

# Entropy Production in Nonlinear, Thermally Driven Hamiltonian Systems

Jean-Pierre Eckmann,<sup>1,2</sup> Claude-Alain Pillet,<sup>3,4</sup> and Luc Rey-Bellet<sup>1,5</sup>

*Received October 29, 1998*

---

We consider a finite chain of nonlinear oscillators coupled at its ends to two infinite heat baths which are at different temperatures. Using our earlier results about the existence of a stationary state, we show rigorously that for arbitrary temperature differences and arbitrary couplings, such a system has a unique stationary state. (This extends our earlier results for small temperature differences.) In all these cases, any initial state will converge (at an unknown rate) to the stationary state. We show that this stationary state continually produces entropy. The rate of entropy production is strictly negative when the temperatures are unequal and is proportional to the mean energy flux through the system.

---

**KEY WORDS:** Open systems; nonequilibrium steady states; control theory; entropy production.

## 1. INTRODUCTION

In a recent paper, [EPR], we have studied the existence of a stationary regime in a non-linear non-equilibrium setup. The model considered was that of a chain of  $n$  non-linear oscillators coupled at each of its two ends to heat baths which are infinite systems at two different temperatures.

---

<sup>1</sup>Section de Mathématiques, Université de Genève, CH-1211 Genève 4, Switzerland.

<sup>2</sup>Département de Physique Théorique, Université de Genève, CH-1211 Genève 4, Switzerland.

<sup>3</sup>PHYMAT, Université de Toulon, F-83957 La Garde Cedex, France.

<sup>4</sup>CPT-CNRS Luminy, F-13288 Marseille Cedex 09, France.

<sup>5</sup>Current address: Department of Mathematics, Rutgers University, Hill Center, Rutgers University, Piscataway, New Jersey 08903.

Under suitable conditions which we sketch below, it has been shown that for sufficiently small temperature differences between the baths, the complete system has a *unique* stationary state, and that every initial state converges to it. Of course, this stationary state is not an equilibrium state but a steady state in which supposedly heat flows (on average) from the hot bath to the cold one. The aim of this paper is to show first that this result extends to arbitrary temperature differences and that the heat flux across the chain is positive.

It should be noted that this is *not* a perturbative result. To prove the existence and uniqueness of the steady state for non-linear, boundary driven problems with arbitrary temperature difference is a difficult problem. See [GLP] and [GKI] for similar results for a gas of particles in a box with thermostatted boundary conditions. For other boundary driven models, see [FGS] (a 1-dimensional hard-core gas).

Our model, whose study was started in [EPR], combines several desirable features, while still allowing for a rather complete set of rigorous results. The main features are the property of being fully Hamiltonian (as those studied in [FGS] and [SL]), with non-linear interactions, and a realistic implementation of the retro-action of the chain on the heat baths. In particular, the system is self-regulating and we do not need any Gaussian thermostats [PH, EM, H, CELS, GC1, GC2, G].

The reader should note that in our model the energy of the chain fluctuates wildly in time and there is no external dissipation term which prevents the energy of the chain from diverging to infinity. The baths can exchange energy with the chain. Also, since the potentials are not monotone, several stationary non-equilibrium states could possibly exist, each corresponding for example to one of the extrema of the potential. We show here that on the contrary, there is exactly one stationary state, no matter how large the temperature difference of the baths is.

Once the uniqueness of the steady state is established, we show that, away from equilibrium, there is a stationary, strictly positive heat flow through the chain and the (thermodynamic) entropy production is strictly negative. We also discuss briefly (heuristically) a suitable version of the Gallavotti–Cohen fluctuation relation [ECM, GC1, GC2, G, K, LS] for the entropy production in the context of our model.

## 2. SETUP AND NOTATIONS

To make this paper accessible without the necessity of referring back to [EPR], we introduce again the model. It deals with an anharmonic chain driven at its ends by two heat baths.

The chain consists of  $n$  particles moving in  $\mathbf{R}^d$ , with  $n$  arbitrary but finite, and its dynamics is described by the following Hamiltonian:

$$H_S(q_1, \dots, q_n, p_1, \dots, p_n) = \sum_{j=1}^n \frac{p_j^2}{2} + V(q_1, \dots, q_n)$$

where the potential  $V$  is of the form<sup>6</sup>

$$V(q) = \sum_{j=1}^n U_j^{(1)}(q_j) + \sum_{i=1}^{n-1} U_i^{(2)}(q_i - q_{i+1})$$

We make the following assumptions on the potential  $V$ :

(H1) Behavior at infinity: We assume that  $V$  is of the form

$$V(q) = \frac{1}{2}(q - a, \mathbf{Q}(q - a)) + F(q)$$

where  $\mathbf{Q}$  is a positive definite  $(dn \times dn)$  matrix,  $a$  is a vector, and  $\partial_{q_i^{(v)}} F \in \mathcal{F}$  for  $i = 1, \dots, n$  and  $v = 1, \dots, d$ . Here,  $\mathcal{F}$  denotes the space of those  $\mathcal{C}^\infty$  functions  $F$  on  $\mathbf{R}^{dn}$  for which  $\partial^\alpha F(q)$  is bounded uniformly in  $q \in \mathbf{R}^{dn}$ , for all multi-indices  $\alpha$ .

(H2) Coupling: Each of the  $(d \times d)$  matrices

$$\mathcal{M}_i(x) \equiv \mathbf{D}^2 U_i^{(2)}(x), \quad i = 1, \dots, n-1$$

of second derivatives, is either uniformly positive or negative definite for  $x \in \mathbf{R}^d$ .

Each heat bath is modeled by an infinite dimensional linear Hamiltonian system, which is a scalar field whose dynamics is governed by a  $d$ -dimensional wave equation:

$$H_{\mathbf{B}}(\phi_i, \pi_i) = \frac{1}{2} \int dx (|\nabla \phi_i|^2 + |\pi_i|^2) \quad (2.1)$$

We will denote the heat baths by the subscripts L and R respectively. We couple the L heat bath to the first particle of the chain and the R heat bath to the  $n$ th particle of the chain. We choose a coupling which is linear both

<sup>6</sup> The two-body potential is slightly more restrictive than in [EPR], since we only take functions of the coordinate differences.

in the field variables and in the particle variables. The total Hamiltonian of the system is then given by

$$\begin{aligned}
 H(q, p, \phi_{\mathbf{L}}, \pi_{\mathbf{L}}, \phi_{\mathbf{R}}, \pi_{\mathbf{R}}) &= \sum_{i \in \{\mathbf{L}, \mathbf{R}\}} H_{\mathbf{B}}(\phi_i, \pi_i) + H_{\mathbf{S}}(q, p) \\
 &+ q_1 \cdot \int dx \nabla \phi_{\mathbf{L}}(x) \rho_{\mathbf{L}}(x) + q_n \cdot \int dx \nabla \phi_{\mathbf{R}}(x) \rho_{\mathbf{R}}(x)
 \end{aligned} \tag{2.2}$$

We consider the heat baths at positive temperatures  $T_{\mathbf{L}}$  and  $T_{\mathbf{R}}$  respectively, i.e., we will assume that the initial conditions of the heat baths are distributed according to the Gaussian measure with mean zero and covariance  $(\cdot, \cdot)_i T_i$ , where  $(\cdot, \cdot)_i$  is the scalar product defined by the quadratic form (2.1),  $i \in \{\mathbf{L}, \mathbf{R}\}$ .

The following reduction to (essentially only) the variables of the small system is explained in detail in [EPR]: We integrate out the variables of the baths and project the dynamics on the variables of the chain. This leads to integro-differential stochastic equations. Under suitable assumptions on the coupling functions  $\rho_{\mathbf{L}}, \rho_{\mathbf{R}}$ , they can be expressed as Markovian equations upon introducing auxiliary variables  $r_{i,m}$  with  $i \in \{\mathbf{L}, \mathbf{R}\}$ , and  $m = 1, \dots, M$ . At the end, one obtains (see [EPR] for details) the following system of stochastic differential equations:

$$\begin{aligned}
 dq_j(t) &= p_j(t) dt, \quad j = 1, \dots, n \\
 dp_1(t) &= -\nabla_{q_1} V(q(t)) dt + \sum_{m=1}^M r_{\mathbf{L},m}(t) dt \\
 dp_j(t) &= -\nabla_{q_j} V(q(t)) dt, \quad j = 2, \dots, n-1 \\
 dp_n(t) &= -\nabla_{q_n} V(q(t)) dt + \sum_{m=1}^M r_{\mathbf{R},m}(t) dt \tag{2.3} \\
 dr_{\mathbf{L},m}(t) &= -\gamma_{\mathbf{L},m} r_{\mathbf{L},m}(t) dt + \lambda_{\mathbf{L},m}^2 \gamma_{\mathbf{L},m} q_1(t) dt - \lambda_{\mathbf{L},m} \sqrt{2\gamma_{\mathbf{L},m} T_{\mathbf{L}}} dw_{\mathbf{L},m}(t) \\
 dr_{\mathbf{R},m}(t) &= -\gamma_{\mathbf{R},m} r_{\mathbf{R},m}(t) dt + \lambda_{\mathbf{R},m}^2 \gamma_{\mathbf{R},m} q_n(t) dt \\
 &- \lambda_{\mathbf{R},m} \sqrt{2\gamma_{\mathbf{R},m} T_{\mathbf{R}}} dw_{\mathbf{R},m}(t), \quad m = 1, \dots, M
 \end{aligned}$$

which defines a Markov diffusion process on  $\mathbf{R}^{2d(n+M)}$ . Each  $w_{\mathbf{L},m}$  and  $w_{\mathbf{R},m}$  is a standard  $d$ -dimensional Brownian motion.

