

Optimized prefactored compact schemes for wave propagation phenomena

Aldo Rona* and Edward Hall†

Department of Engineering, University of Leicester, Leicester, LE1 7RH, England

Ivan Spisso‡

*SuperComputing Applications and Innovation (SCAI) Department,
Cineca, via Magnanelli 6/3, 40033 Casalecchio di Reno, Italy*

A new family of prefactored cost-optimized schemes is developed to minimize the computational cost for a given level of error in linear wave propagation applications, such as aerodynamic sound propagation. This work extends the theory of Pirozzoli¹ to the prefactored compact high-order schemes of Hixon,² which are MacCormack type schemes that use discrete Padé approximations. An explicit multi-step Runge-Kutta scheme advances the states in time. Theoretical predictions for spatial and temporal error bounds are used to drive the optimization process.

Theoretical comparisons of the cost-optimized schemes with a classical benchmark scheme are made. Then, two numerical experiments assess the computational efficiency of the cost-optimised schemes for computational aeroacoustic applications. A polychromatic sinusoidal test-case verifies that the cost-optimized schemes perform according to the design high-order accuracy characteristics for this class of problems. For this test case, upwards of a 50% computational cost-saving at the design level of error is recorded. The final test case shows that the cost-optimized schemes can give substantial cost savings for problems where a fully broadband signal needs to be resolved.

I. Introduction

A. Challenges in modelling wave generation and propagation phenomena

Models for the propagation of waves in a continuum are developed across the full spectrum of physical sciences, including aeroacoustics, where increasingly stringent aircraft noise regulations^{3,4} promote the development of accurate and affordable methods for predicting aerodynamically generated noise. Enhanced Computational Aeroacoustic (CAA) schemes are sought that can model sound generation and its propagation as part of the industrial design process, where predictions of accuracy compatible with the design specifications are required at an affordable cost, produced within specific time-constraints, earlier in the design process, using multi-processor computer clusters.

Computational aeroacoustic schemes that model the sound generation and propagation as one simulation typically face a number of challenges, such as in the direct approach to modelling trailing edge noise, which is a typical exemplary aeroacoustic problem. Small-scale flow structures, of the size of the boundary layer thickness, cross a trailing edge, where they experience large amplitude fluctuations in momentum. The interaction generates small amplitude acoustic pressure perturbations of long wavelength, compared to the boundary layer thickness. These radiate at the speed of sound, which is much larger than the flow structure mean convection speed.

In a direct computational aeroacoustic simulation of the trailing edge problem, the computational domain size scales with the long wavelength of the acoustic pressure perturbation whereas the spatial discretization

*Thermofluids Senior Lecturer, MAIAA

†Marie Curie ITN Technical Project Manager

‡HPC consultant for academic and industrial CFD applications

scales with the small flow structures. This separation of length scales results in large sized computational meshes. Furthermore, supporting large amplitude localised momentum fluctuations and small amplitude acoustic perturbations in the same solution requires numerical schemes with low dispersion and low dissipation characteristics, to prevent the low-amplitude acoustic signal from being distorted and/or suppressed by the flow solver numerical viscosity.

B. Prefactored compact finite-difference schemes

Different aeroacoustic problems exhibit different flow physics. Therefore no single algorithm is available to model all problems with adequate resolution and accuracy. Low Mach number acoustically active flows typically involve capturing complex features that are nevertheless computationally smooth, that is, these features do not involve any sharp change in the flow state, such as from a shock. In these circumstances, it is computationally advantageous to use higher than second-order (high-order) schemes that can achieve exponential (e^{-aN}) convergence rates by increasing the scheme's order as opposed to algebraic (N^{-b}) convergence rates by increasing the spatial mesh refinement; here, N is the number of degrees of freedom and a and n are positive real numbers. Several numerical methods for aeroacoustics have emerged in the last two decades^{5–7} with attendant applications documented in four proceedings of the Computational Aeroacoustic (CAA) workshops on benchmark problems.^{8–10} Lele¹¹ pioneered the use of Padé type compact and explicit optimized schemes in aeroacoustics. This work highlighted the requirement for special near-boundary treatment, driven by the longer finite-difference stencils used in these high-order methods. Hixon² introduced in 1996 a prefactorization method to reduce the non-dissipative central-difference stencil of the compact schemes to two lower-order biased stencils which have easily soluble matrices of reduced size. The advantages of these schemes over traditional compact schemes arise from their reduced stencil size and the independent nature of the resultant factored matrices. By reducing the stencil size of the compact schemes, the prefactorization method reduces the depth of the boundary frame over the perimeter of the computational domain where *ad-hoc* boundary stencils are required. This simplifies the specification of the boundary conditions.¹² Ashcroft and Zhang¹³ has extended the prefactorization method to a broader class of compact schemes using a more general derivation strategy, which combines Fourier analysis with the notion of a numerical wavenumber. This class of optimized prefactored schemes exhibits better wave propagation characteristics than the standard prefactored compact ones.

C. Optimization of the numerical methods

Several finite-difference explicit and compact methods are now available for solving wave propagation problems in a low Mach number flow. It is typical to involve a high performance computer cluster for such a computation, where the cost of the computation is important. It is therefore of interest to develop and implement an optimization strategy on the computational cost, based on an acceptable level of numerical accuracy of the results. The issue of computational efficiency of finite-difference schemes has been investigated in details by Colonius and Lele⁷ and by Spisso and Rona.¹⁴ These authors considered the dispersive and dissipative characteristics of several types of spatial discretization. The error associated with approximate time integration is usually considered separately from the spatial error. Spatial and temporal errors combine in finite difference time-marching schemes to determine the overall numerical accuracy of the solution. Pirozzoli¹ developed in 2007 a general strategy for the analysis of finite-difference schemes for wave propagation problems, trying to involve time integration in the analysis in a natural way. The analysis of the global discretization error has shown the occurrence of two approximately independent sources of error, associated with the discretizations in space and time. The improvement of the performance of the global scheme can be achieved by trying to separately minimize the two contributions. The analysis leads to logical and simple criteria for deriving optimized space- and time-discretization schemes, based on the concepts of spatial and temporal resolving efficiencies. That is achievable by a careful design of the space- and time-discretization schemes, as well as an appropriate choice of the grid spacing and of the time step. This analysis points towards a substantial computer time saving.

D. Obtaining and testing optimized pre-factored compact schemes

The objectives of the present work are to extend the computational cost optimization method of Pirozzoli¹ to prefactored compact schemes, assess by numerical experiment the actual computational cost of the newly

developed optimized schemes, and verify that the desired level of numerical accuracy of the solution is achieved. In completing this aim, several challenges were overcome. For the spatial differentiation, the compact stencils optimized for set levels of error required an analytical pre-factorization that retained the non-dissipative characteristics of the equivalent compact centred scheme and satisfied additional symmetry relationships of pre-factored functions. Secondly, the coupling of the optimized spatial discretization with a separately optimized temporal integration by Bernardini and Pirozzoli¹⁵ for Runge-Kutta schemes was used to obtain a space and time cost-optimized time-marching scheme for a given design level of error.

The numerical tests for verifying the level of numerical accuracy and computational cost of the optimized schemes posed a number of challenges. Each accuracy set-point used for optimizing the scheme produced one set of coefficients for the discrete pre-factored spatial differentiation and temporal integration. Each set required testing over a discrete Courant and wavenumber space of 8030 points to obtain numerically iso-cost and iso-accuracy maps to compare with the a-priori analytical predictions. These computations were performed on the High Performance Computer cluster Alice at the University of Leicester. The cluster is an assembly of 208 pairs of 2.6 GHz CPUs with 64GB RAM per node pair. One segregated node with segregated memory was used for all the numerical tests to prevent other jobs running concurrently on the cluster from affecting the assessment of the computational cost.

E. Paper outline

The paper is organized as follows: in Section II, the theory of the cost-optimization of Pirozzoli,¹ based on the optimization of the computational cost for a given error level, is reviewed and the approximate decoupling of the spatial and temporal error is introduced. The de-coupled errors are used in Section III to obtain a new family of prefactored cost-optimized schemes, with three different levels of numerical accuracy. In this section, the algebraic symmetry properties of the prefactored stencils from Hixon² are used to define a new spatial differentiation for the cost-optimized compact schemes that is prefactored. The spatial differentiation is coupled with a Runge-Kutta time integration that is cost-optimised to a matching level of accuracy by the cost-optimization process of Bernardini and Pirozzoli.¹⁵ The theoretical performance of the optimised schemes are then investigated for a range of cost-levels. Numerical tests in one spatial dimensions are reported in Section IV and compare the numerical accuracy achieved by the computational schemes against the design values from Section III. Section IV also addresses whether the analytic cost function used for the cost-optimization is a good representation of CPU time recorded during the actual computations. Finally, Section IV investigates the performance gain achieved by the optimised schemes versus their non-optimized counterparts. Conclusions from the current work and future perspectives are presented in Section V.

II. Optimization of finite-difference schemes for wave propagation phenomena

A. Compact finite difference schemes

By performing a Taylor's expansion of a smooth function $f : \mathbb{R} \rightarrow \mathbb{R}$, its derivatives can be shown to satisfy the following relations

$$\sum_{j=-P}^Q \alpha_j \frac{\partial f}{\partial x}(x + jh) = \frac{1}{h} \sum_{j=-R}^S a_j f(x + jh) + \mathcal{O}(h^n), \quad (1)$$

where $(S + R) \geq 1$, $Q \leq S$, $P \leq R$, P and Q are non-negative integers, and $h > 0$. Consider now a set of uniformly spaced nodes indexed by i along the \mathbb{R} axis, as shown in Figure 1, where $x_i = (i - 1)h$, for $1 \leq i \leq N$ and $h = L/(N - 1)$. Letting $f_i = f(x_i)$, then finite difference approximations f'_i of the first order derivatives $\frac{\partial f}{\partial x}(x_i)$ can be found using the $(R + S + 1)$ stencil:

$$\sum_{j=-P}^Q \alpha_j f'_{i+j} = \frac{1}{h} \sum_{j=-R}^S a_j f_{i+j}. \quad (2)$$

Selecting $P = Q = 0$ gives an explicit scheme, while other choices yield implicit schemes, otherwise known as Padé or compact schemes. Formally, from Eq. (1), the approximations f'_i have a truncation error of $\mathcal{O}(h^n)$ and standard techniques seek to choose coefficients α_j and a_j , which give the largest possible n for a given stencil width. For brevity, we use $CPQRS$ to denote compact schemes. Of particular interest in CAA is the measure of error in the wave propagation characteristics of a single Fourier component of $f(x)$.

