

MIT Open Access Articles

Unconditional Stability for Multistep ImEx Schemes: Theory

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Rosales, Rodolfo R., Benjamin Seibold, David Shirokoff, and Dong Zhou. "Unconditional Stability for Multistep ImEx Schemes: Theory." *SIAM Journal on Numerical Analysis* 55, no. 5 (January 2017): 2336–2360. © 2017 Society for Industrial and Applied Mathematics.

As Published: <http://dx.doi.org/10.1137/16M1094324>

Publisher: Society for Industrial & Applied Mathematics (SIAM)

Persistent URL: <http://hdl.handle.net/1721.1/115502>

Version: Final published version: final published article, as it appeared in a journal, conference proceedings, or other formally published context

Terms of Use: Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.



UNCONDITIONAL STABILITY FOR MULTISTEP IMEX SCHEMES: THEORY*

RODOLFO R. ROSALES[†], BENJAMIN SEIBOLD[‡], DAVID SHIROKOFF[§], AND
DONG ZHOU[‡]

Abstract. This paper presents a new class of high order linear ImEx (implicit-explicit) multistep schemes with large regions of unconditional stability. Unconditional stability is a desirable property of a time stepping scheme, as it allows the choice of time step solely based on accuracy considerations. Of particular interest are problems for which both the implicit and explicit parts of the ImEx splitting are stiff. Such splittings can arise, for example, in variable coefficient problems, or the incompressible Navier–Stokes equations. To characterize the new ImEx schemes, an unconditional stability region is introduced, which plays a role analogous to that of the stability region in conventional multistep methods. Moreover, computable quantities (such as a numerical range) are provided that guarantee an unconditionally stable scheme for a proposed ImEx matrix splitting. The new approach is illustrated with several examples. Coefficients of the new schemes up to fifth order are provided.

Key words. linear multistep ImEx, unconditional stability, ImEx stability, high order time stepping

AMS subject classifications. 65L04, 65L06, 65L07, 65M12

DOI. 10.1137/16M1094324

1. Introduction. When a stiff differential equation is solved via an explicit time stepping scheme, stability requires time steps that are much smaller than those imposed by accuracy. Implicit schemes can overcome this limitation. Unfortunately, for many practical problems, a fully implicit treatment may be structurally difficult or computationally costly. Implicit-explicit (ImEx) methods are based on splitting the problem into two parts: one to be treated implicitly, and the other explicitly. In many problems, the stiff modes can be conveniently treated implicitly, while the explicitly treated modes are nonstiff. Moreover, for many ImEx schemes a time step restriction is incurred from the explicit part, which is generally acceptable if it is nonstiff.

The study presented here is motivated by a different situation, namely, the case where an ImEx splitting is conducted for which both parts are stiff (see section 1.2 for examples in which this structure arises naturally). In that case, a time step restriction based on the explicit part is not acceptable. We therefore aim for more, namely, that the ImEx time stepping scheme, for the particular splitting, be unconditionally stable; i.e., arbitrarily large time steps can be chosen without losing stability.

At first glance it may sound impossible to achieve unconditional stability if some parts of the problem are treated explicitly. The reason why it *is* possible is that

*Received by the editors September 19, 2016; accepted for publication (in revised form) July 3, 2017; published electronically October 3, 2017.

<http://www.siam.org/journals/sinum/55-5/M109432.html>

Funding: This work was supported by the National Science Foundation through grants DMS–1318942 (Rosales) and DMS–1318709 (Seibold and Zhou) and was partially supported through grants DMS–1719637 (Rosales), DMS–1719693 (Shirokoff), and DMS–1719640 (Seibold and Zhou). D. Shirokoff was supported by a grant from the Simons Foundation (359610).

[†]Department of Mathematics, MIT, Cambridge, MA 02139 (rrr@mit.edu).

[‡]Department of Mathematics, Temple University, Philadelphia, PA 19122 (seibold@temple.edu, dzhou@temple.edu).

[§]Corresponding author. Department of Mathematical Sciences, NJIT, Newark, NJ 07102 (david.g.shirokoff@njit.edu).

the ImEx scheme is applied to problems and splitting choices that possess specific properties, so that the implicit part can stabilize any growing modes produced by the explicit part. This concept goes further than one may think: a properly chosen ImEx scheme can stabilize a large explicit part via a smaller implicit part (see section 5.1).

While the task outlined above is of interest for any time stepping scheme, this paper focuses on ImEx linear multistep methods (LMMs) [8, 15, 49]. These achieve a high order of accuracy by using information from previous time steps. Thus, in each time step, they need a single evaluation of the explicit part, and a single solve with the implicit part [28, Chapter II.3, p. 171]. Because high order multistep methods tend to possess less favorable stability properties than Runge–Kutta methods, the task of achieving unconditional stability is of particular importance.

1.1. Outline of the problem and contributions of this paper. The problem of interest is a linear system of ordinary differential equations

$$(1) \quad \vec{u}_t = \mathbf{L}\vec{u} + \vec{f}(t) \quad \text{with} \quad \vec{u}(0) = \vec{u}_0 ,$$

where $\vec{u}(t), \vec{u}_0, \vec{f}(t) \in \mathbb{R}^N$ and $\mathbf{L} \in \mathbb{R}^{N \times N}$ is a matrix. We assume that \mathbf{L} is stable; i.e., the homogeneous equation $\vec{u}_t = \mathbf{L}\vec{u}$ has solutions that remain bounded for all time (stability is independent of the forcing \vec{f}). The term $\mathbf{L}\vec{u}$ in problem (1) is now split into an implicit part ($\mathbf{A}\vec{u}$) and an explicit part ($\mathbf{B}\vec{u}$), transforming (1) into

$$(2) \quad \vec{u}_t = \mathbf{A}\vec{u} + \mathbf{B}\vec{u} + \vec{f}(t) ,$$

where $\mathbf{B}\vec{u} = \mathbf{L}\vec{u} - \mathbf{A}\vec{u}$.

Of course, the choice of splitting $\mathbf{L} = \mathbf{A} + \mathbf{B}$ is not unique. One approach is to choose \mathbf{A} as the stiff terms in \mathbf{L} (i.e., the terms that would give rise to unnecessarily small time step restrictions if treated explicitly) and \mathbf{B} as the nonstiff terms in \mathbf{L} . In such a case, one can guarantee stability for an ImEx LMM [20] by requiring a time step restriction roughly dictated by an explicit treatment of \mathbf{B} . However, as outlined above, here we are concerned with the situation where such a splitting strategy is not feasible/practical. Hence, we look for ImEx time stepping schemes that are unconditionally stable when applied to (2), where \mathbf{B} can involve stiff terms.

Whether a time stepping scheme (of whatever kind) for (1) or (2) is stable depends on both the scheme and the problem’s right-hand side \mathbf{L} . A classical approach (for non-ImEx schemes) in stability analysis [38, Chapter 7] is to separate stability into a property of the scheme and another property of the problem’s right-hand side, as follows. For a linear scheme, the region of absolute stability $S \subset \mathbb{C}$ is the set of all $z = k\lambda$, where k is the time step, for which the numerical solution remains bounded when applied to the test equation $u_t = \lambda u$. Similarly, one can define a region of unconditional stability $S_u = \{z \in \mathbb{C} : \mu z \in S \ \forall \mu \geq 0\}$ as the largest cone contained within S . If the eigenvalues of \mathbf{L} lie in S_u , then the scheme is unconditionally stable. This concept decouples the scheme stability analysis from the detailed properties of \mathbf{L} , relying on its spectrum $\sigma(\mathbf{L})$ only. Moreover, it allows one to make stability statements about whole classes of problems. For instance, if S_u is the cone $|\theta - \pi| < \alpha$, where $0 < \alpha < \pi/2$ and θ is the polar angle (i.e., the scheme is $A(\alpha)$ stable), then the scheme is unconditionally stable for all problems where \mathbf{L} is negative definite. Conversely, we know that the same scheme is not unconditionally stable if \mathbf{L} is skew-symmetric.

In this paper, an analogous concept is developed for the ImEx framework. This extension is not straightforward, because one now has two right-hand side operators \mathbf{A} and \mathbf{B} that, in general, do not commute and thus do not share a set of common eigenvectors (see section 1.3 for references to the commutative case).

While the fundamental idea of *stability criteria* for ImEx schemes has been presented before (see section 1.3), here we present sufficient criteria for unconditional stability that are less restrictive than prior work. The stability set \mathcal{D} that we introduce depends only on the coefficients of the ImEx schemes, and not on the matrices \mathbf{A} and \mathbf{B} in the splitting (2). Moreover, we devise new high order ImEx schemes with *very large* stability regions that can stabilize splittings of the form (2) which are unstable with current schemes (see section 5).

1.2. Motivating applications. While the ideas developed here apply to an abstract ODE system (2), particular interest lies in systems that arise from a method of lines [38, Chapter 9.2] discretization (e.g., via finite differences, finite elements, or spectral) of linear PDE problems. Two important applications are as follows (letting ∇_h denote the spatial discretization of ∇ in an appropriate basis with smallest length scale h):

- (i) Variable coefficient diffusion with

$$\mathbf{L}\vec{u} = \nabla_h \cdot (d(x) \nabla_h u), \quad \text{where } d(x) > 0.$$

Here \mathbf{L} can be split into a constant coefficient diffusion \mathbf{A} and a variable coefficient diffusion \mathbf{B} . Then fast solvers [24, 46] can treat \mathbf{A} efficiently. However, \mathbf{B} remains stiff, because it scales the same as \mathbf{A} (i.e., like $1/h^2$). See section 5.2 for more details.

- (ii) Nonlocal operators, such as the Stokes operator in the linearized Navier–Stokes equations, whose discretization either yields a dense matrix or requires the addition of extra variables through the introduction of Lagrange multipliers,

$$\mathbf{L}\vec{u} = \nu \nabla_h^2 u - \nabla_h p \quad \text{and constraint} \quad \nabla_h \cdot u = 0.$$

A splitting where $\nu \nabla_h^2 u$ is implicit can create a stiff explicit $\nabla_h p$ [32, 39, 44]. The theory in this paper does not directly apply to cases where $\vec{L}(\vec{u}, t)$ is nonlinear or time dependent, as arises, for instance, with discretizations of the Cahn–Hilliard equation [11]. However, the ideas presented below for linear splitting may nevertheless be useful in stabilizing more general splittings as well.

1.3. Existing results and the new contributions in context. The simplest ImEx scheme that can achieve unconditional stability is a first order in time combination of forward and backward Euler steps. The application to (2) yields

$$(3) \quad \frac{1}{k} (\vec{u}_{n+1} - \vec{u}_n) = \mathbf{A}\vec{u}_{n+1} + \mathbf{B}\vec{u}_n + \vec{f}(nk).$$

Here $k > 0$ is the time step, and \vec{u}_n is the numerical solution at time $t = nk$.

First order in time schemes that achieve unconditional stability originated with Douglas and Dupont [16]. Other first order approaches are (i) iterative schemes for steady state elliptic problems [14]; (ii) variable coefficient diffusion with spectral methods [23, Chapter 9]; (iii) nonlinear convex–concave splittings for the Cahn–Hilliard equation [19]; (iv) nonlocal explicit terms [7]; (v) Hele–Shaw flows [21]; (vi) phase-field models [10, 18, 43, 45]; and (vii) viscosity-pressure splittings in incompressible Navier–Stokes [32, 39].

A disadvantage of first order approaches is that, in addition to the low order, large error constants have been reported for stable splitting choices in dissipative equations [13], as well as dispersive equations [12].

Better accuracy requires higher order ImEx time stepping methods. The following are two of the most commonly used approaches, which can be applied to (2):

- **CN-AB:** Implicit Crank–Nicolson for $\mathbf{A}\vec{u}$, and explicit Adams–Bashforth extrapolation for $\mathbf{B}\vec{u}$.
- **SBDF (semi-implicit backward differentiation formula):** Implicit BDF for $\mathbf{A}\vec{u}$, and explicit Adams–Bashforth extrapolation for $\mathbf{B}\vec{u}$.

For second order schemes, unconditional stability, or at least the absence of a stiff time step restriction, has been reported in practice for the semi-implicit treatment of the incompressible Navier–Stokes equations [34, 35] and the Cahn–Hilliard equation [9]. Rigorous proofs that guarantee unconditional stability for second order ImEx schemes such as SBDF or CN-AB have been given for convex–concave splittings of gradient flow systems [22, 25], a coupled Stokes–Darcy system [37], and a system with an explicit treatment of nonlocal terms [48]. See also [18, 50] for an interpretation of some convex–concave splittings as fully implicit schemes with a rescaled time step.

