



**HAL**  
open science

# Third-order A-stable alternating implicit Runge-Kutta schemes

Alexandre Ern, Jean-Luc Guermond

► **To cite this version:**

Alexandre Ern, Jean-Luc Guermond. Third-order A-stable alternating implicit Runge-Kutta schemes. 2024. hal-04527220

**HAL Id: hal-04527220**

**<https://hal.science/hal-04527220>**

Preprint submitted on 29 Mar 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Third-order A-stable alternating implicit Runge–Kutta schemes\*

Alexandre Ern<sup>†</sup> and Jean-Luc Guermond<sup>‡</sup>

Draft version March 29, 2024

## Abstract

We design pairs of six-stage, third-order, alternating implicit Runge–Kutta (RK) schemes that can be used to integrate in time two stiff operators by an operator-split technique. We also design for each pair a companion explicit RK scheme to be used for a third, nonstiff operator in an IMEX fashion. The main application we have in mind are (non)linear parabolic problems, where the two stiff operators represent diffusion processes (for instance, in two spatial directions) and the nonstiff operator represents (non)linear transport. We identify necessary conditions for linear  $A(\alpha)$ -stability by considering a scalar ODE with two (complex) eigenvalues lying in some fixed cone of the half-complex plane with nonpositive real part. We show numerically that it is possible to achieve  $A(0)$ -stability when combining two operators with negative eigenvalues, irrespective of their relative magnitude. Finally, we show by numerical examples including two-dimensional nonlinear transport problems discretized in space using finite elements that the proposed schemes behave well.

**Keywords.** High-order time integration, Operator splitting, Implicit-explicit time integration, Order barrier.

**MSC.** 35L65, 65M60, 65M12, 65N30

## 1 Introduction

Operator-splitting is a well established and computationally effective approach to design time-integration techniques for a wide class of systems of stiff ordinary differential equations (ODEs) and partial differential equations (PDEs) involving coupled stiff operators. One traditional way to split two stiff operators consists of using methods like Strang splitting [30] at the time-continuous level or the Peaceman–Rachford alternating direction implicit method (ADI) at the time-discrete level [20] (see also Douglas and Rachford [8]). We refer the reader, e.g., to Marchuk [17], Yanenko [33] for early surveys on the subject.

The stiff PDE model we have in mind is that of (non)linear parabolic equations, where two operators are stiff (say (non)linear diffusion in different directions) and a third one is less stiff (say

---

<sup>†</sup>CERMICS, Ecole des Ponts, 77455 Marne-la-Vallée Cedex 2, France and INRIA Paris, 75589 Paris, France

<sup>‡</sup>Department of Mathematics, Texas A&M University 3368 TAMU, College Station, TX 77843, USA.

\*This material is based upon work supported in part by the National Science Foundation grant DMS2110868, the Air Force Office of Scientific Research under contract number FA9550-18-1-0397, the Army Research Office under grant number W911NF-19-1-0431, the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contracts B640889, and INRIA through the International Chair program. The authors are thankful to Ari Rappaport (INRIA Paris) for his instrumental help with the `julia` package.

nonlinear transport). Our objective is to construct a method that is third-order accurate in time when the two stiff operators are split, while the nonstiff operator is treated explicitly in an IMEX fashion. This is a non-trivial task since operator-splitting methods face a second-order accuracy barrier. More precisely, the accuracy of exponential splitting methods is reduced to second-order if one excludes any strategy requiring backward time integration and linear combinations of forward-stepping exponential splitting methods with negative multiplicative coefficients, see Sheng [27], Suzuki [31], Goldman and Kaper [11], and Blanes and Casas [4]. One remedy to break the second-order barrier consists of adopting complex time integration. This idea was suggested by Rosenbrock [24] and Bandrauk and Shen [3]. It was formalized up to fourth-order in Gegechkori et al. [10] and up to order fourteen in Hansen and Ostermann [15] and Castella et al. [5]. A second class of methods also potentially capable of breaking the second-order barrier consists of using defect correction strategies, as shown in Christlieb et al. [6].

The third option, which is the one we consider in the paper, consists of interlacing two implicit Runge–Kutta (RK) schemes. By this, we mean that, at every stage of the method, only one of the two implicit schemes has a nonzero diagonal entry, and this feature alternates at every stage. The resulting RK scheme is called alternating-implicit (in short, AIRK). The prototypical second-order example is actually the Peaceman–Rachford ADI method which is built by combining the implicit midpoint rule with the Crank–Nicolson scheme. This leads to an A-stable, two-stage, second-order AIRK scheme, where only one of the two stiff operators is treated implicitly at each of the two stages. Our ambition here is not to be general, but to demonstrate that the second-order accuracy barrier can be overcome by interlacing two six-stage, third-order implicit RK schemes, while maintaining some form of A-stability. In the paper, we provide two examples of such AIRK schemes. For both examples, the two constitutive implicit RK schemes are singly diagonal and  $A(\alpha)$ -stable, and, in one of the examples, the two schemes are even  $L(\alpha)$ -stable. Moreover, for both examples, we propose a companion explicit RK (ERK) scheme which can be used in conjunction with the AIRK scheme in an IMEX fashion.

The idea of interlacing two (or more) RK schemes has been well explored in the literature. We refer the reader to Cooper and Sayfy [7], Rentrop [21], Rice [22] for early works on the subject, leading in particular to the notion of additive RK (ARK) methods. An importance instance of ARK schemes are the implicit-explicit (IMEX) methods developed by Ascher et al. [1, 2], Kennedy and Carpenter [16], Pareschi and Russo [18, 19], Zhong [34]. The order conditions for ARK schemes are well understood through the concept of P-trees developed by Hairer [13]. A further important development of ARK schemes is the class of generalized-structure ARK (GARK) schemes in Sandu and Günther [25], where several copies of the dependent unknowns are advanced at each stage. We refer the reader, e.g., to González-Pinto et al. [12], Roberts et al. [23], Sarshar et al. [26], Spiteri and Wei [29] for recent developments on the subject. We observe that the present AIRK schemes can be viewed as a particular instance of GARK schemes (see Remark 2.2 for further discussion).

GARK schemes constitute an effective framework to devise high-order operator-split techniques. However, establishing some form of stability for high-order GARK schemes (say, beyond second-order) is still a nontrivial question at the time of this writing. Indeed, even if the implicit RK schemes considered for each operator enjoy some form of linear stability, say  $A(\alpha)$ -stability or even  $L(\alpha)$ -stability, the linear stability of the resulting AIRK scheme generally remains an open question. This question can be studied by considering Dahlquist’s test problem in various settings, whereby a scalar ODE is considered with each operator represented by a complex number in the half-complex plane with nonpositive real part. Quite importantly, the question also needs to be studied numerically on more realistic situations beyond linear stability, e.g., for PDEs modeling nonlinear transport.

One important contribution of the paper is to identify some necessary conditions for  $A(\alpha)$ - and  $L(\alpha)$ -stability when combining two implicit RK schemes into an AIRK scheme, under the

assumption that the spectra of the two split operators lie in some fixed cone around the negative real axis with an acute half angle. This assumption is reasonable for our purposes since the stiff operators represent (non)linear diffusion processes. Moreover, we verify numerically that, when combining two operators with negative eigenvalues, the AIRK schemes we propose are indeed  $A(0)$ -stable, uniformly with respect to the relative magnitude of the eigenvalues. Finally, we assess numerically the performances of the proposed AIRK schemes on a series of challenging test cases resulting from the finite element discretization of two-dimensional nonlinear advection-diffusion problems.

The paper is organized as follows. In Section 2, we make the setting precise and establish useful results to study the linear stability of AIRK schemes. Our main result is Lemma 2.4. In Section 3, we focus on six-stage, implicit schemes and identify sufficient conditions to achieve third-order accuracy as well as necessary conditions to achieve suitable linear stability properties, see, in particular, Lemma 3.3 and Lemma 3.5. We also discuss the design of the companion ERK scheme to be used for the nonstiff operator. In Section 4, we study numerically the properties of the AIRK and ERK schemes obtained in the previous section. We perform a series of tests on two-dimensional advection diffusion equations and nonlinear transport problems discretized in space with finite elements. We close this work with two appendices. In Appendix A, we give two examples of operator-split schemes fulfilling the design conditions identified in Section 3; each example comprises an AIRK scheme and one or two companion ERK scheme(s). In particular, we show numerically that the necessary linear stability conditions identified in Section 3 indeed lead to  $A(0)$ -stability when combining two operators with negative eigenvalues. The  $A(0)$ -stability is uniform with respect to the relative magnitude of the eigenvalues. Finally, in Appendix B, we collect some results on four-stage, third-order and two-stage, second-order AIRK schemes. We show that there is a stability barrier for the former, and that the only possible realization for the latter is essentially the Peacemann–Rachford scheme.

## 2 Setting

In this section, we introduce some useful notions and derive some preliminary results on the linear stability of AIRK schemes.

### 2.1 Model problem

Given a time horizon  $T > 0$ , we want to approximate in time the following nonlinear system of  $I$  coupled ODEs, which consists of seeking  $\mathbf{U} \in C^1([0, T]; \mathbb{R}^I)$  so that

$$\partial_t \mathbf{U}(t) = L_0(t, \mathbf{U}(t)) + L_1(t, \mathbf{U}(t)) + L_2(t, \mathbf{U}(t)), \quad \mathbf{U}(0) = \mathbf{U}^0 \in \mathbb{R}^I, \quad (1)$$

where we make the usual assumption on the Lipschitz continuity with respect to  $\mathbf{U}$  and continuity with respect to  $t$  of  $L_0, L_1, L_2$ . We additionally assume that the Lipschitz constants of  $L_0(t, \cdot) : \mathbb{R}^I \rightarrow \mathbb{R}^I$  and  $L_1(t, \cdot) : \mathbb{R}^I \rightarrow \mathbb{R}^I$  are significantly larger than that of  $L_2(t, \cdot) : \mathbb{R}^I \rightarrow \mathbb{R}^I$ . Our objective is to design a third-order time-stepping method where  $L_2$  is treated explicitly and  $L_0, L_1$  are treated implicitly in an alternating fashion by means of an AIRK scheme.

### 2.2 Butcher tableaux

To achieve the task described above, we want to combine three Butcher tableaux composed of  $s+1$  stages where  $s \geq 2$  is even,

$$\begin{array}{c|c} c & A_0 \\ \hline & b_0 \end{array} \quad \begin{array}{c|c} c & A_1 \\ \hline & b_1 \end{array} \quad \begin{array}{c|c} c & A_2 \\ \hline & b_2 \end{array}. \quad (2)$$

Notice that the three Butcher tableaux share the same time index vector  $c$  (this property is called internal consistency in the context of ARK schemes). We additionally assume that

$$c_1 = 0, \quad c_{s+1} = 1, \quad (3a)$$

$$b_0 = e_{s+1}^\top A_0, \quad b_1 = e_{s+1}^\top A_1, \quad b_2 = e_{s+1}^\top A_2, \quad (3b)$$

$$A_0 U = c \quad A_1 U = c \quad A_2 U = c, \quad (3c)$$

where  $e_{s+1}$  is the last vector of the canonical Cartesian basis of  $\mathbb{R}^{s+1}$  and  $U$  is the column vector in  $\mathbb{R}^{s+1}$  having all its entries equal to one. In (3b), we request that the line vectors  $b_0, b_1, b_2$  be copies of the last row of the matrices  $A_0, A_1, A_2$ , respectively. This property means that the implicit schemes are stiffly accurate. Moreover, the identities (3c) are Butcher's simplifying assumption. Notice that the assumptions (3) imply that  $b_0 U = e_{s+1}^\top A_0 U = e_{s+1}^\top c = c_{s+1} = 1$  and, similarly,  $b_1 U = b_2 U = 1$ .

We assume that the matrices  $A_0, A_1$  are lower triangular with the upper left entry equal to zero, and the matrix  $A_2$  is strictly lower triangular. The scheme associated with  $A_2$  is therefore explicit. The schemes associated with  $A_0, A_1$  are a priori diagonally implicit, but we further simplify the method by requesting that the matrices  $A_0, A_1$  have alternating nonzero coefficients on the diagonal, i.e., we assume that

$$(A_0)_{l,l} = 0, \quad \text{mod}(l, 2) = 1, \quad \forall l \in \{1:s+1\}, \quad (4a)$$

$$(A_1)_{l,l} = 0, \quad \text{mod}(l, 2) = 0, \quad \forall l \in \{1:s+1\}. \quad (4b)$$

We say that the combined RK scheme is alternating-implicit for this reason.

Let  $t^n$  be the current discrete time node and  $\tau^n$  be the current time step. We set  $t^{n+1} := t^n + \tau^n$  and  $t^{n,m} := t^n + c_m \tau^n$  for all  $m \in \{1:s+1\}$ . The IMEX RK scheme associated with (2) consists of marching from  $t^n$  to the next discrete time node  $t^{n+1}$  by performing the following  $s$  stages: Given  $U^n$ , set  $U^{n,1} := U^n$  and compute, for all  $l \in \{2:s+1\}$ ,

$$\begin{aligned} & U^{n,l} - \tau^n \left\{ (A_0)_{ll} L_0(t^{n,l}, U^{n,l}) + (A_1)_{ll} L_1(t^{n,l}, U^{n,l}) \right\} \\ &= U^n + \tau^n \sum_{m \in \{1:l-1\}} \left\{ (A_0)_{lm} L_0(t^{n,m}, U^{n,m}) + (A_1)_{lm} L_1(t^{n,m}, U^{n,m}) \right. \\ & \quad \left. + (A_2)_{lm} L_2(t^{n,m}, U^{n,m}) \right\}, \end{aligned} \quad (5)$$

and finally set  $U^{n+1} := U^{n,s+1}$ . Owing to the assumption (4), we obtain an AIRK scheme since, at every stage, only one of the stiff operators  $L_0, L_1$  is treated implicitly. The operator  $L_2$  is treated explicitly at all stages.