A monochromatic wave of wavelength λ and wavenumber k has a scaled wavenumber $\kappa = kh$ with support $0 \leq \kappa \leq \pi$. Upon taking Fourier coefficients of Eq. (1) and ignoring the contributions from the higher order terms, the scaled pseudo-wavenumber of scheme Eq. (2) can be defined by

$$\bar{\kappa}(\kappa) = \frac{1}{i} \frac{\sum_{j=-R}^S a_j e^{ij\kappa}}{\sum_{j=-P}^Q \alpha_j e^{ij\kappa}}. \quad (3)$$

If $\bar{\kappa} = \kappa$ could be achieved over the entire wavenumber range, then the finite difference estimate of f'_i by Eq. (2) would be exact. In practice, the discrepancies between κ and $\bar{\kappa}$ increase as $\kappa \rightarrow \pi$, driven by the analytical form of the Padé approximation by which $\bar{k} = 0$ at $\kappa = \pi$. $\bar{\kappa}$ is, in general, complex valued and its real part corresponds to dispersive error, while its imaginary part is related to dissipation error. In computational aeroacoustics, where disturbances are transmitted over large distances, it is important to minimise dissipative error. Selecting $R = S$, $Q = P$ and letting a_j be antisymmetric and α_k be symmetric produces a scheme with zero dissipation. This desirable property is included in the current work by restricting the stencil optimisation to tridiagonal compact schemes ($P = Q = 1$) with five point stencil ($R = S = 2$) *i.e.*, C1122 schemes. In this case the scaled pseudo-wavenumber is a real value given by

$$\bar{\kappa}(\kappa) = \frac{\sum_{j=1}^2 2a_j \sin(j\kappa)}{1 + 2\alpha_1 \cos(j\kappa)}. \quad (4)$$

This work shows the optimization of the prefactored compact stencils to the one parameter family of C1122 schemes given by

$$\begin{cases} \alpha_0 = 1, \\ a_0 = 0, \\ a_1 = -a_{-1} = \frac{1}{3}(\alpha_1 + 2), \\ a_2 = -a_{-2} = \frac{1}{12}(4\alpha_1 - 1). \end{cases} \quad (5)$$

The choice $\alpha_1 = 1/3$ yields the sixth-order C1122 scheme. The analysis in this article is, in principle, extensible to more general *CPQRS* compact schemes.

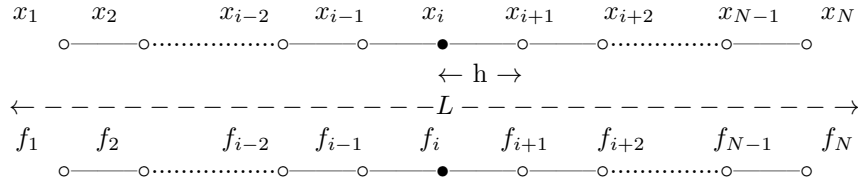


Figure 1. Variation of discrete function $f_i = f(x_i)$ along uniformly discretised streamwise length L .

B. Runge–Kutta time stepping

The semi-discretisation of a general conservation law by means of finite differences gives rise to a non-autonomous system of ODEs of the form

$$\frac{d\mathbf{U}}{dt} = \mathbf{F}(\mathbf{U}(t), t), \quad (6)$$

$$\mathbf{U}(t_0) = \mathbf{U}^0, \quad (7)$$

where \mathbf{U} represents the vector containing the solution values at spatial nodes. The solution can be time-marched from $t = t_n$ to $t = t_n + \Delta t$ by means of an explicit p -stage, two level Runge-Kutta scheme, which is formulated as

$$\mathbf{U}^{(0)} = \mathbf{U}^n, \quad (8)$$

$$\mathbf{U}^{(l)} = \mathbf{U}^n + \beta_l \Delta t \mathbf{F}(\mathbf{U}^{(l-1)}), \quad l = 1, \dots, p, \quad (9)$$

$$\mathbf{U}^{n+1} = \mathbf{U}^{(p)}, \quad (10)$$

where the coefficients α_l can be determined to achieve a given formal order of accuracy by means of a Taylor expansion, or, for example, to minimise dissipation and phase errors.^{16,17} Assuming $\mathbf{F}(\mathbf{U}) = A\mathbf{U}$, where A is linear, the scheme can be rewritten as

$$\mathbf{U}^{n+1} = \mathbf{U}^n + \sum_{j=1}^p (\Delta t)^p \gamma_j A^p \mathbf{U}^n, \quad (11)$$

where $\gamma_j = \Pi_{l=p-j+1}^p \beta_l$. This article uses the 4-stage Runge-Kutta scheme with $\gamma_1 = 1$, $\gamma_2 = 1/2$ and γ_3, γ_4 free parameters. Setting $\gamma_3 = 1/3!$ and $\gamma_4 = 1/4!$ yields a formally 4th order accurate scheme.

C. Linear advection equation analysis

The one dimensional single-scale problem

The one dimensional linear advection equation (LAE)

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0, \quad u(x, 0) = u_0(x), \quad (12)$$

where c is the wave speed, is a simple platform for developing and testing appropriate schemes for more complicated, multi-dimensional, multi-variable problems. Imposing the sinusoidal monochromatic initial condition $u_0(x) = \hat{u}_0 e^{ikx}$ on infinite and periodic domains gives rise to the solution $u(x, t) = \hat{u}_0 e^{i(-\omega t + kx)}$, where the angular frequency ω and the wavenumber k are related by the dispersion relation $\omega = ck$.

For the p -stage Runge-Kutta scheme defined in Eq. (11) and the compact scheme from Eq. (1), the amplification factor $r(\kappa, \sigma)$ from one time level to the next can be obtained¹⁸ as

$$r(\kappa, \sigma) = 1 + \sum_{j=1}^p \gamma_j (-i\sigma \bar{\kappa}(\kappa))^j, \quad (13)$$

where $\sigma = c\Delta t/h$ is the Courant number and $j = 1, \dots, p$. In the case of null spatial error, *i.e.*, when $\bar{\kappa} = \kappa$, the amplification factor is

$$r_t(z, \gamma_j) = 1 + \sum_{j=1}^p \gamma_j (-iz)^j, \quad (14)$$

where $z = \sigma\kappa$. The stability limit z_s is then given by the following condition:

$$z_s = \max\{z : 0 < \bar{z} \leq z, |r_t(\bar{z}, \gamma_j)| \leq 1\}. \quad (15)$$

Let $u_h(\cdot, \cdot)$ denote an approximation to $u(\cdot, \cdot)$. Following the work of Pirozzoli,¹ we consider the relative error in the L_2 -norm across one wavelength defined as

$$E = \frac{\|u_h(\cdot, T) - u(\cdot, T)\|_2}{\|u_0(\cdot)\|_2}, \quad (16)$$

where

$$\|w\|_2 = \left(\frac{1}{\lambda} \int_{x_0}^{x_0+\lambda} |w(x)|^2 dx \right)^{1/2},$$

and $\lambda = 2\pi/k$ is the wavelength. It can then be shown,¹ that, to a good approximation,

$$E \approx (ckT) \frac{|r(k, \sigma) - e^{-i\sigma\kappa}|}{\sigma\kappa}. \quad (17)$$

Multi-dimensional single-scale problem

For a problem in d space dimensions, the cost C_d of its numerical solution can be shown by dimensional arguments⁷ to scale as

$$C_d \propto p N_{\text{op}} \left(\frac{T}{\Delta t} \right) \left(\frac{L}{h} \right)^d, \quad (18)$$

where $(L/h)^d$ is the number of points in the domain, N_{op} is the required number of operation per mesh node, p is the number of RK stages and $T/\Delta t$ is the number of time steps. By introducing non-dimensional groups ckT and kL , which are, respectively, measures of the number of wavelengths travelled by a wave in time T and the number of wavelengths contained in the computational domain L , the normalized error e and normalized cost c_d are obtained

$$e(\kappa, \sigma) \equiv \frac{E}{ckT} = \frac{|r(k, \sigma) - e^{-i\sigma\kappa}|}{\sigma\kappa}, \quad (19)$$

$$c_d(\kappa, \sigma) \equiv \frac{C_d}{(ckT)(kL)^d} = pN_{\text{op}} \frac{1}{\sigma\kappa^{d+1}}. \quad (20)$$

The normalized error represents the numerical error incurred in advecting the wave over one period. Similarly, the normalized cost represents the cost of advecting the wave in one dimension over the same period.

Multi-dimensional multi-scale problems

In practice, aeroacoustic signals are composed of a number of waves of differing wavelengths. It is therefore of interest to extend the error and cost definitions from Eq. (19) and Eq. (20) so they apply to broadband signals, supposing these have a spectrum of finite width $|k| \leq \hat{k}$ and varying propagation velocities $|c| \leq \hat{c}$. Let $\hat{\kappa} = \hat{k}h$, $\hat{\sigma} = \hat{c}\Delta t/h$, then the formulae for the normalized global error and the normalized global cost metric for a broadband signal can be defined as

$$\hat{e}(\hat{\kappa}, \hat{\sigma}) = \frac{1}{\hat{\sigma}\hat{\kappa}} \max_{(\kappa, \sigma) \in [0, \hat{\kappa}] \times [0, \hat{\sigma}]} |r(k, \sigma) - e^{-i\sigma\kappa}|, \quad (21)$$

$$\hat{c}_d(\hat{\kappa}, \hat{\sigma}) = pN_{\text{op}} \frac{1}{\hat{\sigma}\hat{\kappa}^{d+1}}. \quad (22)$$

D. Optimal performance for multi-scale problems

In this work, optimising the performance of a scheme is taken as minimising the computational cost for a given level of error, for given values of p and N_{op} and a given problem, *i.e.*, for specific non-dimensional groups ckT and kL values. To do this, a target level for the relative error ϵ is specified so that

$$\hat{e}(\hat{\kappa}, \hat{\sigma}) = \frac{\epsilon}{ckT} \equiv \tilde{\epsilon},$$

and the optimization consists of finding a pair of values $(\kappa^*(\tilde{\epsilon}), \sigma^*(\tilde{\epsilon}))$ that minimise the cost metric and satisfy the stability condition Eq. (15).