Higher order semi-implicit schemes that guarantee unconditional stability are not as well studied as their first and second order counterparts. Some third order schemes for the Navier–Stokes equations have been found that do not require a diffusion-restricted time step [34, 40]. General sufficient conditions on \mathbf{A} and \mathbf{B} guaranteeing unconditional stability for any order of SBDF have been outlined in [5] and related works [3, 4]. Specifically [3, 4, 5] assume that \mathbf{A} is negative definite and also allow for \mathbf{B} to be nonlinear. The results in [5], applied to the case where \mathbf{B} is a matrix, guarantee unconditional stability for an SBDF scheme of order $1 \leq r \leq 6$ ¹ if

$$(4) \quad \|(-\mathbf{A})^{-1/2} \mathbf{B} (-\mathbf{A})^{-1/2}\|_2 < (2^r - 1)^{-1}.$$

In related work, a set of new second order ImEx coefficients was introduced in [6], allowing for a weaker upper bound in (4)—it can be made arbitrarily close to 1. The unconditional stability criteria devised here are more general than previous bounds such as (4). Instead of prescribing norm bounds, we introduce the concept of unconditional stability diagrams for ImEx schemes. The new diagrams generalize the previous work on ImEx stability regions [20] (see also [36]) to (i) the case of unconditional stability, and (ii) the case where \mathbf{A} and \mathbf{B} do not commute. We then prescribe a set of new ImEx coefficients and show that they can achieve unconditional stability for some problems which violate (4) by orders of magnitude. See also Chapter IV of [28] for an overview of different splitting methods for ODE integration. Other techniques for specific problems are (i) explicit Runge–Kutta schemes with very large stability regions for parabolic problems [1], (ii) semi-implicit deferred correction methods [42], and (iii) semi-implicit schemes when an integration factor (matrix exponential) is easily evaluated [33, 41].

This paper is organized as follows. In sections 2–3 we introduce ImEx LMMs, the new criteria for unconditional stability, and the definition of the unconditional stability region. In section 4 we define new ImEx coefficients, characterize their unconditional stability region, and examine their effect on the approximation error. Finally, section 5 demonstrates how a small implicit term may stabilize a large explicit term. It also provides an example showing how the new coefficients may be used to stabilize splittings (2) that arise from a variable coefficient diffusion problem. We conclude with tables of the new ImEx coefficients in section 7 so that they may be used by practitioners.

2. Mathematical foundations. The purpose of this paper is to examine ImEx LMMs (linear multistep methods) for splittings of the form (2), where $\mathbf{A}\vec{u}$ is treated

¹See [5, equations (1.4)–(1.5), Theorem 2.1, and Remark 2.3].

implicitly, and $\mathbf{B}\vec{u}$ explicitly. Moreover, we are particularly interested in the case where both \mathbf{A} and \mathbf{B} are stiff, i.e., each term alone would result in severely limited time steps (due to stability) when treated explicitly. The goal is to first devise simple sufficient conditions that guarantee unconditionally stability of a time stepping scheme when applied to (2). We will then devise new ImEx schemes that allow one to satisfy the simple unconditional stability conditions, thereby guaranteeing an unconditionally stable scheme.

Here we restrict \mathbf{A} to be real, self-adjoint, and negative definite. Thus, $\mathbf{A}^T = \mathbf{A}$, and $\langle \vec{x}, \mathbf{A}\vec{x} \rangle < 0$ for all $\vec{x} \neq \vec{0}$. We use the notation

$$\langle \vec{x}, \vec{y} \rangle = \vec{x}^T \vec{y} \quad \text{and} \quad \|\vec{x}\|^2 = \langle \vec{x}, \vec{x} \rangle \quad \text{for} \quad \vec{x}, \vec{y} \in \mathbb{C}^N.$$

Note that the restriction above, which is needed for the theoretical presentation in this paper, is not as limiting as it might seem. A self-adjoint, negative definite, matrix \mathbf{A} yields desirable properties for the efficient solution of linear systems [47, Chapter IV, lecture 38] with coefficient matrices of the form $(\mathbf{I} - \gamma\mathbf{A})$, with $\gamma > 0$. The need to solve such linear systems arises in the time stepping of LMMs, as well as for implicit Runge–Kutta schemes. Hence, even if the matrix \mathbf{L} is not symmetric (e.g., the discretization of a dispersive wave problem), it may still be advantageous to take \mathbf{A} to be symmetric and negative definite, with $\mathbf{B} := \mathbf{L} - \mathbf{A}$.

Let \vec{u}_n be the numerical solution of (2) at time $t = nk$, where k is the time step, and let $\vec{f}_n = \vec{f}(nk)$. Then an LMM with $s \geq 1$ steps takes the form

$$(5) \quad \frac{1}{k} \sum_{j=0}^s a_j \vec{u}_{n+j} = \sum_{j=0}^s \left(c_j \mathbf{A} \vec{u}_{n+j} + b_j \mathbf{B} \vec{u}_{n+j} + b_j \vec{f}_{n+j} \right),$$

where (a_j, b_j, c_j) , with $0 \leq j \leq s$, are the time stepping coefficients. Here we will assume that $b_s = 0$ and $a_s, c_s \neq 0$, so that the method is implicit in \mathbf{A} and explicit in \mathbf{B} —i.e., it is an ImEx time stepping scheme. To accompany (5), one must also supply s initial vectors $\vec{u}_0, \vec{u}_1, \dots, \vec{u}_{s-1}$.

We wish to avoid any unnecessarily small time step restriction, and therefore demand that the scheme (5) be *unconditionally stable*. That is, the solutions to (5), with $\vec{f} = 0$, remain bounded for arbitrarily large time steps $k > 0$. This leads to the following.

DEFINITION 1 (unconditional stability). *A scheme (5) is unconditionally stable if, when $\vec{f} = 0$, there exists a constant C such that*

$$\|\vec{u}_n\| \leq C \max_{0 \leq j \leq s-1} \|\vec{u}_j\| \quad \text{for all } n \geq s, k > 0 \quad \text{and} \quad \vec{u}_j \in \mathbb{R}^N, \quad \text{where } 0 \leq j \leq s-1.$$

Note that C may depend on the matrices \mathbf{A}, \mathbf{B} , and the coefficients (a_j, b_j, c_j) , but is independent of the time step k , the time index n , and the initial vectors $\vec{u}_j, 0 \leq j \leq s-1$.

Unconditional stability is a strong requirement for ImEx LMMs that comes with the following caveat: unconditional stability is a coupled property of *both* the set of ImEx coefficients (a_j, b_j, c_j) and the matrices (\mathbf{A}, \mathbf{B}) . Hence, the following hold:

- A given set of coefficients, (a_j, b_j, c_j) , may yield unconditional stability for some splittings (\mathbf{A}, \mathbf{B}) , and not others.
- If the splitting (\mathbf{A}, \mathbf{B}) arises from the spatial discretization of a PDE, then a given set of coefficients (a_j, b_j, c_j) may not yield unconditional stability for all model parameters.

If the matrices \mathbf{A} and \mathbf{B} commute and are diagonalizable, then the stability of (5) can be examined by using the spectra, $\sigma(\mathbf{A})$ and $\sigma(\mathbf{B})$. In this paper we *do not* assume that \mathbf{A} and \mathbf{B} commute. Hence we cannot rely on the existence of common eigenvectors and must develop a different approach to study the stability of (5), as follows:

- We introduce an unconditional stability region/diagram \mathcal{D} , which is computable in terms of the scheme coefficients (a_j, b_j, c_j) *only*.
- We introduce a region in the complex plane that generalizes the notion of spectrum and depends on the matrix splitting (\mathbf{A}, \mathbf{B}) *only*.

This approach opens a pathway to the design of splittings that are guaranteed to be stable for a fixed set of ImEx coefficients, or to the choosing of ImEx coefficients for which a given splitting (\mathbf{A}, \mathbf{B}) yields a stable scheme. In fact, in this paper we introduce a new class of ImEx coefficients that may be chosen to stabilize a given splitting. For these new schemes the coefficients yield diagrams that permit (arbitrarily) large regions of unconditional stability.

3. Stability for linear multistep methods. In this section we review the stability criteria for ImEx LMMs defined by (5). Following a standard procedure [26, Chapter III.4], one may recast the linear recursion relation (5) with matrix coefficients as a single vector recursion on an $s \times N$ vector:

$$(6) \quad \vec{V}^n = \mathbf{W}\vec{V}^{n-1}, \quad \text{where } \vec{V}^n := (\vec{u}_{n+s}, \vec{u}_{n+s-1}, \dots, \vec{u}_{n+1})^T \in \mathbb{R}^{sN}.$$

Here \mathbf{W} is a matrix with block structure:

$$(7) \quad \mathbf{W} = \begin{pmatrix} a_s - kc_s\mathbf{A} & 0 & 0 & \dots & 0 \\ 0 & \mathbf{I} & 0 & \dots & 0 \\ 0 & 0 & \mathbf{I} & \dots & 0 \\ \vdots & & & \ddots & 0 \\ 0 & 0 & 0 & \dots & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{C}_{s-1} & \mathbf{C}_{s-2} & \dots & \mathbf{C}_1 & \mathbf{C}_0 \\ \mathbf{I} & 0 & \dots & 0 & 0 \\ 0 & \mathbf{I} & \dots & 0 & 0 \\ \vdots & & \ddots & 0 & 0 \\ 0 & 0 & \dots & \mathbf{I} & 0 \end{pmatrix},$$

where \mathbf{I} is the $N \times N$ identity matrix, and

$$\mathbf{C}_j = kc_j\mathbf{A} + kb_j\mathbf{B} - a_j\mathbf{I}, \quad 0 \leq j \leq s - 1.$$

Recall [26, Chapter III.4], [27, Chapter V.1] that (6), and hence the scheme (5), is stable for a given k if every semisimple² eigenvalue of \mathbf{W} satisfies $|\zeta| \leq 1$, and every nonsemisimple eigenvalue satisfies $|\zeta| < 1$. In the case when \mathbf{A} and \mathbf{B} do not commute, the eigenvalues of \mathbf{W} depend on both (i) the matrices \mathbf{A} and \mathbf{B} , and (ii) the ImEx time stepping coefficients (a_j, b_j, c_j) . Hence the eigenvalues of \mathbf{W} do not provide a way to characterize unconditional stability in a way analogous to that for non-ImEx schemes: Some set depending on \mathbf{L} only (e.g., its spectrum) must be included within some set that is defined by the scheme coefficients only (the unconditional stability set). In what follows we devise a strategy to get around this problem, so that conditions that guarantee unconditional stability of ImEx schemes can be formulated in a language similar to the one for non-ImEx schemes, or for ImEx schemes with commutative splits (though the set depending on $\mathbf{L} = \mathbf{A} + \mathbf{B}$ is no longer a spectrum).

Let $\vec{V}^* \neq \vec{0}$ be an eigenvector of \mathbf{W} with eigenvalue ζ . Then, due to the structure of the bottom $(s - 1)$ matrix blocks in \mathbf{W} , $\vec{V}^* \in \mathbb{C}^{sN}$ has the form

$$(8) \quad \vec{V}^* = (\zeta^{s-1}\vec{v}, \zeta^{s-2}\vec{v}, \dots, \zeta\vec{v}, \vec{v})^T, \quad \text{where } \vec{v} \neq \vec{0}, \vec{v} \in \mathbb{C}^N.$$

²An eigenvalue ζ is semisimple if its algebraic multiplicity equals its geometric multiplicity.

The characteristic equation for \mathbf{W} can be rewritten in the form

$$\det(\mathbf{W} - \zeta \mathbf{I}) = 0 \iff \det\left(\frac{1}{k}a(\zeta) \mathbf{I} - c(\zeta) \mathbf{A} - b(\zeta) \mathbf{B}\right) = 0,$$

where

$$a(z) = \sum_{j=0}^s a_j z^j, \quad b(z) = \sum_{j=0}^{s-1} b_j z^j, \quad c(z) = \sum_{j=0}^s c_j z^j$$

are polynomials determined by the time stepping coefficients (a_j, b_j, c_j) , $0 \leq j \leq s$.

Hence if ζ is an eigenvalue of \mathbf{W} (with possible algebraic multiplicity greater than one), then there always exists at least one \vec{V}^* from (8) with \vec{v} satisfying

$$(9) \quad \mathbf{T}(\zeta)\vec{v} = 0, \quad \text{where } \mathbf{T}(z) := \left(\frac{1}{k}a(z) \mathbf{I} - c(z) \mathbf{A} - b(z) \mathbf{B}\right).$$

Note that one may also arrive at (9) by substituting the normal mode ansatz $\vec{u}_n = \zeta^n \vec{v}$ into the general linear ImEx time stepping scheme (5).

Clearly if $\mathbf{T}(z)$ is singular for $|z| < 1$, then any eigenvector \vec{V}^* of \mathbf{W} has every eigenvalue (regardless of algebraic multiplicity) $|\zeta| < 1$. Conditions on $\mathbf{T}(z)$ for the stability of (5) can then be stated as follows.

PROPOSITION 2. *If, for a fixed $k > 0$, the matrix $\mathbf{T}(z)$ is nonsingular for all $|z| \geq 1$, i.e., $\det \mathbf{T}(z) \neq 0$ for $|z| \geq 1$, then the scheme (5) is stable.*

Remark 1. Proposition 2 is not sharp as we have omitted the possibility for $\det \mathbf{T}(\zeta) = 0$ with $|\zeta| = 1$.

