

High-Order Taylor–Galerkin Methods for Linear Hyperbolic Systems

A. SAFJAN AND J. T. ODEN

*Texas Institute for Computational and Applied Mathematics, The University of Texas at Austin, 2.400 Taylor Hall
Austin, Texas 78712*

Received August 23, 1994

A new family of high-order Taylor–Galerkin schemes is presented for the analysis of first-order linear hyperbolic systems. The schemes are unconditionally stable which makes them very attractive to use in conjunction with adaptive *hp*-finite element methods for spatial approximation. © 1995 Academic Press, Inc.

1. INTRODUCTION

The Taylor–Galerkin schemes represent generalizations of the Lax–Wendroff algorithm to weak formulations valid for finite element approximations and have been used effectively for producing second- and third-order temporal approximations [3, 6, 2]. However, no procedures for extending these techniques to temporal approximations of arbitrary order appear to be available.

In the present work, we present a new family of stable high order Taylor–Galerkin (TG) methods for the numerical solutions of first-order hyperbolic systems. It is shown that for first-order systems, a multi-stage process can be used that produces schemes of order $2s$ for s -stages, each of which involves the solution of a second-order system of elliptic equations. A detailed stability analysis is provided for the case of linear systems which establishes choices of parameters that result in unconditionally stable schemes.

Advantages of the new family of TG methods developed in this work are enumerated as follows:

1. the TG methods lead to a series of well-posed problems which can be solved using *hp*-adaptive finite element methods with very high accuracy and high rates of convergence.
2. the TG methods have excellent accuracy; some of them possess a built-in temporal error estimate,
3. with a proper choice of parameters, TG methods can be unconditionally stable for certain classes of problems,
4. the TG methods can be classified as semi-implicit schemes, so that they are generally more efficient than fully implicit schemes and provide a natural splitting of the operator that preserves accuracy, and,

5. in the case in which the underlying system of conservation laws is symmetrizable, the resulting system of linear algebraic equations is symmetric and positive definite, which is the best possible setting for most iterative solvers.

The plan of the presentation is as follows. In Section 2, we discuss the properties of the linear hyperbolic system of the form

$$\mathbf{U}_t + \mathbf{A}\mathbf{U} = \mathbf{0} \tag{1.1}$$

in which \mathbf{A} is a self-adjoint operator. For focus, we consider equations of linear acoustics. In Section 3, we formally derive high-order temporal difference schemes for system of conservation laws of the form

$$\mathbf{U}_t + \sum_i \mathbf{F}^i(\mathbf{U})_{,x_i} = \mathbf{0} \tag{1.2}$$

and we derive the corresponding weak formulation. In particular, we do not assume linearity of (1.2). Section 4 is devoted to linear stability analysis and Section 5 gives some *a priori* error estimates. Finally, numerical results are presented in Section 6 and the coefficients of the new TG schemes are listed in the Appendix.

2. MODEL HYPERBOLIC SYSTEM: EQUATIONS OF LINEAR ACOUSTICS

As a starting point we consider the conservation equations of isentropic, compressible inviscid flow in the form

$$\begin{aligned} \rho(u_{k,t} + u_k u_{l,t}) + p_{,k} &= 0 \\ \rho_{,t} + (\rho u_k)_{,k} &= 0 \end{aligned} \tag{2.1}$$

$$p/p_0 = (\rho/\rho_0)^\gamma, \quad k, l = 1, \dots, \nu; \nu = 2 \text{ or } 3,$$

where ρ is the density, $\mathbf{u} = (u_k)$ is the velocity vector, p is the pressure, γ is the ratio of specific heats, ρ_0 and p_0 are the static values of ρ and p , respectively; the standard notation for differentiation is used ($\rho_{,j} = \partial\rho/\partial x_j$, $\rho_{,t} = \partial\rho/\partial t$, etc.) and the summation over repeated indices holds.

Linearizing the equations around the equilibrium state,

$$\rho_0 = \text{const}, \quad p_0 = \text{const}, \quad \mathbf{u} = \mathbf{0}, \quad (2.2)$$

and introducing the small signal sound speed c_0 defined by

$$c_0^2 = \frac{dp}{d\rho}(\rho_0), \quad (2.3)$$

we arrive at the classical equations of linear acoustics in the form

$$\begin{aligned} \rho_0 u_{k,t} + p'_{,k} &= 0 \\ p'_{,t} + \rho_0 u_{k,k} &= 0 \\ p' &= c_0^2 \rho', \end{aligned} \quad (2.4)$$

where ρ' , p' , and \mathbf{u} are perturbed (or acoustical) variables defined by

$$\begin{aligned} \rho &= \rho_0 + \rho', \quad F\rho'F \ll \rho_0 \\ p &= p_0 + p', \quad Fp'F \ll \rho_0 c_0^2 \\ \mathbf{u} &= \mathbf{0} + \mathbf{u}, \quad \|\mathbf{u}\| \ll c_0. \end{aligned} \quad (2.5)$$

Finally, eliminating ρ' and introducing the non-dimensional variables

$$\begin{aligned} u_i^+ &= \frac{u_i}{c_0}, \quad p^+ = \frac{p'}{\rho_0 c_0^2}, \quad t^+ = \frac{t}{t_c}, \\ x_i^+ &= \frac{x_i}{c_0 t_c}, \quad i = 1, 2, 3, \end{aligned} \quad (2.6)$$

we arrive at the non-dimensional version of the equations in the form

$$\begin{aligned} \mathbf{u}_{,t} + \mathbf{grad} p &= \mathbf{0} \\ p_{,t} + \text{div} \mathbf{u} &= 0. \end{aligned} \quad (2.7)$$

Here t_c is some ‘‘characteristic time’’ (e.g., $t_c = l/c_0$, with l being a unit of length), and we omit superscripts ‘‘+’’ on non-dimensional quantities.

Equations (2.7) are to be solved in a domain $\Omega \subset \mathbb{R}^3$, $\nu = 2, 3$. Typically, two particular cases are of interest:

- *interior* problems when Ω is bounded
- *exterior* problems when Ω is a complement of a bounded set.

The initial boundary value problem is further specified by introducing boundary conditions. We consider the following kinds of boundary conditions:

1. Kinematic boundary condition (vibrating boundary)

$$u_n \stackrel{\text{def}}{=} \mathbf{u} \cdot \mathbf{n} = \hat{u}_n \quad \text{on } \Gamma_u, \quad (2.8)$$

where \mathbf{n} is the unit outward normal to the boundary and \hat{u}_n is the prescribed velocity of the vibrating boundary. The special case $\hat{u}_n = 0$ is referred to as the ‘‘solid wall’’ boundary condition.

2. Pressure boundary condition

$$p = \hat{p} \quad \text{on } \Gamma_p, \quad (2.9)$$

where \hat{p} is a prescribed pressure on the Γ_p -part of the boundary ($\partial\Omega = \overline{\Gamma_u} \cup \overline{\Gamma_p}$, $\Gamma_u \cap \Gamma_p = \emptyset$). In particular, the portion of the boundary at which $\hat{p} = 0$ is referred to as the ‘‘pressure release surface.’’

The initial boundary value problem is completed by specifying the initial conditions

$$\mathbf{u} = \mathbf{u}_0, \quad p = p_0 \quad \text{at } t = 0. \quad (2.10)$$

In the case of homogeneous boundary conditions, the problem can be cast into a Hilbert space formulation as follows (see [5] for a detailed discussion). We introduce

- the group variable,

$$\mathbf{U} = (\mathbf{u}^T, p)^T; \quad (2.11)$$

- the complex Hilbert space \mathbf{H} with inner product (\cdot, \cdot) and corresponding norm $\|\cdot\|$,

$$\begin{aligned} \mathbf{H} &= (L^2(\Omega))^\nu \times L^2(\Omega) \\ (\mathbf{U}, \mathbf{V}) &= (\mathbf{U}, \mathbf{V})_{(L^2)^\nu \times L^2} \end{aligned} \quad (2.12)$$

$$\|\mathbf{U}\| = (\mathbf{U}, \mathbf{U})^{1/2};$$

- operator $\mathbf{A} : \mathbf{H} \supset \mathbf{D}(\mathbf{A}) \rightarrow \mathbf{H}$,

$$\mathbf{A}\mathbf{U} \stackrel{\text{def}}{=} -i \begin{pmatrix} \mathbf{0} & \mathbf{grad} \\ \text{div} & 0 \end{pmatrix} \mathbf{U}, \quad (2.13)$$

where the domain of \mathbf{A} , $\mathbf{D}(\mathbf{A})$, consists of all vectors $\mathbf{U} = (\mathbf{u}^T, p)^T$ such that

$$\text{div} \mathbf{u} \in L^2(\Omega), \quad u_n = 0 \quad \text{on } \Gamma_u, \quad (2.14)$$

$$p \in H^1(\Omega), \quad p = 0 \quad \text{on } \Gamma_p, \quad (2.15)$$

and i is the imaginary unit. Note that the boundary condition on p is satisfied in the sense of the trace theorem, whereas the boundary condition on \mathbf{u} is interpreted in the sense of the

generalized Green's formula. For these reasons, the pressure boundary condition is classified as a Dirichlet boundary condition and the velocity boundary condition as a Neumann boundary condition for operator \mathbf{A} . Note also that $k\mathbf{U}k^2$ represents the total (mechanical) energy of the acoustical field:

$$2E = k\mathbf{U}k^2 = (\mathbf{U}, \mathbf{U}) = \int_{\Omega} [u_k u_k + p^2] dx. \quad (2.16)$$

Within the Hilbert space formalism, the initial boundary value problem can be reinterpreted as an abstract Cauchy problem for operator \mathbf{A} ,

$$\begin{aligned} \frac{d}{dt} \mathbf{U} + i\mathbf{A}\mathbf{U} &= \mathbf{0}, & t > 0, \\ \mathbf{U} &= \mathbf{U}_0, & t = 0, \end{aligned} \quad (2.17)$$

where $\mathbf{U}_0 \in \mathbf{H}$ is the initial condition vector $\mathbf{U}_0 = (\mathbf{u}_0^T, p_0)^T$.

An \mathbf{H} -valued function of time $\mathbf{U} = \mathbf{U}(t)$,

$$t \in [0, \infty) \rightarrow \mathbf{U}(t) \in \mathbf{H}, \quad (2.18)$$

is called a weak solution to the Cauchy problem if $\mathbf{U}(t)$ satisfies two conditions:

(i) a regularity assumption

$$\mathbf{U} \in C([0, \infty); \mathbf{H}) \quad (2.19)$$

(ii) a weak form of the equation given by

$$\int_0^{\infty} \int_{\Omega} \mathbf{U}^T (-\overline{\Phi}_{,i} + i\mathbf{A}\overline{\Phi}) dx dt - (\mathbf{U}_0, \Phi(0, \cdot)) = 0 \quad (2.20)$$

for every test function

$$\Phi \in C_0(\mathbb{R}; \mathbf{D}(\mathbf{A})) \cap C^1(\mathbb{R}; \mathbf{H}). \quad (2.21)$$

Note that this definition admits, in particular, solutions in the d'Alembert sense.

We record now some fundamental results concerning operator \mathbf{A} and the existence and uniqueness of weak solutions \mathbf{U} :

1. Operator \mathbf{A} is self-adjoint and, therefore, its spectrum lies on the real line and consists of a point spectrum (eigenvalues) and continuous spectrum (generalized eigenvalues) only.

2. For the interior problems (Ω bounded), the spectrum of \mathbf{A} consists of eigenvalues only. The eigenvalues of \mathbf{A} are symmetrically located on two parts of the real axis and "escape" to infinity (\mathbf{A} is unbounded):

$$\begin{aligned} \sigma(\mathbf{A}) &= \{0, \lambda_1, \lambda_{-1}, \lambda_2, \lambda_{-2}, \dots\} \\ 0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \rightarrow \infty, & \lambda_{-n} \equiv -\lambda_n. \end{aligned} \quad (2.22)$$

Except for the 0-eigenvalue, all eigenvalues are of finite multiplicity and the corresponding eigenspaces $\{\mathbf{U}_n\}$ are orthogonal.

3. For the exterior problems, the discrete spectrum reduces to the single eigenvalue 0 while the rest of the real axis constitutes the continuous spectrum.

4. In both cases, the eigenspace corresponding to zero eigenvalue, i.e., the null space of operator \mathbf{A} , is of infinite dimension and it contains all incompressible velocity fields. More precisely, it is of the form

$$N(\mathbf{A}) = \{(\mathbf{u}^T, p)^T \in \mathbf{H}; \operatorname{div} u = 0\}, \quad \dim N(\mathbf{A}) = \infty, \quad (2.23)$$

where p is an arbitrary constant if $\Gamma_p = \emptyset$ or $p = 0$ otherwise.

5. Operator \mathbf{A} admits a classical spectral decomposition

$$\mathbf{A}\mathbf{U} = \int_{-\infty}^{\infty} \lambda d\mathbf{E}_{\lambda} \mathbf{U}, \quad (2.24)$$

where \mathbf{E}_{λ} is a uniquely defined spectral family of \mathbf{A} (see, e.g., [11]).

6. A weak solution \mathbf{U} exists and is unique. Moreover, it is of the form

$$\mathbf{U}(t) = e^{-i\mathbf{A}t} \mathbf{U}_0 \equiv \int_{-\infty}^{\infty} e^{-i\lambda t} d\mathbf{E}_{\lambda} \mathbf{U}_0. \quad (2.25)$$

In particular, it follows from (2.25) that the energy is conserved

$$\|\mathbf{U}(t)\|^2 = \int_{-\infty}^{\infty} |e^{-i\lambda t}|^2 d(\mathbf{E}_{\lambda} \mathbf{U}_0, \mathbf{U}_0) = \|\mathbf{U}_0\|^2 \quad \forall t \geq 0. \quad (2.26)$$

6. If the initial condition function \mathbf{U}_0 satisfies an additional regularity assumption

$$\mathbf{U}_0 \in \mathbf{D}(\mathbf{A}) \quad (2.27)$$

then the solution $\mathbf{U} \in C^1([0, \infty); \mathbf{H}) \cap C([0, \infty); \mathbf{D}(\mathbf{A}))$. We say then that \mathbf{U} is a *strict solution* to the problem.

3. HIGH-ORDER TAYLOR-GALERKIN (TG) METHODS

In this section we introduce rather formally a class of high-order TG schemes for a system of conservation laws of the form

$$\begin{aligned} \mathbf{U}_{,i} + \mathbf{F}^i(\mathbf{U})_{,i} &= \mathbf{0}, & \mathbf{x} \in \Omega, t \in (0, t^*] \\ \mathbf{U}(\mathbf{x}, 0) &= \mathbf{U}_0(\mathbf{x}), & \mathbf{x} \in \Omega. \end{aligned} \quad (3.1)$$

Here Ω is a bounded domain in \mathbb{R}^2 (or \mathbb{R}^3), $\mathbf{U} = \mathbf{U}(\mathbf{x}, t)$ is a column vector of M unknowns, \mathbf{F}^i , $i = 1, \dots, M$, are vector-valued (linear or non-linear) functions of \mathbf{U} , commas denote the differentiation with respect to time t and spatial variables x_i , and the usual summation convention holds. It is assumed

that system (3.1) is accompanied by appropriate boundary conditions.

Any approximation to (3.1) must involve discretization both in space and time variables. We will adopt the assumption here that the final approximation is obtained by using finite differences in time and finite elements in space variables.

Then, two different approaches are possible. In the *classical method of lines*, an approximation in space variables converts the original initial-boundary-value problem into a system of ordinary differential equations (ODEs), which next is discretized in time using one of many time integration schemes for ODEs.

The TG schemes belong to a different class of methods, known as *the method of discretization in time*, which consists of the same two steps but done in the reverse order. By discretizing in time first, the initial-boundary-value problem (3.1) is converted into a sequence of boundary-value problems (3.2),

$$\begin{aligned} \mathbf{U}_\tau(\mathbf{x}, t_n + \Delta t) &= \mathbf{T}\mathbf{U}_\tau(\mathbf{x}, t_n) \\ \mathbf{U}_\tau(\mathbf{x}, 0) &= \mathbf{U}_0 \\ \Delta t &= t^*/N, \quad t_n = n\Delta t, \quad n = 0, 1, \dots, N, \end{aligned} \tag{3.2}$$

which, in turn, give a basis for a spatial approximation and, consequently, a fully discretized scheme (3.3),

$$\begin{aligned} \mathbf{U}_{\tau h}(\mathbf{x}, t_n + \Delta t) &= \mathbf{T}_h \mathbf{U}_{\tau h}(\mathbf{x}, t_n) \\ \mathbf{U}_{\tau h}(\mathbf{x}, 0) &= \mathbf{U}_{0h} \\ \Delta t &= t^*/N, \quad t_n = n\Delta t, \quad n = 0, 1, \dots, N. \end{aligned} \tag{3.3}$$

Here \mathbf{U}_τ and $\mathbf{U}_{\tau h}$ are the solutions of the semi-discrete problem (3.2) and the fully discrete problem (3.3), respectively, \mathbf{T} is an appropriate transient operator, \mathbf{T}_h denotes its finite-dimensional realization, and \mathbf{U}_{0h} is a suitable approximation of \mathbf{U}_0 .

