

**RESEARCH ARTICLE**

# On the stability of Runge–Kutta methods for arbitrarily large systems of ODEs

**Eitan Tadmor**University of Maryland, College Park,  
Maryland, USA**Correspondence**Eitan Tadmor, University of Maryland,  
College Park, Maryland, USA.Email: [tadmor@umd.edu](mailto:tadmor@umd.edu)**Funding information**ONR, Grant/Award Numbers:  
N00014-2112773, N00014-2412659;  
Fondation Sciences Mathématiques de  
Paris**Abstract**

We prove that Runge–Kutta (RK) methods for numerical integration of arbitrarily large systems of Ordinary Differential Equations are linearly stable. Standard stability arguments—based on spectral analysis, resolvent condition or strong stability, fail to secure the stability of RK methods for arbitrarily large systems. We explain the failure of different approaches, offer a new stability theory based on the numerical range of the underlying large matrices involved in such systems, and demonstrate its application with concrete examples of RK stability for hyperbolic methods of lines.

**CONTENTS**

1. INTRODUCTION . . . . .	822
2. SPECTRAL, RESOLVENT, AND STRONG STABILITY ANALYSIS ARE NOT ENOUGH	828
3. NUMERICAL RANGE AND STABILITY OF COERCIVE RUNGE—KUTTA SCHEMES	833
4. SPECTRAL SETS AND STABILITY OF RUNGE–KUTTA METHODS . . . . .	838
5. EXAMPLES: STABILITY OF TIME-DEPENDENT METHODS OF LINES . . . . .	843
ACKNOWLEDGMENTS . . . . .	852
REFERENCES . . . . .	852
APPENDIX A: THE NUMERICAL RANGE IS $(1 + \sqrt{2})$ -SPECTRAL SET . . . . .	855

-----  
 This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2024 The Author(s). *Communications on Pure and Applied Mathematics* published by Courant Institute of Mathematics and Wiley Periodicals LLC.

# 1 | INTRODUCTION

Runge–Kutta (RK) methods are widely used class of effective methods for numerical integration of systems of Ordinary Differential Equations (ODEs). In particular, such methods are used routinely for integration of large systems of ODEs encountered in various applications. As examples we mention RK integration of large systems of ODEs in molecular dynamics in Chemistry, in many particle systems in Physics, in climate modeling, in cosmology and in spatial discretization of time-dependent PDEs which end up with increasingly large systems of ODEs, so-called “method of lines.” In recent years such problems also arise in integration of high-dimensional data sets/neural networks, for example [4, 12, 22, 45].

## 1.1 | An informal summary of main results

The stability of RK methods encoded in terms of their *region of absolute stability* is well documented [3, 23, 28]. We therefore begin with an informal summary of our results, clarifying the claim made in the title.

We consider linear systems of ODEs,  $\dot{\mathbf{y}} = \mathbb{L}_N \mathbf{y}$ , associated with a general class of  $N \times N$  semi-bounded matrices,  $\mathbb{L}_N$ , and a RK method associated with the polynomial,  $\mathcal{P}(z) = \sum_{k=0}^s a_k z^k$ . The RK method is stable if its computed solution evolved in time, remains comparable to the size of the (initial) data, uniformly in the number of time steps,  $n$ , and the size of the underlying system,  $N$ . Thus, the linearized stability of the RK method in the present context requires

$$\|\mathcal{P}^n(\Delta t \mathbb{L}_N)\| \leq K_{\mathbb{L}}, \quad n = 1, 2, \dots,$$

with a constant  $K_{\mathbb{L}}$  independent of  $n$  and  $N$ . This leads to the classical stability criterion which requires the time-step  $\Delta t$  to be small enough so that

$$\Delta t \sigma(\mathbb{L}_N) \subset \mathcal{A}.$$

Here,  $\sigma(\mathbb{L}_N)$  is the spectrum of  $\mathbb{L}_N$  and  $\mathcal{A} = \{z \in \mathbb{C} : |\mathcal{P}(z)| \leq 1\}$  is the region of absolute stability associated with the RK method under consideration. This classical framework of stability suffices for systems of finite size but *fails* for arbitrarily large systems. There is an extensive literature, going back to the 1980’s, which tried to secure a uniform-in- $N$  stability bound by adapting alternative notions of resolvent stability or strong stability. We discuss the failure of different approaches in §1.2 and further elaborate in §2 below. Alternatively, there were different approaches to secure the uniform-in- $N$  stability bound for restricted classes of  $\mathbb{L}_N$ ’s, satisfying different coercivity restrictions. We mention the strong stability preserving (SSP) theory which necessitates  $\Delta t \sigma(\mathbb{L}_N) \subset \{z : |1 + z| \leq 1\}$  [18] (see §3.4 below), and the wedge condition,  $\Delta t \sigma(\mathbb{L}_N) \subset \mathcal{A} \cap \{z : |\arg z| \geq \pi - \alpha\}$  with  $\alpha < \pi/2$  [57]. The uniform-in- $N$  stability question for general semi-bounded  $\mathbb{L}_N$ ’s remained open. This is addressed in our first main result in Theorem 4.2 below, stating that if

$$\Delta t W(\mathbb{L}_N) \subset \mathcal{A},$$

then stability follows with uniformly bounded  $K_{\mathbb{L}}$ . Here,  $W(\mathbb{L}_N) = \{\langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle : |\mathbf{x}| = 1\}$  is the *numerical range* of  $\mathbb{L}_N$  (associated with general inner product  $\langle \cdot, \cdot \rangle$ ; consult §3.1 below). Concrete examples for applications of this RK stability result are given in §5, in the context of methods of

lines for hyperbolic transport equations, where we recover classical old results and derive some new ones. In this case, we explicitly compute  $W(\mathbb{L}_N)$  for circulant or almost circulant matrices. But in general, the structure of  $W(\mathbb{L}_N)$  as a set in the complex plane is not as accessible as the discrete spectrum  $\sigma(\mathbb{L}_N)$ . This is addressed in our second main Theorem 4.4. Consider RK method whose region of absolute stability contains a non-trivial interval along the imaginary axis,  $\mathcal{A} \supset [-iR, iR]$ . Then, there exists a constant  $C > 0$  (depending on  $R$ ) so that the Courant–Friedrichs–Levy (CFL)-like condition  $\Delta t \|\mathbb{L}_N\| \leq C$  implies  $\Delta t W(\mathbb{L}_N) \subset \mathcal{A}$ , and uniform-in- $N$  stability follows.

### 1.2 | The quest for stability

We consider systems of ODEs,

$$\dot{\mathbf{y}} = \mathbf{F}(t, \mathbf{y}),$$

which govern an  $N$ -vector of unknown solution,  $\mathbf{y}(t) \in \mathbb{R}^N$ , subject to prescribed initial data,  $\mathbf{y}(t_0) = \mathbf{y}_0$ . As a canonical example for one of the most widely used numerical integrators we mention the 4-stage RK method, which computes an approximate solution,  $\{\mathbf{u}_n = \mathbf{u}(t_n)\}_{n>0}$ , at successive time steps  $t_{n+1} := t_n + \Delta t$  [23, §II.1],

$$\mathbf{u}_{n+1} = \mathbf{u}_n + \frac{\Delta t}{6} (\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4) \quad \begin{cases} \mathbf{k}_1 = \mathbf{F}(t_n, \mathbf{u}_n) \\ \mathbf{k}_2 = \mathbf{F}(t_{n+1/2}, \mathbf{u}_n + (\Delta t/2)\mathbf{k}_1) \\ \mathbf{k}_3 = \mathbf{F}(t_{n+1/2}, \mathbf{u}_n + (\Delta t/2)\mathbf{k}_2) \\ \mathbf{k}_4 = \mathbf{F}(t_{n+1}, \mathbf{u}_n + \Delta t\mathbf{k}_3). \end{cases} \quad (1.1)$$

The *linearized stability analysis* examines the behavior of (1.1) for linear systems,  $\mathbf{F}(t, \mathbf{y}) = \mathbb{L}_N \mathbf{y}$ ,

$$\dot{\mathbf{y}} = \mathbb{L}_N \mathbf{y}, \quad (1.2)$$

where (1.1) is reduced to

$$\mathbf{u}_{n+1} = \left( \mathbb{I} + \Delta t \mathbb{L}_N + \frac{1}{2}(\Delta t \mathbb{L}_N)^2 + \frac{1}{6}(\Delta t \mathbb{L}_N)^3 + \frac{1}{24}(\Delta t \mathbb{L}_N)^4 \right) \mathbf{u}_n, \quad n = 0, 1, \dots \quad (\text{RK4})$$

The corresponding iterations for a general  $s$ -stage *explicit* RK method take the form

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathbf{u}_n, \quad n = 0, 1, 2, \dots, \quad \mathcal{P}_s(z) := \sum_{k=0}^s a_k z^k, \quad a_k \in \mathbb{R}, \quad a_s \neq 0. \quad (1.3)$$

Different  $\{a_k\}_{k=0}^s$  dictate different RK methods with emphasis on different aspects of accuracy, efficiency and stability. The resulting  $s$ -stage RK methods (1.3), involve  $N \times N$  matrices, denoted  $\mathbb{L}_N$  to highlight the fact that they are parameterized with respect to  $N$ . As already noted above, such large matrices are often encountered in applications, and we therefore pay particular attention to the question of RK stability that is uniform with respect to the increasingly large dimension  $N$ .

Following [65, §2], we consider (1.2) for the class of *semi-bounded*  $\mathbb{L}_N$ 's, namely —  $\mathbb{L}_N$ 's for which there exist constants  $\eta, K_{\mathbb{H}} > 0$  independent of  $N$ , and *uniformly* positive-definite symmetrizers,  $\mathbb{H}_N$ 's, such that<sup>1</sup>,

$$\mathbb{H}_N \mathbb{L}_N^{\top} + \mathbb{L}_N \mathbb{H}_N \leq 2\eta \mathbb{H}_N, \quad 0 < K_{\mathbb{H}}^{-1} \leq \mathbb{H}_N \leq K_{\mathbb{H}}.$$

It follows that the solutions of the corresponding semi-bounded ODEs (1.2) subject to arbitrary initial data  $\mathbf{y}(0) = \mathbf{y}_0$ , satisfy

$$|\mathbf{y}(t)|_{\ell^2} \leq K_{\mathbb{H}} e^{\eta t} |\mathbf{y}_0|_{\ell^2}.$$

Replacing  $\mathbb{L}_N$  with  $\mathbb{L}_N - \eta \mathbb{I}$ , allows us to consider without loss of generality the case  $\eta = 0$ , corresponding to *negative definite*  $\mathbb{L}_N$ 's,

$$\mathbb{H}_N \mathbb{L}_N^{\top} + \mathbb{L}_N \mathbb{H}_N \leq 0, \quad 0 < K_{\mathbb{H}}^{-1} \leq \mathbb{H}_N \leq K_{\mathbb{H}}. \quad (1.4)$$

Solutions of ODE governed by such negative<sup>2</sup>  $\mathbb{L}_N$ 's remain uniformly bounded in time relative to their initial data  $\mathbf{y}_0$ ,

$$|\mathbf{y}(t)|_{\ell^2} \leq K_{\mathbb{H}} |\mathbf{y}_0|_{\ell^2}. \quad (1.5)$$

**Stability of RK scheme.** The notion of stability of RK schemes requires the numerical solution to satisfy the bound corresponding to (1.5). To this end, one is focused on a family of negative  $\mathbb{L}_N$ 's parametrized by their dimension  $N$ . The  $s$ -stage RK scheme (1.3) is *stable*, if there exist constants,  $K_{\perp} > 0$  and  $C_s > 0$  independent of  $N$ , such that solutions of (1.3) subject to arbitrary initial data  $\mathbf{u}_0$  satisfy, for small enough time step  $\Delta t$ ,

$$\text{Stability of RK scheme: } |\mathbf{u}_n|_{\ell^2} \leq K_{\perp} |\mathbf{u}_0|_{\ell^2}, \quad n = 0, 1, 2, \dots \quad (1.6)$$

The restriction of having small enough time step is encoded in terms of the bound

$$\Delta t \cdot \|\mathbb{L}_N\| \leq C_s; \quad (1.7)$$

in the context of method of lines, the time-step restriction is related to the celebrated CFL condition [5], and we shall therefore often refer to the time-step restriction (1.7) as a CFL condition.

The notion of stability encoded in (1.6) amounts to the question of power-boundedness of  $\mathcal{P}_s(\Delta t \mathbb{L}_N)$ ,

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq K_{\perp}, \quad n = 0, 1, 2, \dots \quad (1.8)$$

*Remark 1.1* (Stability and linearization). The general notion of stability for semi-bounded  $\mathbb{L}_N$ 's, limits the exponential stability bound to a finite time interval,

$$|\mathbf{u}_n|_{\ell^2} \leq K_{\perp} e^{\eta t} |\mathbf{u}_0|_{\ell^2}, \quad n \cdot \Delta t \leq t.$$

<sup>1</sup> Throughout the paper, we use  $K_{\square}$  to denote different constants which are independent of  $N$ .

<sup>2</sup> Throughout the work, we use the term “negative” for short of “negative definite.”

Since we restrict attention to negative  $\mathbb{L}_N$ 's, we may as well let  $n \in \mathbb{N}$ . This notion of stability is invariant against low-order perturbations [30, 59], and therefore allows to recover the stability of RK schemes for smooth solutions of fully nonlinear problems,  $\dot{\mathbf{y}} = \mathbf{F}(t, \mathbf{y})$ . To this end, one can linearize and freeze coefficients at arbitrary  $t = t_*$ , arriving at the linearized system (1.2),

$$\dot{\mathbf{y}} = \mathbb{L}_N \mathbf{y} \quad \text{with} \quad \mathbb{L}_N = \frac{\partial \mathbf{F}(t_*, \mathbf{y}(t_*))}{\partial \mathbf{y}}.$$

We shall not dwell on the details, expect for referring to our discussion on stability in presence of variable coefficients in §5.2 below. This motivates our focus on the question of linearized stability, where  $\mathbb{L}_N$  is a substitute for the  $N \times N$  gradient matrix frozen at arbitrary state.

### 1.3 | Spectral stability analysis

The standard approach to address the question of power-boundedness is *spectral analysis*, in which (1.8) requires  $\max_{1 \leq k \leq N} |\lambda_k(\mathcal{P}_s(\Delta t \mathbb{L}_N))| \leq 1$ . By the spectral mapping theorem,

$$\lambda_k(\mathcal{P}_s(\Delta t \mathbb{L}_N)) = \mathcal{P}_s(\Delta t \lambda_k(\mathbb{L}_N)), \tag{1.9}$$

which leads to the necessary stability condition, requiring small enough time-step dictated by the *region of absolute stability* associated with (1.3),

$$\Delta t \cdot \lambda_k(\mathbb{L}_N) \in \mathcal{A}_s, \quad k = 1, 2, \dots, N, \quad \mathcal{A}_s := \{z \in \mathbb{C} : |\mathcal{P}_s(z)| \leq 1\}. \tag{1.10}$$

Conversely, consider the favorite scenario in which  $\mathbb{L}_N$  is diagonalizable,

$$\mathbb{T}_N \mathbb{L}_N \mathbb{T}_N^{-1} = \Lambda, \quad \Lambda = \begin{bmatrix} \lambda_1(\mathbb{L}_N) & 0 & \dots & \dots & 0 \\ 0 & \lambda_2(\mathbb{L}_N) & \ddots & & \vdots \\ \vdots & \ddots & & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & \dots & 0 & \lambda_N(\mathbb{L}_N) \end{bmatrix}.$$

Then  $\mathcal{P}_s(\Delta t \mathbb{L}_N) = \mathbb{T}_N^{-1} \mathcal{P}_s(\Delta t \Lambda) \mathbb{T}_N$  and (1.10) implies

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| = \|\mathbb{T}_N^{-1} \mathcal{P}_s^n(\Delta t \Lambda) \mathbb{T}_N\| \leq \|\mathbb{T}_N^{-1}\| \cdot \|\mathbb{T}_N\|. \tag{1.11}$$

This guarantees the stability of RK schemes for systems of finite *fixed* dimension.<sup>3</sup> However, here we insist that the stability sought in (1.6) will apply uniformly for increasingly large systems, and since the condition number on the right of (1.11),  $\|\mathbb{T}_N^{-1}\| \cdot \|\mathbb{T}_N\|$ , may grow with  $N$ , spectral condition (1.10) is not enough to secure the desired uniform-in- $N$  stability bound. Indeed, as we elaborate in §2.1 below, the general question of stability, uniformly in  $N$ , cannot be addressed solely in terms of spectral analysis.