An illustration of this optimisation process is depicted in Figure 2(a), which shows iso-lines of the normalised global error $\hat{e}(\hat{\kappa}, \hat{\sigma})$ and normalised cost $\hat{c}_1(\hat{\kappa}, \hat{\sigma})$ in the wavenumber-Courant number $(\hat{\kappa}, \hat{\sigma})$ -plane, where the support is given by $[0 \leq \hat{\sigma} \leq \hat{\sigma}_{\text{max}}, 0 \leq \hat{\kappa} \leq \hat{\kappa}_{\text{max}}]$. Results are shown for the sixth-order compact spatial discretisation scheme coupled with the four-stage, fourth-order RK method, denoted by *C12RK4*, for which $\hat{\sigma}_{\text{max}} = 1.42$ and $\hat{\kappa}_{\text{max}} = \pi$. For the *C12RK4* scheme, the iso-cost curves are convex whereas the iso-error curves are concave.¹ Under such topology, the cost minima occur at the single point of tangency between the iso-lines and the circles in Figure 2(a) show the ‘optimal’ working condition for a number of relative error values.

Let

$$\hat{e}(\hat{\kappa}, \hat{\sigma}) \approx \max(\hat{e}_0(\hat{\kappa}), \hat{e}_t(\hat{\sigma})) \equiv \hat{e}_a(\hat{\kappa}, \hat{\sigma}), \quad (23)$$

be an approximation¹ of $\hat{e}(\hat{\kappa}, \hat{\sigma})$ from Eq. (21) where $\hat{e}_0(\hat{\kappa})$ is the spatial error (assuming no temporal error) and $\hat{e}_t(\hat{\sigma})$ is the temporal error (assuming no spatial error). $\hat{e}_0(\hat{\kappa})$ and $\hat{e}_t(\hat{\sigma})$ are given by

$$\hat{e}_0(\hat{\kappa}) \equiv \frac{1}{\hat{\kappa}} \max_{0 \leq \kappa \leq \hat{\kappa}} |\bar{\kappa} - \kappa|, \quad \hat{e}_t(\hat{\sigma}) \equiv \frac{1}{\hat{\sigma}} \max_{0 \leq z \leq \hat{\sigma}} \left| \sum_{j=0}^P (-iz)^j - e^{-iz} \right|, \quad (24)$$

where $\hat{z} = \hat{\sigma}\hat{\kappa}$. $\hat{e}_a(\cdot, \cdot)$ enables a decoupling of the error into spatial and temporal contributions, which will lead to a simplification of the computations that follow. Figure 2(b) shows a comparison of the optimal working points for the *C12RK4* scheme for the same levels of error as in Figure 2(a), when $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ is used

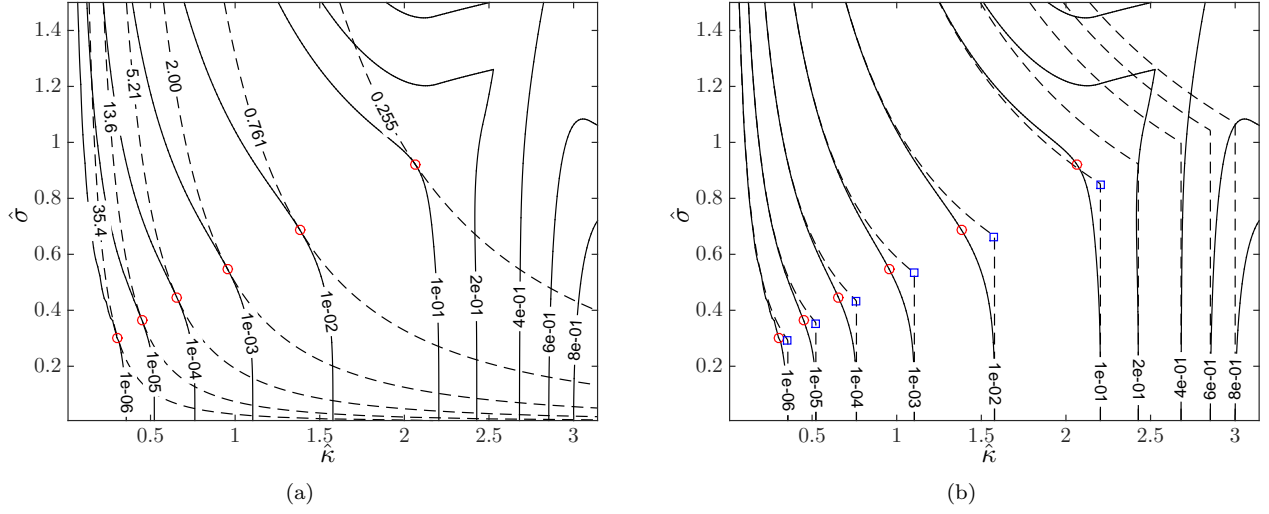


Figure 2. (a) Iso-lines of normalized ‘global’ error function $\hat{e}(\hat{\kappa}, \hat{\sigma})$ (solid lines) and normalized cost function $\hat{c}_1(\hat{\kappa}, \hat{\sigma})$ (dashed lines), for C12RK4 scheme; circles indicate ‘optimal’ working condition. (b) Iso-lines of normalized ‘global’ error function $\hat{e}(\hat{\kappa}, \hat{\sigma})$ (solid lines) and approximate normalized ‘global’ error function $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ (dashed lines); circles and squares indicate the respective ‘optimal’ working conditions.

in place of $\hat{e}(\kappa, \sigma)$. The dashed lines show the maximum of the spatial and temporal errors. The vertical segments are where the spatial error is dominant and the dashed curves are where the temporal error is dominant. The circles correspond to the error $\hat{e}(\hat{\kappa}, \hat{\sigma})$ estimated by Eq. (21) and the squares to $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$. Table 1 shows the optimal values $(\hat{\kappa}^+, \hat{\sigma}^+)$ and $(\hat{\kappa}^*, \hat{\sigma}^*)$, corresponding to $(\hat{\kappa}, \hat{\sigma})$ from Eq. (21) and Eq. (23), respectively, and the relative distance between them defined by

$$D((\hat{\kappa}^+, \hat{\sigma}^+), (\hat{\kappa}^*, \hat{\sigma}^*)) = \left[\left(\frac{\hat{\kappa}^+ - \hat{\kappa}^*}{\hat{\kappa}^+} \right)^2 + \left(\frac{\hat{\sigma}^+ - \hat{\sigma}^*}{\hat{\sigma}^+} \right)^2 \right]^{1/2}. \quad (25)$$

For the errors considered, there is a maximum of 17% difference in the location of the ‘optimal’ working conditions when using the two error definitions. Whereas such a discrepancy has an appreciable magnitude, it still enables operating a finite difference time-marching scheme in the ball-park of its optimal design point, as verified in a propagation of a monochromatic wave.¹⁹ This allows the splitting of the derivation of the cost-optimized coefficients into a spatial component and a temporal component.

ϵ	$\hat{e}(\hat{\kappa}, \hat{\sigma})$		$\hat{e}_a(\hat{\kappa}, \hat{\sigma})$		$D((\hat{\kappa}^+, \hat{\sigma}^+), (\hat{\kappa}^*, \hat{\sigma}^*))$
	$\hat{\kappa}^+$	$\hat{\sigma}^+$	$\hat{\kappa}^*$	$\hat{\sigma}^*$	
1×10^{-1}	2.064	0.922	2.203	0.849	0.10
1×10^{-2}	1.384	0.686	1.577	0.633	0.14
1×10^{-3}	0.954	0.548	1.103	0.534	0.16
1×10^{-4}	0.656	0.446	0.762	0.432	0.16
1×10^{-5}	0.449	0.365	0.522	0.354	0.17
1×10^{-6}	0.307	0.301	0.356	0.291	0.17

Table 1. Comparison between the optimal pair $(\hat{\kappa}^+, \hat{\sigma}^+)$ based on ‘true’ global error $\hat{e}(\hat{\kappa}, \hat{\sigma})$ and the optimal pair $(\hat{\kappa}^*, \hat{\sigma}^*)$ based on ‘approximate’ global error $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ for the classical C12RK4 scheme.

III. Optimized compact, time-marching schemes and extension to prefactored compact schemes

A. Optimized compact schemes

This section describes the optimisation of compact finite difference schemes, by which $\hat{\kappa}$ is maximised for a given normalised error level $\tilde{\epsilon}$, assuming that time integration is exact. This is equivalent to finding the maximum resolving efficiency $\hat{\kappa}_{\text{opt}}$ so that all wavenumbers $\hat{\kappa} \leq \hat{\kappa}_{\text{opt}}$ are resolved with an error less than $\tilde{\epsilon}$. For this purpose, the error $\hat{e}_0(\hat{\kappa})$ from Eq. (24) is used.

The baseline spatial scheme is the compact C1122 scheme of Eq. (5) with the free parameter α_1 . The new class of optimised schemes will be denoted C1122- n , where n represents the exponent in the target $\tilde{\epsilon} = 10^{-n}$. Optimisation in this setting is a non-trivial matter, as can be seen in Figure 3. In fact, the values $(\alpha_1^n, \hat{\kappa}_{\text{opt}}^n)$ at which $\hat{\kappa}$ is highest occur at a corner point of the 10^{-n} error contour. In addition, near this corner point, for each α_1 , there can be multiple values of $\hat{\kappa}$ which lie on the desired error contour. As such, a standard local minimisation routine will not necessarily locate the optimum. The procedure used for estimating $(\alpha_1^n, \hat{\kappa}_{\text{opt}}^n)$ is as follows. For a given α_1 the maximum $\hat{\kappa}$, which gives the prescribed level error, is found by utilising a number of local function solves starting at different initial guesses. A global pattern search algorithm, which does not require the objective function to be smooth, is then used in the α_1^n variable. The optimised values $(\alpha_1^n, \hat{\kappa}_{\text{opt}}^n)$, for $n = 4, 5$ and 6 are shown as the filled symbols in Figure 3 and are presented in Table 2. Table 2 shows that the optimised values for α_1 are above the ‘standard’ value of $1/3$ for all error levels, asymptoting toward $\alpha_1 = 1/3$ as the desired error level is reduced. This trend is indicated by the dotted line in Figure 3. In addition, Table 2 shows the maximum resolving efficiency $\hat{\kappa}^n$ for a normalised error level $\tilde{\epsilon} = 10^{-n}$, when the non-optimised C1122 scheme is used. Using the cost-optimized coefficients increases the spatial resolving efficiency $\hat{\kappa}$ from 47% up to 49% over the range $10^{-6} \leq \epsilon \leq 10^{-4}$, indicating a gain in spatial resolving efficiency from the spatial cost-optimization.