**Remark 2.1** ( $s$ -stage AIRK). *Note that the first stage is trivial ( $U^{n,1} := U^n$ ). The  $(s+2)$ th stage is trivial as well ( $U^{n,s+2} = U^{n,s+1}$ ) owing to the assumption (3b). Hence, the scheme is actually composed of  $s$  stages.*

**Remark 2.2** (Rewriting in GARK format). *Setting  $s' := \frac{s}{2}$ , one can distribute the stage updates  $(U^l)_{l \in \{1:s+1\}}$  (we drop the superscript  $n$  to ease the notation) into the two collections  $(Y^{1,l})_{l \in \{1:s'+1\}}$  and  $(Y^{2,l})_{l \in \{1:s'+1\}}$  so that  $Y^{1,1} = Y^{2,1} = U^1$  and  $Y^{1,l} = U^{2l-2}$ ,  $Y^{2,l} = U^{2l-1}$  for all  $l \in \{2:s'+1\}$ . Then, (5) can be rewritten as follows: For all  $l \in \{2:s'+1\}$ , solve sequentially for  $i \in \{1, 2\}$ ,*

$$Y^{i,l} = U^n + \tau^n \sum_{m \in \{1:l\}} \sum_{p,q \in \{0,1\}} \mathfrak{A}_{lm}^{i,pq} L_p(t^{n,m}, Y^{q,m}), \quad (6)$$

where the eight arrays  $(\mathfrak{A}^{r,pq})_{i,p,q \in \{0,1\}}$  are all of order  $(s'+1)$ , lower triangular, and with upper left diagonal entry equal to zero. Moreover, only the arrays  $\mathfrak{A}^{0,00}$ ,  $\mathfrak{A}^{1,11}$ , and  $\mathfrak{A}^{1,00}$ ,  $\mathfrak{A}^{1,10}$  have nonzero diagonal entries (the latter two do not lead to an implicit treatment owing to the sequential solve in  $i \in \{0,1\}$ ). Notice that GARK schemes are often written by discarding the arrays  $\mathfrak{A}^{i,pq}$  with  $p \neq q$ . These arrays are nonzero in the present AIRK formalism. Another significant difference is that the two variables  $\mathfrak{Y}^{1,l}$  and  $\mathfrak{Y}^{2,l}$  are not synchronized in the present setting. We refer the reader to Section A.4 for an example with a six-stage AIRK scheme.

### 2.3 Linear stability: amplification functions

The classical approach to analyze the linear stability of a single implicit RK scheme consists of considering the scalar ODE  $\partial_t u = \lambda u(t)$  with  $\lambda \in \mathbb{C}^- := \{z \in \mathbb{C} \mid \Re(z) \leq 0\}$  (this ODE is often called Dahlquist's test problem). Separately considering the Butcher tableaux in (2) for  $i = 0$  and  $i = 1$  leads to the following two amplification functions for all  $i \in \{0,1\}$  (which we call single-array amplification functions): For all  $z \in \mathbb{C}^-$ ,

$$R_i(z) := 1 + \frac{\rho_i(z)}{\det(I - zA_i)}, \quad \rho_i(z) = \det(I - zA_i)zb_i(I - zA_i)^{-1}U. \quad (7)$$

(We introduce the function  $\rho_i(z)$  for later use.) Recall that the implicit RK scheme associated with the  $i$ th Butcher tableau is said to be  $A(\alpha)$ -stable if there is an angle  $\alpha_i \in [0, \frac{\pi}{2}]$  such that  $|R_i(z)| \leq 1$  for all  $z \in C(\alpha_i)$ ; see Widlund [32], Hairer and Wanner [14, Def. 3.7&3.9]. Here, for a generic angle  $\beta \in [0, \frac{\pi}{2}]$ , we defined the cone  $C(\beta) := \{z \in \mathbb{C}^- \mid \arg(-z) \leq \beta\}$ . Moreover, the scheme is said to be  $L(\alpha)$ -stable if it is  $A(\alpha)$ -stable and  $\ell_i := \lim_{|z| \rightarrow \infty} R_i(z) = 0$ .

In the present setting with two stiff operators, the natural extension of Dahlquist's test problem is to consider the scalar ODE

$$\partial_t U(t) = \lambda_0 U(t) + \lambda_1 U(t), \quad (8)$$

with  $\lambda_i \in \mathbb{C}^-$  for all  $i \in \{0,1\}$ . This test problem is, however, too general for our present purpose, where the two stiff operators are diffusion operators, so that their spectrum is a discrete subset of the negative real axis in the complex plane. To allow for a bit more generality at this stage, we assume that there is an angle  $\beta \in [0, \frac{\pi}{2}]$  such that  $\lambda_i \in C(\beta)$  for all  $i \in \{0,1\}$ . We have  $\beta = 0$  for diffusion operators. Setting  $z := \lambda_0 + \lambda_1$ ,  $\theta := \frac{\lambda_1}{\lambda_0 + \lambda_1}$ ,  $1 - \theta = \frac{\lambda_0}{\lambda_0 + \lambda_1}$ , (8) reduces to  $\partial_t U(t) = z((1 - \theta)U(t) + \theta U(t))$ . Observe that both  $\theta$  and  $(1 - \theta)$  are in the ball  $B(\beta)$  centered at  $\frac{1}{2}$  and of radius  $\frac{1}{2}(1 + \tan^2(\beta))^{\frac{1}{2}}$ . Therefore, linear stability can be studied by assuming that  $\theta$  and  $(1 - \theta)$  are uniformly bounded.

The amplification function for the scheme (5) applied to the ODE (8) is

$$R_\theta(z) = 1 + \frac{\rho_\theta(z)}{\det(I - zA_\theta)}, \quad \rho_\theta(z) := \det(I - zA_\theta)zb_\theta(I - zA_\theta)^{-1}U. \quad (9)$$

with  $A_\theta := (1 - \theta)A_0 + \theta A_1$  and  $b_\theta := (1 - \theta)b + \theta b_1$ . In the above setting, we can use the following notion of stability for AIRK schemes.

**Definition 2.3** (Sectorial  $A(\alpha)$ -stability and  $L(\alpha)$ -stability for AIRK schemes). *We say that the AIRK scheme (5) is sectorial  $A(\alpha)$ -stable if there is an angle  $\alpha \in [0, \beta]$  s.t. for all  $\theta \in B(\beta)$ ,  $|R_\theta(z)| \leq 1$  for all  $z \in C(\alpha)$ . We say that the scheme is sectorial  $L(\alpha)$ -stable if it is  $A(\alpha)$ -stable and  $\ell_\theta := \lim_{|z| \rightarrow \infty} R_\theta(z) = 0$  for all  $\theta \in B(\beta)$ . In what follows, to ease the terminology, we simply speak of  $A(\alpha)$ - and  $L(\alpha)$ -stability.*

For a lower-triangular matrix  $\Lambda$  of order  $(s+1)$  with diagonal entries  $\{\lambda_i\}_{i \in \{1:s+1\}}$  (the example we have in mind is  $\Lambda = A_\theta$ ), we set

$$\mathrm{tr}_m(\Lambda) := \sum_{\substack{(i_1, \dots, i_m) \in \{1:s\}^m \\ i_1 < \dots < i_m}} \lambda_{i_1} \times \dots \times \lambda_{i_m}, \quad \forall m \in \{1:s+1\}, \quad (10)$$

and we conventionally set  $\mathrm{tr}_0(\Lambda) := 1$ . Notice that  $\mathrm{tr}_1(\Lambda)$  is the usual trace of  $\Lambda$  and  $\mathrm{tr}_{s+1}(\Lambda) = \lambda_1 \times \dots \times \lambda_{s+1}$ . The characteristic polynomial of the matrix  $\Lambda$  is

$$\pi_\Lambda(t) = \det(tI - \Lambda) = \sum_{k \in \{0:s+1\}} (-1)^{s+1-k} \mathrm{tr}_{s+1-k}(\Lambda) t^k. \quad (11)$$

The Hamilton–Cayley theorem gives

$$\pi_\Lambda(\Lambda) = \sum_{k \in \{0:s+1\}} (-1)^{s+1-k} \mathrm{tr}_{s+1-k}(\Lambda) \Lambda^k = 0 \in \mathbb{R}^{s+1, s+1}. \quad (12)$$

Finally, we notice that, whenever the matrix  $\Lambda$  has only  $m$  nonzero diagonal coefficients with  $m \leq s$ , we have  $\mathrm{tr}_k(\Lambda) = 0$  for all  $k \in \{m+1:s+1\}$ . Notice, in particular, that  $\mathrm{tr}_{s+1}(A_\theta) = 0$  and that  $\mathrm{tr}_m(A_0) = \mathrm{tr}_m(A_1) = 0$  for all  $m \geq \frac{s}{2} + 1$  owing to (4).

To gain some insight into the amplification function  $R_\theta(z)$ , we study the function  $\rho_\theta(z)$  defined in (9).

**Lemma 2.4** (Function  $\rho_\theta$ ). *The function  $\rho_\theta$  defined in (9) is a polynomial in  $z$  of degree at most  $s$ ,  $\rho_\theta(z) = \sum_{k \in \{0:s-1\}} \omega_k(\theta) z^{k+1}$ , where for all  $k \in \{0:s-1\}$ ,*

$$\omega_k(\theta) := \sum_{l \in \{0:k\}} \beta_{k-l}(\theta) \tau_l(\theta), \quad (13a)$$

$$\beta_k(\theta) := b_\theta A_\theta^k U, \quad \tau_k(\theta) := (-1)^k \mathrm{tr}_k(A_\theta). \quad (13b)$$

Moreover,  $\omega_k(\theta)$  is a polynomial in  $\theta$  of degree at most  $k$  with real-valued coefficients.

*Proof.* Since  $\Phi_\theta(z) := \det(I - zA_\theta)(I - zA_\theta)^{-1}$  is the transpose of the cofactor matrix of  $(I - zA_\theta)$ . As the matrix  $(I - zA_\theta)$  is lower triangular with upper left entry equal to 1, the entries of the matrix  $\Phi_\theta(z)$  are all polynomials in  $z$  of degree at most  $s$ . Hence,  $\rho_\theta(z)$  is a polynomial of degree at most  $(s+1)$  in  $z$ . To see that the degree of  $\rho_\theta(z)$  is actually at most  $s$  instead of  $(s+1)$ , we compute the coefficients of the matrix-valued polynomial  $\Phi_\theta(z)$ . Since  $\mathrm{tr}_{s+1}(A_\theta) = 0$ , we have

$$\Phi_\theta(z) = \left\{ \sum_{l \in \{0:s\}} (-1)^l \mathrm{tr}_l(A_\theta) z^l \right\} \sum_{m \geq 0} z^m A^m = \sum_{k \in \{0:s\}} \left\{ \sum_{l \in \{0:k\}} (-1)^l \mathrm{tr}_l(A_\theta) A^{k-l} \right\} z^k.$$

Since  $\rho_\theta(z) = zb_\theta \Phi_\theta(z)U$ , we infer using the definitions (13b) that

$$\rho_\theta(z) = \sum_{k \in \{0:s\}} \left\{ \sum_{l \in \{0:k\}} \tau_l(\theta) \beta_{k-l}(\theta) \right\} z^{k+1}.$$

Setting  $\omega_k(\theta) := \sum_{l \in \{0:k\}} \beta_{k-l}(\theta) \tau_l(\theta)$  for all  $k \in \{0:s\}$  as in (13a), and observing that  $\omega_0(\theta) = \beta_0(\theta) \tau_0(\theta) = 1$  (notice that  $\beta_0(\theta) = b_\theta U = (1-\theta) + \theta = 1$ ), we conclude that  $\rho_\theta(z) = \sum_{k \in \{0:s\}} \omega_k(\theta) z^{k+1}$ .

