#### **PINT 2019**



# Toward error estimates for general space-time discretizations of the advection equation

Martin J. Gander<sup>1</sup> · Thibaut Lunet<sup>1</sup>

Received: 8 December 2019 / Accepted: 8 June 2020 / Published online: 18 October 2020 © The Author(s) 2020

#### **Abstract**

We develop new error estimates for the one-dimensional advection equation, considering general space-time discretization schemes based on Runge–Kutta methods and finite difference discretizations. We then derive conditions on the number of points per wavelength for a given error tolerance from these new estimates. Our analysis also shows the existence of synergistic space-time discretization methods that permit to gain one order of accuracy at a given CFL number. Our new error estimates can be used to analyze the choice of space-time discretizations considered when testing Parallel-in-Time methods.

**Keywords** Advection equation · Space-time discretization · Error estimates · Parallel-in-time integration · Parareal

#### 1 Introduction

Over the last decade, Parallel-in-Time (PinT) methods have received sustained attention in the numerical analysis research community, because of the new scientific challenges coming with the ever growing parallel super-computing architectures; for a review, see [12]. Many PinT algorithms are known to work well for parabolic problems (see, e.g. [20,30,37]), but the PinT community still struggles with hyperbolic problems. This motivated many contributions in the literature, some of which considered solutions based on PARAREAL, a well known PinT algorithm proposed in [28], see for instance [3,4,6,29,32,35]. A geometric multigrid interpretation of PARAREAL was given in [20], followed by the genesis of the MGRIT algorithm [10], based on algebraic multigrid methods. MGRIT was first presented during a student paper competition at the Sixteenth Copper Mountain Conference on Multigrid Methods, from which [10] originates. A journal version of this work was also published later in [9]. In contrast to [20], in the first paper for MGRIT the authors consider a Full Approximation Scheme (FAS) interpretation of PARAREAL, and based on this introduce a variant of MGRIT, with additional coarse and fine relax-

Communicated by Robert Speck.

☐ Thibaut Lunet thibaut.lunet@unige.ch

Section de Mathématiques, University of Geneva, 2-4 rue du Liévre, Case postale 64, 1211 Genéve 4, Switzerland ation steps (FCF-relaxation, in contrast with F-relaxation only for the basic MGRIT method). It was then proved in [19] that this variant with FCF-relaxation is identical to PARA-REAL with overlap, where each additional *CF* relaxation adds one coarse time interval more in the overlap in PARAREAL. Many recent contributions focused on applying MGRIT to hyperbolic problems, see e.g. [5,25,26]. There are however also other PinT algorithms especially designed for hyperbolic problems, see e.g. ParaExp [13,14], the diagonalization technique [18], and waveform relaxation [16,17,39,40], and combinations thereof, see e.g. [21], based on earlier work in [43].

The advection equation has proven to be a critical test problem for PinT algorithms ([2,11,26]). Many space-time discretizations can be used to solve the advection equation, some being more diffusive than others, and numerical diffusion has proven to help convergence of PinT algorithms like PARAREAL [34]. It thus seems that a compromise between numerical diffusion and accuracy of the numerical scheme is necessary for effective convergence of such PinT methods applied to the advection equation (and hyperbolic problems in general). Hence, care must be taken when validating PinT algorithms on the advection equation. In particular, when solving such problems with a PinT algorithm, it is important to know how much accuracy the computed solution actually contains, especially since PinT algorithms should be capable of computing over longer time intervals accurate solutions in parallel. It is therefore important to have reliable error estimates that permit informed decisions on which space-time



16 Page 2 of 14 M. J. Gander, T. Lunet

discretization schemes can be chosen for solving the advection problem as a test problem for PinT algorithms.

There are already several error estimates available in the literature for classical space-time discretizations. For instance, in [24, Chapter 2], the authors list many error estimates for high order space discretizations, and derive also conditions on the number of points per wavelength needed to achieve a given level of accuracy. Similar results can also be found in [23,38]. The existing error bounds in the literature have however three major drawbacks:

- Most of them focus on the phase error, and estimate the  $L_1$  error of the numerical solution in Fourier space.
- They involve only specific space discretizations (e.g. centered finite-differences in space), and there is no general methodology for arbitrary space-time discretizations.
- Very few studies also include the impact of the time discretization, and the influence of important numerical parameters, like the Courant-Friedrich-Lewy (CFL) number, are still not completely understood.

Our goal here is to complement the currently available error estimates in the literature, and to provide a general methodology to estimate the error in the approximate solution, depending on the space and time discretization scheme used. We will also derive from these new general error estimates conditions on the minimum number of points per wavelength needed to achieve a given error tolerance in space. Finally, our new error estimates allowed us to discover the existence of numerical schemes with the same order p in space and time that have together a p+1 order accuracy if used with a given CFL number.

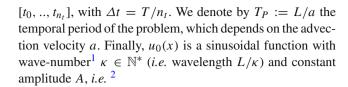
Our manuscript is organized as follows: we state and prove our main theoretical results in Sect. 2. In Sect. 3, we perform numerical experiments to validate the error estimates and illustrate the main consequences they imply. In Sect. 4, we discuss the application of some space-time discretizations for PARAREAL in view of our new error estimates. Finally, we give a conclusion and brief outlook on future work in Sect. 5.

# 2 Discrete fourier analysis

We consider the advection equation

$$\frac{\partial u}{\partial t} = -a \frac{\partial u}{\partial x}, \quad u(x,0) = u_0(x), \tag{1}$$

where a is constant, solved numerically on a periodic spatial domain [0, L], discretized using a uniform mesh  $[0, \Delta x, ..., L - \Delta x] = [x_1, ..., x_{n_x}]$  of  $n_x$  points  $(\Delta x = L/n_x)$ . Numerical time integration is performed decomposing [0, T] using  $n_t + 1$  points in time,  $[0, \Delta t, ..., T] =$ 



$$u_0(x) = A\cos\left(2\pi\kappa\frac{x}{L}\right). \tag{2}$$

#### 2.1 Main theorem

To present our main error estimate for an arbitrary spacetime discretization of the advection equation (1), we need to introduce our measure for the error and two key properties that characterize an arbitrary space and time discretization.

**Definition 1** Let the numerical solution of (1) be represented at time t by

$$\mathbf{u}(t) := [u_1(t), ..., u_{n_x}(t)] \approx [u(x_1, t), ..., u(x_{n_x}, t)].$$
 (3)

We define its discrete norm by

$$\|\mathbf{u}\| := \sqrt{\frac{1}{n_x} \sum_{j=1}^{n_x} u_j(t)^2}.$$
 (4)

The choice of this particular norm will be discussed later in Sect. 2.4, in view of the main results of this paper and their proof.

We consider the use of a finite difference discretization of order p for the space derivative in (1), whose Fourier symbol is<sup>3</sup>

$$z_{adv}(\theta) = i\theta + \alpha_{p+1}\theta^{p+1} + \mathcal{O}(\theta^{p+2}), \quad \alpha_{p+1} \neq 0,$$
 (5)

with  $i = \sqrt{-1}$ . We also consider a general Runge–Kutta type time discretization of order q, whose stability function R satisfies<sup>4</sup>

$$R(z) - e^z = \beta_{q+1} z^{q+1} + \mathcal{O}(z^{q+2}), \quad \beta_{q+1} \neq 0.$$
 (6)

We then have as our main result:



<sup>&</sup>lt;sup>1</sup> Rigorously speaking, the wave-number is the inverse of the wavelength, *i.e.*  $\kappa/L$ . Our  $\kappa$  here is a normalized wave-number, but we decided to call it "wave-number" for simplicity.

<sup>&</sup>lt;sup>2</sup> Any other sinusoidal function of wave-number  $\kappa$  and amplitude *A* can be obtained using a shift of the cosine, which would simply lead to a shift in the *x* coordinate and therefore not change the results.

<sup>&</sup>lt;sup>3</sup> For an example, see (27).

<sup>&</sup>lt;sup>4</sup> For an example, see (30).

$$\sigma = \frac{a\Delta t}{\Delta x} \tag{7}$$

are constant. Then the discrete error of the numerical solution,

$$\epsilon := \| \boldsymbol{u}_{num} - \boldsymbol{u}_{exact} \|, \tag{8}$$

evaluated at time  $T := n_P T_P$  with  $n_P$  a real positive constant, and starting with initial condition (2), depends when  $n_x \to \infty$  on the discretization orders (p, q) as follows:

1. If p < q (space discretization has lower order), then

$$\epsilon = n_P A \frac{|\alpha_{p+1}|}{\sqrt{2}} (2\kappa \pi)^{p+1} \frac{1}{n_x^p} + \mathcal{O}\left(\frac{1}{n_x^{p+1}}\right).$$
 (9)

2. If q < p (time discretization has lower order), then

$$\epsilon = n_P A \frac{|\beta_{q+1}|}{\sqrt{2}} (2\kappa\pi)^{q+1} \sigma^q \frac{1}{n_x^q} + \mathcal{O}\left(\frac{1}{n_x^{q+1}}\right).$$
 (10)

3. If p = q (same order for space and time discretization), then

$$\epsilon = n_P A (2\kappa \pi)^{p+1} \frac{\left| -i^{p+1} \alpha_{p+1} + \beta_{p+1} \sigma^p \right|}{\sqrt{2}} \frac{1}{n_x^p} + \mathcal{O}\left(\frac{1}{n_x^{p+1}}\right). \tag{11}$$

We obtain as a consequence the following result on the required number of discretization points needed for a certain accuracy.