3.1. The stability region \mathcal{D} . The $N \times N$ matrix equation (9) still couples both the matrices (\mathbf{A}, \mathbf{B}) to the scheme coefficients (a_j, b_j, c_j) . To decouple the time stepping stability analysis (i.e., the time stepping coefficients) from the details of the ODE being solved (i.e., the matrices \mathbf{A} and \mathbf{B}), we multiply (9) by the positive definite matrix $(-\mathbf{A})^{p-1}$, where $p \in \mathbb{R} - p$ real is all that is needed for the analysis below to hold. In the examples in section 5, we will eventually focus on $p = 1$, as it is observed that this choice provides sufficient estimates for the test problems we consider. The stability theory obtained with other values of $p \neq 1$ may still, however, be of use in the numerical treatment of other PDEs, distinct from those in section 5. Thus,

$$\frac{1}{k}a(\zeta)(-\mathbf{A})^{p-1}\vec{v} = -c(\zeta)(-\mathbf{A})^p\vec{v} + b(\zeta)(-\mathbf{A})^{p-1}\mathbf{B}\vec{v}.$$

Dotting through with \vec{v} and setting

$$(10) \quad y = -k \frac{\langle \vec{v}, (-\mathbf{A})^p \vec{v} \rangle}{\langle \vec{v}, (-\mathbf{A})^{p-1} \vec{v} \rangle}, \quad \mu = \frac{\langle \vec{v}, (-\mathbf{A})^{p-1} \mathbf{B} \vec{v} \rangle}{\langle \vec{v}, (-\mathbf{A})^p \vec{v} \rangle},$$

we obtain the equation³

$$(11) \quad a(\zeta) = y c(\zeta) - y \mu b(\zeta).$$

Since $(-\mathbf{A})$ is positive definite, y may take any value $y < 0$ as k varies over the allowable values $k > 0$, with any $\vec{v} \neq 0$ fixed. The following definition is then justified by the result in Proposition 2.

³The polynomial (11) with $\mu = 0$ was used in convergence proofs in [3, 4, 5, 15]. A similar equation, $a(\zeta) = \lambda c(\zeta) + \mu b(\zeta)$, was obtained in [8] for commuting matrices \mathbf{A} and \mathbf{B} and studied as a model equation for stability in [20] to estimate explicit time step k restrictions. However, note that here we *do not* assume that \mathbf{A} and \mathbf{B} commute.

DEFINITION 3 (stability). *The polynomial equation (11) is stable, for a given $y < 0$ and $\mu \in \mathbb{C}$, if every solution satisfies $|\zeta| < 1$.*

DEFINITION 4 (unconditional stability region). *We define the region of unconditional stability \mathcal{D} as the values of μ so that (11) is stable for all $y \in \mathbb{R}_{<0} \cup \{-\infty\}$. Formally, define the following sets:*

$$\begin{aligned} \mathcal{D}_y &:= \{\mu \in \mathbb{C} : (11) \text{ is stable for a fixed } y \in \mathbb{R}_{<0}\}, \\ \mathcal{D}_{-\infty} &:= \{\mu \in \mathbb{C} : c(\zeta) - \mu b(\zeta) \text{ has stable roots}\}, \\ \mathcal{D} &= \bigcap_{y \in \mathbb{R}_{<0} \cup \{-\infty\}} \mathcal{D}_y. \end{aligned}$$

Note that \mathcal{D} depends only on the ImEx time stepping coefficients and not on the matrices \mathbf{A}, \mathbf{B} . Moreover, \mathcal{D} may be empty for some schemes.

3.2. Numerical range and sufficient condition for unconditional stability. The exact realizable values of μ defined by the expression in (10), for a given splitting (\mathbf{A}, \mathbf{B}) and time stepping coefficients, are determined through the normal modes \vec{v} . To find these values of μ , which form a discrete, finite set in the complex plane, one must solve the fully coupled eigenvalue problem given by (9). A better and simpler approach is to overestimate the region in the complex plane where the values of μ reside. Specifically, the values of μ belong to the complex set obtained by allowing \vec{v} to vary over all possible vectors. That is,

$$\mu \in W_p, \quad \text{where } W_p := \left\{ \langle \vec{v}, (-\mathbf{A})^{p-1} \mathbf{B} \vec{v} \rangle : \langle \vec{v}, (-\mathbf{A})^p \vec{v} \rangle = 1 \right\}.$$

Using a straightforward change of variables $\vec{v} = (-\mathbf{A})^{\frac{p}{2}} \vec{x}$, and the fact that \mathbf{A} is symmetric, the set W_p can be identified as

$$W_p = W \left((-\mathbf{A})^{\frac{p}{2}-1} \mathbf{B} (-\mathbf{A})^{-\frac{p}{2}} \right).$$

Here $W(\mathbf{X})$ denotes the *numerical range* (also known as the *field of values*) of a matrix $\mathbf{X} \in \mathbb{C}^{N \times N}$ and is defined by

$$(12) \quad W(\mathbf{X}) := \{ \langle \vec{x}, \mathbf{X} \vec{x} \rangle : \|\vec{x}\| = 1, \vec{x} \in \mathbb{C}^N \}.$$

See section SM1 in the supplementary material for a list of standard properties for $W(\mathbf{X})$. One then arrives at a sufficient condition for unconditional stability for (5).

THEOREM 5 (sufficient condition for unconditional stability). *Suppose that a matrix splitting (\mathbf{A}, \mathbf{B}) has sets W_p for $p \in \mathbb{R}$ and that the LMM time stepping coefficients (a_j, b_j, c_j) have an unconditional stability region \mathcal{D} . Then, if there exists a $p \in \mathbb{R}$ such that $W_p \subseteq \mathcal{D}$, the scheme (9) is unconditionally stable.*

Remark 2. Different values of p may modify the size of W_p in the complex plane. The sufficient condition for unconditional stability requires only one value of p to satisfy $W_p \subseteq \mathcal{D}$ (even if other values of p violate $W_p \subseteq \mathcal{D}$).

4. New ImEx coefficients.

4.1. Definition of the new ImEx coefficients. The property of unconditional stability is not limited to LMMs; however, here we focus on LMMs only. Any ImEx LMM where the number of steps equals the order of the scheme $s = r$ is completely

Downloaded 05/11/18 to 18.51.0.240. Redistribution subject to SIAM license or copyright; see http://www.siam.org/journals/ojsa.php

defined by specifying the polynomial $c(z)$. For instance, given $s = r$ and a fixed $c(z)$, the order conditions define the polynomials $a(z), b(z)$ and subsequently all time stepping coefficients. Therefore, the roots⁴ of the polynomial $c(z)$ can also be used to uniquely define any ImEx scheme when $r = s$. The new ImEx coefficients proposed in this paper will be prescribed by the location of the roots of $c(z)$. In particular, regions of unconditional stability \mathcal{D} depend strongly on the location of the roots of $c(z)$ and become large when the roots of $c(z)$ get close to 1 (see also section SM4). Although there are many options for parameterizing how the roots of $c(z)$ approach 1, we choose the simplest approach and *lock* all the roots together.

DEFINITION 6 (new ImEx coefficients). *For orders $1 \leq r \leq 5$ and $0 < \delta \leq 1$, the new ImEx coefficients (a_j, b_j, c_j) , for $0 \leq j \leq r$, are defined as the following polynomial coefficients:*

$$(13) \quad (\text{Implicit coeff.}) \quad c(z) = (z - 1 + \delta)^r,$$

$$(14) \quad (\text{Explicit coeff.}) \quad b(z) = (z - 1 + \delta)^r - (z - 1)^r.$$

The time stepping polynomial $a(z)$ is concisely written as the r th order Taylor polynomial centered at $z = 1$ of the generating function $f(z)$,

$$(15) \quad (\text{Derivative coeff.}) \quad a(z) = \sum_{j=1}^r \frac{f^{(j)}(1)}{j!} (z - 1)^j, \quad f(z) = (\ln z)(z - 1 + \delta)^r.$$

Note that once $c(z)$ is chosen, $a(z)$ and $b(z)$ are uniquely determined. For more on this, see Proposition 7 below. In section 7 we report the ImEx coefficients (a_j, b_j, c_j) as polynomial functions of δ . In the case when $\delta = 1$, the new coefficients recover the combined SBDF—backward differentiation formula (for the implicit $c(z)$) and Adams–Bashforth (for the explicit $b(z)$). For $\delta < 1$ the roots of $c(z)$ shift towards $z = 1$. The new coefficients bear some similarity to the one-parameter, high order, multistep schemes with large absolute stability regions studied in [29, 30]. We stress, however, that our use of the ImEx coefficients in Definition 6 is of a fundamentally different nature than the non-ImEx investigation found in [29, 30]. Specifically, we select a δ value that is strictly bounded away from 0, based on the ImEx splitting (\mathbf{A}, \mathbf{B}) of \mathbf{L} , which yields an unconditionally stable method. Moreover, a subsequent error investigation indicates that δ should be selected as large as possible, while still maintaining unconditional stability.

Remark 3. We limit Definition 6 to orders $r \leq 5$. SBDF schemes ($\delta = 1$) with orders $r \geq 7$ are not zero stable. Furthermore, the characterization of \mathcal{D} for $r = 6$ is not contained within the theory presented in the following subsection. Specifically, Numerical Observation 1 (see section 4.2) fails for $r = 6$ and $\delta = 1$.

PROPOSITION 7. *For all $0 < \delta \leq 1$ and orders $1 \leq r \leq 5$, the ImEx coefficients in Definition 6 are zero stable and satisfy the r th order conditions.*

See section SM2 for the verification of Proposition 7.

4.2. Stability regions for the new ImEx coefficients. The region \mathcal{D} was introduced in the context of the sufficient conditions for unconditional stability. As we will see later (in section 4.3) it also plays a role in the necessary conditions for

⁴Since rescaling the ImEx coefficients (a_j, b_j, c_j) by an overall constant does not modify a scheme, one can take without loss of generality the leading coefficient of $c(z)$ to be 1.

unconditional stability. In this section we characterize the geometry of \mathcal{D} for the ImEx coefficients in Definition 6. This geometry (i.e., the size and shape of \mathcal{D} in the complex plane) fixes classes of splittings (\mathbf{A}, \mathbf{B}) that are, or are not, unconditionally stable. Roughly speaking, for small δ values, \mathcal{D} approaches the union of (i) a large circle with radius $\sim (r\delta)^{-1}$ and center $\sim -(r\delta)^{-1}$, and (ii) a triangular region, symmetric relative to the real axis, with its tip on the positive real axis. See Figure 4.

We first focus on describing the set $\mathcal{D}_{-\infty}$, since by definition the unconditional stability region \mathcal{D} is a subset of $\mathcal{D}_{-\infty}$, i.e., $\mathcal{D} \subseteq \mathcal{D}_{-\infty}$. However, we show later that this subset inclusion is in fact an equality, so that $\mathcal{D} = \mathcal{D}_{-\infty}$. Thus one should keep in mind that statements characterizing $\mathcal{D}_{-\infty}$ are statements about \mathcal{D} . The main result regarding $\mathcal{D}_{-\infty}$ is summarized by the following theorem.

THEOREM 8 (the set $\mathcal{D}_{-\infty}$). *The set $\mathcal{D}_{-\infty}$ is simply connected, contains the origin $\mu = 0$, and has a boundary parameterized by the curve*

$$(16) \quad \partial\mathcal{D}_{-\infty} = \left\{ \frac{(z - 1 + \delta)^r}{(z - 1 + \delta)^r - (z - 1)^r} : |z| = 1, \arg z_0 \leq \arg z \leq 2\pi - \arg z_0 \right\},$$

where

$$(17) \quad \begin{aligned} z_0 &= 1 && \text{for order } r = 1, \text{ and} \\ z_0 &= \frac{2 - \delta - 2(1 - \delta) \cos(\pi/r) e^{i\pi/r}}{2 - \delta - 2 \cos(\pi/r) e^{i\pi/r}} && \text{for orders } 2 \leq r \leq 5. \end{aligned}$$

Moreover, let m_r (resp., m_l) be the rightmost (resp., leftmost) point of $\partial\mathcal{D}_{-\infty}$. Then m_r (resp., m_l) is obtained at the parameter value $z = z_0$ (resp., $z = -1$). Thus

$$\begin{aligned} \text{for } r = 1, \quad m_l &= \frac{-(2-\delta)}{\delta} && \text{and } m_r = 1, \\ \text{for } 2 \leq r \leq 5, \quad m_l &= \frac{-(2-\delta)^r}{2r-(2-\delta)^r} && \text{and } m_r = \frac{(2-\delta)^r}{(2-\delta)^r + 2^r \cos^r(\pi/r)}. \end{aligned}$$

Note that both m_l and m_r are on the real axis.

Proof. For $r = 1$ the proof is straightforward as $\partial\mathcal{D}_{-\infty}$ is a circle for all $0 < \delta \leq 1$. The idea for the proof when $2 \leq r \leq 5$ is to show that $\mathcal{D}_{-\infty} = \varphi^{-1}(\mathcal{T})$ is the preimage of a set \mathcal{T} (which is a triangle for $r \geq 3$ and a strip for $r = 2$) under the mapping of a complex function $\varphi(z)$. The results in the theorem then follow from basic calculus arguments, and the conformal properties of complex mappings.