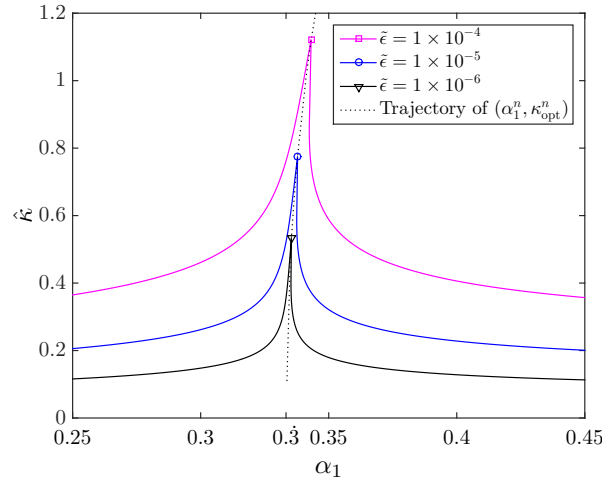


Figure 3. Iso-lines of ‘global’ spatial error $\hat{e}_0(\hat{\kappa})$ for the C1122 scheme from Eq. (5) with varying α_1 . The filled symbols indicate the optimised working conditions. The dotted line indicates the locus of the optimised working conditions that asymptotes to $\alpha_1 = 1/3$ as $\hat{\kappa} \rightarrow 0$.

	C1122- n		C1122		RK4- n				RK4			
n	α_1^n	$\hat{\kappa}_{\text{opt}}^n$	α_1	$\hat{\kappa}^n$	γ_3^n	γ_4^n	z_s^n	\hat{z}_{opt}^n	γ_3	γ_4	z_s	\hat{z}^n
4	0.354740	1.121	1/3	0.762	0.165242	0.0402486	2.826	0.436	1/6	1/24	2.83	0.331
5	0.337838	0.776	1/3	0.522	0.166106	0.0411119	2.828	0.272	1/6	1/24	2.83	0.186
6	0.335419	0.533	1/3	0.357	0.166486	0.0414859	2.829	0.160	1/6	1/24	2.83	0.105

Table 2. Coefficients and resolving efficiencies of optimised spatial C1122- n schemes and temporal RK4- n schemes.

B. Spatial differentiation by prefactored compact finite-difference schemes

In the C1122 scheme of Eq. (5) above, the approximation of the derivatives at the nodal points requires the solution of a tridiagonal system of equations. In this section, following the methodology of Hixon,² the tridiagonal system of equations is recast as two (upper and lower) bidiagonal systems. This gives a *prefactored* compact scheme. Let the derivative f'_i be split into a backward component $f_i'^B$ and a forward component $f_i'^F$ so that

$$f'_i = \frac{1}{2} (f_i'^F + f_i'^B). \quad (26)$$

Suppose that $f_i'^B$ and $f_i'^F$ are found by using, respectively, the following backward and forward finite difference formulae:

$$\alpha_0^F f'_i + \alpha_1^F f'_{i+1} = \frac{1}{h} [a_{i-1}^F f_{i-1} + a_i^F f_i + a_{i+1}^F f_{i+1}], \quad (27)$$

$$\alpha_{-1}^B f_i'^B + \alpha_0^B f'_i = \frac{1}{h} [a_{i-1}^B f_{i-1} + a_i^B f_i + a_{i+1}^B f_{i+1}]. \quad (28)$$

The coefficients (α^F, α^B) and (a^F, a^B) in Eqs. (27) and (28) are determined so that the real (dispersive) components of the scaled pseudo-wavenumbers of the forward and backward stencils are equal to the scaled pseudo-wavenumbers of the original central scheme, and the imaginary (dissipative) components of the scaled pseudo-wavenumbers are equal in magnitude and opposite in sign. Such schemes are termed *MacCormack* schemes²⁰. The scaled pseudo-wavenumber of the generic C1122 scheme of Eq. (5) is given by

$$\bar{\kappa}(\kappa) = \frac{\frac{2(\alpha_1+2)}{3(1+2\alpha_1)} \sin(\kappa) + \frac{(4\alpha_1-1)}{6(1+2\alpha_1)} \sin(2\kappa)}{\frac{1}{(1+2\alpha_1)} + \frac{2\alpha_1}{(1+2\alpha_1)} \cos(\kappa)}. \quad (29)$$

In a similar manner, the scaled pseudo-wavenumber of the generic forward and backward operators are

$$Re(\bar{\kappa}^F(\kappa)) = \frac{(a_1^F \alpha_0^F - a_0^F \alpha_1^F - a_{-1}^F \alpha_0^F) \sin(\kappa) - a_{-1}^F \alpha_1^F \sin(2\kappa)}{(\alpha_0^F)^2 + (\alpha_1^F)^2 + 2\alpha_1^F \alpha_0^F \cos(\kappa)}, \quad (30)$$

$$Im(\bar{\kappa}^F(\kappa)) = \frac{-(a_1^F \alpha_1^F + a_0^F \alpha_0^F) - (a_1^F \alpha_0^F + a_0^F \alpha_1^F + a_{-1}^F \alpha_0^F) \cos(\kappa) - a_{-1}^F \alpha_1^F \cos(2\kappa)}{(\alpha_0^F)^2 + (\alpha_1^F)^2 + 2\alpha_1^F \alpha_0^F \cos(\kappa)}, \quad (31)$$

and for the backward stencil:

$$Re(\bar{\kappa}^B(\kappa)) = \frac{(a_1^B \alpha_0^B + a_0^B \alpha_{-1}^B - a_{-1}^B \alpha_0^B) \sin(\kappa) + a_{-1}^B \alpha_{-1}^B \sin(2\kappa)}{(\alpha_{-1}^B)^2 + (\alpha_0^B)^2 + 2\alpha_{-1}^B \alpha_0^B \cos(\kappa)}, \quad (32)$$

$$Im(\bar{\kappa}^B(\kappa)) = \frac{-(a_0^B \alpha_0^B + a_{-1}^B \alpha_{-1}^B) - (a_1^B \alpha_0^B + a_0^B \alpha_{-1}^B + a_{-1}^B \alpha_0^B) \cos(\kappa) - a_{-1}^B \alpha_{-1}^B \cos(2\kappa)}{(\alpha_{-1}^B)^2 + (\alpha_0^B)^2 + 2\alpha_{-1}^B \alpha_0^B \cos(\kappa)}. \quad (33)$$

Imposing the following conditions

$$\alpha_j^F = \alpha_{-j}^B, \quad a_j^F = -a_{-j}^B, \quad -1 \leq j \leq 1,$$

ensures that the imaginary components of the forward and backward operators are equal in magnitude and opposite in sign and that the real components are equal to one another. Further, by imposing

$$\sum_{j=-1}^1 a_j^F = 0,$$

guarantees that the spatial derivatives vanish in regions of zero gradient. Finally, by matching corresponding terms between Eq. (30) and Eq. (29) the following system of algebraic equations is obtained

$$\left\{ \begin{array}{l} a_1^F \alpha_0^F - a_0^F \alpha_1^F - a_{-1}^F \alpha_0^F = \frac{2(\alpha_1+2)}{3(1+2\alpha_1)}, \\ -a_{-1}^F \alpha_1^F = \frac{4\alpha_1-1}{6(1+2\alpha_1)}, \\ (\alpha_0^F)^2 + (\alpha_1^F)^2 = \frac{1}{(1+2\alpha_1)}, \\ 2\alpha_1^F \alpha_0^F = \frac{2\alpha_1}{(1+2\alpha_1)}, \\ a_1^F + a_0^F + a_{-1}^F = 0. \end{array} \right. \quad (34)$$

The quadratic terms mean that Eq. (34) has two solutions. The lower value solution for α_F is selected to minimize the ratio α_1^F/α_0^F , so that the influence of errors at the boundaries on the interior scheme is minimized, in problems that require closure at the computational domain boundaries.

Substituting α_1^n from Section III.A, Table 2, for α_1 in Eq. (34), the prefactored coefficients for the C1122- n schemes are obtained. The resultant prefactored coefficients are listed in Table 3, in double precision (16 digits) format. It is noticeable that the standard C1122 scheme appears to exhibit a greater central dominance in the stencil, confirmed by the higher valued a_0^F coefficient. The coefficients in Table 3 together with Eqs. (26), (27) and (28) fully define the prefactored compact finite difference approximation f'_i of the first-order spatial derivatives $\partial f/\partial x$ in the f domain interior, $(i+1)h \leq x \leq (N-1)h$.

Figure 4 shows the dispersive characteristics of the classical C1122 and the cost-optimised C1122- n schemes. Figure 4(a) shows the real component of the scaled pseudo-wavenumber $Re(\bar{\kappa}(\kappa))$ of the compact scheme from Eq. (29) over the scaled wavenumber range $0 < \kappa \leq \pi$. The dashed line shows $\bar{\kappa}(\kappa) = \kappa$ for the case of exact differentiation, and the lines with symbols the prefactored compact finite difference approximations using different levels of target error 10^{-n} . The insert in Fig. 4(a) shows that the maximum value of the the real part of the scaled pseudo-wavenumber $Re(\bar{\kappa}(\kappa))$ does not vary significantly among the cost-optimized schemes. Due to the symmetric properties of the MacCormack schemes,²⁰ the real component of the scaled pseudo-wavenumber of the prefactored forward stencil from Eq. (30) is equal to the real component of the scaled pseudo-wavenumber of the prefactored backward stencil from Eq. (32) and is, by Eq. (26), also equal to the real part of the scaled pseudo-wavenumber of the original C1122 scheme of Eq. (29), that is $Re(\bar{\kappa}(\kappa)) = Re(\bar{\kappa}^F(\kappa)) = Re(\bar{\kappa}^B(\kappa))$. Hence, each line with symbols shows $Re(\bar{\kappa}(\kappa)) = Re(\bar{\kappa}^F(\kappa)) = Re(\bar{\kappa}^B(\kappa))$. Additionally, the two imaginary components of scaled pseudo-wavenumber of the prefactored forward and backward stencils, respectively from Eqs. (30) and (32), are equal in magnitude and opposite in sign and they cancel each other out by summing of Eq. (26), to give a dissipation-free scheme.