3.1. Approximation in Time

3.1.1. Direct Taylor–Galerkin (TG) Schemes

Let N be a positive integer. We define a partition in time as $\Delta t = t^*/N$, $t_n = n\Delta t$, $n = 0, 1, \dots, N$, and consider a typical time-step $t_n \rightarrow t_n + \Delta t$. Given the solution $\mathbf{U}_\tau(t_n)$, we seek the next time-step solution $\mathbf{U}_\tau(t_n + \Delta t)$ through an s -stage scheme of the form

$$\begin{aligned} \mathbf{Z}_i - \eta \Delta t^2 \mathbf{Z}_{i,t} &= \mathbf{Z}_0 + \mu_{i0} \Delta t \mathbf{Z}_{0,t} + \nu_{i0} \Delta t^2 \mathbf{Z}_{0,tt} \\ &+ \Delta t \sum_{j=1}^{i-1} \mu_{ij} \mathbf{Z}_{j,t} + \Delta t^2 \sum_{j=1}^{i-1} \nu_{ij} \mathbf{Z}_{j,tt}, \quad i = 1, 2, \dots, s, \end{aligned} \tag{3.4}$$

where

$$\begin{aligned} \mathbf{Z}_j &= \mathbf{U}_\tau(t_n + c_j \Delta t), \quad j = 0, 1, \dots, s, \\ \mu_{ij}, \nu_{ij}, \mu_{i0}, \nu_{i0}, c_i &\in \mathbb{R}, \quad i = 1, 2, \dots, s; j = 1, 2, \dots, i-1, \\ \eta \in \mathbb{R}_+ &\text{ is a stability parameter} \\ s &= \text{number of stages.} \end{aligned} \tag{3.5}$$

Coefficients μ_{ij} , ν_{ij} , μ_{i0} , ν_{i0} , c_i , are to be chosen so as to obtain the highest possible order of accuracy, subject to stability or other constraints. A free parameter η is to be chosen from the stability considerations.

It is convenient to rewrite (3.4) in the compact form

$$\begin{aligned} \begin{pmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \\ \vdots \\ \mathbf{Z}_s \end{pmatrix} - \Delta t^2 \mathbf{N} \otimes \begin{pmatrix} \mathbf{Z}_{1,tt} \\ \mathbf{Z}_{2,tt} \\ \vdots \\ \mathbf{Z}_{s,tt} \end{pmatrix} - \Delta t \mathbf{M} \otimes \begin{pmatrix} \mathbf{Z}_{1,t} \\ \mathbf{Z}_{2,t} \\ \vdots \\ \mathbf{Z}_{s,t} \end{pmatrix} \\ = \mathbf{K} \otimes \begin{pmatrix} \mathbf{Z}_0 \\ \Delta t \mathbf{Z}_{0,t} \\ \Delta t^2 \mathbf{Z}_{0,tt} \end{pmatrix}, \end{aligned} \tag{3.6}$$

where

$$\begin{aligned} \mathbf{N} &= (\nu_{ij}), \quad \mathbf{M} = (\mu_{ij}) \in \mathbb{R}^{s \times s}, \\ \mathbf{K} &= (\kappa_{ij}) \in \mathbb{R}^{s \times 3}, \quad \mathbf{c} = (c_i) \in \mathbb{R}^{s \times 1} \end{aligned}$$

$$\mathbf{N} = \begin{pmatrix} \eta & 0 & 0 & 0 & 0 & 0 \\ \nu_{21} & \eta & 0 & 0 & 0 & 0 \\ \nu_{31} & \nu_{32} & \eta & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \nu_{s1} & \nu_{s2} & \nu_{s3} & \cdots & \nu_{s,s-1} & \eta \end{pmatrix}, \quad \mathbf{c} = (c_1 c_2 \cdots c_{s-1} 1),$$

$$\mathbf{M} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ \mu_{21} & 0 & 0 & 0 & 0 & 0 \\ \mu_{31} & \mu_{32} & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mu_{s1} & \mu_{s2} & \mu_{s3} & \cdots & \mu_{s,s-1} & 0 \end{pmatrix},$$

$$\mathbf{K} = \begin{pmatrix} 1 & \mu_{10} & \nu_{10} \\ 1 & \mu_{20} & \nu_{20} \\ 1 & \mu_{30} & \nu_{30} \\ \vdots & \vdots & \vdots \\ 1 & \mu_{s0} & \nu_{s0} \end{pmatrix}.$$

Matrices \mathbf{N} , \mathbf{M} , and \mathbf{K} , together with vector \mathbf{c} , completely characterize the difference scheme (3.6).

A distinct feature of (3.6) is that the coefficient matrices \mathbf{N} and \mathbf{M} are lower triangular matrices which makes the resulting scheme semi-implicit (i.e., to compute \mathbf{Z}_i it is necessary to know \mathbf{Z}_{i-1} , \mathbf{Z}_{i-2} , ..., \mathbf{Z}_1 , but it is not necessary to know \mathbf{Z}_{i+1}). Moreover, all the diagonal elements of \mathbf{N} are equal ($\nu_{ii} = \eta$, $i = 1, 2, \dots, s$) and so are those of \mathbf{M} ($\mu_{ii} = 0$, $i = 1, 2, \dots, s$). This makes the operator defining the left-hand side of each stage of (3.6) identical for linear (or linearized) problems, and, hence, significantly reduces the cost of applying the method.

A particular choice of zero diagonal elements of \mathbf{M} is made with an eye on a well-posedness of a typical one-stage problem and a possible splitting of the operator defining the left-hand side of each stage.

To make the i th stage solution \mathbf{Z}_i m th order accurate, it is necessary to satisfy the order conditions for \mathbf{Z}_i and to make the previous stage solutions $\mathbf{Z}_{i-1}, \mathbf{Z}_{i-2}, \dots, \mathbf{Z}_1$, to be at least of the order $m - 1$. (Otherwise, some coefficients have to be set to zero, e.g., if \mathbf{Z}_{i-1} is of the order $m - 2$, then, necessarily $\mu_{i,i-1} = 0$.) The order conditions for \mathbf{Z}_i are obtained by expanding it in Taylor series at \mathbf{Z}_0 ,

$$\begin{aligned} \mathbf{Z}_i = & \mathbf{Z}_0 + (c_i \Delta t) \mathbf{Z}_{0,t} + \frac{1}{2!} (c_i \Delta t)^2 \mathbf{Z}_{0,tt} \\ & + \dots + \frac{1}{m!} (c_i \Delta t)^m \frac{\partial^m}{\partial t^m} \mathbf{Z}_0 + O(\Delta t^{m+1}), \end{aligned} \quad (3.8)$$

introducing this expression into the left-hand side of the i th stage equation,

$$\begin{aligned} \mathbf{Z}_i - \sum_{j=1}^s (\mu_{ij} \Delta t \mathbf{Z}_{j,t} + \nu_{ij} \Delta t^2 \mathbf{Z}_{j,tt}) = & \mathbf{Z}_0 \\ & + \kappa_{i2} \Delta t \mathbf{Z}_{0,t} + \kappa_{i3} \Delta t^2 \mathbf{Z}_{0,tt} \end{aligned} \quad (3.9)$$

and equating coefficients of like powers of Δt to zero. This leads to the following system of non-linear algebraic equations:

$$\begin{aligned} c_i^k - k \sum_{j=1}^s c_j^{k-1} \mu_{ij} - k(k-1) \sum_{j=1}^s c_j^{k-2} \nu_{ij} \\ = \begin{cases} \mu_{i0}, & k = 1, \\ 2\nu_{i0}, & k = 2, \\ 0, & \text{otherwise;} \end{cases} \end{aligned} \quad (3.10)$$

$k = 1, 2, \dots, m$.

Equations (3.10) are referred to herein as the *order conditions*.

In addition, we introduce the *commutability constraints*,

$$\begin{aligned} (\nu_{i0} + \eta) \mu_{ij} - (\nu_{ij} + \varepsilon \eta) \mu_{i0} = 0, \\ i = 2, 3, \dots, s; j = 0, 1, \dots, s-1, \end{aligned} \quad (3.11)$$

where

$$\varepsilon = \begin{cases} 1, & j = 0, \\ 0, & j > 0. \end{cases} \quad (3.12)$$

The significance of (3.11) will be discussed in Section 4 in the context of the stability of TG schemes at finite dimensions.

For the sake of illustration, we derive coefficients for *one* particular scheme. The coefficients for several other TG schemes can be found in [8] and [9]. In the sequel, we adopt the following notation:

TG(s, m) = s -stage m th-order scheme which does *not* satisfy (3.11)

$\overline{\text{TG}}(s, m)$ = s -stage m th-order scheme which satisfies constraints (3.11).

2-Stage Scheme of Order 3 with Imposed Commutability Constraints ($\overline{\text{TG}}(2,3)$). We design \mathbf{Z}_1 and \mathbf{Z}_2 to be $O(\Delta t^2)$ and $O(\Delta t^3)$, respectively, and impose constraints (3.11), which leads to the following system of equations:

$$\begin{aligned} \mu_{10} = c_1, & & 2\nu_{10} = c_1^2 - 2\eta \\ \mu_{20} + \mu_{21} = 1, & & 2(\nu_{20} + \nu_{21}) + 2\mu_{21}c_1 = 1 - 2\eta \\ 6\nu_{21}c_1 + 3\mu_{21}c_1^2 = 1 - 6\eta, & & (\nu_{10} + \eta)\mu_{20} - (\nu_{20} + \eta)\mu_{10} = 0 \\ & & (\nu_{10} + \eta)\mu_{21} - \nu_{21}\mu_{10} = 0. \end{aligned} \quad (3.13)$$

The solution of (3.13) reads:

$$\begin{aligned} \mu_{10} = c_1 = \frac{1}{2} [1 \pm (-\frac{1}{8} + 8\eta)^{1/2}], & & \nu_{10} = \frac{1}{2} c_1^2 - \eta, \\ \mu_{20} = \frac{1}{2} (3 - c_1^{-1}), & & \nu_{20} = \frac{1}{4} (3c_1 - 1) - \eta, \\ \mu_{21} = \frac{1}{2} (c_1^{-1} - 1), & & \nu_{21} = \frac{1}{4} (1 - c_1). \end{aligned} \quad (3.14)$$

A nonexistence of solution for $\eta < \frac{1}{24}$ and its non-uniqueness for $\eta \geq \frac{1}{24}$, should be observed.

It is important to realize that the particular forms of the coefficient matrices \mathbf{M} and \mathbf{N} limit the attainable order of (direct) TG schemes. For example, it is not possible to attain order 6 in three stages. To remedy this situation, we introduce a new class of transformed Taylor–Galerkin (TG*) schemes.

3.1.2 Transformed Taylor–Galerkin Schemes (TG*)

The principal idea behind transformed Taylor–Galerkin schemes is to choose \mathbf{M} and \mathbf{N} to be full matrices which can be simultaneously transformed to the lower triangular form which is characteristic for direct TG schemes.

Formally TG* schemes are defined as follows. Given the solution $\mathbf{U}_r(t_n)$ at time $t_n = n\Delta t$, we seek the next time step solution $\mathbf{U}_r(t_n + \Delta t)$ in the form

$$\begin{aligned} \mathbf{Z}_i - \sum_{j=1}^s (\mu_{ij} \Delta t \mathbf{Z}_{j,t} + \nu_{ij} \Delta t^2 \mathbf{Z}_{j,tt}) = & \kappa_{i1} \mathbf{Z}_0 \\ & + \kappa_{i2} \Delta t \mathbf{Z}_{0,t} + \kappa_{i3} \Delta t^2 \mathbf{Z}_{0,tt}, \quad i = 1, 2, \dots, s, \end{aligned} \quad (3.15)$$

where

$$\begin{aligned} \mathbf{Z}_j = \mathbf{U}_r(t_n + c_j \Delta t), & & j = 0, 1, \dots, s, \\ \mu_{ij}, \nu_{ij}, \kappa_{ik}, c_i \in \mathbb{R}; & & i, j = 1, 2, \dots, s; \quad k = 1, 2, 3; \\ \kappa_{i1} = 1; & & i = 1, 2, \dots, s, \\ s = \text{number of stages.} \end{aligned} \quad (3.16)$$

Unknown coefficient matrices $\mathbf{M} = (\mu_{ij})$, $\mathbf{N} = (\nu_{ij})$, $\mathbf{K} = (\kappa_{ik})$, $\mathbf{c} = (c_i)$, and an unknown transformation matrix $\mathbf{R} = (r_{ij})$, are

determined from:

- order conditions ($m_1 + m_2 + \dots + m_s$ equations)

$$\begin{aligned} c_i^k - k \sum_{j=1}^s c_j^{k-1} \mu_{ij} - k(k-1) \sum_{j=1}^s c_j^{k-2} \nu_{ij} \\ = \begin{cases} \kappa_{i2}, & k = 1, \\ 2\kappa_{i3}, & k = 2, \\ 0, & \text{otherwise;} \end{cases} \\ k = 1, 2, \dots, m_i, \quad i = 1, 2, \dots, s, \\ m_i = \text{order of } \mathbf{Z}_i; \end{aligned}$$

- constraints on matrix \mathbf{M} : it is assumed that \mathbf{M} is similar to a strictly lower triangular matrix ($s(s+1)/2$ equations),

$$\sum_{k,l=1}^s r_{ik} \mu_{kl} q_{lj} = 0 \quad \text{for } i \leq j, \quad (3.18)$$

where $\mathbf{R} = (r_{ij})$ is an unknown transformation matrix (finding \mathbf{R} is a part of the problem), and $\mathbf{Q} = (q_{ik}) = \mathbf{R}^{-1}$;

- constraints on matrix \mathbf{N} : it is assumed that \mathbf{N} is similar to a lower triangular matrix with one, s -fold eigenvalue $\eta \in \mathbb{R}_+$ ($s(s+1)/2$ equations),

$$\sum_{k,l=1}^s r_{ik} \nu_{kl} q_{lj} = \begin{cases} 0, & \text{for } i < j, \\ \eta, & \text{for } i = j; \end{cases} \quad (3.19)$$

- constraints on transformation matrix \mathbf{R} : it is assumed that \mathbf{R} is orthonormal ($s(s+1)/2$ equations),

$$\begin{aligned} \mathbf{r}_k \cdot \mathbf{r}_l = \delta_{kl} \\ k = 1, 2, \dots, s; \quad l = 1, 2, \dots, k, \end{aligned} \quad (3.20)$$

where \mathbf{r}_k denotes k th column of \mathbf{R} and “ \cdot ” denotes the usual dot product.

The choice of \mathbf{R} to be orthonormal is justified, since any real matrix with real eigenvalues only is orthogonally similar to a lower triangular matrix (cf. the real Schur decomposition). Obviously, some other choices of \mathbf{R} are possible (e.g., \mathbf{M} and \mathbf{N} can be transformed to the Jordan canonical form).

Introducing the change of variables,

$$\mathbf{Z}_i = \sum_{j=1}^s q_{ij} \mathbf{Z}'_j \quad (3.21)$$

scheme (3.15) can be transformed to the form

$$\begin{aligned} \mathbf{Z}'_i - \sum_{j=1}^s (\mu'_{ij} \Delta t \mathbf{Z}'_{j,i} + \nu'_{ij} \Delta t^2 \mathbf{Z}'_{j,i}) = \kappa'_{i1} \mathbf{Z}_0 \\ + \kappa'_{i2} \Delta t \mathbf{Z}_{0,i} + \kappa'_{i3} \Delta t^2 \mathbf{Z}_{0,i}, \quad i = 1, 2, \dots, s, \end{aligned} \quad (3.22)$$

where

$$\mathbf{M}' = \mathbf{RMR}^{-1}, \quad \mathbf{N}' = \mathbf{RNR}^{-1}, \quad \mathbf{K}' = \mathbf{RK},$$

$$\begin{aligned} \mathbf{M}' &= \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ \mu'_{21} & 0 & 0 & 0 & 0 & 0 \\ \mu'_{31} & \mu'_{32} & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mu'_{s1} & \mu'_{s2} & \mu'_{s3} & \cdot & \mu'_{s,s-1} & 0 \end{pmatrix}, \\ \mathbf{N}' &= \begin{pmatrix} \eta & 0 & 0 & 0 & 0 & 0 \\ \nu'_{21} & \eta & 0 & 0 & 0 & 0 \\ \nu'_{31} & \nu'_{32} & \eta & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \nu'_{s1} & \nu'_{s2} & \nu'_{s3} & \cdot & \nu'_{s,s-1} & \eta \end{pmatrix}. \end{aligned} \quad (3.23)$$

Clearly, it is much easier to solve the “one-step problem” (3.15) by solving the transformed problem (3.22) (which is semi-implicit), rather than to solve (3.15) directly.

As for direct TG schemes, we also introduce the commutability constraints which now take the form

$$\begin{aligned} (\nu'_{i0} + \eta) \mu'_{ij} - (\nu'_{ij} + \varepsilon \eta) \mu'_{i0} = 0, \\ i = 2, 3, \dots, s; \quad j = 0, 1, \dots, s-1, \end{aligned} \quad (3.25)$$

where $\mu'_{i0} \equiv \kappa'_{i2}$ and $\nu'_{i0} \equiv \kappa'_{i3}$, $i = 1, 2, \dots, s$, and ε is defined by (3.12).