**Remark 2.1.** In Eqs. (B.3) the variables  $r_{\mathbf{L},m}, r_{\mathbf{R},m}$  describe both the (random) forces exerted by the heat bath on the chain and the (dissipative) forces due to the retroaction of the heat baths on the chain. The

fact that the variables  $r_{L,m}, r_{R,m}$  obey Markovian stochastic differential equations is a consequence of our choice of coupling functions  $\rho_L, \rho_R$ . In fact, these functions are chosen in such a way that the random forces exerted by the heat baths are not white noises but have covariances which are (sums of) exponentials. Together with the fluctuation theorem relating these random forces with the dissipative forces, one obtains Markovian differential equations on the phase space consisting of the physical variables  $p, q$ , augmented by the auxiliary variables  $r_{L,m}, r_{R,m}$ .

**Remark 2.2.** If the temperatures of the two heat baths are the same, i.e., if  $T_L = T_R$ , the stationary state of the Markov process which solves (2.3) can be written explicitly. It is given by the generalized Gibbs state

$$\mu(dr, dq, dp) = \mu_{T_L, T_L}(dr, dq, dp) = Z^{-1} e^{-G^{(0)}(r, q, p)/T_L} dr dq dp \quad (2.4)$$

where the “energy”  $G^{(0)}$  is given by

$$G^{(0)}(r, q, p) = H_S(q, p) + \sum_{m=1}^M \left( \frac{r_{L,m}^2}{2\lambda_{L,m}^2} + \frac{r_{R,m}^2}{2\lambda_{R,m}^2} - q_1 \cdot r_{L,m} - q_n \cdot r_{R,m} \right) \quad (2.5)$$

The marginal of this measure on the physical phase space is given by

$$\nu(dq, dp) = \int \mu(dr, dq, dp) = \frac{1}{Z'} e^{-H_{\text{eff}}/T_L} dq dp$$

where the effective Hamiltonian  $H_{\text{eff}}$  is given by

$$H_{\text{eff}}(q, p) = H_S(q, p) - \frac{1}{2} q_1^2 \sum_{m=1}^M \lambda_{L,m}^2 - \frac{1}{2} q_n^2 \sum_{m=1}^M \lambda_{R,m}^2 \equiv \frac{1}{2} p^2 + V_{\text{eff}}(q) \quad (2.6)$$

It can be seen from (2.6) that the coupling between the chain and the heat baths induces a renormalization of the potential  $V(q)$ . In particular, because of Condition (H1), if the coupling constants  $\lambda_{im}$  are too large,  $V_{\text{eff}}(q)$  is not confining any more, and the measure  $\mu$  is not a probability measure, but only  $\sigma$ -finite. In the sequel we require the following:

(H3) The coupling constants  $\lambda_{im}, i \in \{L, R\}, m \in \{1, \dots, M\}$  are such that

$$\lim_{|q| \rightarrow \infty} V_{\text{eff}}(q) = \infty$$

### 3. UNIQUENESS OF THE STATIONARY STATE

In [EPR] we proved, under the conditions (H1)–(H3), the existence of an invariant measure for any temperature  $T_L$ ,  $T_R$ , but the uniqueness was shown only for small temperature differences. In this section, we extend the uniqueness result to *arbitrary* temperature differences.

The uniqueness will follow from a dynamical argument: we will show that the Markov process is transitive. This is done using a (well-known) relationship between stochastic differential equations and control theory (see, e.g., [Ku] and references therein).

We explain the method for a general stochastic differential equation of the form

$$dx(t) = b(x(t)) dt + \sigma dw(t) \quad (3.1)$$

where  $x \in \mathbf{R}^k$ ,  $b(x)$  is a  $\mathcal{C}^\infty$  vector field,  $w(t)$  is a standard  $\ell$ -dimensional Brownian motion, and  $\sigma$  is a  $k \times \ell$  matrix. We assume that the vector field  $b(x)$  is such that (3.1) has a unique solution for all  $t > 0$ . One then replaces  $dw(t)$  in (3.1) by  $u(t) dt$ . The function  $u(t) = (u_1(t), \dots, u_\ell(t))$  is called a control. One obtains the system of ordinary differential equations

$$\dot{x} = b(x(t)) + \sigma u(t) \quad (3.2)$$

The correspondence between the two systems is established by the following result of Stroock and Varadhan [SV]. We fix an arbitrary time  $\tau > 0$ . Let  $\mathcal{U}$  denote the set of piecewise constant functions  $u: [0, \tau] \rightarrow \mathbf{R}^\ell$ . Let  $\mathcal{W}$  be the set of all continuous functions  $\varphi$  from  $[0, \tau]$  to  $\mathbf{R}^k$  equipped with the uniform topology and let  $\mathcal{W}_x = \{\varphi \in \mathcal{W} : \varphi(0) = x\}$ . We denote  $\xi_x$  the diffusion process defined by (3.1) with initial condition  $\xi_x(0) = x$ . Then the path  $\xi_x$  belongs almost surely to  $\mathcal{W}_x$ . The support of this diffusion process  $\xi_x$  on  $[0, \tau]$  is, by definition, the smallest closed subset  $\mathcal{S}_x$  of  $\mathcal{W}_x$  such that

$$\mathbf{P}[\xi_x \in \mathcal{S}_x] = 1 \quad (3.3)$$

where  $\mathbf{P}$  is the probability induced by the Brownian motion  $w$ . We denote by  $\varphi_x^u: [0, \tau] \rightarrow \mathbf{R}^k$  the solution of the differential equations (3.2) with control  $u$  and initial condition  $x$ . We next consider the notion of accessibility. Let  $x$  and  $y$  be two points in  $\mathbf{R}^k$ . The point  $y$  is called *accessible* from  $x$  at time  $\tau$  ( $\tau > 0$ ) if there is a control  $u$  such that  $\varphi_x^u(\tau) = y$ . The set of all points which are accessible from  $x$  at time  $\tau$  is denoted  $Y_\tau(x)$ .

**Theorem 3.1** [SV]. One has

$$\mathcal{L}_x = \overline{\{\varphi_x^u : u \in \mathcal{U}\}} \tag{3.4}$$

for all  $x \in \mathbf{R}^k$ , and

$$\text{supp } P(\tau, x, \cdot) = \overline{Y_\tau(x)}$$

for all  $x \in \mathbf{R}^k$  and  $\tau > 0$ , where  $P(\tau, x, dy)$  denotes the transition probability of the process  $\xi_x$ .

**Remark.** The first statement is explicit in [SV] and the second is a straightforward consequence of the first.

The main technical result of this section is the following

**Theorem 3.2.** The control system associated with the stochastic differential equation (2.3) is strongly completely controllable, i.e.,  $\overline{Y_\tau(x)} = \mathbf{R}^{2d(n+M)}$ , for all  $x = (q, p, r_L, r_R)$  and all  $\tau > 0$ .

**Remark 3.3.** We will combine this result with Theorem 3.1 and hypoellipticity to show that the invariant measure has a smooth, strictly positive density.

**Remark 3.4.** It should be noted that strong complete controllability (SCC) can not be deduced from Hörmander’s hypoellipticity condition alone. (See, e.g., [IK] for examples of hypoelliptic diffusions with two invariant measures). Therefore, Theorem 3.2 contains additional information. Various sufficient conditions for SCC have been expressed in terms of differential geometry in [Ku], but these are *not* applicable to Eqs. (2.3).

*Proof of Theorem 3.2.* We will show the strong complete controllability of the control problem associated with (2.3) by an explicit approach using the requirement of effective coupling (condition (H2)) of the chain.