**Corollary 1** Using the same hypotheses as in Theorem 1, for a given error tolerance  $\epsilon_{tol}$  small enough, the number of points required for the whole domain discretization can be estimated by

$$n_{x,min} = A^{1/r} n_0 n_P^{1/r} \kappa^{1+1/r}, \tag{12}$$

where  $n_0$  depends on the tolerance  $\epsilon_{tol}$ , the discretization orders (p, q), and  $r = \min(p, q)$ :

1. If p < q (space discretization has lower order), then

$$n_0 \simeq (2\pi)^{1+1/p} \left[ \frac{|\alpha_{p+1}|}{\epsilon_{tol} \sqrt{2}} \right]^{1/p}. \tag{13}$$

2. If q < p (time discretization has lower order), then

$$n_0 \simeq (2\pi)^{1+1/q} \left[ \frac{|\beta_{q+1}|}{\epsilon_{tol} \sqrt{2}} \right]^{1/q} \sigma. \tag{14}$$

3. If p = q (same order for space and time discretization), then

$$n_0 \simeq (2\pi)^{1+1/p} \left[ \frac{\left| -\alpha_{p+1} + \beta_{p+1} (-i)^{p+1} \sigma^p \right|}{\epsilon_{tol} \sqrt{2}} \right]^{1/p}.$$
 (15)

# 2.2 Preliminary definitions

Performing a non-scaled Discrete Fourier Transform (DFT) of the sinusoidal initial value  $u^0$  corresponding to (2) yields<sup>5</sup>

$$\hat{\boldsymbol{u}}^{0} = \left[0, ..., 0, \frac{An_{x}}{2}, ..., 0, ..., \frac{An_{x}}{2}, 0, ..., 0\right].$$
 (16)

The position of the non-zero components (up to machine precision) is linked to the angular frequency of  $u_0$ ,

$$\omega := \frac{2\kappa\pi}{L}.\tag{17}$$

We consider a space grid, so it is convenient to define the reduced angular frequency

$$\theta := \omega \Delta x = \frac{2\kappa \pi}{n_x} \in [-\pi, \pi]. \tag{18}$$

As the DFT is un-scaled, we have the discrete form of the Parseval relation

$$\sum_{j=1}^{n_x} u_j^2 = \frac{1}{n_x} \sum_{\kappa=0}^{n_x} |\hat{\mathbf{u}}_{\kappa}|^2, \tag{19}$$

which allows us to link the discrete norm of the numerical solution with its Fourier components by

$$\|\boldsymbol{u}\| = \sqrt{\frac{1}{n_x} \sum_{j=1}^{n_x} \boldsymbol{u}_j^2} = \sqrt{\frac{1}{n_x^2} \sum_{\kappa=0}^{n_x} |\hat{\boldsymbol{u}}_{\kappa}|^2}.$$
 (20)

For example, if we apply that to our initial condition, we obtain

$$\|\mathbf{u}^0\| = \sqrt{\frac{1}{n_x^2} 2\left(\frac{An_x}{2}\right)^2} = \frac{A}{\sqrt{2}}$$
 (21)

$$= \sqrt{\frac{1}{L} \int_0^L A^2 \cos^2(2\pi \kappa x/L) dx} = \|u_0\|.$$
 (22)

<sup>&</sup>lt;sup>5</sup> The Fourier transform vector is ordered following increasing wavenumbers.

16 Page 4 of 14 M. J. Gander, T. Lunet

Hence, the (normalized) norms are conserved from discrete to continuous space.

*Space discretization.* We consider any finite-difference discretization that will approximate the first space derivative in (1) using

$$\frac{\partial u}{\partial x}(x_j) \approx \frac{1}{\Delta x} \sum_{d=-s}^{s} c_d u(x_{j+d}),$$
 (23)

where the  $c_d$  are the finite difference coefficients, and s is the stencil half-width. Using the method of lines for (1) allows us to build the system of ordinary differential equations

$$\frac{d\mathbf{u}}{dt} = K\mathbf{u},\tag{24}$$

where K represents the discrete space derivative operator. Because of the periodic boundary conditions and the constant velocity, the matrix K is circulant, and thus becomes a diagonal matrix D in Fourier space. Each coefficient of D is associated with one reduced angular frequency of the solution. The two diagonal coefficients concerning our initial condition are then

$$-a\frac{z_{adv}(\pm\theta)}{\Delta x},\tag{25}$$

where  $z_{adv}$  is the Fourier symbol associated with the space discretization, that is

$$z_{adv}(\theta) = \sum_{d=-s/2}^{s/2} c_d e^{ij\theta}.$$
 (26)

For example, the second order centered finite difference scheme has the Fourier symbol

$$z_{adv,C2}(\theta) = \frac{e^{i\theta} - e^{-i\theta}}{2} = i\sin(\theta)$$
$$= i\theta - \frac{(i\theta)^3}{3!} + \mathcal{O}(\theta^5).$$
 (27)

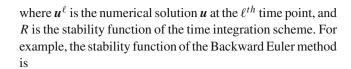
More generally,  $z_{adv}(\theta)$  is an approximation of  $i\theta$  up to order p, that is

$$z_{adv}(\theta) = i\theta + \alpha_{p+1}\theta^{p+1} + \mathcal{O}(\theta^{p+2}), \qquad (28)$$

with  $\alpha_{p+1} \neq 0$ , which is our original hypothesis (5) for Theorem 1 on the finite difference scheme for the space discretization of the advection equation (1).

*Time discretization.* We consider now any Runge–Kutta type time integration method used to solve (24). The numerical solution is then given by

$$\boldsymbol{u}_{num}^{\ell+1} = R(\Delta t K) \boldsymbol{u}^{\ell} = R(\Delta t K)^{\ell+1} \boldsymbol{u}^{0}, \tag{29}$$



$$R_{BE}(z) = \frac{1}{1-z} = 1 + z + z^2 + \mathcal{O}(z^3),$$
 (30)

which implies

$$R_{BE}(z) - e^z = \frac{z^2}{2} + \mathcal{O}(z^3)$$
 (31)

In general, the stability function is an approximation up to some order q of the exponential [27, Lemma 4.4, p.62], that is

$$R(z) - e^z = \beta_{q+1} z^{q+1} + \mathcal{O}(z^{q+2}),$$
 (32)

with  $\beta_{q+1} \neq 0$ , which is our original hypothesis (6) for Theorem 1 on the time integration scheme.

#### 2.3 Proofs of the main results

Estimation of the Fourier error function. For any Runge–Kutta type method, since R is a rational function, its expression is conserved in Fourier space, and we have

$$\hat{\boldsymbol{u}}_{num}^{\ell+1} = R(\Delta t D)^{\ell+1} \hat{\boldsymbol{u}}^0. \tag{33}$$

Using (16) and (25) allows us to get

$$\hat{\boldsymbol{u}}_{num}^{\ell+1} = \left[0, ..., 0, R\left(-\Delta t a \frac{z(-\theta)}{\Delta x}\right)^{\ell+1} \frac{A n_x}{2}, ..., \\ 0, ..., R\left(-\Delta t a \frac{z(\theta)}{\Delta x}\right)^{\ell+1} \frac{A n_x}{2}, 0, ..., 0\right].$$
(34)

Here, we used the short hand notation  $z(\theta) := z_{adv}(\theta)$  and dropped the sub-scripts to be concise in the following computations. Equation (34) can be written using the CFL number as

$$\hat{\boldsymbol{u}}_{num}^{\ell+1} = \left[0, ..., 0, R\left(-\sigma z(-\theta)\right)^{\ell+1} \frac{An_x}{2}, ..., \\ 0, ..., R\left(-\sigma z(\theta)\right)^{\ell+1} \frac{An_x}{2}, 0, ..., 0\right].$$

Let us consider the exact<sup>6</sup> solution of (1), obtained with exact derivation in Fourier space and an exponential time



<sup>&</sup>lt;sup>6</sup> That is, exact up to the given space discretization mesh size. In this case, only a range of wavenumbers for the initial condition can be represented numerically, and higher wave-numbers for the initial condition would require a finer spatial mesh.

integration of (24). The exact solution can be written in Fourier space as

$$\hat{\boldsymbol{u}}_{exact}^{\ell+1} = \left[0, ..., 0, e^{(\sigma i\theta)(\ell+1)} \frac{An_x}{2}, ..., \\ 0, ..., e^{(-\sigma i\theta)(\ell+1)} \frac{An_x}{2}, 0, ..., 0\right].$$