Figure 4(b) shows the spatial error \hat{e}_0 for the classical C1122 and cost-optimized C1122- n schemes against the discrete number of points per wavelength $N_\lambda = 2\pi/\kappa$. The classical scheme C1122 displays a straight line (constant logarithmic roll-off rate), whereas the cost-optimized schemes C1122- n feature local \hat{e}_0 minima of 10^{-n} for selected wavenumbers. At these wavenumbers the spatial error \hat{e}_0 is reduced approximately by one order of magnitude compared to the C1122 scheme. Similar $\hat{e}_0 - N_\lambda$ plots with single and multiple \hat{e}_0 minima are reported in the literature for optimized finite-difference schemes.^{5,11,13,21}

	C1122	C1122-4	C1122-5	C1122-6
α_1^F	0.276393202250021	0.284319733102085	0.280029458740805	0.278074448732061
α_0^F	0.723606797749979	0.715680266897915	0.719970541259195	0.721925551267939
a_1^F	0.87939886704167	0.870514243234187	0.875204935692797	0.877434323134828
a_0^F	-0.758797734083341	-0.741028486468374	-0.750409871385594	-0.754868646269656
a_{-1}^F	-0.12060113295833	-0.129485756765813	-0.124795064307203	-0.122565676865172

Table 3. Prefactored spatial discretization coefficients of the classical C1122 and optimized C1122- n schemes.

C. Cost-optimized Runge-Kutta time integration

The optimisation approach of Section III.A is herein extended to the Runge-Kutta time integration schemes of Eqs. (8–10). Assuming a zero spatial differentiation error, the coefficients γ_m are sought that maximise the temporal resolving efficiency $\hat{z}_{\text{opt}}(\tilde{\epsilon})$ for a given value of the normalised error $\tilde{\epsilon}$

$$\hat{z}_{\text{opt}} = \max\{\tilde{z} : \hat{e}_t(\tilde{z}, \gamma_m) \leq \tilde{\epsilon}\},$$

under the following stability constraint

$$\zeta \hat{z}_{\text{opt}}(\tilde{\epsilon}) \leq z_s, \quad (35)$$

where z_s is the stability limit defined in Eq. (15). The factor $\zeta \geq 1.0$ has been introduced to guarantee a greater stability margin beyond the range of well resolved angular frequencies z . The four-stage, fourth-order RK scheme is selected as the baseline temporal solver. To preserve formal second-order time integration accuracy, let $\gamma_1 = 1$ and $\gamma_2 = 1/2$. This leaves two free parameters γ_3, γ_4 to be determined from $\zeta \hat{z}_{\text{opt}}$.

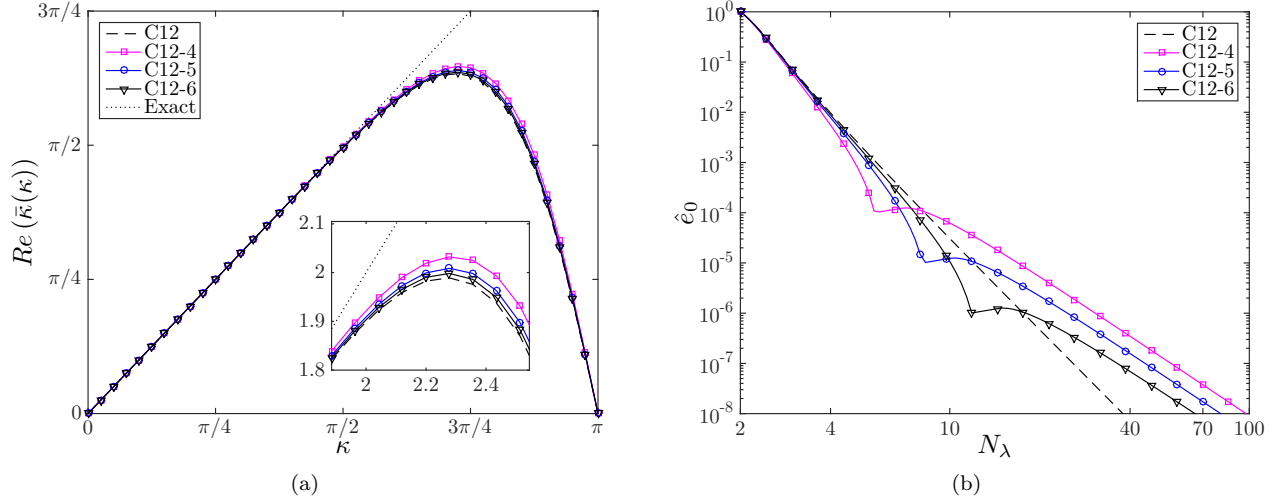


Figure 4. Dispersive characteristics of the prefactored classical C1122 and cost-optimized C1122–4, C1122–5 and C1122–6 schemes. (a) Real component of the prefactored compact scheme from Eq. (29). (b) Spatial error e_0 versus number of points per wavelengths N_λ .

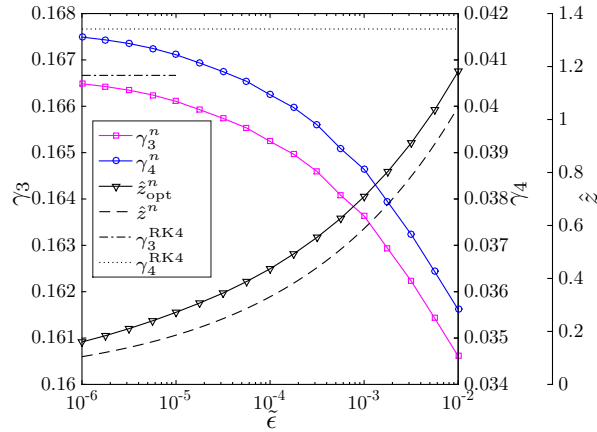


Figure 5. ‘Optimal’ working conditions of the generic RK4 schemes and corresponding temporal resolving efficiency z_{opt}^* for a range of errors $\tilde{\epsilon}$.

To find optimal values for differing levels of target error $\tilde{\epsilon} = 1 \times 10^{-n}$, a global pattern search algorithm similar to the one in Section III.A is used, in which an additional nonlinear constraint is incorporated to satisfy Eq. (35). In this case, $\zeta = 1.1$ was found heuristically to provide a good balance between performance and stability. The results are plotted in Figure 5, these show that the cost-optimised time integration schemes have smaller values of coefficients γ_3 and γ_4 than the classical fourth-order accurate RK scheme, indicated by the discontinuous lines. As the level of normalised error $\tilde{\epsilon}$ decreases, the cost-optimised coefficients γ_3^n , shown by the lines with symbols, tend asymptotically to γ_3 and γ_4 respectively.

The coefficients of the optimised RK4- n schemes, their respective temporal resolving efficiencies and the stability limits z_s of Eq. (15) are listed in the Table 2. The stability limits of the temporal cost-optimized schemes are marginally lower than the limit of the classical RK4 scheme, which is $z_s = 2.83$. Using the cost-optimised coefficients increases the temporal resolving efficiencies between 31% and 52% over the range $10^{-6} \leq \tilde{\epsilon} \leq 10^{-4}$, indicating a gain in temporal resolving efficiency from the temporal optimisation.

D. Theoretical performance of the cost-optimised schemes

Sections III.B and III.C considered the potential benefits of cost optimization in solving ordinary differential equations in space and in time, respectively. This section explores combining C1122- n and RK- n schemes

in a single numerical procedure, denoted C12RK4- n to solve the partial differential equation Eq. (12). The application of the finite difference/Runge-Kutta schemes to a signal composed of waves all with $\kappa < \hat{\kappa}$ can give rise to numerical round-off errors, which manifest as spurious waves with $\kappa > \hat{\kappa}$. In order to counter this, the Courant number, which represents the time step advancement in the RK4 scheme, is restricted to maintain numerical stability so that $\sigma < \sigma_{\max}$, where typically for the schemes considered $\sigma_{\max} \approx 1.5$.

Figure 2(a) has shown that, given a target level of normalised error $\hat{e}(\hat{\kappa}, \hat{\sigma})$, it is possible to determine the optimal working conditions $(\hat{\kappa}^*, \hat{\sigma}^*)$ that give the smallest normalised cost $\hat{c}_1(\hat{\kappa}, \hat{\sigma})$ and Figure 2(b) has shown that the normalised approximate error $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ can be used in place of $\hat{e}(\hat{\kappa}, \hat{\sigma})$ for estimating $(\hat{\kappa}^*, \hat{\sigma}^*)$. Following the same procedure, the normalised cost $c_1(\hat{\kappa}^*, \hat{\sigma}^*)$ of operating the C12RK4 and the C12RK4- n schemes for different levels of normalised approximation of $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ were determined and the results reported in Fig. 6(a). For each C12RK4- n scheme, the filled symbols denotes the normalised cost $\hat{c}_1(\hat{\kappa}^*, \hat{\sigma}^*)$ of operating this scheme at the normalised approximate error $\hat{e}_a(\hat{\kappa}^*, \hat{\sigma}^*) = 10^{-n}$, which is the designed error level of the C12RK4- n scheme. Figure 6(a) confirms that for each C12RK4- n scheme, the designed level of error 10^{-n} of its spatial differentiation C1122- n and temporal integration RK4- n constituents is retained once the spatial differentiation and temporal integration are coupled together in one partial differential equation solver.

Whereas it is encouraging to notice that the C12RK4- n schemes retain their design level based on the approximation $\hat{e}_a(\hat{\kappa}^*, \hat{\sigma}^*)$ at their design point, it is of interest to discuss the error trend away from the $\hat{e}_a = 10^{-n}$ level, in comparison with the conventional C12RK4 scheme. All C12RK4- n schemes exhibit a rate of convergence higher than the C12RK4 scheme up to their respective design error levels. Thereafter, the C12RK4- n schemes exhibit a \hat{e}_a plateau with increasing \hat{c}_1 , followed by a convergence rate that is lower than that of C12RK4. The \hat{e}_a trend is monotonically decreasing for all schemes. The variable error roll-off of the C12RK4- n schemes offers some advantages and limitations in computational physics applications. The monotonicity of the numerical error indicates that the scheme can be operated with confidence in wavenumber bandwidth limited problems knowing the design error will not be exceeded. The variable error roll-off rate, however, prevents the implementation of classical roll-off rate checks for mesh convergence computations as constant error roll-off rates exhibit only at a normalised cost \hat{c}_1 about two orders of magnitude higher than the scheme's optimal operating point. This cost difference is likely to amount to an onerous mesh resolution in a practical computation, the rate of convergence of which requires testing by alternative means. Note, the convergence rates are not linear in any case, but would just appears linear on a log-log plot.

Table 4 shows the computational cost c_1^* required by each of the optimised schemes to achieve their design level of error $\tilde{\epsilon}$. This is then compared with the corresponding cost $\hat{c}_{1(\text{C12RK4})}^*$ required by the standard C12RK4 scheme to obtain the same level of error; the percentage improvement ($\Delta\hat{c}_1$) is between 48% and 56% and increases with increasing error level. In addition, the smallest error $\tilde{\epsilon}_{(\text{C12RK4})}$ achieved by the standard scheme for the same costs c_1^* is shown, and the percentage improvement offered $\Delta\tilde{\epsilon}$ is seen to be in the region of 80% for all the optimised schemes.

