) dy - \frac{1}{2} e^x \int_x^\infty e^{-y} f(y) dy \\ &\geq \frac{C_0}{2} e^{-x} \quad \text{for } x \text{ large enough.} \end{aligned}$$

So a contradiction occurs by combination with (4.27)–(4.30) and (4.35). Thus, $f(x) \equiv 0$ and consequently, $u(x, t) \equiv 0$. The theorem is proved. \square

At the end of this paper, we would like to make a simple comparison between the rod equation and the Camassa–Holm equation. In the whole paper, we do not discuss sufficient conditions to guarantee the global existence of smooth solutions to (1.1) for general γ , which was discussed in [17].

We know for (1.1), $y = u - u_{xx}$ satisfies

$$y_t + \gamma y_x u + 2\gamma y u_x + \frac{3(\gamma - 1)}{2} (u^2)_x = 0.$$

It is a pity that we couldn't obtain the proper particle trajectory line equation to consider this problem just as it was done for the Camassa–Holm equation. Moreover, for $\gamma \neq 1$, few beautiful identities that appeared in the Camassa–Holm equation are obtained for (1.1).

Acknowledgments. The first author is greatly indebted to Prof. Zhou for his constructive suggestions and constant encouragement throughout this work. The authors thank the referees for their constructive and helpful suggestions and comments.

REFERENCES

- [1] T. B. BENJAMIN, J. L. BONA, AND J. J. MAHONY, *Model equations for long waves in nonlinear dispersive systems*, Philos. Trans. R. Soc. Lond. Ser. A, 272 (1972), pp. 47–78.
- [2] A. BRESSAN AND A. CONSTANTIN, *Global dissipative solutions of the Camassa–Holm equation*, Anal. Appl., 5 (2007), pp. 1–27.
- [3] A. BRESSAN AND A. CONSTANTIN, *Global conservative solutions of the Camassa–Holm equation*, Arch. Ration. Mech. Anal., 183 (2007), pp. 215–239.
- [4] R. CAMASSA AND D. HOLM, *An integrable shallow water equation with peaked solitons*, Phys. Rev. Lett., 71 (1993), pp. 1661–1664.
- [5] K. S. CHOU AND C. Z. QU, *Integrable equations arising from the motion of plane curves*, Phys. D, 162 (2002), pp. 9–33.
- [6] A. CONSTANTIN, *Finite propagation speed for the Camassa–Holm equation*, J. Math. Phys., 46 (2005), 023506, 4 pp.
- [7] A. CONSTANTIN, *The trajectories of particles in Stokes waves*, Invent. Math., 166 (2006), pp. 523–535.
- [8] A. CONSTANTIN AND J. ESCHER, *Global existence and blow-up for a shallow water equation*, Ann. Sc. Norm. Super. Pisa Cl. Sci., 26 (1998), pp. 303–328.