We reconsider the stochastic differential equation (2.3). Following the procedure described above we replace the Brownian motions  $w_{im}(t)$  by controls  $u_{im}(t)$  in (2.3), and rewrite the system thus obtained as a system of second (and first) order equations. This leads to:

$$\dot{r}_{L,m} = -\gamma_{L,m} r_{L,m} + \lambda_{L,m}^2 \gamma_{L,m} q_1 + u_{L,m}, \quad m = 1, \dots, M \tag{3.5}$$

$$\ddot{q}_1 = -\nabla_{q_1} V(q) + \sum_{m=1}^M r_{L,m} \tag{3.6}$$

$$\ddot{q}_j = -\nabla_{q_j} V(q), \quad j = 2, \dots, n-1 \quad (3.7)$$

$$\ddot{q}_n = -\nabla_{q_n} V(q) + \sum_{m=1}^M r_{\mathbf{R}, m} \quad (3.8)$$

$$\dot{r}_{\mathbf{R}, m} = -\gamma_{\mathbf{R}, m} r_{\mathbf{R}, m} + \lambda_{\mathbf{R}, m}^2 \gamma_{\mathbf{R}, m} q_n + u_{\mathbf{R}, m}, \quad m = 1, \dots, M \quad (3.9)$$

Here, we have absorbed the constants in front of the Brownian motion in (2.3) into the controls  $u_{im}(t)$ . We will only consider controls  $u$  of class  $\mathcal{C}^\infty$ . Any such function can be uniformly approximated, on any compact interval, by a piecewise constant function. Since a simple Gronwall estimate of Eq. (3.2) shows that

$$\sup_{t \in [0, \tau]} |\varphi_x^u(t) - \varphi_x^v(t)| \leq C(\tau) \sup_{t \in [0, \tau]} |u(t) - v(t)|$$

holds with a constant  $C(\tau)$  depending on the model, but not on  $u$  and  $v$ , we conclude that

$$\{\varphi_x^u(\tau) : u \in \mathcal{C}^\infty(\mathbf{R})\} \subset \overline{Y_\tau(x)}$$

The proof of Theorem 3.2 will now be done in two parts:

*Part 1: Boundary Control of the Chain.* We start by considering the auxiliary problem of controlling a chain of  $n$  oscillators by the motion of the two ends of the chain.

The differential equation reads

$$\ddot{q}_j = f_j(q_{j-1}, q_j, q_{j+1}), \quad j = 1, \dots, n \quad (3.10)$$

where  $q \equiv (q_1, \dots, q_n)$  is the dynamical variable, whereas  $q_0 \equiv u_L$  and  $q_{n+1} \equiv u_R$  are the control variables. The smooth functions  $f_j$  are given by

$$f_j(x, y, z) \equiv -(\nabla U_j^{(1)})(y) - (\nabla U_j^{(2)})(y-z) + (\nabla U_{j-1}^{(2)})(x-y)$$

Here, we define  $U_0^{(2)}(x) = U_n^{(2)}(x) = x^2/2$ .

Note that by Condition (H2), the functions  $\nabla U_j^{(2)}$  are diffeomorphisms. It follows that the equation  $w = f_j(x, y, z)$  can be solved for  $z$ : There exist smooth functions  $g_j$  such that  $w = f_j(x, y, g_j(x, y, w))$  for all  $x, y, w \in \mathbf{R}^d$ . Consequently the differential equation (3.10) is equivalent to the equation

$$q_{j+1} = g_j(q_{j-1}, q_j, \ddot{q}_j), \quad j = 1, \dots, n \quad (3.11)$$



Obviously, for given  $q_0$  and  $q_1$ , this equation can be solved inductively, and in a unique way. To express this solution, let us introduce some notation. For a smooth function  $\varphi$  and an integer  $\alpha$ , we shall denote the collection of the first  $\alpha$  derivatives of  $\varphi$  by  $\varphi^{[\alpha]} \equiv (\varphi, \dot{\varphi}, \dots, \varphi^{(\alpha)})$ . We also set  $p_j \equiv \dot{q}_j$  for  $j = 1, \dots, n$ . For  $q_0 = \zeta$  and  $q_1 = \eta$ , the inductive solution of Eq. (3.11) reads

$$\begin{aligned} u_L &= \zeta && \equiv G_0(\zeta^{[0]}) \\ q_1 &= \eta && \equiv G_1(\eta^{[0]}) \\ p_1 &= \dot{q}_1 && \equiv G_2(\eta^{[1]}) \\ q_2 &= g_1(q_0, q_1, \dot{p}_1) && \equiv G_3(\zeta^{[0]}, \eta^{[2]}) \\ p_2 &= \dot{q}_2 && \equiv G_4(\zeta^{[1]}, \eta^{[3]}) \\ q_3 &= g_2(q_1, q_2, \dot{p}_2) && \equiv G_5(\zeta^{[2]}, \eta^{[4]}) \\ &\vdots && \vdots \\ u_R &= g_n(q_{n-1}, q_n, \dot{p}_n) && \equiv G_{2n+1}(\zeta^{[2n-2]}, \eta^{[2n]}) \end{aligned}$$

We can organize the  $2n + 2$  maps  $G_J$  into a map  $G: \mathbf{R}^{4nd} \rightarrow \mathbf{R}^{4nd}$  in the following way: Denote by  $(a, b)$  a point of  $\mathbf{R}^{4nd}$ , with  $a \equiv (a_0, \dots, a_{2n-2}) \in \mathbf{R}^{(2n-1)d}$  and  $b \equiv (b_0, \dots, b_{2n}) \in \mathbf{R}^{(2n+1)d}$ . With  $a^{[\alpha]} \equiv (a_0, \dots, a_\alpha)$  and  $b^{[\alpha]} \equiv (b_0, \dots, b_\alpha)$ , define  $G(a, b) \equiv (a, \hat{G}(a, b))$  where

$$\hat{G}(a, b) \equiv (G_1(b^{[0]}), G_2(b^{[1]}), G_3(a^{[0]}, b^{[2]}), \dots, G_{2n+1}(a^{[2n-2]}, b^{[2n]}))$$

We have proved that, if  $(u_L, q_1, \dots, q_n, u_R)$  is a solution of Eq. (3.10) on the time interval  $I \subset \mathbf{R}$ , then

$$(u_L^{[2n-2]}, q_1, \dot{q}_1, \dots, q_n, \dot{q}_n, u_R) = G(u_L^{[2n-2]}, q_1^{[2n]}) \tag{3.12}$$

holds on  $I$ . A simple consequence of this fact is that  $G$  is a bijection. Indeed, repeated differentiation of Eq. (3.10) gives

$$\begin{aligned} q_1^{(2)} &= f_1(u_L, q_1, q_2) = F_2(u_L^{[0]}, q_1, q_2) \\ q_1^{(3)} &= \partial_t q_1^{(2)} = F_3(u_L^{[1]}, q_1, \dot{q}_1, q_2, \dot{q}_2) \\ &\vdots \\ q_1^{(2\alpha)} &= \partial_t q_1^{(2\alpha-1)} = F_{2\alpha}(u_L^{[2\alpha-2]}, q_1, \dot{q}_1, \dots, q_{\alpha+1}) \\ q_1^{(2\alpha+1)} &= \partial_t q_1^{(2\alpha)} = F_{2\alpha+1}(u_L^{[2\alpha-1]}, q_1, \dot{q}_1, \dots, q_{\alpha+1}, \dot{q}_{\alpha+1}) \\ &\vdots \\ q_1^{(2n)} &= \partial_t q_1^{(2n-1)} = F_{2n}(u_L^{[2n-2]}, q_1, \dot{q}_1, \dots, u_R) \end{aligned}$$

and thus we find another functional relation  $q_1^{[2n]} = \hat{F}(u_L^{[2n-2]}, q_1, \dot{q}_1, \dots, q_n, \dot{q}_n, u_R)$  for its solutions. It immediately follows that the map  $F: (a, b) \mapsto (a, \hat{F}(a, b))$  satisfies

$$\begin{aligned} F \circ G(u_L^{[2n-2]}, q_1^{[2n]}) &= (u_L^{[2n-2]}, q_1^{[2n]}) \\ G \circ F(u_L^{[2n-2]}, q_1, \dot{q}_1, \dots, u_R) &= (u_L^{[2n-2]}, q_1, \dot{q}_1, \dots, u_R) \end{aligned}$$

on every solution of Eq. (3.10). Since  $u_L^{[2n-2]}(0)$  and either  $q_1^{[2n]}(0)$  or  $(q(0), \dot{q}(0), u_L(0))$  can be prescribed arbitrarily, we conclude that  $F = G^{-1}$ .