Therefore, it only remains to prove that  $\omega_s(\theta) = 0$ . Using (3b), i.e.,  $\beta_m(\theta) = b_\theta A_\theta^m U = e_{s+1}^\top A_\theta^{m+1} U$  for all  $m \geq 0$ , we obtain

$$\begin{aligned} \omega_s(\theta) &= \sum_{l \in \{0:s\}} \pi_l(\theta) \beta_{s-l}(\theta) = e_{s+1}^\top \left( \sum_{l \in \{0:s\}} (-1)^l \operatorname{tr}_l(A_\theta) A_\theta^{s+1-l} \right) U \\ &= e_{s+1}^\top \left( \sum_{l \in \{1:s+1\}} (-1)^{s+1-l} \operatorname{tr}_{s+1-l}(A_\theta) A_\theta^l \right) U \\ &= e_{s+1}^\top \pi_{A_\theta}(A_\theta) U, \end{aligned}$$

where we used that  $\operatorname{tr}_{s+1}(A_\theta) = 0$ . Owing to the Hamilton–Cayley theorem, we conclude that  $\omega_s(\theta) = 0$ . Finally, the expressions (13) show that  $\omega_k(\theta)$  is a polynomial in  $\theta$  of degree at most  $(k+1)$  having real-valued coefficients. Since  $A_\theta U = c$  owing to (3c), the degree is at most  $k$ .  $\square$

### 3 Six-stage, third-order, AIRK schemes

The main focus of the paper is when  $s = 6$ , with both  $A_0$  and  $A_1$  having three nonzero diagonal coefficients interlaced along the diagonal. Thus, we consider two six-stage, implicit RK schemes having the following structure (we omit the vectors  $b_0, b_1$  since the schemes are stiffly accurate, see (3b)):

$$\begin{array}{c|cccccccc} 0 & 0 & & & & & & & 0 & 0 \\ c_2 & A_{21}^0 & A_{22}^0 & & & & & & c_2 & A_{21}^1 & 0 \\ c_3 & A_{31}^0 & A_{32}^0 & 0 & & & & & c_3 & A_{31}^1 & A_{32}^1 & A_{33}^1 \\ c_4 & A_{41}^0 & A_{42}^0 & A_{43}^0 & A_{44}^0 & & & & c_4 & A_{41}^1 & A_{42}^1 & A_{43}^1 & 0 \\ c_5 & A_{51}^0 & A_{52}^0 & A_{53}^0 & A_{54}^0 & 0 & & & c_5 & A_{51}^1 & A_{52}^1 & A_{53}^1 & A_{54}^1 & A_{55}^1 \\ c_6 & A_{61}^0 & A_{62}^0 & A_{63}^0 & A_{64}^0 & A_{65}^0 & A_{66}^0 & & c_6 & A_{61}^1 & A_{62}^1 & A_{63}^1 & A_{64}^1 & A_{65}^1 & 0 \\ 1 & A_{71}^0 & A_{72}^0 & A_{73}^0 & A_{74}^0 & A_{75}^0 & A_{76}^0 & 0 & 1 & A_{71}^1 & A_{72}^1 & A_{73}^1 & A_{74}^1 & A_{75}^1 & A_{76}^1 & A_{77}^1 \end{array}$$

#### 3.1 Third-order conditions

Let  $U$  be the column vector in  $\mathbb{R}^7$  having all its entries equal to 1. Let  $c^2$  be the column vector in  $\mathbb{R}^7$  having all its entries equal  $c_m^2$  for all  $m \in \{1:7\}$ . The single-array third-order conditions are (3c) together with

$$b_0 c = b_1 c = \frac{1}{2}, \quad (14a)$$

$$b_0 c^2 = b_1 c^2 = \frac{1}{3}, \quad (14b)$$

$$b_0 A_0 c = b_1 A_1 c = \frac{1}{6}. \quad (14c)$$

Recall that  $b_0 U = b_1 U = 1$  follows from (3b) and (3c). Moreover, the coupling third-order conditions are

$$b_0 A_1 c = b_1 A_0 c = \frac{1}{6}. \quad (15)$$

**Lemma 3.1** ( $\beta_0(\theta), \beta_1(\theta), \beta_2(\theta)$ ). *Assume (3c), (14) and (15). With the coefficients  $\beta_k(\theta)$  are defined in (13b), the following holds:*

$$\beta_0(\theta) = 1, \quad \beta_1(\theta) = \frac{1}{2}, \quad \beta_2(\theta) = \frac{1}{6}. \quad (16)$$

*Proof.* By linearity, we have  $b_\theta U = 1$ ,  $b_\theta c = \frac{1}{2}$ , and  $A_\theta c = U$ . This shows that  $\beta_0(\theta) = b_\theta U = 1$  and  $\beta_1(\theta) = b_\theta A_\theta U = b_\theta c = \frac{1}{2}$  owing to (3c). Finally, a direct calculation shows that

$$\beta_2(\theta) = b_\theta A_\theta c = (1-\theta)^2 b_0 A_0 c + \theta(1-\theta)(b_0 A_1 c + b_1 A_0 c) + \theta^2 b_1 A_1 c = \frac{1}{6}((1-\theta) + \theta)^2 = \frac{1}{6},$$

where we used (3c), (14c) and (15).  $\square$



### 3.2 Linear stability

This section collects important results concerning the amplification function associated with the combined Butcher tableaux and the amplification functions associated with each tableau individually (which we call single-array amplification function).

**Lemma 3.2** (Function  $\rho_\theta(z)$ ). *The function  $\rho_\theta$  defined in (9) is a polynomial in  $z$  of degree at most 6, of the form  $\rho_\theta(z) = \sum_{k \in \{0:5\}} \omega_k(\theta) z^{k+1}$  with*

$$\omega_5(\theta) = (b_\theta A_\theta^4 c) + (b_\theta A_\theta^3 c) \tau_1(\theta) + (b_\theta A_\theta^2 c) \tau_2(\theta) + \frac{1}{6} \tau_3(\theta) + \frac{1}{2} \tau_4(\theta) + \tau_5(\theta), \quad (17a)$$

$$\omega_4(\theta) = (b_\theta A_\theta^3 c) + (b_\theta A_\theta^2 c) \tau_1(\theta) + \frac{1}{6} \tau_2(\theta) + \frac{1}{2} \tau_3(\theta) + \tau_4(\theta), \quad (17b)$$

$$\omega_3(\theta) = (b_\theta A_\theta^2 c) + \frac{1}{6} \tau_1(\theta) + \frac{1}{2} \tau_2(\theta) + \tau_3(\theta), \quad (17c)$$

$$\omega_2(\theta) = \frac{1}{6} + \frac{1}{2} \tau_1(\theta) + \tau_2(\theta), \quad (17d)$$

$$\omega_1(\theta) = \frac{1}{2} + \text{tr}_1(\theta), \quad (17e)$$

and  $\omega_0(\theta) = 1$ .

*Proof.* Combine Lemma 2.4 with Lemma 3.1 and (3c) to establish (17).  $\square$

**Lemma 3.3** (Necessary condition for  $A(\alpha)$ -stability, AIRK scheme). *A necessary condition for the  $A(\alpha)$ -stability of the AIRK scheme is*

$$\omega_5(\theta) = 0, \quad \forall \theta \in B(\beta). \quad (18)$$

Moreover, under this condition, we have  $\ell_\theta = 1$  for all  $\theta \in B^\circ(\beta) := B(\beta) \setminus \{0, 1\}$ .

*Proof.* We notice that, as  $|z| \rightarrow \infty$ ,  $\rho_\theta(z) \sim \omega_5(\theta) z^6$  for all  $\theta \in B(\beta)$  such that  $\omega_5(\theta) \neq 0$ , whereas  $\det(I - zA_\theta) \sim \theta^3(1 - \theta)^3 \text{tr}_3(A_0) \text{tr}_3(A_1) z^6$  for all  $\theta \in B^\circ(\beta)$ . This implies that  $R_\theta(z) \sim 1 + \frac{\omega_5(\theta)}{\theta^3(1-\theta)^3} (\text{tr}_3(A_0) \text{tr}_3(A_1))^{-1}$  for all  $\theta \in B^\circ(\beta)$  s.t.  $\omega_5(\theta) \neq 0$ . Since  $\omega_5(\theta) \in \mathbb{P}_5[\theta]$ ,  $R_\theta(z)$  can stay bounded as  $|z| \rightarrow \infty$  only if (18) holds true. Finally, the fact that  $\ell_\theta = 1$  for all  $\theta \in B^\circ(\beta)$  readily follows from the above asymptotic expression for  $R_\theta(z)$  and  $\omega_5(\theta) = 0$ .  $\square$

**Remark 3.4** (Barrier on  $L(\alpha)$ -stability). *A striking consequence of (3.3) is that a six-stage, third-order AIRK scheme cannot be  $L(\alpha)$ -stable since  $\ell_\theta = 1 \neq 0$  for all  $\theta \notin \{0, 1\}$ . We shall see though that it is still possible to make the two interlaced implicit RK schemes  $L(\alpha)$ -stable (see Remark 3.7 below for further discussion).*

Let us now consider the single-array amplification functions. Let  $i \in \{0, 1\}$  and set  $\rho_i(z) := \det(I - zA_i) z b_i (I - zA_i)^{-1} U$  (see (7)). We infer from Lemma 3.2 that  $\rho_i(z) = \sum_{k \in \{0:5\}} \omega_k^i z^{k+1}$  with

$$\omega_k^i := \omega_k(i), \quad \forall i \in \{0, 1\}, \forall k \in \{0:5\}. \quad (19)$$

Let us set  $\tau_k^i := \tau_k(i)$  (recall that  $\tau_k(\theta) := (-1)^k \text{tr}_k(A_\theta)$ ).

**Lemma 3.5** (Necessary condition for  $A(\alpha)$ -stability, single RK schemes). *A necessary condition for  $A(\alpha)$ -stability for each single RK scheme is, for all  $i \in \{0, 1\}$ ,*

$$\omega_3^i = \omega_4^i = \omega_5^i = 0, \quad (20a)$$

$$\omega_2^i = (1 - \ell_i) \tau_3^i, \quad \ell_i \in [-1, 1]. \quad (20b)$$

*Proof.* The reasoning is similar to the one in the proof of Lemma 3.3, the only difference being that  $\det(I - zA_i) \sim -\text{tr}_3(A_i)z^3$  as  $|z| \rightarrow \infty$ . Therefore,  $R_i(z)$  can stay bounded as  $|z| \rightarrow \infty$  only if  $\omega_3^i = \omega_4^i = \omega_5^i = 0$ , which gives (20a). Moreover, in this situation, we obtain  $\lim_{|z| \rightarrow \infty} R_i(z) = 1 - \frac{\omega_2^i}{\text{tr}_3(A_i)} = \ell_i \in [-1, 1]$  owing to (20b).  $\square$

Owing to (17) and since  $\tau_4^i = \tau_5^i = 0$  (recall that both matrices  $A_i$  have only three nonzero diagonal coefficients), the conditions (20a) can be rewritten as follows: For all  $i \in \{0, 1\}$ ,

$$(b_i A_i^4 c) + (b_i A_i^3 c) \tau_1^i + (b_i A_i^2 c) \tau_2^i + \frac{1}{6} \tau_3^i = 0, \quad (21a)$$

$$(b_i A_i^3 c) + (b_i A_i^2 c) \tau_1^i + \frac{1}{6} \tau_2^i + \frac{1}{2} \tau_3^i = 0, \quad (21b)$$

$$(b_i A_i^2 c) + \frac{1}{6} \tau_1^i + \frac{1}{2} \tau_2^i + \tau_3^i = 0, \quad (21c)$$

$$\frac{1}{6} + \frac{1}{2} \tau_1^i + \tau_2^i + (1 - \ell_i) \tau_3^i = 0. \quad (21d)$$

**Remark 3.6** (Singly diagonal case). *If the array  $A_i$  is singly diagonal with entry  $a$ , (21d) readily implies that this entry must be a positive root of the cubic equation  $(1 - \ell)x^3 - 3x^2 + \frac{3}{2}x - \frac{1}{6} = 0$ . For  $\ell = 0$ , we obtain  $a = 0.1589\dots$ . For  $\ell = 1$ , the equation becomes quadratic and the positive root is  $a = \frac{1}{6}$ . Notice also that, if both arrays  $A_0$  and  $A_1$  are singly diagonal and such that  $\ell_0 = \ell_1$ , (20b) implies that  $\omega_2^0 = \omega_2^1$ . Since  $\omega_1^0 = \omega_1^1 = \frac{1}{2} + 3a$  by (17e), we infer that the amplification functions  $R_0$  and  $R_1$  are the same.*

**Remark 3.7** (Singular limit). *Recall that  $\ell_\theta = 1$  for all  $\theta \in B^\circ(\beta)$  owing to Lemma 3.3, whereas Lemma 3.5 shows that it is possible to fix  $\ell_i \in [-1, 1]$  for all  $i \in \{0, 1\}$ . There are, therefore, two somewhat natural choices when it comes to fixing the limits  $\ell_i$ . The first one is to select  $\ell_0 = \ell_1 = 0$ , so that the two constitutive implicit RK schemes are  $L(\alpha)$ -stable, but in this case the limits  $\lim_{|z| \rightarrow \infty}$  and  $\lim_{\theta \rightarrow 0}$  (or  $\lim_{\theta \rightarrow 1}$ ) do not commute. The second one is to enforce  $\ell_0 = \ell_1 = 1$ , which leads to two  $A(\alpha)$ -stable implicit RK schemes and the above two limits commute.*

### 3.3 Summary of devising conditions

The devising conditions on the two tableaux composing the AIRK scheme are collected in Table 1. We first collect in the two columns labeled  $i = 0$  and  $i = 1$  the design conditions that are specific to each Butcher tableau. The last four lines of the table (spanning the two columns) collect the design conditions coupling both Butcher tableaux. The design parameters are the column vector  $c \in \mathbb{R}^7$  with  $c_1 = 0$  and  $c_7 = 1$ , the limits  $\ell_0, \ell_1 \in [-1, 1]$ , and a small parameter  $\epsilon \geq 0$ . Since  $\omega_5(\theta)$  is a polynomial of degree at most 5 in  $\theta$  having real coefficients, we infer that  $\omega_5 \equiv 0$  iff  $\omega_5(0) = \omega_5(1) = 0$ ,  $\omega_5'(0) = \omega_5'(1) = 0$  and  $\omega_5''(0) = \omega_5''(1) = 0$ , which are indeed the conditions recorded in Table 1. As  $\omega_4(\theta)z^5$  is the dominating factor in  $\rho(\theta)$ , one can further reduce the magnitude of the amplification function by annihilating  $\rho_4(\theta)$ . This is achieved by setting  $\omega_4(0) = \omega_4(1) = 0$ ,  $\omega_4'(0) = \omega_4'(1) = 0$ , and  $\omega_4(\frac{1}{2}) = 0 =: \epsilon$ . Our numerical experiments have shown that achieving  $\omega_4(\frac{1}{2}) = 0$  is possible if one does not insist on the two tableaux being singly diagonal. But, if one insists on  $A_0$  and  $A_1$  being singly diagonal, then one can only enforce  $\omega_4(\frac{1}{2})$  to be of order  $3.8 \times 10^{-5} \simeq \epsilon$  when  $\ell_0 = \ell_1 = 1$  and  $7.9 \times 10^{-5} \simeq \epsilon$  when  $\ell_0 = \ell_1 = 0$ .