The set $\mathcal{D}_{-\infty}$ consists of the values $\mu \in \mathbb{C}$, which ensure that the solutions $z \in \mathbb{C}$ to the following polynomial equation are stable (see Definition 3):

$$(18) \quad c(z) - \mu b(z) = 0 \iff (z - 1 + \delta)^r - \mu \left((z - 1 + \delta)^r - (z - 1)^r \right) = 0.$$

Note that $0 \in \mathcal{D}_{\infty}$, since $c(z)$ has a single root: $z = 1 - \delta$ (with multiplicity r). As a direct result of the simple structure of the polynomials $c(z)$ and $b(z)$, equation (18) can be solved explicitly to write the solutions $z_j(\mu)$ (for $0 \leq j \leq r - 1$) in terms of μ as

$$(19) \quad z_j(\mu) = 1 + \frac{\delta}{\xi_j \varphi(\mu) - 1}, \quad \text{where } \xi_j = e^{\frac{i2\pi j}{r}}, \quad 0 \leq j \leq r - 1.$$

Here $\varphi(\mu)$ is the complex-valued function defined using a branch cut taken along the negative real axis:

$$(20) \quad \varphi(\mu) := \left(\frac{\mu}{\mu - 1} \right)^{1/r}, \quad \text{where } (Re^{i\theta})^{1/r} := R^{1/r} e^{\frac{i\theta}{r}}, \quad (-\pi < \theta \leq \pi, R \geq 0).$$

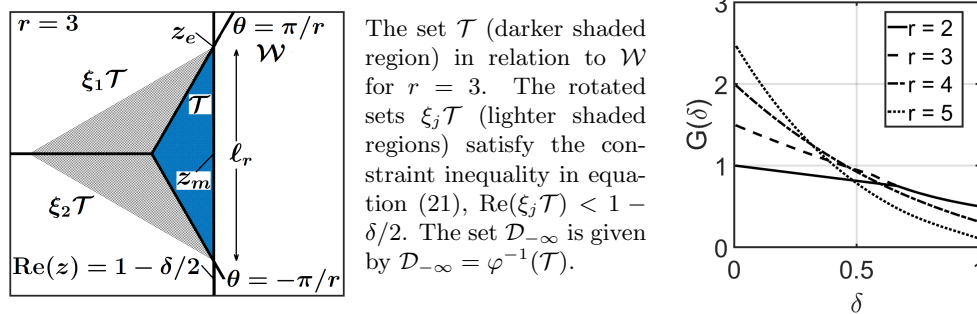


FIG. 1. Left: the set \mathcal{T} . Right: plot of $G(\delta)$, as defined by (24).

Observe that $\varphi(\mu)$ is the composition of a Möbius transformation (which has the property that it is a one-to-one mapping of the compactified complex plane to itself, with the identification that the point $1 \rightarrow \infty$ and $\infty \rightarrow 1$), with the r th root function. Hence, $\varphi(\mu) : \mathbb{C} \rightarrow \mathcal{W}$, where

$$\mathcal{W} = \left\{ z \in \mathbb{C} : z = 0, \text{ or } -\frac{\pi}{r} < \arg z \leq \frac{\pi}{r} \right\}.$$

Next, we note that the modulus constraints $|z_j| < 1$ restrict the range of $\varphi(\mu)$ to the intersection of r half-planes given by the following inequalities:

$$(21) \quad \left| 1 + \frac{\delta}{\xi_j \varphi(\mu) - 1} \right| < 1 \iff \operatorname{Re}(\xi_j \varphi(\mu)) < 1 - \frac{\delta}{2}.$$

Clearly, inequality (21) must be satisfied by all roots $0 \leq j \leq r-1$. Satisfying inequality (21) for $j=0$, however, will automatically guarantee the satisfaction of the remaining $1 \leq j \leq r-1$ inequalities. To make this correspondence precise, we introduce the set \mathcal{T} (which is a triangle for $r \geq 3$, a strip for $r=2$, and half-plane for $r=1$), obtained by taking the intersection of \mathcal{W} with the $j=0$ inequality in (21),

$$(22) \quad \mathcal{T} = \left\{ z \in \mathcal{W} : \operatorname{Re}(z) < 1 - \frac{\delta}{2} \right\}.$$

Figure 1 (left) shows the triangle \mathcal{T} , as well as the rotated triangles $\xi_j \mathcal{T}$, for $r=3$. A simple use of inequalities,⁵ whose geometric interpretation is highlighted in Figure 1 (left), shows that if $w \in \mathcal{T}$, then $\operatorname{Re}(\xi_j w) < 1 - \frac{\delta}{2}$. Hence, if $\varphi(\mu) \in \mathcal{T}$, then $\mu \in \mathcal{D}_{-\infty}$. That is, $\mathcal{D}_{-\infty} = \varphi^{-1}(\mathcal{T})$ is the preimage of \mathcal{T} under the mapping $\varphi(z)$. The sets $\mathcal{D}_{-\infty}$, for the parameter value $\delta=1$ and orders $1 \leq r \leq 3$, are shown in Figure 2.

The properties of $\mathcal{D}_{-\infty}$ now follow by observing that the set $\varphi^{-1}(\mathcal{T}) = M(\mathcal{T}^r)$ is the image under the Möbius transformation $M(z) = z/(z-1)$ of the set \mathcal{T}^r , where $\mathcal{T}^r = \{z^r : z \in \mathcal{T}\}$ is the r th power of \mathcal{T} . Below, we will use the following simple properties [2, Chapter 3] of the Möbius transformation $M(z)$ in the Riemann sphere, with the understanding that $M(1) = \infty$ and $M(\infty) = 1$:

(M1) The real axis is invariant under $M(z)$.

(M2) If D is a closed disk centered on the real axis, with $\operatorname{Re}(D) < 1$, then $M(D)$ is also a disk centered on the real axis with $\operatorname{Re}(M(D)) < 1$.

⁵Specifically, if $w = Re^{i\theta}$ with $R < (1 - \delta/2) \sec(\theta)$ so that $\operatorname{Re}(w) < 1 - \delta/2$, then $\operatorname{Re}(\xi_j w) = R \cos(\theta + 2\pi j/r) < (1 - \delta/2)$, since $\cos(\theta + 2\pi j/r) \leq \cos(\theta)$ for $|\theta| \geq \pi/r$.

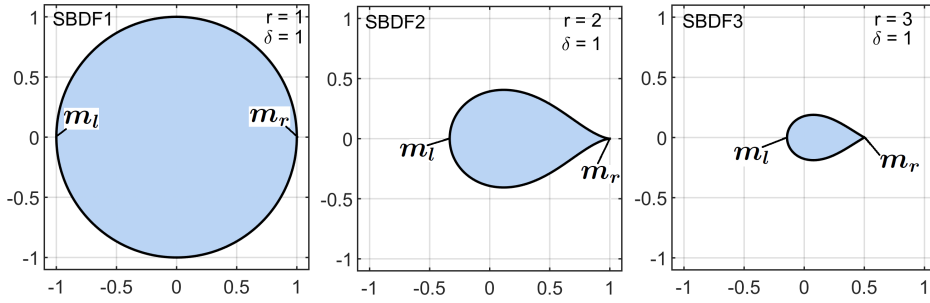


FIG. 2. The sets $\mathcal{D}_{-\infty}$ (which by virtue of Proposition 9 equal \mathcal{D}) are shown shaded. The parameters are $\delta = 1$ (SBDF schemes) and orders $r = 1, 2, 3$ (left to right). Formulas for the boundary are given by Theorem 8.

(M3) The half-plane $\text{Re}(z) \leq 1$ is invariant under $M(z)$. Any half-plane $\text{Re}(z) \leq \alpha < 1$ ($\alpha \in \mathbb{R}$) is mapped to a disk D with center on the real axis and $\text{Re}(D) < 1$.

(M4) M is a continuous map on the Riemann sphere, and $M = M^{-1}$.

Note that $\mathcal{D}_{-\infty} = \varphi^{-1}(\mathcal{T}) = M(\mathcal{T}^r)$ is simply connected, since M is continuous and \mathcal{T}^r is simply connected. To obtain the formula for the boundary $\partial\mathcal{D}_{-\infty}$, we observe that the line segments $\theta = \pm\pi/r$ on $\partial\mathcal{T}$ are mapped (under the r th power, $\mathcal{T} \rightarrow \mathcal{T}^r$) to identical line segments along the negative real axis. Further, these segments are contained in the interior of \mathcal{T}^r . Hence the boundary of \mathcal{T}^r , and subsequently the boundary $\partial\mathcal{D}_{-\infty} = \varphi^{-1}(\ell_r)$, is the preimage of the line or line segment which is the right side of \mathcal{T} . Here ℓ_r is defined as follows:

$$\text{For } r = 2, \quad \ell_2 = \left\{ \text{Re}(z) = 1 - \delta/2 \right\},$$

$$\text{For } r \geq 3, \quad \ell_r = \left\{ (1 - \tau)\bar{z}_e + \tau z_e : 0 < \tau \leq 1, z_e = (1 - \delta/2) \sec(\pi/r) e^{i\pi/r} \right\}.$$

Substituting $\varphi(\ell_r)$ into (19) for $j = 0$ yields the root locus parameterization of the boundary $\partial\mathcal{D}_{-\infty}$ stated in the theorem. The value z_0 in the theorem statement corresponds to substituting the endpoint \bar{z}_e of ℓ_r for $\mu = \varphi^{-1}(z_e)$ into the formula for $z_0(\mu)$ in (19),

$$z_0 = \frac{2 - \delta - 2(1 - \delta) \cos(\pi/r) e^{i\pi/r}}{2 - \delta - 2 \cos(\pi/r) e^{i\pi/r}} \quad \text{for } 2 \leq r \leq 5.$$

In the above expression, and for our subsequent calculations below, it is understood that for $r = 2$, z_e is taken as $z_e = (1 - \delta/2) + i\infty$.

Finally, to verify the result for the right- and leftmost endpoints of $\partial\mathcal{D}_{-\infty}$, our goal is to show that \mathcal{T}^r is contained in a suitably chosen disk ($r \geq 3$) or half-plane ($r = 2$) and to use properties (M1)–(M3). First denote the midpoint of ℓ_r as $z_m = (1 - \delta/2)$. Then the only values of $\partial\mathcal{T}^r$ along the real axis are z_m^r and z_e^r . Hence by property (M1), $m_l := \varphi^{-1}(z_m)$ and $m_r := \varphi^{-1}(z_e)$ are the only values of $\partial\mathcal{D}_{-\infty}$ along the real axis. To show that m_l and m_r are the leftmost and rightmost points of $\mathcal{D}_{-\infty}$ for $r = 2$, note that \mathcal{T}^r is contained within the half-plane $\text{Re}(z) \leq z_m^2$ and contains the point along the negative real axis $-\infty \in \mathcal{T}^r$. Hence, by property (M3), $m_r = 1$ is the rightmost point, and by combining properties (M1) and (M3), m_l is the leftmost point of $\partial\mathcal{D}_{-\infty}$. For $r \geq 3$, it is sufficient to show that \mathcal{T}^r is contained in

the disk $D = \{z \in \mathbb{C} : |z - z_d| \leq R_d\}$ centered at $z_d = \frac{1}{2}(z_e^r + z_m^r)$ with a radius $R_d = \frac{1}{2}(z_m^r - z_e^r)$, and right and left endpoints z_m^r and z_e^r , respectively. This is because properties (M1) and (M2) imply that $m_r = M(z_m^r)$ and $m_l = M(z_e^r)$ will be preserved as the right- and leftmost points of $\partial\mathcal{D}_{-\infty}$ under the transformation $M(z)$. To show $\mathcal{T}^r \subseteq D$, write the boundaries $\partial\mathcal{T}^r$ and ∂D in polar coordinates $r e^{i\theta}$, with $r = f(\theta)$ and $r = g(\theta)$, respectively. Then, with $\beta_r = \sec^r(\pi/r)$,

$$f(\theta) = (1 - \delta/2)^r \sec^r(\theta/r) \quad \text{and}$$

$$g(\theta) = (1 - \delta/2)^r \left(\frac{1}{2}(1 - \beta_r) \cos(\theta) + \sqrt{\beta_r + \left(\frac{1}{2}(1 - \beta_r) \cos(\theta) \right)^2} \right).$$

By symmetry across the real axis, it is sufficient to show that $f(\theta) \leq g(\theta)$ for $0 \leq \theta \leq \pi$. This is true (i.e., after manipulating $f(\theta) \leq g(\theta)$), provided that the following inequality holds for $0 \leq \theta \leq \pi$:

$$h_r(\theta) := \beta_r - \sec^{2r}(\theta/r) + \sec^r(\theta/r) \cos(\theta)(1 - \beta_r) \geq 0.$$

Expanding $\cos(\theta)$ in powers of $\cos(\theta/r)$ via the binomial series, a direct computation of $h_r(\theta)$ (on $0 \leq \theta \leq \pi$) yields

$$h_3(\theta) = (\sec^2(\theta/3) - 1)(4 - \sec^2(\theta/3))(5 + \sec^2(\theta/3)) \geq 0,$$

$$h_4(\theta) = (\sec^2(\theta/4) - 1)(2 - \sec^2(\theta/4))(10 + 3\sec^2(\theta/4) + \sec^4(\theta/4)) \geq 0.$$

For $h_5(\theta)$ we write

$$h_5(\theta) = (\sec^2(\theta/5) - 1) \tilde{h}_5(\sec^2(\theta/5)),$$

where $\tilde{h}_5(x) = -x^4 - x^3 - x^2 - (5\beta_5 - 4)x - 16 + 15\beta_5$.