We now construct an error estimating method for TG* schemes. Let \mathbf{M} , \mathbf{N} , \mathbf{K} , \mathbf{c} , and \mathbf{R} be the coefficient matrices for TG* (s, m), and let $\bar{\mathbf{M}}$, $\bar{\mathbf{N}}$, $\bar{\mathbf{K}}$, $\bar{\mathbf{c}}$, and $\bar{\mathbf{R}}$ be the coefficient matrices for TG* ($s+1, m+1$), where TG* (s, m) and TG* ($s+1, m+1$) are referred to as the basic and the error estimating methods, respectively. The following choice of $\bar{\mathbf{M}}$, $\bar{\mathbf{N}}$, $\bar{\mathbf{K}}$, $\bar{\mathbf{c}}$, and $\bar{\mathbf{R}}$ guarantees that TG* (s, m) is embedded in TG* ($s+1, m+1$), i.e., that the m th-order method TG* (s, m) is a part of the $(m+1)$ th-order method TG* ($s+1, m+1$) (we symbolically denote embedding by “TG* (s, m) \subset TG* ($s+1, m+1$)”):

$$\begin{aligned} \bar{\mathbf{M}}_{(s+1) \times (s+1)} &= \left(\begin{array}{cccc|c} & & & & 0 \\ & \mathbf{M}_{s \times s} & & & 0 \\ & & & & \cdot \\ & & & & 0 \\ \hline & & & & 0 \\ \mu'_{s+1,1} & \mu'_{s+1,2} & \cdot & \mu'_{s+1,s} & 0 \end{array} \right) \\ \bar{\mathbf{N}}_{(s+1) \times (s+1)} &= \left(\begin{array}{cccc|c} & & & & 0 \\ & \mathbf{N}_{s \times s} & & & 0 \\ & & & & \cdot \\ & & & & 0 \\ \hline & & & & 0 \\ \nu'_{s+1,1} & \nu'_{s+1,2} & \cdot & \nu'_{s+1,s} & \eta \end{array} \right), \end{aligned}$$

$$\begin{aligned} \bar{\mathbf{K}}_{(s+1) \times 3} &= \begin{pmatrix} \mathbf{K}_{s \times 3} & & \\ 1 & \kappa_{s+1,2} & \kappa_{s+1,3} \end{pmatrix} \\ \bar{\mathbf{R}}_{(s+1) \times (s+1)} &= \left(\begin{array}{cccc|c} & & & & 0 \\ & \mathbf{R}_{s \times s} & & & 0 \\ & & & & \cdot \\ & & & & 0 \\ 0 & 0 & \cdot & 0 & 1 \end{array} \right), \\ \bar{\mathbf{c}}_{(s+1) \times 1} &= (\mathbf{c}_{s \times 1} \mid 1). \end{aligned} \quad (3.26)$$

As \mathbf{Z}_s and \mathbf{Z}_{s+1} are $O(\Delta t^{m+1})$ (resp. $O(\Delta t^{m+2})$) approximations to the solution $\mathbf{U}(t_n + \Delta t)$, the difference $\|\mathbf{Z}_{s+1} - \mathbf{Z}_s\|$ defines a *relative* error for \mathbf{Z}_s . In addition, the special form of coefficient matrices (3.26) allows us to compute \mathbf{Z}_{s+1} very economically: the extra computation involves the solution for one additional stage only.

The coefficients for several TG* schemes were derived numerically for some optimal values of stability parameter η and are listed in an Appendix. The following notation is adopted:

TG*(s, m) = s -stage m th-order transformed scheme which does *not* satisfy (3.25)

$\overline{\text{TG}}^*(s, m)$ = s -stage m th-order transformed scheme which satisfies (3.25).

3.2. Approximation in Space

First, using the original equations (3.1), we calculate the time derivatives in terms of spatial derivatives as

$$\mathbf{U}_{,i} = -\mathbf{F}^k(\mathbf{U})_{,k} \quad (3.27)$$

$$\mathbf{U}_{,ii} = (\mathbf{A}^k \mathbf{A}^l \mathbf{U}_{,l})_{,k}, \quad (3.28)$$

where $\mathbf{A}^k = \mathbf{F}^k_{,i}$ are the Jacobian matrices corresponding to fluxes \mathbf{F}^k . It is important to notice that $\mathbf{U}_{,ii}$ can be expressed in terms of spatial derivatives of \mathbf{U} of order at most 2, and, therefore, they can be effectively handled by C^0 continuous finite elements.

Next, replacing the time derivatives in (3.6) (or (3.15)) by formulas (3.27) and (3.28), we arrive at the system of equations

$$\begin{aligned} \mathbf{Z}_i - \sum_{j=1}^s (\mu_{ij} \Delta t [-\mathbf{F}^k(\mathbf{Z}_j)_{,k}] + \nu_{ij} \Delta t^2 [(\mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{j,l})_{,k}]) &= \mathbf{Z}_0 \\ + \kappa_{i2} \Delta t [-\mathbf{F}^k(\mathbf{Z}_0)_{,k}] + \kappa_{i3} \Delta t^2 [(\mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{0,l})_{,k}], & \quad (3.29) \\ i &= 1, 2, \dots, s, \end{aligned}$$

where the indices i and j are used to denote a particular stage, the indices k and l refer to the axes of a Cartesian coordinate system, a comma denotes partial differentiation, and the summation convention for k and l holds.

Finally, multiplying (3.29) by a vector-valued test function \mathbf{W} , integrating over Ω , and integrating by parts, we arrive at a variational formulation of the form

Given $\mathbf{Z}_0 \in \mathbf{X}$

Find $\mathbf{Z}_i \in \mathbf{X}$, $i = 1, 2, \dots, s$, such that

$$\begin{aligned} \int_{\Omega} \mathbf{W}^T \mathbf{Z}_i \, dx & \\ - \sum_{j=1}^s \mu_{ij} \Delta t \left(\int_{\Omega} \mathbf{W}^T_k \mathbf{F}^k(\mathbf{Z}_j) \, dx - \int_{\partial\Omega} \mathbf{W}^T \mathbf{F}^k(\mathbf{Z}_j) n_k \, ds \right) & \\ - \sum_{j=1}^s \nu_{ij} \Delta t^2 \left(- \int_{\Omega} \mathbf{W}^T_k \mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{j,l} \, dx + \int_{\partial\Omega} \mathbf{W}^T \mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{j,l} n_k \, ds \right) & \\ = \int_{\Omega} \mathbf{W}^T \mathbf{Z}_0 \, dx & \quad (3.30) \\ + \kappa_{i2} \Delta t \left(\int_{\Omega} \mathbf{W}^T_k \mathbf{F}^k(\mathbf{Z}_0) \, dx - \int_{\partial\Omega} \mathbf{W}^T \mathbf{F}^k(\mathbf{Z}_0) n_k \, ds \right) & \\ + \kappa_{i3} \Delta t^2 \left(- \int_{\Omega} \mathbf{W}^T_k \mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{0,l} \, dx + \int_{\partial\Omega} \mathbf{W}^T \mathbf{A}^k \mathbf{A}^l \mathbf{Z}_{0,l} n_k \, ds \right) & \end{aligned}$$

for all test functions $\mathbf{W} \in \mathbf{X}$, where $\mathbf{n} = (n_k)$ is the unit outward normal. The weak form (3.30) is the basis for finite element approximations.

4. A LINEAR STABILITY ANALYSIS

Let \mathbf{H} be a Hilbert space (with inner product (\cdot, \cdot) and the corresponding norm $\|\cdot\|$), and let $\mathbf{A}: \mathbf{H} \supset \mathbf{D}(\mathbf{A}) \rightarrow \mathbf{H}$ be a self-adjoint operator in \mathbf{H} . Being self-adjoint, \mathbf{A} admits the spectral decomposition of the form [11]

$$\mathbf{A}\mathbf{U} = \int_{-\infty}^{\infty} \lambda \, d\mathbf{E}_{\lambda} \mathbf{U} \quad (4.1)$$

$$\mathbf{D}(\mathbf{A}) = \left\{ \mathbf{U} \in \mathbf{H}: \int_{-\infty}^{\infty} \lambda^2 \, d\|\mathbf{E}_{\lambda} \mathbf{U}\|^2 < \infty \right\}, \quad (4.2)$$

where \mathbf{E}_{λ} is a uniquely defined spectral family. Let

$$\mathbf{X} = \mathbf{D}(\mathbf{A}) \quad (4.3)$$

and let $\|\cdot\|_{\mathbf{A}}$ denote the graph norm

$$\|\mathbf{U}\|_{\mathbf{A}}^2 = \|\mathbf{U}\|^2 + \|\mathbf{A}\mathbf{U}\|^2. \quad (4.4)$$

Equipped with the graph norm, \mathbf{X} is a Hilbert space.

In this section we consider stability properties of the TG schemes for solving abstract Cauchy problems of the form

$$\begin{aligned} \frac{d}{dt} \mathbf{U} + i\mathbf{A}\mathbf{U} &= \mathbf{0}, \quad 0 < t \leq t^*, \\ \mathbf{U} &= \mathbf{U}_0, \quad t = 0. \end{aligned} \quad (4.5)$$

For our focus, we consider \mathbf{A} to be the operator of linear acoustics introduced in Section 2. It should be noted, however, that (4.1)–(4.5) describe a more general situation; e.g., it may describe the propagation of stress waves in elastic solids.

4.1. Semidiscrete Schemes

We start by specifying the general formula (3.30) to the case of Cauchy problem (4.5). The rigid solid wall boundary conditions (2.8) take the form

$$u_n \stackrel{\text{def}}{=} \mathbf{u} \cdot \mathbf{n} = 0, \quad \dot{u}_n \stackrel{\text{def}}{=} -\partial p / \partial n = 0 \quad \text{on } \Gamma_u \quad (4.6)$$

and the pressure boundary conditions (2.9) read

$$p = 0, \quad \dot{p} \stackrel{\text{def}}{=} -\text{div } \mathbf{u} = 0 \quad \text{on } \Gamma_p. \quad (4.7)$$

Thus, both (4.6) and (4.7) result in vanishing of the boundary integrals in (3.30), which reduces to a sequence of s linear variational boundary-value problems of the form

$$\left[\begin{array}{l} \text{Given } \mathbf{Z}_0 \in \mathbf{X} \\ \text{Find } \mathbf{Z}_i \in \mathbf{X}, i = 1, 2, \dots, s, \text{ such that} \\ A(\mathbf{Z}_i, \mathbf{W}) - \Delta t \sum_{j=1}^s \mu_{ij} C(\mathbf{Z}_j, \mathbf{W}) + \Delta t^2 \sum_{j=1}^s \nu_{ij} B(\mathbf{Z}_j, \mathbf{W}) \quad (4.8) \\ = A(\mathbf{Z}_0, \mathbf{W}) + \kappa_{i2} \Delta t C(\mathbf{Z}_0, \mathbf{W}) - \kappa_{i3} \Delta t^2 B(\mathbf{Z}_0, \mathbf{W}) \\ \text{for all test functions } \mathbf{W} \in \mathbf{X} \end{array} \right.$$

where $\mathbf{X} = \mathbf{D}(\mathbf{A})$, and the bilinear (sesquilinear) forms A , B , and C are defined by

$$\begin{aligned} A, B, C: \mathbf{X} \times \mathbf{X} &\rightarrow \mathbb{C}, \\ A(\mathbf{U}, \mathbf{W}) &= (\mathbf{U}, \mathbf{W}) \\ B(\mathbf{U}, \mathbf{W}) &= (\mathbf{A}\mathbf{U}, \mathbf{A}\mathbf{W}) \\ C(\mathbf{U}, \mathbf{W}) &= -i(\mathbf{A}\mathbf{U}, \mathbf{W}). \end{aligned} \quad (4.9)$$

Applying transformation (3.21), (4.8) can be reduced to the semi-implicit form (for direct TG schemes, (4.8) is already in that form, i.e., $q_{ij} = \delta_{ij}$, $i, j = 1, 2, \dots, s$, where $\mathbf{Q} = (q_{ij})$ is the inverse of transformation matrix \mathbf{R})

$$\left[\begin{array}{l} \text{Given } \mathbf{Z}_0 \in \mathbf{X} \\ \text{Find } \mathbf{Z}'_i \in \mathbf{X}, i = 1, 2, \dots, s, \text{ such that} \\ A(\mathbf{Z}'_i, \mathbf{W}) + \eta \Delta t^2 B(\mathbf{Z}'_i, \mathbf{W}) \quad (4.10) \\ = \kappa'_{i1} A(\mathbf{Z}_0, \mathbf{W}) + \kappa'_{i2} \Delta t C(\mathbf{Z}_0, \mathbf{W}) - \kappa'_{i3} \Delta t^2 B(\mathbf{Z}_0, \mathbf{W}) \\ + \Delta t \sum_{j=1}^{i-1} \mu'_{ij} C(\mathbf{Z}'_j, \mathbf{W}) - \Delta t^2 \sum_{j=1}^{i-1} \nu'_{ij} B(\mathbf{Z}'_j, \mathbf{W}) \\ \text{for all test functions } \mathbf{W} \in \mathbf{X} \end{array} \right.$$

and the solution at time $t_n + \Delta t$ is given by

$$\mathbf{U}_r(t_n + \Delta t) = \sum_{j=1}^s q_{sj} \mathbf{Z}'_j. \quad (4.11)$$

It should be noted that bilinear forms A and B are symmetric (Hermitian)

$$\begin{aligned} \overline{A(\mathbf{W}, \mathbf{U})} &= A(\mathbf{U}, \mathbf{W}) \\ \overline{B(\mathbf{W}, \mathbf{U})} &= B(\mathbf{U}, \mathbf{W}) \end{aligned} \quad (4.12)$$

and that bilinear form C is skew-symmetric (skew-Hermitian)

$$\overline{C(\mathbf{W}, \mathbf{U})} = -C(\mathbf{U}, \mathbf{W}). \quad (4.13)$$

It is convenient to introduce an equivalent operator form of (4.10),

$$\mathbf{U}_r(t_n + \Delta t) = \mathbf{T}\mathbf{U}_r(t_n). \quad (4.14)$$

Here \mathbf{T} is a transient operator defined by

$$\mathbf{T} = \sum_{i=1}^s q_{si} \mathbf{T}'_i \quad (4.15)$$

with \mathbf{T}'_i given by the recurrence relation

$$\begin{aligned} \mathbf{T}'_1 &= \mathbf{T}'_{10} \\ \mathbf{T}'_i &= \mathbf{T}'_{i0} + \sum_{j=1}^{i-1} \mathbf{T}'_{ij} \mathbf{T}'_j, \quad i = 2, 3, \dots, s, \end{aligned} \quad (4.16)$$

operators $\mathbf{T}'_{ij}: \mathbf{X} \rightarrow \mathbf{X}$ defined by

$$\begin{aligned} B_1(\mathbf{T}'_{i0} \mathbf{U}, \mathbf{W}) &= L_{i0}(\mathbf{U}, \mathbf{W}) \\ \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}, i &= 1, 2, \dots, s, \\ B_1(\mathbf{T}'_{ij} \mathbf{U}, \mathbf{W}) &= L_{ij}(\mathbf{U}, \mathbf{W}) \\ \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}, j &= 1, 2, \dots, s-1, \\ & \quad i = 2, 3, \dots, s, i > j, \end{aligned} \quad (4.17)$$

and $B_1, L_{ij}: \mathbf{X} \times \mathbf{X} \rightarrow \mathbb{C}$, defined by

$$\begin{aligned} B_1(\mathbf{U}, \mathbf{W}) &= A(\mathbf{U}, \mathbf{W}) + \eta \Delta t^2 B(\mathbf{U}, \mathbf{W}) \\ L_{i0}(\mathbf{U}, \mathbf{W}) &= \kappa'_{i1} A(\mathbf{U}, \mathbf{W}) + \kappa'_{i2} \Delta t C(\mathbf{U}, \mathbf{W}) - \kappa'_{i3} \Delta t^2 B(\mathbf{U}, \mathbf{W}) \\ L_{ij}(\mathbf{U}, \mathbf{W}) &= \mu'_{ij} \Delta t C(\mathbf{U}, \mathbf{W}) - \nu'_{ij} \Delta t^2 B(\mathbf{U}, \mathbf{W}). \end{aligned} \quad (4.18)$$

It is easy to check that bilinear form B_1 is continuous and coercive,

$$B_1(\mathbf{U}, \mathbf{W}) \leq M \|\mathbf{U}\|_{\mathbf{A}} \|\mathbf{W}\|_{\mathbf{A}} \quad \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}, M > 0, \quad (4.19)$$

$$B_1(\mathbf{U}, \mathbf{U}) \geq \gamma \|\mathbf{U}\|_{\mathbf{A}}^2 \quad \forall \mathbf{U} \in \mathbf{X}, \gamma > 0, \quad (4.20)$$

and that bilinear forms L_{ij} are continuous,

$$L_{ij}(\mathbf{U}, \mathbf{W}) \leq M_{ij} \|\mathbf{U}\|_{\mathbf{A}} \|\mathbf{W}\|_{\mathbf{A}} \quad \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}, M_{ij} > 0, \quad (4.21)$$

and, therefore, by virtue of the Lax–Milgram theorem, operators \mathbf{T}'_{ij} are well defined.

Remark 4.1. The solution strategy for the “one-stage problem” deserves special attention. Since the formal operator equivalent to the bilinear form B_1 has zero off-diagonal terms,

$$\mathbf{I} + \eta \Delta t^2 \mathbf{A}^2 = \mathbf{I} - \eta \Delta t^2 \begin{pmatrix} \mathbf{grad} \operatorname{div} & \mathbf{0} \\ 0 & \operatorname{div} \mathbf{grad} \end{pmatrix}, \quad (4.22)$$

velocity equations and pressure equations can be solved independently and in parallel (the coupling between velocity and pressure is only through the right-hand side). Moreover, for the most common case of the potential velocity field, we have $\mathbf{grad} \operatorname{div} \mathbf{u} = \nabla^2 \mathbf{u}$ and, hence, decoupling of velocity equations in the interior of Ω . (Velocity components are coupled at the boundary, though). For this case, (4.22) reads

$$\mathbf{I} + \eta \Delta t^2 \mathbf{A}^2 = \mathbf{I} - \eta \Delta t^2 \nabla^2, \quad (4.23)$$

where ∇^2 is the Laplace operator. Note also that the matrix defining the left-hand side of the resulting algebraic problem is symmetric and positive definite, which is a very convenient setting for iterative solvers. The above two properties are the result of the absence of the first-order terms at the left-hand side in TG schemes.