There are altogether 48 unknowns (24 for each Butcher tableau), and there are altogether 35 design conditions in Table 1. Moreover, we restrict ourselves to singly diagonal arrays, i.e., we additionally require that

$$A_{22}^0 = A_{44}^0 = A_{66}^0, \quad A_{33}^1 = A_{55}^1 = A_{77}^1, \quad (22)$$

giving four additional devising conditions. The above undetermined system of 39 nonlinear equations can be solved. The results reported in Appendix A have been obtained by using the nonlinear

Table 1: Design conditions for six-stage, third-order, AIRK schemes

#cdts.	$i = 0$	$i = 1$	ref
12	$A_0 U = c$	$A_1 U = c$	(3c)
2	$b_0 c = \frac{1}{2}$	$b_1 c = \frac{1}{2}$	(14a)
2	$b_0 c^2 = \frac{1}{3}$	$b_1 c^2 = \frac{1}{3}$	(14b)
2	$b_0 A_0 c = \frac{1}{6}$	$b_1 A_1 c = \frac{1}{6}$	(14c)
6	$\omega_3^0 = \omega_4^0 = \omega_5^0 = 0$	$\omega_3^1 = \omega_4^1 = \omega_5^1 = 0$	(20a)
2	$\omega_2^0 = (1 - \ell_0) \tau_3^0$	$\omega_2^1 = (1 - \ell_1) \tau_3^1$	(20b)
2	$b_1 A_0 c = b_0 A_1 c = \frac{1}{6}$		(15)
4	$\omega_5'(0) = \omega_5''(0) = \omega_5'(1) = \omega_5''(1) = 0$		(18)
2	$\omega_4'(0) = \omega_4'(1) = 0$		–
1	$\omega_4(\frac{1}{2}) = \epsilon$		–

solver `nlsolve` in `julia`. As the problem is highly nonlinear, the algorithm is first run with  $\epsilon = 0$  without enforcing (22). Then, one uses this solution as initialization to run the algorithm again with (22) but ignoring the constraint  $\omega_4(\frac{1}{2}) = 0$ . We refer the reader to Appendix A for two examples and some implementation details.

### 3.4 Companion ERK scheme

We now design a companion ERK scheme that can be used in combination with the above AIRK scheme in the IMEX setting. Therefore, we consider a third Butcher array in the form (we again omit the vector  $b_2$ )

$$\begin{array}{c|cccccc}
 0 & 0 & & & & & \\
 c_2 & A_{21}^2 & 0 & & & & \\
 c_3 & A_{31}^2 & A_{32}^2 & 0 & & & \\
 c_4 & A_{41}^2 & A_{42}^2 & A_{43}^2 & 0 & & \\
 c_5 & A_{51}^2 & A_{52}^2 & A_{53}^2 & A_{54}^2 & 0 & \\
 c_6 & A_{61}^2 & A_{62}^2 & A_{63}^2 & A_{64}^2 & A_{65}^2 & 0 \\
 1 & A_{71}^2 & A_{72}^2 & A_{73}^2 & A_{74}^2 & A_{75}^2 & A_{76}^2 \quad 0
 \end{array}$$

To obtain a third-order scheme, we enforce

$$A_2 U = c, \quad b_2 c = \frac{1}{2}, \quad b_2 c^2 = \frac{1}{3}, \quad b_2 A_2 c = \frac{1}{6}, \quad (23)$$

together with the coupling conditions

$$b_2 A_0 c = b_0 A_2 c = b_2 A_1 c = b_1 A_2 c = \frac{1}{6}. \quad (24)$$

This gives altogether 13 conditions for 21 unknowns. In some cases, we enforce the following three conditions to achieve linear order four:

$$b_2 c^3 = \frac{1}{4}, \quad b_2 A_2 c^2 = \frac{1}{12}, \quad b_2 A_2 c = \frac{1}{24}. \quad (25)$$

The resulting undetermined set of 13 or 16 nonlinear equations can be solved. We refer the reader to Appendix A for two examples obtained by using the nonlinear solver `nlsolve` in `julia`.

## 4 Numerical experiments

In this section, we illustrate numerically the performance of the method described §3 using the Butcher tableaux given in Appendix A. All the tests reported in this section are done in double precision.

### 4.1 ODEs

We start illustrating the proposed method by solving the following  $2 \times 2$  system of ODEs:

$$\partial_t \mathbf{U}(t) = L(\mathbf{U}(t)) + \mathbf{F}(t), \quad \mathbf{U}(0) = \mathbf{U}^0 \in \mathbb{R}^2, \quad (26)$$

where  $L := L_0 + L_1$  with  $L_0 := -P_0 D_0 P_0^{-1}$ ,  $L_1 := -P_1 D_1 P_1^{-1}$ , and

$$P_0 := \begin{pmatrix} 1 & 3 \\ 3 & -1 \end{pmatrix}, \quad D_0 := \begin{pmatrix} 0.023 & 0 \\ 0 & 0.073 \end{pmatrix}, \quad (27a)$$

$$P_1 := \begin{pmatrix} 2 & -3 \\ -1 & -1 \end{pmatrix}, \quad D_1 := \begin{pmatrix} 0.024 & 0 \\ 0 & 0.1345 \end{pmatrix}. \quad (27b)$$

The two matrices  $L_0$  and  $L_1$  do not commute. More precisely, denoting  $\|\cdot\|_{\text{Fr}}$  the Frobenius norm, we have  $2\|L_0 L_1 - L_1 L_0\|_{\text{Fr}} / \|L_0 + L_1\|_{\text{Fr}} \simeq 0.74$ . The matrix  $L$  is diagonalizable, and its two eigenvalues are approximately  $\lambda_0 \approx -0.085$ ,  $\lambda_1 \approx -0.17$ . Denoting  $L = P D P^{-1}$  the diagonal decomposition of  $L$ , and  $\mathbf{P}_0, \mathbf{P}_1$  the two columns of  $P$ , we initialize the system with  $\mathbf{U}^0 := \mathbf{P}_0 + 3\mathbf{P}_1$ . When  $\mathbf{F} \equiv 0$ , the exact solution to the autonomous system is  $\mathbf{U}_{\text{auto}}(t) = \mathbf{P}_0 e^{\lambda_0 t} + 3\mathbf{P}_1 e^{\lambda_1 t}$ . We also construct a solution with a nonzero source by setting  $\mathbf{F}(t) := \partial_t \mathbf{W} - L(\mathbf{W}(t))$  with  $\mathbf{W}(t) := (\cos(t), \sin(2t))^T$ . In this case, the exact solution is  $\mathbf{U}_{\text{auto}}(t) + \mathbf{W}(t)$ .

**Remark 4.1** (Sources). *Notice that there is variety of choices to approximate the source term in (26). For instance, one can regroup  $L_0$  and  $\mathbf{F}$  or regroup  $L_1$  and  $\mathbf{F}$ . One can also consider a convex combination by regrouping  $L_0$  and  $\alpha\mathbf{F}$  and regrouping  $L_1$  and  $(1 - \alpha)\mathbf{F}$  for all  $\alpha \in [0, 1]$ . Finally, one can also treat  $\mathbf{F}$  by using the companion matrix  $A_2$  for the ERK scheme. The tests reported below are done by regrouping  $L_0$  and  $\mathbf{F}$ . No significant difference is observed when using any of the other choices (not shown here for brevity).*

We test the method using the decomposition  $L = L_0 + L_1$  and the Butcher tableaux from Section A.1. The problem is solved over the time interval  $[0, T]$  with  $T := 10$ . The  $\ell^2$ -norm of the error divided by the  $\ell^2$ -norm of  $\mathbf{U}^0$  is measured at  $T$  for various time steps  $\tau_i = 2^{-i}$ ,  $i \in \{0:9\}$ . The results are reported in Table 4.1 for the two solutions. Up to machine accuracy, we observe third-order convergence rates as expected.

### 4.2 Heat equation

We continue with the two-dimensional heat equation

$$\partial_t u(\mathbf{x}, t) - \mu \Delta u(\mathbf{x}, t) = f(\mathbf{x}, t), \quad (\mathbf{x}, t) \in D \times (0, T), \quad u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \mathbf{x} \in D, \quad (28)$$

supplemented with either Dirichlet or Neuman boundary conditions and  $\mu := 1$ .

Table 2:  $\ell^2$ -errors and convergence rates for the ODE system (26). Butcher tableaux from Section A.1

$i$	autonomous sol.		non-autonomous sol.	
	error	rate	error	rate
0	0.1381E-05	–	0.2062E-02	–
1	0.1690E-06	3.03	0.2119E-03	3.28
2	0.2090E-07	3.02	0.2522E-04	3.07
3	0.2598E-08	3.01	0.3112E-05	3.02
4	0.3239E-09	3.00	0.3875E-06	3.01
5	0.4043E-10	3.00	0.4837E-07	3.00
6	0.5054E-11	3.00	0.6043E-08	3.00
7	0.6673E-12	2.92	0.7552E-09	3.00
8	0.1222E-12	2.45	0.9437E-10	3.00
9	0.7246E-14	4.08	0.1181E-10	3.00

#### 4.2.1 The setting

The tests are done in the unit square  $D := (0, 1)^2$ . We test homogeneous Dirichlet and homogeneous Neumann boundary conditions. Using the notation  $\mathbf{x} := (x, y)$ , the two exact solutions we use are

$$u_{\text{Dir}}(\mathbf{x}, t) = (2 + \sin(t)) \sin(2\pi x) \sin(3\pi y) + 64x(1-x)y(1-y) \sin(x+y+t), \quad (29)$$

$$u_{\text{Neu}}(\mathbf{x}, t) = (2 + \sin(t)) \cos(2\pi x) \cos(3\pi y) + 4x^2(1.5-x)y^2(1.5-y)(2 + \sin(\pi t)). \quad (30)$$

We apply the operator-splitting method by using the directional decomposition  $\Delta = \partial_{xx} + \partial_{yy}$ , i.e.,  $L_0(v) = \partial_{xx}v$  and  $L_1(v) = \partial_{yy}v$ . Although, in this case, it is traditional to use finite differences to realize the approximation in space, we illustrate the method by using continuous finite elements. Let  $V_h$  be the said finite element space and  $\{\varphi_i\}_{i \in \mathcal{V}}$  be the associated shape functions. The set  $\mathcal{V}$  is used to enumerate the shape functions with  $\#\mathcal{V} = I$ . Let  $(g, h)_{L^2(D)} := \int_D g(\mathbf{x})h(\mathbf{x}) \, dx$  be the canonical inner product in  $L^2(D)$ . We define the bilinear forms  $a_0(u_h, v_h) := (\mu \partial_x u_h \partial_x v_h)_{L^2(D)}$  and  $a_1(u_h, v_h) := (\mu \partial_y u_h \partial_y v_h)_{L^2(D)}$ . Then we consider the semi-discrete problem consisting of seeking  $u_h \in C^1([0, T]; V_h)$  such that, for all  $t \in [0, T]$ ,

$$\partial_t(u_h(t), \varphi_i)_{L^2(D)} + a_0(u_h(t), \varphi_i) + a_1(u_h(t), \varphi_i) = (f(t), \varphi_i)_{L^2(D)}, \quad \forall i \in \mathcal{V}, \quad (31)$$

and  $u_h(\cdot, 0) = u_{0h}$ , where  $u_{0h}$  is some quasi-optimal approximation of  $u_0$  in  $V_h$ . Let  $\mathcal{M}$  be the mass matrix associated with the  $L^2(D)$ -inner product and  $\mathcal{S}_0, \mathcal{S}_1$  be the stiffness matrices associated with the bilinear forms  $a_0$  and  $a_1$ , respectively. Let  $\mathbf{F}(t)$  be the vector in  $\mathbb{R}^I$  with entries  $(f(t), \varphi_i)_{L^2(D)}$ . Then, setting  $u_h(\mathbf{x}, t) = \sum_{i \in \mathcal{V}} \mathbf{U}_i(t) \varphi_i(\mathbf{x})$ , the system (31) reduces to solving the ODE system

$$\mathcal{M} \partial_t \mathbf{U}(t) = \mathcal{S}_0 \mathbf{U}(t) + \mathcal{S}_1 \mathbf{U}(t) + \mathbf{F}(t). \quad (32)$$

We solve (32) using the method presented in the paper. We use continuous finite elements of degree 2 to match the third-order accuracy in time of the method. We recall that the theoretical convergence rate for quadratic elements is cubic in the  $L^2$ -norm, quadratic in the  $H^1$ -seminorm, and the Riesz projection of the solution to (28) is superconvergent in the  $H^1$ -seminorm up to third-order. We run the simulations up to  $T := \frac{1}{2}$  on six consecutively refined meshes.