We claim now that $\tilde{h}_5(x) \geq 0$ for $1 \leq x \leq \sec^2(\pi/5)$. For this, note that $\beta_5 > \sec^4(\pi/6) = 16/9$, which shows that $\tilde{h}_5(1) > 10(16/9 - 3/2) > 0$. By construction, we also know that the boundaries $\partial\mathcal{T}^5$ and D touch at $\theta = \pi$, which implies $f(\pi) = g(\pi)$. This can then be used to show that $\tilde{h}_5(\sec^2(\pi/5)) = 0$. Finally, applying Descartes' rule of signs to the derivative $\tilde{h}'_5(x)$ shows that $\tilde{h}'_5(x)$ has no roots for $x > 0$. Hence, $\tilde{h}_5(x)$ is decreasing, and thus $\tilde{h}_5(x) \geq 0$ on $1 \leq x \leq \sec^2(\pi/5)$. \square

Figure 2 illustrates Theorem 8 by plotting the sets $\mathcal{D}_{-\infty}$ for the well-known SBDF schemes. Using the characterization of $\mathcal{D}_{-\infty}$ in Theorem 8, we are now in a position to show not only that $\mathcal{D} \subseteq \mathcal{D}_{-\infty}$, but that this inclusion is also an equality: $\mathcal{D} = \mathcal{D}_{-\infty}$.

To first illustrate that $\mathcal{D} = \mathcal{D}_{-\infty}$, in Figure 3 we plot \mathcal{D}_y for different values of y , using the *boundary locus* [38, Chapter 7.6] method. Specifically, \mathcal{D}_y is a region whose boundary is a subset of the locus

$$(23) \quad \Gamma_y := \left\{ \frac{1}{b(z)}(c(z) - y^{-1}a(z)) : |z| = 1 \right\}, \quad \Gamma_{-\infty} := \left\{ \frac{c(z)}{b(z)} : |z| = 1 \right\}.$$

Equation (23) is obtained by isolating μ in (11) and letting z vary over the unit circle. Figure 3 shows the nested stability regions \mathcal{D}_y for orders $r = 3, 4, 5$ and fixed parameter value $\delta = 1$. In the figure, the solid curve traces out Γ_y corresponding to the boundary locus for \mathcal{D}_y . The dashed curves show as a reference Γ_y for different y values. Although the plots are only for one value of δ , the limiting behavior $\mathcal{D} = \mathcal{D}_{-\infty}$ is observed for all $0 < \delta \leq 1$.

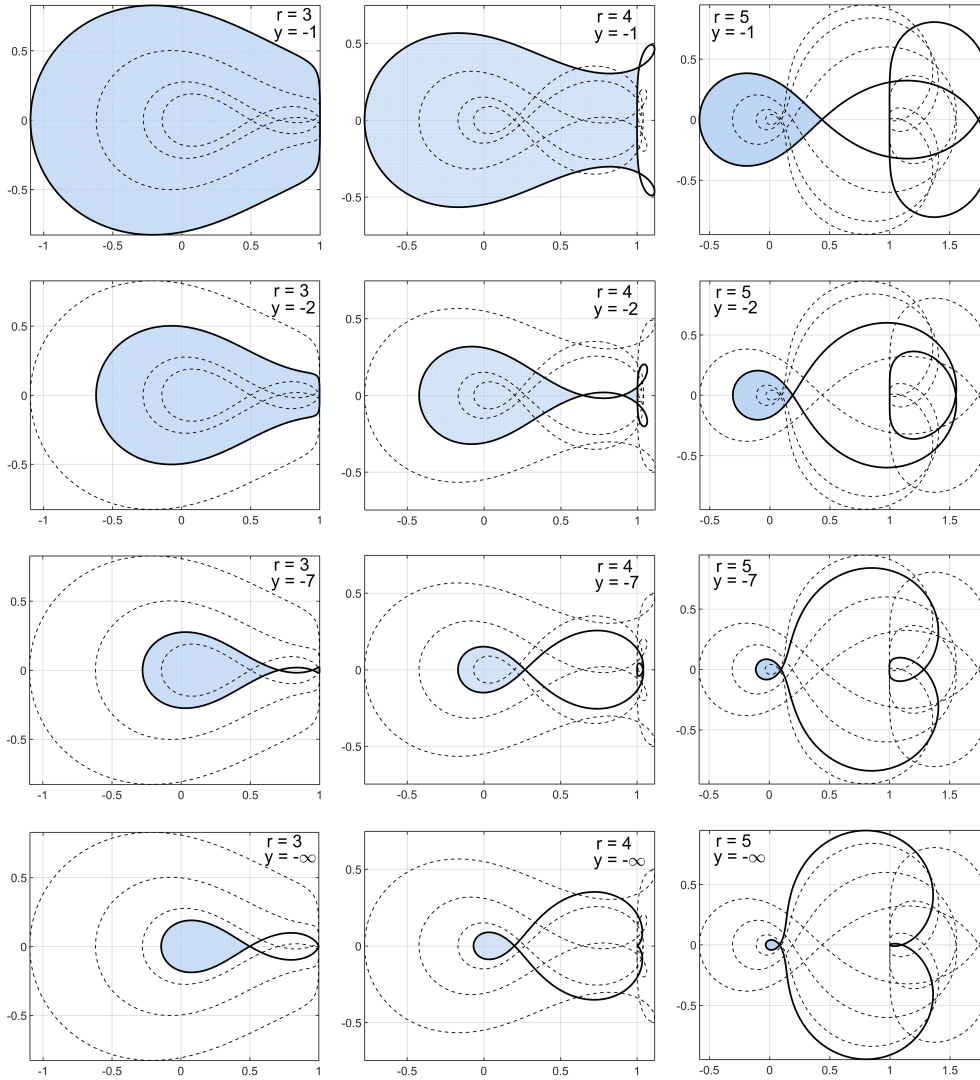


FIG. 3. Visualization of Proposition 9: $\mathcal{D}_{-\infty}$ is contained in \mathcal{D}_y for all $y < 0$. Plot of the boundary locus Γ_y (black curves) and the stability regions \mathcal{D}_y (shaded regions) for $y = -1, -2, -7, -\infty$ (top to bottom), orders $r = 3, 4, 5$ (left to right), and fixed parameter value $\delta = 1$. In each plot, the dashed lines $\Gamma_{-1}, \Gamma_{-2}, \Gamma_{-7}, \Gamma_{-\infty}$ are shown for reference. Note that the inclusion $\mathcal{D}_{-\infty} \subseteq \mathcal{D}_y$ is valid for all $y \in \mathbb{R}_{<0}$.

We now show that the set equality $\mathcal{D} = \mathcal{D}_{-\infty}$ is a direct consequence of the fact that the function $G(\delta)$ (defined below for the ImEx schemes in Definition 6) is positive. Note that $G(\delta)$, roughly speaking, is a measure of the distance of Γ_y to the set $\mathcal{D}_{-\infty}$ —and it is the key to showing that $\mathcal{D} = \mathcal{D}_{-\infty}$:

$$(24) \quad G(\delta) := \inf_{y < 0} \min_{w \in \Gamma_y} \left[\left(\operatorname{Re}(\varphi(w)) - (1 - \delta/2) \right) (1 - y)\delta^{-2} \right].$$

This function may be numerically computed, which leads to the following.

Numerical Observation 1. Numerical computations (shown in Figure 1, right) in-

dicates that for $0 < \delta \leq 1$ and $2 \leq r \leq 5$, $G(\delta) > 0$.

This fact is introduced below as an assumption in Proposition 9.

The positive factor $(1-y)\delta^{-2}$ in (24) is included to rescale the difference between $\operatorname{Re}(\varphi(w))$ and $(1-\delta/2)$, which vanishes as $y \rightarrow -\infty$ or $\delta \rightarrow 0$. This rescaling helps to visually verify that $G(\delta)$ does not change sign, even as $y \rightarrow -\infty$ or $\delta \rightarrow 0$. To computationally handle the infinite interval $-\infty < y < 0$, we introduce the change of variables $\tilde{y} = (1-y)^{-1}$, so that $0 < \tilde{y} < 1$. For each fixed value of \tilde{y} , we parameterize Γ_y as the image of the unit circle, which then allows us to compute $G(\delta)$ as a double minimization over two real variables on bounded intervals.

PROPOSITION 9 (the set $\mathcal{D} = \mathcal{D}_{-\infty}$). (i) For $r = 1$ and $0 < \delta \leq 1$, $\mathcal{D} = \mathcal{D}_{-\infty}$. (ii) For $2 \leq r \leq 5$, assume that $G(\delta) > 0$, $0 < \delta \leq 1$. Then

$$(25) \quad \mu \in \mathcal{D}_{-\infty} \implies \mu \in \mathcal{D}_y \quad \text{for any } y \in \mathbb{R}_{<0}.$$

In other words, for every $y \in \mathbb{R}_{<0}$ the set \mathcal{D}_y contains the limiting set $\mathcal{D}_{-\infty}$. As a result, the unconditional stability region is $\mathcal{D} = \mathcal{D}_{-\infty}$.

Proof. For (i), the proof is straightforward as \mathcal{D}_y is a disk centered at $1 - (\delta^{-1} - y^{-1})$ with radius $\delta^{-1} - y^{-1}$. For (ii) the proof involves two steps. First, we use a standard continuity argument to show that if $\mu \in \mathcal{D}_{-\infty}$, but $\mu \notin \mathcal{D}_{y_0}$ for some $y_0 < 0$, then there is an intermediate y -value ($-\infty < y < y_0$) where μ must lie on the boundary locus $\mu \in \Gamma_y$. Next we show that Γ_y is bounded away from $\mathcal{D}_{-\infty}$ when $y < 0$. It then follows that $\mu \in \mathcal{D}_y$ whenever $\mu \in \mathcal{D}_{-\infty}$.

To proceed with the first step, we define the following polynomial function based on (11):

$$(26) \quad P(z; \tilde{y}) := c(z) - \mu b(z) + \frac{\tilde{y}}{1-\tilde{y}} a(z).$$

Here $y = 1 - \tilde{y}^{-1}$, so that $0 < \tilde{y} < 1$ (resp., $\tilde{y} = 0$) corresponds to $y < 0$ (resp., $y = -\infty$), which will be useful in the subsequent continuity argument. To minimize additional notation, we will continue to use \mathcal{D}_y and Γ_y as sets, and \tilde{y} as the parameter in the polynomials, with the understanding that $y = 1 - \tilde{y}^{-1}$. Then \mathcal{D}_y is defined as $\mu \in \mathbb{C}$ such that $P(z; \tilde{y})$ has r roots inside the unit circle or, alternatively, (i) $P(z; \tilde{y}) \neq 0$ on the unit circle $|z| = 1$, and (ii) the function $F(\tilde{y}) = r$, where $F(\tilde{y})$ counts the number of roots $|z| < 1$ via the Cauchy integral formula:

$$F(\tilde{y}) := \frac{1}{2\pi i} \oint_{|z|=1} \frac{P_z(z; \tilde{y})}{P(z; \tilde{y})} dz.$$

Now, $F(\tilde{y})$ is continuous as a function of \tilde{y} , and also a constant, as long as it is defined. The only way $F(\tilde{y})$ may change values is if $P(z; \tilde{y}) = 0$ vanishes for some $|z| = 1$ on the unit circle, which implies $\mu \in \Gamma_y$. Hence, if for a given μ , $F(0) = r$ and $F(\tilde{y}_0) \neq r$, then there must exist a point $0 < \tilde{y} < \tilde{y}_0$ such that $\mu \in \Gamma_y$.

To show that Γ_y does not intersect $\mathcal{D}_{-\infty}$ for $0 < \tilde{y} < 1$, we exploit the fact that the mapping $\varphi(z)$, defined in Theorem 8, simplifies the shape of $\varphi(\mathcal{D}_{-\infty}) = \mathcal{T}$. Specifically, $\varphi(z)$ is a one-to-one mapping of \mathbb{C} to the wedge \mathcal{W} , so that it is sufficient to show that the mappings of $\mathcal{D}_{-\infty}$ and Γ_y under $\varphi(z)$ do not intersect, i.e., $\varphi(\Gamma_y)$ does not intersect \mathcal{T} , for $0 < \tilde{y} < 1$. Since \mathcal{T} is contained within the half-plane $\operatorname{Re}(z) < 1 - \delta/2$, we arrive at the following observation: if

$$(27) \quad \operatorname{Re}(\varphi(w)) - (1 - \delta/2) > 0 \quad \text{for all } 0 < \tilde{y} < 1, w \in \Gamma_y,$$

then $\varphi(\Gamma_y)$ and \mathcal{T} do not intersect. Multiplying the left-hand side of the inequality (27) by the positive factor $\delta^{-2}\tilde{y}^{-1} = \delta^{-2}(1 - y) > 0$, and minimizing over $0 < \tilde{y} < 1$, $w \in \Gamma_y$, yields the function $G(\delta)$. Hence, we arrive at the conclusion that $\varphi(\Gamma_y)$ and \mathcal{T} do not intersect whenever $G(\delta) > 0$, which, together with the first step of the proof, implies $D_{-\infty} \subseteq \mathcal{D}_y$ for all $y < 0$. \square

With the exact boundary locus description in Theorem 8 and the subsequent result that $\mathcal{D} = \mathcal{D}_{-\infty}$, one may provide an asymptotic description of \mathcal{D} in the limit $\delta \ll 1$.

Remark 4 (asymptotic \mathcal{D}). Define the circle C as

$$C = \left\{ z \in \mathbb{C} : \left| z + \frac{1}{r\delta} - \frac{r+1}{2r} \right| \leq \frac{1}{r\delta} \right\}.$$

Taking the asymptotic limit $\delta \ll 1$ and values of $|z - 1| \gg \delta$ in formula (16) for $\partial\mathcal{D}$ (which correspond to points in \mathcal{D} away from the rightmost values along the real axis), the exact boundary $\partial\mathcal{D}$ approaches the circle ∂C : $\mathcal{D} \approx C + \mathcal{O}(\delta)$. For $r = 1$, the domain $\mathcal{D} = C$ is a circle for all $0 < \delta \leq 1$.