To investigate stability properties of (4.10), we need to estimate the eigenvalues of transient operator \mathbf{T} , or, equivalently, the eigenvalues of the following eigenproblem:

$$\left\{ \begin{array}{l} \text{Find an eigenpair } (\Lambda, \mathbf{0} \neq \mathbf{U} \in \mathbf{X}) \\ \text{and } \mathbf{Z}_i \in \mathbf{X}, i = 1, 2, \dots, s, \mathbf{Z}_s = \Lambda \mathbf{U}, \text{ such that} \\ A(\mathbf{Z}_i, \mathbf{W}) - \Delta t \sum_{j=1}^s \mu_{ij} C(\mathbf{Z}_j, \mathbf{W}) + \Delta t^2 \sum_{j=1}^s \nu_{ij} B(\mathbf{Z}_j, \mathbf{W}) \\ = A(\mathbf{U}, \mathbf{W}) + \kappa_{i2} \Delta t C(\mathbf{U}, \mathbf{W}) - \kappa_{i3} \Delta t^2 B(\mathbf{U}, \mathbf{W}) \\ \text{for all test functions } \mathbf{W} \in \mathbf{X} \end{array} \right. \quad (4.24)$$

Toward this end, we recall that the underlying operator \mathbf{A} is self-adjoint (in the complex sense) and, in the case of an interior problem, has a pure point spectrum. In this case, spectral decomposition (2.24) reduces to

$$\mathbf{A} = \sum_{n=-\infty}^{\infty} \lambda_n \mathbf{P}_n \quad (4.25)$$

$$\mathbf{I} = \sum_{n=-\infty}^{\infty} \mathbf{P}_n \quad (4.26)$$

$$A(\mathbf{U}_m, \mathbf{U}_n) = \delta_{mn}, \quad (4.27)$$

where $(\lambda_n, \mathbf{U}_n)$ is the n th eigenpair of \mathbf{A} , $\mathbf{P}_n(\cdot) = A(\cdot, \mathbf{U}_n) \mathbf{U}_n$, for $n \neq 0$, \mathbf{P}_0 is the projector on $\mathbf{N}(\mathbf{A})$, $\lambda_{-n} \equiv -\lambda_n$, and $\mathbf{U}_{-n} \equiv \overline{\mathbf{U}_n}$.

By introducing (4.25)–(4.27) into (4.24) and taking

$$\mathbf{W} = \mathbf{U}_n \perp \mathbf{N}(\mathbf{A}) \quad (4.28)$$

we obtain the equation for $\Lambda = \Lambda_n$,

$$\begin{aligned} A(\mathbf{Z}_i, \mathbf{U}_n) & - i \Delta t \lambda_n \sum_{j=1}^s \mu_{ij} A(\mathbf{Z}_j, \mathbf{U}_n) \\ & + \Delta t^2 \lambda_n^2 \sum_{j=1}^s \nu_{ij} A(\mathbf{Z}_j, \mathbf{U}_n) = A(\mathbf{U}, \mathbf{U}_n) \\ & + i \kappa_{i2} \Delta t \lambda_n A(\mathbf{U}, \mathbf{U}_n) - \kappa_{i3} \Delta t^2 \lambda_n^2 A(\mathbf{U}, \mathbf{U}_n), \quad i = 1, 2, \dots, s, \end{aligned} \quad (4.29)$$

where $\mathbf{Z}_s = \Lambda_n \mathbf{U}$. From (4.29) we get Λ_n as a function of $\Delta t \lambda_n$, η , μ_{ij} , ν_{ij} :

$$\Lambda_n = \Lambda_n(\Delta t \lambda_n, \eta, \mu_{ij}, \nu_{ij}). \quad (4.30)$$

By choosing $\mathbf{U} = \mathbf{W} \in \mathbf{N}(\mathbf{A})$ it is easily seen that the corresponding eigenvalue Λ_0 is equal to unity.

To show the unconditional stability of TG methods, we need to show the existence of an η such that

$$|\Lambda_n| \leq 1 \quad \forall \Delta t \in \mathbb{R}_+. \quad (4.31)$$

To this end, we introduce the coefficients $\mu_{ij}(\eta)$ and $\nu_{ij}(\eta)$ derived in Section 3 into (4.30) and compute $|\Lambda_n|$ as a function of η and $\Delta t \lambda_n$. To illustrate, we consider $\overline{\text{TG}}(2, 3)$.

$\overline{\text{TG}}(2, 3)$. For the 2-stage third-order scheme satisfying commutability conditions and defined by (3.14), we obtain

$$\begin{aligned} |\Lambda_n|^2 & = 1 + \{[(\frac{1}{2} - 2\eta)^2 - 2(\frac{1}{6} - 2\eta) - 6\eta^2 + 2\theta](\Delta t \lambda_n)^4 \\ & + [(\frac{1}{6} - 2\eta)^2 - 4\eta^2 - 2(\frac{1}{2} - 2\eta)\theta](\Delta t \lambda_n)^6 + [\theta^2 \\ & - \eta^4](\Delta t \lambda_n)^8\} / \{1 + \eta(\Delta t \lambda_n)^2\}^4, \end{aligned} \quad (4.32)$$

where

$$\theta = -[\frac{1}{4}(3c_1 - 1) - \eta]\eta + \frac{1}{4}(1 - c_1)(\frac{1}{2}c_1^2 - \eta) \quad (4.33)$$

$$c_1 = \frac{1}{2}[1 \pm (-\frac{1}{3} + 8\eta)^{1/2}].$$

By analyzing the irrational functions of η appearing in the numerator of the fraction in (4.32), it can be shown that the stable scheme corresponds to the plus sign in the equation defining c_1 and that an optimal value of η guaranteeing unconditional stability is given by a zero of the quartic polynomial,

$$1024\eta^4 - 672\eta^3 + \frac{304}{3}\eta^2 - \frac{56}{9}\eta + \frac{4}{27} = 0, \quad (4.34)$$

the numerical value of which is given in Table A.VII.

We now analyze the stability properties of transformed TG schemes. To this end, we observe that both direct and transformed TG schemes, when applied to Cauchy problem (4.5), with $t = \Delta t$, generate the following eigenvalues Λ_n (which are approximations to $e^{-i\lambda_n \Delta t}$),

$$\Lambda_n = [1 + \eta(\Delta t \lambda_n)^2]^{-s} [\chi_0 + \chi_1(i\lambda_n \Delta t) + \chi_2(i\lambda_n \Delta t)^2 + \cdots + \chi_{2s}(i\lambda_n \Delta t)^{2s}], \quad (4.35)$$

where the coefficients χ_i , $i = 0, 1, 2, \dots, 2s$, can be deduced from (4.30). On the other hand, for an s -stage m th-order scheme ($m \leq 2s$), the coefficients χ_j , $j = 0, 1, 2, \dots, m$, can be obtained directly by expanding the function $[1 + \eta(\Delta t \lambda_n)^2]^s e^{-i\lambda_n \Delta t}$ in Taylor series at $i\lambda_n \Delta t = 0$ and equating the like powers of $i\lambda_n \Delta t$. Accordingly, identity (4.36),

$$\begin{aligned} [1 + \eta(\Delta t \lambda_n)^2]^s \left[1 + (i\lambda_n \Delta t) + \frac{1}{2!} (i\lambda_n \Delta t)^2 \right. \\ \left. + \cdots + \frac{1}{(m)!} (i\lambda_n \Delta t)^m \right] = [\chi_0 + \chi_1(i\lambda_n \Delta t) + \chi_2(i\lambda_n \Delta t)^2 \\ + \cdots + \chi_{2s}(i\lambda_n \Delta t)^{2s}] + O(|\lambda_n \Delta t|^{m+1}), \end{aligned} \quad (4.36)$$

leads to the relations for $\chi_i(\eta)$,

$$\begin{aligned} \chi_0(\eta) = \chi_1(\eta) = 1, \quad \chi_2(\eta) = \frac{1}{2!} - \binom{s}{1} \eta \\ \chi_3(\eta) = \frac{1}{3!} - \binom{s}{1} \eta, \quad \chi_4(\eta) = \frac{1}{4!} - \frac{1}{2!} \binom{s}{1} \eta + \binom{s}{2} \eta^2, \end{aligned} \quad (4.37)$$

etc., where $\binom{s}{j} = s!/j!(s-j)!$ is the Newton binomial symbol.

Since the expansion is unique, the coefficients χ_j , $j = 0, 1, 2, \dots, m$, are necessarily the same for both direct and transformed TG schemes (with the same number of stages s and the same order m). If, in addition, the remaining coefficients, i.e., χ_j , $j = m+1, m+2, \dots, 2s$, are also identical, then the stability results proven for direct schemes are also valid for transformed schemes. Thus, e.g., it follows immediately that $\overline{TG}^*(2, 3)$ is unconditionally stable for the η given by the zero of (4.34). We now consider some of the remaining transformed schemes.

$\overline{TG}^*(3, 4)$ -I and $\overline{TG}^*(3, 4)$ -II. We compute $|\Lambda_n|$ in terms of coefficients χ_j ,

$$\begin{aligned} |\Lambda_n|^2 = 1 + [(\chi_3^2 - 2\chi_2\chi_4 - \chi_5 + \chi_6) - 20\eta^3](\lambda_n \Delta t)^6 \\ + (\chi_4^2 + 2(\chi_2\chi_6 - \chi_3\chi_5) - 15\eta^4)(\lambda_n \Delta t)^8 \\ + (\chi_5^2 - 2\chi_4\chi_6 - 6\eta^5)(\lambda_n \Delta t)^{10} \\ + (\chi_6^2 - \eta^6)(\lambda_n \Delta t)^{12} / [1 + \eta(\Delta t \lambda_n)^2]^6, \end{aligned} \quad (4.38)$$

where χ_j , $j \leq 4$ are given by (4.37), and χ_5 and χ_6 are implicitly defined by (4.30) and are functions of coefficient matrices

$\mathbf{M} = (\mu_{ij})$, $\mathbf{N} = (\nu_{ij})$, $\mathbf{K} = (\kappa_{ik})$, $\mathbf{c} = (c_i)$, and $\mathbf{Q} = (q_{ij})$, defining the scheme. Since the analytical expressions for χ_5 and χ_6 as a function of η are unknown, we compute coefficient matrices for a given value of stability parameter η and check the sign of each term appearing in the numerator of (4.38). The results of numerical computation of the ranges of η guaranteeing unconditional stability are given in Table A.VII.

4.2. Fully Discrete Schemes

Thus far we investigated the stability properties of TG methods for \mathbf{X} being infinite-dimensional and we showed the existence of ranges of the stability parameter η guaranteeing satisfaction of (4.31). We now introduce an arbitrary conforming finite element approximation on domain Ω , and we denote the resulting finite dimensional finite element space by \mathbf{X}_h ($\mathbf{X}_h \subset \mathbf{X}$, $\dim \mathbf{X}_h < \infty$). The finite dimensional counterpart of eigenvalue problem (4.24) (i.e., when bilinear forms A , B , and C , (4.9), are restricted to $\mathbf{X}_h \times \mathbf{X}_h$), reads

$$\left[\begin{array}{l} \text{Find an eigenpair } (\Lambda_h, \mathbf{0} \neq \mathbf{U}_h \in \mathbf{X}_h \subset \mathbf{X}) \\ \text{and } \mathbf{Z}_{h,i} \in \mathbf{X}_h, i = 1, 2, \dots, s; \mathbf{Z}_{h,s} = \Lambda_h \mathbf{U}_h, \text{ such that} \\ A(\mathbf{Z}_{h,i}, \mathbf{W}) - \Delta t \sum_{j=1}^s \mu_{ij} C(\mathbf{Z}_{h,j}, \mathbf{W}) + \Delta t^2 \sum_{j=1}^s \nu_{ij} B(\mathbf{Z}_{h,j}, \mathbf{W}) \\ = A(\mathbf{U}_h, \mathbf{W}) + \kappa_{i2} \Delta t C(\mathbf{U}_h, \mathbf{W}) - \kappa_{i3} \Delta t^2 B(\mathbf{U}_h, \mathbf{W}) \quad (4.39) \\ \text{for all test functions } \mathbf{W} \in \mathbf{X}_h \end{array} \right.$$

As before, it is convenient to introduce an equivalent operator form of (4.39) (cf. (4.14)–(4.18)),

$$\mathbf{T}_h \mathbf{U}_h = \Lambda_h \mathbf{U}_h, \quad \mathbf{U}_h \neq \mathbf{0}, \quad (4.40)$$

where $\mathbf{T}_h: \mathbf{X}_h \rightarrow \mathbf{X}_h$, is defined as

$$\mathbf{T}_h = \sum_{i=1}^s q_{si} \mathbf{T}'_{h,i}, \quad (4.41)$$

$\mathbf{T}'_{h,i}$ is given by the recurrence relation

$$\begin{aligned} \mathbf{T}'_{h,1} = \mathbf{T}'_{h,10} \\ \mathbf{T}'_{h,i} = \mathbf{T}'_{h,i0} + \sum_{j=1}^{i-1} \mathbf{T}'_{h,ij} \mathbf{T}'_{h,j}, \quad i = 2, 3, \dots, s, \end{aligned} \quad (4.42)$$

and operators $\mathbf{T}'_{h,ij}: \mathbf{X}_h \rightarrow \mathbf{X}_h$ are defined by

$$\begin{aligned} B_1(\mathbf{T}'_{h,i0} \mathbf{U}, \mathbf{W}) = L_{i0}(\mathbf{U}, \mathbf{W}) \\ \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}_h, i = 1, 2, \dots, s, \\ B_1(\mathbf{T}'_{h,ij} \mathbf{U}, \mathbf{W}) = L_{ij}(\mathbf{U}, \mathbf{W}) \\ \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}_h, j = 1, 2, \dots, s-1; \\ i = 2, 3, \dots, s, i > j. \end{aligned} \quad (4.43)$$

There is a fundamental difference between (4.41) and (4.15), however, namely, that, in general, the operators $\mathbf{T}'_{h,ij}$ do *not* commute with each other,

$$\mathbf{T}'_{h,ij}\mathbf{T}'_{h,kl} \neq \mathbf{T}'_{h,kl}\mathbf{T}'_{h,ij}, \quad i \neq k \text{ or } j \neq l, \quad (4.44)$$

and, therefore, the eigenvectors of $\mathbf{T}'_{h,ij}$ differ from that of $\mathbf{T}'_{h,kl}$, and, consequently, from that of \mathbf{T}_h . Moreover, due to the presence of skew-symmetric form C , the eigenvectors of $\mathbf{T}'_{h,ij}$ do not form an orthogonal system. These facts together imply that the technique used to prove the unconditional stability of TG methods at the infinite dimensional level, where the operators \mathbf{T}'_{ij} (and, consequently \mathbf{T}) have the same (orthonormal) system of eigenvectors, namely, the eigenvectors of the underlying operator \mathbf{A} , cannot be used at finite dimensional levels.

We shall now restrict ourselves to TG schemes satisfying commutability conditions, (3.11) or (3.25). We first put the operator $\mathbf{T}'_{h,ij}$ into the following form:

$$\begin{aligned} B_1(\mathbf{T}'_{h,ij}\mathbf{U}, \mathbf{W}) &= a_{ij}A(\mathbf{U}, \mathbf{W}) + c_{ij}\Delta t C(\mathbf{U}, \mathbf{W}) \\ &\quad - b_{ij}B_1(\mathbf{U}, \mathbf{W}) \quad \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}_h, \\ j &= 0, 1, 2, \dots, s-1; i = 1, 2, \dots, s; i > j; \quad (4.45) \\ a_{ij} &= \begin{cases} \kappa'_{i1} + \eta^{-1}\kappa'_{i3}, & j = 0, \\ \eta^{-1}\nu'_{ij}, & j > 0; \end{cases} \quad b_{ij} = \begin{cases} \eta^{-1}\kappa'_{i3}, & j = 0, \\ \eta^{-1}\nu'_{ij}, & j > 0; \end{cases} \\ c_{ij} &= \begin{cases} \kappa'_{i2}, & j = 0, \\ \mu'_{ij}, & j > 0. \end{cases} \end{aligned}$$

Next, we consider another operator, $\mathbf{T}'_{h,kl}$, say, with $i \neq k$ or $j \neq l$, and express it in the form

$$\begin{aligned} B_1(\mathbf{T}'_{h,kl}\mathbf{U}, \mathbf{W}) &= a_{ij}^{-1}a_{kl}[a_{ij}A(\mathbf{U}, \mathbf{W}) + c_{ij}\Delta t C(\mathbf{U}, \mathbf{W})] \\ &\quad + a_{ij}^{-1}(a_{ij}c_{kl} - a_{kl}c_{ij})\Delta t C(\mathbf{U}, \mathbf{W}) \quad (4.46) \\ &\quad - b_{ij}B_1(\mathbf{U}, \mathbf{W}) \quad \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}_h. \end{aligned}$$

It is immediate from (4.46) that $\mathbf{T}'_{h,ij}$ and $\mathbf{T}'_{h,kl}$ will have the same system of eigenvectors (and, hence, will commute), if

$$a_{ij}c_{kl} - a_{kl}c_{ij} = 0. \quad (4.47)$$

One easily recognizes in (4.47) the commutability constraints introduced in Section 3. Thus, commutability constraints (3.25) (or (3.11)) guarantee that *all* operators $\mathbf{T}'_{h,ij}$ have the common system of eigenvectors.