#### 4.2.2 Approximation of source term

As mentioned in Remark 4.1, the source  $\mathbf{F}(t)$  in the ODE system (32) can be handled in a variety of ways. We investigate in this section the three methods discussed in Remark 4.1 to handle this situation. We show three series of tests using the Dirichlet solution (29). In the first series of tests, we treat  $\mathbf{F}(t)$  using the companion Butcher tableau  $A_2$ , i.e., we set  $L_2(t) := \mathbf{F}(t)$ . In the second

series, we regroup  $F(t)$  and  $\mathcal{S}_1 U(t)$  (i.e., we set  $L_1(t, U(t)) := \mathcal{S}_1 U(t) + F(t)$ ), and in the third series, we combine  $F(t)$  and  $\mathcal{S}_0 U(t)$  (i.e., we set  $L_0(t, U(t)) := \mathcal{S}_0 U(t) + F(t)$ ). In all the tests, we use the L-stable pair  $(A_0, A_1)$  from Section A.1. The Dirichlet solution (29) has been manufactured to amplify the phenomenon we are about to discuss now.

The results are reported in Table 3. We show both the relative  $L^2$ -norm and  $H^1$ -seminorm of the solution at the final time  $T = \frac{1}{2}$ . We observe a loss of convergence as the mesh is refined for the first and the second methods (see the orange columns in the table). The asymptotic convergence rate in the  $L^2$ -norm and  $H^1$ -seminorm for these two methods is  $\mathcal{O}(h^{2.25})$  and  $\mathcal{O}(h^{1.5})$ , respectively, instead of the optimal rates  $\mathcal{O}(h^3)$  and  $\mathcal{O}(h^2)$ . Visual inspection of the solutions reveals the formation of spurious boundary layers as often observed for many splitting methods when enforcing Dirichlet boundary conditions. On the other hand, we observe that the third method does not suffer from any order reduction (see the green column in the table). The convergence rate in the  $H^1$ -seminorm is even superconvergent, which is a clear indication that no spurious boundary layer appears.

Table 3: Source approximation.  $\mathbb{P}_2$  approximation of (28) with the Dirichlet solution (29)

$L_2(t) := F(t)$			$L_1(t, U(t)) := S_1(U(t)) + F(t)$			$L_0(t, U(t)) := S_0(U(t)) + F(t)$		
$I$	$L^2$ -err	rate	$L^2$ -err	rate	$L^2$ -err	rate	$L^2$ -err	rate
441	2.78E-03	–	3.32E-03	–	2.19E-03	–	2.19E-03	–
1681	2.58E-04	3.55	4.42E-04	3.01	1.23E-04	4.31	1.23E-04	4.31
6561	4.13E-05	2.69	8.67E-05	2.39	9.55E-06	3.76	9.55E-06	3.76
25921	8.45E-06	2.31	1.84E-05	2.26	1.30E-06	2.90	1.30E-06	2.90
103041	1.78E-06	2.26	3.91E-06	2.24	1.83E-07	2.84	1.83E-07	2.84
410881	3.76E-07	2.25	8.29E-07	2.24	2.44E-08	2.91	2.44E-08	2.91
$I$	$H^1$ -err	rate	$H^1$ -err	rate	$H^1$ -err	rate	$H^1$ -err	rate
441	1.21E-02	–	1.35E-02	–	1.16E-02	–	1.16E-02	–
1681	1.88E-03	2.79	2.99E-03	2.25	1.54E-03	3.01	1.54E-03	3.01
6561	3.95E-04	2.29	9.33E-04	1.71	1.90E-04	3.08	1.90E-04	3.08
25921	1.14E-04	1.81	3.17E-04	1.57	2.28E-05	3.08	2.28E-05	3.08
103041	3.68E-05	1.64	1.10E-04	1.54	2.85E-06	3.01	2.85E-06	3.01
410881	1.22E-05	1.59	3.80E-05	1.53	4.01E-07	2.84	4.01E-07	2.84

The possible loss of convergence is only observed for the Dirichlet problem. Systematic tests have shown that this phenomenon does not occur for the Neuman problem (not shown here). While an explanation of this phenomenon is still lacking, this series of test shows that sources that do not depend on the solution should be combined with the operator  $L_0$  (i.e., one should use the Butcher tableau  $A_0$  for the source). This latter approach is systematically used in the tests reported in the rest of the paper.

#### 4.2.3 L-stable vs. A-stable tableaux

Our next objective is to compare the performances of the two AIRK methods, i.e., the one using the  $L(\alpha)$ -stable tableaux (see Section A.1) and the one using the  $A(\alpha)$ -stable tableaux (see Section A.2). We report in Tables 4 and 5 the relative error in the  $L^2$ -norm and the relative error in the  $H^1$ -seminorm for the Dirichlet and the Neumann solutions, respectively.

Table 4:  $\mathbb{P}_2$  approximation of (28) with the Dirichlet solution (29)

$I$	A-stable		L-stable		A-stable		L-stable	
	$L^2$ -err	rate	$L^2$ -err	rate	$H^1$ -err	rate	$H^1$ -err	rate
441	2.02E-03	–	2.19E-03	–	1.14E-02	–	1.16E-02	–
1681	1.30E-04	4.10	1.23E-04	4.31	1.52E-03	3.01	1.54E-03	3.01
6561	1.65E-05	3.03	9.55E-06	3.76	1.91E-04	3.05	1.90E-04	3.08
25921	2.34E-06	2.85	1.30E-06	2.90	2.53E-05	2.95	2.28E-05	3.08
103041	3.15E-07	2.91	1.83E-07	2.84	4.30E-06	2.57	2.85E-06	3.01
410881	4.41E-08	2.84	2.44E-08	2.91	2.21E-06	0.96	4.01E-07	2.84

Table 5:  $\mathbb{P}_2$  approximation of (28) with the Neumann solution (30)

$I$	A-stable		L-stable		A-stable		L-stable		
	$L^2$ -err	rate	$L^2$ -err	rate	$H^1$ -err	rate	$H^1$ -err	rate	
441	3.50E-03	–	3.75E-03	–	441	1.80E-02	–	1.80E-02	–
1681	2.38E-04	4.02	2.33E-04	4.15	1681	2.81E-03	2.77	2.85E-03	2.75
6561	2.79E-05	3.15	1.84E-05	3.73	6561	4.53E-04	2.68	4.50E-04	2.71
25921	3.95E-06	2.85	2.23E-06	3.07	25921	9.35E-05	2.30	7.39E-05	2.63
103041	5.01E-07	2.99	2.98E-07	2.92	103041	1.56E-05	2.60	1.26E-05	2.56
410881	6.39E-08	2.98	3.91E-08	2.94	410881	2.74E-06	2.51	2.19E-06	2.53

We observe third-order accuracy in the  $L^2$ -norm for both the L( $\alpha$ )-stable and the A( $\alpha$ )-stable methods and for both the Dirichlet and the Neumann problems. The approximation is again superconvergent in the  $H^1$ -seminorm. We notice a slight loss of convergence in the  $H^1$ -seminorm on the finest meshes for the Dirichlet problem using the method with the A( $\alpha$ )-stable tableaux. This effect is not observed for the method with the L( $\alpha$ )-stable tableaux. Overall, the method with the two L( $\alpha$ )-stable tableaux is slightly more accurate than that with the two A( $\alpha$ )-stable tableaux. In the remainder of the paper, we only report the results obtained with the L( $\alpha$ )-stable tableaux for brevity.

### 4.3 Heat equation coupled with (non)linear transport

Here, we consider the heat equation augmented with a transport term treated explicitly. This term can be either linear or nonlinear.

#### 4.3.1 Linear transport

We start by characterizing the convergence properties of the companion tableaux  $A_2$  presented in Appendix A by solving the linear transport equation

$$\partial_t u + \mathbf{v} \cdot \nabla u = 0, \quad (\mathbf{x}, t) \in D \times (0, T), \quad u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \mathbf{x} \in D, \quad (33)$$

supplemented with Dirichlet boundary conditions at the inflow boundary  $\partial D^- := \{\mathbf{x} \in \partial D \mid \mathbf{v}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) < 0\}$ . We consider  $D := (0, 1)^2$  and  $\mathbf{v}(\mathbf{x}) := (1, 1)^\top$ . The initial data is  $u_0(\mathbf{x}) := \exp((r(\mathbf{x})^2 + 2a^2)/(r(\mathbf{x})^2 - a^2))$  for  $r(\mathbf{x}) \leq a$  and  $u_0(\mathbf{x}) = 0$  otherwise, with  $r(\mathbf{x}) := \|\mathbf{x} - \mathbf{x}_0\|_{\ell^2}$ ,  $\mathbf{x}_0 := (\frac{1}{4}, \frac{1}{4})^\top$  and  $\mathbf{a} := 0.2$ . The exact solution is  $u(\mathbf{x}, t) = u_0(\mathbf{x} - \mathbf{v}t)$ .

Table 6:  $\mathbb{P}_1$  and  $\mathbb{P}_3$  approximation of (33) using the  $A_2$  companion tableaux.

$\mathbb{P}_1$ , A-st			$\mathbb{P}_3$ , A-st			$\mathbb{P}_3$ , L-st, 3rd		$\mathbb{P}_3$ , L-st, 4-th	
$I$	$L^2$ -err	rate	$I$	$L^2$ -err	rate	$L^2$ -err	rate	$L^2$ -err	rate
121	5.55E-01	–	961	9.17E-02	–	9.96E-02	–	9.65E-02	–
441	1.58E-01	1.94	3721	1.41E-02	2.76	1.54E-02	2.76	1.52E-02	2.73
1681	3.09E-02	2.44	14641	1.43E-03	3.34	1.87E-03	3.08	1.47E-03	3.41
6561	4.76E-03	2.75	58081	1.01E-04	3.85	2.13E-04	3.15	9.97E-05	3.91
25921	4.79E-04	3.34	231361	4.64E-06	4.45	2.53E-05	3.08	4.36E-06	4.53
103041	3.46E-05	3.81	923521	2.30E-07	4.34	3.16E-06	3.01	2.06E-07	4.41

We run the simulations using continuous finite elements up to the final time  $T := \frac{1}{2}$ . We test the three tableaux  $A_2$  from Appendix A using  $\mathbb{P}_1$  and  $\mathbb{P}_3$  finite elements. Recall that we have three tableaux  $A_2$  at our disposal. One for the A-stable pair (see Section A.2), which is fourth-order accurate, and two for the L-stable pair (see Section A.1), one which is third-order accurate and the other which is fourth-order accurate but with a smaller stability region. The results are shown in Table 6. We report the relative  $L^2$ -norm of the error at  $T = \frac{1}{2}$ . The results reported in the first and second tables are obtained with the tableau  $A_2$  associated with the A-stable pair. The

results shown in the third table are obtained with the third-order tableau  $A_2$  associated with the L-stable pair, and the results in the fourth table are obtained with the fourth-order tableau  $A_2$  also associated with the L-stable pair. The expected convergence rate is observed in all the cases.

### 4.3.2 Burgers-like nonlinear transport equation

We now focus our attention on a variation of the viscous Burgers equation

$$\partial_t u - \mu \Delta u + \nabla \cdot \mathbf{f}(u) = 0, \quad \mathbf{x} \in D_\infty, \quad t > 0, \quad (34)$$

in the semi-infinite domain  $D_\infty := \mathbb{R} \times (0, 1)$ , with the flux  $\mathbf{f}(u) := (u(1 - u), 0)^\top$ . We enforce homogeneous Neumann boundary conditions on the top and bottom boundaries of the domain. Setting  $\mathbf{x} := (x, y)$ , we also enforce  $\lim_{x \rightarrow -\infty} u(\mathbf{x}, t) = u_L$  and  $\lim_{x \rightarrow +\infty} u(\mathbf{x}, t) = u_R$ . We use the initial data

$$u_0(\mathbf{x}) := \bar{u} + \delta \tanh\left(\frac{\delta}{\mu}(x - x_0)\right), \quad \bar{u} := \frac{1}{2}(u_L + u_R), \quad \delta := \frac{1}{2}(u_R - u_L). \quad (35)$$

The solution to this Cauchy problem is a wave moving at speed  $s := 1 - 2\bar{u}$ ,

$$u(\mathbf{x}, t) = u_0(\mathbf{x} - \mathbf{s}t) \quad \text{with} \quad \mathbf{s} := (s, 0). \quad (36)$$

We set  $u_L := -1$  and  $u_R := 1$  in the tests reported below so that  $s = 1$ . We also set  $\mu = 0.01$ .

Table 7:  $\mathbb{P}_1$  and  $\mathbb{P}_3$  approximation of (34) using the L-stable tableaux  $(A_0, A_1)$  and the third-order companion tableau  $A_2$ .