<sup>3</sup>The precise necessary and sufficient characterization for power-boundedness of a single matrix,  $\|\mathcal{P}^n\| \leq K$ , requires that the eigenvalues  $\lambda_k(\mathcal{P})$  are inside the unit disc and those on the unit circle are simple or non-defective in the sense of having fully diagonalizable eigenspace; the constant  $K$  may still depend on the dimension of  $\mathcal{P}$ .

## 1.4 | Resolvent stability

We now appeal to a stronger notion of stability of RK method. An  $s$ -stage RK method  $\mathcal{P}_s(\cdot)$  is *stable* if the corresponding RK schemes (1.3) are stable for all negative  $\mathbb{L}_N$ 's,

$$\text{Stability of RK method: } \|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq K_{\mathbb{L}} \text{ for all negative } \mathbb{L}_N' \text{'s.} \quad (1.12)$$

Observe that we are making a distinction between the stability of RK *scheme* — which examines the boundedness of RK protocol  $\mathcal{P}_s^n(\Delta t \mathbb{L}_N)$  for a specific family of negative  $\mathbb{L}_N$ 's, vs. the stability of RK *method*—which examines the behavior of RK protocol  $\mathcal{P}_s^n(\Delta t \cdot)$ , for *all* negative  $\mathbb{L}_N$ 's.

This stronger notion of stability restricts the class of stable RK methods. In particular, their stability question should apply to the scalar ODEs,  $\dot{y} = \lambda y$ , for all negative  $\text{Re } \lambda \leq 0$ , which in turn implies that (1.10) must hold for purely imaginary  $\lambda = i\sigma$ , so that  $|\mathcal{P}_s(i\Delta t\sigma)| \leq 1$ , for small enough step-size,  $\Delta t$ . In other words, a stable RK method *must* satisfy the following interval condition.

**Definition 1.2** Imaginary interval condition<sup>4</sup>. A Runge–Kutta method is said to satisfy the *imaginary interval condition* if there exists a constant  $R_s > 0$  such that

$$|\mathcal{P}_s(i\sigma)| \leq 1, \quad -R_s \leq \sigma \leq R_s. \quad (1.13)$$

In other words, the region of absolute stability of a stable RK method must contain a non-trivial interval along the imaginary axis  $[-iR_s, iR_s] \subset \mathcal{A}_s$ . This secures the stability of RK method for scalar hyperbolic ODEs,  $\dot{y} = i\sigma y$ , with small enough step-size  $\Delta t\sigma < R_s$ .

The interval condition excludes the standard 1-stage forward Euler method (for historical perspective of Euler's method which dates back to 1768 see [69, §1]),

$$\text{Forward Euler: } \mathbf{u}_{n+1} = (\mathbb{I} + \Delta t \mathbb{L}_N) \mathbf{u}_n, \quad (\text{RK1})$$

for which  $\mathcal{P}_1(z) = 1 + z \rightsquigarrow |\mathcal{P}_1(i\sigma)| > 1$  for all  $\sigma \neq 0$ . The imaginary interval condition (1.13) also excludes the 2-stage Heun's method [9, §8.3.3] (also known as modified Euler method),

$$\text{Heun method: } \mathbf{u}_{n+1} = \left( \mathbb{I} + \Delta t \mathbb{L}_N + \frac{1}{2}(\Delta t \mathbb{L}_N)^2 \right) \mathbf{u}_n, \quad (\text{RK2})$$

since  $\mathcal{P}_2(z) = 1 + z + \frac{1}{2}z^2 \rightsquigarrow |\mathcal{P}_2(i\sigma)| > 1$  for all  $\sigma \neq 0$ .

On the other hand, the 3-stage Kutta method,

$$\text{Kutta method: } \mathbf{u}_{n+1} = \left( \mathbb{I} + \Delta t \mathbb{L}_N + \frac{1}{2}(\Delta t \mathbb{L}_N)^2 + \frac{1}{6}(\Delta t \mathbb{L}_N)^3 \right) \mathbf{u}_n, \quad (\text{RK3})$$

as well as the 4-stage Runge–Kutta method, (RK4), and its higher-order embedded version RK45 of Dormand-Prince method [11, 23, 35], do satisfy the interval condition with  $R_3 = \sqrt{3}$ , and respectively,  $R_4 = 2\sqrt{2}$ ; this is depicted in Figure 1. A precise characterization of general  $s$ -stage RK methods satisfying the interval condition was given in [35, Theorem 3.1] and will be recalled in (4.8b) below.

<sup>4</sup> So-called “local stability along the imaginary line” in [35, Definition 2.1].

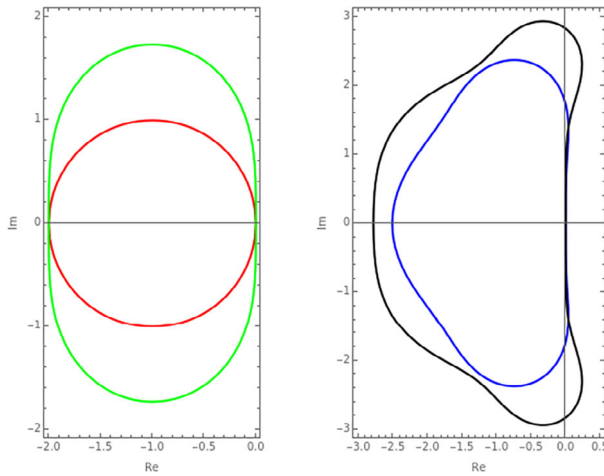


FIGURE 1 Regions of absolute stability,  $\mathcal{A}_s$ ,  $s = 1, 2$  (left), and  $s = 3, 4$  (right).

The interval condition 1.2 is necessary for stability of a RK method. Kreiss and Wu [36], proved the converse in the sense that the interval condition is sufficient for *resolvent stability*, namely—(1.13) implies that the following holds.

**Definition 1.3** (Resolvent stability). The RK method is *resolvent stable* if there exist constants  $K_R > 0$  and  $0 < C_s < R_s$ , independent of  $N$ , such that for small step-size,

$$\|(z\mathbb{I} - \mathcal{P}_s(\Delta t \mathbb{L}_N))^{-1}\| \leq \frac{K_R}{|z| - 1}, \quad \forall |z| > 1, \quad \Delta t \cdot \|\mathbb{L}_N\| \leq C_s. \tag{1.14}$$

So the interval condition implies resolvent stability which in turn guarantees the stability of RK schemes for systems of finite *fixed* dimension, in view of the Kreiss matrix theorem [30, 51, §4.9]. Indeed, in [62] and its improvement [42], it was proved that (1.14) implies

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq 2eK_R N, \quad n = 1, 2, \dots \tag{1.15}$$

However, as we shall elaborate in §2.2 below, the  $N$ -dependent bound on the right cannot be completely removed and hence resolvent stability does not secure the desired stability uniformly for arbitrarily large  $N$ .

### 1.5 | Strong stability

A Runge–Kutta scheme (1.3) is *strongly stable* if there exists  $K_{\mathcal{T}} > 0$  independent of  $N$  such that  $\mathcal{P}_s(\Delta t \mathbb{L}_N)$  is uniformly similar to a contraction,

$$\|\mathcal{T}_N \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathcal{T}_N^{-1}\| \leq 1, \quad \|\mathcal{T}_N^{-1}\| \cdot \|\mathcal{T}_N\| \leq K_{\mathcal{T}}. \tag{1.16}$$

A strongly stable RK scheme is clearly stable, for

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| = \|\mathcal{T}_N^{-1} (\mathcal{T}_N \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathcal{T}_N^{-1})^n \mathcal{T}_N\| \leq \|\mathcal{T}_N^{-1}\| \cdot \|\mathcal{T}_N\| \leq K_{\mathcal{T}} \tag{1.17}$$

The choice  $\mathcal{T}_N = \mathbb{T}_N$  recovers (1.11) as a special case of (1.17).

To secure strong stability it remains to construct a uniformly bounded symmetrizer  $\mathcal{H}_N := \mathcal{T}_N^* \mathcal{T}_N$  with  $0 < K_{\mathcal{T}}^{-1} \leq \mathcal{H}_N \leq K_{\mathcal{T}}$ . We addressed this issue in [65], proving the strong stability of the 3-stage RK method (RK3) with symmetrizer  $\mathcal{H}_N = \mathbb{H}_N$  and  $C_3 = 1$ , thus providing the first example of a RK method which is stable uniformly for arbitrarily large system of ODEs. It was later extended to all  $s$ -stage RK methods of order  $s = 3[\text{mod}4]$ , [61]. The question arises whether strong stability can be extended using proper symmetrizers,  $\mathcal{H}_N$ , for other  $s$ -stage RK methods for arbitrary  $s$ ? In [65] we conjectured that the 4-stage (RK4) fails strong stability in the sense that it is not uniformly similar to a contraction, or equivalently — as outlined in §2.3 below, that there is no symmetrizer  $\mathcal{H}_N := \mathcal{T}_N^* \mathcal{T}_N$  such that (1.16) <sub>$s=4$</sub>  holds. This was confirmed in [60, Proposition 1.1] and was later extended in [1, Theorem 2], where it was shown that strong stability *fails* for all  $s$ -stage  $p$ -order accurate<sup>5</sup> RK methods with  $s = r \in 4\mathbb{N}$ .

*Remark 1.4.* In fact, the issue of RK4 stability is more subtle as  $\|\mathcal{P}_s^{2n}(\Delta t \mathbb{L}_N)\| \leq 1$ , which can be interpreted to say that the 8-stage RK4 is strongly stable. We refer the interested reader to [60, 61] and the references therein.

**The stability question for RK schemes.** We come out from the above discussion, lacking a definitive answer to the question of stability of RK schemes/methods for arbitrarily large systems of ODEs. Thus, for example, the stability question for the widely used RK4 remains open. At this stage, the three different approaches — spectral analysis, resolvent condition and strong stability failed to determine whether RK4 method for example, is stable uniformly in  $N$  for the general class of negative  $\mathbb{L}_N$ 's. We therefore raise the question:

Are the Runge–Kutta methods (1.3) stable for arbitrarily large semi-bounded systems?

The title of the paper is an affirmative answer to this question. The answer is given in §4 in terms of the numerical range of  $\mathbb{L}_N$ .

## 2 | SPECTRAL, RESOLVENT, AND STRONG STABILITY ANALYSIS ARE NOT ENOUGH

In this section, we further elaborate with specific counterexamples, on the failure of spectral analysis, resolvent condition and strong stability to capture the uniform-in- $N$  stability of general  $s$ -stage RK schemes/methods. Spectral and resolvent analysis are shown to be too weak to secure stability, while strong stability argument is too restrictive.

### 2.1 | Spectral analysis is not enough

We recall the spectral analysis led to the necessary stability condition (1.10)

$$\Delta t \cdot \lambda_k(\mathbb{L}_N) \subset \mathcal{A}_s, \quad k = 1, 2, \dots, N.$$

<sup>5</sup>The RK method (1.3) is  $r$ -order accurate if  $|e^z - \mathcal{P}_s(z)| = \mathcal{O}(|z|^{r+1})$ ,  $|z| \ll 1$ ; see (4.8a) below. Thus,  $r$  is the largest index for which  $a_k = 1/k!$  for  $k = 1, 2, \dots, r$ .



As noted above, this spectral condition is not sufficient to secure stability in case of ill-conditioned eigensystems,  $\|\mathbb{T}_N^{-1}\| \cdot \|\mathbb{T}_N\|$ , which grows with  $N$ . An alternative approach, trying to circumvent this difficulty of ill-conditioning is to use a unitary triangulation

$$\mathcal{P}_s(\Delta t \mathbb{L}_N) = \mathbb{U}_N^* (\Lambda_{\mathcal{P}} + \mathbb{R}_N) \mathbb{U}_N, \quad \Lambda_{\mathcal{P}} := \mathcal{P}_s(\Delta t \Lambda),$$

where  $\Lambda$  and  $\Lambda_{\mathcal{P}}$  are the diagonals made of the eigenvalues of  $\mathbb{L}_N$  and, respectively,  $\mathcal{P}_s(\Delta t \mathbb{L}_N)$ , and  $\mathbb{R}_N$  is a nilpotent upper triangular matrix,  $(\mathbb{R}_N)_{ij} = 0, j \leq i$ . Since  $\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| = \|(\Lambda_{\mathcal{P}} + \mathbb{R}_N)^n\|$ , it remains to study the power-boundedness of the triangular matrix  $\Lambda_{\mathcal{P}} + \mathbb{R}_N$ . But we claim that even a most favorable scenario, in which the spectral stability analysis (1.10) secures the eigenvalues *strictly* inside the unit disc,

$$\theta := \max_{1 \leq k \leq N} |\mathcal{P}_s(\Delta t \lambda_k(\mathbb{L}_N))| < 1, \tag{2.1}$$

will not suffice to guarantee the stability of RK method. Indeed, we may assume without restriction that  $\mathbb{R}_N$  is arbitrarily small by its further re-scaling,<sup>6</sup> so that

$$\begin{aligned} \|\mathbb{S}_{\delta} \mathbb{R}_N \mathbb{S}_{\delta}^{-1}\| &= \|\{\mathbb{R}_{ij} \delta_{\epsilon}^{i-j}\}_{j>i}\| \leq \epsilon, \\ \mathbb{S}_{\delta} &= \begin{bmatrix} \delta_{\epsilon} & 0 & \dots & 0 \\ 0 & \delta_{\epsilon}^2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & \delta_{\epsilon}^N \end{bmatrix}, \quad \delta_{\epsilon} := \frac{\epsilon}{\|\mathbb{R}_N\|_F}. \end{aligned}$$

Here, an arbitrary  $\epsilon > 0$  is at our disposal to be determined below. It follows that

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| = \|\mathbb{U}_N^* \mathbb{S}_{\delta}^{-1} (\Lambda_{\mathcal{P}} + \mathbb{S}_{\delta} \mathbb{R}_N \mathbb{S}_{\delta}^{-1})^n \mathbb{S}_{\delta} \mathbb{U}_N\| \leq \|\mathbb{S}_{\delta}^{-1}\| \times (\|\Lambda_{\mathcal{P}}\| + \epsilon)^n \times \|\mathbb{S}_{\delta}\|.$$

By assumption,  $\|\Lambda_{\mathcal{P}}\| = \theta < 1$ . Set  $\epsilon := 1/2(1 - \theta)$ , we then end up with the stability bound

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq \delta_{\epsilon}^{1-N} \left(\frac{1 + \theta}{2}\right)^n = \left(\frac{2\|\mathbb{R}_N\|_F}{1 - \theta}\right)^{N-1} \left(\frac{1 + \theta}{2}\right)^n. \tag{2.2}$$

This bound secures the stability of finite dimensional systems – in fact, it recovers the well-known fact that matrices of finite *fixed* dimension with eigenvalues strictly inside the unit disc have exponentially decreasing iterates. But the argument breaks down when we examine the dependence on  $N$ , since the bound (2.2) is *not* uniform in  $N$ : for  $n = N - 1$ , for example, we find that unless  $\mathbb{R}_N$  is sufficiently small,<sup>7</sup> then there is an *exponential growth* in  $N$ ,

$$\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq \left(\frac{2\|\mathbb{R}_N\|_F}{1 - \theta}\right)^{N-1} \left(\frac{1 + \theta}{2}\right)_{n=N-1}^n = \left(\|\mathbb{R}_N\|_F \frac{1 + \theta}{1 - \theta}\right)^{N-1}. \tag{2.3}$$

<sup>6</sup>  $\|\cdot\|_F$  refers to Frobenius norm,  $\|A\|_F^2 = \text{trace}(A^T A)$ .

