

**Stability analysis of the
numerical Method of characteristics
applied to a class of
energy-preserving hyperbolic systems.
Part I: Periodic boundary conditions**

T.I. Lakoba*, Z. Deng

Department of Mathematics and Statistics, 16 Colchester Ave.,
University of Vermont, Burlington, VT 05401, USA

February 23, 2019

Abstract

We study numerical (in)stability of the Method of characteristics (MoC) applied to a system of non-dissipative hyperbolic partial differential equations (PDEs) with periodic boundary conditions. We consider three different solvers along the characteristics: simple Euler (SE), modified Euler (ME), and Leap-frog (LF). The two former solvers are well known to exhibit a mild, but unconditional, numerical instability for non-dissipative ordinary differential equations (ODEs). They are found to have a similar (or stronger, for the MoC-ME) instability when applied to non-dissipative PDEs. On the other hand, the LF solver is known to be stable when applied to non-dissipative ODEs. However, when applied to non-dissipative PDEs within the MoC framework, it was found to have by far the strongest instability among all three solvers. We also comment on the use of the fourth-order Runge–Kutta solver within the MoC framework.

Keywords: Method of characteristics, Coupled-wave equations, Numerical instability.

*tlakoba@uvm.edu, 1 (802) 656-2610

1 Introduction

In this series of two papers we address numerical stability of the Method of characteristics (MoC) applied to a class of non-dissipative hyperbolic partial differential equations (PDEs) in 1 time + 1 space dimension. For reference purposes, we will now present the idea of the MoC using the following system as an example:

$$\begin{aligned}u_{1t} + u_{1x} &= f_1(u_1, u_2) \\ u_{2t} - u_{2x} &= f_2(u_1, u_2),\end{aligned}\tag{1}$$

where $f_{1,2}$ are some differentiable functions. In the MoC, each of the equations in (1) is transformed to an ordinary differential equation (ODE) in its respective variable, $\eta_{\pm} = x \mp t$, and then is solved by an ODE numerical solver.

The MoC is widely used and is described in most textbooks on numerical solution of PDEs. As for any numerical method, in order to converge to a true solution, the method should be both consistent and stable. It is, therefore, quite surprising that stability of the MoC has not been investigated in as much detail as that of most other numerical methods. For example, in [1], it was considered by the von Neumann analysis (which implies periodic boundary conditions (BC)) for a simple scalar model

$$u_t + u_x = 0,\tag{2a}$$

with the ODE solver being the simple Euler method. Clearly, that analysis could be straightforwardly generalized for a vector model

$$\mathbf{u}_t + \mathbf{A}\mathbf{u}_x = \mathbf{0},\tag{2b}$$

where \mathbf{A} is any constant diagonalizable matrix with real eigenvalues; an analysis for a similar model in three spatial dimensions was presented in [2]. However, we have been unable to find a systematic stability analysis — even for periodic BC — for a system of the form

$$\mathbf{u}_t + \mathbf{A}\mathbf{u}_x = \mathbf{B}\mathbf{u},\tag{3}$$

where \mathbf{A} , \mathbf{B} are constant matrices. A partial exception is a conference paper [3], where a model with a somewhat more complicated right-hand side (rhs), describing a specific engineering application, was studied by the von Neumann analysis. The focus of that paper was on the examination of the effect of various parameters of that particular model rather than on the impact on stability of the MoC by the ODE solver used. (In fact, the ODE solver used in [3] was the simple Euler method.)

In this paper we will examine the effect of the ODE solver on the stability of the MoC applied to the model problem (3) with *periodic* BC. In most applications, other types of BC are more appropriate for hyperbolic systems (3); e.g., fully or partially nonreflecting BC should be assigned in problems of electromagnetic wave propagation in distributed-feedback lasers ([4], Sec. 3.8) or fiber Bragg gratings [5]. However, studying stability of the MoC with periodic BC is a reasonable first step in the direction of more general stability studies, as it allows one to use the standard

analytical tool, the von Neumann analysis. A study of stability of the MoC with nonreflecting BC is carried out — by a different analytical technique — in a companion paper [6].

We will further limit the scope of the problem as follows. First, and most importantly, we will consider only non-dissipative and stable systems. This implies that matrix \mathbf{B} possesses only imaginary eigenvalues. Second, we will consider only three ODE solvers: simple Euler (SE), modified Euler (ME), and the Leap-frog (LF). Let us note that for a non-dissipative stable ODE, the numerical solutions obtained by the SE and ME are known to be mildly, but unconditionally, unstable, with the growth rates of the numerical instability being $O(\Delta t)$ and $O(\Delta t^3)$, respectively; see, e.g., Secs. 3 and 4 below. On the contrary, the LF solver is known to quasi-preserve (i.e., not shift systematically) at least some of the conserved quantities of non-dissipative ODEs.

Our third assumption is that the physical system described by (3) comprises only two groups of quantities (waves) that propagate with two constant and different velocities, c_{\pm} . One well-known and general example of such a system is the (linear) Klein–Gordon equation:

$$u_{tt} - c^2 u_{xx} = g_1(x, t)u + g_2(x, t), \quad (4)$$

where $c_{\pm} \equiv \pm c$ and $g_{1,2}$ are arbitrary continuous functions. The fact that (4) can be written in the form (3) is demonstrated, e.g., in [7]. A representative list of specific physical systems, leading to (3), is given in Section 2. With our third assumption, the constant matrix \mathbf{A} in (3) can be diagonalized into the form

$$\mathbf{A} = c_1 \mathbf{I}_N + c_2 \mathbf{\Sigma}, \quad \mathbf{\Sigma} = \text{diag}(\mathbf{I}_{N_1}, -\mathbf{I}_{N_2}), \quad (5a)$$

where $c_{\pm} = c_1 \pm c_2$, \mathbf{I}_N is the $N \times N$ identity matrix, and N is the dimension of vector \mathbf{u} . Without loss of generality (i.e., by the change of variables: $t_{\text{new}} = t_{\text{old}}$, $x_{\text{new}} = (x_{\text{old}} - c_1 t_{\text{old}})/c_2$) one can set

$$c_1 = 0, \quad c_2 = 1. \quad (5b)$$

We will also let $N_1 = N_2 = N/2$. Particular conclusions of the forthcoming analysis may change if the diagonalization of matrix \mathbf{A} is different from (5a). However, the methodology of the analysis will not be affected.

Our analysis reveals two facts which, to our knowledge, have not been previously pointed out. First, in contrast to the situation with most numerical methods, the most numerically unstable Fourier harmonics may occur not at the edges of the spectrum but in the “middle” (see the footnote in Section 4) of it. In fact, the (in)stability of the highest Fourier harmonics is the same as that of the lowest ones. This fact easily follows from the foregoing analysis and holds for the MoC employing any ODE solver considered in this work, as long as matrix \mathbf{A} has the form (5).

Second, and quite unexpectedly, the LF, which outperforms the SE and ME when applied to non-dissipative ODEs, performs much worse than those methods when applied to non-dissipative stable PDEs within the MoC framework.

Let us now comment on the main advantage and disadvantage of the MoC over other methods that are used to solve the Klein–Gordon equation (4) and the more general system (3). There

exists a vast literature on this subject, and we will mention only three groups of such methods. First, there are finite-difference methods, where the x - and t -partial derivatives are discretized on the (x, t) -grid with a desired accuracy. Many of those methods, e.g., Lax–Wendroff, MacCormack, etc., are covered in textbooks (see, e.g., [8]). More accurate variants of them often use higher-order Runge–Kutta-type discretization in time and higher-order discretization in space (see, e.g., [9, 10, 11] and a review [12]). A variation (for the purpose of this grouping) of finite-difference methods are collocation-type methods, where spatial derivatives are not approximated explicitly but rather evaluated via inner product with certain basis functions, e.g., B-splines [13] or radial basis functions [14]. The second group comprises split-step methods, which handle the \mathbf{A} - and \mathbf{B} -terms in (3) in consecutive substeps (see, e.g., [15, 16, 17]). The third group is that of exponential-time differencing and integrating factor methods, which treat the \mathbf{A} -part of (3) “exactly” (i.e., with accuracy inherent to that of discretization of \mathbf{u}_x) while approximating the effect of the \mathbf{B} -term by finite differencing in time (see, e.g., [18, 19, 17]). On the other hand, the MoC has been used for equations of the type (3)–(5) primarily in the context of electromagnetic wave propagation in: optical fiber Bragg gratings [20], distributed-feedback semiconductor lasers [21, 22, 23] and Raman lasers [24], periodic two-level resonant media [25, 26], birefringent optical fibers [27], and stimulated Brillouin scattering [28, 29].

The main advantage of the MoC over methods of the first group is that it preserves the linear dispersion relation for high wavenumbers:

$$\omega = \pm ck, \quad |k| \gg 1, \quad (6)$$

of the plane wave, i.e., $\exp[i(kx - \omega t)]$, solution of (3) (where c is an eigenvalue of \mathbf{A}) or (4). We will demonstrate this statement in Section 3 for the MoC-SE; it can be shown similarly for other MoC schemes. In contrast, of all methods of the first group, we know of only one, Ref. [30], that preserves the dispersion relation (6) for high wavenumbers¹; all others significantly distort it. A distortion of the dispersion relation that makes ω a non-monotone function of k can lead to nonphysical reflection of high-wavenumber modes from the boundaries into the computational domain (see, e.g., [32]) and subsequent contamination of the solution; e.g., in the relativistic field theory this is referred to as “fermion doubling”. A substantial effort has been invested into construction of numerical boundary schemes that suppress those reflections; see, e.g., a review [33] and a recent paper [34]. Even if the dispersion relation $\omega(k)$ of a finite-difference method is monotone, its deviation from the analytical expression (6) leads to numerical dispersion (see, e.g., [8]) and can sometimes lead to physically erroneous conclusions about the solution [7]. A series of studies [9, 10, 35] was devoted to construction of so-called low-dispersion and low-dissipation finite-difference methods in computational acoustics to minimize such distortions. On the other hand, the MoC schemes provide those desirable properties inherently.

¹However, the method of [30] uses a leap-frog time discretization and therefore cannot be used for long-term computation of dissipative or non-dissipative but nonlinear hyperbolic systems, unless a certain low-wavenumber filtering is used to remove parasitic modes (see, e.g., [31]).