Let space-only optimised schemes be denoted by C12RK4- S_n , where the standard RK4 scheme is used with the space optimised schemes of Section III.A, and time-only optimised schemes be denoted by C12RK4- T_n where the standard C1122 spatial discretisation is combined with the temporal optimised schemes of Section III.C; consistently n represents the exponent of the error level. Figure 6(b) shows a comparison between the fully-optimised C12RK4-5 scheme and C12RK4-S5 and C12RK4-T5. The filled symbols and dotted lines show the smallest normalised cost \hat{c}_1 required by each scheme to give a normalised error of $\hat{e}_a = 10^{-5}$. To achieve an error of 10^{-5} , the fully optimised scheme shows a cost advantage of 32.9% and 37.4%, over the C12RK4-T5 and C12RK4-S5 schemes, respectively. Similar results hold for other error levels n , but these are omitted from Figure 6(b) for clarity. These results highlight the advantages achieved when a scheme fully optimised in both space and time is used.

scheme	$\tilde{\epsilon}$	\hat{c}_1^*	$\hat{c}_{1(\text{C12RK4})}^*$	$\Delta\hat{c}_1(\%)$	$\tilde{\epsilon}_{(\text{C12RK4})}$	$\Delta\tilde{\epsilon}$
C12RK4-4	10^{-4}	73.65	142.67	48.45	4.98×10^{-4}	79.92
C12RK4-5	10^{-5}	170.63	370.19	53.91	6.48×10^{-5}	84.57
C12RK4-6	10^{-6}	422.36	963.55	56.17	7.28×10^{-6}	86.26

Table 4. Theoretical performance of cost-optimized schemes C12RK4- n for different target errors in one dimensional space and comparisons with standard C12RK4 scheme.

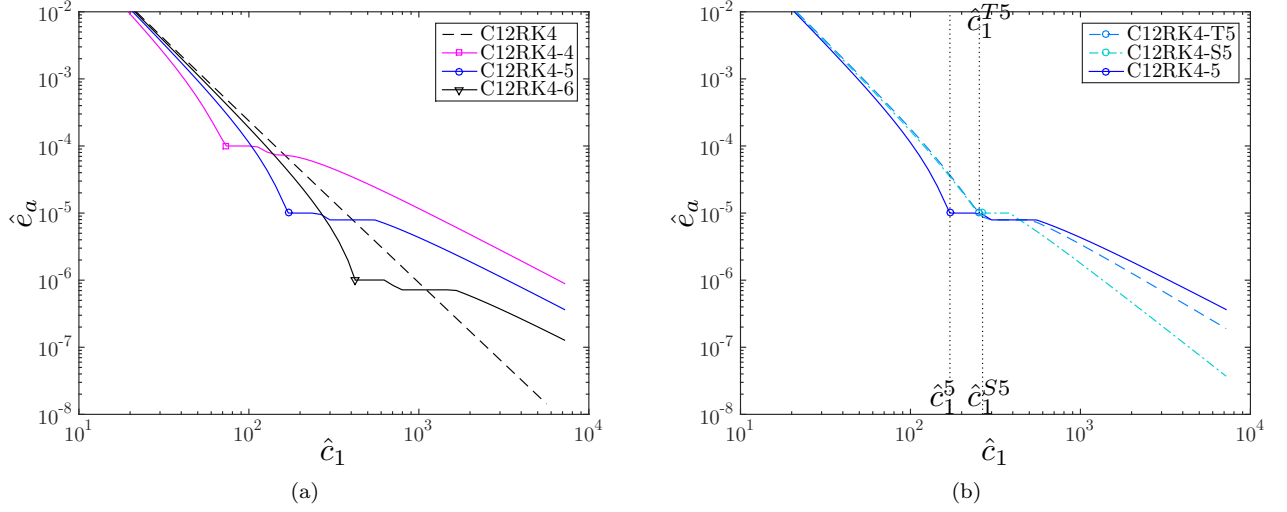


Figure 6. (a) Error estimated from Eq. (23) as a function of c_1 for C12RK4- n schemes, one-dimensional implementation. (b) Comparison of the estimated error among the C12RK4-5, C12RK4-T5 and C12RK4-S5 schemes.

IV. Numerical Experiments

This section presents a series of numerical examples to highlight the improvement in computational efficiency that is attainable by using cost-optimised prefactored schemes.

A. Polychromatic wave

The first example solves numerically the one-dimensional LAE equation Eq. (12), on the interval $(0, 1)$ with $c = 1$, periodic boundary conditions $u(0, t) = u(1, t)$, $t > 0$ and initial condition $u(x, 0) = \sum_{j=1}^4 \sin(2^{(j+1)}\pi x)$. The solution is a polychromatic wave with $\bar{k} = 32\pi$. For a given theoretical cost level \hat{c}_1 , a simulation is run with the optimal Courant numbers $\hat{\sigma}^*$ and optimal pseudo-wavenumber $\hat{\kappa}^*$ found in Section III.D. Once an approximate solution $u_h(x, t)$ has been found, the normalised relative numerical error \bar{e} can be calculated. Here, \bar{e} is given by

$$\bar{e}^2 = \frac{\sum_{i=1}^{N-1} (u_h(x_i, T) - u(x_i, T))^2}{\sum_{i=1}^{N-1} (u(x_i, 0))^2}, \quad (36)$$

where N is the number of nodes in the computational domain.

Figure 7(a) and Figure 7(b) show respectively the iso-lines of the theoretical error function $\hat{e}_a(\cdot, \cdot)$ and the iso-lines of the numerical error $\bar{e}(\cdot, \cdot)$ for the C12RK4-4 scheme. The locus of optimal points $(\hat{\sigma}^*, \hat{\kappa}^*)$ and the design point $(\hat{\kappa}_{\text{opt}}, \hat{\sigma}_{\text{opt}})$ for the C12RK4-4 scheme are overlaid on both figures. The *a posteriori* computed numerical error $\bar{e}(\hat{\kappa}, \hat{\sigma})$ is dependent on the wavenumber spectrum of the problem being solved. Still, the comparison between Figure 7(a) and Figure 7(b) shows that there is a strong correlation between the theoretical error $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ and the numerical error $\bar{e}(\hat{\kappa}, \hat{\sigma})$ for this polychromatic wave propagation problem. Figure 7(b) shows that, in the range $1 \leq \hat{\kappa} \leq 1.2$, there is a valley of error lower than the theory predicts. There are two branches to the locus of optimal points. The design point $(\hat{\kappa}_{\text{opt}}, \hat{\sigma}_{\text{opt}})$ lies on the right hand branch and is located close to the valley of low error. The discontinuity in the locus of optimal points is due to the non-monotonicity of the approximate spatial error function $\hat{e}_a(\cdot, \cdot)$ as \bar{k} increases from 0.75 to 1.2 for a fixed $\bar{\sigma} < 0.5$. This is shown in Figure 7(b) by the contour error labels over the range $0.75 < \hat{\kappa} < 1.2$. The non-monotonic error variation is lower in magnitude than the error change outside this region, as indicated by the packing of the contour lines in Figure 7(b), so that this region can be regarded as featuring an error plateau. In practice, by increasing the normalised cost \hat{c}_1 above that given by the design point, the theoretical error initially plateaus until the jump from the left to right branch. By increasing the cost above that given by the design point, the theoretical error initially plateaus, until the jump from the right to the left branch occurs. Thereafter, the error monotonically decays with the reduction of both $\hat{\kappa}$ and $\hat{\sigma}$ as shown in Figure 6(a).

Scheme	\bar{e}	\hat{c}_1^*	$\hat{c}_{1(\text{C12RK4})}^*$	$\Delta\hat{c}_1(\%)$	$\bar{e}_{(\text{C12RK4})}$	$\Delta e(\%)$
C12RK4-4	6.443×10^{-5}	74.376	170.56	56.39	4.706×10^{-4}	86.31
C12RK4-5	7.702×10^{-6}	171.7987	412.19	58.32	6.332×10^{-5}	87.84
C12RK4-6	9.565×10^{-7}	424.136	981.3	56.78	7.191×10^{-6}	86.71

Table 5. Example A. Performance of cost-optimised schemes at their operating points $(\hat{\kappa}^*, \hat{\sigma}^*)$ and comparison with the standard C12RK4 scheme.

Figure 8(a) plots the numerical error \bar{e} against cost \hat{c}_1 for the optimised schemes C12RK4- n and the standard C12RK4 scheme at a non-dimensional time $T = 1$ and overlays the theoretical results from Figure 8(a). The normalised numerical errors $\bar{e}(\hat{\kappa}^*, \hat{\sigma}^*)$ computed *a posteriori* compare very favourably with the approximate error estimates $\hat{e}_a(\hat{\kappa}^*, \hat{\sigma}^*)$. Each of the optimised schemes gives a significant improvement over the C12RK4 scheme close to their respective design level of error, although the maximum cost reduction is seen at error levels slightly below the design level of error. The approximate error $\hat{e}_a(\hat{\kappa}^*, \hat{\sigma}^*)$ in Figure 8(a) is shown to be consistently of the same order of magnitude of the *a posteriori* numerical error $\bar{e}(\hat{\kappa}^*, \hat{\sigma}^*)$ over the normalised cost range $10^1 \leq \hat{c}_1 \leq 10^4$. More specifically $\hat{e}_a \geq \bar{e}$. This is a desirable feature, as it indicates that \hat{e}_a can be used as an upper bound indicator for the normalised error for the application of a cost optimized prefactored compact scheme to linear advection problems. It implies that the numerical solution will not exceed an acceptable level of numerical error $\bar{e} = 10^{-n}$ that is arbitrarily set by the numerical modeller selecting a specific C12RK4- n scheme. The *a posteriori* numerical error trends shown in Figure 8(a) can be explained by considering the iso-error maps in the $(\hat{\kappa}, \hat{\sigma})$ plane. For instance, in the case of the C12RK4-4 scheme, the initial rapid decrease in error corresponds to moving along the right hand branch of the locus of optimal points in Figure 7(b), until the design point is reached. The *a posteriori* normalized numerical error would be then expected to plateau with increasing cost, based on the trend from its approximation \hat{e}_a , however, the numerical error continues to decrease, due to the locus of optimal points $(\hat{\kappa}^*, \hat{\sigma}^*)$ just beginning to enter the valley of low numerical error. As the locus of points ‘climbs’ back out of the valley, the *a posteriori* normalised error \bar{e} increases with cost until $\hat{e}_a \approx \bar{e}$ then begins to plateau until the jump from the left to the right branch, where the *a posteriori* normalised numerical error converges at increasing \hat{c}_1 and decreasing $\hat{\kappa}$. A similar discussion applies to both the C12RK4-5 and C12RK4-6 curves in Figure 6(a). Table 5 reports the performance of the C12RK4- n schemes benchmarked against the C12RK4 scheme at their respective design points. For each optimised scheme operating at their design point, the numerical error \bar{e} is close to the corresponding design level of error; in fact, the error is slightly below the design level. For example, the C12RK4-4 scheme has an error of 6.443×10^{-5} at its design point, compared with the theoretical error 1×10^{-4} .