To solve our control problem it suffices to remark that the set of solutions  $(u_L, q, u_R)$  satisfying  $(u_L^{[2n-2]}(t_0), q_1(t_0), \dot{q}_1(t_0), \dots, u_R(t_0)) = (a, b)$  is identical with the set of solutions satisfying  $(u_L^{[2n-2]}(t_0), q_1^{[2n]}(t_0)) = F(a, b)$ . Since for  $\tau > 0$  and arbitrary  $(a, b), (a', b') \in \mathbf{R}^{4nd}$  one can find functions  $u_L$  and  $q_1$  for which

$$\begin{aligned} (u_L^{[2n-2]}(0), q_1^{[2n]}(0)) &= F(a, b) \\ (u_L^{[2n-2]}(\tau), q_1^{[2n]}(\tau)) &= F(a', b') \end{aligned}$$

we see that the system (3.10) is strongly controllable.

*Part 2: Completion of the Proof of Theorem 3.2.* We reduce the problem of Eqs. (3.5)–(3.9) to the case dealt with in Part 1, by introducing the auxiliary variables  $q_0$  and  $q_{n+1}$ . Recalling the definition  $U_0^{(2)}(x) = U_n^{(2)}(x) = x^2/2$ , we can rewrite the control problem associated with our stochastic differential equation as

$$\begin{aligned} \ddot{q}_j &= f_j(q_{j-1}, q_j, q_{j+1}), \quad j = 1, \dots, n \\ \sum_{m=1}^M r_{L,m} &= q_1 - q_0 \\ \sum_{m=1}^M r_{R,m} &= q_{n+1} - q_n \end{aligned} \tag{3.13}$$

$$\begin{aligned} u_{L,m} &= \dot{r}_{L,m} + \gamma_{L,m} r_{L,m} - \lambda_{L,m}^2 \gamma_{L,m} q_1, \quad m = 1, \dots, M \\ u_{R,m} &= \dot{r}_{R,m} + \gamma_{R,m} r_{R,m} - \lambda_{R,m}^2 \gamma_{R,m} q_n, \quad m = 1, \dots, M \end{aligned}$$

with the boundary conditions

$$\begin{aligned} (r_L(0), q_1(0), \dot{q}_1(0), \dots, q_n(0), \dot{q}_n(0), r_R(0)) &= x, \\ (r_L(\tau), q_1(\tau), \dot{q}_1(\tau), \dots, q_n(\tau), \dot{q}_n(\tau), r_R(\tau)) &= y \end{aligned}$$

The equation  $\sum_{m=1}^M r_{L,m} = q_1 - q_0$  serves to compensate the term  $q_1 - q_0$  produced when differentiating  $f_1$  in (3.13). Given the boundary data for  $r_L$ ,  $q_1$ ,  $q_n$  and  $r_R$  we obtain boundary values for  $q_0$  and  $q_{n+1}$ . We can thus control the first equation by our previous result. This gives us  $q_0, \dots, q_{n+1}$ . Selecting arbitrary functions  $r_{L,2}, \dots, r_{L,M}$  and  $r_{R,2}, \dots, r_{R,M}$  satisfying the corresponding boundary data, we define

$$r_{L,1} \equiv q_1 - q_0 - \sum_{m=2}^M r_{L,m},$$

$$r_{R,1} \equiv q_{n+1} - q_n - \sum_{m=2}^M r_{R,m}$$

These two functions will also satisfy the boundary conditions. Finally we use the last two sets of equations to determine the control variables  $u_{L,m}$  and  $u_{R,m}$ . This concludes the proof of Theorem 3.2.

**Remark 3.5.** It is obvious from the proof of Theorem 3.2 that this theorem is valid under much weaker conditions than those given in (H1). It is enough to require that the stochastic differential equation (2.3) has a unique solution for all  $t > 0$ . In particular we do not need to restrict ourselves to potentials which are “quadratic at infinity” as required in the proof of the existence of the invariant measure.

The main result of this section is:

**Theorem 3.6.** If Conditions (H1)–(H3) are satisfied, the Markov process which solves (2.3) has a unique invariant measure  $\mu = \mu_T$ . The measure  $\mu$  has a  $\mathcal{C}^\infty$  density  $\rho(r, q, p)$ . This density is an exponentially decaying, strictly positive function of  $r$ ,  $q$ , and  $p$ . The invariant measure is ergodic and mixing.

**Remark 3.7.** In fact, combining this result with information from [EPR], one can show that

$$\rho(r, q, p) = f(r, q, p) \exp(-G^{(0)}(r, q, p)/T^*)$$

where  $G^{(0)}$  was defined in (2.5), and  $T^* = \max(T_L, T_R)$ . The function  $f$  is in the Schwartz space  $\mathcal{S}$  when  $T_L \neq T_R$ , (and is a constant otherwise).

*Proof.* The proof is a combination of Theorem 3.1 and Theorem 3.2 with results in [EPR]. The existence of the invariant measure  $\mu$  is proven in [EPR, Theorem 2.1]. Furthermore, using hypoellipticity, we showed

that the density  $\rho$  of  $\mu$  is  $\mathcal{C}^\infty$ . Also, the transition probabilities  $P(t, x, dy)$  have a smooth density,  $p$ , defined by  $P(t, x, dy) = p(t, x, y) dy$  with  $p(t, x, y) \in \mathcal{C}^\infty((0, \infty), \mathbf{R}^{2d(n+M)}, \mathbf{R}^{2d(n+M)})$ .

We next show that the support of  $\mu$  is all of the extended phase space  $X \equiv \mathbf{R}^{2d(n+M)}$ . In Theorem 3.2, we have seen that (2.3) is strongly completely controllable. By Theorem 3.1, we conclude that the support of  $P(\tau, x, \cdot)$  is the whole phase space for every  $x \in X$  and all  $\tau > 0$ . Therefore, we have, for all  $t > 0$ , all  $x$ , and all open sets  $Y$ , the inequality

$$P(t, x, Y) > 0$$

Since  $\mu(Y) = \int \mu(dx) P(t, x, Y)$  (because  $\mu$  is invariant), we conclude that  $\text{supp } \mu = X$  and thus the density  $\rho$  is Lebesgue almost everywhere positive.

We next show that  $\rho(x) > 0$ , for all  $x$ , by assuming the contrary and deriving a contradiction. Assume that there is a  $y$  for which  $\rho(y) = 0$ . By the invariance of the measure we have, for any  $t > 0$ ,

$$0 = \rho(y) = \int dx \rho(x) p(t, x, y)$$

This implies  $p(t, x, y) = 0$  for Lebesgue almost all  $x$ . Since the transition kernel  $p$  is smooth, we conclude that the function  $p(t, \cdot, y)$  is identically zero every  $t > 0$ . On the other hand, since  $p$  is the kernel of a strongly continuous semigroup, we also have  $p(t, x, y) \rightarrow \delta(x - y)$  as  $t \rightarrow 0$ . This is a contradiction, and we have shown  $\rho(y) > 0$ , for all  $y \in X$ .

We next show uniqueness. We have just shown that every invariant measure must have a smooth, strictly positive density. Since every ergodic component is mutually singular to any other, the invariant measure is unique (and ergodic). The property of mixing of the invariant measure has been deduced from uniqueness in the proof of [EPR, Theorem 3.9]. This concludes the proof of Theorem 3.6.

**Remark 3.8.** We proved in [EPR, Lemma 3.7] that the density  $\rho = \rho_T$  is a real analytic function of  $\zeta = (T_L - T_R)/(T_L + T_R)$ . In particular, this yields the standard perturbative results near equilibrium ( $\zeta = 0$ ).

#### 4. TIME-REVERSAL, ENERGY FLUX, AND ENTROPY PRODUCTION

In this section, we ask questions which are intimately related to the Hamiltonian nature of our model. After introducing appropriate notation,

we introduce time reversal, and draw some consequences. In particular, we are able to show that the system exhibits non-zero mean energy flux as soon as  $T_L \neq T_R$ , and we relate the flux to the entropy production.