The split-step methods of group two can preserve the dispersion relation (6) if they use collocation with nonlocal basis functions for spatial discretization, e.g., the discrete Fourier transform. However, the latter imposes periodic BC, which often do not correspond to the actual BC of the problem at hand. To overcome this problem and be able to impose any type of BC, one can use collocation with Chebyshev polynomials. However, this would severely restrict the stability threshold (from $\Delta t_{\text{thresh}} = O(1/M_{\text{colloc}})$ to $\Delta t_{\text{thresh}} = O(1/M_{\text{colloc}}^2)$, where M_{colloc} is the number of collocation nodes) of the numerical scheme compared to the case where local finite differences or the discrete Fourier transform are used [36]. Incidentally, if the spatial derivative substep in the split-step method is computed by the MoC, this both preserves the dispersion relation and allows one to use a BC prescribed by the physics of the problem; see, e.g., [37, 38]. Finally, the situation with methods of group three is similar to that with methods of group two: the dispersion relation is preserved only by a nonlocal spatial discretization, which either imposes periodic BC or severely reduces the stability threshold of the method. Also, unlike the situation with split-step methods, here it is not known how spatial discretization could be implemented with a MoC scheme.

The main disadvantage of the MoC compared to all the methods of the above three groups, as applied to Eqs. (3) or (4), is that only first- and second-order accurate explicit MoC schemes are known. We briefly comment on the absence of higher-order explicit MoC schemes in Section 6. There exist implicit fourth-order MoC schemes based on Runge–Kutta [20] and multistep [23] solvers. As any implicit solvers, they are slower and more difficult to implement than explicit ones. Moreover, while no numerical instability for them has been reported, neither a stability analysis nor a systematic numerical study of stability was done for them, to the best of our knowledge. We present the first such an analysis, with the simplifying assumption of periodic BC and for three low-order explicit MoC schemes, in this paper.

The main part of this work is organized as follows. In Section 2 we specify the physical model to be considered below. Our analysis will be developed for the linearized version of that model. In Section 3 we present the von Neumann analysis and its verification by direct numerical simulations of the PDE for the case when the MoC employs the SE solver. In Sections 4 and 5 we repeat steps of Section 3 for the cases where the ODE solvers are the ME and LF, respectively.

It should be noted that in all three cases, the von Neumann analysis leads one to 4×4 generalized eigenvalue problems which involve both matrices \mathbf{B} and \mathbf{A} (or, as per (5a), $\mathbf{\Sigma}$). The largest-magnitude eigenvalue, which determines the stability of the numerical method, of those problems cannot be analytically found or even related to the *individual* eigenvalues of the “physical” matrices \mathbf{B} and \mathbf{A} . However, for given \mathbf{B} and \mathbf{A} , that eigenvalue can be easily and quickly found numerically by standard built-in commands in software like Matlab, and thus we use this semi-analytical approach.

In Section 6 we comment on the case when the ODE solver is the fourth-order classical Runge–Kutta (cRK) method. We will not, however, analyze it in any detail because, as we will demonstrate, there is an uncertainty about using the cRK solver within the MoC framework. A detailed investi-

gation of this issue is outside the scope of this paper. In Section 7 we present our conclusions and demonstrate the validity of our analysis for a broader class of systems than the particular system considered in Sections 2–5. First, in Appendix A we apply our analysis of the MoC-LF scheme to a class of models of which the model considered in the main text is a special case. These models still result in linearized equations with constant coefficients. Then, in Appendix B, we present a different model with a *spatially localized* solution and demonstrate by direct numerical simulations that essential conclusions of our von Neumann analysis remain valid for (in)stability of the MoC applied to that non-constant-coefficient system.

2 Physical model

While our study will focus on a linear problem of a rather general form (see Eq. (12a) below), we will begin by stating a specific nonlinear problem which had originally motivated our study and whose linearization leads to (12a). We consider the system

$$\underline{\mathbf{S}}_t^\pm \pm \underline{\mathbf{S}}_x^\pm = \underline{\mathbf{S}}^\pm \times \hat{\mathbf{J}} \underline{\mathbf{S}}^\mp, \quad (7)$$

where $\underline{\mathbf{S}}^\pm \equiv [S_1^\pm, S_2^\pm, S_3^\pm]^T$, $\hat{\mathbf{J}} = \text{diag}(1, -1, -2)$, and superscript ‘T’ denotes the transposition. This system is a representative of a class of models that arise in studying propagation of light in birefringent optical fibers with Kerr nonlinearity [39]–[42]; see Appendix A for more detail. The nonlinear system (7) has a rather special form and arises in a specific application. However, its *linearized* form (see Eq. (12a) below), for which we will analyze the stability of the MoC, is quite general. A wide and diverse range of physical problems lead to the same equation, as we will explain after Eqs. (13).

In the component form, system (7) is:

$$(\partial_t + \partial_x)S_1^+ = S_3^+ S_2^- - 2S_2^+ S_3^-, \quad (8a)$$

$$(\partial_t + \partial_x)S_2^+ = 2S_1^+ S_3^- + S_3^+ S_1^-, \quad (8b)$$

$$(\partial_t + \partial_x)S_3^+ = -(S_1^+ S_2^- + S_2^+ S_1^-), \quad (8c)$$

$$(\partial_t - \partial_x)S_1^- = S_3^- S_2^+ - 2S_2^- S_3^+, \quad (8d)$$

$$(\partial_t - \partial_x)S_2^- = 2S_1^- S_3^+ + S_3^- S_1^+, \quad (8e)$$

$$(\partial_t - \partial_x)S_3^- = (\partial_t + \partial_x)S_3^+. \quad (8f)$$

It has four families of soliton/kink solutions, a special case of one of which is (see, e.g., [41]):

$$S_1^\pm = \pm \frac{1}{\sqrt{3}} \text{sech}(\sqrt{2}x), \quad S_2^\pm = \mp \tanh(\sqrt{2}x), \quad S_3^\pm = \sqrt{2} S_1^\pm. \quad (9)$$

(The other three families differ from (9) by combinations of signs of the solution’s components.) However, we will focus on the analysis of the (in)stability of the MoC applied to a simpler, *constant* solution:

$$S_{1,3}^\pm = 0, \quad S_2^\pm = \pm 1. \quad (10)$$

This solution corresponds to the asymptotic value of (9) as $x \rightarrow -\infty$. It can be shown (see, e.g., [43]) that this solution is stable in the sense to be defined below. However, we have found that when simulated by certain “flavors” of the MoC, it can be numerically unstable. If that asymptotic, constant solution is numerically unstable, then so will be any “physically interesting” non-constant solutions, like (9), possessing the asymptotics (10). Thus, numerical stability of the constant solution (10) is necessary for the successful performance of the MoC on the soliton/kink solution (9) and related ones. Moreover, considering the numerical stability of the MoC simulating the simpler solution (10) will allow us to focus on the behavior of the numerical method without being distracted by the complexity of the physical model. The methodology of our analysis, as well as at least some of our general conclusions, will hold for the MoC applied to other non-dissipative hyperbolic PDEs; we will discuss this in Section 7 and provide details in the Appendices.

For future reference, we linearize Eqs. (8) on the background of solution (10):

$$S_j^\pm = S_{j0}^\pm + s_j^\pm, \quad j = 1, 2, 3; \quad (11)$$

here S_{j0}^\pm are the components of the exact solution (10) and s_j^\pm are small perturbations. The linearized system (8) reduces to:

$$\mathbf{s}_t + \mathbf{\Sigma} \mathbf{s}_x = \mathbf{P} \mathbf{s}, \quad (12a)$$

$$(s_2^\pm)_t \pm (s_2^\pm)_x = 0, \quad (12b)$$

where $\mathbf{s} = [s_1^+, s_3^+, s_1^-, s_3^-]^T$, the 4×4 matrices in (12a) are:

$$\mathbf{\Sigma} = \text{diag}(I, -I), \quad \mathbf{P} = \begin{pmatrix} -A & B \\ -B & A \end{pmatrix}, \quad (13a)$$

and the 2×2 matrices in (13a) are:

$$I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \quad B = -\begin{pmatrix} 0 & 2 \\ 1 & 0 \end{pmatrix}. \quad (13b)$$

Non-dissipative linearized equations of the form (12a), with the same or similar $\mathbf{\Sigma}$ but a more general form of \mathbf{P} , arise in a wide range of physical applications involving the interaction of waves propagating with distinctly different velocities. In addition to the optics-related applications mentioned in Section 1 [20]–[29], other examples include: acousto-optical interactions, known as Stimulated Brillouin Scattering ([44], Sec. 8.3; [45, 46]); Stimulated Raman Scattering ([44], Sec. 9.4; [47]) and, more generally, parametric interaction among two or more electromagnetic waves [48] in semi-classical optics [49, 50], quantum optics [51, 52], photonics [53], and plasma physics applications [54, 55]; interaction of crossed beams in photorefractive materials ([44], Sec. 10.6); interaction of counter- and co-propagating light waves in a medium with a periodic refractive index ([4], Sec. 3.4; [5, 56, 57]) and in laser resonator cavities [58, 59]; and the relativistic field theory [60, 61] (also see [62] for a general form of coupled-wave-like field models with a cubic nonlinearity and [37] for more recent references on these models’ solutions). Interestingly, models similar to those in the

relativistic field theory also arise in nonlinear fiber optics [63, 64]. In Section 7 we will demonstrate that our results for system (12), (13) are also applicable to the one-dimensional Gross–Neveu model [61] of the relativistic field theory.

In Eqs. (12) and (13) and in what follows we will adopt the following notations. Boldfaced quantities with an underline, $\underline{\mathbf{S}}^\pm$, and with a hat, $\hat{\mathbf{J}}$, will continue to denote 3×1 vectors and 3×3 matrices, respectively, as in (7). Boldfaced quantities *without* an underline or a hat will denote 4×4 matrices or 4×1 vectors, as in (12a); the ambiguity of the same notations for matrices and vectors here will not cause any confusion. Finally, underlined letters in regular (not boldfaced) font will denote 2×1 vectors; e.g.:

$$\underline{\mathbf{s}}^\pm \equiv [s_1^\pm, s_3^\pm]^T.$$

Clearly then, $\mathbf{s} \equiv [(\underline{\mathbf{s}}^+)^T, (\underline{\mathbf{s}}^-)^T]^T$.