The circle C in Remark 4 is obtained via an asymptotic computation, i.e., $\delta \rightarrow 0$, of (16). Specifically, note that the starting value of the locus description for $\mathcal{D}_{-\infty}$ in Theorem 8 satisfies $|z_0 - 1| = \mathcal{O}(\delta)$, so that the locus parameter z almost traces through an entire circle. Consider points $|z - 1| \gg \delta$ and expand $c(z)/b(z)$ in a Laurent series in powers of δ about $z = 1$:

$$\begin{aligned} \frac{c(z)}{b(z)} &= \frac{(z - 1)^r + \delta r(z - 1)^{r-1} + \dots}{\delta r(z - 1)^{r-1} + \delta^2 \frac{r(r-1)}{2}(z - 1)^{r-2} + \dots} \\ (28) \quad &= \frac{1}{\delta r} \left(\frac{(z - 1) + \delta r + \mathcal{O}(\delta^2)}{1 + \delta \frac{(r-1)}{2}(z - 1)^{-1} + \mathcal{O}(\delta^2)} \right) = \frac{1}{r\delta}(z - 1) + \frac{r+1}{2r} + \mathcal{O}(\delta). \end{aligned}$$

For values $|z| = 1$, equation (28) describes the boundary of the circle C defined in Remark 4 with radius $\frac{1}{r\delta}$ and center $\frac{r+1}{2r} - \frac{1}{r\delta}$. Hence $\partial\mathcal{D} \approx \frac{1}{r\delta}(z - 1) + \frac{r+1}{2r} + \mathcal{O}(\delta)$ for $|z - 1| \gg \mathcal{O}(\delta)$. Figure 4 shows the regions \mathcal{D} for different parameter values δ and orders $2 \leq r \leq 5$. In particular, the figure illustrates how the regions \mathcal{D} grow larger with decreasing δ values, and also approach the asymptotic circle C .

Having precise estimates for the geometric properties of \mathcal{D} , such as the formulas for m_r , m_l , and C , is very useful for the design of unconditionally stable schemes. Specifically the design of an unconditionally stable scheme requires a simultaneous choice of matrix splitting (\mathbf{A}, \mathbf{B}) and time stepping coefficients (a_j, b_j, c_j) . If one knows, either through direct numerical computation or analytic estimates, W_p for a matrix splitting (\mathbf{A}, \mathbf{B}) , then the estimates for m_r , m_l , and C can be used to choose a δ value large enough to guarantee that $W_p \subseteq \mathcal{D}$. Such a choice of δ will then provide the suitable time stepping coefficients that guarantee unconditional stability. We highlight such an approach in several numerical examples in section 5, as well as in greater detail in a companion paper on the practical aspects of unconditional stability for multistep ImEx schemes.

4.3. Necessary conditions for unconditional stability. The sufficient conditions for unconditional stability $W_p \subseteq \mathcal{D}$ are not sharp, and we supplement them with additional necessary conditions. Let

$$\sigma((-\mathbf{A})^{-1}\mathbf{B}) = \{ \mu \in \mathbb{C} : \mu(-\mathbf{A})\vec{u} = \mathbf{B}\vec{u}, \vec{u} \neq \vec{0} \}$$

be the generalized eigenvalues of $(-\mathbf{A}), \mathbf{B}$.

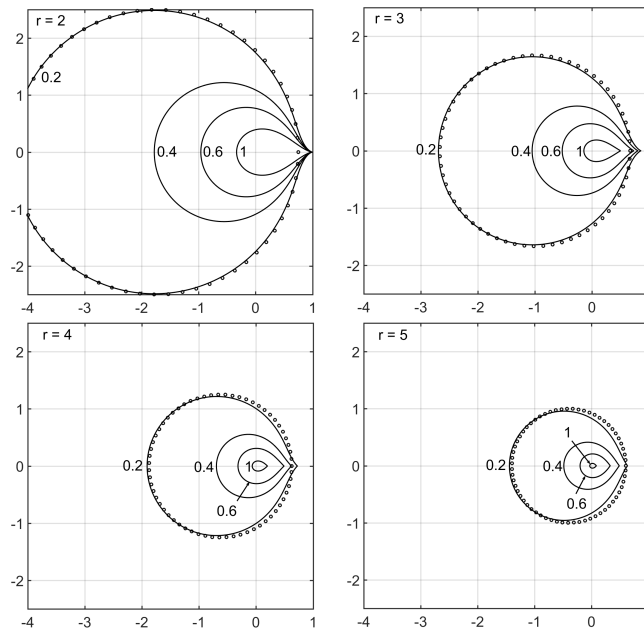


FIG. 4. Region of unconditional stability \mathcal{D} for orders $r = 2, 3, 4, 5$. In each subfigure, the boundary $\partial\mathcal{D}$ (solid line) is shown for parameter values $\delta = 0.2, 0.4, 0.6, 1$ ($\delta = 1$ corresponds to SBDF). For small $\delta \ll 1$, the stability region becomes arbitrarily large. With the exception of points near the positive real axis, it approaches the asymptotic circle C defined in Theorem 8. The dots (\circ) show C for $\delta = 0.2$.

PROPOSITION 10 (necessary condition for unconditional stability). *Given a set of ImEx time stepping coefficients (Definition 6) and the corresponding stability diagram \mathcal{D} , a necessary condition for unconditional stability of the scheme in (5) is that the eigenvalues satisfy $\sigma((-\mathbf{A})^{-1}\mathbf{B}) \subseteq \mathcal{D} \cup \Gamma_{-\infty}$.*

Proof. The idea behind the necessary condition is that in the limit of large time steps $k \rightarrow \infty$, the nonlinear eigenvalue problem (9) governing stability can be solved using the eigenvectors of the matrix $(-\mathbf{A})^{-1}\mathbf{B}$. As a result, a necessary condition for unconditional stability may be placed on the eigenvalue spectrum $\mu \in \sigma((-\mathbf{A})^{-1}\mathbf{B})$.

We first prove a slightly stronger statement. Let

$$\mathcal{A} := \{\mu \in \mathbb{C} : (18) \text{ has a solution } |z| > 1\}.$$

Then $\sigma((-\mathbf{A})^{-1}\mathbf{B}) \subseteq \mathcal{A}^c$ is a necessary condition for unconditionally stability. This is because, in the limit $k \rightarrow \infty$, the nonlinear eigenvalue problem (9) becomes

$$(29) \quad \mathbf{T}(z)\vec{u} = -c(z)\mathbf{A}\vec{u} - b(z)\mathbf{B}\vec{u} = 0.$$

Hence, an eigenvector \vec{u}_μ to $(-\mathbf{A})^{-1}\mathbf{B}$ with eigenvalue $\mu \in \sigma((-\mathbf{A})^{-1}\mathbf{B})$ becomes an eigenvector of (29):

$$(30) \quad \mathbf{T}(z)\vec{u}_\mu = -(c(z) - \mu b(z))\mathbf{A}\vec{u}_\mu = 0.$$

Thus the eigenvalues z satisfy (18), since $\mathbf{A}\vec{u}_\mu \neq 0$ because \mathbf{A} is invertible. If μ is also in \mathcal{A} , then at least one solution to (30) satisfies $|z| > 1$. Finally, we note that any nonlinear eigenvalue $|z| > 1$, arising in the limit $k \rightarrow \infty$, will yield a slightly

perturbed eigenvalue when $0 < k^{-1} \ll 1$. Thus, for any k sufficiently large (but finite) an unstable eigenvalue satisfying $|z| > 1$ will exist.

Finally we observe that $\mathcal{A}^c \subseteq \mathcal{D} \cup \Gamma_{-\infty}$. The reason is that every $\mu \in \mathcal{A}^c$ has one of the following properties: (i) all solutions to (18) have $|z| < 1$, implying $\mu \in \mathcal{D}$, or (ii) at least one solution to (18) has $|z| = 1$ (while all the others have $|z| < 1$), implying $\mu \in \Gamma_{-\infty}$. \square

Remark 5. Numerical experiments (such as the diagrams in Figure 3) suggest that the set $\mathcal{D} \cup \Gamma_{-\infty}$ in Proposition 10 can be further reduced to include only the portion of $\Gamma_{-\infty}$ that is the boundary $\partial\mathcal{D}$ and the single point $\{1\}$.

Remark 6. In the limit $\delta \rightarrow 0$, \mathcal{D} approaches the circle C which encompasses an entire complex half-plane:

$$(31) \quad \left\{ \mu \in \mathbb{C} : \operatorname{Re}(\mu) < \frac{r+1}{2r}, 1 \leq r \leq 5 \right\} \subseteq \lim_{\delta \rightarrow 0} \mathcal{D}.$$

The limiting \mathcal{D} also contains the real half-line $(-\infty, (1 + \cos^r(\pi/r))^{-1})$ for $2 \leq r \leq 5$.

4.4. Numerical error dependence on δ for the new ImEx coefficients.

Up to now, the results appear to indicate that one should choose $\delta \ll 1$ (extremely small) to yield a large unconditional stability region. In this section we describe why this is not a good strategy. In particular, we investigate the dependence of the *global truncation error* (GTE) on δ for the new ImEx coefficients. We do so by running numerical tests and computing the *error constants* which characterize the leading order asymptotic GTE behavior in k .

The GTE at time $t_n = nk$ is defined by $\|\vec{u}_n - \vec{u}^*(nk)\|_{\ell^\infty}$ and depends on \mathbf{L} , the time stepping coefficients, and the forcing $\vec{f}(t)$. Here $\vec{u}^*(t)$ is the exact ODE solution to (1) at time t . Formally, the new ImEx schemes given in Definition 6 achieve r th order accuracy, so that the $\text{GTE} = \mathcal{O}(k^r)$. The leading order constant in the GTE depends on $\mathbf{A}, \mathbf{B}, \vec{f}$ and the time stepping coefficients (for error constants in an LMM, see [26, equation (2.3), p. 373]). In ImEx schemes one may examine two separate error constants, an implicit $C_{I,r}$ (resp., explicit $C_{E,r}$) constant characterizing the error of a purely implicit (resp., explicit) scheme where $\mathbf{B} = 0$ (resp., $\mathbf{A} = 0$):

$$(32) \quad C_{I,r} := \frac{R_{I,r}}{c(1)} = \delta^{-r} R_{I,r}, \quad C_{E,r} := \frac{R_{E,r}}{b(1)} = \delta^{-r} R_{E,r}.$$

Here we have used the fact that $c(1) = b(1) = \delta^r$ for the new ImEx schemes, while the constants $R_{I,r}, R_{E,r}$ quantify how much the r th order coefficients (when $r = s$) fail to satisfy the $(r + 1)$ th order conditions (SM3):

$$R_{I,r} = \frac{1}{(r+1)!} \sum_{j=0}^r (a_j j^{r+1} - (r+1)c_j j^r), \quad R_{E,r} = \frac{1}{(r+1)!} \sum_{j=0}^r (a_j j^{r+1} - (r+1)b_j j^r).$$

Even though $R_{I,r}, R_{E,r}$ depend on δ , both constants satisfy $R_{I,r} = \mathcal{O}(1), R_{E,r} = \mathcal{O}(1)$ for all values of $0 < \delta \leq 1$. As a result, the asymptotic $\delta \ll 1$ behavior on the GTE for the new ImEx coefficients is $\text{GTE} = \mathcal{O}(\delta^{-r} k^r)$. The numerical tests in section 5, as well as those in section SM3, confirm the estimate $\text{GTE} \sim \delta^{-r} k^r$. As a result of this scaling, we adopt the following general philosophy: given a splitting (\mathbf{A}, \mathbf{B}) , choose δ as large as possible while maintaining unconditional stability.

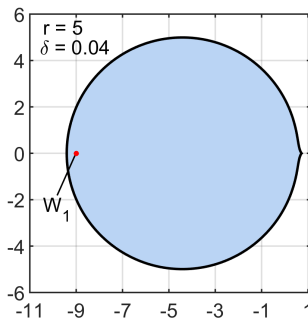


FIG. 5. ODE example (33). Note that W_1 (\circ) is contained in \mathcal{D} (shaded region) for the parameter value $\delta = 0.04$ and order $r = 5$. Using the new ImEx coefficients with $\delta = 0.04$ yields an unconditionally stable scheme.

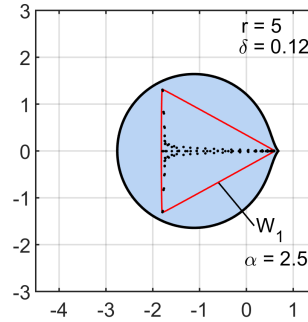


FIG. 6. PDE example for variable diffusion coefficient problem. The stability diagram \mathcal{D} (shaded region) is for the parameter $\delta = 0.12$ and order $r = 5$, and contains W_1 (red curve shows the boundary). The black dots show the generalized eigenvalues $\sigma((-\mathbf{A})^{-1}\mathbf{B})$.