#### 4.1. Notation

It will be useful to streamline the notation. It is convenient to introduce first  $\rho_{L,m} = (\lambda_{L,m} \gamma_{L,m}^{1/2})^{-1} r_{L,m}$  and similarly for the  $r_{R,m}$ . Then the equations of motion are

$$\begin{aligned} dq_j(t) &= p_j(t) dt, \quad j = 1, \dots, n \\ dp_1(t) &= -\nabla_{q_1} V(q(t)) dt + \sum_{m=1}^M \lambda_{L,m} \gamma_{L,m}^{1/2} \rho_{L,m}(t) dt \\ dp_j(t) &= -\nabla_{q_j} V(q(t)) dt, \quad j = 2, \dots, n-1 \\ dp_n(t) &= -\nabla_{q_n} V(q(t)) dt + \sum_{m=1}^M \lambda_{R,m} \gamma_{R,m}^{1/2} \rho_{R,m}(t) dt \end{aligned} \quad (4.1)$$

$$\begin{aligned} d\rho_{L,m}(t) &= -\gamma_{L,m} \rho_{L,m}(t) dt + \lambda_{L,m} \gamma_{L,m}^{1/2} q_1(t) dt - \sqrt{2} T_L^{1/2} dw_{L,m}(t) \\ d\rho_{R,m}(t) &= -\gamma_{R,m} \rho_{R,m}(t) dt + \lambda_{R,m} \gamma_{R,m}^{1/2} q_n(t) dt \\ &\quad - \sqrt{2} T_R^{1/2} dw_{R,m}(t), \quad m = 1, \dots, M \end{aligned}$$

We can write this system in vector notation: We write the Hamiltonian of the chain (the small system) as

$$H_S(q, p) = \frac{p^2}{2} + V(q)$$

with  $q, p \in \mathbf{R}^{nd}$ . The two reservoirs, L and R, are described by the variables  $\rho = (\rho_L, \rho_R) \in \mathbf{R}^{Md} \oplus \mathbf{R}^{Md}$ . The “energy” of the complete system, i.e., chain and reservoirs, is then given by  $G^{(0)}(r, q, p) = G^{(1)}(\rho, q, p)$ , where now

$$G^{(1)}(\rho, q, p) = H_S(q, p) + \frac{1}{2} \rho \cdot \Gamma \rho - q \cdot \lambda \Gamma^{1/2} \rho$$

Here,  $\Gamma = \Gamma_L \oplus \Gamma_R$ , where  $\Gamma_i$  is the diagonal  $(M \times M)$  matrix  $\text{diag}(\gamma_{i,1}, \dots, \gamma_{i,M})$ , with  $i \in \{L, R\}$ . Note that by assumption, the  $\gamma$ 's are all strictly positive. We also define  $A$  as the  $(2Md \times nd)$  matrix given by

$$q \cdot \lambda \rho = q_1 \cdot A_L \rho_L + q_n \cdot A_R \rho_R = q_1 \sum_{m=1}^M \lambda_{L,m} \rho_{L,m} + q_n \sum_{m=1}^M \lambda_{R,m} \rho_{R,m}$$

With these notations, the equations of motion can be written as:

$$dq = \nabla_p G^{(1)} dt = p dt$$

$$dp = -\nabla_q G^{(1)} dt = -(\nabla_q V(q) - \Lambda \Gamma^{1/2} \rho) dt$$

$$d\rho = -\nabla_\rho G^{(1)} dt - (2T^{1/2}) dw = -(\Gamma \rho - \Gamma^{1/2} \Lambda^T q) dt - (2T^{1/2}) dw$$

Here  $w = w_L \oplus w_R = (w_{L,1}, \dots, w_{L,M}, w_{R,1}, \dots, w_{R,M})$  is a  $2Md$ -dimensional standard Brownian motion, and  $T$  is the  $(2M \times 2M)$  diagonal temperature matrix

$$T = \text{diag}(T_L, \dots, T_L, T_R, \dots, T_R)$$

It is useful to introduce the (final!) change of variables  $s = \rho - F^T q$ , where  $F = \Lambda \Gamma^{-1/2}$ . In terms of these variables, one can introduce the effective potential

$$V_{\text{eff}}(q) = V(q) - \frac{1}{2} q \cdot \Lambda \Lambda^T q \quad (4.2)$$

and the “energy” is now  $G(s, q, p) = G^{(1)}(\rho, q, p)$  with

$$G(s, q, p) = \frac{1}{2} p^2 + V_{\text{eff}} + \frac{1}{2} s \cdot \Gamma s \quad (4.3)$$

Finally, with the adjoint change in the derivatives  $\nabla_q \rightarrow \nabla_q - F \nabla_s$ , the equations of motion become

$$dq = \nabla_p G dt = p dt$$

$$dp = -(\nabla_q - F \nabla_s) G dt = -(\nabla_q V_{\text{eff}}(q) - F \Gamma s) dt \quad (4.4)$$

$$ds = -(\nabla_s + F^T \nabla_p) G dt - (2T^{1/2}) dw = -(\Gamma s + F^T p) dt - (2T^{1/2}) dw$$

**Notation.** In the sequel, we shall write  $G_p$  for  $\nabla_p G$  and  $G_q$  for  $\nabla_q G$  (these are vectors with  $nd$  components), and  $G_s$  for  $\nabla_s G$  (this is a vector with  $2Md$  components).

The generator  $L$  of the diffusion process takes, in the variables  $s, q, p$ , the form

$$L = \nabla_s \cdot T \nabla_s - G_s \cdot \nabla_s + (G_p \cdot \nabla_q - G_q \cdot \nabla_p) + ((F G_s) \cdot \nabla_p - G_p \cdot F \nabla_s) \quad (4.5)$$

If  $f$  is a function on the phase space  $X$ , we let

$$S^t f(x) = (e^{Lt} f)(x) = \int f(\xi_x(t)) d\mathbf{P}(w)$$

The associated Fokker–Planck operator  $L^T$  is the adjoint of  $L$  in the space  $L^2(\mathbf{R}^{d(2M+2n)}, dx)$ , i.e.,

$$L^T = \nabla_s \cdot T \nabla_s + \nabla_s \cdot G_s - (G_p \cdot \nabla_q - G_q \cdot \nabla_p) - ((FG_s) \cdot \nabla_p - G_p \cdot F \nabla_s) \quad (4.6)$$

**Remark 4.1.** The density  $\rho$  of the invariant measure is the (unique) normalized solution of the equation

$$L^T \rho = 0$$

## 4.2. The Entropy Production $\sigma$

We now establish the relation between the energy flux and the entropy production. Since we are dealing with a Hamiltonian setup, the energy flux is naturally defined by the time derivative of the mean evolution  $S^t$  of the effective energy,  $H_{\text{eff}}(q, p) = p^2/2 + V_{\text{eff}}(q)$ . Differentiating, we get from the equations of motion

$$\begin{aligned} \partial_t S^t H_{\text{eff}} &= S^t L H_{\text{eff}} \\ L H_{\text{eff}} &= p \cdot (-\nabla_q V_{\text{eff}} + F \Gamma s) + \nabla_q V_{\text{eff}} \cdot p = p \cdot F \Gamma s \end{aligned}$$

We define the total flux by  $\Phi = p \cdot F \Gamma s$ , and inspection of the definition of  $F$  and  $\Gamma$  leads to the identification of the flux at the left and right ends of the chain:

$$\Phi = \Phi_L + \Phi_R$$

with

$$\begin{aligned} \Phi_L &= p_1 \cdot A_L \Gamma_L^{1/2} s_L \\ \Phi_R &= p_n \cdot A_R \Gamma_R^{1/2} s_R \end{aligned}$$

Note that  $A_L \Gamma_L^{1/2} s_L$  is the net force exerted by the left bath on the chain. Therefore,  $\Phi_L = p_1 \cdot A_L \Gamma_L^{1/2} \rho_L - L q_1 \cdot A_L^2 q_1 / 2$  is, up to a time-derivative which vanishes in the stationary state, the total power dissipated by the left bath. A similar interpretation holds for  $\Phi_R$ . Furthermore, observe that

$$\langle \Phi \rangle_\mu = 0 \quad (4.7)$$

where, generally,

$$\langle f \rangle_\mu \equiv \int \mu(dx) f(x)$$

The Equation (4.7) holds because  $\Phi = LH_{\text{eff}}$  and  $L^T\mu = 0$ .

We next proceed to define the entropy production in the setting of our model. Since we have been able to identify the energy flux on the ends of the chain, we *define* the (thermodynamic) entropy production  $\sigma$  by

$$\sigma = \frac{\Phi_L}{T_L} + \frac{\Phi_R}{T_R} = p \cdot FT^{-1}\Gamma s \quad (4.8)$$

We refer to [CL] and references therein for a detailed discussion of the various types of entropy production in non-equilibrium stationary states. In Subsection 4.4, we will explain, in the context of our model, the relationship between the entropy production  $\sigma$  and the Gibbs entropy.

### 4.3. Time-Reversal, Generalized Detailed Balance Condition, and Negativity of the Entropy Production

**Definition.** We define the “time-reversal” map  $J$  by  $(Jf)(s, q, p) = f(s, q, -p)$ . This map is the projection onto the space of the  $s, q, p$  of the time-reversal of the Hamiltonian flow (on the full phase space of chain plus baths) defined by the original problem (2.2).