Seeking the solution of (12a) to be proportional to $e^{ikx-i\omega t}$, one can show that $\omega \in \mathbb{R}$ for all $k \in \mathbb{R}$. This means that solution (10) is stable on the infinite line, as mentioned above. In particular, $\omega(k=0) \in \mathbb{R}$, which is equivalent to the statement that all nonzero eigenvalues of \mathbf{P} are purely imaginary. Indeed, from (13) one finds:

$$\lambda_{\mathbf{P}} = 0, 0, \pm i\sqrt{6}; \quad \text{so } \lambda_{\mathbf{P}} \in i\mathbb{R}. \quad (14)$$

Returning to the full set of equations (7) and scalar-multiplying each of them by its respective $\underline{\mathbf{S}}^+$ or $\underline{\mathbf{S}}^-$, we notice that they admit the following “conservation” relations:

$$(\partial_t \pm \partial_x) |\underline{\mathbf{S}}^\pm|^2 = 0, \quad (15a)$$

where $|\dots|$ stands for the length of the vector. Two other conservation relations can be obtained in a similar fashion. Note that for periodic BC, considered in this paper, these relations become conservation laws. E.g., (15a) yield:

$$\partial_t \int_0^L |\underline{\mathbf{S}}^\pm|^2 dx = 0, \quad (15b)$$

where L is the length of the spatial domain. Moreover, a Hamiltonian of system (7) can be constructed in certain action-angle variables [65].

Therefore, it is desirable that the numerical method also conserve or quasi-conserve (i.e., not lead to a systematic shift in) at least some of these quantities. A simple, explicit such a method for ODEs is LF, whereas explicit Euler methods are known to lead to mild numerical instability for conservative ODEs. For that reason we had expected that using the LF as the ODE solver in the MoC would produce better results than explicit Euler solvers. However, we found that, on the contrary, it produces the worst results. We will now turn to the description and analysis of the SE, ME, and LF solvers for the MoC with periodic BC. We will refer to the corresponding “flavors” of the MoC as MoC-SE, MoC-ME, and MoC-LF, respectively.

3 MoC with the SE solver

The SE is a first-order method and, moreover, is well-known to lead to a mild yet conspicuous numerical instability when applied to conservative ODEs (see, e.g., [66]). The reason that we consider this method is that it is simple enough to illustrate the approach. Thus, describing sufficient details in this Section will allow us to skip them in subsequent sections, devoted to second-order methods, where such details are more involved.

3.1 Analysis

The form of the MoC-SE equations for system (8) is:

$$(S_j^\pm)_{m+1}^{n+1} = (S_j^\pm)_{m\mp 1}^n + h f_j^\pm((\underline{\mathbf{S}}^+)_{m\mp 1}^n, (\underline{\mathbf{S}}^-)_{m\mp 1}^n), \quad j = 1, 2, 3; \quad (16)$$

where f_j^\pm are the nonlinear functions on the rhs of (8), and $m = 0, \dots, M$, with $(M+1)$ being the number of grid points. Note that in this paper we have set the temporal and spatial steps equal, $\Delta t = \Delta x = h$, to ensure having a regular grid. To impose periodic BC, we make the following identifications in (16):

$$(\underline{\mathbf{S}}^\pm)_{-1}^n \equiv (\underline{\mathbf{S}}^\pm)_M^n, \quad (\underline{\mathbf{S}}^\pm)_{M+1}^n \equiv (\underline{\mathbf{S}}^\pm)_0^n. \quad (17)$$

Linearizing Eqs. (16), one arrives at:

$$\begin{pmatrix} \underline{s}^+ \\ \underline{s}^- \end{pmatrix}_m^{n+1} = \begin{pmatrix} \underline{s}^+ \\ \underline{0} \end{pmatrix}_{m-1}^n + \begin{pmatrix} \underline{0} \\ \underline{s}^- \end{pmatrix}_{m+1}^n + h \begin{pmatrix} P^{++} & P^{+-} \\ \mathcal{O} & \mathcal{O} \end{pmatrix} \begin{pmatrix} \underline{s}^+ \\ \underline{s}^- \end{pmatrix}_{m-1}^n + h \begin{pmatrix} \mathcal{O} & \mathcal{O} \\ P^{-+} & P^{--} \end{pmatrix} \begin{pmatrix} \underline{s}^+ \\ \underline{s}^- \end{pmatrix}_{m+1}^n, \quad (18)$$

where \mathcal{O} is the 2×2 zero matrix and

$$\begin{pmatrix} P^{++} & P^{+-} \\ P^{-+} & P^{--} \end{pmatrix} \equiv \begin{pmatrix} -A & B \\ -B & A \end{pmatrix} = \mathbf{P}. \quad (19)$$

Seeking the solution of (18) with periodic BC in the form $(\mathbf{s})_m^n = (\mathbf{s})^n e^{ikmh}$ reduces (18) to:

$$(\mathbf{s})^{n+1} = \mathbf{N}(z) (\mathbf{s})^n, \quad (20a)$$

$$\mathbf{N}(z) = \mathbf{Q}(\mathbf{I} + h\mathbf{P}), \quad \mathbf{Q} = \exp[-iz\boldsymbol{\Sigma}], \quad (20b)$$

where $z = kh$, and $\boldsymbol{\Sigma}$ is defined in (13).

Below we will show that eigenvalues of \mathbf{N} satisfy $|\lambda_{\mathbf{N}}| > 1$, which implies that the MoC-SE with periodic BC is unstable. One can easily see this for the longest- and shortest-period Fourier harmonics ($z = 0$ and $z = \pi$, respectively), where one has $\mathbf{Q}(z = 0) = \mathbf{I}$ and $\mathbf{Q}(z = \pi) = -\mathbf{I}$ and invokes relation (14). Then one has:

$$|\lambda_{\mathbf{N}}(z = 0, \pi)| = |1 + ih|\lambda_{\mathbf{P}}|| > 1. \quad (21a)$$

However, for other Fourier harmonics of the numerical error, i.e., for $z \neq 0$ or π , an analytical relation between the eigenvalues $\lambda_{\mathbf{N}}$ of the numerical method and the eigenvalues $\lambda_{\mathbf{P}}$ of the “physical” matrix \mathbf{P} cannot be established. In particular,

$$|\lambda_{\mathbf{N}}(z \neq 0, \pi)| \neq |1 + h\lambda_{\mathbf{P}}|, \quad (21b)$$

other than by accident. Therefore, to determine $|\lambda_{\mathbf{N}}|$, one has to find these eigenvalues numerically. Let us stress that this limits only the result, but *not the methodology*, of our von Neumann analysis to a specific matrix \mathbf{P} . Indeed, for any \mathbf{P} , the relation between the amplification matrix \mathbf{N} of the numerical scheme and the “physical” matrix \mathbf{P} is given by (20b). The eigenvalues of \mathbf{N} are found within a second by modern software, whereas direct numerical simulations of scheme (16) (and other MoC schemes considered in Sections 4 and 5) take much longer.

The largest (in magnitude) eigenvalue of \mathbf{N} , found by Matlab, is shown as a function of the normalized wavenumber z in Fig. 1(a). The two curves illustrate the general trend as h is varied. They also show that the strongest instability of the MoC-SE applied to (7) occurs in the ODE- and “anti-ODE” limits, i.e. for $z = 0$ and $z = \pi$. In particular, since $\mathbf{Q}(z = \pi) = -\mathbf{Q}(z = 0)$, one concludes that $|\lambda_{\mathbf{N}}(z = \pi)| = |\lambda_{\mathbf{N}}(z = 0)|$. In other words, the instability of the MoC-SE for the highest and lowest Fourier harmonics is the same, in contrast to that of other numerical methods.

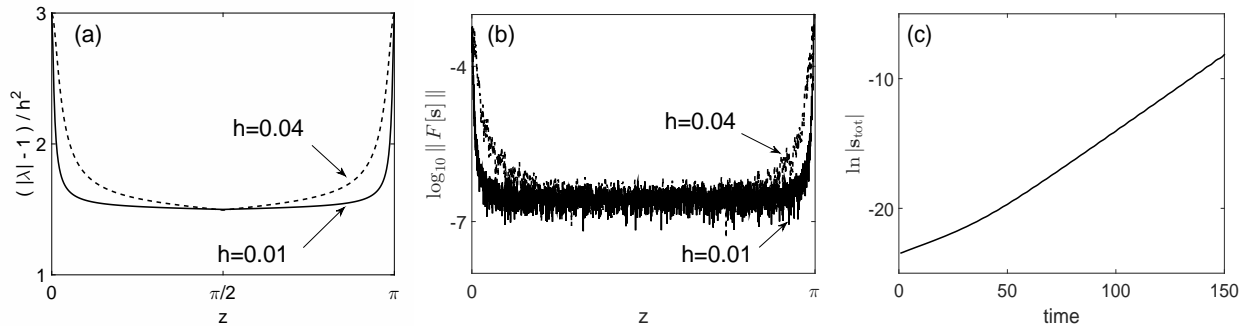


Figure 1: (a) Eigenvalue of matrix \mathbf{N} in (20) for $h = 0.01$ and $h = 0.04$. (b) Spectrum of the numerical error obtained by simulating scheme (16), (17) for $L = 100$ and $t = 520$ (for $h = 0.01$) and $t = 130$ (for $h = 0.04$). (c) Time evolution of the error (24) for $h = 0.04$.

Equations (20) illustrate the statement, made in Section 1, that MoC schemes preserve the high- k dispersion relation (6) for system (12a). Indeed, (20a) implies that the frequency ω_{numer} of the numerical solution is an eigenvalue of $-(\ln \mathbf{N})/(ih)$, where the minus sign is used due to the definition of ω after relation (6). Then, (20b) implies that for $|k| \gg 1$,

$$\omega_{\text{numer}} = -(\ln \mathbf{Q} + O(h))/(ih) = \pm k(1 + O(1/|k|)), \quad (22)$$

which agrees with (6).