5. Two illustrative examples. In this section we highlight the potential of the new ImEx coefficients to obtain unconditionally stable schemes. The first example (section 5.1) illustrates that a small implicit term can stabilize a larger explicit term. The second example (section 5.2) represents the numerical discretization of a variable coefficient diffusion equation. For this stiff problem, unconditional stability for orders $r > 2$ is beyond the capabilities of classical SBDF schemes; however, the new coefficients achieve the goal.

5.1. A single variable ODE. Consider the ODE

$$(33) \quad u_t = -10u = -u - 9u,$$

with splitting $\mathbf{A}u := -u$ and $\mathbf{B}u := -9u$. For this simple case, (\mathbf{A}, \mathbf{B}) are numbers, or 1×1 matrices. An important observation is that $|\mathbf{A}| = 1$, while $|\mathbf{B}| = 9$; i.e., the implicit term is 9 times smaller than the explicit term.

The set $W_1 = \{-9\}$ consists of one element, and it is also equal to the generalized eigenvalue $\sigma((-\mathbf{A})^{-1}\mathbf{B}) = \{-9\}$. Therefore unconditional stability requires that $\{-9\} \subseteq \mathcal{D}$. Using the fact that the leftmost endpoint of \mathcal{D} is given by m_l in the formula in Remark 8, one obtains unconditional stability for an r th order scheme, provided that

$$\frac{-(2-\delta)^r}{2^r - (2-\delta)^r} < -9 \iff \delta < 2 \left[1 - \left(\frac{9}{10} \right)^{1/r} \right].$$

Note that for a fixed δ value, the unconditional stability regions \mathcal{D} become smaller with increasing r . Setting $r = 5$ inside the inequality yields $\delta < 0.0417$. Therefore, a choice of the parameter value $\delta = 0.04$ inside the new ImEx coefficients guarantees that $W_1 \subseteq \mathcal{D}$ for $r = 5$ (as seen in Figure 5), and hence subsequently for all $1 \leq r \leq 5$. Hence, the smaller implicit term stabilizes the instabilities generated by the explicit term, thus achieving unconditional stability.

5.2. A PDE example: Variable coefficient diffusion. This example demonstrates how one might use the new ImEx coefficients, in conjunction with the sufficient conditions for unconditional stability, to avoid a stiff time step restriction in the spatial discretization of a PDE. Specifically, we numerically solve the variable coefficient

diffusion equation on the domain $\Omega = (-1, 1)$:

$$u_t = (d(x)u_x)_x + f(x, t) \quad \text{on } \Omega \times (0, T],$$

with Dirichlet boundary conditions, $u = 0$, on $x \in \{-1, 1\}$. Here $d(x) > 0$ is a spatially dependent diffusion coefficient.

For the spatial discretization, we adopt a Chebyshev spectral method [46, Chapters 5–7] using the $N + 2$ Chebyshev collocation points⁶

$$x_j = \cos\left(\frac{j\pi}{N+1}\right), \quad 0 \leq j \leq N+1, \quad \text{with} \quad \vec{u} = (u(x_1), \dots, u(x_N))^T \in \mathbb{R}^N.$$

We also use the boundary conditions to set $u(x_0) = u(x_{N+1}) = 0$ so that there are only N independent variables. Let \mathbf{D}_N be the spectral differentiation matrix, so that $\mathbf{D}_N \vec{u} \approx u_x(x)$. The matrix \mathbf{L} is then built using the Dirichlet boundary conditions by constructing

$$\mathbf{L} = \mathbf{D}_N \text{diag}(d(x_j)) \mathbf{D}_N.$$

Here \mathbf{L} acts on \vec{u} at the N grid points x_1, x_2, \dots, x_N (see [46, Chapter 7, p. 62] for details). Note that due to collocation of the boundary conditions, the matrix \mathbf{L} as well as the Laplacian $(\mathbf{D}_N)^2$ are not symmetric. However, the spectrum of \mathbf{L} and \mathbf{D}_N^2 are still purely real, in contrast with the situation in truly asymmetric problems, such as advection-diffusion.

In practice, the semi-implicit time stepping of (5), using the schemes defined by Definition 6, requires both a choice of splitting (\mathbf{A}, \mathbf{B}) and a set of new ImEx coefficients fixed by a choice of δ . For this example we consider a splitting where \mathbf{A} is a scalar multiple of the symmetrized part of the discrete, spectral Laplacian:

$$\mathbf{A} = \frac{\alpha}{2} \left((\mathbf{D}_N)^2 + (\mathbf{D}_N^2)^T \right), \quad \mathbf{B} = \mathbf{L} - \mathbf{A}$$

with an $\alpha > 0$. Here the choice of \mathbf{A} is negative definite and symmetric.

It is worth noting that, in general, \mathbf{A} and \mathbf{B} do not commute, therefore motivating the use of the new unconditional stability criteria. For this class of splittings, we focus on using the generalized numerical range W_1 . The reason is that the size and shape of W_1 depends only very weakly on N for large N .

There are now two free variables to choose: (i) α , which fixes the relative splitting of the (symmetric) implicit Laplacian to the explicit variable diffusion, and (ii) δ , which fixes the ImEx coefficients. Ideally, one would like to simultaneously choose α and δ to obtain unconditional stability and also minimize the overall error in the scheme. For a detailed discussion on how one may minimize the error, we defer to a companion paper on practical aspects of unconditional stability. Here we state briefly how one may first choose α , followed by δ , to obtain unconditional stability.

Decreasing α moves the set W_1 left in the complex plane—into a region that may be stabilized by the new ImEx coefficients. Specifically, we choose α small enough so that the rightmost point of W_1 is pushed to the left of the rightmost point of the limiting set \mathcal{D} (see Remark 6 for the rightmost point of \mathcal{D}). Once W_1 is sufficiently far left, we choose a sufficiently small δ value to ensure that $W_1 \subseteq \mathcal{D}$. To compute W_1 , we first build the matrix $\mathbf{X} = (-\mathbf{A})^{-\frac{1}{2}} \mathbf{B} (-\mathbf{A})^{-\frac{1}{2}}$, followed by using the MATLAB

⁶Note that here the points $1 = x_0 > x_1 > \dots > x_{N+1} = -1$ are in *reverse* order, following the usage in [46].

TABLE 1

Errors for variable coefficient diffusion test case $\alpha = 2.5$, $\delta = 0.12$, $t_f = 1$, $N = 100$. Exact solution $u^* = \sin(20t) \sin(2\pi x) e^{\sin(2\pi x)}$. Note that an explicit scheme, such as explicit Euler, would require a time step restriction $\mathcal{O}(N^{-2}) \sim 10^{-4} \sim 2^{-13}$. Here unconditional stability allows one to choose a time step based solely on accuracy considerations.

k	$r = 1$		$r = 2$		$r = 3$		$r = 4$		$r = 5$	
	Error	Rate	Error	Rate	Error	Rate	Error	Rate	Error	Rate
2^{-6}	2.1e+00	0.4	1.4e+00	0.5	1.0e+00	1.8	1.9e+00	4.4	4.0e+00	5.9
2^{-7}	1.3e+00	0.7	7.6e-01	0.9	4.4e-01	1.2	4.2e-01	2.2	6.8e-01	2.6
2^{-8}	7.0e-01	0.9	1.8e-01	2.1	2.4e-01	0.9	1.5e-01	1.5	1.9e-02	5.2
2^{-9}	3.6e-01	1.0	7.3e-02	1.3	5.1e-02	2.2	3.8e-03	5.3	4.8e-03	2.0
2^{-10}	1.8e-01	1.0	3.0e-02	1.3	5.8e-03	3.1	5.5e-04	2.8	1.8e-04	4.7
2^{-11}	8.2e-02	1.1	8.8e-03	1.8	6.0e-04	3.3	5.4e-05	3.4	4.7e-06	5.3
2^{-12}	3.9e-02	1.1	2.3e-03	1.9	6.7e-05	3.2	3.9e-06	3.8	1.2e-07	5.3
2^{-13}	1.9e-02	1.0	6.0e-04	2.0	7.9e-06	3.1	2.6e-07	3.9	3.7e-09	5.0

Chebfun routine [17] to compute $W_1 = W(\mathbf{X})$, based on a classical algorithm due to Johnson [31].

Finally, we perform a convergence test using the variable diffusion coefficient

$$d(x) = 4 + 3 \cos(2\pi x).$$

Figure 6 shows the set W_1 for a variable coefficient $d(x)$ and a value of $\alpha = 2.5$. In addition, the figure also shows a plot of the enclosing stability region \mathcal{D} for order $r = 5$ and the parameter value $\delta = 0.12$. Note that the unconditional stability region \mathcal{D} becomes smaller as the order r increases, so that $\delta = 0.12$ automatically guarantees unconditional stability for all orders $1 \leq r \leq 5$. For a convergence test, we use a manufactured solution approach and prescribe a forcing function $f(x, t)$ to yield an exact solution:

$$u^*(x, t) = \sin(20t) \sin(2\pi x) e^{\sin(2\pi x)}.$$

The numerical test case is also chosen to satisfy the *exact* initial data: $\vec{u}_j = \vec{u}^*(x, jk)$ evaluated at the grid points, for $j = 0, -1, \dots, -r + 1$. Table 1 shows the absolute $L^\infty(\Omega)$ errors for an integration time $t_f = 1$ and grid $N = 100$. Convergence rates for $1 \leq r \leq 5$ are observed as expected. Computations are done using MATLAB with double precision floating point arithmetic. Errors are limited to 10^{-9} for $r = 5$ due to machine precision and round-off errors.

Remark 7. An important observation is that the set W_1 remains bounded as $N \rightarrow \infty$. This result is of great practical relevance: one fixed value of δ can yield a stability region that contains W_1 for arbitrary N . For instance, the convergence results in Table 1 are all computed using the same value of δ . Therefore, the new time stepping schemes can be advantageous in PDE applications where the parameter δ can be chosen for a particular splitting of the differential operators, and can hold uniformly for any level of discretization of those operators (i.e., for a whole family of matrix splittings).

This example can be seen as a blueprint for many practical applications: the implicit part is simple and efficient to solve for (symmetric, constant coefficient), and the new ImEx coefficients enable one to obtain a numerical approximation that is unconditionally stable, thus avoiding diffusive-type time step restriction associated with explicit methods.

6. Discussion and conclusions. We have introduced a stability region \mathcal{D} , along with a generalized numerical range, as a way to guarantee unconditional stability for ImEx LMMs with a negative definite implicit term. It should be stressed that this type of study of unconditional stability is, structurally, not limited to ImEx LMMs and can also be examined in the context of any other time stepping scheme, such as Runge–Kutta methods, exponential integrators, deferred correction, or Richardson extrapolation. Moreover, unconditional stability (and further generalizations of \mathcal{D}) can in principle be examined also when the implicit term is not symmetric negative definite, such as for stiff wave problems.

In addition to sufficient criteria for unconditional stability, we have also introduced a family of ImEx LMM coefficients, parameterized by $0 < \delta \leq 1$ (which reduce to classical SBDF when $\delta = 1$). This parameter δ incurs crucial implications for stability, and the examples in section 5 highlight how the new ImEx coefficients can yield highly efficient time stepping schemes.

In light of these substantial advantages, three points of caution have to be stressed:

- (a) The error constant for an r th order method scales as δ^{-r} .
- (b) Computations with $\delta \ll 1$ may substantially amplify round-off errors.
- (c) L-stability, or small growth factors, are desirable properties for stiff equations, and lost for $\delta < 1$. If one uses the new ImEx coefficients as a fully implicit scheme (i.e., choosing $\mathbf{A} := \mathbf{L}$, $\mathbf{B} = 0$), then stability of the test equation $u_t = \lambda u$ is characterized by roots of the polynomial $a(z) - k\lambda c(z) = 0$. In the limit $k \rightarrow \infty$, the roots approach $\zeta := 1 - \delta$ (repeated r times). L-stability is only attained when the roots ζ have $\delta = 1$, corresponding to SBDF. Moreover, if $\delta \ll 1$, then the growth factor $1 - \delta$ is close to 1, implying that stiff modes may require many time steps to decay.

To conclude, major drawbacks of the new ImEx schemes are incurred only if $\delta \ll 1$. In practice, a moderate δ value (for instance, $\delta \sim 0.1$) is frequently sufficient to stabilize a matrix splitting. In such a case the debilitating drawbacks of the new coefficients pale in comparison to the alternative of having to use a stiff time step restriction.

7. Tables of new ImEx coefficients. This section presents the new ImEx coefficients (a_j, b_j, c_j) for $0 \leq j \leq r$ as a function of $0 < \delta \leq 1$. To use the coefficients in practice, (i) choose a small enough value of δ that guarantees unconditional stability, (ii) substitute the chosen value of δ into the tables in this section to obtain the time stepping coefficients at the required order.