**Notation.** To obtain simple formulas for the entropy production  $\sigma$  we write the strictly positive density  $\rho$  of the invariant measure  $\mu$  as

$$\rho = J e^{-R} e^{-\varphi} \quad (4.9)$$

where we have introduced the quantity

$$R = R(s) = \frac{1}{2}s \cdot \Gamma T^{-1}s \quad (4.10)$$

Let  $L^*$  denote the adjoint of  $L$  in the space  $\mathcal{H}_\mu = L^2(X, d\mu)$  associated with the invariant measure  $\mu$  with density  $\rho$ . In terms of the adjoint  $L^T$  on  $L^2(X, ds dq dp)$ , we have

$$L^* = \rho^{-1} L^T \rho \quad (4.11)$$



Let  $L_\lambda = L + \lambda\sigma$ , where  $\lambda \in \mathbf{R}$ . (This definition is suggested by the paper [K], see below.) We have the following important symmetry property:

**Theorem 4.2.** One has the operator identity

$$J e^{-J\varphi} (L_\lambda)^* e^{J\varphi} J = L_{1-\lambda} \tag{4.12}$$

In particular, one has

$$J e^{-J\varphi} L^* e^{J\varphi} J - L = \sigma \tag{4.13}$$

**Remark 4.3.** This relation may be viewed as a generalization to non-equilibrium of the detailed balance condition (at equilibrium, one has  $JL^*J - L = 0$ ).

**Remark 4.4.** Recently, a lot of interest has been generated in the wake of papers by Gallavotti and Cohen, [GC1, GC2, G, and references therein], in which intriguing relations for the fluctuations of the entropy production have been found. These papers dealt first with numerical experiments by [ECM], which were then abstracted to the general context of dynamical systems. In further work, these ideas have been successfully applied to thermostatted systems modeling non-equilibrium problems. In the papers [K] and [LS] these ideas have been further extended to non-equilibrium models described by stochastic dynamics. In the context of our model, the setup is as follows: One considers the observable

$$W(t) = \int_0^t dt \sigma(\xi_x(\tau))$$

By ergodicity,  $\lim_{t \rightarrow \infty} t^{-1} W(t) = \langle \sigma \rangle_\mu$ , for  $\mu$ -almost all  $x$ . We are interested in the rate function  $\hat{e}$  for the large deviations of  $W(t)/(t \langle \sigma \rangle_\mu)$ , and want to argue (heuristically) that it satisfies

$$\hat{e}(w) - \hat{e}(-w) = -w \langle \sigma \rangle_\mu \tag{4.14}$$

In particular this means that at equal temperatures, when  $\langle \sigma \rangle_\mu = 0$ , the fluctuations are symmetric around the mean 0, while at unequal temperatures, the odd part is linear in  $w$  and proportional to the mean entropy production. This is the celebrated Gallavotti–Cohen fluctuation theorem.

The rate function  $\hat{e}$  is characterized by the relation

$$\inf_{w \in I} \hat{e}(w) = - \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbf{P} \left( \frac{W(t)}{t \langle \sigma \rangle_\mu} \in I \right)$$

Under suitable conditions it can be expressed as the Legendre transform of the function

$$e(\lambda) \equiv - \lim_{t \rightarrow \infty} t^{-1} \log \langle e^{-\lambda W(t)} \rangle_{\mu}$$

Formally,  $-e(\lambda)$  can be represented as the maximal eigenvalue of  $L_{\lambda}$ . Observing now the relation (4.12), one sees immediately that

$$e(\lambda) = e(1 - \lambda) \quad (4.15)$$

This in turn implies (4.14). A rigorous derivation of the program outlined above lacks several technical ingredients, in particularly more spectral information about  $L_{\lambda}$  seems to be necessary.

The relation (4.12) has a number of other consequences which we enumerate now, before going to the proof of Theorem 4.2. It allows to prove that the entropy production is negative in our model and it yields an interesting symmetry relation (see Theorem 4.6).

**Proposition 4.5.** One has the following identities (between functions):

$$L\varphi = \sigma + |T^{1/2} \nabla_s \varphi|^2 \quad (4.16)$$

$$L^* J\varphi = -\sigma - |T^{1/2} \nabla_s J\varphi|^2 \quad (4.17)$$

Here,  $|f|^2 \equiv f \cdot f$ .

**Theorem 4.6.** In the stationary state  $\mu$  the entropy production satisfies the identity:

$$\langle \sigma \rangle_{\mu} = -\langle |T^{1/2} \nabla_s \varphi|^2 \rangle_{\mu} = -\langle |T^{1/2} \nabla_s J\varphi|^2 \rangle_{\mu} \leq 0 \quad (4.18)$$

**Remark 4.7.** In Subsection 4.5, we will show that the heat flux is non-zero for unequal temperatures by showing that the entropy production in the stationary state satisfies:

$$\langle \sigma \rangle_{\mu} = 0$$

if and only if  $T_L = T_R$ .

The remainder of this subsection is devoted to the proofs of Theorem 4.2, Proposition 4.5, and Theorem 4.6.

*Proof of Theorem 4.2.* We show the identity

$$e^R J L^T J e^{-R} = L + \sigma \tag{4.19}$$

Starting with the relation

$$e^R \nabla_s e^{-R} = \nabla_s - (\nabla_s R) = \nabla_s - T^{-1} G_s$$

we get, using the definition of  $L^T$ ,

$$\begin{aligned} e^R J L^T J e^{-R} &= \nabla_s \cdot T \nabla_s - G_s \cdot \nabla_s + (G_p \cdot \nabla_q - G_q \cdot \nabla_p) \\ &\quad + (F G_s \cdot \nabla_p - p \cdot F \nabla_s) + G_p \cdot F T^{-1} G_s \end{aligned}$$

Note that the sum of all the terms except the last equals  $L$ , while the last equals

$$G_p \cdot F T^{-1} G_s = p \cdot F T^{-1} G_s = \frac{p_1 F_L \Gamma_L^S L}{T_L} + \frac{p_n F_R \Gamma_R^S R}{T_R} = \sigma$$

We have thus shown (4.19). Combining (4.19) with the expressions (4.11) and (4.9) for  $L^*$  and  $\rho$  we obtain the identity:

$$L + \sigma = e^R J L^T J e^{-R} = e^R J e^{-R} e^{-J\varphi} L^* e^{J\varphi} e^R J e^{-R} = J e^{-J\varphi} L^* e^{J\varphi} J \tag{4.20}$$

which is (4.12) for  $\lambda = 0$ , i.e., Eq.(4.13). Observing now that  $J\sigma J = -\sigma$ , we obtain

$$J e^{-J\varphi} L_\lambda^* e^{J\varphi} J = L_{1-\lambda}$$

and thus conclude the proof of Theorem 4.2.

*Proof of Proposition 4.5.* From (4.12) we obtain the identity

$$J L^* J = e^\varphi (L + \sigma) e^{-\varphi} \tag{4.21}$$

A straightforward computation shows that, for any smooth function  $f$ , we have the following operator identity

$$e^f L e^{-f} = L - 2(\nabla_s f) \cdot T \nabla_s - (L f) + |T^{1/2}(\nabla_s f)|^2 \tag{4.22}$$

Applying (4.22) with  $f = \varphi$  we obtain from Eq. (4.21) the operator identity

$$J L^* J = L - 2(\nabla_s \varphi) \cdot T \nabla_s - (L \varphi) + |T^{1/2}(\nabla_s \varphi)|^2 + \sigma \tag{4.23}$$

Since  $L^* = \rho^{-1}L^T\rho$  and  $L^T\rho = 0$  we have  $L^*1 = 0$ . Applying the operator identity (4.23) to the function 1 and noting that  $J1 = 1$  we get

$$0 = JL^*J1 = -L\varphi + |T^{1/2}(\nabla_s\varphi)|^2 + \sigma \quad (4.24)$$

and this is the identity (4.16). With this, (4.23) simplifies to

$$JL^*J = L - 2(\nabla_s\varphi) \cdot T\nabla_s \quad (4.25)$$

Applying the operator identity (4.25) to the function  $\varphi$ , and using (4.24) we get

$$JL^*J\varphi = L\varphi - 2|T^{1/2}(\nabla_s\varphi)|^2 = \sigma - |T^{1/2}(\nabla_s\varphi)|^2$$

or, equivalently,

$$L^*J\varphi = -\sigma - |T^{1/2}(\nabla_sJ\varphi)|^2$$

which proves (4.17). With this we have concluded the proof of Proposition 4.5.

*Proof of Theorem 4.6.* Theorem 4.6 is a simple consequence of Proposition 4.5. From (4.16), using the invariance of the measure  $\mu$ , we get

$$\begin{aligned} \langle \sigma \rangle_\mu &= \langle L\varphi \rangle_\mu - \langle |T^{1/2}(\nabla_s\varphi)|^2 \rangle_\mu \\ &= -\langle |T^{1/2}(\nabla_s\varphi)|^2 \rangle_\mu \end{aligned}$$

which yields the first equality in (4.18). The second inequality is obtained in the same way using (4.17). We have

$$\begin{aligned} \langle \sigma \rangle_\mu &= -\langle L^*J\varphi \rangle_\mu - \langle |T^{1/2}(\nabla_sJ\varphi)|^2 \rangle_\mu \\ &= -\langle |T^{1/2}(\nabla_sJ\varphi)|^2 \rangle_\mu \end{aligned} \quad (4.26)$$

where the last equality in (4.26) follows from the identity

$$\langle L^*J\varphi \rangle_\mu = \int dx (L^T\rho J\varphi) = \int dx \rho J\varphi (L1) = 0$$

It is obvious from (4.26) that the entropy production in the stationary state is a non-positive quantity and this concludes the proof of Theorem 4.6.