3.2 Numerical verification

First, in Fig. 1(b) we show the logarithm of the Fourier spectrum of the numerical error obtained by applying the MoC-SE to (8). Namely, we simulated scheme (16) with the initial condition

being (10) plus white (in space) noise of magnitude 10^{-12} and then subtracted from the numerical solution the exact solution (10). The curves in Fig. 1(b) have the same qualitative shapes as the corresponding curves in Fig. 1(a), which confirms the validity of the analysis in Section 3.1. Indeed, it follows from (20a) that (individual) Fourier harmonics of the numerical error satisfy

$$\|(\mathbf{s})^n\| \propto |\lambda_{\mathbf{N}}|^n \approx \exp [((|\lambda_{\mathbf{N}}| - 1)/h) t], \quad \Rightarrow \quad \ln \|(\mathbf{s})^n\| \propto ((|\lambda_{\mathbf{N}}| - 1)/h) t, \quad (23)$$

where $\|\dots\|$ denotes the Euclidean norm of the vector.

Second, to verify our analytical results from yet another perspective, we will compare the growth rate measured for a “total” numerical error (see below) with the growth rate predicted by the analysis of Section 3.1. The “total” numerical error is computed as

$$|\mathbf{s}_{\text{tot}}| = \left(\sum_{m=0}^M \|(\mathbf{s})_m\|^2 \right)^{1/2}. \quad (24)$$

Note that the summation over m smooths out the noisy spatial profile of the error. The perturbations s_2^\pm were indeed found to be much smaller than $s_{1,3}^\pm$ (as long as the latter are themselves sufficiently small), as predicted by (12b), and hence they are not included in the “total” error (24). Now, according to Section 3.1, the highest growth rate occurs for harmonics with $z = 0, \pi$. It follows from (21a) that

$$|\lambda_{\mathbf{N}}(z = 0)|^n = (|1 + ih \max |\lambda_{\mathbf{P}}| |)^{t/h} \approx \exp \left[\frac{h}{2} \max |\lambda_{\mathbf{P}}|^2 t \right]. \quad (25)$$

Thus, the theoretical growth rate is

$$\gamma_{\text{theor}} = 3h, \quad (26)$$

where we have used (14).

On the other hand, the growth rate can be measured from the numerical solution:

$$\gamma_{\text{meas}} = \frac{\ln |\mathbf{s}_{\text{tot}}(t_2)| - \ln |\mathbf{s}_{\text{tot}}(t_1)|}{t_2 - t_1}, \quad (27)$$

where $t_{1,2}$ are some times when the dependence of $\ln |\mathbf{s}_{\text{tot}}|$ on t appears to be linear (as, e.g., in Fig. 1(c) for $t > 40$). Using (27), we have found for the parameters listed in the caption to Fig. 1 that γ_{meas} agrees with its theoretical value (26) to two significant figures.

4 MoC with the ME solver

The MoC-ME applied to system (8) yields the following scheme:

$$\overline{S}_j^\pm = (S_j^\pm)_{m\mp 1}^n + h f_j^\pm ((\underline{\mathbf{S}}^+)_{m\mp 1}^n, (\underline{\mathbf{S}}^-)_{m\mp 1}^n), \quad j = 1, 2, 3; \quad (28a)$$

$$(S_j^\pm)_{m+1}^{n+1} = \frac{1}{2} \left[(S_j^\pm)_{m\mp 1}^n + \overline{S}_j^\pm + h f_j^\pm (\underline{\mathbf{S}}^+_m, \underline{\mathbf{S}}^-_m) \right], \quad (28b)$$

where the notations are the same as in Section 3. Note that the rhs of (28a) is the same as that of (16). Following the analysis of Section 3.1, one obtains a counterpart of (20) where now

$$\mathbf{N} = \frac{1}{2} [\mathbf{Q} + (\mathbf{I} + h\mathbf{P}) \mathbf{Q} (\mathbf{I} + h\mathbf{P})]. \quad (29)$$

As we discussed in Section 3.1, the eigenvalues of \mathbf{N} cannot be analytically related to the eigenvalues of \mathbf{P} and hence need to be found numerically. (As for the MoC-SE before, this takes fractions of a second for modern software, unlike the much longer direct simulations of scheme (28).) The magnitude of the largest-in-magnitude such eigenvalue is shown in Fig. 2(a). The logarithm of the Fourier spectrum of $\|\mathbf{s}\|$ is plotted in Fig. 2(b) and is seen to agree qualitatively with the result in Fig. 2(a); see also (23).

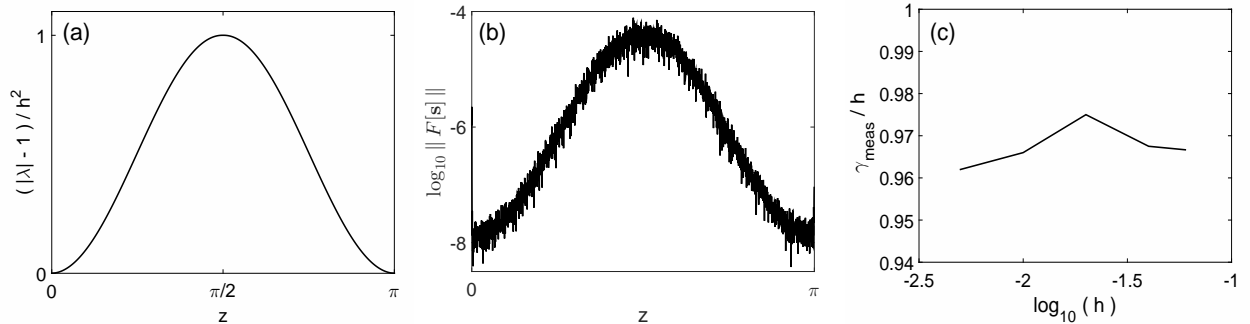


Figure 2: (a) Eigenvalue of matrix \mathbf{N} in (29) for $h = 0.01$. The curves for $h \leq 0.05$ all look qualitatively similar to the one shown here; in particular, their maximum value is 1. (b) Spectrum of the numerical error obtained by simulating scheme (28), (17) for $L = 100$ and $t = 800$ with $h = 0.01$. (c) Comparison of the growth rate of MoC-ME measured from scheme (28), (17) with that established by formula (30b). Note that the value 1 of the ratio would correspond to the perfect agreement between the numerics and analysis.

Let us note that the numerical instability of the MoC-ME scheme in the ODE and anti-ODE limits is much weaker than that of the MoC-SE, because from (29),

$$|\lambda_{\mathbf{N}}(z = 0, \pi)| = \left(1 + \frac{h^4}{4} \max |\lambda_{\mathbf{P}}|^4\right)^{1/2},$$

and hence similarly to (25), one has

$$\gamma_{\text{theor}}(z = 0, \pi) = \frac{h^3}{8} \max |\lambda_{\mathbf{P}}|^4. \quad (30a)$$

However, Fig. 2(a) shows that *unlike* the MoC-SE, the strongest instability of the MoC-ME occurs in the “*middle*”² of the Fourier spectrum. Moreover, this growth rate has the same order of magnitude as the growth rate of the MoC-SE’s instability. Therefore, it is the growth rate of the harmonics with $k \sim \pi/(2h)$ that one would measure in the numerical experiment. Using the information from Fig. 2(a) and its caption, one has

$$|\lambda_{\mathbf{N}}(z = \pi/2)| \approx 1 + h^2,$$

and hence

$$\gamma_{\text{theor}}(z = \pi/2) \approx h. \quad (30b)$$

²We used the quotes because, strictly speaking, this is the middle of the right *half* of the spectrum.

Figure 2(c) shows that the growth rates computed from the (semi-)analytical formula (30b) and measured, as explained in Section 3.2, in direct numerical simulations of scheme (28), agree reasonably well.

In Section 7 we will demonstrate that these conclusions about the numerical instability of the MoC-ME remain qualitatively valid for a different system of PDEs, which has spatially localized (as opposed to constant) coefficients; see Appendix B.

5 MoC with the LF solver

Unlike the Euler ODE solvers considered in the previous two sections, the LF involves the solution at three time levels:

$$(S_j^\pm)^{n+1} = (S_j^\pm)^{n-1} + 2h f_j^\pm((\underline{\mathbf{S}}^+)^n_{m\mp 1}, (\underline{\mathbf{S}}^-)^n_{m\mp 1}), \quad j = 1, 2, 3. \quad (31)$$

Given the initial condition, one can find the solution at time level $n = 1$ with, e.g., the MoC-SE. To enforce periodic BC, convention (17) needs to be supplemented by

$$(\underline{\mathbf{S}}^\pm)^n_{-2} \equiv (\underline{\mathbf{S}}^\pm)^n_{M-1}, \quad (\underline{\mathbf{S}}^\pm)^n_{M+2} \equiv (\underline{\mathbf{S}}^\pm)^n_1. \quad (32)$$

The counterpart of (20) is now

$$(\mathbf{s})^{n+1} - 2h\mathbf{Q}\mathbf{P}(\mathbf{s})^n + \mathbf{Q}^2(\mathbf{s})^{n-1} = \mathbf{0}. \quad (33)$$

The (in)stability of the MoC-LF is determined by whether the magnitude of the largest eigenvalue of the quadratic eigenvalue problem obtained from (33) and leading to

$$\det(\lambda^2\mathbf{I} - 2\lambda h\mathbf{Q}\mathbf{P} + \mathbf{Q}^2) = 0 \quad (34)$$

exceeds unity. As previously for the MoC-SE and MoC-ME, this 4×4 eigenvalue problem has to be solved numerically for a given \mathbf{P} . The largest magnitude of its eigenvalue, computed with Matlab's command `polyeig` for \mathbf{P} given by (13), is plotted in Fig. 3(a).

It should be noted that the quantity $(|\lambda| - 1)$ for the MoC-LF scales as h , in contrast to the cases of MoC-SE and MoC-ME, where it scales as h^2 . Therefore, the instability growth rate of the MoC-LF is: $\gamma = O(1)$, which is much greater than those rates of the MoC-SE and MoC-ME. The reason for this will come out as a by-product when we qualitatively explain, in the next paragraph, why the strongest instability of the MoC-LF for the PDE system in question occurs near $z = \pi/2$.