<sup>7</sup> To avoid an exponential growth of the upper-bound in (2.2) requires  $\|\mathbb{R}_N\|_F \leq \frac{1-\theta}{1+\theta}$ ; a more delicate tuning of the scaling parameter  $\delta_{\epsilon}$  shows that uniform bound is achieved for  $\|\mathbb{R}_N\|_F < 1 - \theta$ .

This bound is sharp in the sense that the power-growth hinted on the right of (2.3) is realized by the powers of the increasingly large  $N \times N$  Jordan blocks

$$\|\mathbb{J}_q^n\| \sim \left(\frac{2}{1-q}\right)^N \left(\frac{1+q}{2}\right)^n, \quad \mathbb{J}_q := \begin{bmatrix} -q & 1+q & \dots & \dots & 0 \\ 0 & -q & 1+q & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & -q & 1+q \\ 0 & \dots & \dots & \dots & -q \end{bmatrix}. \quad (2.4)$$

Although  $|\lambda_k(\mathbb{J}_q)| < 1$  for  $-1 < q < 1$ , there is a non-uniform growth of  $\|\mathbb{J}_q^n\|$  with  $0 < q < 1$ , corresponding to  $q = \theta$  in (2.3), when  $n \sim N \uparrow \infty$ . These increasingly large Jordan blocks realize the extreme case of ill-conditioning warned in (1.11).

### 2.1.1 | Instability of forward Euler scheme

The extremal example (2.4) is not just of academic interest. The following classical example [51, §6.6], [36, §3], [65, §5.1] sheds light on what can go wrong with spectral analysis. Consider the transport equation with fixed speed  $a > 0$

$$\begin{cases} y_t(x, t) = ay_x(x, t), & (t, x) \in \mathbb{R}_+ \times (0, 1) \\ y(1, t) = 0. \end{cases} \quad (2.5)$$

Its spatial part is discretized using one-sided spatial differences on equi-spaced grid,  $\{x_\nu := \nu \Delta x\}_{\nu=0}^N$ ,  $\Delta x = 1/N$ , covering the interval  $[0, 1]$ ,

$$\begin{cases} \frac{d}{dt} y(x_\nu, t) = a \frac{y(x_{\nu+1}, t) - y(x_\nu, t)}{\Delta x}, & \nu = 0, 1, \dots, N-1, \\ y(x_N, t) = 0. \end{cases} \quad (2.6)$$

This amounts to method of lines for the  $N$ -vector of unknowns,  $\mathbf{y}(t) := (y(x_0, t), \dots, y(x_{N-1}, t))^T$ , governed by the  $N \times N$  semi-discrete system in terms of the forward-difference operator  $\mathbb{D}_N^+$ ,

$$\dot{\mathbf{y}}(t) = a \mathbb{D}_N^+ \mathbf{y}, \quad \mathbb{D}_N^+ := \frac{1}{\Delta x} \begin{bmatrix} -1 & 1 & \dots & \dots & 0 \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & -1 & 1 \\ 0 & \dots & \dots & \dots & -1 \end{bmatrix}. \quad (2.7)$$

Observe that  $\mathbb{D}_N^+$  is semi-bounded — in fact it is *strictly dissipative* in the sense that

$$(\mathbb{D}_N^+)^T + \mathbb{D}_N^+ \leq -2 \left(1 - \cos\left(\frac{\pi}{N+1}\right)\right) \mathbb{I}_{N \times N}.$$

This system (2.7) is integrated using one-stage Forward Euler method, (RK1), augmented with boundary condition  $u(x_N, t) = 0$ ,

$$\mathbf{u}_{n+1} = \mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+) \mathbf{u}_n, \quad \mathbf{u}_n := (u(x_0, t^n), \dots, u(x_{N-1}, t^n))^T, \quad n = 0, 1, 2, \dots, \quad (2.8)$$

which encodes the fully discrete finite difference scheme

$$\begin{cases} \frac{u(x_\nu, t^{n+1}) - u(x_\nu, t^n)}{\Delta t} = a \frac{u(x_{\nu+1}, t^n) - u(x_\nu, t^n)}{\Delta x}, & \nu = 0, 1, \dots, N-1, \\ u(x_N, t^{n+1}) = 0. \end{cases} \quad (2.9)$$

The computation proceeds with hyperbolic scaling of fixed mesh ratio,  $\Delta t / \Delta x$ . This is precisely the regime  $N \sim n$  indicated in (2.3), in which case it is known that the forward Euler scheme (2.9) is *unstable*, if it violates the CFL condition  $0 < a\Delta t / \Delta x < 1$ . Observe that  $\mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+)$  amounts to a Jordan block,

$$\mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+) = \mathbb{I} + \Delta t \cdot a \mathbb{D}_N^+ = \mathbb{J}_q, \quad q = a\delta - 1, \quad \delta := \frac{\Delta t}{\Delta x}.$$

Therefore, the instability of  $\mathbb{J}_q$  with  $q \in (0, 1]$  follows, corresponding to  $1 < a\delta < 2$ , which was already claimed by the bound (2.4). In particular, the RK1 scheme (2.8) is unstable, despite having  $|\lambda_k(\mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+))| = |q| < 1$ .

Now consider integration of (2.7) using 4-stage (RK4). Spectral stability analysis

$$|\lambda_k(\mathcal{P}_4(\Delta t \cdot a \mathbb{D}_N^+))| = |\mathcal{P}_4(-a\delta)| \leq 1,$$

leads to the CFL condition,  $0 < a\delta \leq R_4 = 2\sqrt{2}$ , which *fails* to guarantee stability, since it does not capture the power-growth of the increasingly large Jordan block  $a\delta \mathbb{D}_N^+$ . We conclude that even in the most favorable scenario (2.1), spectral analysis is not enough to secure a uniform-in- $N$  stability of RK methods for increasingly large systems.

## 2.2 | Resolvent stability is not enough

Recall that the imaginary interval condition (1.13) is necessary for the stability of RK method. Kreiss and Wu [36, Theorem 3.6] proved that the converse holds in the sense of *resolvent stability*. Here, resolvent stability is interpreted in the sense that there exists a constant  $K_R > 0$  independent of  $N$ , such that for all negative  $\mathbb{L}_N$ 's, if the time step is small enough,  $\Delta t \cdot \|\mathbb{L}_N\| \leq C_s$ , then the corresponding  $s$ -stage RK method satisfies

$$\|(z\mathbb{I} - \mathcal{P}_s(\Delta t \mathbb{L}_N))^{-1}\| \leq \frac{K_R}{|z| - 1}, \quad \forall |z| > 1. \quad (2.10)$$

The size of the time step is dictated by region of absolute stability,  $\mathcal{A}_s$ , specifically  $-C_s \leq R_s$  is the radius of largest half disc inscribed inside  $\mathcal{A}_s$ ,

$$B_{C_s}^-(0) := \{z : \operatorname{Re} z < 0, |z| < C_s\} \subset \mathcal{A}_s, \quad \mathcal{A}_s = \{z \in \mathbb{C} : |\mathcal{P}_s(z)| \leq 1\}.$$

The notion of stability in the sense of power-boundedness, (1.8), implies that the resolvent condition holds with  $K_R = K_{\perp}$ . The Kreiss Matrix Theorem [30, 51, §4.9], states that the converse holds for a family of matrices with a *fixed* dimension. Yet this does not enable us to conclude the uniform-in- $N$  power-boundedness stability of RK method sought in (1.12), since the resolvent bound (2.10) may still allow growth  $\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \lesssim NK_R$ . In [62] we conjectured that this linear dependence on  $N$  is the best possible. This was confirmed in [42] proving that

$$\sup_{A \in M_N(\mathbb{C})} \frac{\sup_{|z| > 1} (|z| - 1) \|(z\mathbb{I} - A)^{-1}\|}{\sup_{n \geq 1} \|A^n\|} \sim eN.$$

The above linear-growth-in- $N$  behavior was exhibited by a sequence of increasingly large  $N \times N$  Jordan blocks,  $A_N = N\mathbb{J}_0$ . We observe that the  $A_N$ 's in this case are not resolvent bounded uniformly in  $N$ ; it is only the ratio on the left that exhibits the sharp linear bound in  $N$ . A concrete example of a family of matrices in  $M_N(\mathbb{C})$  which are resolvent stable yet their powers admit logarithmic growth in  $N$  was constructed in [44] and was improved to linear growth in  $N$  [27].

*Remark 2.1* (Dissipative resolvent condition). In [63] we considered a stronger resolvent condition of the form

$$\|(z\mathbb{I} - \mathcal{P}_s(\Delta t \mathbb{L}_N))^{-1}\| \leq \frac{K_R}{|z - 1|}, \quad \forall \{z : |z| \geq 1, z \neq 1\}. \quad (2.11)$$

In [52] it was proved that (2.11) implies  $n^{-1} \|\mathcal{P}_s^n\| \xrightarrow{n \rightarrow \infty} 0$ . In [63] we stated the improved logarithmic bound  $\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \lesssim \log(n)$ ; this was proved in [66]. More on the dissipative resolvent (2.11) and related notions of stability can be found in [53, 67, 68]. The dissipative resolvent bound (2.11) reflects a flavor of coercivity condition which will be visited in §3.3 below; however, it does not secure uniform-in- $N$  power-boundedness. A more precise notion of a dissipative resolvent condition of order  $2r > 0$  requires the existence of  $\eta_r > 0$  such that

$$\|(z\mathbb{I} - \mathcal{P}_s(\Delta t \mathbb{L}_N))^{-1}\| \leq \frac{K_R}{\text{dist}\{z, \Omega_r\}} \quad \forall z \notin \Omega_r, \quad := \{w : |w| + \eta_r |w - 1|^{2r} \leq 1\}. \quad (2.12)$$

The resolvent bound (2.12) reflects the classical notion of “dissipativity of order  $2r$ ” due to Kreiss [31]. It remains an open question whether (2.12) implies uniform-in- $N$  power-boundedness.

### 2.3 | Strong stability is not enough

The contractivity stated in (1.16),  $\|\mathcal{T}_N \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathcal{T}_N^{-1}\| \leq 1$  with uniformly bounded  $\|\mathcal{T}_N^{-1}\| \cdot \|\mathcal{T}_N\| \leq K_{\mathcal{T}}$ , is equivalent to strong stability in the sense that there exist uniformly positive definite symmetrizer  $\mathcal{H}_N$  and  $K_{\mathcal{H}} > 0$ , such that

$$\mathcal{P}_s^{\top}(\Delta t \mathbb{L}_N) \mathcal{H}_N \mathcal{P}_s(\Delta t \mathbb{L}_N) \leq \mathcal{H}_N, \quad 0 < \frac{1}{K_{\mathcal{H}}} \leq \mathcal{H}_N \leq K_{\mathcal{H}}. \quad (2.13)$$

Just set  $\mathcal{H}_N = \mathcal{T}_N^* \mathcal{T}_N$  with uniformly bounded  $K_{\mathcal{H}} = K_{\mathcal{T}}$ . In other words, (2.13) tells us that<sup>8</sup>

$$\|\mathcal{P}_s(\Delta t \mathbb{L}_N)\|_{\mathcal{H}_N} \leq 1, \quad \Delta t \cdot \|\mathbb{L}_N\| \leq C_s. \tag{2.14}$$

This coincides with the usual notion of strong stability,<sup>9</sup> for example [48, 65]. It follows that a strongly stable RK scheme,  $\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N)\mathbf{u}_n$ , satisfies

$$|\mathbf{u}_{n+1}|_{\mathcal{H}_N} = |\mathcal{P}_s(\Delta t \mathbb{L}_N)\mathbf{u}_n|_{\mathcal{H}_N} \leq |\mathbf{u}_n|_{\mathcal{H}_N} \leq \dots \leq |\mathbf{u}_0|_{\mathcal{H}_N},$$

and hence the RK iterations satisfy the uniform-in- $N$  stability bound,  $|\mathbf{u}(t_n)|_{\ell^2} \leq K_{\mathcal{H}} |\mathbf{u}_0|_{\ell^2}$ . The strong stability of the 3-stage RK method (RK3) with symmetrizer  $\mathcal{H}_N = \mathbb{H}_N$  and  $C_3 = 1$ , was proved in [65] and was later extended in [61, Theorem 4.2] to all  $s$ -stage RK methods of order  $s = 3[\text{mod}4]$ , namely—for small enough time step,  $\Delta t \cdot \|\mathbb{L}_N\| \leq C_s$ , there holds,

$$\|\mathcal{P}_s(\Delta t \mathbb{L}_N)\|_{\mathbb{H}_N} \leq 1, \quad \mathcal{P}_s(z) = \sum_{k=0}^s \frac{z^k}{k!}, \quad s = 3[\text{mod}4]. \tag{2.15}$$

As mentioned above, this line of arguing stability by construction of the strong stability symmetrizer, fails to extend to  $s$ -stage RK methods with  $s \in 4\mathbb{N}$ , [1, 49, 60]. But this does not mean that the latter RK methods are necessarily unstable. Indeed, the general question whether stable methods are necessarily strongly stable was addressed in [13]—they are not. It leaves open the possibility that the question stability can be pursued by other approaches—other than strong stability. This will be addressed in the next section.

### 3 | NUMERICAL RANGE AND STABILITY OF COERCIVE RUNGE—KUTTA SCHEMES

#### 3.1 | Numerical range

We let  $\ell^2_H(\mathbb{C}^N)$  denote the weighted Euclidean space associated with a given positive definite matrix  $H > 0$ , and equipped with

$$\langle \mathbf{x}, \mathbf{y} \rangle_H := \mathbf{x}^* H \mathbf{y}, \quad |\mathbf{x}|_H^2 := \langle \mathbf{x}, H \mathbf{x} \rangle, \quad H > 0.$$

Let  $A \in M_N(\mathbb{C})$  be an  $N \times N$  matrix with possibly complex-valued entries. The  $H$ -weighted numerical range,  $W_H(A)$ , is the set in the complex plane

$$W_H(A) := \{ \langle A \mathbf{x}, \mathbf{x} \rangle_H : \mathbf{x} \in \mathbb{C}^N, |\mathbf{x}|_H = 1 \}.$$

In the case of the standard Euclidean framework corresponding to  $H = \mathbb{I}$ , we drop the subscript  $\mathbb{H} = \mathbb{I}$  and remain with the usual  $|\cdot|_{\ell^2} = \langle \cdot, \cdot \rangle$ , and the corresponding numerical range denoted

<sup>8</sup> We let  $|\cdot|_H$  denote the weighted norm,  $|\mathbf{w}|_H^2 = \langle \mathbf{w}, H \mathbf{w} \rangle$ , and  $\|\cdot\|_H$  denote the corresponding induced matrix norm,  $\|\mathcal{P}\|_H := \max_{\mathbf{w} \neq 0} |\mathcal{P} \mathbf{w}|_H / |\mathbf{w}|_H$ .

<sup>9</sup> also called monotonicity in the literature on Runge–Kutta methods.