Then, using (20) we can write

$$\epsilon^{2} = \|\mathbf{u}_{num}^{n_{t}} - \mathbf{u}_{exact}^{n_{t}}\|^{2}$$

$$= \frac{1}{n_{x}^{2}} \left| R \left( -\sigma z(-\theta) \right)^{n_{t}} - e^{(\sigma i \theta) n_{t}} \right|^{2} \frac{A^{2} n_{x}^{2}}{4}$$

$$+ \frac{1}{n_{x}^{2}} \left| R \left( -\sigma z(\theta) \right)^{n_{t}} - e^{(-\sigma i \theta) n_{t}} \right|^{2} \frac{A^{2} n_{x}^{2}}{4}$$

$$= \frac{A^{2}}{4} \left| R \left( -\sigma z(-\theta) \right)^{n_{t}} - e^{(\sigma i \theta) n_{t}} \right|^{2}$$

$$+ \frac{A^{2}}{4} \left| R \left( -\sigma z(\theta) \right)^{n_{t}} - e^{(-\sigma i \theta) n_{t}} \right|^{2}.$$
(35)

Hence, in order to compute the discrete error, we have to estimate what we call the Fourier error function

$$\hat{E}(\theta) := R \left( -\sigma z(\theta) \right)^{n_t} - e^{(-\sigma i\theta)(n_t)}. \tag{36}$$

Asymptotic approximation. We first define

$$\zeta := -\sigma i\theta = -\sigma i \frac{2\kappa \pi}{n_x}.$$
(37)

We use the hypothesis that  $\kappa$  and  $n_P$  are constant, and from the definitions we have

$$\zeta n_t = -\frac{a\Delta t}{\Delta x} i \frac{2\kappa \pi}{n_x} n_t = -\frac{a}{L} 2i\kappa \pi T_P n_P$$

$$= -2i\kappa \pi n_P =: c \in \mathbb{C}.$$
(38)

This shows the equivalence

$$\zeta \sim 1/n_t.$$
 (39)

It is worth mentioning that the CFL number  $\sigma$  can be set arbitrarily large or small, depending on the computation<sup>7</sup>. But the hypothesis in Theorem 1 that  $\sigma$  is constant mainly states that it does not vary with  $n_x$ ,  $n_t$  or  $\theta$ . This implies that  $\theta$  and  $\zeta$  scale with one-another, and similarly for  $n_x$  and  $n_t$ , *i.e.* 

$$\theta \sim \zeta \sim 1/n_t \sim 1/n_x. \tag{40}$$

Next, we define

$$\tilde{\zeta} := -\sigma z(\theta) = \zeta - \sigma \alpha_{p+1} \theta^{p+1} + \mathcal{O}\left(\theta^{p+2}\right), 
= \zeta - \alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right),$$
(41)

and then (36) becomes

$$\hat{E}(\theta) = R(\tilde{\zeta})^{\frac{c}{\zeta}} - e^{c} =: \tilde{E}(\zeta). \tag{42}$$

Since we also have the equivalence  $\zeta \sim \tilde{\zeta}$ , we get for any integer r that

$$\mathcal{O}(\zeta^r) = \mathcal{O}(\tilde{\zeta}^r). \tag{43}$$

To study  $\tilde{E}(\zeta)$ , we first perform a Taylor expansion of  $R(\tilde{\zeta})$ : using (6) and (41) we have

$$\begin{split} R(\tilde{\zeta}) &= e^{\tilde{\zeta}} + \beta_{q+1} \tilde{\zeta}^{q+1} + \mathcal{O}\left(\tilde{\zeta}^{q+2}\right) \\ &= \exp\left[\zeta - \alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right)\right] \\ &+ \beta_{q+1} \left(\zeta - \alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right)\right)^{q+1} \\ &+ \mathcal{O}\left(\zeta^{q+2}\right) \\ &= e^{\zeta} \exp\left[-\alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right)\right] \\ &+ \beta_{q+1} \zeta^{q+1} + \mathcal{O}\left(\zeta^{q+2}\right) \\ &= e^{\zeta} \left[1 - \alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right)\right] \\ &+ \beta_{q+1} \zeta^{q+1} + \mathcal{O}\left(\zeta^{q+2}\right) \\ &= e^{\zeta} \left[1 - \alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \mathcal{O}\left(\zeta^{p+2}\right)\right] \\ &+ \beta_{q+1} \zeta^{q+1} + \mathcal{O}\left(\zeta^{q+2}\right)\right]. \end{split}$$

Then, using the notation<sup>8</sup>

$$\psi(\zeta) := -\alpha_{p+1} \frac{i^{p+1}}{\sigma^p} \zeta^{p+1} + \beta_{q+1} \zeta^{q+1},$$
  

$$r := \min(p, q).$$

<sup>&</sup>lt;sup>7</sup> Generally, one looks to maximize  $\sigma$ , in order to get larger time steps and reduce computation time. In practice for advection dominated problems,  $\sigma = \mathcal{O}(1)$  for explicit methods, and for implicit method  $\sigma = \mathcal{O}(10)$  or maybe rarely  $\mathcal{O}(100)$ .

 $<sup>^{8}</sup>$   $\psi$  is actually a local truncation error estimate in Fourier space for the space-time discretization.

16 Page 6 of 14 M. J. Gander, T. Lunet

we obtain

$$\begin{split} R(\tilde{\zeta})^{\frac{1}{\zeta}} &= \exp\left[\frac{1}{\zeta}\ln\left(R(\tilde{\zeta})\right)\right] \\ &= \exp\left[\frac{1}{\zeta}\ln\left[e^{\zeta}\left(1 + \psi(\zeta) + \mathcal{O}\left(\zeta^{r+2}\right)\right)\right]\right] \\ &= \exp\left[\frac{\zeta + \ln\left[1 + \psi(\zeta) + \mathcal{O}\left(\zeta^{r+2}\right)\right]}{\zeta}\right] \\ &= \exp\left[\frac{\zeta + \psi(\zeta) + \mathcal{O}\left(\zeta^{r+2}\right)}{\zeta}\right] \\ &= e \cdot \exp\left[\frac{\psi(\zeta)}{\zeta} + \mathcal{O}\left(\zeta^{r+1}\right)\right] \\ &= e\left(1 - \alpha_{p+1}\frac{i^{p+1}}{\sigma^p}\zeta^p + \beta_{q+1}\zeta^q + \mathcal{O}\left(\zeta^{r+1}\right)\right). \end{split}$$

Taking the expression to the power c, a last Taylor expansion yields

$$R(\tilde{\zeta})^{\frac{c}{\zeta}} = e^{c} \left( 1 - c\alpha_{p+1} \frac{i^{p+1}}{\sigma^{p}} \zeta^{p} + \mathcal{O}(\zeta^{p+1}) + c\beta_{q+1} \zeta^{q} + \mathcal{O}(\zeta^{q+1}) \right),$$

which gives us

$$\tilde{E}(\zeta) = e^{c} \left( -c\alpha_{p+1} \frac{i^{p+1}}{\sigma^{p}} \zeta^{p} + \mathcal{O}(\zeta^{p+1}) + c\beta_{q+1} \zeta^{q} + \mathcal{O}(\zeta^{q+1}) \right).$$

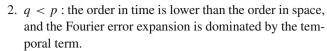
Finally, using the definition in (37) and (38), we get

$$\tilde{E}(\zeta) = e^{c} \left[ -n_{P} \alpha_{p+1} i^{p+1} \frac{(2\kappa \pi)^{p+1}}{n_{x}^{p}} + \mathcal{O}\left(\frac{1}{n_{x}^{p+1}}\right) + n_{P} \beta_{q+1} \frac{(2\kappa \pi)^{q+1}}{n_{x}^{q}} \sigma^{q} + \mathcal{O}\left(\frac{1}{n_{x}^{q+1}}\right) \right]$$

$$=: E(n_{x}).$$
(44)

Origins of discretization errors. We identify two main terms in the Taylor expansion of  $E(n_x)$ : the first one with  $n_x^{-p}$  is the contribution of the space discretization error; the second one with  $n_x^{-q}$  is the contribution of the time discretization error. Hence depending on the values of p and q, three cases need to be considered:

1. *p* < *q* : the order in space is lower than the order in time, and the Fourier error expansion is dominated by the spatial term.



3. q = p: same order in time and space, and the Fourier error expansion is dominated by a mixed term, including error in time and space.

We note from (38) that c is a purely imaginary number, so  $|e^c| = 1$ , which indicates that this term can be ignored when evaluating the absolute value of the Fourier error function in (35). We also recall for what follows that

$$\hat{E}(\theta) = \tilde{E}(\zeta) = E(n_x). \tag{45}$$

Space discretization with lower order. When p < q, the Fourier error becomes

$$E(n_x) = e^c \left[ \alpha_{p+1} i^{p+1} \frac{(2\kappa \pi)^{p+1}}{n_x^p} + \mathcal{O}\left(\frac{1}{n_x^{p+1}}\right) \right], \quad (46)$$

which implies that

$$|\hat{E}(\theta)|^2 = |\hat{E}(-\theta)|^2 + \mathcal{O}\left(\frac{1}{n_r^{p+1}}\right).$$
 (47)

Then, combining (35) with (46), we get

$$\epsilon = n_P A \frac{|\alpha_{p+1}|}{\sqrt{2}} (2\kappa \pi)^{p+1} \frac{1}{n_x^p} + \mathcal{O}\left(\frac{1}{n_x^{p+1}}\right),$$
 (48)

which proves the first part of Theorem 1.