Figure 8(b) shows the normalized numerical error \bar{e} against computational time, at the same $(\hat{\kappa}^*, \hat{\sigma}^*)$ as in Figure 8(a). Computations were carried out on a machine with an Intel Xeon Ivy Bridge CPU running at 2.6GHz and the computational time was averaged over 100 runs per $(\bar{\kappa}^*, \bar{\sigma}^*)$ pair to ensure the consistency of the results. Qualitatively, Figure 8(b) is very similar to Figure 8(a). This indicates that \hat{c}_1 is a good pseudo-variable for the CPU time; this is confirmed in Figure 8(c), where the CPU time, t , is plotted against \hat{c}_1 for the C12RK4-5 scheme and shows $\hat{c}_1 \propto t$. The computations were repeated for the C12RK4-4 and C12RK4-5 schemes and the constant of proportionality was found to be the same for all three schemes.

Table 5 that the percentage cost reduction $\Delta\hat{c}_1$ over the standard C12RK4 scheme in achieving the same numerical error level \bar{e} attained by the optimised schemes operated at their design points. This cost reduction varies from 50% to 60% and it increases with the scheme number n . Similarly, the percentage error reduction Δe achieved by the optimised schemes over the C12RK4 scheme for the same cost \hat{c}_1^* is over 80% for each of the optimised schemes. Using the $t \propto \hat{c}_1$ relationship from Figure 8(c), each of the optimised schemes offers the same (percentage) reduction in CPU time over the C12RK4 scheme as the $\Delta\hat{c}(\%)$ reported in Table 5.

B. Gaussian pulse

The second sample application of the C12RK4- n schemes is derived from Test case B from.¹⁵ A numerical solution is sought for the one-dimensional LAE equation Eq. (12), over the domain $(-100, 100)$, with $c = 1$, periodic boundary conditions, and the initial condition $u(x, 0) = \frac{1}{2}e^{-(x/3)^2}$. The solution has a broadband Fourier decomposition and provides a good test of the C12RK- n schemes applied to a broadband signal. The

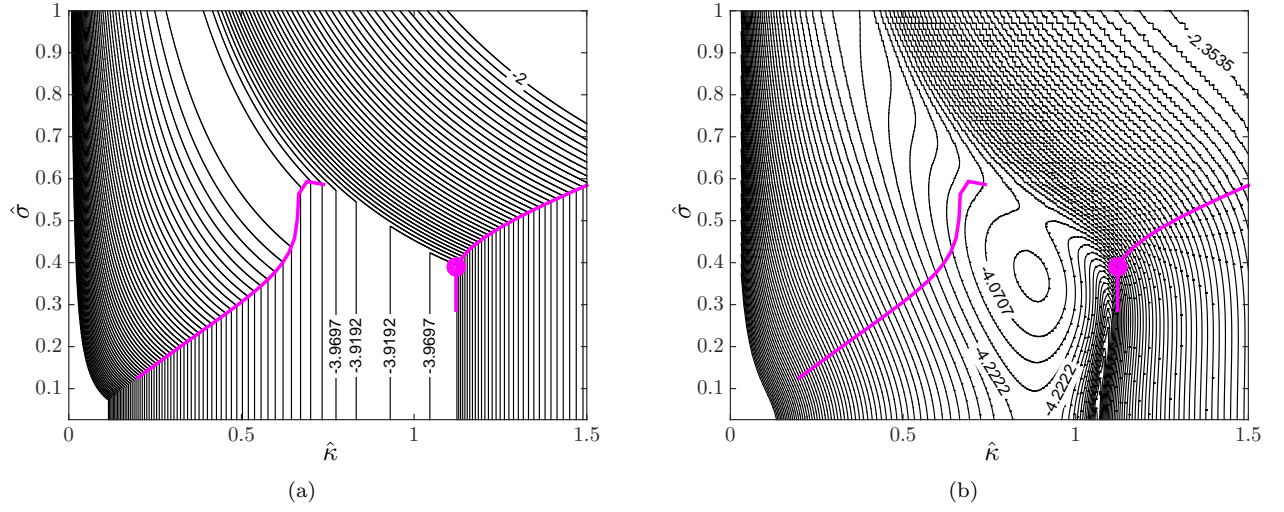


Figure 7. (a) Iso-lines of global approximate error $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ for the C12RK4-4 scheme and (b) iso-lines of polychromatic wave numerical error. The iso-lines are equally spaced on a log-scale and the levels shown are the exponent of 10. Bold (magenta) lines indicate the locus of optimal (κ^*, σ^*) for the C12RK4-4 scheme, while the filled circle corresponds to $(\hat{\kappa}_{\text{opt}}, \hat{\sigma}_{\text{opt}})$.

characteristic length scale for this problem is $\lambda = 6$, which corresponds to $\bar{k} = \pi/3$. For a given cost level \hat{c}_1 , the optimal Courant number $\hat{\sigma}^*$ and optimal pseudo-wavenumber $\hat{\kappa}^*$ have been chosen based on those computed in Section III.D. Figure 9(a) plots the normalised numerical error \bar{e} , computed using Eq. (36), against cost \hat{c}_1 for each of the optimised schemes C12RK4- n and the standard C12RK4 scheme at the final time $T = 100$. In Figure 9(a), all optimal operating points of the C12RK4- n schemes, denoted by the filled symbols, lie to the left of the continuous line, therefore, numerical solutions with the same level of normalised numerical error \bar{e} can be obtained for a lower computational cost than when using the C12RK4 scheme. In this case the reduction in computational cost \hat{c}_1 obtained by using the optimised schemes is not as large as in the polychromatic wave test case, but there is still a clear advantage in using the optimised schemes at their respective design points. This is better appreciated in the enlargement of Figure 9(b) in the vicinity of $\bar{e} = 1 \times 10^{-6}$ and $\bar{e} = 1 \times 10^{-5}$. Table 6 shows a comparison between the optimised schemes C12RK4- n , $n = 4, 5, 6$, operating at their design points and the standard C12RK4 scheme. Table 6 confirms that, when operating at their design points, all optimised schemes require a 15% to 25% lower cost than the C12RK4 scheme to achieve the same normalised numerical error \bar{e} . Similarly, when the optimised schemes are run at their design point, they offer between a 30% and 50% reduction in normalized numerical error for the same normalized cost \hat{c}_1^* when compared with C12RK4 scheme.

Notwithstanding the positive outcome from this broadband wave propagation test, which points to a consistent performance advantage of the C12RK4- n schemes over the C12RK4 scheme, this performance advantage is significantly lower than that predicted by the theory in Table 4. The background of this lower performance gain is explored by reporting in Figure 9(a) by the dotted lines the approximate normalised error $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ versus the normalised cost \hat{c}_1 for the C12RK4- n and C12RK4 schemes. The approximate normalised error in an *a priori* estimate that assumes equal contributions from all wavenumbers below $\bar{k} = \pi/3$. Recall that Figure 4(a) indicates an increase in dispersion with increasing wavenumber. Given that the Gaussian pulse has a definite roll-off at high wavenumbers, $\hat{e}_a(\hat{\kappa}, \hat{\sigma})$ overestimates \bar{e} , as shown by the dotted curves lying above the continuous and discontinuous line in Figure 9(a). The approximate errors can be considered to be upper bound estimates for \bar{e} because \hat{e}_a assumes the equal weighting of wavenumbers over the whole wavenumber spectrum, which is not necessarily the case in a physical problem. Given that the baseline *a posteriori* normalised numerical error is lower than \hat{e}_a , a more modest reduction in $\bar{e} = 10^{-n}$ appears to be a reasonable outcome by the use of a C12RK4- n scheme.