$W(A)$ . If  $A$  is real symmetric then  $W(A)$  is an interval on the real line (and conversely—if  $W(A)$  is a real interval then  $A$  is symmetric [29, Problem 3.9]); if  $A$  is skew-symmetric then  $W(A)$  is an interval on the imaginary line. For general  $A$ 's, the Hausdorff–Toeplitz theorem asserts that  $W(A)$  is a convex set in  $\mathbb{C}$ . As an example, we compute the numerical range of the  $N \times N$  translation matrix,  $J_0$ ,

$$J_0 := \begin{bmatrix} 0 & 1 & \dots & \dots & 0 \\ 0 & 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & 0 & 1 \\ 0 & \dots & \dots & \dots & 0 \end{bmatrix}_{N \times N}. \tag{3.1}$$

For any unit vector  $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$  we set a new unit vector  $x_j(\xi) := e^{ij\xi} x_j$  to find

$$\langle J_0 \mathbf{x}(\xi), \mathbf{x}(\xi) \rangle = \sum_{j=1}^{N-1} x_{j+1}(\xi) \overline{x_j(\xi)} = e^{i\xi} \langle J_0 \mathbf{x}, \mathbf{x} \rangle, \quad x_j(\xi) := e^{ij\xi} x_j,$$

which proves that  $W(J_0)$  is a disc centered at the origin,  $B_\rho(0)$ ; its radius,  $\rho = \rho_N$ , is found by considering the eigenvalues  $\lambda_k(Re J_0) = \cos(\frac{k\pi}{N+1})$ ,  $k = 1, 2, \dots, N$ : since for any  $A$ ,  $Re W(A) = W(Re A)$ , we find  $\rho_N = \lambda_1(Re J_0) = \cos(\frac{\pi}{N+1})$ , and we conclude that  $W(J)$  is the disc  $B_{\rho_N}(0)$ ,

$$W(J_0) = \{z : |z| \leq \rho_N\}, \quad \rho_N = \cos\left(\frac{\pi}{N+1}\right). \tag{3.2}$$

### 3.2 | The numerical radius

The numerical radius of  $A \in M_N(\mathbb{C})$  is given by

$$r_H(A) := \max_{\|\mathbf{x}\|_H=1} |\langle A\mathbf{x}, \mathbf{x} \rangle_H|.$$

The role of the numerical radius in addressing the question of stability was pioneered in the celebrated work of Lax and Wendroff, [40], in which they proved the stability of their 2D Lax–Wendroff (LW) scheme, that is, power-boundedness of a family amplification matrices,  $\|G^n\| \leq Const.$ , by securing  $r(G) \leq 1$ . The original proof, by induction on  $N$  (!), was later replaced by Halmos inequality, [24, 47]

$$r(G^n) \leq r^n(G). \tag{3.3}$$

Note that although the numerical radius is not sub-multiplicative, that is — although  $r(AB) \leq r(A)r(B)$  may fail for general  $A, B \in M_N(\mathbb{R})$  [15], Halmos' inequality states that it holds whenever  $A = B$ .

Since for all  $A$ 's there holds  $\|A\| \leq 2r(A)$ , (3.3) immediately yields the stability asserted by Lax and Wendroff

$$r(G) \leq 1 \rightsquigarrow \|G^n\| \leq 2r(G^n) \leq 2, \tag{3.4}$$

and more important for our purpose—power-boundedness is secured uniformly in  $N$ . It is straightforward to extend these arguments to the weighted framework [62, §3]

$$r_{\mathbb{H}}(G^n) \leq r_{\mathbb{H}}^n(G), \text{ and therefore } r_{\mathbb{H}}(G) \leq 1 \rightsquigarrow \|G^n\| \leq 2K_{\mathbb{H}}, \quad 0 < K_{\mathbb{H}}^{-1} \leq \mathbb{H} \leq K_{\mathbb{H}}. \quad (3.5)$$

*Remark 3.1.* H.-O. Kreiss proved the LW stability (3.4) by linking it to a (strict) resolvent condition

$$r(A) \leq 1 \rightsquigarrow \|(z\mathbb{I} - A)^{-1}\| \leq \frac{1}{|z| - 1}, \quad \forall |z| > 1$$

and conversely, the numerical range is the smallest set  $S = W(A)$ , which induces the strict resolvent condition [56],

$$\|(z\mathbb{I} - A)^{-1}\| \leq \frac{1}{\text{dist}(z, S)}, \quad \forall z \in S^c.$$

### 3.3 | Coercivity and RK stability

We turn to verify the stability of the 1-stage forward Euler scheme (RK1),

$$\mathbf{u}_{n+1} = (\mathbb{I} + \Delta t \mathbb{L}_N) \mathbf{u}_n.$$

There are two regions of interest in the complex plane that we need to consider: the weighted numerical range,  $W_{\mathbb{H}_N}(\mathbb{L}_N)$ , and the region of absolute stability associated with forward Euler,  $\mathcal{A}_1 = \{z : |1 + z| \leq 1\}$ . We make the assumption that the time step  $\Delta t$  is small enough so that

$$\Delta t W_{\mathbb{H}_N}(\mathbb{L}_N) \subset \mathcal{A}_1, \quad \mathcal{A}_1 = \{z : |1 + z| \leq 1\}, \quad (3.6)$$

then

$$\begin{aligned} r_{\mathbb{H}_N}(\mathcal{P}_1(\Delta t \mathbb{L}_N)) &= \max_{|\mathbf{x}|_{\mathbb{H}_N} = 1} |1 + \langle \Delta t \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N}| \\ &= \max_{z \in \Delta t W_{\mathbb{H}_N}(\mathbb{L}_N)} |1 + z| \leq \max_{z \in \mathcal{A}_1} |\mathcal{P}_1(z)| = 1. \end{aligned} \quad (3.7)$$

We summarize by stating the following.

**Theorem 3.2** (Numerical range stability of RK1). *Consider the forward Euler scheme associated with 1-stage forward Euler method (RK1),*

$$\mathbf{u}_{n+1} = (\mathbb{I} + \Delta t \mathbb{L}_N) \mathbf{u}_n, \quad n = 0, 1, 2, \dots,$$

*and assume the CFL condition (3.6) holds. Then the scheme is stable, and the following stability bound holds*

$$|\mathbf{u}_n|_{\ell^2} \leq 2K_{\mathbb{H}} |\mathbf{u}_0|_{\ell^2}, \quad \forall n \geq 1.$$

**Example 3.1.** As an example for theorem 3.2 we consider the one-sided differences (2.7),

$$\Delta t \cdot a \mathbb{D}_N^+ = a \frac{\Delta t}{\Delta x} \begin{bmatrix} -1 & 1 & \dots & \dots & 0 \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & -1 & 1 \\ 0 & \dots & \dots & \dots & -1 \end{bmatrix} = a\delta(-\mathbb{1} + \mathbb{J}_0), \quad (3.8)$$

$$a > 0, \quad \delta = \frac{\Delta t}{\Delta x}.$$

By translation and dilation,  $W(\Delta t \cdot a \mathbb{D}_N^+) = a\delta(-1 \oplus W(\mathbb{J}_0))$ , where (3.2) tells us that  $W(\mathbb{J}_0)$  is the ball  $B_{\rho_N}(0)$ . Hence  $W(\Delta t \cdot a \mathbb{D}_N^+)$  is given by the shifted ball,

$$W(\Delta t \cdot a \mathbb{D}_N^+) = \left\{ z : |z + a\delta| \leq a\delta\rho_N \right\}, \quad \delta = \frac{\Delta t}{\Delta x}, \quad \rho_N = \cos\left(\frac{\pi}{N+1}\right). \quad (3.9)$$

In particular,  $W(\Delta t \cdot a \mathbb{D}_N^+) \subset B_1(-1)$  uniformly in  $N$  if and only if the CFL condition  $a\delta \leq 1$  holds, which in turn secures the stability of the 1-stage forward Euler method, (RK1), for one-sided the transport equation(2.7),  $\mathbf{u}_{n+1} = (\mathbb{1} + \Delta t \cdot a \mathbb{D}_N^+)\mathbf{u}_n$ .

**Corollary 3.3** (Stability of forward Euler scheme). *Consider the forward Euler scheme (2.9) associated with 1-stage RK method (RK1),*

$$\mathbf{u}_{n+1} = \mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+)\mathbf{u}_n, \quad n = 0, 1, 2, \dots$$

*The scheme is stable under the CFL condition,  $0 < a\delta \leq 1$ , and the following stability bound holds  $\|\mathbf{u}_n\|_{\ell^2} \leq 2\|\mathbf{u}_0\|_{\ell^2}$ ,  $\forall n \geq 1$ .*

The last corollary can be recast in terms of a stability statement for  $\mathbb{J}_q = \mathcal{P}_1(\Delta t \cdot a \mathbb{D}_N^+)$ ,

$$\|\mathbb{J}_q^2\| \leq 2, \quad q \in (-1, 0).$$

This complements the statement of instability of  $\mathbb{J}_q$  in the range  $q \in (0, 1]$ , discussed in §2.1.1.

We note that the stability of  $\mathbb{J}_q$ ,  $q \in [-1, 0)$  can be independently verified by its induced  $\ell^1$ -norm,

$$\|\mathbb{J}_q\|_{\ell^1} = |-q| + |1+q| = 1 \rightsquigarrow \|\mathbb{J}_q^2\|_{\ell^1} \leq 1, \quad q \in [-1, 0). \quad (3.10)$$

However, the  $\ell^2$ -stability  $\|\mathbb{J}_q\|_{\ell^2} \leq 2$  stated in corollary 3.3 and the  $\ell^1$ -stability (3.10) are not equivalent uniformly in  $N$ . Also,  $\mathbb{J}_q$  is subject to  $\ell^2$  von-Neumann stability analysis [51, §4.7]

$$\max_{\varphi} |-q + (1+q)e^{i\varphi}| = 1, \quad q \in [-1, 0).$$

However, since the underlying problem (2.9) is not periodic, von Neumann stability analysis may not suffice: it requires the normal mode analysis [32] to prove  $\ell^2$ -stability. Thus, the numerical



range argument summarized in corollary 3.3 offers a genuinely different approach of addressing the question of stability, at least for 1-stage RK1.

*Remark 3.4 (Coercivity).* The CFL restriction encoded in (3.6),  $|\langle \Delta t \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N} + 1| \leq 1$ , leads to the sub-class of negative  $\mathbb{L}_N$ 's which satisfy the *coercivity bound*

$$2\text{Re} \langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N} \leq -\beta |\langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N}|^2, \quad \forall \mathbf{x} \in \{\mathbb{C}^N : |\mathbf{x}|_{\mathbb{H}_N} = 1\}. \tag{3.11}$$

Indeed, if  $\mathbb{L}_N$  is  $\beta$ -coercive in the sense that (3.11) holds with  $\beta > 0$ , then (3.6) is satisfied under the CFL condition  $\Delta t \leq \beta$ , and stability follows,  $r_{\mathbb{H}_N}(\mathbb{I} + \Delta t \mathbb{L}_N) \leq 1$ . We note that (3.11) places a weaker coercivity condition than the stronger notion of coercivity introduced in [43]

$$\mathbb{L}_N^T \mathbb{H}_N + \mathbb{H}_N \mathbb{L}_N \leq -\beta \mathbb{L}_N^T \mathbb{H}_N \mathbb{L}_N, \quad \beta > 0. \tag{3.12}$$

Indeed, the latter implies (3.11), for

$$2\text{Re} \langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N} \leq -\beta \langle \mathbb{L}_N^T \mathbb{H}_N \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle = -\beta |\mathbb{L}_N \mathbf{x}|_{\mathbb{H}_N}^2 \leq -\beta |\langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N}|^2, \quad |\mathbf{x}|_{\mathbb{H}_N} = 1.$$

One can then revisit the coercivity-based examples for stable RK methods in [43] using the relaxed coercivity (3.11). The notion of  $\beta$ -coercivity is related to the dissipative resolvent condition (2.11) but we shall not dwell on this point in this work.

### 3.4 | Numerical range stability of SSP RKs

We extend theorem 3.2 to multi-stage RK methods using their SSP format [18, §3]. We demonstrate the first three cases of RKs,  $s = 2, 3, 4$ .

Assume that the numerical range stability (3.7) holds. For example, the CFL condition  $\Delta t \leq \beta$  for  $\beta$ -coercive  $\mathbb{L}_N$ 's, (3.11), implies  $r_{\mathbb{H}_N}(\mathbb{I} + \Delta t \mathbb{L}_N) \leq 1$ . Then, for the 2-stage RK method, (RK2), we have by Halmos inequality (3.3)

$$\begin{aligned} r_{\mathbb{H}_N}(\mathcal{P}_2(\Delta t \mathbb{L}_N)) &\leq 1/2 + 1/2 r_{\mathbb{H}_N}^2(\mathbb{I} + \Delta t \mathbb{L}_N) \leq 1/2 + 1/2 = 1, \\ \mathcal{P}_2(\Delta t \mathbb{L}_N) &\equiv 1/2 \mathbb{I} + 1/2 (\mathbb{I} + \Delta t \mathbb{L}_N)^2. \end{aligned}$$

Similarly, the 3-stage RK method (RK3) can be expressed as

$$\mathcal{P}_3(\Delta t \mathbb{L}_N) \equiv 1/3 \mathbb{I} + 1/2 (\mathbb{I} + \Delta t \mathbb{L}_N) + 1/6 (\mathbb{I} + \Delta t \mathbb{L}_N)^3,$$

and hence if (3.7) holds, then the stability of (RK3) follows from Halmos inequality,

$$\begin{aligned} r_{\mathbb{H}_N}(\mathcal{P}_3(\Delta t \mathbb{L}_N)) &\leq 1/3 + 1/2 r_{\mathbb{H}_N}(\mathbb{I} + \Delta t \mathbb{L}_N) + 1/6 r_{\mathbb{H}_N}^3(\mathbb{I} + \Delta t \mathbb{L}_N) \\ &\leq 1/3 + 1/2 + 1/6 = 1. \end{aligned}$$

A similar argument applies to the 4-stage RK (RK4),

$$\mathcal{P}_4(\Delta t \mathbb{L}_N) \equiv 3/8 \mathbb{I} + 1/3(\mathbb{I} + \Delta t \mathbb{L}_N) + 1/4(\mathbb{I} + \Delta t \mathbb{L}_N)^2 + 1/24(\mathbb{I} + \Delta t \mathbb{L}_N)^4;$$

the numerical stability (3.7),  $r_{\mathbb{H}_N}(\mathbb{I} + \Delta t \mathbb{L}_N) \leq 1$  implies the stability of RK4,

$$\begin{aligned} r_{\mathbb{H}_N}(\mathcal{P}_4(\Delta t \mathbb{L}_N)) &\leq 3/8 + 1/3 r_{\mathbb{H}_N}(\mathbb{I} + \Delta t \mathbb{L}_N) + 1/4 r_{\mathbb{H}_N}^2(\mathbb{I} + \Delta t \mathbb{L}_N)^2 + 1/24 r_{\mathbb{H}_N}^4(\mathbb{I} + \Delta t \mathbb{L}_N) \\ &\leq 3/8 + 1/3 + 1/4 + 1/24 = 1. \end{aligned}$$

We summarize by stating

**Corollary 3.5** (Coercivity implies stability of RKs,  $s = 2, 3, 4$ ). *Consider the RK schemes*

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathbf{u}_n, \quad n = 0, 1, 2, \dots, \quad s = 2, 3, 4.$$

*Assume the numerical range stability (3.7) holds. In particular if  $\mathbb{L}_N$  is  $\beta$ -coercive in the sense of (3.11), and that the CFL condition,  $\Delta t \leq \beta$ , is satisfied. Then these  $s$ -stage RK schemes are stable,*

$$|\mathbf{u}(t_n)|_{\mathbb{H}_N} \leq 2|\mathbf{u}(0)|_{\mathbb{H}_N} \rightsquigarrow |\mathbf{u}(t_n)|_{\ell^2} \leq 2K_{\mathbb{H}}|\mathbf{u}(0)|_{\ell^2}.$$

The building block of corollary 3.5 is the condition of numerical range stability (3.7) originated with (RK1). While this argument is sharp for the 1-stage forward Euler, this SSP-based argument is too restrictive for multi-stage RKs. In particular, corollary 3.5 rules out the large sub-class of negative yet non-coercive  $\mathbb{L}_N$ 's, due to a numerical range which has non-trivial intersection with the imaginary axes. In particular, this includes the important sub-class of skew-symmetric (hyperbolic)  $\mathbb{L}_N$ 's with purely imaginary numerical range. For example, if the one-sided differences in (2.7) are replaced by centered-differences

$$\mathbf{u}_{n+1} = (\mathbb{I} + \Delta t \cdot a \mathbb{D}_N^0) \mathbf{u}_n, \quad \mathbb{D}_N^0 := \frac{1}{\Delta x} \begin{bmatrix} 0 & 1 & \dots & \dots & 0 \\ -1 & 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & -1 & 0 & 1 \\ 0 & \dots & \dots & -1 & 0 \end{bmatrix}_{N \times N}. \quad (3.13)$$

The numerical range lies on the imaginary interval,  $W(\Delta t \cdot a \mathbb{D}_N^0) = [-iR, iR]$  with  $R = R_N = a\delta \cos(\frac{\pi}{N+1})$ . The 1-stage forward Euler (3.13) fails to satisfy the imaginary interval condition, and therefore, corollary 3.5 fails to capture the stability of the corresponding RKs,  $\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \cdot a \mathbb{D}_N^0) \mathbf{u}_n$  for  $s = 3, 4$ .