- [9] A. CONSTANTIN AND J. ESCHER, *On the blow-up rate and the blow-up set of breaking waves for a shallow water equation*, Math. Z., 233 (2000), pp. 75–91.
- [10] A. CONSTANTIN AND J. ESCHER, *Well-posedness, global existence and blow-up phenomena for a periodic quasi-linear hyperbolic equation*, Comm. Pure Appl. Math., 51 (1998), pp. 475–504.
- [11] A. CONSTANTIN AND J. ESCHER, *Wave breaking for nonlinear nonlocal shallow water equations*, Acta Math., 181 (1998), pp. 229–243.
- [12] A. CONSTANTIN AND J. ESCHER, *Particle trajectories in solitary water waves*, Bull. Amer. Math. Soc., 44 (2007), pp. 423–431.
- [13] A. CONSTANTIN AND D. LANNES, *The hydrodynamical relevance of the Camassa-Holm and Degasperis-Procesi equations*, Arch. Ration. Mech. Anal., (2008), DOI: 10.1007/s00205-008-0128-2.
- [14] A. CONSTANTIN AND L. MOLINET, *Global weak solutions for a shallow water equation*, Comm. Math. Phys., 211 (2000), pp. 45–61.
- [15] A. CONSTANTIN AND L. MOLINET, *Orbital stability of solitary waves for a shallow water equation*, Phys. D, 157 (2001), pp. 75–89.
- [16] A. CONSTANTIN AND W. STRAUSS, *Stability of peakons*, Comm. Pure Appl. Math., 53 (2000), pp. 603–610.
- [17] A. CONSTANTIN AND L. MOLINET, *Stability of a class of solitary waves in compressible elastic rods*, Phys. Lett. A, 270 (2000), pp. 140–148.
- [18] H. H. DAI, *Model equations for nonlinear dispersive waves in a compressible Mooney-Rivlin rod*, Acta Mech., 127 (1998), pp. 193–207.
- [19] H. H. DAI AND Y. HUO, *Solitary shock waves and other travelling waves in a general compressible hyperelastic rod*, R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci., 456 (2000), pp. 331–363.
- [20] Z. GUO, *Blow up, global existence, and infinite propagation speed for the weakly dissipative Camassa-Holm equation*, J. Math. Phys., 49 (2008), 033516.
- [21] D. HENRY, *Compactly supported solutions of the Camassa-Holm equation*, J. Nonlinear Math. Phys., 12 (2005), pp. 342–347.
- [22] D. HENRY, *Compactly supported solutions of a family of nonlinear partial differential equations*, Dyn. Contin. Discrete Impuls. Syst. Ser. A Math. Anal., 15 (2008), pp. 145–150.
- [23] D. HENRY, *Persistence properties for a family of nonlinear partial differential equations*, Nonlinear Anal., DOI:10.1016/J.NA.2008.02.104.
- [24] A. HIMONAS, G. MISIOLK, G. PONCE, AND Y. ZHOU, *Persistence properties and unique continuation of solutions of the Camassa-Holm equation*, Comm. Math. Phys., 271 (2007), pp. 511–512.
- [25] R. S. JOHNSON, *Camassa-Holm, Korteweg-de Vries and related models for water waves*, J. Fluid Mech., 455 (2002), pp. 63–82.
- [26] Y. LI AND P. OLVER, *Well-posedness and blow-up solutions for an integrable nonlinear dispersive model wave equation*, J. Differential Equations, 162 (2000), pp. 27–63.
- [27] Y. LIU AND Y. ZHOU, *Blow-up phenomenon for a periodic rod equation*, J. Phys. A: Math. Theor., 41 (2008), 344013.
- [28] Y. LIU, P. WITTEW, AND Y. ZHOU, *Blow-up phenomenon for a periodic rod equation revisited*, preprint, East China Normal University, Shanghai, 2008.
- [29] H. P. MCKEAN, *Breakdown of a shallow water equation*, Asian J. Math., 2 (1998), pp. 867–874.
- [30] L. MOLINET, *On well-posedness results for Camassa-Holm equation on the line: A survey*, J. Nonlinear Math. Phys., 11 (2004), pp. 521–533.
- [31] O. G. MUSTAFA, *Global conservative solutions of the hyperelastic rod equation*, Int. Math. Res. Not. 2007, (2007), Article id rnm040, 26 pp.
- [32] G. RODRIGUEZ-BLANCO, *On the Cauchy problem for the Camassa-Holm equation*, Nonlinear Anal., 46 (2001), pp. 309–327.
- [33] R. SELIGER, *A note on the breaking of waves*, Proc. Roy. Soc. Lond. Ser. A., 303 (1968), pp. 493–496.
- [34] S. SHKOLLER, *Geometry and curvature of diffeomorphism groups with H^1 metric and mean hydrodynamics*, J. Funct. Anal., 160 (1998), pp. 337–365.
- [35] J. F. TOLAND, *Stokes waves*, Topol. Methods Nonlinear Anal., 7 (1996), pp. 1–48.
- [36] E. WAHLÉN, *A blow-up result for the periodic Camassa-Holm equation*, Arch. Math., 84 (2005), pp. 334–340.
- [37] E. WAHLÉN, *On the blow-up of solutions to a nonlinear dispersive rod equation*, J. Math. Anal. Appl., 323 (2006), pp. 1318–1324.
- [38] Z. XIN AND P. ZHANG, *On the weak solution to a shallow water equation*, Comm. Pure Appl. Math., 53 (2000), pp. 1411–1433.

- [39] Z. YIN, *On the blow-up of solutions of a periodic nonlinear dispersive wave equation in compressible elastic rods*, J. Math. Anal. Appl., 288 (2003), pp. 232–245.
- [40] Y. ZHOU, *Wave breaking for a periodic shallow water equation*, J. Math. Anal. Appl., 290 (2004), pp. 591–604.
- [41] Y. ZHOU, *Wave breaking for a shallow water equation*, Nonlinear Anal., 57 (2004), pp. 137–152.
- [42] Y. ZHOU, *Blow-up phenomenon for a periodic rod equation*, Phys. Lett. A, 353 (2006), pp. 479–486.
- [43] Y. ZHOU, *Local well-posedness and blow-up criteria of solutions for a rod equation*, Math. Nachr., 278 (2005), pp. 1726–1739.
- [44] Y. ZHOU, *Blow-up of solutions to a nonlinear dispersive rod equation*, Calc. Var. Partial Differential Equations, 25 (2006), pp. 63–77.
- [45] Y. ZHOU, *Stability of solitary waves for a rod equation*, Chaos Solitons Fractals, 21 (2004), pp. 977–981.