Since  $n_x \to \infty$ , we have  $\epsilon \to 0$  such that we can neglect the term  $\mathcal{O}\left(1/n_x^{p+1}\right)$ , and for a given error tolerance  $\epsilon_{tol}$  small enough, we have

$$n_x \ge A^{1/p} (2\pi)^{1+1/p} \left(\frac{|\alpha_{p+1}|}{\epsilon_{tol}\sqrt{2}}\right)^{1/p} n_P^{1/p} \kappa^{1+1/p},$$
 (49)

which proves the first part of Corollary 1.

Time discretization with lower order. When q < p, the Fourier error becomes

$$E(n_x) = e^c \left[ n_P \beta_{q+1} \frac{(2\kappa \pi)^{q+1}}{n_x^q} \sigma^q + \mathcal{O}\left(\frac{1}{n_x^{q+1}}\right) \right]. \quad (50)$$

Again, the Fourier error function verifies a similar relation as in (47), so we get

$$\epsilon = n_P A \frac{|\beta_{q+1}|}{\sqrt{2}} \sigma^q \frac{1}{n_x^q} + \mathcal{O}\left(\frac{1}{n_x^{q+1}}\right),\tag{51}$$



which proves the second part of Theorem 1, and we obtain similarly as from the first part

$$n_x \ge A^{1/q} (2\pi)^{1+1/q} \left( \frac{|\beta_{q+1}|}{\epsilon_{tol} \sqrt{2}} \right)^{1/q} \sigma n_P^{1/q} \kappa^{1+1/q},$$
 (52)

which proves the second part of Corollary 1.

Space and time discretization with the same order. When p = q, the Fourier error function becomes

$$E(n_x) = e^c \left[ n_P (2\kappa \pi)^{p+1} \left( -i^{p+1} \alpha_{p+1} + \beta_{p+1} \sigma^p \right) + \mathcal{O}\left(\theta^{p+2}\right) \right].$$
(53)

Again, the Fourier error function verifies a similar relation as in (47), so we get

$$\epsilon = n_P A \frac{|-i^{p+1} \alpha_{p+1} + \beta_{p+1} \sigma^p|}{\sqrt{2}} \frac{1}{n_x^p} + \mathcal{O}\left(\frac{1}{n_x^{p+1}}\right), (54)$$

which concludes the proof of Theorem 1, and for Corollary 1 we obtain proceeding as before

$$n_{x} \ge A^{1/p} (2\pi)^{1+1/p} \left( \frac{\left| -\alpha_{p+1} + \beta_{p+1} (-i)^{p+1} \sigma^{p} \right|}{\epsilon_{tol} \sqrt{2}} \right)^{1/q} \times n_{p}^{1/p} \kappa^{1+1/p}.$$

#### 2.4 Discussion of the norm and the error measure

Choice of the norm. We comment now on the definition of the discrete norm in (4). The Parseval theorem (19) used in the proof of Theorem 1 requires to select an  $L^2$ -type norm. Several solutions exist, and we give here the three most natural candidates:

$$\|\mathbf{u}\| := \sqrt{\frac{1}{n_x} \sum_{j=1}^{n_x} u_j(t)^2},$$

$$\|\mathbf{u}\|_{L^2} := \sqrt{\Delta x \sum_{j=1}^{n_x} u_j(t)^2}, = \sqrt{L} \|\mathbf{u}\|,$$

$$\|\mathbf{u}\|_2 := \sqrt{\sum_{j=1}^{n_x} u_j(t)^2}. = \sqrt{n_x} \|\mathbf{u}\|.$$
(55)

 $\|u\|_2$  is the standard Euclidean vector norm, but this norm grows with the number of elements in the vector, and thus can not converge when  $n_x \to \infty$  to a continuous norm of a given continuous solution. Thus, a scaled norm is necessary to build a coherent error estimate, a condition that is satisfied by both  $\|u\|$  and  $\|u\|_{L^2}$ . The latter is the counterpart of the continuous  $L^2$  norm  $\sqrt{\int_L |u(x)|^2 dx}$ , but it varies with the domain length

L. Using it for the error estimate would then induce a dependence on the domain length, which is generally avoided by setting L=1 implicitly. In that case,  $\|u\|_{L^2}$  simply becomes the  $\|u\|$  norm, that has no dependence on L. Note that our discrete norm  $\|u\|$  is similar to the norm used to compute the Mean Square Error in statistics and the root-mean-square velocity in Computational Fluid Dynamics. And finally, our discrete norm  $\|u\|$  is more convenient to apply the discrete form of the Parseval theorem (19). Thus, we chose  $\|u\|$  so our error estimate does not depend on the dimensionality of the problem. Also, note that error estimates based on  $\|u\|_{L^2}$  or  $\|u\|_2$  can be easily retrieved by multiplying the results of Theorem 1 by the appropriate factor.

Choice of the error measure. Once a discrete norm is chosen, defining a measure for the error is not straightforward. First, we defined a norm "in space", but our solution depends also on time. To handle this, there are many approaches in the scientific community, one commonly used consists in looking at the largest  $L^2$  error in space over the simulation time interval. This can be interpreted as an  $L^\infty$  measure in time combined with an  $L^2$  measure in space. For the advection equation, it is already known that the error necessarily grows with time, such that the maximum  $L^2$  error is reached at the end of the simulation time interval (i.e. for t=T). This motivates our choice when defining  $\epsilon$  in (8). We used an absolute error measure, but a relative error measure could also have been used,

$$\epsilon^{rel} := \frac{\|\boldsymbol{u}_{num} - \boldsymbol{u}_{exact}\|}{\|\boldsymbol{u}_{ref}\|},\tag{56}$$

with  $u_{ref}$  being either  $u_0$  or  $u_{exact}$ . Note that choosing one or the other would have no impact on our error measurements if they are computed after an integer number of periods, and for those times,  $u_0 = u_{exact}$ .

Looking at the proof of Theorem 1, it is easy to see that

$$\epsilon^{rel} = \epsilon \text{ when } A = \sqrt{2},$$
(57)

with the same property for the results of Corollary 1. The relative error measure does not depend on the amplitude of the wave, which is an interesting property when the error is dominated by numerical diffusion and contains only one frequency:  $\epsilon^{rel}$  then indicates the percentage of decrease in the solution values<sup>9</sup>. But if one considers several wave-number components of the solution, then the absolute error measure  $\epsilon$  is more appropriate, since it allows a comparison of errors with different weights, related to each wave-number amplitude. Thus, we decided to express our main results in the absolute error measure  $\epsilon$ .

 $<sup>\</sup>overline{{}^{9}}$  If we build an artificial solution  $u_{num} = cu_{exact}$  with  $0 \le c < 1$ , then we have  $\epsilon^{rel} = (1 - c)$ .



16 Page 8 of 14 M. J. Gander, T. Lunet

## 3 Numerical experiments

#### 3.1 Illustration of the error estimates

We set a = L = 1, so all numerical integration parameters are fully determined by  $\sigma$ ,  $n_P$  and  $n_x$ . In particular,  $n_t$  is computed using the fact that

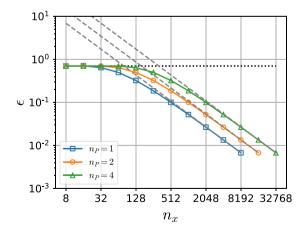
$$\Delta t = \frac{T}{n_t} = \frac{T_P \sigma}{n_x},\tag{58}$$

with  $T_P = L/a = 1$ . Test values for  $\sigma$ ,  $n_P$  and  $n_x$  are set to have  $n_t \in \mathbb{N}$ . We compare the error estimates given when neglecting the big  $\mathcal{O}$  term in the formulas of Theorem 1 to the effective error in numerically computed solutions of (1) in Fig. 1.

On the left, a Backward Euler scheme is used in combination with a 6<sup>th</sup> order centered finite difference discretization. In this case, the time discretization is of lower order, and  $\epsilon$  can be estimated with the dominant term in (10) (dashed gray lines). We observe that the estimates are sharp for large values of  $n_x$ , and this does not depend on the simulation time interval, represented by various values of  $n_P$ . Similar results were also obtained for other configurations with lower order time discretization schemes, and also when the space discretization is of lower order.

In Fig. 1 on the right, the SDIRK3 scheme of [1, Theorem 5] is used in combination with a  $3^{\rm rd}$  order upwind finite difference scheme. Here, since the time and space discretization have the same order, (11) was used to estimate  $\epsilon$  by neglecting the big  $\mathcal{O}$  term, considering a value of the CFL number  $\sigma = 8$ . Again, the theoretical error estimates are sharp for large values of  $n_x$ , with various choices of  $n_P$ .