First note that in the ODE limit, $z = 0$, one has $\mathbf{Q}(0) = \mathbf{I}$, and hence from (34) one recovers the well-known relation between the amplification factor λ of the ODE-LF and eigenvalues of \mathbf{P} :

$$\lambda(z = 0) = h\lambda_{\mathbf{P}} \pm \sqrt{1 + (h\lambda_{\mathbf{P}})^2}. \quad (35)$$

Given $\lambda_{\mathbf{P}} \in i\mathbb{R}$, as would be the case for any stable non-dissipative system, this yields $|\lambda(z = 0)| = 1$ (here and below we assume that $h \max |\lambda_{\mathbf{P}}| < 1$). A similar situation holds in the anti-ODE limit,

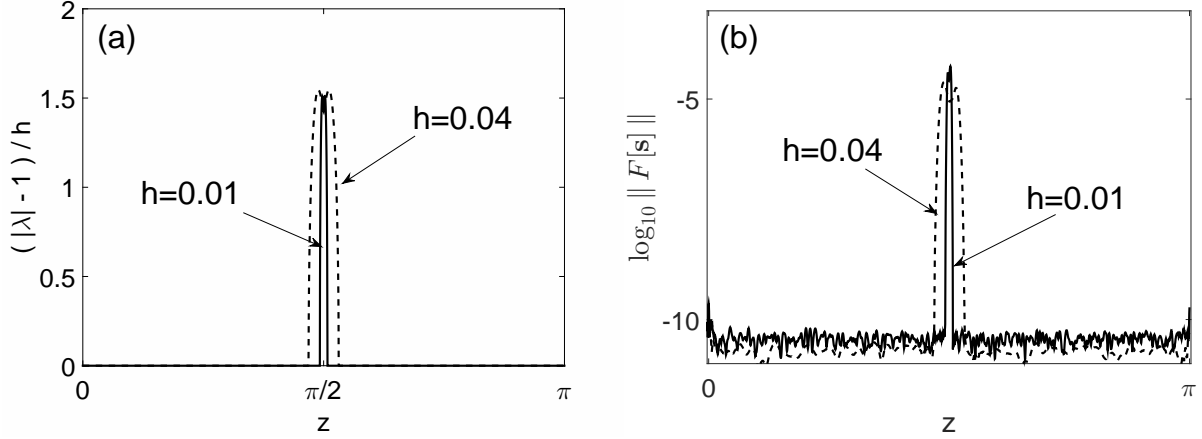


Figure 3: (a) Eigenvalue of problem (34) for $h = 0.01$ and $h = 0.04$. (b) Spectrum of the numerical error obtained by simulating scheme (31), (17), (32) for $L = 100$ and $t = 10$.

where $\mathbf{Q} = -\mathbf{I}$. However, for $z = \pi/2$, $\mathbf{Q}(\pi/2) = -i \text{diag}(I, -I)$ can be said to be the most dissimilar from $\mathbf{I} \equiv \text{diag}(I, I)$, and the counterpart of (35) is:

$$\lambda(z = \pi/2) = -h\lambda_{\mathbf{P}_{\text{mod}}} \pm \sqrt{-1 + (h\lambda_{\mathbf{P}_{\text{mod}}})^2}, \quad (36)$$

where $\mathbf{P}_{\text{mod}} = \mathbf{Q}(\pi/2)\mathbf{P}$. For $\lambda_{\mathbf{P}_{\text{mod}}} \in i\mathbb{R}$, this yields $\max |\lambda(z = \pi/2)| > 1$, i.e., a numerical instability. In the case of matrix \mathbf{P} defined in (13), $\lambda_{\mathbf{P}_{\text{mod}}} = 0, 0, \pm i\sqrt{2}$, and thus the above statements explain our result shown in Fig. 3(a). In Section 7 we will discuss how this generalizes to some other energy-preserving hyperbolic PDEs.

Moreover, it also follows from (36) that the growth rate of this instability is about $\max |\lambda_{\mathbf{P}_{\text{mod}}}| = O(1)$. (For the specific \mathbf{P} in (13), we will show in the second paper of this work that the growth rate is $3/2$ instead of $\sqrt{2} = \max |\lambda_{\mathbf{P}_{\text{mod}}}|$; this minor difference is responsible for the subtly double-peaked profile of the curves in Fig. 3(a).) We confirmed this by direct numerical simulations: representative spectra of the numerical error are shown in Fig. 3(b) and are seen to agree with the shapes of the curves in Fig. 3(a); see (23). The solution obtained by the MoC-LF becomes completely destroyed by the numerical instability around $t = 25$ regardless of h . In contrast, for the MoC-ME, the solution remains essentially unaffected by the growing numerical error over times $O(1/h) \gg 10$.

6 Comments about the MoC with the cRK solver

The fourth-order cRK method may appear as a natural choice of a higher-order ODE solver to be combined with the MoC. Indeed, not only is its accuracy for ODEs $O(\Delta t^4)$, but also the systematic error that it introduces into the numerically computed conserved quantities of energy-preserving ODEs is usually negligible for sufficiently small time steps.

However, a straightforward use of the cRK solver in the MoC framework produces an unsatisfactory numerical method. The corresponding scheme is:

$$(\underline{\mathbf{S}}^\pm)_{m+1}^{n+1} = (\underline{\mathbf{S}}^\pm)_{m\mp 1}^n + \frac{1}{6} ((\underline{\boldsymbol{\kappa}}_1^\pm)_m + 2(\underline{\boldsymbol{\kappa}}_2^\pm)_m + 2(\underline{\boldsymbol{\kappa}}_3^\pm)_m + (\underline{\boldsymbol{\kappa}}_4^\pm)_m); \quad (37a)$$

$$(\underline{\boldsymbol{\kappa}}_1^\pm)_m = h \underline{\mathbf{f}}^\pm (\underline{\mathbf{S}}_{m\mp 1}^+, \underline{\mathbf{S}}_{m\mp 1}^-); \quad (37b)$$

$$(\underline{\boldsymbol{\kappa}}_\alpha^\pm)_m = h \underline{\mathbf{f}}^\pm \left(\underline{\mathbf{S}}_{m\mp 1}^+ + \frac{1}{2}(\underline{\boldsymbol{\kappa}}_{\alpha-1}^+)_m, \underline{\mathbf{S}}_{m\mp 1}^- + \frac{1}{2}(\underline{\boldsymbol{\kappa}}_{\alpha-1}^-)_m \right), \quad \alpha = 2, 3; \quad (37c)$$

$$(\underline{\boldsymbol{\kappa}}_4^\pm)_m = h \underline{\mathbf{f}}^\pm (\underline{\mathbf{S}}_{m\mp 1}^+ + (\underline{\boldsymbol{\kappa}}_3^+)_m, \underline{\mathbf{S}}_{m\mp 1}^- + (\underline{\boldsymbol{\kappa}}_3^-)_m). \quad (37d)$$

A von Neumann analysis analogous to that presented in Section 4 yields a dependence $\max |\lambda_{\mathbf{N}}(z)|$ similar to the one shown in Fig. 2(a), and direct numerical simulations confirm that result for system (8). Thus, scheme (37) has the instability growth rate $\gamma = O(h)$, just as the MoC-ME (28). In fact, for the specific PDE system (8), this γ is twice that of the MoC-ME's, and hence scheme (37) offers no advantage over the MoC-ME.

Perhaps, a proper MoC-cRK needs to ensure that the arguments of $\underline{\mathbf{f}}^\pm$ in (37c) and (37d) be taken on the same characteristics. E.g., (37d) would then need to be replaced with

$$(\underline{\boldsymbol{\kappa}}_4^\pm)_m = h \underline{\mathbf{f}}^\pm (\underline{\mathbf{S}}_{m-1}^+ + (\underline{\boldsymbol{\kappa}}_3^+)_m, \underline{\mathbf{S}}_{m+1}^- + (\underline{\boldsymbol{\kappa}}_3^-)_m),$$

which mimics the last term in (28). The node numbers of some of the arguments of $\underline{\mathbf{f}}^\pm$ in (37c) would then also need to be modified in some way. This may result in a more stable MoC-cRK. However, a detailed exploration of this topic will have to begin with an analysis of the accuracy of the resulting method away from the ODE limit. This is outside the scope of this paper, and therefore we will not pursue it here. Let us mention that we have been unable to find any published research which would study the problem of setting up the proper MoC-cRK for the case of crossing characteristics.

7 Summary and discussion

We have presented results of the von Neumann analysis of the MoC combined with three ODE solvers: SE, ME, and LF, as applied to an energy-preserving hyperbolic PDE system, (7) or (8). Our main findings are as follows.

First, we have found that the numerical instability of the highest and lowest Fourier harmonics in the MoC schemes has the same growth rate. This is in contrast with the situation of other numerical methods, where it is the highest Fourier harmonics that become numerically unstable before the lower ones. Moreover, we have found that the harmonics in the ‘‘middle’’ of the spectrum, i.e., with $|k| \approx \pi/(2h)$, can become the most numerically unstable. This occurs, for the PDE system in question, for the MoC-ME and MoC-LF schemes. While such harmonics for the MoC-LF occupy a narrow spectral band and hence can, in principle, be filtered out, for the MoC-ME they occupy a substantial portion of the spectrum and hence cannot be filtered out without destroying the solution. It should also be pointed out that at least for some energy-preserving *PDEs*, as the one considered here, the growth rate of the most unstable harmonics has the same order of magnitude, $O(h)$, for the MoC-SE and MoC-ME. This should be contrasted with the situation for energy-preserving *ODEs*, where the instability growth rate of the ME method, $O(h^3)$, is much lower than that, $O(h)$, of the SE method.

Second, we have found that the MoC based on the LF solver, which (the solver) is known to considerably outperform the Euler schemes for energy-preserving ODEs, performs the worst for the energy-preserving PDE that we considered. Namely, the Fourier harmonics with $|k| \approx \pi/(2h)$ grow at the rate $\gamma = O(1)$, which by far exceeds the growth rate $\gamma = O(h)$ of the numerically most unstable harmonics in both the MoC-SE and MoC-ME.

In Section 6 we have pointed out that, for hyperbolic PDE systems with crossing characteristics, the formulation of a MoC based on the cRK ODE solver and having numerical stability (or, at worst, a milder instability than the MoC-ME) for $|kh| \sim \pi/2$, remains an open problem.