## 4 | SPECTRAL SETS AND STABILITY OF RUNGE-KUTTA METHODS

We now turn our attention to the stability of multi-stage RK methods,  $\mathcal{P}_s(\Delta t \mathbb{L}_N)$ . Clearly, spectral analysis is not enough. On the other hand, direct computation based on  $\ell^1$  or  $\ell^2$ -von Neumann analysis is not accessible: even the entries in the example of one-sided differences,  $\mathcal{P}_s(\Delta t \cdot a \mathbb{D}_N^+)$ ,

for  $s = 3, 4$ , become excessively complicated to write down. Instead, we suggest to pursue a stability argument based on numerical radius along the lines of (3.7), starting with

$$r(\mathcal{P}_s(\Delta t \mathbb{L}_N)) = \max_{\substack{|\mathbf{x}|=1 \\ \mathbf{x} \in \mathbb{C}^N}} \left| \sum_{k=0}^s a_k \langle (\Delta t \mathbb{L}_N)^k \mathbf{x}, \mathbf{x} \rangle \right|.$$

This requires a proper functional calculus of numerical range, relating the sets  $W(\mathcal{P}_s(\Delta t \mathbb{L}_N))$  and  $\{|\mathcal{P}_s(z)|, : z \in W(\Delta t \mathbb{L}_N)\}$ , similar to the role of the spectral mapping theorem (1.9) as the centerpiece of spectral stability analysis. To this end we recall the notion of a  $K$ -spectral set developed in [8, 10], which dates back to von Neumann [46]; we refer to [54] for a most recent overview.

**Definition 4.1** ( $K$ -spectral sets). Given  $A \in M_N(\mathbb{C})$ , we say that  $\Omega \subset \mathbb{C}$  is a  $K$ -spectral set of  $A$  if there exists a finite  $K > 0$  such that for all analytic  $f$ 's bounded on  $\Omega$ , there holds

$$\|f(A)\|_H \leq K \max_{z \in \Omega} |f(z)|. \tag{4.1}$$

In a remarkable work [6], Crouzeix proved that for every matrix  $A$ , the numerical range  $W_H(A)$  is a  $K$ -spectral set of  $A$  with  $K \leq 11.08$ ; this was later improved to  $K = 1 + \sqrt{2}$  [7]. An elegant proof of Crouzeix and Palencia  $(1 + \sqrt{2})$ -bound [50] is included in an appendix. It follows, in particular, that for all polynomials  $p$ ,

$$\|p(A)\|_H \leq (1 + \sqrt{2}) \max_{z \in W_H(A)} |p(z)|. \tag{4.2}$$

**Theorem 4.2** (Stability of Runge–Kutta schemes). Consider the  $s$ -stage explicit RK method  $\mathcal{P}_s(z) = \sum_{k=0}^s a_k z^k$ , associated with region of absolute stability  $\mathcal{A}_s = \{z : |\mathcal{P}_s(z)| \leq 1\}$ . Then, the RK scheme

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathbf{u}_n, \quad n = 0, 1, 2, \dots$$

is stable under the CFL condition  $\Delta t W_{\mathbb{H}_N}(\mathbb{L}_N) \subset \mathcal{A}_s$ ,

$$\Delta t W_{\mathbb{H}_N}(\mathbb{L}_N) \subset \mathcal{A}_s \iff \|\mathbf{u}_n\|_{\ell^2} \leq (1 + \sqrt{2}) K_{\mathbb{H}} \|\mathbf{u}_0\|_{\ell^2}, \quad n = 1, 2, \dots \tag{4.3}$$

For proof we apply (4.2) with  $p = \mathcal{P}_s^n$ :

$$\begin{aligned} \|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\|_{\mathbb{H}_N} &\leq (1 + \sqrt{2}) \max_{z \in \Delta t W_{\mathbb{H}_N}(\mathbb{L}_N)} |\mathcal{P}_s^n(z)| \\ &\leq (1 + \sqrt{2}) \max_{z \in \mathcal{A}_s} |\mathcal{P}_s^n(z)| \leq 1 + \sqrt{2}, \end{aligned}$$

and hence  $\|\mathcal{P}_s^n(\Delta t \mathbb{L}_N)\| \leq (1 + \sqrt{2}) K_{\mathbb{H}}$ .

*Remark 4.3* (Implicit RK methods). The argument above makes a critical use of the striking fact that the spectral set bound,  $K = 1 + \sqrt{2}$ , is independent of neither the increasing degree,  $\deg(\mathcal{P}_s^n) = sn$ , nor of the increasingly large dimension,  $\dim(\mathbb{L}_N) = N$ . In fact, since (4.2) applies

to the larger algebra of *rational functions* bounded on  $W_H(A)$ , theorem 4.2 can be equally well formulated to general implicit RK methods [23, II.7].

We recall the spectral stability analysis (1.10) which is quantified in terms of the spectrum  $\sigma(\mathbb{L}_N)$

$$\Delta t \sigma(\mathbb{L}_N) \subset \mathcal{A}_S, \quad \sigma(A) := \{\lambda_k(A) : k = 1, 2, \dots, N\}.$$

In the terminology of (4.1), the spectrum  $\sigma(\mathbb{L}_N)$  is not a spectral set for  $\mathbb{L}_N$ . Theorem 4.2 tells us that replacing the spectrum with the larger set of  $H$ -weighted numerical range,  $W_{\mathbb{H}_N}(\mathbb{L}_N) \supset \sigma(\mathbb{L}_N)$ , provides a very general framework for the stability of any Runge–Kutta scheme, in conjunction with any  $\mathbb{L}_N$ . For example, the forward Euler (RK1) applies to the one-sided difference (3.8) which was covered in Corollary 3.3. Observe that for *normal* matrices,<sup>10</sup>  $\mathbb{L}_N$ , there holds  $\text{conv}\{\sigma(\mathbb{L}_N)\} = W(\mathbb{L}_N)$ , for example [25]. Thus, the gap  $W_{\mathbb{H}_N}(\mathbb{L}_N) \setminus \text{conv}\{\sigma(\mathbb{L}_N)\}$  comes into play in the stability statement (4.3) when normality uniform-in- $N$  fails — precisely the scenario described in §2.1 for failure of spectral analysis to secure stability (in this context we remark that since  $\cap_{H>0} W_H(\mathbb{L}_N) = \text{conv}\{\sigma(\mathbb{L}_N)\}$ , it is essential to restrict attention to uniformly bounded symmetrizers  $\mathbb{H}_N$  in (1.4)).

A main drawback of the CFL condition (4.3) is its formulation in terms of a weighted numerical range which is not always easily accessible. Here comes the imaginary interval condition, (1.13), which provides an accessible sufficient condition for stability of multi-stage RK methods.

**Theorem 4.4** (Stability of Runge–Kutta methods). *Consider the  $s$ -stage explicit RK method and assume it satisfies the imaginary interval condition (1.13), namely—there exists  $R_s > 0$  such that*

$$\max_{-R_s \leq \sigma \leq R_s} |\mathcal{P}_s(i\sigma)| \leq 1, \quad \mathcal{P}_s(z) = 1 + z + a_2 z^2 + \dots + a_s z^s. \tag{4.4}$$

Then, there exists a constant  $0 < C_s < R_s$  such that for all negative  $\mathbb{L}_N$ 's (1.4), the RK method

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathbf{u}_n, \quad n = 0, 1, 2, \dots,$$

is stable under the CFL condition  $\Delta t \cdot r_{\mathbb{H}_N}(\mathbb{L}_N) \leq C_s$ ,

$$\Delta t \cdot r_{\mathbb{H}_N}(\mathbb{L}_N) \leq C_s \iff \|\mathbf{u}_n\|_{\ell^2} \leq (1 + \sqrt{2}) K_{\mathbb{H}} \|\mathbf{u}_0\|_{\ell^2}, \quad n = 1, 2, \dots \tag{4.5}$$

*Proof.* Recall  $B_\alpha^-$  denotes the semi-disc,  $B_\alpha^- := \{z : \text{Re } z \leq 0, |z| \leq \alpha\}$ . Consider an arbitrary negative  $\mathbb{L}_N$ ,

$$2\text{Re} \langle \mathbb{L}_N \mathbf{x}, \mathbf{x} \rangle_{\mathbb{H}_N} = \langle (\mathbb{L}_N^\top \mathbb{H}_N + \mathbb{H}_N \mathbb{L}_N) \mathbf{x}, \mathbf{x} \rangle \leq 0$$

The negativity of  $\mathbb{L}_N$  states that the weighted numerical range  $W_{\mathbb{H}_N}(\mathbb{L}_N)$  lies on the left side of complex plane, and in fact, inside the left semi-disc

$$W_{\mathbb{H}_N}(\mathbb{L}_N) \subset B_{r_{\mathbb{H}_N}(\mathbb{L}_N)}^- := \{z : \text{Re } z \leq 0, |z| \leq r_{\mathbb{H}_N}(\mathbb{L}_N)\}.$$

<sup>10</sup>  $\mathbb{L}_N^* \mathbb{L}_N = \mathbb{L}_N \mathbb{L}_N^*$  where  $\mathbb{L}_N^*$  is the  $\ell^2$ -adjoint of  $\mathbb{L}_N$ .

Next, we make use of [35, Theorem 3.2] which asserts<sup>11</sup> that for an  $s$ -stage RK method satisfying the imaginary interval condition, its region of absolute stability contains a non-trivial semi-disc  $B_{C_s}^-$  with  $C_s \leq R_s$ , so that

$$\mathcal{A}_s \supset B_{C_s}^- := \{z : \operatorname{Re} z \leq 0, |z| \leq C_s\}, \quad C_s \leq R_s. \tag{4.6}$$

We conclude that for small step-size (4.5)

$$\Delta t W_{\mathbb{H}_N}(\mathbb{L}_N) \subset \Delta t B_{r_{\mathbb{H}_N}(\mathbb{L}_N)}^- = B_{\Delta t \cdot r_{\mathbb{H}_N}(\mathbb{L}_N)}^- \subset B_{C_s}^- \subset \mathcal{A}_s.$$

Theorem 4.2 implies stability (1.6) with  $K_{\mathbb{L}} = (1 + \sqrt{2})K_{\mathbb{H}}$ . □

*Remark 4.5.* We note that theorem 4.4 makes use of the semi-disc  $B_{C_s}^-$  as a spectral set for  $\mathcal{P}_s(\Delta t \mathbb{L}_N)$ . In this case, one expects a sharper bound, compared with (4.2) [54, §3.2],  $\|p(A)\|_H \leq 2 \max_{z \in W_H(A)} |p(z)|$ . The constant 2—corresponding to (3.4) with  $p(z) = z^n$ , agrees with *Crouzeix’s conjecture* [6] regarding the optimality of the numerical range as 2-spectral set.

### 4.1 | Optimality of the numerical radius-based CFL condition

We observe that the CFL condition quoted in (4.5),

$$\Delta t \cdot r_{\mathbb{H}_N}(\mathbb{L}_N) \leq C_s, \tag{4.7}$$

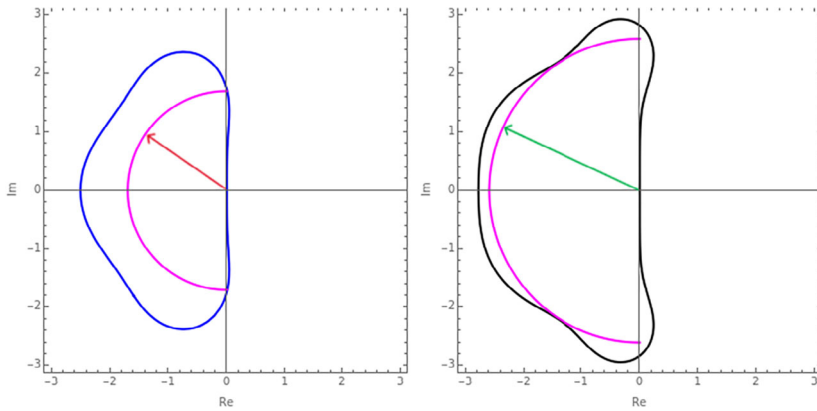
offers a refinement of the CFL condition (1.7). Indeed, since  $\mathbb{H}_N$  is uniformly bounded  $0 < K_{\mathbb{H}}^{-1} \leq \mathbb{H}_N \leq K_{\mathbb{H}}$ , we have

$$r_{\mathbb{H}_N}(\mathbb{L}_N) \leq \|\mathbb{L}_N\|_{\mathbb{H}_N} \leq K_{\mathbb{H}} \|\mathbb{L}_N\|,$$

and hence, the CFL condition—compare with (1.7),  $\Delta t \cdot \|\mathbb{L}_N\| \leq C'_s$  with  $C'_s := C_s/K_{\mathbb{H}}$ , implies that (4.5) holds, and stability follows.

In fact, we claim that (4.7) offers an optimal CFL condition in the following sense. The proof of theorem 4.4 compares two semi-discs: on one hand we identified  $B_{C_s}^-$  as the largest semi-disc inscribed inside  $\mathcal{A}_s$  (this is a property of the RK method under consideration); on the other hand, we identified  $B_{r_{\mathbb{H}_N}(\mathbb{L}_N)}^-$  as the smallest semi-disc which contains  $W_{\mathbb{H}_N}(\mathbb{L}_N)$ . The CFL condition (4.7) secures the dilation of the latter semi-disc inside the former, and there, we seek the smallest semi-disc associated with  $\mathbb{L}_N$  which satisfies a set of desired requirements. We claim that we cannot find a smaller semi-disc which will secure this line of argument. Indeed, let  $\|\cdot\|$  denote an arbitrary (vector) norm on  $M_N(\mathbb{C})$ , with a semi-disc  $B_{\|\mathbb{L}_N\|}^-$  which would be a candidate for a better CFL condition, that is, an even smaller semi-disc  $B_{\|\mathbb{L}_N\|}^- \subset B_{r_{\mathbb{H}_N}(\mathbb{L}_N)}^-$ . Clearly, by the necessity encoded in (1.10), the CFL condition requires that  $\|A\|$  is *spectrally dominant* in the sense that  $\|A\| \geq |\lambda_{\max}(A)|$  for all  $A \in M_N(\mathbb{C})$ . Moreover, since power-boundedness is invariant under unitary transformations,  $\|(UAU^*)^n\|_{\mathbb{H}_N} = \|A^n\|_{\mathbb{H}_N}$ , we ask that the semi-disc associated with  $\|\cdot\|$

<sup>11</sup> Note that this requires  $\mathcal{P}_s(0) = \mathcal{P}'_s(0) = 1$  in (4.4).



**FIGURE 2** The semi-circles  $B_{C_3}^-(0)$  inscribed inside  $\mathcal{A}_3$  (left) and  $\mathcal{A}_4$  (right),

be unitarily invariant,

$$UB_{\|\mathbb{L}_N\|}^- U^* = B_{\|\mathbb{L}_N\|}^- \text{ for all } U\text{'s such that } |U\mathbf{x}|_{\mathbb{H}_N} = |\mathbf{x}|_{\mathbb{H}_N}.$$

It follows from the main theorem of [14] that the semi-disc  $B_{\|\mathbb{L}_N\|}^-$  must contain  $B_{r_{\mathbb{H}_N}(\mathbb{L}_N)}^-$ . That is, the corresponding CFL condition (4.7) is optimal in the sense that it is the smallest, spectrally dominant, unitarily invariant semi-disc which makes the argument of theorem 4.4 work.