Furthermore, we observe that (4.47) implies the relation between the eigenvalues of the two operators,

$$\lambda'_{h,kl} = a_{ij}^{-1}a_{kl}[\lambda'_{h,ij} + b_{ij}] - b_{kl}, \quad (4.48)$$

where $\lambda'_{h,kl}$ is an eigenvalue of $\mathbf{T}'_{h,kl}$ and $\lambda'_{h,ij}$ is that of $\mathbf{T}'_{h,ij}$. This, together with (4.41) and (4.42), enables us to express the eigenvalues of transient operator \mathbf{T}_h in terms of the eigenvalues of any of the operators $\mathbf{T}'_{h,ij}$ (for focus, we choose $\mathbf{T}'_{h,10}$):

$$\Lambda_h = \Lambda_h(\lambda'_{h,10}). \quad (4.49)$$

We now consider eigenproblem (4.50), defining eigenvalues $\lambda'_{h,10}$:

$$\left[\begin{array}{l} \text{Find an eigenpair } (\lambda'_{h,10}, \mathbf{0} \neq \mathbf{U}_h \in \mathbf{X}_h) \text{ such that} \\ \lambda'_{h,10}[A(\mathbf{U}_h, \mathbf{W}) + \eta\Delta t^2 B(\mathbf{U}_h, \mathbf{W})] \\ \quad = \kappa'_{i1}A(\mathbf{U}_h, \mathbf{W}) + \kappa'_{i2}\Delta t C(\mathbf{U}_h, \mathbf{W}) - \kappa'_{i3}\Delta t^2 B(\mathbf{U}_h, \mathbf{W}) \\ \text{for all test function } \mathbf{W} \in \mathbf{X}_h \end{array} \right. \quad (4.50)$$

By taking $\mathbf{W} = \mathbf{U}_h$, $\lambda'_{h,10}$ can be expressed as

$$\lambda'_{h,10} = \frac{\kappa'_{i1} - \kappa'_{i3}\Delta t^2(b/a) - i\kappa'_{i2}\Delta t(c/a)}{1 + \eta\Delta t^2(b/a)}, \quad (4.51)$$

where $a, b, c \in \mathbb{R}$ are defined as

$$\begin{aligned} a &= A(\mathbf{U}_h, \mathbf{U}_h) > 0 \\ b &= B(\mathbf{U}_h, \mathbf{U}_h) \geq 0 \\ ic &= C(\mathbf{U}_h, \mathbf{U}_h) \end{aligned} \quad (4.52)$$

and we used the fact that A and B are Hermitian and C is skew-Hermitian (cf. (4.12) and (4.13)). We also note that A is positive-definite and that B is semi positive-definite on \mathbf{X}_h . By introducing (4.52) into (4.49), we obtain the eigenvalue of transient operator \mathbf{T}_h as a function of *two* variables $(b/a)\Delta t^2$ and $(c/a)\Delta t$:

$$\Lambda_h = \Lambda_h\left(\frac{c}{a}\Delta t, \frac{b}{a}\Delta t^2\right). \quad (4.53)$$

In addition, we record the relation between b/a and c/a ,

$$\left|\frac{c}{a}\right|^2 \leq \frac{b}{a} \quad (4.54)$$

which follows from the Cauchy-Schwarz inequality

$$|c|^2 = |(A\mathbf{U}, \mathbf{U})|^2 \leq (\mathbf{U}, \mathbf{U})(A\mathbf{U}, A\mathbf{U}) = ab. \quad (4.55)$$

Thus, the stability analysis is reduced to the study of (4.53), subject to constraint (4.54). To illustrate the above ideas, we consider the stability of $\overline{\text{TG}}(2, 3)$. (The proofs of the unconditional stability of $\overline{\text{TG}}^*(2, 3)$, $\overline{\text{TG}}^*(3, 4)$ -I, and $\overline{\text{TG}}^*(3, 4)$ -II can be found in [9].)

$\overline{\text{TG}}(2, 3)$. We introduce coefficients (3.14) into (4.49) and obtain Λ_h as a quadratic polynomial of a complex variable $\lambda_{h,10}$,

$$\Lambda_h = b_0 + b_1\lambda_{h,10} + b_2\lambda_{h,10}^2, \quad (4.56)$$

where $b_i \in \mathbb{R}$, $i = 0, 1, 2$, are given by

$$b_0 = \frac{1}{2}c_1^{-2}(1 - 3c_1), \quad b_1 = c_1^{-2}(2c_1 - 1), \quad b_2 = \frac{1}{2}c_1^{-2}(1 - c_1), \quad (4.57)$$

and $c_1 = \frac{1}{2}[1 + (-\frac{1}{3} + 8\eta)^{1/2}]$. Since the coefficients of the quadratic polynomial are real and its discriminant, D , is non-negative for all values of η for which the scheme exists,

$$D = c_1^{-2}(2c_1 - 1) = c_1^{-2}(-\frac{1}{3} + 8\eta)^{1/2} \geq 0, \quad (4.58)$$

Λ_h can be factored as

$$\Lambda_h = b_2 \prod_{j=1}^2 (\lambda_{h,10} - \eta_j), \quad (4.59)$$

where η_j , $j = 1, 2$, are necessarily *real* zeros of the quadratic polynomial. Then,

$$\begin{aligned} |\Lambda_h| &= \left| \Lambda_h \left(\frac{c}{a} \Delta t, \frac{b}{a} \Delta t^2 \right) \right| \\ &= |b_2| \prod_{j=1}^2 \left| \frac{1 - \nu_{10} \Delta t^2 (b/a) - i\mu_{10} \Delta t (c/a)}{1 + \eta \Delta t^2 (b/a)} - \eta_j \right| \\ &\leq |b_2| \prod_{j=1}^2 \left| \frac{1 - \nu_{10} \Delta t^2 (b/a) - i\mu_{10} \Delta t \sqrt{b/a}}{1 + \eta \Delta t^2 (b/a)} - \eta_j \right| \\ &= |\Lambda_h(\sqrt{b/a} \Delta t, (b/a) \Delta t^2)|, \end{aligned} \quad (4.60)$$

where we used inequality (4.54). Finally, we note that $|\Lambda_h(\sqrt{b/a} \Delta t, (b/a) \Delta t^2)|^2$ has necessarily the same functional form as (4.32) with $\Delta t \lambda_n$ replaced with $\sqrt{b/a} \Delta t$, and, hence, $|\Lambda_h| \leq 1$ for η defined implicitly by (4.34).

Remark 4.2. The proofs of unconditional stability of $\overline{\text{TG}}$ and $\overline{\text{TG}}^*$ schemes at finite dimensions are also valid at the infinite dimensional level.

Remark 4.3. The stability analysis presented in Section 4.2 is valid for periodic boundary conditions and for a general *symmetrizable* linear conservation laws of the form (3.1).

Finally, we note that the formal proof of stability of TG schemes which do not satisfy commutability conditions, is not complete. On the other hand, an extensive benchmarking on a variety of problems and meshes, including problems with singularities and arbitrary *hp*-adapted meshes, suggests that unconditional stability at the infinite dimensional level may imply unconditional stability at finite dimensions.

5. A PRIORI ERROR ESTIMATION FOR TG METHODS

5.1. Temporal Approximation Error

Let $\mathbf{U}(t)$ be the solution of the Cauchy problem (4.5) and let $\mathbf{U}_r(t)$ be its semidiscrete approximation

$$\begin{aligned} \mathbf{U}_r(t_n + \Delta t) &= \mathbf{T} \mathbf{U}_r(t_n) \\ \mathbf{U}_r(0) &= \mathbf{U}_0 \end{aligned} \quad (5.1)$$

$$\Delta t = t^*/N, \quad t_n = n\Delta t, \quad n = 0, 1, \dots, N,$$

where \mathbf{T} is a transient operator corresponding to particular TG scheme and defined by (4.15)–(4.18).

From the definition of \mathbf{T} and spectral decomposition (4.1), it follows that \mathbf{T} can be represented as a rational function of the underlying operator \mathbf{A} ,

$$\mathbf{T} \mathbf{U} = r_s(i\Delta t \mathbf{A}) \mathbf{U}, \quad (5.2)$$

where $r_s: \mathbb{Z} \rightarrow \mathbb{Z}$ is given by

$$r_s(z) = P_\mu(z)/(1 - \eta z^2)^s, \quad z \neq \pm \eta, \mu \leq 2s, \quad (5.3)$$

P_μ is a μ th-order polynomial (with real coefficients), s denotes the number of stages of the TG scheme, and $\eta \in \mathbb{R}_+$ is the stability parameter. $r_s(z)$ is analytic for $|\text{Re}(z)| < \eta$ and, for an m th-order scheme,

$$|r_s(iy) - e^{-iy}| = O(|y|^{m+1}) \quad \text{as } y < y_0, y \in \mathbb{R}. \quad (5.4)$$

Moreover, for each scheme there exists an η_0 such that

$$|r_s(iy)| \leq 1 \quad \forall \eta \geq \eta_0, \forall y \in \mathbb{R}. \quad (5.5)$$

We prove the above assertions for the case of a 2-stage direct TG scheme (the proof for other schemes is similar). In this case $\mathbf{T}: \mathbf{X} \rightarrow \mathbf{X}$, $\mathbf{X} = \mathbf{D}(\mathbf{A})$, is given by

$$\mathbf{T} \mathbf{U} = (\mathbf{T}_{20} + \mathbf{T}_{21} \mathbf{T}_{10}) \mathbf{U} \quad (5.6)$$

and \mathbf{T}_{20} , \mathbf{T}_{21} , $\mathbf{T}_{10}: \mathbf{X} \rightarrow \mathbf{X}$, are defined by

$$\mathbf{B}_1(\mathbf{T}_{ij} \mathbf{U}, \mathbf{W}) = L_{ij}(\mathbf{U}, \mathbf{W}) \quad \forall \mathbf{W}, \mathbf{U} \in \mathbf{X}, i = 1, 2; j = 0, 1, \quad (5.7)$$

where B_i and L_{ij} are defined by (4.16). Assume that $\mathbf{U} \in \mathbf{D}(\mathbf{A}^4)$. Then, integration by parts yields

$$\begin{aligned} \mathbf{T}_{i0}\mathbf{U} &= (\mathbf{I} + \eta\Delta t^2 \mathbf{A}^2)^{-1}(\mathbf{I} - i\mu_{i0}\Delta t \mathbf{A} - \nu_{i0}\Delta t^2 \mathbf{A}^2)\mathbf{U}, \\ i &= 1, 2, \\ \mathbf{T}_{21}\mathbf{U} &= (\mathbf{I} + \eta\Delta t^2 \mathbf{A}^2)^{-1}(-i\mu_{21}\Delta t \mathbf{A} - \nu_{21}\Delta t^2 \mathbf{A}^2)\mathbf{U} \end{aligned} \quad (5.8)$$

and consequently,

$$\begin{aligned} \mathbf{T}\mathbf{U} &= [(\mathbf{I} + \eta\Delta t^2 \mathbf{A}^2)^{-1}(\mathbf{I} - i\mu_{20}\Delta t \mathbf{A} - \nu_{20}\Delta t^2 \mathbf{A}^2) \\ &\quad + (\mathbf{I} + \eta\Delta t^2 \mathbf{A}^2)^{-1}(-i\mu_{21}\Delta t \mathbf{A} - \nu_{21}\Delta t^2 \mathbf{A}^2) \\ &\quad \times (\mathbf{I} + \eta\Delta t^2 \mathbf{A}^2)^{-1}(\mathbf{I} - i\mu_{10}\Delta t \mathbf{A} - \nu_{10}\Delta t^2 \mathbf{A}^2)]\mathbf{U} \\ &= \int_{-\infty}^{\infty} r_s(i\Delta t\lambda) d\mathbf{E}_\lambda \mathbf{U} = r_s(i\Delta t\mathbf{A})\mathbf{U}, \end{aligned} \quad (5.9)$$

where $r_s(\cdot)$ is given by

$$\begin{aligned} r_s(z) &= [1 - (\mu_{20} + \mu_{21})z - (\nu_{20} - \eta + \mu_{21}\mu_{10} + \nu_{21})z^2 \\ &\quad + (-\mu_{20}\eta + \mu_{21}\nu_{10} + \nu_{21}\mu_{10})z^3 \\ &\quad + (-\nu_{20}\eta + \nu_{21}\nu_{10})z^4](1 - \eta z^2)^{-2}. \end{aligned} \quad (5.10)$$

Hence, (5.2) follows.

To show (5.3), we note that if \mathbf{A} has a point spectrum with eigenvalues λ_n , then $r_s(\Delta t\lambda_n)$ are nothing but eigenvalues of transient operator \mathbf{T} , which were analysed in Section 4. Finally, order condition (5.4) follows from order conditions (3.10). (Alternatively, (5.4) can be checked directly.)

We now define \mathbf{E}_τ , the temporal approximation error,

$$\mathbf{E}_\tau \stackrel{\text{def}}{=} \mathbf{U}(t^*) - \mathbf{U}_\tau(t^*). \quad (5.11)$$

To estimate $\|\mathbf{E}_\tau\|$ we need the following result. Let $\mathbf{U} \in \mathbf{D}(\mathbf{A}^{m+1})$, $m + 1 \geq \mu$, and (5.4) holds. Then there exists a constant C_1 such that the one-step error satisfies

$$\|e^{-i\Delta t\mathbf{A}}\mathbf{U} - r_s(i\Delta t\mathbf{A})\mathbf{U}\| \leq C_1\Delta t^{m+1} \|\mathbf{A}^{m+1}\mathbf{U}\|. \quad (5.12)$$

Indeed,

$$\begin{aligned} &\|e^{-i\Delta t\mathbf{A}}\mathbf{U} - r_s(i\Delta t\mathbf{A})\mathbf{U}\|^2 \\ &= \int_{-\infty}^{\infty} |e^{-i\lambda\Delta t} - r_s(i\lambda\Delta t)|^2 d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) \\ &\leq (C_1\Delta t^{m+1})^2 \int_{-\infty}^{\infty} |\lambda^{m+1}|^2 d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) \\ &= (C_1\Delta t^{m+1})^2 \|\mathbf{A}^{m+1}\mathbf{U}\|^2. \end{aligned} \quad (5.13)$$

Now, assume that (5.13) and (5.5) hold. Then, the temporal approximation error is bounded as follows:

$$\begin{aligned} \|\mathbf{E}_\tau\| &= \|e^{-i\Delta t\mathbf{A}}\mathbf{U}_0 - r_s^N(i\Delta t\mathbf{A})\mathbf{U}_0\| \\ &\leq \sum_{j=0}^{N-1} \|e^{-i(N-j-1)\Delta t\mathbf{A}}\| \|r_s^j(i\Delta t\mathbf{A})\| \|e^{-i\Delta t\mathbf{A}}\mathbf{U}_0 - r_s(i\Delta t\mathbf{A})\mathbf{U}_0\| \\ &\leq N \|e^{-i\Delta t\mathbf{A}}\mathbf{U}_0 - r_s(i\Delta t\mathbf{A})\mathbf{U}_0\| \\ &\leq NC_1\Delta t^{m+1} \|\mathbf{A}^{m+1}\mathbf{U}_0\| \\ &= C_1 t^* \Delta t^m \|\mathbf{A}^{m+1}\mathbf{U}_0\|, \quad \mathbf{U}_0 \in \mathbf{D}(\mathbf{A}^{m+1}), m + 1 \geq \mu. \end{aligned} \quad (5.14)$$

We also estimate temporal approximation error in the energy norm $\|\cdot\|_E$ defined by the bilinear form appearing at the left-hand side of TG schemes,

$$\|\cdot\|_E \stackrel{\text{def}}{=} B_1(\cdot, \cdot)^{1/2}, \quad (5.15)$$

where B_1 is defined in (4.18). We first observe that, if $\mathbf{U} \in \mathbf{D}(\mathbf{A}^{m+2})$, $m + 2 \geq \mu$, and (5.4) holds, then the one-step error is estimated by

$$\begin{aligned} &\|e^{-i\Delta t\mathbf{A}}\mathbf{U} - r_s(i\Delta t\mathbf{A})\mathbf{U}\|_E^2 \\ &= \int_{-\infty}^{\infty} |e^{-i\lambda\Delta t} - r_s(i\lambda\Delta t)|^2 (1 + \eta\Delta t^2\lambda^2) d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) \\ &\leq (C_1\Delta t^{m+1})^2 \int_{-\infty}^{\infty} |\lambda^{m+1}|^2 (1 + \eta\Delta t^2\lambda^2) d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) \\ &= (C_1\Delta t^{m+1})^2 \|\mathbf{A}^{m+1}\mathbf{U}\|_E^2. \end{aligned} \quad (5.16)$$

Next, we note that $e^{-i\Delta t\mathbf{A}}$ is an isometry in the energy norm:

$$\begin{aligned} &\|e^{-i\Delta t\mathbf{A}}\mathbf{U}\|_E^2 \\ &= \int_{-\infty}^{\infty} |e^{-i\Delta t\lambda}|^2 (1 + \eta\Delta t^2\lambda^2) d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) = \|\mathbf{U}\|_E^2. \end{aligned} \quad (5.17)$$

Finally, if (5.5) holds, then $r_s(i\Delta t\mathbf{A})$ is a contraction in the energy norm

$$\begin{aligned} &\|r_s(i\Delta t\mathbf{A})\mathbf{U}\|_E^2 \\ &= \int_{-\infty}^{\infty} |r_s(i\lambda\Delta t)|^2 (1 + \eta\Delta t^2\lambda^2) d(\mathbf{E}_\lambda \mathbf{U}, \mathbf{U}) \leq \|\mathbf{U}\|_E^2. \end{aligned} \quad (5.18)$$

Consequently, assuming that $\mathbf{U}_0 \in \mathbf{D}(\mathbf{A}^{m+2})$, estimate (5.14) is also valid in the energy norm:

$$\|\mathbf{E}_\tau\|_E \leq C_1 t^* \Delta t^m \|\mathbf{A}^{m+1}\mathbf{U}_0\|_E, \quad \mathbf{U}_0 \in \mathbf{D}(\mathbf{A}^{m+2}), m + 2 \geq \mu. \quad (5.19)$$

5.2. Spatial Approximation Error

We introduce a family of finite-dimensional subspaces \mathbf{X}_h of \mathbf{X} , indexed by a parameter $h \in (0, 1)$, and a corresponding family of finite dimensional operators \mathbf{T}_h , approximating tran-

sient operator \mathbf{T} . In particular, we consider

$$\mathbf{X}_h = [X_{hp}]^{\nu+1} \cap \mathbf{X}, \quad \nu = 2 \text{ or } 3, \quad (5.20)$$

where X_{hp} is a standard finite element space of piecewise polynomials of degree p and we assume that the following approximation property holds:

$$\inf_{\chi \in \mathbf{X}_h} \{ \|\mathbf{U} - \chi\| + h \|\mathbf{A}(\mathbf{U} - \chi)\| \} \leq Ch^s \|\mathbf{U}\|_{\mathbf{H}^s(\Omega)} \quad (5.21)$$

$$\forall \mathbf{U} \in \mathbf{X} \cap \mathbf{H}^s(\Omega), \quad 1 \leq s \leq p + 1.$$

Here $\mathbf{H}^s(\Omega) = [H^s(\Omega)]^{\nu+1}$, $H^s(\Omega)$ denotes the usual Sobolev space of order s , and C is a constant independent of h and \mathbf{U} . We assume that $\Omega \subset \mathbf{R}^\nu$, $\nu = 2, 3$, and that $\partial\Omega$ is sufficiently smooth.