There remains a question as to whether the above conclusions of our analysis, which necessarily had to be obtained for a *specific* system (as we explained in the main text), apply to other energy-preserving hyperbolic PDEs. Again, for the same reasons, this cannot be answered in general, as no analytical relation between the eigenvalues of the “physical” matrix \mathbf{P} and the numerical amplification factor $|\lambda_{\mathbf{N}}|$ can be obtained for an arbitrary wavenumber k . Therefore, in Appendices A and B we consider several specific PDEs with the same linearized form (12a), where matrix \mathbf{P} has a different form than (13). In Appendix A we limit ourselves to systems with spatially constant coefficients and focus on how qualitatively different eigenvalues (see below) of \mathbf{P} affect the amplification factor for the MoC-LF. In Appendix B we will demonstrate that our conclusions about numerical instability of both MoC-LF and MoC-ME remain qualitatively valid for a different system whose linearization has spatially-*localized* entries in matrix \mathbf{P} .

Thus, we will now explore what “patterns” of numerical instability can be expected in simulations of energy-preserving PDEs by the MoC-LF. To that end, we consider a number of systems³ representing a broader class of such PDEs, which include our system (8), (10) as a special case. The total of 18 systems, including it, were considered, and a summary of results is presented in Figs. 3(a) and 4 and explained below. The description of the systems is found in Appendix A.

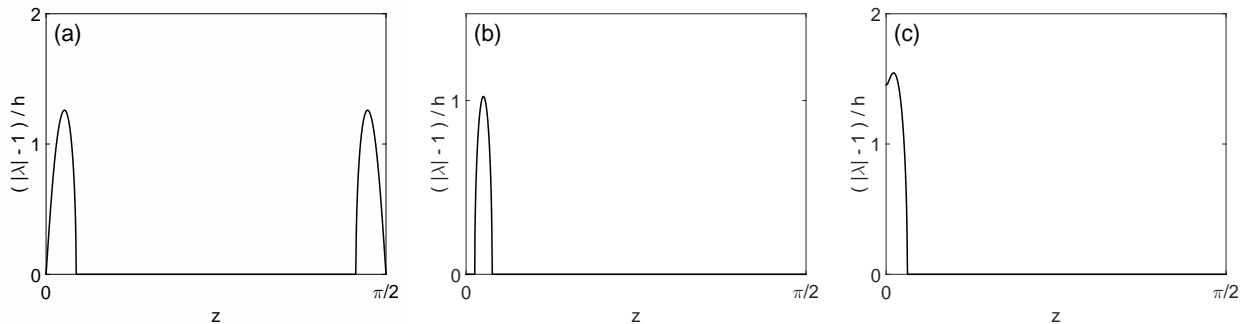


Figure 4: Eigenvalue of problem (34) for $h = 0.04$ and for the matrices \mathbf{P} corresponding to the physically unstable systems considered in Appendix A. The “other half” of the spectrum, $z \in [\pi/2, \pi]$, is reflectionally symmetric about the point $z = \pi/2$ to the one shown here. (a) $\lambda_{\mathbf{P}} = \pm i\alpha, \pm i\beta$; (b) $\lambda_{\mathbf{P}} = \pm i\alpha$ (a pair of double roots); (c) $\lambda_{\mathbf{P}} = \pm\alpha, \pm\beta$. See Section 7 for more details.

³By a “system” we will refer here to the PDE plus the stationary solution on whose background we perform the linearization. Both of them are needed to define entries in matrix \mathbf{P} .

Before we proceed, let us remind the reader that the stability of the *physical* problem (12a), as opposed to the stability of the numerical method, was defined before Eq. (14). For example, system (8), (10) is stable, which can be seen from Fig. 3(a). Indeed, the MoC-LF accurately approximates the solution for $k = O(1)$ (i.e., for $z \equiv kh = O(h)$ in Fig. 3(a)), and the fact that $|\lambda| = 1$ there corresponds to $\omega \in \mathbb{R}$ in the text before Eq. (14). Let us stress that our focus is on the presence of numerical instability of the MoC-LF for physically *stable* systems. Indeed, even if the numerical method is found to be free of numerical instability for a physically unstable system, then the corresponding solution will eventually (typically within $t = O(1)$) evolve towards the physically stable solution.

We now continue with the summary of stability results for the MoC-LF. These were obtained by solving the quadratic eigenvalue problem (34) with Matlab's command `polyeig` for 18 specific matrices \mathbf{P} corresponding to each of the systems listed in Appendix A. For all the seven of those systems which are physically stable, the dependence $\max |\lambda(z)|$ was found to be similar to that shown in Fig. 3(a). Thus, for all of the stable systems considered, the MoC-LF has numerically unstable modes with $|kh| \approx \pi/2$ and growth rate $\gamma = O(1)$. Thus, this method should be deemed unsuitable for simulation of these, and possibly other, stable energy-preserving PDEs with crossing characteristics.

Let us mention, for completeness, that for the physically unstable systems considered in Appendix A, we found two possibilities with respect to the numerical (in)stability of the MoC-LF. In the first group, there are unstable systems such that: $\lambda_{\mathbf{P}} = \pm i\alpha, \pm i\beta$ for $\alpha, \beta \in \mathbb{R}$ (i.e., the mode with $k = 0$ is stable), with $\alpha \neq \beta$ and at most one of α or β could be zero. We have found that such systems have $\omega \notin \mathbb{R}$ in the immediate vicinity of $k = 0$. For systems in this group, the MoC-LF has also numerically unstable modes with $|k|h$ near (but not exactly at) $\pi/2$. This is illustrated in Fig. 4(a). In the second group there are unstable systems with all other possibilities of $\lambda_{\mathbf{P}}$ compatible with the energy-preserving nature of the problem: $\lambda_{\mathbf{P}} = \pm i\alpha$ (two pairs of repeated imaginary eigenvalues), including the case $\alpha = 0$; and $\lambda_{\mathbf{P}} = \pm\alpha, \pm\beta$ (in this case, even the $k = 0$ mode is unstable). This is illustrated in Figs. 4(b,c), which show *no numerical* instability in the MoC-LF. However, as we have mentioned above, once the initial, physically unstable, solution evolves sufficiently near a stable one, the MoC-LF will be invalidated by the numerically unstable harmonics, which are seen in Figs. 3(a) and 4(a).

Let us now demonstrate the validity of qualitative conclusions of our von Neumann analysis for a system whose linearization (12a) has a spatially dependent matrix \mathbf{P} . This system is the soliton (see Fig. 5(a)) of the Gross–Neveu model [61] in the relativistic field theory; its details are presented in Appendix B. It has received considerable attention in the past decade both from the analytical and numerical perspectives: see [37], [17], [67]–[71], and references therein. The plots of the spectra of the numerical errors, obtained for this soliton with the MoC-LF (Fig. 5(b)) and MoC-ME (Fig. 5(c)), are qualitatively similar to such plots in Sections 5 and 4, respectively. In particular, the result in panel (b) agrees with our conclusion that the MoC-LF is unsuitable for simulations of

a physically stable hyperbolic PDE system with crossing characteristics: the numerical instability at $k \approx \pi/(2h)$ will destroy the solution soon after $t = 20$. In contrast, panel (c) shows that the MoC-ME would be a suitable method for such systems, even for relatively long times. Moreover, we also verified our conclusion (see (30b) and Fig. 2(c)) that the growth rate of the most unstable Fourier harmonics, with $k \approx \pi/(2h)$, of the MoC-ME scales as $O(h)$ (which is considerably greater than the growth rate of that method's error for energy-preserving ODEs). To that end, we repeated the simulations for several values of $h = L/M$, where $M = 2^{11}, 2^{12}, 2^{13}, 2^{14}$ and computed γ_{meas} using (27) and (44). The results fit the dependence $\gamma_{\text{meas}} = 6.4 \cdot 10^{-3} \cdot (2^{11}/M) \cdot \ln(10)$, which is equivalent to $\gamma_{\text{meas}} = O(h)$, to two significant figures.

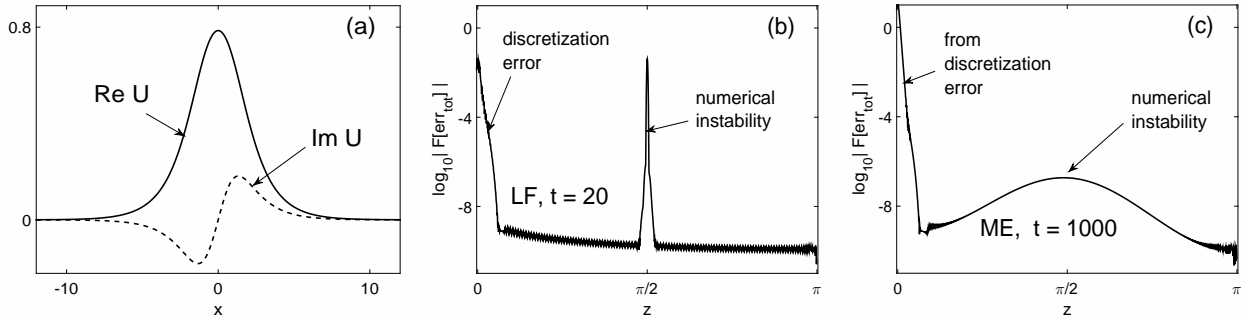


Figure 5: (a): Soliton solution (43b) for $\Omega = 0.7$. Note that $\text{Re}(V)=\text{Re}(U)$, $\text{Im}(V)=-\text{Im}(U)$. (b) and (c): Spectra of the total error (44) obtained by the MoC-LF at $t = 20$ (b) and by the MoC-ME (c) at $t = 1000$. See Section 7 and Appendix B for more details; in particular, the origin of the peaks near $z = 0$ in panels (b) and (c) is explained at the end of Appendix B.

Acknowledgement

This work was supported in part by the NSF grant DMS-1217006.