A main aspect of theorem 4.4 is going beyond any specific coercivity requirement which was sought in the SSP-based arguments in §3.4. It applies to *all* negative  $\mathbb{L}_N$ 's, thus addressing the question sought in [43, §3.5]. A precise characterization for RK methods satisfying the imaginary interval condition was given in [35, Theorem 3.1]. Consider an explicit  $s$  stage RK method, accurate of order  $r \geq 1$ ,

$$P_s(z) = \sum_{k=0}^r \frac{z^k}{k!} + \sum_{k=r+1}^s a_k z^k, \quad r \geq 1. \tag{4.8a}$$

It satisfies the imaginary interval condition (1.13) if and only if

$$\begin{cases} (-1)^{\frac{r+1}{2}} (a_{r+1} - 1) < 0, & r \text{ is odd,} \\ (-1)^{\frac{r+2}{2}} (a_{r+2} - (r+2)a_{r+1} + r+1) < 0, & r \text{ is even.} \end{cases} \tag{4.8b}$$

In the particular case of  $s = r = 3, 4$  we find that the 3-stage RK method (RK3) and 4-stage RK method RK4 satisfy the imaginary interval condition and hence the existence of semi-discs with radii  $C_3 = \sqrt{3}$  and  $C_4 = 2.61$ , shown in Figure 2 which imply stability under the respective CFL conditions,

$$\Delta t \cdot \|\mathbb{L}_N\| \leq C'_s, \quad C'_s = C_s/K_{\mathbb{H}}.$$

In particular, this extends the strong stability statement of 3-stage (RK3) in [65, Theorem 2] and provides a stability proof for the 4-stage RK (RK4) for arbitrarily large systems with negative  $\mathbb{L}_N$ 's.

Comparing the RK4 stability requirement,  $C_4 = 2.61$ , vs. the RK4 interval condition mentioned earlier,  $R_4 = 2\sqrt{2}$ , reflects the stricter stability requirement associated with the larger numerical range  $W_{\mathbb{H}_N}(\mathbb{L}_N) \supset \sigma(\mathbb{L}_N)$ .

Condition (4.8b) becomes more restrictive for higher order methods; instead, one can increase  $r$  and use  $s$ -stage protocol,  $s > r$  to form a dissipative term  $\sum_{k=r+1}^s a_k z^k$  which enforces the imaginary interval condition. In particular, the 7-stage Dormand-Prince method [11], with embedded fourth- and fifth-order accurate RK45,  $(r, s) = (5, 7)$  which is used in MATLAB, does satisfy the imaginary interval condition (4.8b) [55]. See the example of the 10-stage explicit RK method SSPRK(10,4) in [49, Fig. 2].

## 5 | EXAMPLES: STABILITY OF TIME-DEPENDENT METHODS OF LINES

We demonstrate application of the new stability results for arbitrarily large systems in the context of methods of lines for difference approximation of the scalar hyperbolic equation

$$y_t = a(x)y_x, \quad (t, x) \in \mathbb{R}_+ \times [0, 1],$$

augmented with proper boundary conditions. The stability results extend, mutatis mutandis,<sup>12</sup> to multi-dimensional hyperbolic problems,  $\mathbf{y}_t = \sum_{j=1}^d A_j(x)\mathbf{y}_{x_j}$ . Stability theories for such difference approximations were developed in the classical works in the 50s–70s, for example [20, 31, 32, 37, 38, 40, 51] and can be found in the more recent texts of [21, 26, 41]. Our aim here is to revisit the question of stability for RK time-discretizations of such difference approximations, from a perspective of the stability theory developed in §4. A central part of this approach requires computation of the (weighted) numerical range of the large matrices that arise in the context of such difference approximations. The development of full stability theory along these lines is beyond the scope of this paper, and is left for future work.

### 5.1 | Periodic problems. Constant coefficients

We consider the 1-periodic problem

$$\begin{cases} y_t(x, t) = ay_x(x, t), & (t, x) \in \mathbb{R}_+ \times [0, 1] \\ y(0, t) = y(1, t). \end{cases}$$

Its spatial part is discretized using finite-difference method with constant coefficients (depending on  $a$ ),  $\{q_\alpha\}$ , and acting on a discrete grid,  $x_\nu = \nu\Delta x$ ,  $\Delta x = 1/N$ ,

$$\frac{d}{dt}y(x_\nu, t) = Q(E)y(x_\nu, t), \quad \nu = 0, 1, \dots, N - 1, \quad Q(E) := \frac{1}{\Delta x} \sum_{\alpha=-\ell}^r q_\alpha E^\alpha.$$

<sup>12</sup> In particular,  $\ell^2$ -stability needs to be adjusted to weighted  $H$ -stability, weighted by the smooth symmetrizer  $H = H(x, \xi)$  so that  $H(x, \xi) \sum_j A_j(x)e^{ij\xi}$  is symmetric.

Here  $E$  is the 1-periodic translation operator,  $Ey_\nu = y_{(\nu+1)[\text{mod}N]}$ . The resulting scheme amounts to a system of ODEs for the  $N$ -vector of unknowns,  $\mathbf{y}(t) = (y(x_0, t), \dots, y(x_{N-1}, t))^\top$ , which admits the *circulant* matrix representation

$$\dot{\mathbf{y}}(t) = Q(\mathbb{E}_N)\mathbf{y}, \quad Q(\mathbb{E}_N) = \frac{1}{\Delta x} \sum_{\alpha=-\ell}^r q_\alpha \mathbb{E}^\alpha, \tag{5.1}$$

$$\mathbb{E}_N := \begin{bmatrix} 0 & 1 & \dots & \dots & 0 \\ 0 & 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & 0 & 1 \\ 1 & \dots & \dots & 0 & 0 \end{bmatrix}_{N \times N}.$$

The numerical range of circulant matrices is given by convex polytopes. Indeed, let  $\mathbb{F}$  denote the unitary Fourier matrix,  $\mathbb{F}_{jk} = \left\{ \frac{1}{\sqrt{N}} e^{2\pi ijk/N} \right\}_{j,k=1}^N$ . Then  $\mathbb{F}$  diagonalizes  $\mathbb{E}_N$ ,

$$\langle \mathbb{E}_N \mathbf{x}, \mathbf{x} \rangle = \langle \widehat{\mathbb{E}}_N \widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle,$$

$$\widehat{\mathbb{E}}_N := \mathbb{F}^* \mathbb{E}_N \mathbb{F} = \begin{bmatrix} e^{\frac{2\pi i}{N}} & 0 & \dots & \dots & 0 \\ 0 & e^{2\frac{2\pi i}{N}} & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \ddots & e^{(N-1)\frac{2\pi i}{N}} & 0 \\ 0 & \dots & \dots & \dots & 1 \end{bmatrix}, \quad \widehat{\mathbf{x}} := \mathbb{F}^* \mathbf{x}.$$

and hence  $W(\mathbb{E}_N) = \left\{ \sum_{j=1}^N |\widehat{x}_j|^2 e^{2\pi i j/N} : \sum_j |\widehat{x}_j|^2 = 1 \right\}$  is the regular  $N$ -polytope with vertices at  $\{e^{2\pi i j/N}\}_{j=1}^N$ . This should be compared with the numerical range of the Jordan block (3.2).

It follows that  $\widehat{Q(\mathbb{E}_N)} = Q(\widehat{\mathbb{E}}_N)$  and hence the action of the  $N \times N$  circulant  $Q(\mathbb{E}_N)$  is encoded in terms of its *symbol*,  $\widehat{q}(\xi) := \frac{1}{\Delta x} \sum_\alpha q_\alpha e^{i\alpha\xi}$ ,

$$\langle Q(\mathbb{E}_N) \mathbf{x}, \mathbf{x} \rangle = \langle \widehat{Q(\mathbb{E}_N)} \widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle,$$

$$\widehat{Q(\mathbb{E}_N)} = \begin{bmatrix} \widehat{q}\left(\frac{2\pi}{N}\right) & 0 & \dots & \dots & 0 \\ 0 & \widehat{q}\left(2\frac{2\pi}{N}\right) & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \widehat{q}\left((N-1)\frac{2\pi}{N}\right) & 0 \\ 0 & \dots & \dots & 0 & \widehat{q}(2\pi) \end{bmatrix}.$$

**Lemma 5.1** (Numerical range of circulant matrices). *The numerical range of the circulant matrix  $Q(\mathbb{E}_N)$  is given by the convex polytope with vertices at  $\{\widehat{q}(2\pi j/N)\}_{j=1}^N$ ,*

$$W(Q(\mathbb{E}_N)) = \left\{ \sum_j |\widehat{x}_j|^2 \widehat{q}\left(\frac{2\pi j}{N}\right) : |\widehat{\mathbf{x}}| = 1 \right\}.$$



We now appeal to theorem 3.2 which secures the stability of forward Euler time discretization for  $\mathbb{L}_N = Q(\mathbb{E}_N)$ , provided the CFL condition  $\Delta t W(Q(\mathbb{E}_N)) \subset B_1(-1)$  holds.

**Proposition 5.2** (Stability—difference schemes with constant coefficients. I). *Consider the fully-discrete finite difference scheme*

$$\mathbf{u}_{n+1} = \mathbf{u}_n + \frac{\Delta t}{\Delta x} \sum_{\alpha} q_{\alpha} \mathbb{E}_N^{\alpha} \mathbf{u}_n, \quad n = 0, 1, 2, \dots$$

The scheme is stable under the CFL condition,

$$\max_{1 \leq j \leq N} \left| 1 + \Delta t \cdot \hat{q}(2\pi j/N) \right| \leq 1, \quad \hat{q}(\xi) := \frac{1}{\Delta x} \sum_{\alpha} q_{\alpha} e^{i\alpha \xi}, \tag{5.2}$$

and the following stability bound holds  $\|\mathbf{u}_n\|_{\ell^2} \leq 2\|\mathbf{u}_0\|_{\ell^2}, \forall n \geq 1$ .

Since the CFL condition (5.2) guarantees that  $\mathbb{L}_N = \mathbb{I} + \Delta t \cdot Q(\mathbb{E}_N)$  is coercive, the result goes over to SSP-based multi-stage RK time differencing. In fact, theorem 4.4 applies for multi-stage RK time differencing and for all negative  $Q(\mathbb{E}_N)$ 's.

**Proposition 5.3** (Stability—difference schemes with constant coefficients. II). *Consider the fully-discrete finite difference scheme*

$$\begin{aligned} \mathbf{u}_{n+1} &= \mathcal{P}_s(\Delta t \cdot Q(\mathbb{E}_N)) \mathbf{u}_n, \quad n = 0, 1, 2, \dots, \\ \mathcal{P}_s(z) &= \sum_{k=0}^s a_k z^k, \quad Q(\mathbb{E}_N) = \frac{1}{\Delta x} \sum_{\alpha} q_{\alpha} \mathbb{E}_N^{\alpha}. \end{aligned} \tag{5.3}$$

Here,  $\mathcal{P}_s$  is an  $s$ -stage RK stencil satisfying the imaginary interval condition, so that (4.6) holds with  $C_s > 0$ . If the spatial discretization is negative,  $\text{Re} \hat{q}(2\pi j/N) \leq 0$ , then the scheme (5.3) is stable under the CFL condition

$$\max_{1 \leq j \leq N} |\Delta t \cdot \hat{q}(2\pi j/N)| \leq C_s, \quad \hat{q}(\xi) = Q(e^{i\xi}), \tag{5.4}$$

and the following stability bound holds,

$$\|\mathbf{u}_n\|_{\ell^2} \leq (1 + \sqrt{2})\|\mathbf{u}_0\|_{\ell^2}, \quad n = 1, 2, \dots$$

Propositions 5.2 and 5.3 recover von-Neumann stability analysis for difference schemes with constant coefficients [21, §4.2]. We shall consider three examples.

**Example 5.1** (One-sided differences). Consider the periodic setup of the one-sided difference (2.7),

$$\mathbf{u}_{n+1} = (\mathbb{I} + \Delta t \cdot Q(\mathbb{E}_N)) \mathbf{u}_n, \quad Q(\mathbb{E}_N) = \frac{a}{\Delta x} (\mathbb{E}_N - \mathbb{I}),$$

with spatial symbol  $\hat{q}(\xi) = \frac{a}{\Delta x}(e^{i\xi} - 1)$ . This amounts to the  $N \times N$  system

$$\mathbf{u}_{n+1} = \mathbb{L}_N \mathbf{u}_n,$$

$$\mathbb{L}_N = \begin{bmatrix} 1-\delta a & \delta a & 0 & \dots & 0 \\ 0 & 1-\delta a & \delta a & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & \ddots & 1-\delta a & \delta a \\ \delta a & 0 & \dots & 0 & 1-\delta a \end{bmatrix}_{N \times N}, \quad \delta = \frac{\Delta t}{\Delta x}.$$

Using proposition 5.2 we secure stability under the usual CFL condition  $\delta a \leq 1$ ,

$$\delta a \leq 1 \rightsquigarrow \max_{1 \leq j \leq N} |1 + \delta a(e^{2\pi i j/N} - 1)|^2 = |1 - \delta a|^2 + 2|1 - \delta a|\delta a + (\delta a)^2 \leq 1.$$

This extends to multi-stage time differencing, RKs,  $s = 3, 4$

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \cdot Q(\mathbb{E}_N))\mathbf{u}_n, \quad \mathcal{P}_s(z) = \sum_{k=0}^s \frac{z^k}{k!}, \quad s = 3, 4.$$

Clearly,  $Re \hat{q} \leq 0$ , and we can appeal to proposition 5.3 which secures stability under CFL condition  $\delta a \leq C_s$ ; indeed,

$$\delta a \leq C_s \rightsquigarrow \delta a(e^{2\pi i j/N} - 1) \in B_{C_s}^-, \quad j = 1, 2, \dots, N.$$

**Example 5.2** (Centered differences). Consider the periodic setup of the centered spatial difference scheme (3.13), combined with multi-stage RK time differencing, RKs,  $s = 3, 4$ ,

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \cdot Q(\mathbb{E}_N))\mathbf{u}_n, \quad Q(\mathbb{E}_N) = \frac{a}{2\Delta x}(\mathbb{E}_N - \mathbb{E}_N^{-1}).$$

Spatial differencing has purely imaginary symbol  $\hat{q}(\xi) = \frac{a}{\Delta x}i \sin(\xi)$ , and we invoke proposition 5.3 which secures stability under the CFL condition (5.4),

$$\delta a = \max_{1 \leq j \leq N} |\delta a i \sin(2\pi i j/N)| \leq C_s, \quad \delta = \frac{\Delta t}{\Delta x}.$$

This line of argument extends to higher order centered differences [65, §5.2], for example the fourth-order difference

$$Q(\mathbb{E}_N) = \frac{a}{12\Delta x}(-\mathbb{E}_N^2 + 8\mathbb{E}_N - 8\mathbb{E}_N^{-1} + \mathbb{E}_N^{-2})$$

or the fourth-order finite-element difference

$$Q(\mathbb{E}_N) = \begin{bmatrix} 4/6 & 1/6 & 0 & \dots & 1/6 \\ 1/6 & 4/6 & 1/6 & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 4/6 & 1/6 \\ 1/6 & 0 & \dots & 1/6 & 4/6 \end{bmatrix}^{-1} \times \frac{1}{2\Delta x} \begin{bmatrix} 0 & 1 & \dots & \dots & -1 \\ -1 & 0 & 1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 & 1 \\ 1 & 0 & \dots & -1 & 0 \end{bmatrix}_{N \times N}.$$

**Example 5.3** (LW differencing). We use the LW protocol for second-order spatial difference [40] (observe that the mesh ratio,  $\delta = \Delta t/\Delta x$ , is kept fixed),

$$Q_{LW}(\mathbb{E}_N) = \frac{a}{2\Delta x}(\mathbb{E}_N - \mathbb{E}_N^{-1}) + \frac{\delta a^2}{2\Delta x}(\mathbb{E}_N - 2\mathbb{I} + \mathbb{E}_N^{-1}),$$

with symbol

$$\hat{q}_{LW}(\xi) = \frac{a}{\Delta x}i \sin(\xi) + \frac{\delta a^2}{\Delta x}(\cos(\xi) - 1).$$

Stability of the LW scheme

$$\mathbf{u}_{n+1} = (\mathbb{I} + \Delta t Q_{LW}(\mathbb{E}_N))\mathbf{u}_n \tag{5.5}$$

follows provided CFL condition (5.2) holds, namely,  $|1 + \Delta t \hat{q}_{LW}(2\pi j/N)| \leq 1$ . Noting that

$$\hat{q}_{LW}(\xi) = \frac{2a}{\Delta x}i \sin(\xi/2) \cos(\xi/2) - \frac{2\delta a^2}{\Delta x} \sin^2(\xi/2),$$

it is a standard argument, for example [21, §1.2] to conclude that  $\delta a \leq 1$  secures the desired CFL condition,

$$\delta a \leq 1 \rightsquigarrow \max_{1 \leq j \leq N} |1 + \Delta t \hat{q}_{LW}(2\pi j/N)|^2 \leq 1.$$

We note that LW differencing has a negative symbol  $Re \hat{q}_{LW}(\xi) \leq 0$ , and therefore theorem 4.4 secures the stability of higher-order time discretizations of LW scheme

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t Q_{LW}(\mathbb{E}_N))\mathbf{u}_n, \quad s = 3, 4,$$

under the relaxed CFL condition,  $2\delta a \leq C_s$ . Indeed,

$$2\delta a \leq C_s \rightsquigarrow \max_{1 \leq j \leq N} |\Delta t \hat{q}_{LW}(2\pi j/N)|^2 \leq \max_{\xi} \{4(\delta a \sin(\xi/2) \cos(\xi/2))^2 + 4(\delta a \sin(\xi/2))^4\} \leq C_s^2.$$

The constant coefficient case in the period setup involves the algebra of circulant matrices, all of which are uniformly diagonalizable by the Fourier matrix  $F$ . This is a rather special case, in which von Neumann spectral stability analysis prevails for arbitrarily large systems. Clearly, the numerical range-based stability results of Sections 3 and 4 offer a more general framework for studying stability of general non-periodic cases. Examples are outlined below.