We now consider the following approximation of (5.6),

$$\begin{aligned} \mathbf{U}_{th}(t_n + \Delta t) &= \mathbf{T}_h \mathbf{U}_{th}(t_n) \\ \mathbf{U}_{th}(0) &= \mathbf{U}_{0h} \\ \Delta t &= t^*/N, \quad t_n = n\Delta t, \quad n = 0, 1, \dots, N, \end{aligned} \quad (5.22)$$

where $\mathbf{T}_h: \mathbf{X}_h \rightarrow \mathbf{X}_h$ is defined by (4.41),

$$\mathbf{U}_{0h} = \mathbf{P}_h \mathbf{U}_0, \quad (5.23)$$

and $\mathbf{P}_h: \mathbf{X} \rightarrow \mathbf{X}_h$ is the projection operator defined by

$$B_1(\mathbf{P}_h \mathbf{U}, \mathbf{W}) = B_1(\mathbf{U}, \mathbf{W}) \quad \forall \mathbf{U} \in \mathbf{X}, \forall \mathbf{W} \in \mathbf{X}_h. \quad (5.24)$$

Since \mathbf{P}_h is the orthogonal projection with respect to the inner product $B_1(\cdot, \cdot)$, then

$$\|\mathbf{U} - \mathbf{P}_h \mathbf{U}\|_E = \inf_{\chi \in \mathbf{X}_h} \|\mathbf{U} - \chi\|_E, \quad \mathbf{U} \in \mathbf{X}. \quad (5.25)$$

We also have the relation between operators $\mathbf{T}_{h,ij}$ and \mathbf{T}_{ij} :

$$\mathbf{T}_{h,ij} = \mathbf{P}_h \mathbf{T}_{ij}|_{\mathbf{X}_h}. \quad (5.26)$$

We now define \mathbf{E}_h , the spatial approximation error,

$$\mathbf{E}_h \stackrel{\text{def}}{=} \mathbf{U}_\tau(t^*) - \mathbf{U}_{th}(t^*), \quad (5.27)$$

and we estimate it in the energy norm. (For simplicity, we consider a 2-stage direct TG scheme.)

We first record the following inequality:

$$\begin{aligned} & \|\mathbf{P}_h \mathbf{T} \mathbf{U} - \mathbf{T}_h \mathbf{P}_h \mathbf{U}\|_E \\ &= \|\mathbf{P}_h(\mathbf{T}_{20} + \mathbf{T}_{21} \mathbf{T}_{10}) \mathbf{U} - (\mathbf{P}_h \mathbf{T}_{20} + \mathbf{P}_h \mathbf{T}_{21} \mathbf{P}_h \mathbf{T}_{10}) \mathbf{P}_h \mathbf{U}\|_E \\ &\leq (\|\mathbf{T}_{20}\|_E + \|\mathbf{T}_{21}\|_E \|\mathbf{T}_{10}\|_E) \|\mathbf{U} - \mathbf{P}_h \mathbf{U}\|_E \\ &\quad + \|\mathbf{T}_{21}\|_E \|\mathbf{T}_{10} \mathbf{U} - \mathbf{P}_h \mathbf{T}_{10} \mathbf{U}\|_E \\ &\leq C_2 \max\{\|\mathbf{U} - \mathbf{P}_h \mathbf{U}\|_E, \|\mathbf{T}_{10} \mathbf{U} - \mathbf{P}_h \mathbf{T}_{10} \mathbf{U}\|_E\}, \end{aligned} \quad (5.28)$$

where

$$C_2 = \max\{\|\mathbf{T}_{20}\|_E + \|\mathbf{T}_{21}\|_E \|\mathbf{T}_{10}\|_E, \|\mathbf{T}_{21}\|_E\}. \quad (5.29)$$

Next, we observe that transient operator \mathbf{T} (as well as operators \mathbf{T}_{ij}) preserve the regularity of \mathbf{U} in the sense that $\mathbf{U} \in \mathbf{H}^m(\Omega) \cap \mathbf{X}$ implies that $\mathbf{T}\mathbf{U} \in \mathbf{H}^m(\Omega) \cap \mathbf{X}$. Indeed, for $\mathbf{U} \in \mathbf{N}(\mathbf{A})$ we have $\mathbf{T}\mathbf{U} = \mathbf{U}$, so it is sufficient to consider $\mathbf{U} \perp \mathbf{N}(\mathbf{A})$. But for $\mathbf{U} \perp \mathbf{N}(\mathbf{A})$, $\exists q$ such that

$$\mathbf{u} = \mathbf{grad} \, q \quad (5.30)$$

and a typical one-stage problem is reduced to scalar problems of the form $w - a \nabla^2 w = f$, for which the usual regularity results apply. To illustrate, we consider the case $\partial\Omega = \Gamma_u$, $\Gamma_p = \emptyset$; i.e., the solid rigid wall boundary conditions defined by (4.6), and a typical one-stage problem, say $\mathbf{Z}_1 = \mathbf{T}_{10} \mathbf{Z}_0$:

$$\mathbf{Z}_1 - \eta \Delta t^2 \begin{pmatrix} \mathbf{grad} \, \text{div} & \mathbf{0} \\ \mathbf{0} & \text{div} \, \mathbf{grad} \end{pmatrix} \mathbf{Z}_1 = \mathbf{F}(\mathbf{Z}_0). \quad (5.31)$$

If $\mathbf{Z}_0 \in \mathbf{D}(\mathbf{A}^m) \cap \mathbf{X}$, $m \geq 2$, then $\mathbf{F} \in \mathbf{H}^{m-2}(\Omega)$. Let $\mathbf{Z}_1 = [\mathbf{u}^\top, p]^\top$ and $\mathbf{F} = [\mathbf{f}_u^\top, f_p]^\top$. We note that (5.31) decouples into velocity and pressure equations. Accordingly, the velocity equations under condition (5.30) take the form

$$\mathbf{grad} \, (q - \eta \Delta t^2 \nabla^2 q) = \mathbf{f}_u, \quad \partial q / \partial n|_{\partial\Omega} = 0, \quad (5.32)$$

which imply that

$$q - \eta \Delta t^2 \nabla^2 q = f_q, \quad \partial q / \partial n|_{\partial\Omega} = 0, \quad (5.33)$$

with $f_q \in H^{m-1}(\Omega)$. Assuming that $\partial\Omega$ is a C^{m+1} -manifold, it may be shown (see, e.g., [10]) that $q \in H^{m+1}(\Omega)$, and, hence, $\mathbf{u} = \mathbf{grad} \, q \in [H^m(\Omega)]^\nu$. Likewise, the pressure equation reads

$$p - \eta \Delta t^2 \nabla^2 p = f_p, \quad \partial p / \partial n|_{\partial\Omega} = 0, \quad (5.34)$$

with $f_p \in H^{m-2}(\Omega)$. Assuming that $\partial\Omega$ is a C^m -manifold, we have from standard regularity results that $p \in H^m(\Omega)$.

Next, we make an assumption that (asymptotically) the time-step Δt and the mesh parameter h satisfy

$$\Delta t/h = \sigma, \quad (5.35)$$

where σ is a constant. Condition (5.35) should be understood in a sense of regularity of partition of $[0, t^*] \times \Omega$ and *not* as a stability condition. With (5.35) in force, we have

$$\begin{aligned}
& \|U - P_h U\|_E \\
& \leq \inf_{\chi \in X_h} \{ \|U - \chi\| + \sqrt{\eta} \sigma h \|A(U - \chi)\| \} \\
& \leq C_3 h^s \|U\|_{H^s(\Omega)} \quad \forall U \in X \cap H^s(\Omega), \\
& \quad 1 \leq s \leq p + 1,
\end{aligned} \tag{5.36}$$

where C_3 is a constant independent of h and U . Hence, from (5.28)–(5.36), we get the following estimate:

$$\begin{aligned}
& \|P_h T U - T_h P_h U\|_E \leq C_2 C_4 h^s \|U\|_{H^s(\Omega)} \\
& U \in X \cap H^s(\Omega), \quad 1 \leq s \leq p + 1,
\end{aligned} \tag{5.37}$$

where $C_4 = C_3 \max\{1, \|T_{10}\|_{H^s(\Omega)}\}$. Finally, we note the uniform quasi-boundedness of T_h^N ,

$$\|T_h^N\|_{Eh} \stackrel{\text{def}}{=} \sup_{\substack{U \in X_h \\ U \neq 0}} \frac{\|T_h^N U\|_E}{\|U\|_E} \leq M \exp(Bt^*) \quad \text{independent of } N \tag{5.38}$$

which follows from unconditional stability of TG schemes at finite dimensions.

Now, the spatial approximation error $\|E_h\|_E$ is bounded as

$$\begin{aligned}
\|E_h\|_E &= \|T^N U_0 - T_h^N P_h U_0\|_E \\
&\leq \sum_{i=1}^N \|T_h^{N-i}\|_{Eh} \| (P_h T - T_h P_h) T^{i-1} U_0 \|_E \\
&\quad + \| (T^N U_0 - P_h T^N U_0) \|_E \\
&\leq M \exp(Bt^*) N \| (P_h T - T_h P_h) T^{i-1} U_0 \|_E \\
&\quad + \| (T^N U_0 - P_h T^N U_0) \|_E \\
&\leq M \exp(Bt^*) N C_2 C_4 h^s \|T^{i-1} U_0\|_{H^s(\Omega)} \\
&\quad + C_3 h^s \|T^N U_0\|_{H^s(\Omega)} \\
&\leq M \exp(Bt^*) \sigma t^* C_2 C_4 h^{s-1} \|T^{i-1} U_0\|_{H^s(\Omega)} \\
&\quad + C_3 h^s \|T^N U_0\|_{H^s(\Omega)}, \\
&U_0 \in H^s(\Omega) \cap X, \quad 2 \leq s \leq p + 1,
\end{aligned} \tag{5.39}$$

where

$$\begin{aligned}
& \| (P_h T - T_h P_h) T^{i-1} U_0 \|_E \\
&= \max_{i=1,2,\dots,N} \{ \| (P_h T - T_h P_h) T^{i-1} U_0 \|_E \},
\end{aligned} \tag{5.40}$$

and C_2, C_3, C_4, M, B , and σ are independent of h and Δt .

Finally, combining (5.19) and (5.39), we get the estimate of the total approximation error $\|E\|_E$,

$$\begin{aligned}
\|E\|_E &= \|U(t^*) - U_{th}(t^*)\|_E \leq \|E_\tau\|_E + \|E_h\|_E \\
&\leq C(t^*) [\Delta t^m \|A^{m+1} U_0\|_E + h^{s-1} \|U_0\|_{H^s(\Omega)}],
\end{aligned} \tag{5.41}$$

$$U_0 \in \{D(A^{m+2}), m + 2 \geq \mu\} \cap \{H^s(\Omega) \cap X, 2 \leq s \leq p + 1\}.$$

6. NUMERICAL EXAMPLES

EXAMPLE 1 (*The vibrating cylinder problem (2D)*). The purpose of this example is to illustrate the convergence properties of TG schemes. In particular, the focus is on temporal approximation error. The test problem is defined as follows (see Fig. 1).

1. Governing equations (2.13) are to be solved in $\Omega \times (0, t^*]$,

$$\Omega = \{(r, \theta): r_i < r < r_o, \theta_1 < \theta < \theta_2\} \tag{6.1}$$

with $r_i = 1, r_o = 1.6, \theta_1 = -\theta_2 = 5^\circ$, and $t^* = 0.5$.

2. Boundary conditions,

$$u_n = \begin{cases} A \exp(|\tau|^2 - \varepsilon^2)^{-1}, & |\tau| < \varepsilon \\ 0, & |\tau| \geq \varepsilon \end{cases} \quad \text{at } r = r_i; \tag{6.2}$$

$$\begin{aligned}
u_n &= 0 \quad \text{at } \theta = \theta_1, \theta = \theta_2, \\
p &= 0 \quad \text{at } r = r_o,
\end{aligned} \tag{6.3}$$

where $A = -0.005, \varepsilon = 0.25$, and $\tau = t - 0.25$.

3. Initial condition,

$$U_0 = 0. \tag{6.4}$$

We choose a fixed finite element mesh consisting of one layer of elements of order $p = 4$, with the mesh-size parameter in

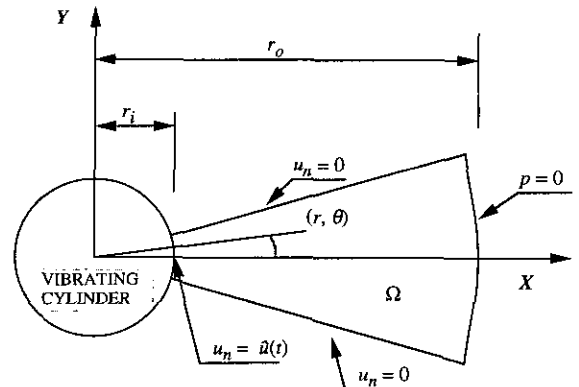


FIG. 1. The vibrating cylinder problem.

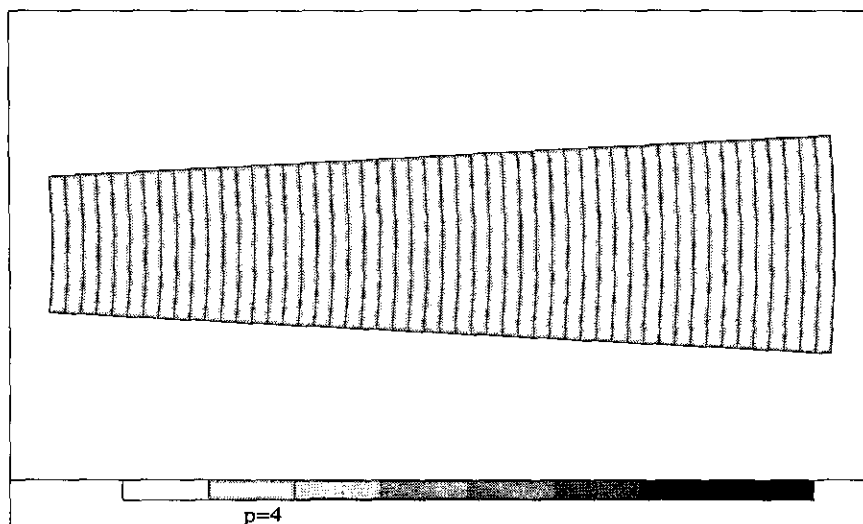


FIG. 2. The vibrating cylinder problem. Mesh elements are of fourth order.

radial direction $h = 0.012$ (Fig. 2). Note that, since the temporal approximation error can be observed at the continuous level only (i.e., when no spatial approximation is introduced), it is necessary to select a rich enough finite dimensional finite element space, so that the spatial approximation error can be neglected.

Figure 3 shows the graph of the relative error

$$\|E\|_R = \|U(t^*) - U_h(t^*)\| \|U(t^*)\|^{-1} \quad (6.5)$$

as a function of the time-step Δt for various TG schemes. (Note that the exact solution U can be found by means of the Fourier-Bessel series.) The graph, which is plotted in a log-log scale, is a straight line and thus

$$\|E\|_R = C \Delta t^\alpha, \quad (6.6)$$

where α is the slope of the line.

For comparison, we also plot the results for diagonally implicit Runge-Kutta (RK) schemes (DIRKs), which, together with singly implicit RKs (SIRKs), are the only schemes comparable with TGs. In particular, DIRKs and SIRKs enjoy the properties of unconditional stability for certain classes of problems, high-order accuracy, and semi-implicitness. (For more details on implicit RK methods for solving ODEs see the monograph of Butcher [1]; see also [7] for details on the application of implicit RK methods for problems in linear acoustics.) We use the following notation for DIRKs: $\text{DIRK}(s, m) = s$ -stage m th-order diagonally implicit RK scheme. In addition, $\text{TG}(2,4)$ -

I and $\text{TG}(2,4)$ -II denote direct TG schemes defined by (3.13) and (3.14) of [8].