Generally, we observed with all numerical experiments that estimating  $\epsilon$  by neglecting the big  $\mathcal{O}$  terms in Theorem 1



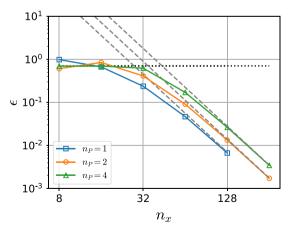
**Fig. 1** Comparison of the theoretical error estimates in Theorem 1 to the effective numerical error. Left: Backward Euler and  $6^{th}$  order centered finite differences,  $\kappa = 2$ , A = 1,  $\sigma = 1$ . Right: SDIRK3 and  $3^{rd}$ 

was a reasonably good approximation for error levels equal or lower than  $\epsilon = 1e^{-2}$ . For low values of  $n_x$ , the error reaches a plateau, with an error around  $\epsilon = A/\sqrt{2}$  (black dotted lines in Fig. 1), which corresponds to the magnitude of the exact solution, hence the numerical approximation has no accuracy whatsoever.

## 3.2 Mesh requirements for a given error tolerance

The results given in Sect. 3.1 have shown the accuracy of the error estimates given in Theorem 1 for general sinusoidal initial conditions as soon as  $n_x$  is large enough to achieve a relatively small value for  $\epsilon$ . This does not apply directly to general initial conditions, but since any solution of (1) can be represented by a linear combination of sinusoidal functions, knowing the highest relevant wavenumber in the initial condition indicates the mesh size requirement to achieve a given error tolerance on  $\epsilon_{tol}$  for the numerical solution to be accurate. If the spatial mesh is fine enough to get an error tolerance  $\epsilon_{tol}$  for the wave-number  $\kappa_{max}$ , then this will automatically apply to all lower wave-numbers.

This motivates Corollary 1, which estimates the minimum number of mesh points  $n_{x,min}$  required to get a given error tolerance for a wave-number  $\kappa$ , after simulating  $n_P$  temporal periods of the advection problem (1). In general,  $n_{x,min}$  depends on the limiting order of the space-time discretization  $r = \min(p,q)$ , as indicated in (12). But three parameters also influence the number of mesh points  $n_{x,min}$  required: the wave-number  $\kappa$ , its associated amplitude A and the number  $n_P$  of simulated temporal periods for (1). The term  $A^{1/r}n_P^{1/r}\kappa^{1+1/r}$  shows that increasing  $n_P$ , the wave-number  $\kappa$  or the amplitude A will naturally increase  $n_{x,min}$  or a given absolute error tolerance, but if higher order methods are used, this increase will be slower. In particular, one can



order upwind finite differences,  $\kappa = 1$ , A = 1,  $\sigma = 8$ . The analytical error estimates are indicated by gray dashed lines, black dotted lines indicate the magnitude of the exact solution



**Table 1** Standard finite difference space disretizations with their  $\alpha_{p+1}$ , and corresponding  $n_0$  when the space discretization has the lower order, for a given relative error tolerance  $\epsilon_{rol}^{rel} = 0.01$ 

Name (ID)	$\alpha_{p+1}$	$n_0(\epsilon_{tol}^{rel} = 0.01)$	
1st order upwind (U1)	$\frac{1}{2}$	1973.9	
2nd order upwind (U2)	$\frac{1}{3}$	90.9	
2nd order centered (C2)	$-\frac{i}{6}$	64.3	
3rd order upwind (U3)	$\frac{1}{12}$	23.5	
4th order upwind (U4)	$\frac{1}{20}$	14.9	
4th order centered (C4)	$-\frac{i}{30}$	13.4	
5th order upwind (U5)	$\frac{1}{60}$	10.0	
6th order centered (C6)	$-\frac{i}{140}$	8.1	
8th order centered (C8)	$-\frac{i}{630}$	6.3	

never be better than "infinite" order of accuracy that would induce no dependency on  $n_P$  and A, and a linear dependency on  $\kappa$ .

Finally, the factor  $n_0$  in (12) depends only on the spacetime discretization schemes. It is particularly important to evaluate the accuracy of numerical schemes, as it indicates the number of mesh points required to get an accuracy of  $\epsilon_{tol}$  for the largest unitary sinusoidal component of the solution ( $\kappa = A = 1$ ), after only one temporal period of (1). This is the focus of the next paragraphs, where values for  $n_0$ are computed for different discretization schemes. Because we focus on one given wavenumber, we consider the relative error tolerance  $\epsilon_{tol}^{rel}$  which represents a percent of error relative to the initial solution (cf. definition of  $\epsilon^{rel}$  in (56), Sect. 2.4, § 2).

Lower order for the space discretization. We consider the case when the space discretization has the lower order, and compute  $n_0(\epsilon_{rol}^{rel}=0.01)$  for several classical finite difference schemes of various order. More details on those schemes are given in Appendix A. The corresponding  $n_0$  values are given in Tab. 1, along with the  $\alpha_{p+1}$  coefficients for each scheme.

We observe a huge  $n_0$  value for the first order method (almost two thousand), and increasing the order rapidly lowers  $n_0$ . But after 4<sup>th</sup> order, the reduction of  $n_0$  when increasing p is not as important any more as before. These results strongly suggest that low order space discretizations should never be considered in practical numerical integrations of smooth solutions of the advection equation (1), given their poor ratio between accuracy and efficiency. Furthermore, higher order difference schemes up to a certain order can easily bring good accuracy at the expense of a small increase in the computational cost of the numerical spatial derivative evaluation, and thus should be favored when solving the advection problem (1) for smooth solutions. Finally, we can also compare upwind and centered schemes with equal

**Table 2** Classical Runge–Kutta time integration schemes with their  $\beta_{q+1}$  coefficient, and corresponding  $n_0$  when the time discretization is of lower order, for a given relative error tolerance  $\epsilon_{tol}^{rel}=0.01$  and CFL number  $\sigma=1$ . The methods are classified into explicit (upper part) and implicit (lower part)

Name	$\beta_{q+1}$	$n_0(\epsilon_{tol}^{rel} = 0.01)$
Forward Euler	$-\frac{1}{2}$	1973.9
Heun	$-\frac{1}{6}$	64.3
RK32	$-\frac{1}{24}$	32.2
RK33	$-\frac{1}{24}$	18.7
RK53	$-\frac{1}{312}$	7.9
RK4	$-\frac{1}{24}$	9.5
Backward Euler	$\frac{1}{2}$	1973.9
Trapezoidal	$\frac{1}{12}$	45.6
SDIRK2	$\simeq 4.04e^{-2}$	31.7
SDIRK2-2	$\simeq -1.37$	184.6
SDIRK3	$\simeq -2.59e^{-2}$	15.9
Gauss-Legendre	$-\frac{1}{720}$	6.1
SDIRK54	$\simeq -8.46e^{-4}$	5.4

order (U2-C2 and U4-C4). For both cases, centered schemes required less mesh points than upwind ones ( $\simeq +41\%$  for U2 compared to C2), but this gap becomes smaller for higher order ( $\simeq +11\%$  for U4 compared to C4). This comparison can be completed considering computational cost (see the coefficients in Appendix A): upwind schemes require more floating point operations than centered ones, which makes them also more costly.

Lower order for the time discretization. We now consider the case when the time discretization has the lower order, and compute  $n_0(\epsilon_{tol}^{rel}=0.01)$  for several classical time integration schemes of various orders. The corresponding  $n_0$  values are given in Table 2, along with the  $\beta_{p+1}$  coefficient for each scheme.

Among the methods considered some are classical (Forward and Backward Euler, Heun, Runge–Kutta of 4th order (RK4), Trapezoidal Rule, Gauss-Legendre) and can be found in many academic textbooks. The others come from the research literature: the Runge–Kutta methods of order 3 with 3 and 5 stages (RK33, RK53) and of order 2 with 3 stages (RK32) from [41], the Singly Diagonally Implicit Runge Kutta methods of order 2 (SDIRK2 and SDIRK2-2) and order 3 (SDIRK3) from [1], and of order 4 with 5 stages (SDIRK54) from [42, Sec. IV.6]. Each  $n_0$  value was computed for a unitary CFL number  $\sigma = 1$ . In practice, given one effective  $\sigma$  value, one should consider the given  $n_0$  multiplied by  $\sigma$ .

Like when we used lower order for the space discretization, first order methods in time require a huge number of mesh points to reach the desired accuracy, while going to higher order quickly reduces  $n_0$  to the order of 10. Also,



methods of the same order do not require necessarily the same number of mesh points for the given error tolerance. For instance, RK33 and RK53 are both 3<sup>rd</sup> order methods, but RK53 has an  $n_0$  less than half the one of RK33; the two additional stages of RK53 compared to RK33 provide not only a better accuracy than RK33, but also a better accuracy than the 4<sup>th</sup> order method RK4. This is due to the chosen error tolerance ( $\epsilon_{tol}^{rel}=0.01$ ), which is not small enough to get the beneficial effects of the 4<sup>th</sup> order. This would appear for lower  $\epsilon_{tol}^{rel}$  values (e.g.  $n_{0,RK4}(\epsilon_{tol}^{rel}=1.0e^{-3})\simeq 16.9$  compared to  $n_{0,RK53}(\epsilon^{rel}=1.0e^{-3})\simeq 17.1$ ).