$\mathbb{P}_1$			$\mathbb{P}_2$			$\mathbb{P}_3$		
$I$	$L^2$ -err	rate	$I$	$L^2$ -err	rate	$I$	$L^2$ -err	rate
121	4.11E-01	–	441	1.30E-01	–	961	6.73E-02	–
441	1.34E-01	1.73	1681	3.15E-02	2.12	3721	4.92E-03	3.87
1681	4.41E-02	1.66	6561	2.95E-03	3.48	14641	1.78E-03	1.49
6561	8.90E-03	2.35	25921	7.06E-04	2.08	58081	1.83E-04	3.30
25921	2.23E-03	2.02	103041	5.55E-05	3.69	231361	1.09E-05	4.08

We run the simulations in the truncated domain  $D := (0, 1)^2$  up to the final time  $T := \frac{1}{2}$  using continuous finite elements of degree 1, 2, and 3. We also use the decomposition  $L_0(u) := \mu \partial_{xx} u$ ,  $L_1(u) := \mu \partial_{yy} u$ , and  $L_2(u) := -\partial_x(\frac{1}{2}u^2)$ . We compute the relative  $L^2$ -norm of the error at  $T = \frac{1}{2}$ . The results are reported in Table 7. For the sake of brevity, we show the results only for the L-stable pair  $(A_0, A_1)$  with the third-order companion tableau  $A_2$  from Section A.1. We observe again that the expected convergence rates are achieved for all the polynomial degrees. The rate is close to 2 for the  $\mathbb{P}_1$  approximation and ranges between 2 and 3.5 for the  $\mathbb{P}_2$  and  $\mathbb{P}_3$  approximations.

### 4.3.3 Nonconservative nonlinear transport equation

We finally consider a nonlinear advection-diffusion equation with a nonconservative transport term. We use the Cole–Hopf transformation to manufacture the solution. We first set

$$w(\mathbf{x}, t) := 2 + \mu + \sin(m\pi x) \sin(n\pi y) e^{-kt}, \quad (37)$$

with  $m := 3$ ,  $n := 2$ ,  $k := \mu(m^2 + n^2)\pi^2$ . Notice that the function  $w$  solves the heat equation  $\partial_t w - \mu \Delta w = 0$  and that  $w(\mathbf{x}, t) \geq 1 + \mu > 1$  for all  $\mathbf{x}$  and all  $t$ . We then set  $u = -\mu \log(w)$ . The scalar field  $u(\mathbf{x}, t)$  solves the nonlinear transport equation

$$\partial_t u - \mu \Delta u + \mathbf{v} \cdot \nabla (\frac{1}{2}u^2) = 0, \quad (\mathbf{x}, t) \in D \times (0, T), \quad (38)$$



Table 8:  $\mathbb{P}_1$  and  $\mathbb{P}_3$  approximation of (37) using the L-stable tableaux  $(A_0, A_1)$  and the third-order companion tableau  $A_2$ .

$\mathbb{P}_1$			$\mathbb{P}_2$			$\mathbb{P}_3$		
$I$	$L^2$ -err	rate	$I$	$L^2$ -err	rate	$I$	$L^2$ -err	rate
121	1.80E-02	–	441	4.95E-04	–	961	4.44E-05	–
441	5.08E-03	1.96	1681	3.39E-05	4.01	3721	2.76E-06	4.10
1681	1.31E-03	2.03	6561	2.17E-06	4.04	14641	1.76E-07	4.02
6561	3.29E-04	2.03	25921	1.37E-07	4.03	58081	1.25E-08	3.85
25921	8.24E-05	2.02	103041	8.60E-09	4.01	231361	1.49E-09	3.07
103041	2.06E-05	2.01	410881	5.90E-10	3.87	923521	2.86E-10	2.39

with the space-time-dependent velocity  $\mathbf{v} := \frac{1}{w \log(w)} \nabla w$ .

We solve (37) in the unit square  $D := (0, 1)^2$  using the decomposition  $L_0(u) := \mu \partial_{xx} u$ ,  $L_1(u) := \mu \partial_{yy} u$ , and  $L_2(t, u) := -\mathbf{v}(\cdot, t) \cdot \nabla (\frac{1}{2} u^2)$ . We run the simulations with  $\mu := 0.01$  up to  $T := \frac{1}{2}$ . The results are reported in Table 8. For the sake of brevity, we only show the results for the L-stable pair  $(A_0, A_1)$  with the third-order companion tableau  $A_2$  from Section A.1. Here again, we observe the expected convergence rates. We also observe that the  $\mathbb{P}_2$  and  $\mathbb{P}_3$  approximations are superconvergent.

## A Two examples of six-stage, third-order schemes

In this section, we present two examples of six-stage, third-order RK schemes. Each example comprises an AIRK scheme (based on two implicit, singly diagonal RK schemes) and a companion ERK scheme. In the first example, the two constitutive implicit schemes are L( $\alpha$ )-stable (i.e.,  $\ell_0 = \ell_1 = 0$ ), whereas they are only A( $\alpha$ )-stable in the second example with  $\ell_0 = \ell_1 = 1$ . We focus on the equidistributed choice  $c_m = \frac{m-1}{6}$  for all  $m \in \{1:7\}$  for the time index array. This has the advantage of maximizing the efficiency of the ERK scheme; see Shu and Osher [28] and the discussion in [9].

For both examples, we solve first the design conditions identified in Section 3.3 to obtain the AIRK scheme. Recall that there are 39 conditions for 48 unknowns. Then, we solve the conditions identified in Section 3.4 to obtain the companion ERK scheme. Recall that there are 13 or 16 conditions for 21 unknowns depending on the accuracy one wants to reach for the ERK scheme. It turns out that for the L( $\alpha$ )-stable schemes, the third-order ERK array leads to a larger stability region than the fourth-order one. This is why we present the two possibilities. On the other hand, for the A( $\alpha$ )-stable schemes, the tableau  $A_2$  can be computed to ensure either third- or fourth-order accuracy, both with a rather large stability region.

In all cases, the resulting sets of coupled nonlinear equations are solved using the nonlinear solver `nlsolve` in `julia`. The optimization is done in quadruple precision for the L-stable tableaux (i.e., `BigFloat` numbers), and double precision for the A-stable tableaux. The residuals associated with the design conditions are less than  $10^{-22}$  for the L-stable tableaux and  $10^{-17}$  for the A-stable tableaux. For the AIRK scheme, solving from scratch the coupled nonlinear equations is somewhat challenging. Thus, the solution procedure employs an iterative fixed-point strategy, where the array  $A_0$  is designed given an array  $A_1$  and vice versa, until the prescribed tolerance is achieved.

The resulting Butcher arrays are reported in the following two sections. We only give the arrays  $A_0, A_1, A_2$  since the line vectors  $b_0, b_1, b_2$  are the last row of the associated array and are never used; see (5). To facilitate the reading, we also indicate for each row  $m \in \{1:7\}$  the value of the coefficient  $c_m$

### A.1 Example 1: $L(\alpha)$ -stable schemes

In this section, we give the  $L(\alpha)$ -stable arrays  $A_0$  and  $A_1$ , together with two possibilities for the companion array  $A_2$  mentioned above (one giving third-order and one giving linear order four). All the arrays are obtained using quadruple precision in `julia`. The accuracy on the design conditions is  $10^{-22}$ . The half-angle of the cone for  $A(\alpha)$ -stability is  $\alpha \approx 75^\circ$ . The amplification functions are illustrated in Figures 1 and A.1 (recall that  $R_0(z) = R_1(z)$  for singly diagonal tableaux, see Remark 3.6).

(i) Array  $A_0$

0	0				
1	0.007682766677990120	0.158983899988676547			
1	0.015365533395673803	0.317967799937659530	0		
1	0.067134743376864802	0.338274603424258278	-0.064393246789799627	...	
1	0.179050077617480914	0.169386371595552944	-0.216637439810267733	...	
1	0.201408968898570210	-0.018586441143895167	0.081249411695151912	...	
1	0.055256411220552875	-0.205127582453523036	1.186467117918441255	...	
1	...	0.158983899988676547			
1	...	0.534867657263900542	0		
1	...	0.477549665944474862	-0.067272172049645030	0.158983899988676547	
1	...	-0.381199971239714302	-0.252773137564567394	0.597377162118810602	0

(ii) Array  $A_1$

0	0				
1	0.16666666666666667	0			
1	0.08798574877573975	0.086363684567082812	0.158983899988676547		
1	0.148272588694077508	0.123809962338217855	0.227917448967704637	0	
1	0.092684091881748154	0.127270401977042040	0.162221507266258003	...	
1	0.166157946222573266	0.125070105123173022	0.124434611239232582	...	
1	0.048973226160787361	0.171916361228143705	0.213459859384815078	...	
2	...	0.125506765552941923	0.158983899988676547		
1	...	0.184260860904362666	0.233409809843991798	0	
1	...	0.179406092880142377	0.227260560357434931	0	0.158983899988676547

(iii) Array  $A_2$ , third-order

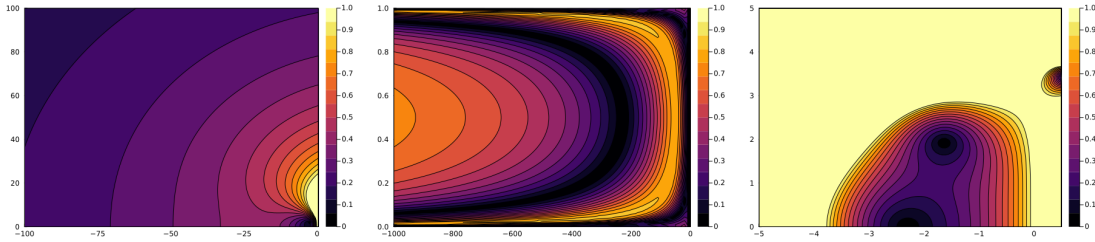
0	0				
1	0.16666666666666667	0			
1	-0.050619531693917875	0.383952865027251208	0		
1	0.115313313956073817	0.099138194215039115	0.285548491828887068	0	
1	0.065658564993170963	0.094245074373801537	0.202738372713947835	...	
1	0.062680510743166078	0.208831301672964596	0.168457244447138580	...	
1	0.187538570996657661	0.031430875635301389	0.109386484984970433	...	
2	...	0.304024654585746332	0		
1	...	0.182720713146197586	0.210643563323866492	0	
1	...	0.107869581266703755	0.392685024987187330	0.171089462129179432	0

(iv) Array  $A_2$ , linear order four

0	0				
1	0.16666666666666667	0			
2	−0.002065923995011051	0.335399257328344385	0		
3	0.009076043244499938	0.095774428321976104	0.395149528433523958	0	
4	0.268333342495086566	−0.084075704836160660	0.076139507867936172	...	
5	0.176995156036447256	0.003750298725649624	0.079363041718674150	...	
6	0.119787399084949175	−0.089727659939499215	0.661036648908505113	...	
7	...	0.406269521139804589	0		
8	...	0.337529406250193346	0.235695430602368957	0	
9	...	−0.142617977938011797	0.062099653483759240	0.389421936400297484	0

We show in the left panel of Figure 1 the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  (recall that  $R_0(z) = R_1(z)$  because the tableaux are singly diagonal). We show in the center panel the absolute value of the amplification function  $R_\theta(x)$  along the real negative x-axis, for  $x \leq 0$  and  $\theta \in [0, 1]$ ; see (9) for the definition of  $R_\theta(z)$ . The modulus of the amplification function  $R_2(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  for the explicit tableau giving third-order accuracy is shown in the right panel of the figure.

Figure 1: L-stable pair. Left: modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . Center: absolute value of the amplification function  $R_\theta(x)$  along the real negative x-axis, for  $x \leq 0$  and  $\theta \in [0, 1]$ . Right: modulus of the amplification function  $R_2(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  for the explicit tableau giving third-order accuracy.

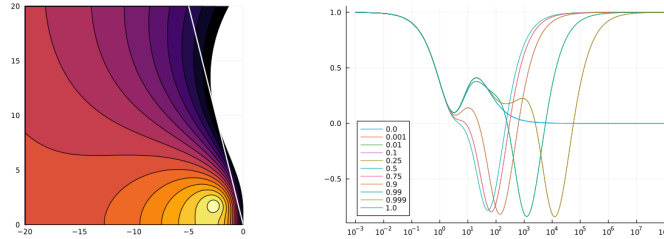


We show in the left panel of Figure A.1 a zoom close to the origin of the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . The modulus is larger than 1 only in the white region. We observe that  $A(\alpha)$ -stability holds for  $\alpha \approx 75^\circ$ . The white line materializes the limit of the stability cone. We show in the right panel of the figure the amplification function  $R_\theta(-x)$  for  $x \in [0, 10^8]$  and  $\theta \in \{0, 0.001, 0.01, 0.1, 0.25, 0.5, 0.75, 0.9, 0.99, 0.999, 1.0\}$ . We observe L-stability for the two extreme tableaux (i.e.,  $\theta \in \{0, 1\}$ ), and we observe  $A(0)$ -stability for all the intermediate values of  $\theta$ , as stated in Remark 3.7.

### A.2 Example 2: $A(\alpha)$ -stable schemes with $\ell_0 = \ell_1 = 1$

In this section, we give the A-stable arrays  $A_0$  and  $A_1$ , together with the companion array  $A_2$  giving linear order four. (Increasing the order from three to four does not affect the stability region of  $A_2$ ). All the arrays are obtained using double precision in `julia`. The accuracy on the design conditions is  $10^{-17}$ . The half-angle of the cone for  $A(\alpha)$ -stability is  $\alpha \approx 50^\circ$ . The amplification functions are illustrated in Figures A.2 and A.2 (recall that  $R_0(z) = R_1(z)$  for singly diagonal tableaux, see Remark 3.6).