From the analysis of Fig. 3 we can draw the following conclusions:

- all the schemes exhibit rates of convergence very close to those theoretically predicted;
- the high-order methods give results which are orders-of-magnitude better than those given by low-order methods, cf., e.g., $\text{TG}(1, 2)$ and $\text{TG}(3, 6)$;
- the TG schemes perform better than DIRKs (at least, for this particular problem);

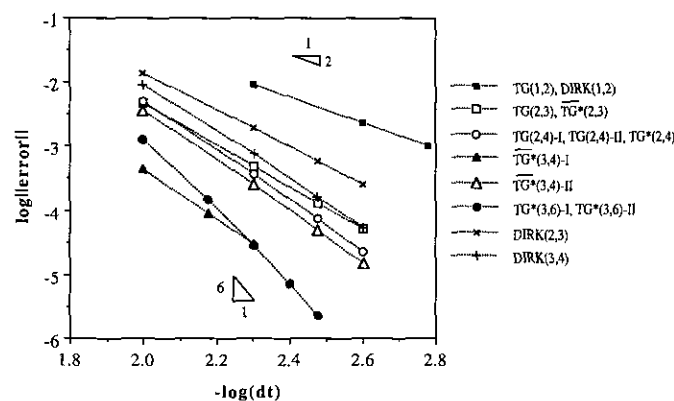


FIG. 3. The vibrating cylinder problem. The relative approximation error measured in L^2 -norm as a function of time step Δt .

—the results for certain schemes are indistinguishable and, therefore, are marked by the same symbol, cf., e.g., TG(2,3) and $\overline{\text{TG}}^*(2, 3)$; this confirms theoretical results that these schemes give *identical* results for linear problems at the infinite-dimensional level (i.e., they yield identical approximations to e^{-iAt}).

Finally, Fig. 4 shows the graph of $\|\mathbf{E}\|_R$ as a function of the number of (re)solutions of the resulting linear system, which is a measure of the cost of the solution. Recall, that, e.g., one step of the size Δt of TG(2, 4) is equivalent to two steps of the size $\Delta t/2$ of TG(1, 2), etc. The results depicted in Fig. 4 indicate that the high order schemes are cost-effective as compared with low-order schemes; cf., TG(1, 2) and TG(3, 6). Again, TGs perform better than DIRKs. In fact, the differences are still larger in favor of TGs, as they possess a built-in operator splitting and the resulting linear system is symmetric and positive definite. On the other hand, DIRKs yield an unsymmetrical linear system with all the equations coupled.

EXAMPLE 2 (Diffraction of a plane wave from a rigid circular cylinder (2D)). As the second example we consider the problem of diffraction of a plane wave consisting of a short train of pulses by a rigid cylinder (Fig. 5). The purpose of this example is to demonstrate the results of a more complicated finite element simulation; the computational and algorithmic details are skipped (they can be found in [8]).

The problem is of interest as its solution is known for the time-harmonic case, which, in turn, may be used to construct the solution of a transient problem. The time-harmonic problem is formulated as

$$\begin{aligned} \Phi_{s,t} - c_0^2 \nabla^2 \Phi_s &= 0 \quad \text{in } \mathbb{R}^2 - \{(x, y): x^2 + y^2 \leq a^2\} \\ (\Phi_p + \Phi_s)_{,r} |_{r=a} &= 0 \\ \sqrt{r}(\Phi_{,r} - ik\Phi) &= 0 \quad \text{as } r \rightarrow \infty, \end{aligned} \tag{6.7}$$

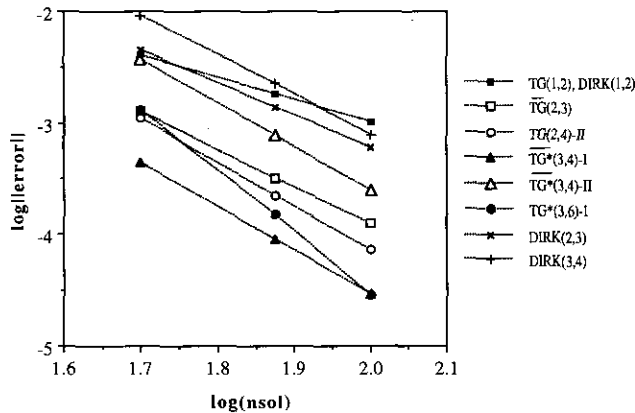


FIG. 4. The vibrating cylinder problem. The relative approximation error measured in L^2 -norm as a function of the number of solutions.

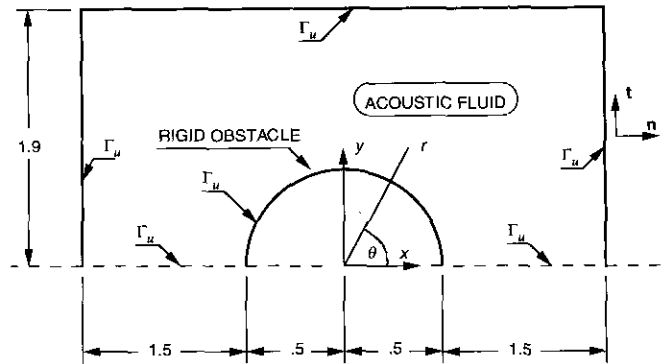


FIG. 5. Scattering of a plane wave by a rigid cylinder. Problem definition.

where

- $\Phi_p = A_p e^{i(kx - i\omega t)}$ = velocity potential of the incoming wave
- Φ_s = velocity potential of the scattered wave
- $\Phi = (\Phi_p + \Phi_s)$ = velocity potential
- ∇^2 = the Laplace operator
- $k = \omega/c_0$ = the wave number
- A_p = amplitude of the incident wave
- a = radius of the cylinder.

The solution of (6.7) reads (see, e.g., [4]):

$$\begin{aligned} \Phi_s(r, \theta, t) &= e^{-i\omega t} \sum_{m=0}^{\infty} B_m H_m^{(1)}(kr) \cos(m\theta) \\ B_0 &= -A_p \frac{J_0'(ka)}{H_0^{(1)'}(ka)}, \quad B_m = -2A_p i^m \frac{J_m'(ka)}{H_m^{(1)'}(ka)}, \quad m > 0, \end{aligned} \tag{6.8}$$

where J_m ($H_m^{(1)}$), $m = 0, 1, \dots$, are the Bessel (resp., Hankel) functions of the first kind. The time-harmonic velocity and pressure, $\mathbf{u}_\omega = (u_{\omega x}, u_{\omega y})$ and p_ω , are computed from the velocity potential Φ as follows:

$$\mathbf{u}_\omega = -\text{grad} \Phi, \quad p_\omega = -\rho_0 c_0 i \Phi. \tag{6.9}$$

The solution of a transient problem is given formally by superposition of (6.9) for various frequencies. In particular, if the initial condition for the transient problem is representable in the form of Fourier integral,

$$\begin{aligned} g(x) \equiv u_x(x, y, 0) &= (\rho_0 c_0)^{-1} p(x, y, 0) \\ &= \int_0^\infty a(\omega) \sin(\omega x/c_0) d\omega + \int_0^\infty b(\omega) \cos(\omega x/c_0) d\omega, \end{aligned} \tag{6.10}$$

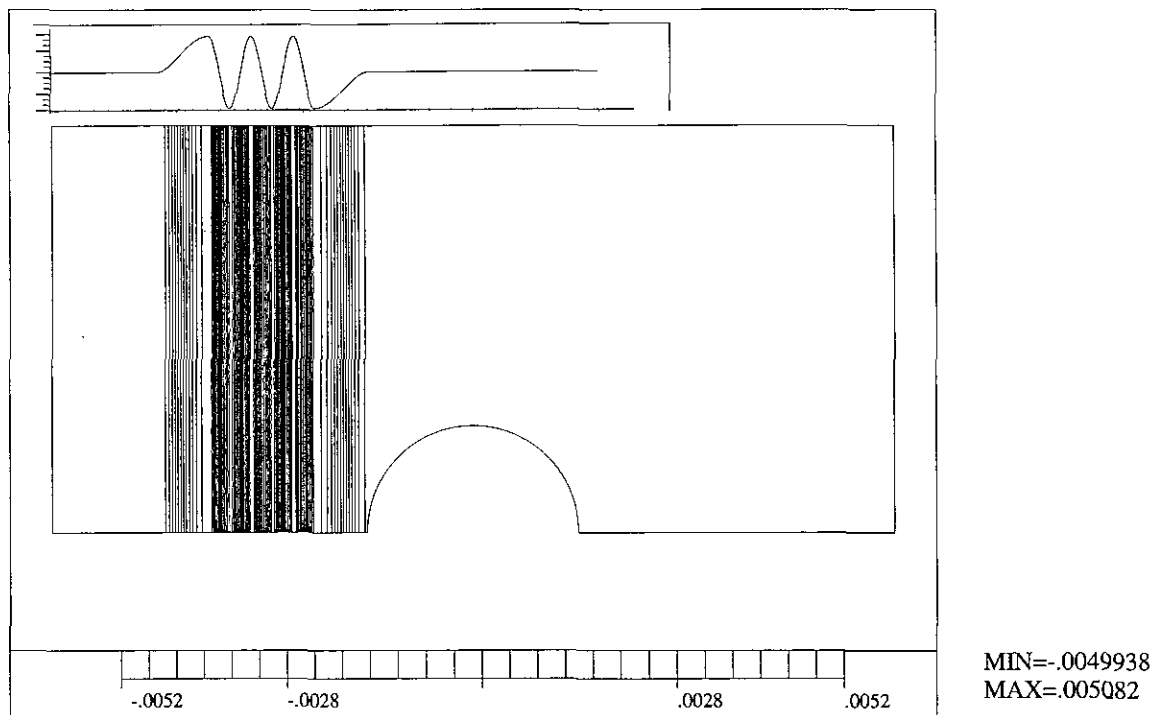


FIG. 6. Rigid scattering problem. Initial condition, acoustical pressure.

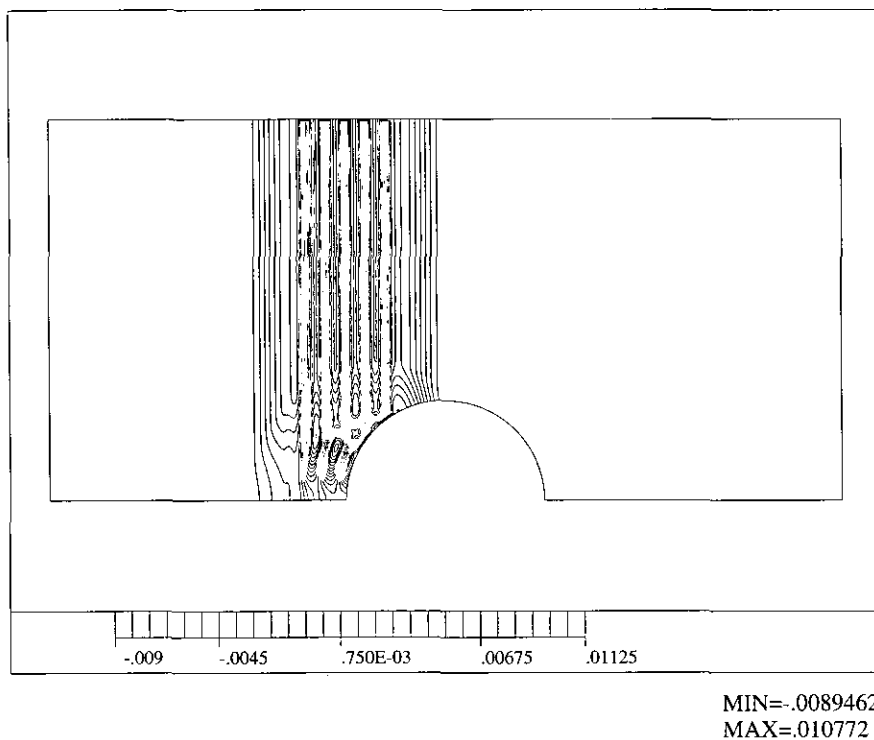


FIG. 7. Rigid scattering problem. The pressure at time $t = 0.5$.

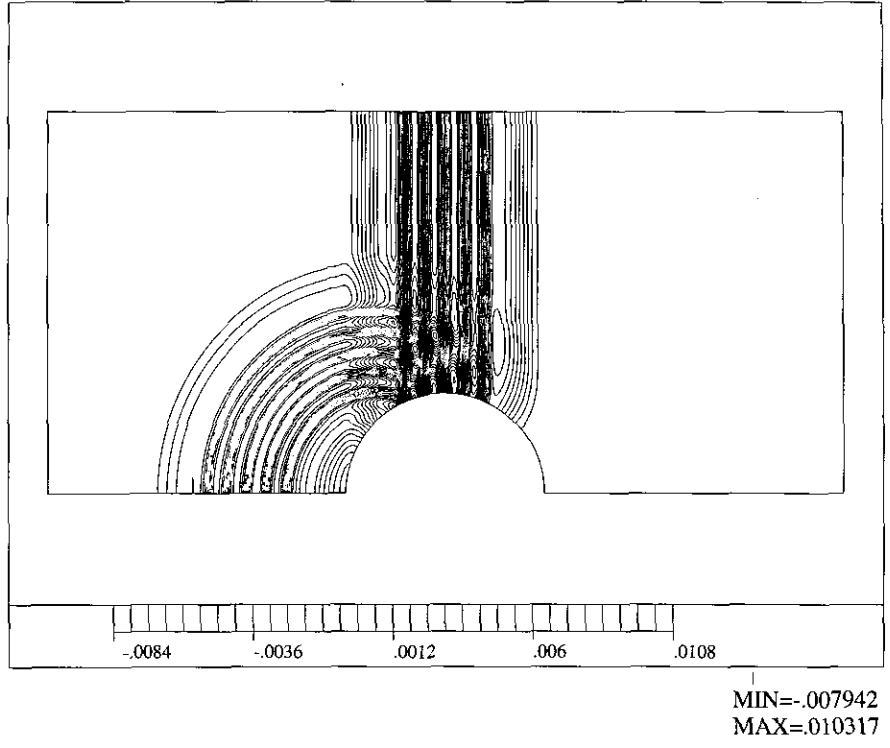


FIG. 8. Rigid scattering problem. The pressure at time $t = 1.0$.

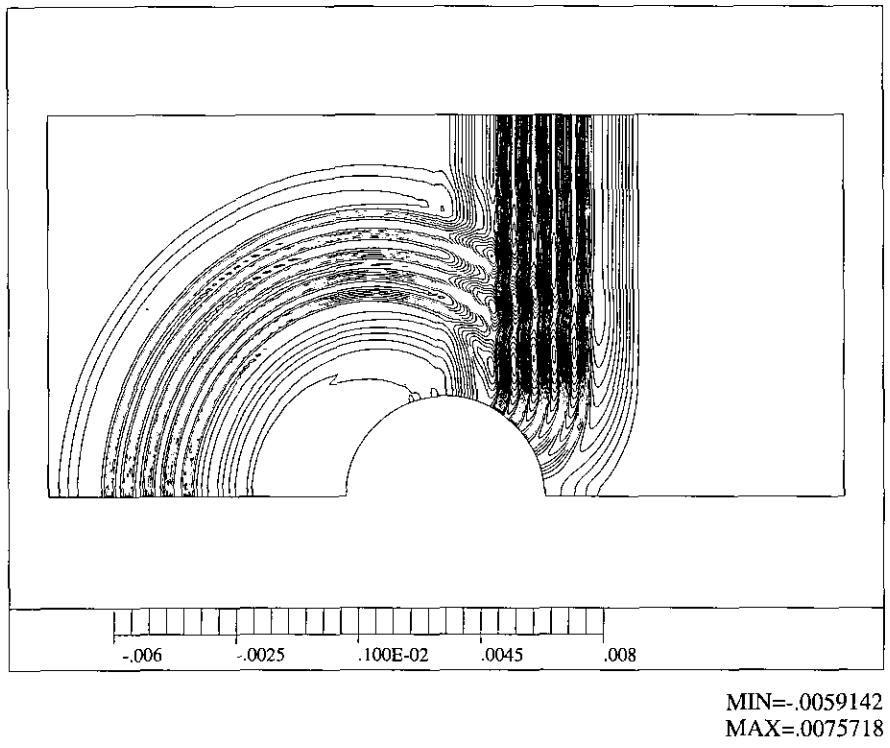


FIG. 9. Rigid scattering problem. The pressure at time $t = 1.5$.

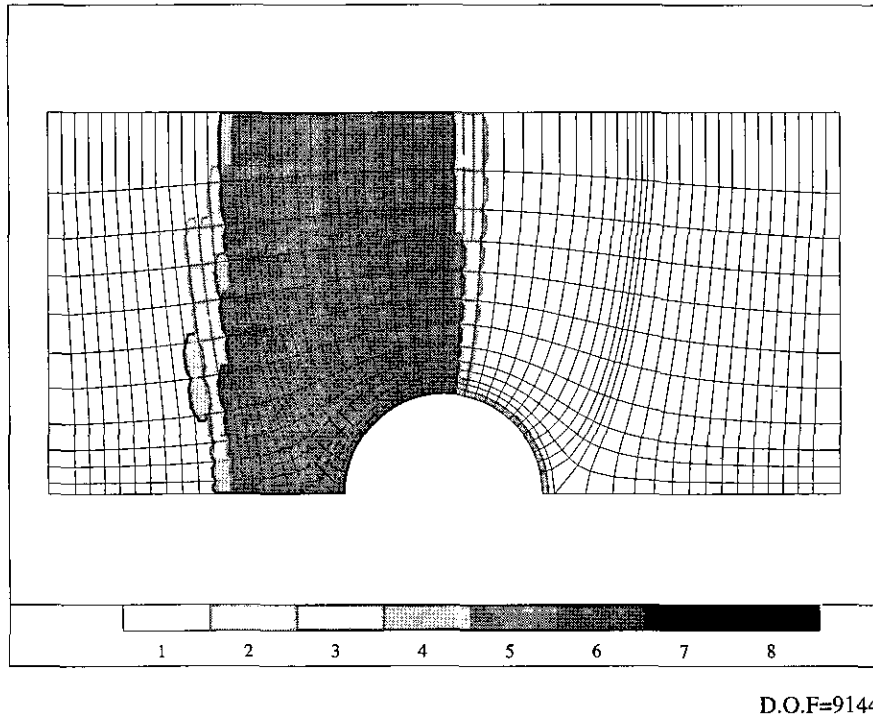


FIG. 10. Rigid scattering problem. Finite element mesh at time $t = 0.5$. Different shades correspond to different element spectral orders p .