Same order for space and time discretizations. Finally, we consider the case when p = q. In this case, one should consider (15) of Corollary 1, which depends on the following term governed by the CFL number  $\sigma$ :

$$C_{\alpha,\beta} = \left| -\alpha_{p+1} i^{p+1} + \beta_{p+1} \sigma^p \right|^{1/p}. \tag{59}$$

On the one hand, for low CFL numbers  $\sigma$ , the term multiplying  $\beta_{p+1}$  can become negligible, and we would have

$$C_{\alpha,\beta} \simeq \left| \alpha_{p+1} \right|^{1/p},\tag{60}$$

which would reduce to the case where the space discretization has the lower order. On the other hand, large values of the CFL number  $\sigma$  would make the  $\alpha_{p+1}$  term negligible, and

$$C_{\alpha,\beta} \simeq \left|\beta_{p+1}\right|^{1/p} \sigma,$$
 (61)

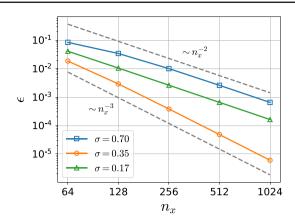
which would reduce to the case where the time discretization has the lower order. This indicates a result that one could naturally expect: for low CFL numbers  $\sigma$ , the error will be dominated by the error of the space discretization, and for large CFL numbers  $\sigma$ , it will be dominated by the error of the time discretization.

### 3.3 Synergistic space-time discretizations

In the case where the time and space discretization have the same order, the Taylor expansion of  $\epsilon$  in (11) has its leading order term multiplied by  $C_{\alpha,\beta}^p$  defined in (59). Under certain sign conditions for  $\alpha_{p+1}$  and  $\beta_{q+1}$ , there exists a CFL number  $\sigma$  which makes the error term of order p vanish, leading to a higher order error term of at most order p+1. If this choice exists, the optimal CFL number is

$$\sigma^{\star} = \left| \frac{\alpha_{p+1}}{\beta_{p+1}} \right|^{1/p}. \tag{62}$$

In that case, we call the space-time discretization "synergistic", since the combination of the two schemes of order p allow us to obtain a scheme of at least global order p + 1.



**Fig. 2** Numerical error of SDIRK2-2 combined with a C2 space discretization, when varying the CFL number  $\sigma$ . Gaussian initial condition, and simulation during one temporal period of (1). Dashed lines represent error reduction of order 2  $(n_x^{-2}, \text{top})$  and of order 3  $(n_x^{-3}, \text{bottom})$ 

Note that in addition, the  $p^{th}$  order error term is canceled for all wave-numbers  $\kappa$ , which leads to the increase of order of accuracy for any initial condition represented on the spatial grid.

A special case of this phenomenon is well known: if we consider Forward Euler (FE,  $\beta_2 = -1/2$ ) combined with first order Upwind finite differences (U1,  $\alpha_2 = 1/2$ ), this combination is synergistic for  $\sigma_{FE/U1}^{\star} = 1$ , and in fact not only the error term of order 1 is canceled, but all the error terms of higher order are also canceled, and the method becomes an exact solver.

There are other combinations of space and time discretization schemes which become synergetic. For instance, we can consider the SDIRK2-2 method from [1], combined with  $2^{\rm nd}$  order centered finite differences (C2) in space. This combination is synergistic and the second order error term is canceled for  $\sigma_{SDIRK2-2/C2}^{\star} \simeq 0.35$ . We illustrate this property by considering an initial condition of the form

$$u_0(x) = e^{-100(x/L - 1/2)^2},$$
 (63)

and simulating up to  $T \simeq T_P$ .

We plot in Fig. 2 the numerical error in the solution at t = T, as a function of  $n_x$ , for three CFL values  $\sigma \in \{\sigma^*/2, \sigma^*, 2\sigma^*\}$ . We observe second order convergence for  $\sigma \in \{0.70, 0.17\}$ , and the error is naturally lower for the smaller CFL number. For  $\sigma = \sigma^*$  however, not only the errors are smaller than for the other two CFL numbers, but we also clearly observe the third order convergence in  $1/n_x$ . Note that we showed this example only for illustration purposes: considering such a low value  $\sigma^*$  for the CFL number defies the purpose of using an implicit method. Furthermore, the SDIRK2-2 method is in practice way less accurate than its counterpart SDIRK2 (see the  $n_0$  values in Table 2), which is probably the reason why it is rarely used in the literature, in



comparison to other SDIRK methods. However, this example is interesting as it illustrates the need to investigate the construction of new space-time discretizations intended to be synergistic, with a  $\sigma^*$  value that would meet some basic efficiency requirements, a subject we intend to further study.

# **4 Application to PinT methods**

We consider now the application of the PARAREAL algorithm to the advection problem (1), for which one needs to define two solvers with different levels of accuracy: a fine solver  $\mathcal{F}$ , which should compute a solution with the desired target level of accuracy, and a coarse solver  $\mathcal{G}$ , which needs to be much faster than  $\mathcal{F}$ , used to achieve time parallel computations. For more information on the PARAREAL algorithm, see [15,20, 28]. We consider N=50 time sub-intervals for PARAREAL (which corresponds to the number of processors that can be used for time parallelization) and  $\Delta T=0.25T_P$  the length of those time sub-intervals, which gives a total computation time  $T=12.5T_P$ . We set the domain length to L=2. At each iteration k of the PARAREAL algorithm, we measure its error by

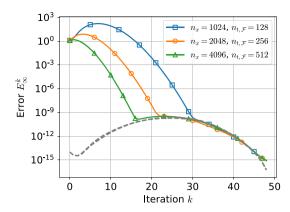
$$E_{\infty}^{k} := \sup_{n>0} \left\| \boldsymbol{u}_{n}^{\mathcal{F}} - \boldsymbol{u}_{n}^{k} \right\|, \tag{64}$$

where  $u_n^{\mathcal{F}}$  is the solution obtained applying  $\mathcal{F}$  sequentially at the end of the  $n^{\text{th}}$  time sub-interval, and  $u_n^k$  is the corresponding PARAREAL solution after k iterations.

We use first BE-C6 as our space-time discretization for the fine solver  $\mathcal{F}$ , with  $\sigma_{\mathcal{F}}=1$ . The coarse solver  $\mathcal{G}$  is using the same BE-C6 scheme, but with a coarsening factor m=32 for the time step, which induces  $\sigma_{\mathcal{G}}=32$ . Note that the choice of implicit time integration simplifies greatly the coarsening for  $\mathcal{G}$ , since the time integration is unconditionally stable for any time-step length. On the other hand, we chose C6 to have a non-dissipative space discretization scheme. We consider the sinusoidal function in (2) with  $A=\kappa=1$ .

The error  $E_{\infty}^k$  is plotted in Fig. 3 for the first iterations of PARAREAL, where we chose different levels of accuracy for  $\mathcal{F}$  by using different values for  $n_x$  and  $n_{t,\mathcal{F}}$  (number of fine time steps on each PARAREAL time sub-interval), keeping the CFL number constant for both solvers ( $\sigma_{\mathcal{F}} = 1$ ,  $\sigma_{\mathcal{G}} = 32$ ).

We see two convergence behaviors on this plot: for early iterations, increasing the accuracy of the fine (and coarse solver) improves the convergence behavior of PARAREAL. For  $n_x=4096$ , PARAREAL even achieves an accuracy of  $E_\infty^k=0.01$  compared to the fine solver after only 7 iterations. This result might at first glance encourage the use of PARAREAL for advection problems, if a high order non dissipative scheme is used for the space discretization, since apparently rapid convergence can be achieved. We see how-



**Fig. 3** Convergence of PARAREAL for the advection problem, using the BE-C6 space-time discretization, and varying the level of accuracy by using different numbers of mesh points  $n_x$  and  $n_{I,\mathcal{F}}$ . Dashed gray lines represent  $E_{\infty}^k$ , using a zero initial condition and a random initial guess with unitary norm, scaled down to  $1e^{-14}$ 

ever also a second convergence behavior in Fig. 3, which the three experiments have in common: after a certain number of iterations, their convergence behavior changes, and follows a new, common curve. One can make this convergence behavior appear right from the start if one starts PARAREAL iterations with a random initial guess, which contains all frequencies 10, in contrast to our experiment setup where we only solve for one frequency in the initial condition and use the coarse propagator to obtain the initial guess. To illustrate this, we use a zero initial condition and start the algorithm with an initial guess containing random values uniformly distributed in [0, 1], scaled such that  $\|\boldsymbol{u}_n^0\| = 1$ . We represent the associated  $E_{\infty}^{k}$  errors with the dashed gray curves in Fig. 3, scaled down to machine precision. We see that the PARAREAL method follows these convergence curves after the initial iterations, due to the (random) roundoff error in the convergence curves in solid lines, amplified by the method. So starting with a random initial guess of order one instead of the coarse solve, none of the methods would have converged before reaching the maximum of 50 parareal iterations in this setting with N = 50 time sub-intervals, as it was predicted in [20, Theorem 5.6, Remark 5.9 and Figure 5.1 on the left] with an analysis at the continuous level for the advection equation. We next use our new results from Sect. 2 to investigate the accuracy of the solutions computed in this example, and to study what happens if one is using the truly needed higher order methods in time.