Order		$j = 3$	$j = 2$	$j = 1$	$j = 0$
1	a_j	.	.	δ	$-\delta$
	c_j	.	.	1	$(\delta-1)$
	b_j	.	.	0	δ
2	a_j	.	$2\delta - \frac{1}{2}\delta^2$	$-4\delta + 2\delta^2$	$2\delta - \frac{3}{2}\delta^2$
	c_j	.	1	$2(\delta - 1)$	$(\delta - 1)^2$
	b_j	.	0	2δ	$(\delta - 1)^2 - 1$
3	a_j	$3\delta - \frac{3}{2}\delta^2 + \frac{1}{3}\delta^3$	$-9\delta + \frac{15}{2}\delta^2 - \frac{3}{2}\delta^3$	$9\delta - \frac{21}{2}\delta^2 + 3\delta^3$	$-3\delta + \frac{9}{2}\delta^2 - \frac{11}{6}\delta^3$
	c_j	1	$3(\delta - 1)$	$3(\delta - 1)^2$	$(\delta - 1)^3$
	b_j	0	3δ	$-6\delta + 3\delta^2$	$(\delta - 1)^3 + 1$

Order			$j = 4$
4	a_j	.	$4\delta - 3\delta^2 + \frac{4}{3}\delta^3 - \frac{1}{4}\delta^4$
	c_j	.	1
	b_j	.	0
		$j = 3$	$j = 2$
	a_j	$-16\delta + 18\delta^2 - \frac{22}{3}\delta^3 + \frac{4}{3}\delta^4$	$24\delta - 36\delta^2 + 18\delta^3 - 3\delta^4$
	c_j	$4(\delta - 1)$	$6(\delta - 1)^2$
	b_j	4δ	$-12\delta + 6\delta^2$
		$j = 1$	$j = 0$
	a_j	$-16\delta + 30\delta^2 - \frac{58}{3}\delta^3 + 4\delta^4$	$4\delta - 9\delta^2 + \frac{22}{3}\delta^3 - \frac{25}{12}\delta^4$
	c_j	$4(\delta - 1)^3$	$(\delta - 1)^4$
	b_j	$12\delta - 12\delta^2 + 4\delta^3$	$(\delta - 1)^4 - 1$

Order		$j = 5$	$j = 4$
5	a_j	$5\delta - 5\delta^2 + \frac{10}{3}\delta^3 - \frac{5}{4}\delta^4 + \frac{1}{5}\delta^5$	$-25\delta + 35\delta^2 - \frac{65}{3}\delta^3 + \frac{95}{12}\delta^4 - \frac{5}{4}\delta^5$
	c_j	1	$5(\delta - 1)$
	b_j	0	5δ
		$j = 3$	$j = 2$
	a_j	$50\delta - 90\delta^2 + \frac{190}{3}\delta^3 - \frac{65}{3}\delta^4 + \frac{10}{3}\delta^5$	$-50\delta + 110\delta^2 - \frac{280}{3}\delta^3 + 35\delta^4 - 5\delta^5$
	c_j	$10(\delta - 1)^2$	$10(\delta - 1)^3$
	b_j	$-20\delta + 10\delta^2$	$30\delta + 10\delta^3 - 30\delta^2$
		$j = 1$	$j = 0$
	a_j	$25\delta - 65\delta^2 + \frac{200}{3}\delta^3 - \frac{365}{12}\delta^4 + 5\delta^5$	$-5\delta + 15\delta^2 - \frac{55}{3}\delta^3 + \frac{125}{12}\delta^4 - \frac{137}{60}\delta^5$
	c_j	$5(\delta - 1)^4$	$(\delta - 1)^5$
	b_j	$-20\delta + 30\delta^2 - 20\delta^3 + 5\delta^4$	$(\delta - 1)^5 + 1$

REFERENCES

- [1] A. ABDULLE AND A. A. MEDOVNIKOV, *Second order Chebyshev methods based on orthogonal polynomials*, Numer. Math., 90 (2001), pp. 1–18.
- [2] L. AHLFORS, *Complex Analysis*, 3rd ed., McGraw-Hill, 1979.
- [3] G. AKRIVIS, *Implicit-explicit multistep methods for nonlinear parabolic equations*, Math. Comp., 82 (2012), pp. 45–68.
- [4] G. AKRIVIS, M. CROUZEIX, AND C. MAKRIDAKIS, *Implicit-explicit multistep finite element methods for nonlinear parabolic problems*, Math. Comp., 67 (1998), pp. 457–477.
- [5] G. AKRIVIS, M. CROUZEIX, AND C. MAKRIDAKIS, *Implicit-explicit multistep methods for quasi-linear parabolic equations*, Numer. Math., 82 (1999), pp. 521–541.
- [6] G. AKRIVIS AND F. KARAKATSANI, *Modified implicit-explicit BDF methods for nonlinear parabolic equations*, BIT, 43 (2003), pp. 467–483.
- [7] M. ANITESCU, W. LAYTON, AND F. PAHLEVANI, *Implicit for local effects, explicit for nonlocal is unconditionally stable*, Electron. Trans. Numer. Anal., 18 (2004), pp. 174–187.
- [8] U. M. ASCHER, S. J. RUUTH, AND B. T. R. WETTON, *Implicit-explicit methods for time-dependent partial differential equations*, SIAM J. Numer. Anal., 32 (1995), pp. 797–823, <https://doi.org/10.1137/0732037>.
- [9] V. BADALASSI, H. CENICEROS, AND S. BANERJEE, *Computation of multiphase systems with phase field models*, J. Comput. Phys., 190 (2003), pp. 371–397.
- [10] A. BERTOZZI, N. JU, AND J.-W. LU, *A biharmonic-modified forward time stepping method for fourth order nonlinear diffusion equations*, Discrete Contin. Dyn. Syst., 29 (2011), pp. 1367–1391.
- [11] J. CAHN AND J. HILLIARD, *Free energy of a nonuniform system I. Interfacial free energy*, J. Chem. Phys., 28 (1958), pp. 258–267.
- [12] H. CENICEROS, *A semi-implicit moving mesh method for the focusing nonlinear Schrödinger equation*, Commun. Pure Appl. Anal., 1 (2002), pp. 1–14.

- [13] A. CHRISTLIEB, J. JONES, K. PROMISLOW, B. WETTON, AND M. WILLOUGHBY, *High accuracy solutions to energy gradient flows from material science models*, J. Comput. Phys., 257 (2014), pp. 193–215.
- [14] P. CONCUS AND G. H. GOLUB, *Use of fast direct methods for the efficient numerical solution of nonseparable elliptic equations*, SIAM J. Numer. Anal., 10 (1973), pp. 1103–1120, <https://doi.org/10.1137/0710092>.
- [15] M. CROUZEIX, *Une méthode multipas implicite-explicite pour l'approximation des équations d'évolution paraboliques*, Numer. Math., 35 (1980), pp. 257–276.
- [16] J. DOUGLAS AND T. DUPONT, *Alternating-direction Galerkin methods on rectangles*, in Numerical Solution of Partial Differential Equations, Vol. II, SYNSPADE-1970 (University of Maryland, College Park, MD), B. Hubbard, ed., Academic Press, New York, 1971, pp. 133–213.
- [17] T. A. DRISCOLL, N. HALE, AND L. N. TREFETHEN, EDS., *Chebfun Guide*, 2014, Pafnuty Publications, Oxford, 2014.
- [18] M. ELSEY AND B. WIRTH, *A simple and efficient scheme for phase field crystal simulation*, ESAIM Math. Model. Numer. Anal., 47 (2013), pp. 1413–1432.
- [19] D. EYRE, *Unconditionally gradient stable time marching the Cahn-Hilliard equation*, in Computational and Mathematical Models of Microstructural Evolution, J. W. Bullard, R. Kalia, M. Stoneham, and L. Chen, eds., MRS Proc. 529, Cambridge University Press, Cambridge, UK, 1998, pp. 1686–1712, <https://doi.org/10.1557/PROC-529-39>.
- [20] J. FRANK, W. HUNSDORFER, AND J. VERWER, *On the stability of IMEX LM methods*, Appl. Numer. Math., 25 (1997), pp. 193–205.
- [21] K. GLASNER, *A diffuse interface approach to Hele-Shaw flow*, Nonlinearity, 16 (2003), pp. 49–66.
- [22] K. GLASNER AND S. ORIZAGA, *Improving the accuracy of convexity splitting methods for gradient flow equations*, J. Comput. Phys., 315 (2016), pp. 52–64.
- [23] D. GOTTLIEB AND S. A. ORSZAG, *Numerical Analysis of Spectral Methods*, CBMS-NSF Regional Conf. Ser. in Appl. Math. 26, SIAM, Philadelphia, 1977, <https://doi.org/10.1137/1.9781611970425>.
- [24] L. GREENGARD AND V. ROKHLIN, *A fast algorithm for particle simulations*, J. Comput. Phys., 73 (1987), pp. 325–348.
- [25] Z. GUAN, J. LOWENGRUB, C. WANG, AND S. WISE, *Second-order convex splitting schemes for periodic nonlocal Cahn-Hilliard and Allen-Cahn equations*, J. Comput. Phys., 277 (2014), pp. 48–71.
- [26] E. HAIRER, S. P. NØRSETT, AND G. WANNER, *Solving Ordinary Differential Equations I: Nonstiff Problems*, 2nd revised ed., Springer-Verlag, Berlin, 1987.
- [27] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*, Vol. 1, Springer-Verlag, Berlin, 1991.
- [28] W. HUNSDORFER AND J. VERWER, *Numerical Solution of Time-Dependent Advection-Diffusion-Reaction Equations*, Springer Ser. Comput. Math. 33, Springer, 2003.
- [29] R. JELTSCH AND O. NEVANLINNA, *Stability of explicit time discretizations for solving initial value problems*, Numer. Math., 37 (1981), pp. 61–91.
- [30] R. JELTSCH AND O. NEVANLINNA, *Stability and accuracy of time discretizations for initial value problems*, Numer. Math., 40 (1982), pp. 245–296.
- [31] C. R. JOHNSON, *Numerical determination of the field of values of a general complex matrix*, SIAM J. Numer. Anal., 15 (1978), pp. 595–602, <https://doi.org/10.1137/0715039>.
- [32] H. JOHNSTON AND J.-G. LIU, *Accurate, stable and efficient Navier-Stokes solvers based on explicit treatment of the pressure term*, J. Comput. Phys., 199 (2004), pp. 221–259.
- [33] L. JU, J. ZHANG, L. ZHU, AND Q. DU, *Fast explicit integration factor methods for semilinear parabolic equations*, J. Sci. Comput., 62 (2015), pp. 431–455.
- [34] G. KARNIADAKIS, M. ISRAELI, AND S. A. ORSZAG, *High-order splitting methods for the incompressible Navier-Stokes equations*, J. Comput. Phys., 97 (1991), pp. 414–443.
- [35] J. KIM AND P. MOIN, *Application of a fractional step method to incompressible Navier-Stokes equations*, J. Comput. Phys., 59 (1985), pp. 308–323.
- [36] T. KOTO, *Stability of implicit-explicit linear multistep methods for ordinary and delay differential equations*, Front. Math. China, 4 (2009), pp. 113–129.
- [37] W. LAYTON AND C. TRENCH, *Stability of two IMEX methods, CNLF and BDF2-AB2, for uncoupling systems of evolution equations*, Appl. Numer. Math., 62 (2012), pp. 112–120.
- [38] R. J. LEVEQUE, *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-Dependent Problems*, SIAM, Philadelphia, 2007, <https://doi.org/10.1137/1.9780898717839>.
- [39] J.-G. LIU, J. LIU, AND R. L. PEGO, *Stability and convergence of efficient Navier-Stokes solvers via a commutator estimate*, Comm. Pure Appl. Math., 60 (2007), pp. 1443–1487.

- [40] J.-G. LIU, J. LIU, AND R. L. PEGO, *Stable and accurate pressure approximation for unsteady incompressible viscous flow*, J. Comput. Phys., 229 (2010), pp. 3428–3453.
- [41] P. A. MILEWSKI AND E. G. TABAK, *A pseudospectral procedure for the solution of nonlinear wave equations with examples from free-surface flows*, SIAM J. Sci. Comput., 21 (1999), pp. 1102–1114, <https://doi.org/10.1137/S1064827597321532>.
- [42] M. L. MINION, *Semi-implicit spectral deferred correction methods for ordinary differential equations*, Commun. Math. Sci., 1 (2003), pp. 471–500.
- [43] G. SHENG, T. WANG, Q. DU, K. WANG, Z. LIU, AND L. Q. CHEN, *Coarsening kinetics of a two phase mixture with highly disparate diffusion mobility*, Commun. Comput. Phys., 8 (2010), pp. 249–264.
- [44] D. SHIROKOFF AND R. R. ROSALES, *An efficient method for the incompressible Navier-Stokes equations on irregular domains with no-slip boundary conditions, high order up to the boundary*, J. Comput. Phys., 230 (2011), pp. 8619–8646.
- [45] P. SMEREKA, *Semi-implicit level set methods for curvature and surface diffusion motion*, J. Sci. Comput., 19 (2003), pp. 439–456.
- [46] L. N. TREFETHEN, *Spectral Methods in MATLAB*, SIAM, Philadelphia, 2000, <https://doi.org/10.1137/1.9780898719598>.
- [47] L. N. TREFETHEN AND D. BAU, *Numerical Linear Algebra*, SIAM, Philadelphia, 2000.
- [48] C. TRENCHIA, *Second order implicit for local effects and explicit for nonlocal effects is unconditionally stable*, Romai J., 12 (2016), pp. 163–178.
- [49] J. M. VARAH, *Stability restrictions on second order, three level finite difference schemes for parabolic equations*, SIAM J. Numer. Anal., 17 (1980), pp. 300–309, <https://doi.org/10.1137/0717025>.
- [50] J. XU, Y. LI, AND S. WU, *On the Accuracy of Partially Implicit Schemes for Phase Field Modeling*, preprint, <https://arxiv.org/abs/1604.05402>, 2016.