**Other Observables for the Entropy Production.** The analysis done for the entropy production  $\sigma$  can be repeated for other observables, (see also [LS] for a similar generalization). A family of such observables can be obtained by replacing the conjugation operator  $e^{J\varphi}J$  of (4.12) by any conjugation operator of the form  $e^f e^{J\varphi}J$ , where  $f=f(q, p)$  satisfies  $Jf=f$ .<sup>7</sup> Interesting examples are obtained when one considers the energy flux between position  $j$  and the  $j+1$  on the chain,  $j=1, \dots, n-1$ . To this end we write the Hamiltonian  $H_{\text{eff}}$  as follows:

$$H_{\text{eff}}(q, p) = \sum_{i=1}^n H_i(q, p)$$

where

$$H_1(q, p) = \frac{p_1^2}{2} + U_1^{(1)}(q_1) - \frac{1}{2} q_1^2 \sum_{m=1}^M \lambda_{L,m}^2 + \frac{1}{2} U_1^{(2)}(q_1 - q_2)$$

$$H_i(q, p) = \frac{p_i^2}{2} + U_i^{(1)}(q_i) + \frac{1}{2} U_{i-1}^{(2)}(q_{i-1} - q_i) + \frac{1}{2} U_i^{(2)}(q_i - q_{i+1}), \quad i=2, \dots, n-1$$

$$H_n(q, p) = \frac{p_n^2}{2} + U_n^{(1)}(q_n) - \frac{1}{2} \sum_{m=1}^M \lambda_{R,m}^2 q_n^2 + \frac{1}{2} U_{n-1}^{(2)}(q_{n-1} - q_n)$$

For any  $j=1, \dots, n-1$ , we choose  $f = -S_j$ , where

$$S_j(q, p) = \frac{1}{T_L} \sum_{i=1}^j H_i(q, p) + \frac{1}{T_R} \sum_{i=j+1}^n H_i(q, p)$$

We write now the invariant density  $\rho$  as

$$\rho = J e^{-R} e^{-S_j} e^{-\psi_j}$$

i.e.,  $\psi_j = \varphi - S_j$ . Variants of computations done above show that we have the operator identity, similar to (4.20):

$$e^{S_j} e^R J L^T J e^{-R} e^{-S_j} = L + \sigma_j$$

<sup>7</sup> The operators  $e^f e^{J\varphi}J$  are all formally selfadjoint on  $\mathcal{H}_\mu$ .

where  $\sigma_j$  is given by the relation

$$\sigma_j = \sigma - LS_j \quad (4.27)$$

since  $S_j$  does not depend on the variable  $s$ . Our choice of  $S_j$  has been made in such a way that

$$\sigma_j = \left( \frac{1}{T_L} - \frac{1}{T_R} \right) (p_j - p_{j+1}) \cdot \nabla U^{(2)}(q_j - q_{j+1})$$

i.e.,  $\sigma_j$  is the energy flux between position  $j$  and  $j+1$  on the chain multiplied by the difference between the inverse temperatures of the heat baths. Using next that  $JS_j = S_j$  one derives easily a relation corresponding to (4.12), namely,

$$J e^{-J\psi_j} (L + \lambda \sigma_j)^* e^{J\psi_j} J = L + (1 - \lambda) \sigma_j \quad (4.28)$$

We have thus found  $n-1$  “entropy productions”  $\sigma_j$ , which, under the assumptions made for  $\sigma$ , satisfy a fluctuation theorem. Note that these entropy productions are all different observables, but, because of Eq. (4.27), the expectations of  $\sigma$  and  $\sigma_j$  in the stationary state  $\mu$  coincide.

#### 4.4. Relation with the Gibbs Entropy

We give now a second proof of the negativity of entropy production in our model using the Gibbs entropy.

Let  $\nu_0$  be a probability measure in the variables  $x = (s, q, p)$  and let  $\nu_t$  denote the time evolution of  $\nu_0$  given by

$$\nu_t(dx) = \int \nu_0(dx') P(t, x', x)$$

Because of the properties of the transition probabilities  $P(t, x', x)$  proven in Sect. 3, we see that  $\nu_t$  is a probability measure (for any  $t > 0$ ) with a smooth positive density denoted  $f_t$  in the sequel. The time evolution of  $f$  is then given by the equation

$$\partial_t f_t = L^T f_t$$

We define the Gibbs entropy as

$$S(f) = - \int dx f(x) \log f(x)$$

and we compute next the entropy change in time. We get:

$$\begin{aligned} \partial_t S(f_t) &= - \int dx (\partial_t f_t) (1 + \log f_t) \\ &= -(L^T f_t, (1 + \log f_t)) \\ &= -(f_t, L \log f_t) \\ &= -(f_t, f_t^{-1} L f_t) + (f_t, |T^{1/2} \nabla_s \log f_t|^2) \equiv X_1 \end{aligned}$$

The last term is the (additional) contribution from the second order derivative (in  $s$ ) appearing in  $L$  when it acts on  $\log f$ . We can transform  $X_1$  further by writing it as

$$X_1 = -(L^T 1, f_t) + (f_t, |T^{1/2} \nabla_s \log f_t|^2) \tag{4.29}$$

Since  $L^T 1 = \text{Tr } \Gamma$  the first term in (4.29) is equal to  $-\text{Tr } \Gamma$ . We use the definition (4.10) of  $R$  and the analog of (4.9) to define  $\varphi_t$ :

$$J e^{-R} e^{-\varphi_t} = f_t \tag{4.30}$$

Since  $J R J = R$ , we see that  $-\log f_t = R + J \varphi_t$ . Expanding the second term in (4.29) we obtain

$$\begin{aligned} &(f_t, |T^{1/2} \nabla_s \log f_t|^2) \\ &= (f_t, |T^{-1/2} \Gamma s|^2) + (f_t, |T^{+1/2} \nabla_s J \varphi_t|^2) + 2(f_t, \Gamma s \cdot \nabla_s J \varphi_t) \end{aligned} \tag{4.31}$$

Since we have the relation

$$\nabla_s f_t = -f_t (\nabla_s R + \nabla_s J \varphi_t)$$

we obtain

$$f_t \nabla_s J \varphi_t = -f_t \nabla_s R - \nabla_s f_t$$

Using this and integrating by parts, we rewrite the third term in (4.31) as

$$\begin{aligned} 2(f_t, \Gamma s \cdot \nabla_s J \varphi_t) &= -2 \int dx \Gamma s \cdot (f_t \nabla_s R + \nabla_s f_t) \\ &= -2(f_t, |T^{-1/2} \Gamma s|^2) + 2 \text{Tr } \Gamma \end{aligned}$$

Altogether we obtain

$$\partial_t S(f_t) = \text{Tr } \Gamma - (f_t, |T^{-1/2} \Gamma s|^2) + (f_t, |T^{1/2} \nabla_s J \varphi_t|^2)$$

Using the identity,

$$LR = \text{Tr } \Gamma - |T^{-1/2} \Gamma_s|^2 - \sigma$$

we obtain finally,

$$\partial_t S(f_t) = \int dx f_t \sigma + \int dx f_t LR + \int dx f_t |T^{1/2} \nabla_s J \varphi_t|^2 \quad (4.32)$$

In line with the ideas of [CL], we can write this last identity in the form:

$$\partial_t S(f_t) - \langle \sigma \rangle_{v_t} = \langle |T^{1/2} \nabla_s J \varphi_t|^2 \rangle_{v_t} + \partial_t \langle R \rangle_{v_t}$$

This shows that the (rearrangement) entropy produced in addition to the thermodynamic entropy  $\sigma$  is a positive quantity, up to a (time-) boundary term. Also note that if  $v_t = \mu$ , i.e., if the system is in the stationary state, then we get the identity

$$0 = \langle \sigma \rangle_\mu + \langle |T^{1/2} J \nabla_s \varphi|^2 \rangle_\mu \quad (4.33)$$

which we already found in Theorem 4.6.

#### 4.5. Strict Positivity of the Heat Flux

We first show that the thermodynamic entropy production, as defined in (4.8), is negative in our model. As an immediate consequence we will show that, in the stationary state, energy is flowing from the hotter heat bath to the colder one.

**Theorem 4.8.** The entropy production  $\sigma$  satisfies:

$$\langle \sigma \rangle_\mu = 0$$

if and only if  $T_L = T_R$ .