Limiting discretization scheme. As it was shown in Theorem 1, the error of a space-time discretization scheme is dominated by its lower order component. In our example, the BE-C6 scheme (q = 1, p = 6) is equivalent to any other discretization using a space discretization with p > 1 (e.g. C2,

 $<sup>^{10}</sup>$  For the importance of testing iterative solvers with random initial guess, see [22, end of Section 5.1].



**16** Page 12 of 14 M. J. Gander, T. Lunet

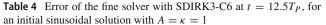
**Table 3** Error of the fine solver with BE-C6 at  $t=12.5T_P$ , for an initial sinusoidal solution with  $A=\kappa=1$ 

$n_{t,\mathcal{F}}$	128	256	512
$n_X$	1024	2048	4096
$\epsilon$	$1.70e^{-1}$	$8.52e^{-2}$	$4.26e^{-2}$

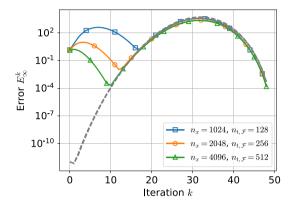
U3, ...). If we use PARAREAL with those lower order space discretizations, we obtain the same convergence curves as in Fig. 3. The important point is that in Fig. 3, the difference between the PARAREAL iterate and the fine solver is shown, and the fine solver accuracy is not at all questioned. The use of the high order "non-dissipative" scheme has however no impact on the global accuracy of the numerical solution, obtained by both applying  $\mathcal F$  sequentially and PARAREAL in parallel. It is very important to study the error of the fine solver before trying to solve the space-time discretized problem in a time parallel fashion.

Error of the fine solver. We thus now investigate the accuracy of the fine solver for the values of  $n_x$  and  $n_{t,\mathcal{F}}$  we used. Since the error estimate increases linearly with the simulation time, it is important to check the error  $\epsilon$  at the final simulation time, that is  $t = 12.5T_P$ . We show in Table 3 the values of  $\epsilon$  for various values of  $n_x$  and  $n_t \neq \infty$ . We see that the error levels are quite high, convergence is first order due to BE in time, and comparing to the error levels of PARAREAL with the fine solver in Fig. 3, we see that the accuracy achieved by PARAREAL compared to the fine solver is not really of interest for the numerical solution of the advection problem in this case. Furthermore, the number of mesh points to achieve this error level with BE-C6 is quite high, which raises the question whether BE should be used at all to solve the advection problem. Our results in Sect. 3.2 clearly indicate the need of higher order time integration schemes for smooth solutions of the advection equation.

The issue of using higher order time-discretization. The simplest coarsening strategy used in PARAREAL consists in taking larger time steps of the fine solver scheme to obtain the coarse solver. In our example, we used a coarsening factor m=32, making the coarse solver  $\mathcal{G}$  thirty two times cheaper than  $\mathcal{F}$ , which leads to the parallel speed-up of PARAREAL. This also induces an error ratio of m between the fine and coarse solver, as we can see from (10) in Theorem 1, and the fact that BE has q = 1. We now keep the same configuration, but use a higher order time discretization for  $\mathcal{F}$  and  $\mathcal{G}$ , that is SDIRK3. This allows us to increase the accuracy of the fine solver, as illustrated in Table 4. But this increase in accuracy brings some drawbacks. In particular, we have an error ratio of  $m^3 = 32^3$  between the fine and coarse solver, while the cost ratio still stays the same. This will make PARAREAL convergence more difficult, since the coarse solver accuracy is more degraded, compared to  $\mathcal{F}$ . We illustrate this by the plot of  $E_{\infty}^{k}$  for the SDIRK3-C6 case in Fig. 4.



$n_{t,\mathcal{F}}$	128	256	512
$n_{\chi}$	1024	2048	4096
$\epsilon$	$3.32e^{-7}$	$4.15e^{-8}$	$5.19e^{-9}$



**Fig. 4** Convergence of PARAREAL for the advection problem, using the SDIRK3-C6 space-time discretization, and varying the level of accuracy, with different number of mesh points  $n_x$ . Dashed gray lines represent  $E_{\infty}^k$ , using a zero initial condition and a random initial guess with unitary norm, scaled down to  $1e^{-12}$ 

We see that not only the initial convergence is degraded compared to the result in Fig. 3, but also the error level of the fine solver can not be reached any more by the PARAREAL iteration before reaching the final 50th iteration, since again convergence for the one target frequency in the experiment changes to the non convergence of parareal with all frequencies present due to roundoff, as indicated by the gray dashed lines computed as for Fig. 3. This shows that increasing the accuracy also makes the true non convergence behavior worse when all frequencies are present.

The difficult choice for time integration. In practice, time integration will often have the lower order, since increasing the space-discretization order does not increase the computational cost too much, and can easily be achieved by increasing the space discretization stencil. In that case however, choosing an implicit time integration scheme will not be advisable for the advection problem, since the benefit of using very large time steps with implicit unconditional stability will be canceled by the  $\sigma^q$  factor in (10) of Theorem 1. This leaves us with explicit time integrators for the fine solver in PARA-REAL, which will then make it difficult to use an implicit scheme for the coarse solver, because it will be hardly possible to have a large cost ratio between an explicit  $\mathcal F$  and an implicit  $\mathcal G$ .

On the other hand, building a coarse solver with explicit time integration cannot be based on an increase of the time step, because of their intrinsic CFL limitation. One can use then coarsening in space, but previous analyses in the lit-



erature [29,33] show that one needs necessarily high order interpolation between the coarse and fine grids, and the coarse solver could then also suffer major efficiency loss in the context of massive space parallelization, not even talking about further high frequency error components introduced due to this process, which the iteration can not eliminate. This makes it very tricky for actual PinT methods based on multilevel techniques, and it seems unavoidable that one has to develop new techniques for the coarser levels when applying PARAREAL like methods to the advection problem, or more generally, hyperbolic problems. In that spirit, interesting ideas have emerged recently in the scientific community for the design of coarse level propagators more suited for PinT integration of hyperbolic problems, see e.g. [5,31,36], and testing these new methods under the stress of a random initial guess with precise accuracy requirements of the computed solution is an important task. The difficulty of good coarse corrections is reminiscent of the tremendous difficulties faced when trying to solve Helmholtz problems using multilevel techniques, see for example [7,8].

#### **5 Conclusion**

We developed general error estimates for the linear advection equation which complement estimates already available in the literature. In particular, our new results allow us to estimate the discrete  $L^2$  norm of the discretization error in the numerical solution for any general space-time discretization based on Runge-Kutta type time integration methods and finite difference space discretizations. We then used our new estimates to derive conditions on the minimum number of points per wavelength needed, which clearly show the interest of using high order discretization methods for smooth solutions of the advection equation.

We also discovered with our new error estimates the existence of synergistic space-time discretizations that allow us to gain one order of accuracy at a given CFL number. We have only shown one example for illustrative purposes, but this discovery opens up the path for developing new spacetime discretization methods constructed specifically to have this property.

We have finally also shown that our new error estimates can be used to analyze the choice of space-time discretizations used in PinT applications, in order to validate them on the advection problem, which is a critical first step toward a more generalized application of PinT methods to larger scale hyperbolic problems. Our preliminary analysis for applying PARAREAL to the advection equation indicates limitations of current classical numerical schemes, and the need of developing new coarse time integrations adapted to hyperbolic problems.

**Acknowledgements** The authors would like to thank Scott MacLachlan and Adrien Laurent for many useful discussions that led to this work. The authors also acknowledge the financial support of the Swiss National Science Foundation project 200020-168999. Finally, the authors thank two anonymous reviewers for the time they took for making helpful comments that allowed us to improve the quality of this manuscript.

Funding Open access funding provided by University of Geneva.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecomm ons.org/licenses/by/4.0/.

# A Space discretization schemes for advection

We give in Table 5 the coefficients of the space discretization schemes investigated in Sect. 3.2. Most of them belong to common literature on finite-difference schemes. The U2 and U4 scheme are not very common, but used in [5].

Table 5 Standard finite difference space disretizations coefficients for advection. Half of the coefficients for C8 are not indicated to reduce table width, as they can be easily deduced from the first coefficients (see others centered schemes)

ID	Finite differences formula
U1	$\frac{u_i - u_{i-1}}{\Delta x}$
U2	$\frac{3u_i - 4u_{i-1} + u_{i-2}}{2\Delta x}$
C2	$\frac{u_{i+1}-u_{i-1}}{2\Delta x}$
U3	$\frac{2u_{i+1} + 3u_i - 6u_{i-1} + u_{i-2}}{6\Delta x}$
U4	$\frac{3u_{i+1} + 10u_i - 18u_{i-1} + 6u_{i-2} - u_{i-1}}{12\Delta x}$
C4	$\frac{-u_{i+2} + 8u_{i+1} - 8u_{i-1} + u_{i-2}}{12\Delta x}$
U5	$\frac{-3u_{i+2} + 30u_{i+1} + 20u_i - 60u_{i-1} + 15u_{i-2} - 2u_{i-3}}{60\Delta x}$
C6	$\frac{u_{i+3} - 9u_{i+2} + 45u_{i+1} - 45u_{i-1} + 9u_{i-2} - u_{i-3}}{60\Delta x}$
C8	$\frac{-3u_{i+4} + 32u_{i+3} - 168u_{i+2} + 672u_{i+1} - \dots}{840\Delta x}$

# References

- 1. Alexander, R.: Diagonally implicit Runge-Kutta methods for stiff ODEs. SIAM J. Numer. Anal. 14(6), 1006-1021 (1977)
- 2. Bal G.: On the convergence and the stability of the Parareal algorithm to solve partial differential equations. In: Kornhuber, R., et al.