Appendix A: A broader class of PDE systems with constant coefficients

The class of models that we referred to in Section 7 generalizes Eqs. 7:

$$\underline{\mathbf{S}}_t^\pm \pm \underline{\mathbf{S}}_x^\pm = \underline{\mathbf{S}}^\pm \times \hat{\mathbf{J}}_c \underline{\mathbf{S}}^\mp + \underline{\mathbf{S}}^\pm \times \hat{\mathbf{J}}_s \underline{\mathbf{S}}^\pm, \quad (38)$$

where $\hat{\mathbf{J}}_c$ and $\hat{\mathbf{J}}_s$ are matrices accounting, respectively, for cross- and self-interaction among components of the Stokes vectors $\underline{\mathbf{S}}^\pm$. This class of models describes propagation of electromagnetic waves in optical fibers with various types of birefringence. The model considered in the main text is a special case of model (38) with

$$\hat{\mathbf{J}}_c = \alpha \text{diag}(1, -1, -2), \quad \hat{\mathbf{J}}_s = (2 - 3\alpha) \text{diag}(0, 0, 1), \quad (39a)$$

which corresponds to a highly spun birefringent fiber [41]. The parameter α , which in model (7) is set to $2/3$, accounts for the ellipticity of the fiber's core. Two other models are: that of a randomly birefringent fiber [42], with

$$\hat{\mathbf{J}}_{\mathbf{c}} = \text{diag}(-1, 1, -1), \quad \hat{\mathbf{J}}_{\mathbf{s}} = \mathcal{O}, \quad (39b)$$

and that of an isotropic (i.e., non-birefringent) fiber [40]:

$$\hat{\mathbf{J}}_{\mathbf{c}} = \text{diag}(-2, -2, 0), \quad \hat{\mathbf{J}}_{\mathbf{s}} = \text{diag}(-1, -1, 0). \quad (39c)$$

Let us stress that only the above forms of matrices $\hat{\mathbf{J}}_{\mathbf{c},\mathbf{s}}$ correspond to physically meaningful situations in the context of birefringent fibers, and thus it would not make sense to consider model (38) with arbitrary $\hat{\mathbf{J}}_{\mathbf{c},\mathbf{s}}$.

Models (38), (39) each have many stationary constant solutions. We will consider only six of them for each of those three models:

$$S_j^+ = 1, \quad S_j^- = \pm S_j^+ \quad \text{for one of } j = 1, 2, \text{ or } 3, \text{ with} \quad (40)$$

the other two components of $\underline{\mathbf{S}}^\pm$ being 0.

For brevity, we will refer to these solutions as $(j\pm)$, where j and \pm correspond to the particular choice of the component and the sign in (40). For example, solution (10) of the main text corresponds to $(2-)$ in (40).

Of these 18 systems, the following 7 are stable on the infinite line (in the sense specified before Eq. (14)): model (39a) with solutions $(1+)$, $(2-)$, $(3-)$; model (39b) with solutions $(1-)$, $(2+)$, $(3-)$; model (39c) with solution $(3+)$. The following 3 systems exhibit instability for $k = 0$: model (39a) with solutions $(1-)$, $(2+)$ and model (39c) with solution $(3-)$. (In other words, $\lambda_{\mathbf{P}} \in \mathbb{R} \setminus \{0\}$ for these systems.) The remaining 8 systems exhibit instability for perturbations with $k \neq 0$.

Appendix B: The soliton solution of the Gross–Neveu model

In physical variables, this model has the form [61, 37]:

$$\begin{aligned} i(\psi_t + \chi_x) + (|\psi|^2 - |\chi|^2)\psi - \psi &= 0, \\ i(\chi_t + \psi_x) + (|\chi|^2 - |\psi|^2)\chi + \chi &= 0. \end{aligned} \quad (41)$$

A change of variables $u = (\psi + \chi)/\sqrt{2}$, $v = (\psi - \chi)/\sqrt{2}$ transforms (41) to the form (1):

$$\begin{aligned} u_t + u_x &= i(|v|^2 u + v^2 u^*) - iv, \\ v_t - v_x &= i(|u|^2 v + u^2 v^*) - iv. \end{aligned} \quad (42)$$

The standing soliton solution of this system is (see, e.g., [37] and references therein):

$$\{u, v\} = \{U(x), V(x)\} \exp[-i\Omega t], \quad \Omega \in (0, 1); \quad (43a)$$

$$\{U(x), V(x)\} = \sqrt{1 - \Omega} \frac{\cosh(\beta x) \pm i\mu \sinh(\beta x)}{\cosh^2(\beta x) - \mu^2 \sinh^2(\beta x)}; \quad (43b)$$

with $\beta = \sqrt{1 - \Omega^2}$ and $\mu = \sqrt{(1 - \Omega)/(1 + \Omega)}$. The physical stability of this soliton for sufficiently small Ω has been an issue of recent controversy [37, 71] between analytical and numerical results. However, for Ω sufficiently away from 0, the soliton has been found to be stable by both analysis and numerics (see references in [71]). In Fig. 5(a) we show this solution for $\Omega = 0.7$.

The linearized system (42) on the background of soliton (43) has the form (12a), where the entries of matrix \mathbf{P} are localized functions of x . Therefore, a von Neumann analysis cannot be rigorously applied to it. However, one can expect that its predictions could be valid qualitatively, based on the principle of frozen coefficients. To demonstrate that this is indeed the case, we simulated system (42) with the initial condition shown in Fig. 5(a), to which a small white noise was added. The length of the computational domain was $L = 64$. The spectra of the numerical errors obtained by schemes (31), (32) (MoC-LF) and (28), (17) (MoC-ME) for $h = L/2^{12} \approx 0.016$ are shown in Figs. 5(b,c). The numerical error was defined similarly to (24):

$$\text{err}_{\text{tot}} = \left(\sum_{m=0}^M |u_m^n - U(x_m)e^{-i\Omega t_n}|^2 + |v_m^n - V(x_m)e^{-i\Omega t_n}|^2 \right)^{1/2}. \quad (44)$$

Note that the peaks near $z = 0$, seen in Figs. 5(b,c), do *not* correspond to any numerical instability. In panel (b) that peak is caused directly by the discretization error of the scheme. In panel (c), a much higher such peak is also caused by the discretization error, but indirectly. Indeed, a slight error in the propagation constant Ω , which inevitably occurs due to a limited accuracy of the numerical scheme, causes $\{u_m^n, v_m^n\}$ in (44) to differ from their respective $\{U(x_m), V(x_m)\} \exp[-i\Omega t_n]$ by $O(1)$ for $t_n \gg 1$.

References

- [1] B. Gustafsson, H.-O. Kreiss, J. Olinger, Time-dependent problems and difference methods, 2nd Ed., John Wiley & Sons, Hoboken, NJ, 2013; Sec. 7.4.
- [2] Yu.Ya. Mikhailov, Stability of some numerical schemes of the three-dimensional method of characteristics, USSR Comp. Math. Math. Phys. 8 (1968) 312–315.
- [3] A.C. Zecchin, A.R. Simpson, M.F. Lambert, von Neumann stability analysis of a method of characteristics visco-elastic pipeline model, 10th Int. Conf. on Pressure Surges (Edinburgh, Scotland), 2008.
- [4] G. Morthier, P. Vankwikelberge, Handbook of distributed feedback laser diodes, Artech House, Boston, 1997.
- [5] R. Kashyap, Fiber Bragg gratings, Academic, Burlington, MA, 2010; Chap. 4.
- [6] T.I. Lakoba, Z. Deng, Stability analysis of the numerical Method of characteristics applied to a class of energy-preserving hyperbolic systems. Part II: Nonreflecting boundary conditions, J. Comput. Appl. Math. <https://doi.org/10.1016/j.cam.2019.01.042>.

- [7] C.M. Schober, T.H. Wlodarczyk, Dispersion, group velocity, and multisymplectic discretizations, *Math. Comput. Simul.* 80 (2009) 741–751.
- [8] J.D. Hoffman, Numerical methods for engineers and scientists, McGraw-Hill, New York, 1992; Chap. 15.
- [9] K. Tselios, T.E. Simos, Runge–Kutta methods with minimal dispersion and dissipation for problems arising from computational acoustics, *J. Comput. Appl. Math.* 175 (2005) 173–181.
- [10] J. Berland, C. Bogey, C. Bailly, Low-dissipation and low-dispersion fourth-order Runge–Kutta algorithm, *Comp. Fluids* 35 (2006) 1459–1463.
- [11] J. Martin-Vaquero, A.H. Encinas, A. Queiruga-Dios, V. Gayoso-Martinez, A. Martin del Rey, Numerical schemes for general Klein–Gordon equations with Dirichlet and nonlocal boundary conditions, *Nonlin. Anal.: Model. Control* 23 (2018) 50–62.
- [12] W. Bao, X. Dong, Analysis and comparison of numerical methods for the Klein–Gordon equation in the nonrelativistic limit regime, *Numer. Math.* 120 (2012) 189–229.
- [13] J. Rashidinia, M. Ghasemi, R. Jalilian, Numerical solution of the nonlinear Klein–Gordon equation, *J. Comput. Appl. Math.* 233 (2010) 1866–1878.
- [14] M. Dehghan, A. Shokri, Numerical solution of the nonlinear Klein–Gordon equation using radial basis functions, *J. Comput. Appl. Math.* 230 (2009) 400–410.
- [15] R.J. LeVeque, J. Olinger, Numerical methods based on additive splittings for hyperbolic partial differential equations, *Math. Comp.* 40 (1983) 469–497.
- [16] Z. Toroker, M. Horowitz, Optimized split-step method for modeling nonlinear pulse propagation in fiber Bragg gratings, *J. Opt. Soc. Am. B* 25 (2008) 448–457.
- [17] W.Z. Bao, Y.Y. Cai, X.W. Jia, J. Yin, Error estimates of numerical methods for the nonlinear Dirac equation in the nonrelativistic limit regime, *Science China Math.* 59 (2016) 1461–1494.
- [18] F. de la Hoz, F. Vaddillo, An integrating factor for nonlinear Dirac equations, *Comp. Phys. Commun.* 181 (2010) 1195–1203.
- [19] A. Mohebbi, Z. Asgari, A. Shahrezaee, Fast and high accuracy numerical methods for the solution of nonlinear Klein–Gordon equations, *Z. Naturforsch.* 66a (2011) 735–744.
- [20] C.M. de Sterke, K.R. Jackson, B.D. Robert, Nonlinear coupled-mode equations on a finite interval: a numerical procedure, *J. Opt. Soc. Am. B* 8 (1991) 403–412.
- [21] N.G.R. Broderick, C.M. de Sterke, K.R. Jackson, Coupled mode equations with free carrier effects: a numerical solution, *Opt. Quant. Electron* 26 (1994) S219–S234.