## 5.2 | Periodic problems. Variable coefficients

We consider the 1-periodic problem with  $C^2$ -variable coefficient  $a(\cdot)$

$$\begin{cases} y_t(x, t) = a(x)y_x(x, t), & (t, x) \in \mathbb{R}_+ \times [0, 1] \\ y(0, t) = y(1, t). \end{cases} \tag{5.6}$$

The spatial part is discretized using finite-difference method with  $a(x)$ -dependent variable coefficients,  $\{q_\alpha(x)\}$ , and acting on a discrete grid,  $x_\nu = \nu\Delta x$ ,  $\Delta x = 1/N$ ,

$$\frac{d}{dt}y(x_\nu, t) = Q(E)y(x_\nu, t), \quad \nu = 0, 1, \dots, N-1, \quad (5.7)$$

$$Q(E) := \frac{1}{\Delta x} \sum_{\alpha=-\ell}^r q_\alpha(x)E^\alpha.$$

The accuracy requirement places the restriction  $\sum_\alpha q_\alpha(x) = 0$ ,  $\sum_\alpha \alpha q_\alpha(x) = a(x)$  and so on. The difference scheme (5.7) amounts to an  $N \times N$  system of ODEs with "slowly varying" circulantcy, that is,  $Q(x, \mathbb{E}_N)_{ij}$  changes smoothly in the sense that  $|Q(x, \mathbb{E}_N)_{i+1, j+1} - Q(x, \mathbb{E}_N)_{ij}|$  is bounded independent of  $1/\Delta x$ .

$$\Delta x \sum_\alpha \alpha^2 |q_\alpha(x)|_{C^2} \leq K_q. \quad (5.8)$$

Let  $\widehat{Q}$  denote the formal symbol associated with (5.7)

$$\widehat{Q}(x, \xi) := \frac{1}{\Delta x} \sum_{\alpha=-\ell}^r q_\alpha(x)e^{i\alpha\xi}.$$

Assume that the symbol is negative  $Re \widehat{Q}(x, \xi) \leq 0$ . Then by the sharp Gårding inequality [37, Theorem 1.1], see also [39], the corresponding difference operator is semi-bounded,<sup>13</sup> namely—there exists a constant  $\eta > 0$  depending on  $K_q$  but otherwise independent of  $N$ , such that

$$Re Q(x, \mathbb{E}_N) \leq 2\eta \mathbb{1}_{N \times N}. \quad (5.9)$$

Theorem 4.4 applies to  $Q(x, \mathbb{E}_N) - \eta \mathbb{1}$ , implying its power-boundedness under the CFL condition (1.7),

$$\|\mathcal{P}_s^n(\Delta t(Q(x, \mathbb{E}_N) - \eta \mathbb{1}))\| \leq 1 + \sqrt{2}, \quad \Delta t \cdot r(Q(x, \mathbb{E}_N)) \leq C_s.$$

Next, we note that the shift  $-\eta \mathbb{1}$  produces only a finite bounded perturbation  $B$ , namely

$$\begin{aligned} \mathcal{P}_s(\Delta t \cdot Q(x, \mathbb{E}_N)) &= \mathcal{P}_s(\Delta t \cdot (Q(x, \mathbb{E}_N) - \eta \mathbb{1}) + \Delta t \cdot \eta \mathbb{1}) \\ &= \mathcal{P}_s(\Delta t \cdot (Q(x, \mathbb{E}_N) - \eta \mathbb{1})) + \Delta t \cdot B, \\ B &= \eta \sum_{k=1}^s a_k k (\Delta t \cdot Q(x, \mathbb{E}_N))^{k-1}, \end{aligned}$$

where  $\|B\| \leq \eta K_B$  with  $K_B = \sum_{k=1}^s |a_k| k C_s^{k-1}$ . We now invoke the fact (due to [30, 59]) that bounded perturbations of power-bounded matrices remain power bounded,<sup>14</sup>

$$\|A^n\| \leq K_A \rightsquigarrow (A + \Delta t \cdot B)^n \leq K_A e^{K_A \|B\| t_n}, \quad t_n = n\Delta t.$$

<sup>13</sup> Note that  $Q(x, \mathbb{E}_N)$  is unbounded,  $\|Q(x, \mathbb{E}_N)\| = \mathcal{O}(1/\Delta x)$ .

<sup>14</sup> This follows from the identity  $(X + Y)^n \equiv X^n + \sum_{k=0}^{n-1} X^{n-k-1} Y (X + Y)^k$ ,  $n = 1, 2, \dots$  and using induction with  $(X, Y) = (A, \Delta t \cdot B)$ .

This implies the desired stability bound

$$|\mathbf{u}(t_n)| \leq \| \mathcal{P}_s^n(\Delta t \cdot Q(x, \mathbb{E}_N)) \| \cdot |\mathbf{u}_0| \leq (1 + \sqrt{2})e^{(1+\sqrt{2})\eta K_B t_n} |\mathbf{u}_0|.$$

We summarize by stating

**Proposition 5.4** (Stability—finite difference schemes with variable coefficients). *Consider the fully-discrete finite difference scheme*

$$\mathbf{u}_{n+1} = \mathcal{P}_s(Q(x, \mathbb{E}_N))\mathbf{u}_n, \quad n = 0, 1, 2, \dots, \tag{5.10}$$

where  $Q(x, \mathbb{E}_N) = \frac{1}{\Delta x} \sum_{\alpha} q_{\alpha}(x) \mathbb{E}_N^{\alpha}$  is a local difference operator, (5.8), and  $\mathcal{P}_s$  is an  $s$ -stage RK stencil satisfying the imaginary interval condition, (4.6). If the spatial symbol is negative,

$$\operatorname{Re} \hat{Q}(x, \xi) \leq 0, \quad \hat{Q}(x, \xi) := \frac{1}{\Delta x} \sum_{\alpha} q_{\alpha}(x) e^{i\alpha\xi}, \tag{5.11}$$

then the scheme (5.10) is stable under the CFL condition

$$\max_{\xi} |\Delta t \cdot \hat{Q}(x, \xi)| \leq C_s, \tag{5.12}$$

and the following stability bound holds with  $K_B := \sum_{k=1}^s |a_k| k C_s^{k-1}$ ,

$$|\mathbf{u}_n|_{\ell^2} \leq (1 + \sqrt{2})e^{(1+\sqrt{2})\eta K_B t_n} |\mathbf{u}_0|_{\ell^2}, \quad n = 1, 2, \dots, \quad \Delta t \cdot r(Q(x, \mathbb{E}_N)) \leq C_s.$$

*Remark 5.5.* The stability analysis of difference schemes with variable coefficients in [31, 37] bounds the norm of  $\| \mathcal{P}_s(\Delta t Q(x, \mathbb{E}_N)) \| \leq 1 + \mathcal{O}(\Delta t)$ . However, the result is limited to one-step forward difference in time,  $\mathbb{I} + \Delta t Q(x, \mathbb{E}_N)$ . The essence of proposition 5.4 is extension to RK time-differentiating of higher orders  $s \geq 1$ .

**Stability of Fourier method.** There are two approaches to handle the stability of difference approximations of problems with variable coefficients: the von-Neumann spectral analysis based on sharp Gårding inequality (5.9), or the energy method for example [64, §2]; both approaches requires local stencils (5.8). An alternative approach for stability with variable coefficients in based on numerical dissipation [31]. As an extreme example for using our RK stability result, we consider the Fourier method [33, §4], [17], which is neither local nor dissipative. Set  $\Delta x = 1/(2N+1)$  with an odd number of  $(2N + 1)$  gridpoints. The Fourier method for (5.6) amounts to  $(2N + 1) \times (2N + 1)$  system of ODEs

$$\begin{aligned} \dot{\mathbf{y}}(t) &= Q(\mathbb{D}_N^{\mathbb{F}})\mathbf{y}(t), \\ Q(\mathbb{D}_N^{\mathbb{F}}) &= A\mathbb{D}_N^{\mathbb{F}}, \quad A = \begin{bmatrix} a(x_0) & & & \\ & a(x_1) & & \\ & & \ddots & \\ & & & a(x_{2N}) \end{bmatrix}, \end{aligned} \tag{5.13}$$

where the diagonal matrix  $A$  encodes  $a(x)$  and  $\mathbb{D}_N^F$  is the  $(2N+1) \times (2N+1)$  Fourier differencing matrix

$$\mathbb{D}_N^F = \mathbb{F} \begin{bmatrix} -iN & 0 & \dots & 0 \\ 0 & -i(N-1) & 0 & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & i(N-1) & 0 \\ 0 & \dots & \dots & 0 & iN \end{bmatrix} \mathbb{F}^*, \quad \mathbb{F}_{jk} = \left\{ \frac{e^{ijk\Delta x}}{\sqrt{2N+1}} \right\}_{j,k=1}^{2N+1}.$$

The Fourier difference method is neither local,  $(\mathbb{D}_N^F)_{jk} = \frac{(-1)^{j-k}}{2 \sin((k-j)\Delta x/2)}$  fails (5.8), nor dissipative, and the method is unstable in presence of variable coefficients [16]. However, there is a different weighted-stability. Specifically — for the prototypical case  $a(x) = \sin(x)$ , there exists a symmetrizer  $\mathbb{H}_N$  such that [16, Theorem 2.1]

$$Q(\mathbb{D}_N^F)^T \mathbb{H}_N + \mathbb{H}_N Q(\mathbb{D}_N^F) \leq \mathbb{H}_N,$$

where the  $\mathbb{H}_N$ -norm corresponds to the  $H^1$ -norm

$$|\mathbf{u}|_{\mathbb{H}_N}^2 = |\mathbf{u}|_{H^1}^2, \quad |\mathbf{u}|_{H^s}^2 := \sum_{k=-N}^N (1 + k^2)^{\frac{s}{2}} |\hat{u}_k|^2.$$

**Proposition 5.6** (Stability—Fourier method). *Consider the time discretization of the Fourier method,*

$$\dot{\mathbf{y}}(t) = Q(\mathbb{D}_N^F)\mathbf{y}(t),$$

$$Q(\mathbb{D}_N^F) = A\mathbb{D}_N^F, \quad A = \begin{bmatrix} \sin(x_0) & & & \\ & \sin(x_1) & & \\ & & \ddots & \\ & & & \sin(x_{2N}) \end{bmatrix},$$

using RK methods which satisfy the imaginary interval condition,

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \cdot Q(\mathbb{D}_N^F))\mathbf{u}_n, \quad n = 1, 2, \dots, \quad \Delta t \cdot N \leq C_s.$$

The Fourier method is  $H^1$ -stable

$$|\mathbf{u}_n|_{\mathbb{H}_N} \leq (1 + \sqrt{2})e^{tn/2} |\mathbf{u}_0|_{\mathbb{H}_N}. \tag{5.14}$$

We note that the symmetrizer  $\mathbb{H}_N$  is not uniformly bounded from below,  $N^{-2}\mathbb{1} \leq \mathbb{H}_N \leq 4\mathbb{1}$ , so  $\ell^2$ -stability fails. Converted to  $\ell^2$ -framework, (5.14) yields

$$|\mathbf{u}_n|_{\ell^2} \leq N|\mathbf{u}_n|_{\mathbb{H}_N} \leq N(1 + \sqrt{2})e^{tn/2} |\mathbf{u}_0|_{\mathbb{H}_N} = 2N(1 + \sqrt{2})e^{tn/2} |\mathbf{u}_0|_{\ell^2}.$$

### 5.3 | Initial-boundary value problems

We consider the problem (2.5) in the strip

$$\begin{cases} y_t(x, t) = ay_x(x, t), & a > 0, & (t, x) \in \mathbb{R}_+ \times [0, 1] \\ y(1, t) = 0. \end{cases}$$

A general stability theory for difference approximations of initial-boundary value problems was developed in [20, 32]. It is based on normal mode analysis and secures the resolvent-type stability of such approximations. The following example shows how to utilize the framework offered in theorem 4.2, to study the stability of difference approximations of initial-boundary value problems.

**Example 5.4** (One-sided difference). Consider an interior centered differencing augmented with one-sided difference at the outflow boundary  $x = 0$ ,

$$\begin{cases} \frac{d}{dt}y(x_0, t) = a \frac{y(x_1, t) - y(x_0, t)}{\Delta x} \\ \frac{d}{dt}y(x_\nu, t) = a \frac{y(x_{\nu+1}, t) - y(x_{\nu-1}, t)}{2\Delta x}, & \nu = 1, 2, \dots, N - 1 \\ y(x_N, t) = 0. \end{cases} \tag{5.15}$$

We emphasize that we treat the semi-infinite problem, which amounts to method of lines for the infinite-vector of unknowns,  $\mathbf{y}(t) := (y(x_0, t), y(x_1, t), \dots, y(x_{N-1}, t))^T$ , governed by the semi-discrete system

$$\dot{\mathbf{y}}(t) = \mathbb{L}_N \mathbf{y}(t), \quad \mathbb{L}_N = \frac{a}{\Delta x} \begin{bmatrix} -1 & 1 & 0 & \dots & \dots & 0 \\ -1/2 & 0 & 1/2 & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & & \vdots \\ \vdots & \ddots & -1/2 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & -1/2 & 0 & 1/2 \\ 0 & \dots & \dots & 0 & -1/2 & 0 \end{bmatrix}. \tag{5.16}$$

Although the matrix  $\mathbb{L}_N$  is not negative,  $\mathbb{L}_N^T + \mathbb{L}_N = \frac{a}{\Delta x} \begin{bmatrix} -2 & 1/2 \\ 1/2 & 0 \end{bmatrix} \oplus 0_{(N-2) \times (N-2)}$ , it is weighted negative with the simple symmetrizer  $\mathbb{H}_N$ :

$$\mathbb{L}_N^T \mathbb{H}_N + \mathbb{H}_N \mathbb{L}_N = \frac{a}{\Delta x} \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix} \oplus 0_{(N-2) \times (N-2)} \leq 0, \quad \mathbb{H}_N := \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix} \oplus \mathbb{I}_{(N-2) \times (N-2)}.$$

Using theorem 4.4, we conclude the stability of time discretization of (5.16) using any RK method satisfying the imaginary interval condition, (4.8). In particular, the fully-discrete schemes based

on the  $s$ -stage RK time discretization

$$\mathbf{u}_{n+1} = \mathcal{P}_s(\Delta t \mathbb{L}_N) \mathbf{u}_n, \quad s = 3, 4, \quad n = 1, 2, \dots,$$

are stable under the CFL condition  $\Delta t \cdot r_{\mathbb{H}_N}(\mathbb{L}_N) \leq C_s$ ,

$$|\mathbf{u}(t_n)| \leq 4(1 + \sqrt{2})|\mathbf{u}_0|.$$

Observing the simple bound,  $r_{\mathbb{H}_N}(\mathbb{L}_N) \leq \frac{a}{\Delta x} K_{\mathbb{H}}$  with  $K_{\mathbb{H}} = 2$ , we end with CFL condition sufficient for stability,  $\delta a \leq C_s/2$ .