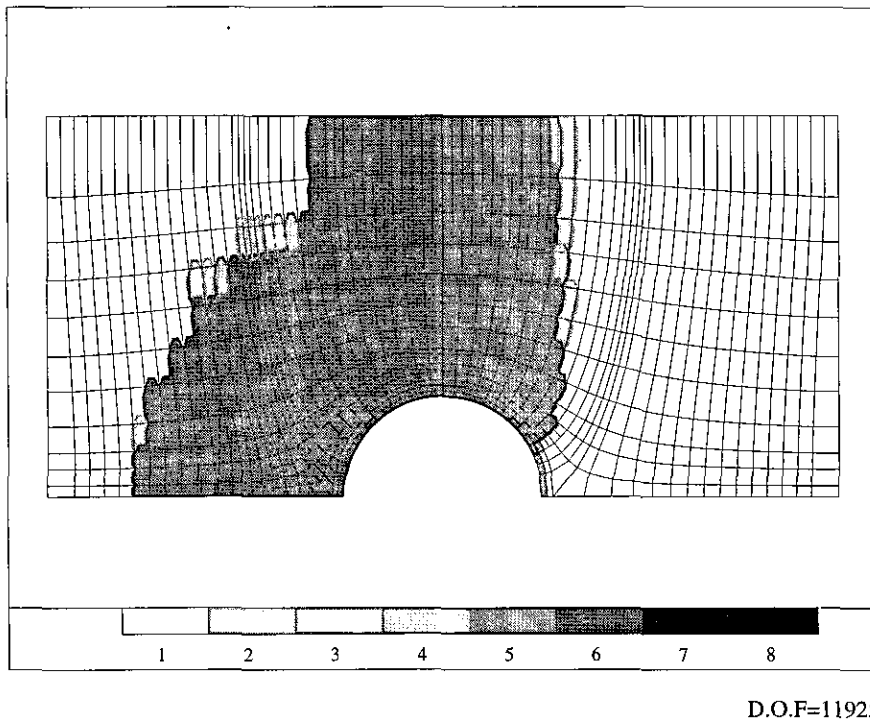
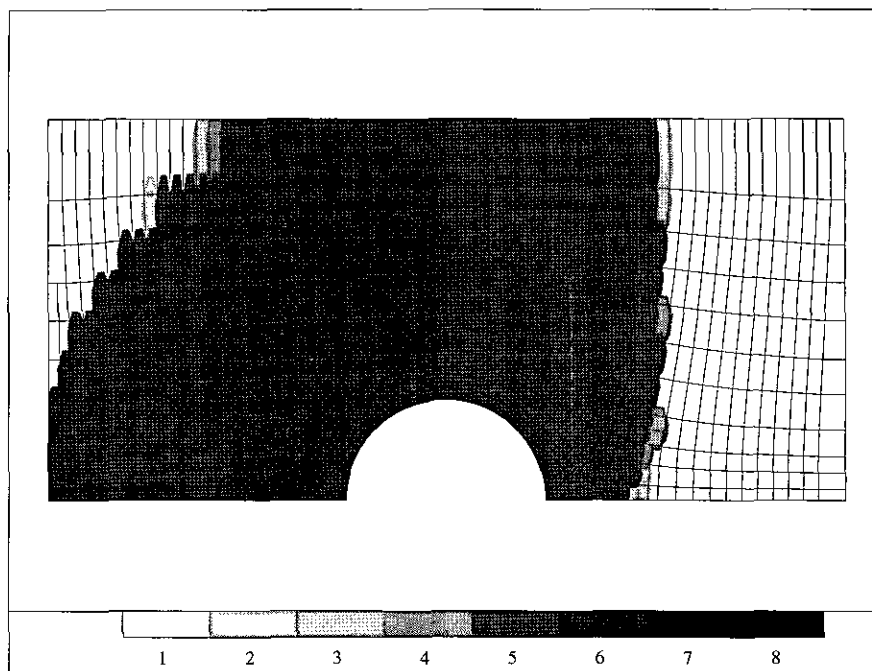


FIG. 11. Rigid scattering problem. Finite element mesh at time $t = 1.0$. Different shades correspond to different element spectral orders p .



D.O.F.=17521

FIG. 12. Rigid scattering problem. Finite element mesh at time $t = 1.5$. Different shades correspond to different element spectral orders p .

where

$$\begin{aligned} a(\omega) &= (2\pi c_0)^{-1} \int_{-\infty}^{\infty} g(\xi) \sin(\omega \xi / c_0) d\xi \\ b(\omega) &= (2\pi c_0)^{-1} \int_{-\infty}^{\infty} g(\xi) \cos(\omega \xi / c_0) d\xi, \end{aligned} \quad (6.11)$$

then the solution of the transient problem reads

$$\begin{aligned} u_x(x, y, t) &= \int_0^{\infty} [a(\omega) \operatorname{Re}(u_{\omega x}) + b(\omega) \operatorname{Im}(u_{\omega x})] d\omega \\ u_y(x, y, t) &= \int_0^{\infty} [a(\omega) \operatorname{Re}(u_{\omega y}) + b(\omega) \operatorname{Im}(u_{\omega y})] d\omega \\ p_x(x, y, t) &= \rho_0 c_0 \int_0^{\infty} [a(\omega) \operatorname{Re}(p_{\omega}) + b(\omega) \operatorname{Im}(p_{\omega})] d\omega. \end{aligned} \quad (6.12)$$

(In (6.12), it is understood that Φ_p should be scaled such that $|\Phi_{p,x}| = 1$.)

We assume that $\rho_0 = c_0 = 1$, $t \in (0, 1.5]$, and we solve the transient problem for the following initial condition function $g(x)$,

$$g(x) = A[H(x + \frac{3}{2}) - H(x + \frac{1}{2})] \times \begin{cases} 1 - 32(\frac{5}{4} + x)^2 + 256(\frac{5}{4} + x)^4, & -\frac{3}{2} \leq x < -\frac{5}{4} \\ -\cos[10\pi(x + \frac{3}{4})], & -\frac{5}{4} \leq x < -\frac{3}{4} \\ -1 + 32(\frac{3}{4} + x)^2 - 256(\frac{3}{4} + x)^4, & -\frac{3}{4} \leq x \leq -\frac{1}{2} \end{cases}, \quad (6.13)$$

where $A = 0.005$. Note, that $ka \approx 15.7$ for the fundamental frequency of $g(x)$, (6.13).

The problem was solved by using TG(2,4)-II method with stability parameter $\eta = 0.471$ and constant time step $\Delta t = 0.01$. Figure 6 shows the initial condition function $p(x, 0)$ in the form of contour map, and Figs. 7, 8, and 9 show pressure distribution at time $t = 0.5, 1.0$, and 1.5 , respectively. In general, we observe a complicated pattern of acoustical field with a non-zero value at the "shadowed" part of the cylindrical obstacle. Figures 10–12 show the corresponding finite element meshes. Finally, the relative difference $\|\mathbf{E}\|_R$ between the finite element solution at time $t = 1.5$ and the "exact" solution, at time $t = 1.5$, was 0.008. ("Exact" means that (6.12) was approximated by using 200 harmonics).

APPENDIX: COEFFICIENTS FOR TAYLOR GALERKIN SCHEMES

TABLE A.I

Coefficients for $TG^*(2, 4) \subset TG^*(3, 5)$ Schemes

$TG^*(2,4) \subset TG(3,5) , \eta = 0.47048000000000000000$		
Matrix R		
0.58666465485616976956	-0.80982997150172899736	0.
0.80982997150172899736	0.58666465485616976956	0.
0.	0.	1.
Matrix N		
0.75470195223744144069	-0.39233905362347625774	0.
0.20589874341489391717	0.18625804776255855930	0.
-0.00054376374825531749	-0.54497970609218374031	0.47048000000000000000
Matrix M		
0.40625593944633613718	-0.56079437058449297368	0.
0.29430375373313155419	-0.40625593944633613718	0.
0.	0.48013870337838156715	0.
Matrix K		
1.	-1.86292057229752427631	3.05310658997760529390
1.	1.11195218571320458299	1.10784490598270736102
1.	0.51986129662161843284	0.09490476646205749064
Vector c		
-2.01745900343568111280	1.	1.

TABLE A.II

Coefficients for $TG^*(3, 6)$ -I Scheme

$TG^*(3,6)-I , \eta = 0.69038000000000000000$		
Matrix R		
0.26393836138048191070	-0.95416090162990922853	-0.14111525499599545264
0.08633218817953910556	0.16908590957365693888	-0.98181297020755524192
0.96066815013089654495	0.24695531778936810990	0.12700305641847137595
Matrix N		
-1.51477676690670500019	10.55956378899615263605	-3.85279357060094799420
-0.60998123986938262505	3.56073337017473914750	-0.96737419432330487490
-0.00003694734022275365	0.34223800736008426286	0.02518339673196585268
Matrix M		
-2.59847697113301497696	9.39372782597028278750	1.38928171890125801151
-0.71878626830570126784	2.59847697113301497696	0.38430073979795441312
0.	0.	0.
Matrix K		
1.	-10.83328594134298975378	-13.83679053292063218969
1.	-1.85087132591624429169	-5.25971521870792121198
1.	1.	0.13261554324817263809
Vector c		
-2.64875336760446393172	0.41312011670902383054	1.

TABLE A.III

Coefficients for TG*(3, 6)-II \subset TG*(4, 7) Schemes

TG*(3,6)-II \subset TG(4,7), $\eta = 0.69038000000000000000$			
Matrix R			
0.02224184480109202572	-0.98125984384904011953	0.19140119955005908246	0.
0.07721500339363600078	-0.18919064667067411723	-0.97889976119276529601	0.
0.99676634352829395506	0.03655158083707934575	0.07156003317177499355	0.
0.	0.	0.	1.
Matrix N			
2.73658233602343014522	-78.99547051648442543859	11.84772656147499912431	0.
0.04565302849603750074	-0.91014957802886358317	0.18161651682403524061	0.
-0.00372950210283742613	0.97423418370868410104	0.24470724200543343795	0.
-0.00000060450236298743	0.18901866121515345568	-0.71205506269416832456	η
Matrix M			
2.21375357137569219934	-76.46444931464393297412	8.22108964178280425021	0.
0.04921903580116693922	-1.39397566614026097193	0.02644123841848942350	0.
-0.00491794189357611907	1.73905947691862811394	-0.81977790523543122740	0.
0.	0.	0.26926492895827408295	0.
Matrix K			
1.	62.67439990893028457962	90.97277191213729001705	
1.	1.60245261047011668961	1.25802615598729407174	
1.	0.08563637021037923253	-0.40606625013563028600	
1.	0.73073507104172591704	0.06339207702310377335	
Vector c			
-3.35520619255615194494	0.28413721854951208041	1.	1.

TABLE A.IV

Coefficients for TG*(2, 3) \subset TG*(3, 4) Schemes

TG*(2,3) \subset TG(3,4), $\eta = 0.47275000000000000000$			
Matrix R			
0.07006051682826013583	0.99754274293473614459	0.	
-0.99754274293473614459	0.07006051682826013583	0.	
0.	0.	1.	
Matrix N			
0.59709758822264003580	1.77049841834590787948	0.	
-0.00873331630040819613	0.34840241177735996419	0.	
-0.01109705484164502740	-0.26962298411765853121	0.47275000000000000000	
Matrix M			
0.17409175710303003393	2.47877087930407633877	0.	
-0.01222700336859269107	-0.17409175710303003393	0.	
0.	0.	0.	
Matrix K			
1.	0.63272527281421274541	-0.02081670349477487591	
1.	1.18631876047162272501	0.37459555605993475163	
1.	1.	0.30797003895930355862	
Vector c			
3.28558790922131911812	1.	1.	

TABLE A.V

Coefficients for $\overline{TG}^*(3, 4) \subset TG^*(4, 5)$ Schemes

$\overline{TG}^*(3, 4) - I \subset TG(4, 5), \eta = 0.10729000000000000000$			
Matrix R			
0.98399373314097109822	0.00936879186603509431	0.17795662077671139925	0.
0.00807236752521729137	0.99993535454537160828	-0.00800772207058889966	0.
-0.17802013937145755943	-0.00644301708773878444	-0.98400575074995061340	0.
0.	0.	0.	1.
Matrix N			
0.04729050669560846788	-0.00418195803292290401	0.01088211382826766613	0.
-0.56078289881371093275	0.11245728915328477198	0.10141948539993317719	0.
-0.30223807862492379408	-0.02339576765995930211	0.16212220415110676012	0.
-0.05219757041704398092	0.00471124028614898069	-0.19005412977308306273	η
Matrix M			
-0.21805739510262578416	-0.01519857626899199782	0.03954909056588031224	0.
-2.03806484686541434172	0.01877958546734692211	0.36859092603936268782	0.
-1.09843007825783911204	-0.08502772059216003849	0.19927780963527886205	0.
0.	0.	0.47070873531283376150	0.
Matrix K			
1.	0.05667232432254755363	-0.09169634979549908529	
1.	3.12027717368014897627	0.75126858867780009859	
1.	1.98417998921472028848	0.43866623286046816173	
1.	0.52929126468716623849	0.15954172459114430145	
Vector c			
-0.13703455648318991611	1.46958283832144424448	1.	1.

TABLE A.VI

Coefficients for $\overline{TG}^*(3, 4) \subset TG^*(4, 5)$ Schemes

$\overline{TG}^*(3, 4) - II \subset TG(4, 5), \eta = 0.57591000000000000000$			
Matrix R			
-0.21477634527149241550	0.27873156080796127890	0.93604478446353113660	0.
0.66545848012066667074	-0.65974399604188925659	0.34914591637619348548	0.
0.71486791279740130668	0.69788722347972340343	-0.04378687653224016150	0.
0.	0.	0.	1.
Matrix N			
1.05718156130944617817	-0.43363433767584146248	0.94589374433850670507	0.
0.36789979523718828505	0.28487109597861427048	1.36770262760438198847	0.
0.00087652096072221507	-0.01283343538248216781	0.38567734271193955133	0.
0.00705442445906670772	-0.00350543708780584948	-0.68967738115452365280	η
Matrix M			
0.61231152634559182099	-0.55170370436127152925	1.20344040437547900096	0.
0.46807105026318115092	-0.37028258030127063372	1.74009883571086130921	0.
0.00111517889374858822	-0.01632770568434363405	-0.24202894604432118727	0.
0.00000000000000000000	0.00000000000000000000	0.54573544698770473614	0.
Matrix K			
1.	1.23114114324622188687	0.39175628539340639937	
1.	1.02503687189268280971	0.22976011155218597312	
1.	1.25724147283491623309	0.41227091860113140593	
1.	0.45426455301229526385	0.06448294679555805841	
Vector c			
2.49518936960602117957	2.86292417756545463613	1.	1.

TABLE A.VII

Ranges of η Guaranteeing Unconditional Stability

TG(2,4)-I, TG(2,4)-II, TG*(2,4)	$\eta \cong \frac{1}{4} + \frac{1}{2} \sqrt{\frac{7}{36}}$
TG*(3,6)-I, TG*(3,6)-II	$\eta \cong \frac{1}{4} + \frac{\sqrt{7}}{6} \cos \left[\frac{1}{3} \arccos \left(\frac{183}{10\sqrt{343}} \right) \right]$
$\overline{\text{TG}}(2,3), \overline{\text{TG}}^*(2,3)$	$\eta \cong 0.4727411766\dots$
$\overline{\text{TG}}^*(3,4)$ -I	$\eta \cong 0.107280333\dots$
$\overline{\text{TG}}^*(3,4)$ -II	$\eta \cong 0.5759095682\dots$

ACKNOWLEDGMENTS

Support of this work by the Army Research Office (ARO) under Grant DAAL03-92-G-0253 and the Office of Naval Research (ONR) under Grant N00014-89-J-1451 is gratefully acknowledged.

REFERENCES

- Butcher, *The Numerical Analysis of Ordinary Differential Equations* (Wiley, New York, 1987).
- L. Demkowicz, J. T. Oden, W. Rachowicz, and O. Hardy, *Comput. Methods Appl. Mech. Engrg.* **88**, 363 (1991).
- J. Donea, *Int. J. Numer. Meth. Eng.* **20**, 101 (1984).
- M. C. Junger and D. Feit, *Sound, Structures and Their Interaction in Mathematical Physics*, 2nd ed. (MIT Press, Cambridge, MA, 1986).
- R. Leis, *Initial Boundary Value Problems in Mathematical Physics* (Wiley, New York, 1986).
- J. T. Oden, "Formulation and Application of Certain Primal and Mixed Finite Element Models of Finite Deformations of Elastic Bodies," in *Computing Methods in Applied Sciences and Engineering*, edited by R. Glowinski and J. L. Lions (Springer Verlag, New York/Berlin, 1974).
- A. Safjan, L. Demkowicz, and J. T. Oden, *Int. J. Numer. Methods Eng.* **32**, 677 (1991).
- A. Safjan and J. T. Oden, *Comput. Methods Appl. Mech. Engrg.* **103**, 187 (1993).
- A. Safjan and J. T. Oden, TICAM Report, 94-09, August 1994, (unpublished).
- R. E. Showalter, *Hilbert Space Methods for Partial Differential Equations* (Pitman, London, 1977).
- K. Yosida, *Functional Analysis* (Springer-Verlag, Berlin/Heidelberg/New York, 1974).