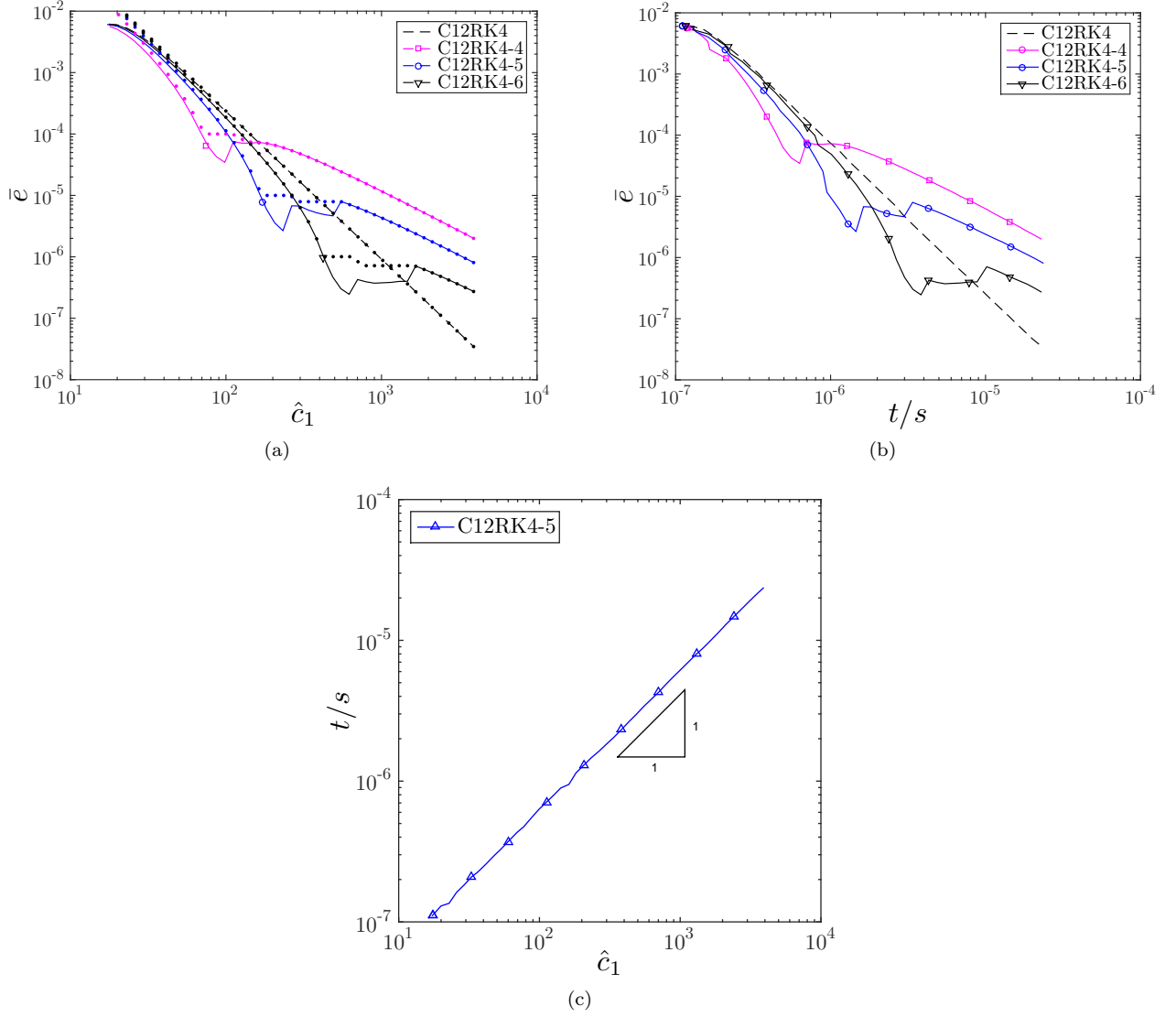


Figure 8. Example A: (a) Numerical error as a function of cost function \hat{c}_1 (b) Numerical error as a function of CPU time for the classical C12RK4 and optimised C12RK4- n schemes. (c) CPU time compared with theoretical cost.

V. Conclusion

In this paper, cost-optimized prefactored schemes C12RK4- n have been developed based on both spatial and temporal *a priori* error estimates. Numerical experiments have verified that, for representative test cases, these cost-optimized schemes show significant reductions in theoretical cost compared to the non-optimized scheme when an *a posteriori* numerical error is computed. In addition, the numerical experiments verified that the theoretical cost used to derive the cost-optimized schemes is a good representative for CPU time. For a problem with a polychromatic wave solution, each of the cost-optimized C12RK4- n schemes achieve upwards of a 50% reduction in cost to achieve the same error level as the classical non-optimized C12RK4 scheme, when operated at their design points. For the Gaussian pulse test case, where the solution spans the full wavenumber spectrum, a cost reduction of between 15% and 25% is obtained. Such levels of improvement in performance are attractive for both computationally expensive aeroacoustic applications and for those applications requiring a high resolution both in frequency and in wavenumber space. Further studies are required to investigate how the computational advantage of the cost-optimized schemes is affected by the use of different computational boundary conditions and by the application of these schemes to model multi-dimensional and non-linear wave propagation phenomena.

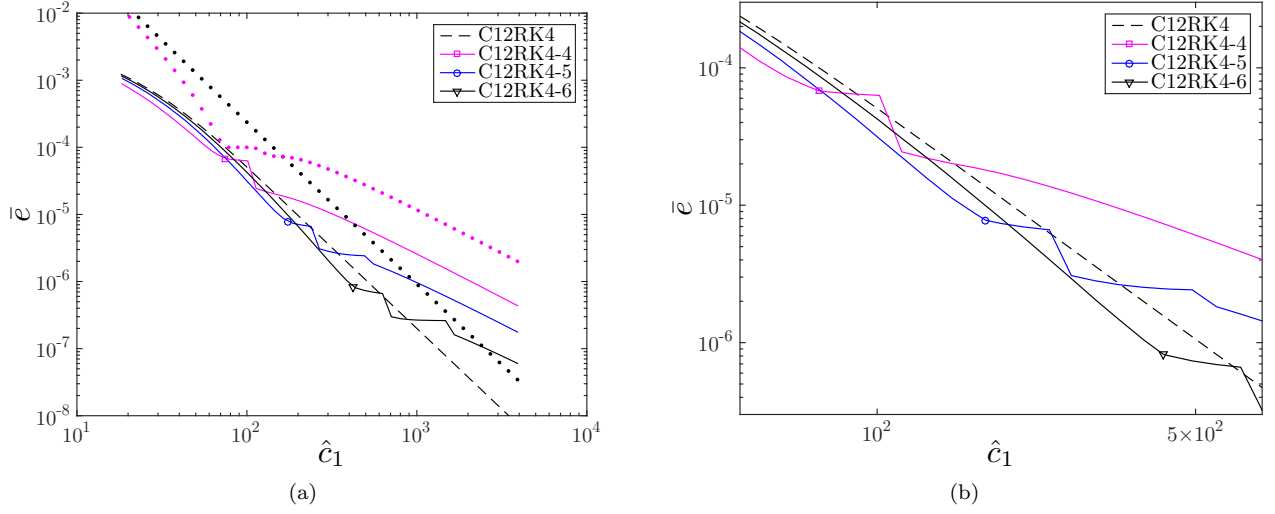


Figure 9. Example B: (a) Numerical error as a function of one-dimensional cost function c_1 with theoretical results overlaid (dotted lines) and (b) zoom around the C12RK4-4, C12RK4-5 and C12RK4-6 design points.

Scheme	\bar{e}	\hat{c}_1^*	$\hat{c}_1^*_{((C12RK4))}$	$\Delta\hat{c}_1(\%)$	$\bar{e}_{(C12RK4)}$	$\Delta e(\%)$
C12RK4-4	6.784×10^{-5}	74.677	88.29	15.42	9.954×10^{-5}	31.85
C12RK4-5	7.7647×10^{-6}	173.227	218.78	20.82	1.365×10^{-5}	43.12
C12RK4-6	8.2525×10^{-7}	424.245	554.45	23.48	1.573×10^{-6}	47.54

Table 6. Example B: Performance of cost-optimised schemes at their operating points (κ^*, σ^*) and comparison with standard C12RK4 scheme.

References

- ¹Pirozzoli, S., “Performance analysis and optimization of finite-difference schemes for wave propagation problems,” *Journal of Computational Physics*, Vol. 222, 2007, pp. 809–831.
- ²Hixon, R., “Prefactored small-stencil compact schemes,” *Journal of Computational Physics*, Vol. 165, No. 2, 2000, pp. 522–41.
- ³UK, G., “The Future of Air Transport, White Paper and the Civil Aviation Bill,” Tech. rep., UK Department of Transport, 2003.
- ⁴ACARE, “The Strategic Research Agenda SRA-1, SRA-2 and the 2008 Addendum to the Strategic Research Agenda,” Tech. rep., <http://www.acare4europe.org/>, 2009.
- ⁵Rock, W. D., Desmet, W., Baelmans, M., and Sas, P., editors, *An overview of high-order finite difference schemes for computational aeroacoustics*, Leuven, Belgium, 2004, Proceedings of the 2004 International Conference on Noise and Vibration Engineering, ISMA.
- ⁶Kurbatskii, K. A. and Mankbadi, R. R., “Review of computational aeroacoustics algorithms,” *International Journal of Computational Fluid Dynamics*, Vol. 18, No. 6, 2004, pp. 533–546.
- ⁷Colonius, T. and Lele, S., “Computational aeroacoustics: progress on nonlinear problems of sound generation,” *Progress in Aerospace Sciences*, Vol. 40, 2004, pp. 365–416.
- ⁸Hardin, J., Ristorcelli, J. R., and Tam, C., editors, *ICASE/LaRC Workshop on Benchmark Problems in Computational Aeroacoustics (CAA)*, Hampton, Virginia, October 1995, NASA Conference Publication 3300.
- ⁹Hardin, J. and Tam, C., editors, *Second Computational Aeroacoustics (CAA) Workshop on Benchmark Problems*, Langley Research Center, November 1996, NASA Conference Publication, Hampton, Virginia.
- ¹⁰Hardin, J., Huff, D., and Tam, C., editors, *Third Computational Aeroacoustics (CAA) Workshop on Benchmark Problems*, Langley Research Center, August 2000, NASA Conference Publication 2000-209790, Cleveland, Ohio.
- ¹¹Lele, S. K., “Compact finite difference schemes with spectral-like resolution,” *Journal of Computational Physics*, Vol. 103, No. 1, 1992, pp. 16–42.
- ¹²Hixon, R., “On Increasing the accuracy of MacCormack schemes for Aeroacoustic Applications,” Nasa Contract Report 202311 ICOMP-96-11, NASA, December 1996.
- ¹³Ashcroft, G. and Zhang, X., “Optimized prefactored compact schemes,” *Journal of Computational Physics*, Vol. 190, No. 2, 2003, pp. 459–477.
- ¹⁴Spisso, I. and Rona, A., “A selective overview of high-order finite difference schemes for aeroacoustic applications,” *International Conference on Sound and Vibration*^{14th}, Cairns, Australia, 2007.

- ¹⁵Bernardini, M. and Pirozzoli, S., “A general strategy for the optimization of Runge-Kutta schemes for wave propagation phenomena,” *Journal of Computational Physics*, Vol. 228, 2009, pp. 4182–4199.
- ¹⁶Bogey, C. and Bailly, C., “A family of low dispersive and low dissipative explicit schemes for flow and noise computations,” *Journal of Computational Physics*, Vol. 194, No. 1, 2004, pp. 194–214.
- ¹⁷Hu, F. Q., Hussaini, M. Y., and Manthey, J. L., “Low-dissipation and low-dispersion Runge-Kutta schemes for computational acoustics,” *Journal of Computational Physics*, Vol. 124, No. 1, 1996, pp. 177–191.
- ¹⁸Hirsch, C., *Numerical Computation of Internal and External flow*, Vol. 1, Fundamental of Computational Fluid Dynamics, New York, 2nd ed., 2007.
- ¹⁹Spisso, I., *Development of a Prefactored High-Order Compact Scheme for Low-Speed Aeroacoustics*, Ph.D. thesis, University of Leicester, optional, optional 2013, optional.
- ²⁰Hixon, R. and Turkel, E., “Compact Implicit MacCormack-Type Schemes with High Accuracy,” *Journal of Computational Physics*, Vol. 158, No. 1, 2000, pp. 51–70.
- ²¹Tam, C. K. W. and Web, J. C., “Dispersion-Relation-Preserving finite difference schemes for Computational Acoustics,” *Journal of Computational Physics*, Vol. 107, 1993, pp. 262–281.