The last example depends on verifying weighted negativity,  $\mathbb{L}_N^\top \mathbb{H}_N + \mathbb{H}_N \mathbb{L}_N \leq 0$ , which requires the construction of a proper symmetrizer on a case by case basis. A systematic approach for studying the weighted negativity for properly designed boundary treatment augmenting centered difference schemes was developed in [2, 19, 34, 58]. To extend our RK stability framework to larger classes of difference approximations of initial-boundary values problems requires a more precise characterization of the *weighted* numerical range of Teoplitz-like spatial discretizations. This is left for future study.

## ACKNOWLEDGMENTS

Research was supported in part by ONR grants N00014-2112773 and N00014-2412659 and the Fondation Sciences Mathématiques de Paris (FSMP) while being hosted by LJLL at Sorbonne University.

## REFERENCES

1. F. Achleitner, A. Arnold, and A. Jüngel, Necessary and sufficient conditions for strong stability of explicit Runge-Kutta methods. In: *From Particle Systems to Partial Differential Equations* (Carlen, E., Gonalves, P., Soares, A.J., eds), *Proceedings in Mathematics & Statistics*, vol 465, pp. 1–21, Springer, 2024.
2. E. Burman, A. Ern, and M. A. Fernández, *Explicit Runge-Kutta Schemes and Finite Elements with Symmetric Stabilization for First-Order Linear PDE Systems*, *SIAM J. Numer. Anal.* **48** (2010), no. 6, 2019–2042.
3. J. C. Butcher, *Numerical Methods for Ordinary Differential Equations*, John Wiley & Sons, 2008. ISBN 978-0-470-72335-7.
4. R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, *Neural Ordinary Differential Equations*, *Adv. Neural Inf. Process. Syst.* **31** (2018). doi: [10.1007/978-3-030-04167-0](https://doi.org/10.1007/978-3-030-04167-0)
5. R. Courant, K. Friedrichs, and H. Lewy, *On the partial difference equations of mathematical physics*, *Math. Ann.* **100** (1928), 32–74, doi: [10.1147/rd.112.0215](https://doi.org/10.1147/rd.112.0215)
6. M. Crouzeix, *Numerical range and functional calculus in Hilbert space*, *J. Funct. Anal.* **244** (2007), 668–690.
7. M. Crouzeix and C. Palencia, *The numerical range is a  $(1 + \sqrt{2})$ -spectral set*, *SIAM J. Matrix Anal. Appl.* **38** (2017), 649–655.
8. M. Crouzeix and A. Greenbaum, *Spectral sets: Numerical range and beyond*, *SIAM J. Matrix Anal. Appl.* **40** (2019), no. 3, 1087–1101.
9. G. Dahlquist and Å. Björck, *Numerical Methods*, Prentice-Hall Series in Automatic Computation, 1973.
10. B. Delyon and F. Delyon, *Generalization of Von Neumann's spectral sets and integral representation of operators*, *Bull. Soc. Math. France* **127** (1999), 25–41.
11. J. R. Dormand and P. J. Prince, *A family of embedded Runge-Kutta formulae*, *J. Computational and Appl. Math.* **6** (1980), no. 1, 19–26.
12. W. E, *A proposal on machine learning via dynamical systems*, *Commun. Math. Stat.* **5** (2017), 1–11.
13. S. R. Foguel, *A counterexample to a problem of Sz-Nagy*, *Proc. AMS* **15** (1964), 788–790.
14. S. Friedland and E. Tadmor, *Optimality of the Lax-Wendroff condition*, *Linear Algebra Appl.* **56** (1984), 121–129.



15. M. Goldberg and E. Tadmor, *On the numerical radius and its applications*, Linear Algebra Appl. **42** (1982), 263–284.
16. J. Goodman, T. Hou, and E. Tadmor, *On the stability of the unsmoothed Fourier method for hyperbolic equations*, Numer. Math **67** (1994), no. 1, 93–129.
17. D. Gottlieb and S. A. Orszag, *Numerical analysis of spectral methods: Theory and applications*, SIAM, Philadelphia, 1977.
18. S. Gottlieb, C.-W. Shu, and E. Tadmor, *Strong stability preserving high order time discretization methods*, SIAM Rev. **43** (2001), 89–112.
19. B. Gustafsson, *On the implementation of boundary conditions for the method of lines*, BIT Numer. Math. **38** (1998), 293–314.
20. B. Gustafsson, H.-O. Kreiss, and A. Sundström, *Stability theory of difference approximations for mixed initial boundary value problems. II*, Math. Comp. **26** (1972), no. 119, 649–686.
21. B. Gustafsson, H.-O. Kreiss, and J. Olinger, *Time dependent problems and difference methods*, Wiley, Hoboken, New Jersey, 2013.
22. E. Haber and L. Ruthotto, *Stable architectures for deep neural networks*, Inverse Probl. **34** (2017), no. 1.
23. E. Hairer, S. Nørsett, and G. Wanner, *Solving ordinary differential equations I: Nonstiff problems*, Springer-Verlag, Berlin, 1993. ISBN 978-3-540-56670-0.
24. P. R. Halmos, *A Hilbert space problem book*, Van Nostrand, New York, 1967.
25. P. Henrici, *Bounds for iterates, inverses, spectral variation and fields of values of non-normal matrices*, Numer. Math. **4** (1962), 24–40.
26. J. S. Hesthaven, *Numerical Methods for Conservation Laws: From Analysis to Algorithms*, Soc. Ind. Appl. Math. (SIAM) (2017).
27. K. J. In’T Hout and M. N. Spijker, *Analysis of error growth via stability regions in numerical initial value problems*, BIT Numer. Math. **43** (2003), 363–385.
28. A. Iserles, *A First Course in the Numerical Analysis of Differential Equations*, Cambridge University Press, Cambridge, UK, 1996. ISBN 978-0-521-55655-2.
29. T. Kato, *Perturbation Theory for Linear Operators*, Classic in Mathematics, Springer, Berlin, 1995.
30. H.-O. Kreiss, *Über die Stabilitätsdefinition für differenzgleichungen die partielle differential-gleichungen approximieren*, BIT Numer. Math. **2** (1962), 153–181.
31. H. Kreiss, *On difference approximations of the dissipative type for hyperbolic differential equations*, Comm. Pure Applied Math. **17** (1964), no. 3, 335–353.
32. H.-O. Kreiss, *Stability theory for difference approximations of mixed initial boundary value problems. I*, Math. Comp. **22** (1968), no. 104, 703–714.
33. H.-O. Kreiss and J. Olinger, *Comparison of accurate methods for the integration of hyperbolic equations*, Tellus **24** (1972), no. 3, 199–215.
34. H.-O. Kreiss and G. Scherer, *Finite element and finite difference methods for hyperbolic partial differential equations*, C. de Boor (ed.), Mathematical Aspects of Finite Elements in Partial Differential Equations, Academic Press, New York, 1974.
35. H.-O. Kreiss and G. Scherer, *Method of lines for hyperbolic differential equations*, SIAM J. Numer. Anal. **29** (1992), no. 3, 640–646.
36. H.-O. Kreiss and L. Wu, *On the stability definition of difference approximations for the initial boundary value problem*, Appl. Numerical Math. **12** (1993), no. 1-3, 213–227.
37. P. D. Lax and L. Nirenberg, *On stability for difference schemes; A sharp form of Gårding’s inequality*, Comm. Pure Applied Math. **19** (1966), no. 4, 473–492.
38. P. D. Lax and R. D. Richtmyer, *Survey of stability of linear finite difference equations*, CPAM **9** (1956), 267–293.
39. P. D. Lax and B. Wendroff, *On the stability of difference schemes*, Comm. Pure Appl. Math. **15** (1962), 363–371.
40. P. D. Lax and B. Wendroff, *Difference schemes with high order of accuracy for solving hyperbolic equations*, Comm. Pure Appl. Math. **17** (1964), no. 3, 381–392.
41. R. J. LeVeque, *Finite Difference Methods for Ordinary and Partial Differential eQuations: Steady-State and Time-Dependent Problems*, Society for Industrial and Applied Mathematics (2007).
42. R. LeVeque and N. Trefethen, *On the resolvent condition in the Kreiss Matrix Theorem*, BIT Numer. Math. **24** (1984), 584–591.

43. D. Levy and E. Tadmor, *From semi-discrete to fully-discrete: Stability of Runge-Kutta schemes by the energy method*, SIAM Rev. **40** (1998), 40–73.
44. C. A. McCarthy and J. Schwartz, *On the norm of a finite Boolean algebra of projections and applications to theorems of Kreiss and Morton*, Comm. Pure Appl. Math. **18** (1965), 191–201.
45. S. Mishra, *A machine learning framework for data driven acceleration of computations of differential equations*, Math. Engineering 2019 **1** (2018), no. 1, 118–146.
46. J. von Neumann, *Erne spektraltheorie für allgememe operatoren ernes unitären raumes*, Math. Nachr. **4** (1951), 258–281.
47. C. Pearcy, *An elementary proof of the power inequality for the numerical radius*, Mich. Math. J. **13** (1966), 289–291.
48. H. Ranocha, *On strong stability of explicit Runge-Kutta methods for nonlinear semibounded operators*, IMA J. Numerical Anal. **41** (2021), no. 1, 654–682.
49. H. Ranocha and P. Öffner,  *$L_2$  Stability of explicit Runge-Kutta schemes*, J. Sci. Comput. **75** (2018), 1040–1056.
50. T. Ransford and F. L. Schwenninger, *Remarks on the Crouzeix–Palencia proof that the numerical range is a  $(1 + \sqrt{2})$ -spectral set*, SIAM J. Matrix Anal. and Appl. **39** (2018), no. 1, 342–345.
51. R. D. Richtmyer and K. W. Morton, *Difference methods for initial-value problems*, 2nd ed., Interscience, New York, 1967.
52. R. K. Ritt, *A condition that  $\lim_{n \rightarrow \infty} n^{-1}T^n = 0$* , Proc. Amer. Math. Soc. **4** (1953), 898–899.
53. F. L. Schwenninger, *Functional calculus estimates for Tadmor-Ritt operators*, J. Math. Anal. Appl. **439** (2016), no. 1, 103–124.
54. F. L. Schwenninger and J. de Vries, *On abstract spectral constants*. In: Operator and Matrix Theory, Function Spaces, and Applications (Ptak, M., Woerdeman, H.J., Wojtylak, M., eds), Operator Theory: Advances and Applications, vol 295, Birkhuser (2024). doi: [https://doi.org/10.1007/978-3-031-50613-0\\_15](https://doi.org/10.1007/978-3-031-50613-0_15)
55. L. F. Shampine and M. W. Reichelt, *The MATLAB ode suite*, SISSC **18** (1997), no. 1, 1–22.
56. M. N. Spijker, *Numerical ranges and stability estimates*, Applied Numer. Math. **13** (1993), 241–249.
57. M. N. Spijker and F. A. J. Straetemans, *Error growth analysis via stability regions for discretizations of initial value problems*, BIT Numer. Math. **37** (1997), 442–464.
58. B. Strand, *Summation by parts for finite difference approximations for  $d/dx$* , J. Comp. Phys. **110** (1994), no. 1, 47–67.
59. G. Strang, *Accurate partial difference methods*, Numer. Math. **6** (1964), no. 1, 37–46.
60. Z. Sun and C.-W. Shu, *Stability of the fourth order Runge-Kutta method for time-dependent partial differential equations*, Ann. Math. Sci. Appl. **2** (2017), no. 2, 255–284.
61. Z. Sun and C. W. Shu, *Strong stability of explicit Runge-Kutta time discretizations*, SIAM J. Numer. Anal. **57** (2019), no. 3, 1158–1182.
62. E. Tadmor, *The equivalence of  $L^2$ -stability, the resolvent condition and strict  $H$ -stability*, Linear Algebra Appl. **41** (1981), 151–159.
63. E. Tadmor, *The resolvent condition and uniform power boundedness*, Linear Algebra Appl. **80** (1986), 250–252.
64. E. Tadmor, *Stability analysis of finite-difference, pseudospectral and Fourier-Galerkin approximations for time-dependent problems*, SIAM Rev. **29** (1987), 525–555.
65. E. Tadmor, *From semi-discrete to fully discrete: Stability of Runge-Kutta schemes by the energy method. II “Collected Lectures on the Preservation of Stability under Discretization”*, D. Estep and S. Tavener, (eds.), , Proceedings in Applied Mathematics, vol. 109, SIAM, 2002, pp. 25–49.
66. P. Vitse, *Functional calculus under the Tadmor-Ritt condition, and free interpolation by polynomials of a given degree*, J. Funct. Anal. **210** (2004), 43–72.
67. P. Vitse, *The Riesz turndown collar theorem giving an asymptotic estimate of the powers of an operator under the Ritt condition*, Rend. Circ. Mat. Palermo **53** (2004), 283–312.
68. P. Vitse, *A band limited and Besov class functional calculus for Tadmor-Ritt operators*, Arch. Math. **85** (2005), 374–385.
69. G. Wanner, *Kepler, Newton and numerical analysis*, Acta Numer. **19** (2010), 561–598.

**APPENDIX A: THE NUMERICAL RANGE IS  $(1 + \sqrt{2})$ -SPECTRAL SET**

In his remarkable work [6], Crouzeix proved that  $W_H(A)$  is a  $K$ -numerical set with  $K = 11.08$  which was later improved by Crouzeix & Palencia to  $K = 1 + \sqrt{2}$ . We quote here the elegant proof of Ransford & Schwenninger [50] for Crouzeix & Palencia  $(1 + \sqrt{2})$ -bound, based on the following lemma. In particular, we refer to the recent review [54].

**Lemma A.1** (Ransford & Schwenninger  $(1 + \sqrt{2})$ -spectral set). *Let  $T$  be a Hilbert space bounded operator  $\|T\| < \infty$ , and let  $\Omega$  be a bounded open set containing the spectrum of  $T$ . Suppose that for each  $f$  analytic on  $\Omega$ , there exists an analytic  $g$  on  $\Omega$  such that the following holds (here and below,  $\|f\|_\Omega := \sup_\Omega |f|$ ):*

$$\|g\|_\Omega \leq \|f\|_\Omega \text{ and } \|f(T) + g(T)^*\| \leq 2\|f\|_\Omega. \tag{A.1}$$

Then

$$\|f(T)\| \leq (1 + \sqrt{2})\|f\|_\Omega$$

*Proof.* Let  $K := \sup_{\|f\|_\Omega=1} \|f(T)\|$ . By assumption, for each  $f$ ,  $\|f\|_\Omega \leq 1$ , there exists  $g$  such that (A.1) holds. Ransford & Schwenninger invoked the identity

$$f(T)f(T)^*f(T)f(T)^* \equiv f(T)(f(T) + g(T)^*)^*f(T)f(T)^* - (fgf)(T)f(T)^*.$$

A simple exercise shows that the norm of the quantity on the left equals  $\|f(T)\|^4$ . Since by (A.1)<sub>1</sub>,  $\|(fgf)\|_\Omega \leq 1$  hence  $\|fgf(T)\| \leq K$ , and since by (A.1)<sub>2</sub>,  $\|f(T) + g(T)^*\| \leq 2$ , then the expression on the right does not exceed

$$\begin{aligned} \|f(T)\|^4 &= \|f(T)f(T)^*f(T)f(T)^*\| \\ &\leq \|f(T)\|\|f(T) + g(T)^*\|\|f(T)\|\|f(T)^*\| + \|(fgf)(T)\|\|f(T)^*\| \\ &\leq 2K^3 + K^2. \end{aligned}$$

Hence,  $K^4 = \sup_{\|f\|_\Omega=1} \|f(T)\|^4 \leq 2K^3 + K^2$  which implies  $K \leq 1 + \sqrt{2}$ . □

Note that the lemma does not involve the numerical range of  $T$  — this comes into play in the construction of  $g = g_\Omega$  satisfying (A.1), in terms of Cauchy transform,

$$g_\Omega(z) := \frac{1}{2\pi i} \int_{\partial\Omega} \frac{\overline{f(\zeta)}}{\zeta - z} d\zeta, \quad z \in \Omega.$$

The main thrust of the work, originated in [46] and then developed in [6, 10] and finally [7], is to show that such  $g_\Omega$  with  $\Omega = W_H(T)$  satisfies (A.1).