Figure 2: L-stable pair. Left: zoom on the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . The modulus is larger than 1 in the white region only. Here,  $A(\alpha)$ -stability holds for  $\alpha \approx 75^\circ$ ; see the white dashed line. Right: amplification function  $R_\theta(-x)$  for  $x \in [0, 10^8]$  and  $\theta \in \{0, 0.001, 0.01, 0.1, 0.25, 0.5, 0.75, 0.9, 0.99, 0.999, 1.0\}$ .



(i) Array  $A_0$

0	0				
1	0	0.166666666666667			
2	0	0.333333333333333	0		
3	0.0881690356651937	0.2077230531651217	0.0374412445030180	...	
4	0.1912570743416719	0.0339232115988989	0.0809855895872098	...	
5	0.2217555743144974	-0.1981876469320450	0.4032535763162587	...	
6	-0.0181549513013415	-0.0576199238642526	1.1548881877024293	...	
7	...	0.166666666666667			
8	...	0.3605007911388862	0		
9	...	0.3112596743406823	-0.0714145113727266	0.166666666666667	
10	...	-0.4373955069083602	-0.2686190973268506	0.6269012916983754	0

(ii) Array  $A_1$

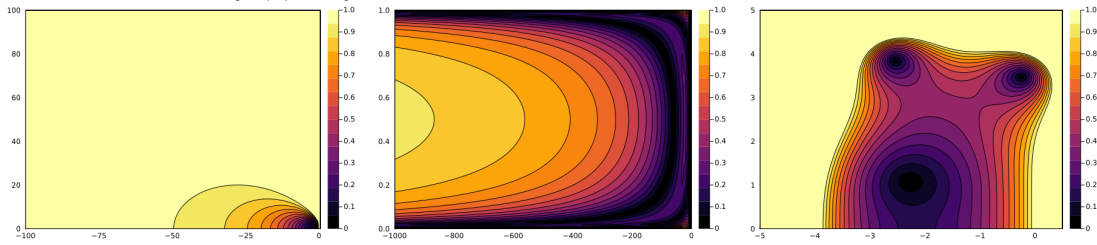
0	0			
1	0.166666666666667	0		
2	0.0961730695098136	0.0704935971568530	0.166666666666667	
3	0.3873667070462485	0.0334791581520742	0.0791541348016774	0
4	0.0482618178342044	0.0808153322470430	0.2741288261693861	...
5	0.3340345537873168	-0.0091489895287693	0.1060064658492590	...
6	0.0633044277927422	0.0951956813187544	0.3345863892872825	...
7	...	0.0967940237493665	0.166666666666667	
8	...	0.1479737995151694	0.2544675037103578	0
9	...	0.1253557996315356	0.2148910353030186	0
10	...			0.166666666666667

(iii) Array  $A_2$  (linear order four)

0	0			
1	0.166666666666667	0		
2	-0.0164974824288459	0.3498308157621792	0	
3	0.1757799381308423	0.0540524791927349	0.2701675826764229	0
4	-0.0229059377360897	0.1748847700986353	0.2836095136036662	...
5	0.0866385339448006	0.3019999712813553	0.1537929988619701	...
6	0.0471394455060848	0.1524277686616651	0.4188944702924878	...
7	...	0.2310783207004548	0	
8	...	-0.2072244075470651	0.4981262367922724	0
9	...	-0.1426444779083035	0.1831972427620590	0.3409855506860067
10	...			0

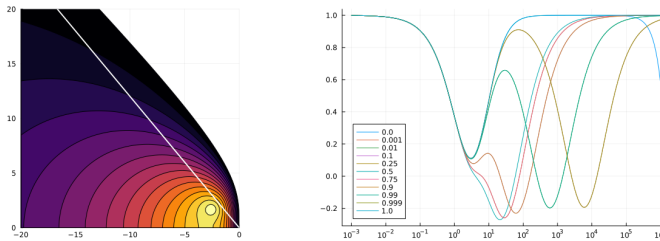
We show in the left panel of Figure A.2 the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  (recall that  $R_0(z) = R_1(z)$  because the tableaux are singly diagonal). We show in the center panel the absolute value of the amplification function  $R_\theta(x)$  along the real negative x-axis, for  $x \leq 0$  and  $\theta \in [0, 1]$ . The modulus of the amplification function  $R_2(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  for the explicit tableau is shown in the right panel of the figure.

Figure 3: A-stable pair. Left: modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . Center: absolute value of the amplification function  $R_\theta(x)$  along the real negative x-axis, for  $x \leq 0$  and  $\theta \in [0, 1]$ . Right: modulus of the amplification function  $R_2(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$  for the explicit tableau.



We show in the left panel of Figure A.2 a zoom close to the origin of the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . The modulus is larger than 1 only in the white region. We observe that  $A(\alpha)$ -stability holds for  $\alpha \approx 50^\circ$ . The white line materializes the limit of the stability cone. We show in the right panel of the figure the amplification function  $R_\theta(-x)$  for  $x \in [0, 10^6]$  and  $\theta \in \{0, 0.001, 0.01, 0.1, 0.25, 0.5, 0.75, 0.9, 0.99, 0.999, 1.0\}$ . We observe  $A(0)$ -stability for all the values of  $\theta$ . Since the tableaux  $A_0$  and  $A_1$  have been computed in double precision only,  $A$ -stability is numerically lost on the tableau  $A_0$  for  $x \geq 10^6$ . This technical problem can be resolved by using quadruple precision as we did for the L-stable tableaux. We have verified that  $A$ -stability still holds for all the other tableaux up to  $x = 10^{10}$ .

Figure 4: A-stable pair. Left: zoom on the modulus of the amplification function  $R_0(z)$  in the half complex plane  $\{\Re(z) \leq 0\}$ . The modulus is larger than 1 in the white region. Here,  $A(\alpha)$ -stability holds for  $\alpha \approx 50^\circ$ ; see the white dashed line. Right: amplification function  $R_\theta(-x)$  for  $x \in [0, 10^6]$  and  $\theta \in \{0, 0.001, 0.01, 0.1, 0.25, 0.5, 0.75, 0.9, 0.99, 0.999, 1.0\}$ .



### A.3 Some implementation details

In this section, we give some details on how the conditions on  $\omega_4(\theta)$  and  $\omega_5(\theta)$  can be implemented. We first observe that

$$\begin{aligned}\tau_1(\theta) &= \zeta_{10}(\theta)\tau_1^0 + \zeta_{01}(\theta)\tau_1^1, \\ \tau_2(\theta) &= \zeta_{20}(\theta)\tau_2^0 + \zeta_{11}(\theta)\tau_1^0\tau_1^1 + \zeta_{02}(\theta)\tau_2^1, \\ \tau_3(\theta) &= \zeta_{30}(\theta)\tau_3^0 + \zeta_{21}(\theta)\tau_2^0\tau_1^1 + \zeta_{12}(\theta)\tau_1^0\tau_2^1 + \zeta_{03}(\theta)\tau_3^1, \\ \tau_4(\theta) &= \zeta_{31}(\theta)\tau_3^0\tau_1^1 + \zeta_{22}(\theta)\tau_2^0\tau_2^1 + \zeta_{13}(\theta)\tau_1^0\tau_3^1, \\ \tau_5(\theta) &= \zeta_{32}(\theta)\tau_3^0\tau_2^1 + \zeta_{23}(\theta)\tau_2^0\tau_3^1, \\ \tau_6(\theta) &= \zeta_{33}(\theta)\tau_3^0\tau_3^1,\end{aligned}$$

with  $\zeta_{mn}(\theta) = (1-\theta)^m\theta^n$ . Furthermore, we give  $\frac{d^p}{d\theta^p}\beta_k(\theta)$ , for all  $k \in \{2, 3, 4\}$  and all  $p \in \{1, 2\}$  using the shorthand notation  $\delta b := b_1 - b_0$  and  $\delta A := A_1 - A_0$ :

$$\beta_3'(\theta) = \delta b A_\theta^2 c + b_\theta (A_\theta^2)' c, \quad (39a)$$

$$\beta_3''(\theta) = 2\delta b (A_\theta^2)' c + b_\theta (A_\theta^2)'' c, \quad (39b)$$

$$\beta_4'(\theta) = \delta b A_\theta^3 c + b_\theta (A_\theta^3)' c, \quad (39c)$$

$$\beta_4''(\theta) = 2\delta b (A_\theta^3)' c + b_\theta (A_\theta^3)'' c, \quad (39d)$$

$$\beta_5'(\theta) = \delta b A_\theta^4 c + b_\theta (A_\theta^4)' c, \quad (39e)$$

$$\beta_5''(\theta) = \delta b (A_\theta^4)' c + b_\theta (A_\theta^4)'' c. \quad (39f)$$

with

$$(A_\theta^2)' = \delta A A_\theta + A_\theta \delta A, \quad (40a)$$

$$(A_\theta^3)' = \delta A A_\theta^2 + A_\theta \delta A A_\theta + A_\theta^2 \delta A, \quad (40b)$$

$$(A_\theta^4)' = \delta A A_\theta^3 + A_\theta \delta A A_\theta^2 + A_\theta^2 \delta A A_\theta + A_\theta^3 \delta A, \quad (40c)$$

$$(A_\theta^2)'' = 2\delta A^2, \quad (40d)$$

$$(A_\theta^3)'' = 2(\delta A^2 A_\theta + \delta A A_\theta \delta A + A_\theta \delta A^2), \quad (40e)$$

$$\begin{aligned}(A_\theta^4)'' &= 2(\delta A^2 A_\theta^2 + \delta A A_\theta \delta A A_\theta + \delta A A_\theta^2 \delta A + A_\theta \delta A^2 A_\theta \\ &\quad + A_\theta \delta A A_\theta \delta A + A_\theta^2 \delta A^2).\end{aligned} \quad (40f)$$

Putting everything together gives

$$\begin{aligned}\omega_5'(\theta) &= \beta_5'(\theta) + \beta_4'(\theta)\tau_1(\theta) + \beta_4(\theta)\tau_1'(\theta) + \beta_3'(\theta)\tau_2(\theta) + \beta_3(\theta)\tau_2'(\theta) \\ &\quad + \frac{1}{6}\tau_3'(\theta) + \frac{1}{2}\tau_4'(\theta) + \tau_5'(\theta),\end{aligned} \quad (41a)$$

$$\begin{aligned}\omega_5''(\theta) &= \beta_5''(\theta) + \beta_4''(\theta)\tau_1(\theta) + 2\beta_4'(\theta)\tau_1'(\theta) + \beta_4(\theta)\tau_1''(\theta) \\ &\quad + \beta_3''(\theta)\tau_2(\theta) + 2\beta_3'(\theta)\tau_2'(\theta) + \beta_3(\theta)\tau_2''(\theta) \\ &\quad + \frac{1}{6}\tau_3''(\theta) + \frac{1}{2}\tau_4''(\theta) + \tau_5''(\theta),\end{aligned} \quad (41b)$$

$$\omega_4'(\theta) = \beta_4'(\theta) + \beta_3'(\theta)\tau_1(\theta) + \beta_3(\theta)\tau_1'(\theta) + \frac{1}{6}\tau_2'(\theta) + \frac{1}{2}\tau_3'(\theta) + \tau_4'(\theta). \quad (41c)$$

### A.4 GARK rewriting

In this section, we illustrate the rewriting of the above seven-stage AIRK schemes as combinations of four-stage schemes using the GARK formalism. Specifically, the AIRK scheme with the above

Butcher arrays rewrites in the format (6) upon setting

$$\begin{aligned}
\mathfrak{A}^{0,00} &= \begin{pmatrix} 0 & & & & \\ 0 & A_{22}^0 & & & \\ 0 & A_{42}^0 & A_{44}^0 & & \\ 0 & A_{62}^0 & A_{64}^0 & A_{66}^0 & \\ & & & & \end{pmatrix}, & \mathfrak{A}^{0,01} &= \begin{pmatrix} 0 & & & & \\ A_{21}^0 & 0 & & & \\ A_{41}^0 & A_{43}^0 & 0 & & \\ A_{61}^0 & A_{63}^0 & A_{65}^0 & 0 & \\ & & & & \end{pmatrix}, \\
\mathfrak{A}^{0,10} &= \begin{pmatrix} 0 & & & & \\ 0 & 0 & & & \\ 0 & A_{42}^1 & 0 & & \\ 0 & A_{62}^1 & A_{64}^1 & 0 & \\ & & & & \end{pmatrix}, & \mathfrak{A}^{0,11} &= \begin{pmatrix} 0 & & & & \\ A_{21}^1 & 0 & & & \\ A_{41}^1 & A_{43}^1 & 0 & & \\ A_{61}^1 & A_{63}^1 & A_{65}^1 & 0 & \\ & & & & \end{pmatrix}, \\
\mathfrak{A}^{1,00} &= \begin{pmatrix} 0 & & & & \\ 0 & A_{32}^0 & & & \\ 0 & A_{52}^0 & A_{54}^0 & & \\ 0 & A_{72}^0 & A_{74}^0 & A_{76}^0 & \\ & & & & \end{pmatrix}, & \mathfrak{A}^{1,01} &= \begin{pmatrix} 0 & & & & \\ A_{31}^0 & 0 & & & \\ A_{51}^0 & A_{53}^0 & 0 & & \\ A_{71}^0 & A_{73}^0 & A_{75}^0 & 0 & \\ & & & & \end{pmatrix}, \\
\mathfrak{A}^{1,10} &= \begin{pmatrix} 0 & & & & \\ 0 & A_{32}^1 & & & \\ 0 & A_{52}^1 & A_{54}^1 & & \\ 0 & A_{72}^1 & A_{74}^1 & A_{76}^1 & \\ & & & & \end{pmatrix}, & \mathfrak{A}^{1,11} &= \begin{pmatrix} 0 & & & & \\ A_{31}^1 & A_{33}^1 & & & \\ A_{51}^1 & A_{53}^1 & A_{55}^1 & & \\ A_{71}^1 & A_{73}^1 & A_{75}^1 & A_{77}^1 & \\ & & & & \end{pmatrix}.
\end{aligned}$$

## B Further remarks on AIRK schemes

In this appendix, we collect two results on four-stage, third-order and two-stage, second-order AIRK schemes, respectively.