- [22] Y.-H. Liao, H.G. Winful, Extremely high-frequency self-pulsations in chirped grating distributed-feedback semiconductor lasers, *Appl. Phys. Lett.* 69 (1996) 2989–2991.
- [23] J. Chi, A. Fernandez, L. Chao, Comprehensive modeling of wave propagation in photonic devices, *IET Commun.* 6 (2012) 473–477.
- [24] P.S. Westbrook, K.S. Abedin, T. Kremp, Distributed Feedback Raman and Brillouin Fiber Lasers, In: *Raman Fiber Lasers*, Y. Feng Ed., Springer Series in Optical Sciences 207, 2017; pp. 235–272.
- [25] B.I. Mantsyzov, Gap 2π pulse with an inhomogeneously broadened line and an oscillating solitary wave, *Phys. Rev. A* 51 (1995) 4939–4943.
- [26] R.A. Vlasov, A.M. Lemeza, Bistable moving optical solitons in resonant photonic crystals, *Phys. Rev. A* 84 (2011) 023828.
- [27] V.V. Kozlov, S. Wabnitz, Instability of optical solitons in the boundary value problem for a medium of finite extension, *Lett. Math. Phys* 96 (2011) 405–413.
- [28] M. Merklein, I.V. Kabakova, T.F.S. Büttner, D.-Y. Choi, B. Luther-Davies, S.J. Madden, B.J. Eggleton, Enhancing and inhibiting stimulated Brillouin scattering in photonic integrated circuits, *Nat. Commun.* 6 (2015) 6396.
- [29] M. Matsumoto, G. Miyashita, Efficiency and stability of pulse compression using SBS in a fiber with frequency-shifted loopback, *IEEE Photon. Technol. Lett.* 29 (2017) 3–6.
- [30] R. Hammer, W. Pötz, A. Arnold, A dispersion and norm preserving finite difference scheme with transparent boundary conditions for the Dirac equation in (1+1)D, *J. Comput. Phys.* 256 (2014) 728–747.
- [31] T.I. Lakoba, Long-time simulations of nonlinear Schrödinger-type equations using step size exceeding threshold of numerical instability, *J. Sci. Comp.* 72 (2016) 14–48.
- [32] T. Colonius, Numerically nonreflecting boundary and interface conditions for compressible flow and aeroacoustic computations, *AIAA J.* 35 (1997) 1126–1133.
- [33] D. Givoli, High-order local non-reflecting boundary conditions: a review, *Wave Motion* 39 (2004) 319–326.
- [34] I. Alonso-Mallo, A.M. Portillo, Time exponential splitting technique for the Klein–Gordon equation with Hagstrom–Warburton high-order absorbing boundary conditions, *J. Comput. Phys.* 311 (2016) 196–212.
- [35] J. Berland, C. Bogey, O. Mardsen, C. Bailly, High-order, low dispersive and low dissipative explicit schemes for multiple-scale and boundary problems, *J. Comput. Phys.* 224 (2007) 637–662.

- [36] J.A.C. Weideman, L.N. Trefethen, The eigenvalues of second-order spectral differentiation matrices, *SIAM J. Numer. Anal.* 25 (1988) 1279–1298.
- [37] S. Shao, N.R. Quintero, F.G. Mertens, F. Cooper, A. Khare, A. Saxena, Stability of solitary waves in the nonlinear Dirac equation with arbitrary nonlinearity, *Phys. Rev. E* 90 (2014) 032915.
- [38] T.I. Lakoba, Numerical study of solitary wave stability in cubic nonlinear Dirac equations in 1D, *Phys. Lett. A* 382 (2018) 300–308.
- [39] V.E. Zakharov, A.V. Mikhailov, Polarization domains in nonlinear media, *JETP Lett.* 45 (1987) 349–352.
- [40] S. Pitois, G. Millot, S. Wabnitz, Nonlinear polarization dynamics of counterpropagating waves in an isotropic optical fiber: theory and experiments, *J. Opt. Soc. B* 18 (2001) 432–443.
- [41] S. Wabnitz, Chiral polarization solitons in elliptically birefringent spun optical fibers, *Opt. Lett.* 34 (2009) 908–910.
- [42] V.V. Kozlov, J. Nuno, S. Wabnitz, Theory of lossless polarization attraction in telecommunication fibers, *J. Opt. Soc. B* 28 (2011) 100–108.
- [43] Z. Deng, Uncommon numerical instability in the method of characteristics applied to hyperbolic equations, M.S. Thesis, University of Vermont, 2016.
- [44] R.W. Boyd, *Nonlinear optics*, Academic, San Diego, 1992.
- [45] D.J. Kaup, The first-order perturbed SBS equations, *J. Nonlin. Sci.* 3 (1993) 427–443.
- [46] C.E. Mungan, S.D. Rogers, N. Satyan, J.O. White, Time-dependent modeling of Brillouin scattering in optical fibers excited by a chirped diode laser, *IEEE J. Quant. Electron.* 48 (2012) 1542–1546.
- [47] F.Y.F. Chu, A.C. Scott, Inverse scattering transform for wave-wave scattering, *Phys. Rev. A* 12 (1975) 2060–2064.
- [48] D.J. Kaup, A. Rieman, A. Bers, Space-time evolution of nonlinear three-wave interactions. I. Interactions in an homogeneous medium, *Rev. Mod. Phys.* 51 (1979) 275–310.
- [49] D.J. Kaup, Simple harmonic generation: an exact method of solution, *Stud. Appl. Math.* 59 (1978) 25–35.
- [50] E. Ibragimov, A. Struthers, D.J. Kaup, Parametric amplification of chirped pulses in the presence of a large phase mismatch, *J. Opt. Soc. Am. B* 18 (2001) 1872–1876.
- [51] C.J. McKinstrie, L. Mejling, M.G. Raymer, K. Rottwitt, Quantum-state-preserving optical frequency conversion and pulse reshaping by four-wave mixing, *Phys. Rev. A* 85 (2012) 053829.

- [52] D.V. Reddy, M.G. Raymer, C.J. McKinstrie, Efficient sorting of quantum-optical wave packets by temporal-mode interferometry, *Opt. Lett.* 39 (2014) 2924–2927.
- [53] C.J. McKinstrie, D.S. Cargill, Simultaneous frequency conversion, regeneration and reshaping of optical signals, *Opt. Expr.* 20 (2012) 6881–6886.
- [54] C.J. McKinstrie, E.J. Turano, Spatiotemporal evolution of parametric instabilities driven by short laser pulses: One-dimensional analysis, *Phys. Plasmas* 3 (1996) 4683–4696.
- [55] C.J. McKinstrie, V.A. Smalyuk, R.E. Giacone, H.X. Vu, Power exchange between crossed laser beams and the associated frequency cascade, *Phys. Rev. E* 55 (1997) 2044–2047.
- [56] J.E. Sipe, C.M. de Sterke, B.J. Eggleton, Rigorous derivation of coupled mode equations for short, high-intensity grating-coupled, co-propagating pulses, *J. Mod. Opt.* 49 (2002) 1437–1452.
- [57] S.A.M.S. Chowdhury, J. Atai, Stability of Bragg grating solitons in a semilinear dual core system with dispersive reflectivity, *IEEE J. Quant. Electron.* 50 (2014) 458–465.
- [58] M. Homar, J.V. Moloney, M. San Miguel, Traveling wave model of a multimode Fabry–Pérot laser in free running and external cavity configurations, *IEEE J. Quant. Electron.* 32 (1996) 553–566.
- [59] A.M. Yacomotti, L. Furfaro, X. Hachair, F. Pedaci, M. Giudici, J. Tredicce, J. Javaloyes, S. Balle, E.A. Viktorov, P. Mandel, Dynamics of multimode semiconductor lasers, *Phys. Rev. A* 69 (2004) 053816.
- [60] W.E. Thirring, A soluble relativistic field model, *Ann. Phys.* 3 (1958) 91–112.
- [61] D.J. Gross, A. Neveu, Dynamical symmetry breaking in asymptotically free field theories, *Phys. Rev. D* 10 (1974) 3235–3253.
- [62] M. Chugunova, D. Pelinovsky, Block-diagonalization of the symmetric first-order coupled-mode system, *SIAM J. Appl. Dyn. Syst.* 5 (2006) 55–83.
- [63] A.B. Aceves, S. Wabnitz, Self-induced transparency solitons in nonlinear refractive periodic media, *Phys. Lett. A* 141 (1989) 37–42.
- [64] M. Romagnoli, S. Trillo, S. Wabnitz, Soliton switching in nonlinear couplers, *Opt. Quantum Electron.* 24 (1992) S1237–1267.
- [65] E. Assemat, A. Picozzi, H.-R. Jauslin, D. Sugny, Hamiltonian tools for the analysis of optical polarization control, *J. Opt. Soc. B* 29 (2012) 559–571.
- [66] D.F. Griffiths, D.J. Higham, Numerical methods for ordinary differential equations, Springer-Verlag, London, 2010; Chaps. 13 and 15.

- [67] G. Berkolaiko, A. Comech, On spectral stability of solitary waves of nonlinear Dirac equation in 1D, *Math. Model. Nat. Phenom.* 7 (2012) 13–31.
- [68] D.E. Pelinovsky, A. Stefanov, Asymptotic stability of small gap solitons in the nonlinear Dirac equations, *J. Math. Phys.* 53 (2012) 073705.
- [69] N. Boussaïd, A. Comech, On spectral stability of the nonlinear Dirac equation, *J. Funct. Anal.* 271 (2016) 1462–1524.
- [70] J. Xu, S. Shao, H. Tang, Numerical methods for nonlinear Dirac equation, *J. Comput. Phys.* 245 (2013) 131–149.
- [71] J. Cuevas-Maraver, P.G. Kevrekidis, A. Saxena, F. Cooper, F.G. Mertens, Solitary waves in the nonlinear Dirac equation at the continuum limit: Stability and dynamics, In: *Ord. Part. Diff. Eqs.*, Nova Science, Boca Raton, 2015; Chap. 4.