*Proof.* Note that if  $T_L = T_R$ , then  $\sigma = \partial_t S^i H_{\text{eff}} / T_L |_{t=0}$  and therefore  $\langle \sigma \rangle_\mu = 0$ . We will show that if  $T_L \neq T_R$ , then  $\langle \sigma \rangle_\mu \neq 0$ . We will proceed by assuming the converse, namely  $\langle \sigma \rangle_\mu = 0$ , and produce a contradiction. The assumption implies by (4.18) that  $\langle |T^{1/2} \nabla_s \varphi|^2 \rangle_\mu = 0$ . Since  $\rho$  is positive, this means that  $\nabla_s \varphi = 0$ , and therefore  $\varphi$  does not depend on the  $s$  variables. From (4.16) we obtain

$$0 = -L\varphi + |T^{1/2} \nabla_s \varphi|^2 + \sigma = -L\varphi + \sigma$$



Using the definition of  $L$  and  $\sigma$  and the fact that  $\varphi$  does not depend on  $s$ , we obtain the equation

$$0 = (p \cdot \nabla_q \varphi - (\nabla_q V_{\text{eff}}) \cdot \nabla_p \varphi) + FTS \cdot (\nabla_p \varphi - T^{-1}p)$$

Since  $\varphi$  does not depend on  $s$  we get

$$\begin{aligned} p \cdot \nabla_q \varphi - (\nabla_q V_{\text{eff}}) \cdot \nabla_p \varphi &= 0 \\ \nabla_{p_1} \varphi &= T_L^{-1} p_1 \\ \nabla_{p_n} \varphi &= T_R^{-1} p_n \end{aligned} \tag{4.34}$$

We will show that this system of linear equations has no solution unless  $T_L = T_R$ . To see this we consider the system of equations

$$\begin{aligned} p \cdot \nabla_q \varphi - (\nabla_q V_{\text{eff}}) \cdot \nabla_p \varphi &= 0 \\ \nabla_{p_1} \varphi &= T_L^{-1} p_1 \end{aligned} \tag{4.35}$$

This system has a solution which is given by  $H_{\text{eff}}(q, p)/T_L$ . We claim that this is the unique solution (up to an additive constant) of (4.35).

If this holds true, then the only solution of (4.34) is given by  $H_{\text{eff}}(q, p)/T_L$  and this is incompatible with the third equation in (4.34) when  $T_L \neq T_R$ .

Since (4.35) is a linear inhomogeneous equation, it is enough to show that the only solutions of the homogeneous equation

$$\begin{aligned} p \cdot \nabla_q \varphi - (\nabla_q V_{\text{eff}}) \cdot \nabla_p \varphi &= 0 \\ \nabla_{p_1} \varphi &= 0 \end{aligned} \tag{4.36}$$

are the constant functions. Since  $\nabla_{p_1} \varphi = 0$ ,  $\varphi$  does not depend on  $p_1$ , we conclude that the first equation in (4.36) reads

$$p_1 \cdot \nabla_{q_1} \varphi + f_1(q_1, \dots, q_n, p_2, \dots, p_n) = 0$$

where  $f_1$  does not depend on the variable  $p_1$ . Thus we see that  $\nabla_{q_1} \varphi = 0$  and therefore  $\varphi$  does not depend on the variable  $q_1$  either. By the first equation in (4.36) we now get

$$-\nabla_{q_1} U_1^{(2)}(q_1 - q_2) \cdot \nabla_{p_2} \varphi + f_2(q_2, \dots, q_n, p_2, \dots, p_n) = 0$$

where  $f_2$  does not depend on  $p_1$  and  $q_1$ . By condition (H2) we see that  $\nabla_{p_2} \varphi = 0$  and hence  $f$  does not depend on  $p_2$ . Iterating the above procedure

we find that the only solutions of (4.36) are the constant functions. This concludes the proof of Theorem 4.8.

**Corollary 4.9.** The stationary state  $\mu = \mu_{T_L, T_R}$  produces a non-vanishing mean heat flux in the direction from the hotter to the colder reservoir. The mean heat flux vanishes only if  $T_L = T_R$ .

*Proof.* The entropy production  $\sigma$  is given by

$$\sigma = \frac{\Phi_L}{T_L} + \frac{\Phi_R}{T_R}$$

where  $\Phi_L$  is the energy flow from the left heat bath to the chain and similarly for  $\Phi_R$ . In the stationary state we have, by (4.7),

$$\langle \Phi_L + \Phi_R \rangle_\mu = 0$$

and therefore

$$\langle \Phi_L \rangle_\mu = -\langle \Phi_R \rangle_\mu$$

We obtain from Theorem 4.8, for  $T_L \neq T_R$ :

$$0 > \langle \sigma \rangle_\mu = \left( \frac{1}{T_L} - \frac{1}{T_R} \right) \langle \Phi_L \rangle_\mu$$

If, say,  $T_L > T_R$ , we get  $\langle \Phi_L \rangle_\mu > 0$  and thus energy flows from the hotter to the cooler heat bath.

## ACKNOWLEDGMENTS

We would like to thank G. Ben Arous, Ch. Mazza, H. Spohn, D. Stroock, A.-S. Sznitman, and L. E. Thomas for encouragements and helpful discussions, and V. Zagrebnov for useful questions and comments. This work was supported in part by the Fonds National Suisse.

## REFERENCES

- [CELS] N. I. Chernov, G. L. Eyink, J. L. Lebowitz, and Ya. G. Sinai, Steady-state electric conduction in the periodic Lorentz gas, *Commun. Math. Phys.* **154**:569–601 (1993).
- [CL] N. I. Chernov and J. L. Lebowitz, Stationary nonequilibrium states in boundary driven Hamiltonian systems: Shear flow. Preprint (1997).
- [ECM] D. J. Evans, E. G. D. Cohen, and G. P. Morriss, Probability of second Law violations in shearing steady flows, *Phys. Rev. Lett.* **71**:2401–2404 (1993).

- [EM] D. Evans and G. Morriss, *Statistical Mechanics of Nonequilibrium Liquids* (Academic Press, New York, 1990).
- [EPR] J.-P. Eckmann, C.-A. Pillet, and L. Rey-Bellet, Non-equilibrium statistical mechanics of anharmonic chains coupled to two heat baths at different temperatures, *Commun. Math. Phys.* **201**:657–697 (1995).
- [FGS] J. Farmer, S. Goldstein, and E. R. Speer, Invariant states of a thermally conducting barrier, *J. Stat. Phys.* **34**:263–277 (1984).
- [G] G. Gallavotti, Chaotic hypothesis and universal large deviations properties, *Doc. Math. J. DMV Extra Volume ICM I*:65–93 (1998).
- [GC1] G. Gallavotti and E. G. D. Cohen, Dynamical ensembles in nonequilibrium statistical mechanics, *Phys. Rev. Lett.* **74**:2694–2697 (1995).
- [GC2] G. Gallavotti and E. G. D. Cohen, Dynamical ensembles in stationary states, *J. Stat. Phys.* **80**:931–970 (1995).
- [GKI] S. Goldstein, C. Kipnis, and N. Ianiro, Stationary states for a mechanical system with stochastic boundary conditions, *J. Stat. Phys.* **41**:915–939 (1985).
- [GLP] S. Goldstein, J. L. Lebowitz, and E. Presutti, Stationary states for a mechanical system with stochastic boundaries, in *Random Fields* (Colloquia Mathematicae Societatis János Bolyai, Amsterdam, North-Holland 27, 1981).
- [H] W. G. Hoover, *Computational Statistical Mechanics* (Elsevier, 1991).
- [IK] K. Ishihara and H. Kunita, A classification of the second order degenerate elliptic operators and its probabilistic characterization, *Z. Wahrsch. Verw. Geb.* **39**:235–254 (1974).
- [K] J. Kurchan, Fluctuation theorem for stochastic dynamics, *J. Phys. A* **31**:3719–3729 (1998).
- [Ku] H. Kunita, Supports of diffusion processes and controllability problems, in *Proc. Intern. Symp. SDE Kyoto.*, pp. 163–185 (1976).
- [LS] J. L. Lebowitz and H. Spohn, The Gallavotti-Cohen fluctuation theorem for stochastic dynamics. Preprint (1998).
- [PH] H. A. Posch and W. G. Hoover, Equilibrium and non equilibrium Lyapunov spectra for dense fluids and solids, *Phys. Rev. A* **39**:2175–2188 (1989).
- [SL] H. Spohn and J. L. Lebowitz, Stationary non-equilibrium: States of infinite harmonic systems, *Commun. Math. Phys.* **54**:97 (1977).
- [SV] D. W. Stroock and S. R. S. Varadhan, On the support of diffusion processes with applications to the strong maximum principle, in *Proc. 6th Berkeley Symp. Math. Stat. Prob. III*, pp. 333–368 (1972).