16 Page 14 of 14 M. J. Gander, T. Lunet

- (eds.) Domain Decomposition Methods in Science and Engineering, Lecture Notes in Computational Science and Engineering, vol. 40, pp. 426–432. Springer (2005)
- Chen, F., Hesthaven, J. S., Maday, Y., Nielsen, A. S.: An adjoint approach for stabilizing the Parareal method. Tech. rep., EPFL-ARTICLE-211097 (2015)
- Chen, F., Hesthaven, J.S., Zhu, X.: On the use of reduced basis methods to accelerate and stabilize the parareal method. In: Reduced Order Methods for Modeling and Computational Reduction, pp. 187–214. Springer (2014)
- De Sterck, H., Falgout, R.D., Friedhoff, S., Krzysik, O.A., MacLachlan, S.P.: Optimizing MGRIT and Parareal coarse-grid operators for linear advection. arXiv preprint arXiv:1910.03726 (2019)
- Eghbal, A., Gerber, A.G., Aubanel, E.: Acceleration of unsteady hydrodynamic simulations using the parareal algorithm. J. Comput. Sci. 19, 57–76 (2017)
- Ernst, O.G., Gander, M.J.: Why it is difficult to solve Helmholtz problems with classical iterative methods. In: Numerical analysis of multiscale problems, pp. 325–363. Springer (2012)
- 8. Ernst, O.G., Gander, M.J.: Multigrid methods for Helmholtz problems: a convergent scheme in 1D using standard components. Direct Inverse Prob. Wave Propag. Appl. 14, 135–186 (2013)
- Falgout, R.D., Friedhoff, S., Kolev, T.V., MacLachlan, S.P., Schroder, J.B.: Parallel time integration with multigrid. SIAM J. Sci. Comput. 36(6), C635–C661 (2014)
- Friedhoff, S., Falgout, R., Kolev, T., MacLachlan, S., Schroder, J.:
   A multigrid-in-time algorithm for solving evolution equations in parallel. In: Sixteenth Copper Mountain Conference on Multigrid Methods, Copper Mountain, CO, United States (2013)
- Gander, M.J.: Analysis of the parareal algorithm applied to hyperbolic problems using characteristics. Bol. Soc. Esp. Mat. Apl. 42, 21–35 (2008)
- Gander, M.J.: 50 years of time parallel time integration. In: Carraro, T., Geiger, S.K., Rannacher R. (eds.) Multiple Shooting and Time Domain Decomposition Methods, pp. 69–114. Springer (2015)
- Gander, M.J., Güttel, S.: ParaExp: a parallel integrator for linear initial-value problems. SIAM J. Sci. Comput. 35(2), C123–C142 (2013)
- Gander, M.J., Güttel, S., Petcu, M.: A nonlinear ParaExp algorithm.
   In: International Conference on Domain Decomposition Methods, pp. 261–270. Springer (2017)
- Gander, M.J., Hairer, E.: Nonlinear convergence analysis for the Parareal algorithm. In: Widlund, O.B., Keyes, D.E. (eds.) Domain Decomposition Methods in Science and Engineering, Lecture Notes in Computational Science and Engineering, vol. 60, pp. 45– 56. Springer (2008)
- Gander, M.J., Halpern, L.: Absorbing boundary conditions for the wave equation and parallel computing. Math. Comput. 74(249), 153–176 (2005)
- Gander, M.J., Halpern, L., Nataf, F.: Optimal Schwarz waveform relaxation for the one dimensional wave equation. SIAM J. Numer. Anal. 41(5), 1643–1681 (2003)
- Gander, M.J., Halpern, L., Rannou, J., Ryan, J.: A direct time parallel solver by diagonalization for the wave equation. SIAM J. Sci. Comput. 41(1), A220–A245 (2019)
- Gander, M.J., Kwok, F., Zhang, H.: Multigrid interpretations of the parareal algorithm leading to an overlapping variant and MGRIT. Comput. Vis. Sci. 19(3–4), 59–74 (2018)
- Gander, M.J., Vandewalle, S.: Analysis of the parareal time-parallel time-integration method. SIAM J. Sci. Comput. 29(2), 556–578 (2007)
- Gander, M.J., Wu, S.L.: Convergence analysis of a periodic-like waveform relaxation method for initial-value problems via the diagonalization technique. Numer. Math. 143(2), 489–527 (2019)

- Gander, M.J., et al.: Schwarz methods over the course of time. Electron. Trans. Numer. Anal 31(5), 228–255 (2008)
- Gustafsson, B.: High Order Difference Methods for Time Dependent PDE, vol. 38. Springer, Berlin (2007)
- 24. Gustafsson, B., Kreiss, H.O., Oliger, J.: Time Dependent Problems and Difference Methods, vol. 24. Wiley, New Jersey (1995)
- Hessenthaler, A., Nordsletten, D., Röhrle, O., Schroder, J.B., Falgout, R.D.: Convergence of the multigrid reduction in time algorithm for the linear elasticity equations. Numer. Linear Algebra Appl. 25(3), e2155 (2018)
- Howse, A.J., Sterck, H.D., Falgout, R.D., MacLachlan, S., Schroder, J.: Parallel-in-time multigrid with adaptive spatial coarsening for the linear advection and inviscid burgers equations. SIAM J. Sci. Comput. 41(1), A538–A565 (2019)
- Iserles, A.: A First Course in the Numerical Analysis of Differential Equations., vol. 44. Cambridge university press, Cambridge (2009)
- Lions, J.L., Maday, Y., Turinici, G.: A "Parareal" in time discretization of PDE's. C. R. Math. Acad. Sci. Paris 332(7), 661–668 (2001)
- 29. Lunet, T., Bodart, J., Gratton, S., Vasseur, X.: Time-parallel simulation of the decay of homogeneous turbulence using parareal with spatial coarsening. Comput. Vis. Sci. 19(1–2), 31–44 (2018)
- Neumüller, M.: Space-Time Methods: Fast Solvers and Applications. Monographic Series TU Graz: Computation in Engineering and Science (2013)
- Nguyen, H., Tsai, R.: A stable parareal-like method for the second order wave equation. J. Comput. Phys. 405, 109156 (2020)
- Nielsen, A.S., Brunner, G., Hesthaven, J.S.: Communication-aware adaptive parareal with application to a nonlinear hyperbolic system of partial differential equations. J. Comput. Phys. 371, 483–505 (2018)
- Ruprecht, D.: Convergence of parareal with spatial coarsening. PAMM 14(1), 1031–1034 (2014)
- Ruprecht, D.: Wave propagation characteristics of parareal. Comput. Vis. Sci. 19(1–2), 1–17 (2018)
- Ruprecht, D., Krause, R.: Explicit parallel-in-time integration of a linear acoustic-advection system. Comput. Fluids 59, 72–83 (2012)
- Schmitt, A., Schreiber, M., Peixoto, P., Schäfer, M.: A numerical study of a semi-Lagrangian parareal method applied to the viscous burgers equation. Comput. Vis. Sci. 19(1–2), 45–57 (2018)
- Speck, R., Ruprecht, D., Emmett, M., Bolten, M., Krause, R.: A space-time parallel solver for the three-dimensional heat equation.
   In: Parallel Computing: Accelerating Computational Science and Engineering (CSE). Advances in Parallel Computing, pp. 263–272 (2014)
- Swartz, B., Wendroff, B.: The relative efficiency of finite difference and finite element methods. i: hyperbolic problems and splines. SIAM J. Numer. Anal. 11(5), 979–993 (1974)
- 39. Ta'asan, S., Zhang, H.: Fourier-Laplace analysis of the multigrid waveform relaxation method for hyperbolic equations. BIT Numer. Math. **36**(4), 831–841 (1996)
- 40. Taasan, S., Zhang, H.: On the multigrid waveform relaxation method. SIAM J. Sci. Comput. 16(5), 1092–1104 (1995)
- 41. Wang, R., Spiteri, R.J.: Linear instability of the fifth-order WENO method. SIAM J. Numer. Anal. 45(5), 1871–1901 (2007)
- 42. Wanner, G., Hairer, E.: Solving Ordinary Differential Equations II. Springer, Berlin Heidelberg (1996)
- 43. Wu, S.L.: Toward parallel coarse grid correction for the parareal algorithm. SIAM J. Sci. Comput. **40**(3), A1446–A1472 (2018)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