### B.1 Four-stage, third-order, implicit RK schemes

In this section, we show that there is a barrier to design four-stage, third-order AIRK schemes. Indeed, the single-array RK schemes cannot be A-stable. We set  $s = 4$  since we consider four-stage schemes. Since our result concerns any single-array implicit RK scheme having only two nonzero diagonal coefficients, we drop in this section the subscripts and simply write  $A$  for the Butcher array and set  $b = e_5^T A$ .

**Lemma B.1** (Stability barrier on four-stage, third-order implicit RK schemes). *Assume that the matrix  $A \in \mathbb{R}^{5,5}$  is lower-triangular with only two nonzero diagonal entries, and that the RK scheme is of order three. Then,  $\lim_{|z| \rightarrow \infty} |R(z)| \geq 1 + \sqrt{3}$ .*

*Proof.* Adapting the arguments in the proof of Lemma 2.4, we infer that

$$\rho(z) := \det(I - zA)zb(I - zA)^{-1}U = \sum_{k \in \{0:3\}} \omega_k z^{k+1},$$

with  $\omega_0 = 1$  and (recall that  $\tau_l(A) = (-1)^l \text{tr}_l(A)$ )

$$\begin{aligned}
\omega_1 &= \frac{1}{2} + \tau_1(A), \\
\omega_2 &= \frac{1}{6} + \frac{1}{2}\tau_1(A) + \tau_2(A), \\
\omega_3 &= (bA^2c) + \frac{1}{6}\tau_1(A) + \frac{1}{2}\tau_2(A).
\end{aligned}$$

Moreover, reasoning as in the proof of Lemma 3.5, a necessary condition to achieve  $A(\alpha)$ -stability is

$$\omega_1 = (\ell - 1)\tau_2(A), \quad \ell \in [-1, 1], \quad \omega_2 = 0, \quad \omega_3 = 0.$$

The conditions on  $\omega_1$  and  $\omega_2$  determine  $\tau_1(A)$  and  $\tau_2(A)$ :

$$\tau_1(A) = \frac{1}{3} \frac{2 + \ell}{1 + \ell}, \quad \tau_2(A) = \frac{1}{6(1 + \ell)}.$$

The standard inequality  $\tau_1(A)^2 \geq 4\tau_2(A)$  gives  $(2 + \ell)^2 \geq 6(1 + \ell)$ , i.e.,  $\ell^2 - 2\ell - 2 \geq 0$ . This, in turn, requires  $\ell \geq 1 + \sqrt{3}$ , which contradicts  $\ell \in [-1, 1]$ .  $\square$

## B.2 Two-stage, second-order, implicit RK schemes

In this section, we show that any two-stage, second-order, implicit RK scheme having only one nonzero diagonal coefficient, say  $a$ , must satisfy  $a = \frac{1}{2}$  and  $\lim_{|z| \rightarrow \infty} R(z) = -1$ . We set  $s = 2$  since we consider two-stage schemes, and, as above, we simply write  $A$  for the Butcher array and set  $b = e_3^\top A$ .

**Lemma B.2.** *Assume that the matrix  $A \in \mathbb{R}^{3,3}$  is lower-triangular with only one nonzero diagonal entry, say  $a$ , and that the RK scheme is of order two. Then,  $a = \frac{1}{2}$ , and the amplification function is given by  $R(z) = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}$ , so that  $\lim_{|z| \rightarrow \infty} R(z) = -1$ .*

*Proof.* Reasoning as above shows that

$$\rho(z) := \det(I - zA)zb(I - zA)^{-1}U = \sum_{k \in \{0:1\}} \omega_k z^{k+1} = z + \left(\frac{1}{2} - a\right)z^2.$$

Since  $R(z) = 1 + \frac{\rho(z)}{1 - az}$ , a necessary condition for A-stability is  $\omega_2 = 0$ , i.e.,  $a = \frac{1}{2}$ . This readily gives  $R(z) = 1 + \frac{z}{1 - \frac{1}{2}z} = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}$ , so that  $\lim_{|z| \rightarrow \infty} R(z) = -1$ .  $\square$

**Remark B.3** (Combined amplification function). *Consider two two-stage, second-order, implicit RK schemes, one having the diagonal entry  $\frac{1}{2}$  on the second line and the other having the diagonal entry  $\frac{1}{2}$  on the third line. Reasoning as above, we readily obtain*

$$\rho_\theta(z) = zb_\theta \left( 1 + (A_\theta + \tau_1(A_\theta)I)z \right) U = z.$$

Hence,

$$R_\theta(z) = 1 + \frac{z}{(1 - \frac{1}{2}\theta z)(1 - \frac{1}{2}(1 - \theta)z)}.$$

We immediately recover that  $\ell_\theta = 1$  when  $\theta \notin \{0, 1\}$ , whereas  $\ell_0 = \ell_1 = -1$ .

Using the second-order conditions (namely (3c) together with  $bc = \frac{1}{2}$ ), we infer that the two implicit RK schemes take the form

$$\begin{array}{c|ccc} 0 & 0 & & \\ \gamma & \gamma - \frac{1}{2} & \frac{1}{2} & \\ 1 & 1 - \frac{1}{2\gamma} & \frac{1}{2\gamma} & 0 \\ \hline & 1 - \frac{1}{2\gamma} & \frac{1}{2\gamma} & 0 \end{array} \quad \begin{array}{c|ccc} 0 & 0 & & \\ \gamma & \gamma & 0 & \\ 1 & \frac{1}{2} & 0 & \frac{1}{2} \\ \hline & \frac{1}{2} & 0 & \frac{1}{2} \end{array} \quad (42)$$

with parameter  $\gamma \in (0, 1)$ . The most natural choice is  $\gamma = \frac{1}{2}$ , which leads, as expected, to the midpoint and Crank–Nicolson schemes.

## References

- [1] U. M. Ascher, S. J. Ruuth, and B. T. R. Wetton. Implicit-explicit methods for time-dependent partial differential equations. *SIAM J. Numer. Anal.*, 32(3):797–823, 1995.



- [2] U. M. Ascher, S. J. Ruuth, and R. J. Spiteri. Implicit-explicit Runge-Kutta methods for time-dependent partial differential equations. *Appl. Numer. Math.*, 25(2-3):151–167, 1997.
- [3] A. D. Brandrauk and H. Shen. Improved exponential split operator method for solving the time-dependent Schrödinger equation. *Chem. Phys. Lett.*, 176(5):428–432, 1991.
- [4] S. Blanes and F. Casas. On the necessity of negative coefficients for operator splitting schemes of order higher than two. *Appl. Numer. Math.*, 54(1):23–37, 2005.
- [5] F. Castella, P. Chartier, S. Descombes, and G. Vilmart. Splitting methods with complex times for parabolic equations. *BIT*, 49(3):487–508, 2009.
- [6] A. J. Christlieb, Y. Liu, and Z. Xu. High order operator splitting methods based on an integral deferred correction framework. *J. Comput. Phys.*, 294:224–242, 2015.
- [7] G. J. Cooper and A. Sayfy. Additive Runge-Kutta methods for stiff ordinary differential equations. *Math. Comp.*, 40(161):207–218, 1983.
- [8] J. Douglas, Jr. and H. H. Rachford, Jr. On the numerical solution of heat conduction problems in two and three space variables. *Trans. Amer. Math. Soc.*, 82:421–439, 1956.
- [9] A. Ern and J.-L. Guermond. Invariant-domain-preserving high-order time stepping: I. Explicit Runge-Kutta schemes. *SIAM J. Sci. Comput.*, 44(5):A3366–A3392, 2022.
- [10] Z. Gegechkori, J. Rogava, and M. Tsiklauri. The fourth order accuracy decomposition scheme for an evolution problem. *M2AN Math. Model. Numer. Anal.*, 38(4):707–722, 2004.
- [11] D. Goldman and T. J. Kaper. Nth-order operator splitting schemes and nonreversible systems. *SIAM J. Numer. Anal.*, 33(1):349–367, 1996.
- [12] S. González-Pinto, D. Hernández-Abreu, M. S. Pérez-Rodríguez, A. Sarshar, S. Roberts, and A. Sandu. A unified formulation of splitting-based implicit time integration schemes. *J. Comput. Phys.*, 448:Paper No. 110766, 22, 2022.
- [13] E. Hairer. Order conditions for numerical methods for partitioned ordinary differential equations. *Numer. Math.*, 36(4):431–445, 1980/81.
- [14] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations. II. Stiff and Differential-algebraic Problems*, volume 14 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2010. Second revised edition, paperback.
- [15] E. Hansen and A. Ostermann. High order splitting methods for analytic semigroups exist. *BIT*, 49(3):527–542, 2009.
- [16] C. A. Kennedy and M. H. Carpenter. Additive Runge-Kutta schemes for convection-diffusion-reaction equations. *Appl. Numer. Math.*, 44(1-2):139–181, 2003.
- [17] G. I. Marchuk. Splitting and alternating direction methods. In *Handbook of numerical analysis, Vol. I*, Handb. Numer. Anal., I, pages 197–462. North-Holland, Amsterdam, 1990.
- [18] L. Pareschi and G. Russo. Implicit-explicit Runge-Kutta schemes for stiff systems of differential equations. In *Recent trends in numerical analysis*, volume 3 of *Adv. Theory Comput. Math.*, pages 269–288. Nova Sci. Publ., Huntington, NY, 2001.

- [19] L. Pareschi and G. Russo. Implicit-Explicit Runge-Kutta schemes and applications to hyperbolic systems with relaxation. *J. Sci. Comput.*, 25(1-2):129–155, 2005.
- [20] D. W. Peaceman and H. H. Rachford, Jr. The numerical solution of parabolic and elliptic differential equations. *J. Soc. Indust. Appl. Math.*, 3:28–41, 1955.
- [21] P. Rentrop. Partitioned Runge-Kutta methods with stiffness detection and stepsize control. *Numer. Math.*, 47(4):545–564, 1985.
- [22] J. R. Rice. Split Runge-Kutta method for simultaneous equations. *J. Res. Nat. Bur. Standards Sect. B*, 64B:151–170, 1960.
- [23] S. Roberts, J. Loffeld, A. Sarshar, C. S. Woodward, and A. Sandu. Implicit multirate GARK methods. *J. Sci. Comput.*, 87(1):Paper No. 4, 32, 2021.
- [24] H. H. Rosenbrock. Some general implicit processes for the numerical solution of differential equations. *The Computer Journal*, 5(4):329–330, 01 1963.
- [25] A. Sandu and M. Günther. A generalized-structure approach to additive Runge-Kutta methods. *SIAM J. Numer. Anal.*, 53(1):17–42, 2015.
- [26] A. Sarshar, S. Roberts, and A. Sandu. Alternating directions implicit integration in a general linear method framework. *J. Comput. Appl. Math.*, 387:Paper No. 112619, 13, 2021.
- [27] Q. Sheng. Solving linear partial differential equations by exponential splitting. *IMA J. Numer. Anal.*, 9(2):199–212, 1989.
- [28] C.-W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock-capturing schemes. *J. Comput. Phys.*, 77(2):439 – 471, 1988.
- [29] R. J. Spiteri and S. Wei. Fractional-step Runge-Kutta methods: representation and linear stability analysis. *J. Comput. Phys.*, 476:Paper No. 111900, 18, 2023.
- [30] G. Strang. On the construction and comparison of difference schemes. *SIAM J. Numer. Anal.*, 5:506–517, 1968.
- [31] M. Suzuki. General theory of fractal path integrals with applications to many-body theories and statistical physics. *J. Math. Phys.*, 32(2):400–407, 1991.
- [32] O. B. Widlund. A note on unconditionally stable linear multistep methods. *Nordisk Tidskr. Informationsbehandling (BIT)*, 7:65–70, 1967.
- [33] N. N. Yanenko. *The method of fractional steps. The solution of problems of mathematical physics in several variables.* Springer-Verlag, New York-Heidelberg, 1971. Translated from the Russian by T. Cheron. English translation edited by M. Holt.
- [34] X. Zhong. Additive semi-implicit Runge-Kutta methods for computing high-speed nonequilibrium reactive flows. *J. Comput. Phys.*, 128(1):19–31, 1996.