

Contents lists available at ScienceDirect

Journal of Computational Physics

journal homepage: www.elsevier.com/locate/jcp



Schur complement IMplicit-EXplicit formulations for discontinuous Galerkin non-hydrostatic atmospheric models



Sohail Reddy ^{a,*,1}, Maciej Waruszewski ^b, Felipe A.V. de Braganca Alves ^a, Francis X. Giraldo ^{a,*}

- ^a Department of Applied Mathematics, Naval Postgraduate School, Monterey, CA, United States of America
- ^b Sandia National Laboratories, Albuquerque, NM, United States of America

ARTICLE INFO

Article history: Received 3 December 2022 Received in revised form 2 June 2023 Accepted 7 July 2023 Available online 22 July 2023

Keywords:
Semi-implicit time integration
IMplicit-EXplicit integration
Discontinuous Galerkin method
Atmospheric modeling
Euler equations
Navier-Stokes equations

ABSTRACT

This work presents IMplicit-EXplicit (IMEX) formulations for discontinuous Galerkin (DG) discretizations of the compressible Euler equations governing non-hydrostatic atmospheric flows. In particular, we show two different IMEX formulations that not only treat the stiffness due to the governing dynamics but also the domain discretization. We present these formulations for two different equation sets typically employed in atmospheric modeling. For both equation sets, efficient Schur complements are derived and the challenges and remedies for deriving them are discussed. The performance of these IMEX formulations of different orders are investigated on both 2D (box) and 3D (sphere) test problems and shown to achieve their theoretical rates of convergence and their efficiency with respect to both mesoscale and global applications are presented.

Published by Elsevier Inc.

1. Introduction

The discontinuous Galerkin (DG) method has been used successfully for developing compressible Euler and Navier-Stokes models. However, one of the main drawbacks of the method is that it is not competitive with either low-order or continuous (CG) methods due to the smaller time-step required for stability. It is well-known that high-order methods impose a more stringent Courant-Friedrichs-Lewy (CFL) condition compared to low-order methods (see, e.g., [1]). It is less clearly understood that most DG methods in the literature impose a more stringent CFL condition than CG methods; the reason is twofold: (1) upwind-based numerical fluxes move the eigenvalues from the imaginary to the real axis (since they add dissipation), requiring the time-integration method to be stable for such eigenspectra (e.g., Runge-Kutta methods satisfy this condition but neutrally stable methods do not); (2) most Godunov methods in use today (e.g., FV and DG methods) rely on one-dimensional numerical fluxes - this means that if the flow moves diagonally across a gridpoint then it must cross two element edges. If the numerical method is not constructed in such a manner to handle this situation (as in the evolutionary DG method described in [2]) then smaller time-steps are required. However, our experience has shown that using standard (one-dimensional fluxes that require a smaller time-step) give faster time-to-solution than more complicated fluxes such as those used in [2].

^{*} Corresponding authors.

E-mail addresses: sredd001@fiu.edu, reddy6@llnl.gov (S. Reddy), mwarusz@sandia.gov (M. Waruszewski), felipe.alves.br@nps.edu (F.A.V. de Braganca Alves), fxgirald@nps.edu (F.X. Giraldo).

¹ Also with Center for Applied Scientific Computing (CASC) at Lawrence Livermore National Laboratory, Livermore, CA, USA.

In order to circumvent the stringent CFL condition imposed by explicit methods, researchers have explored implicitexplicit (IMEX) methods. The first two such works found in the literature are by [3] and [4] where the authors begin with standard double Butcher tableaux whereby the stiff portion is handled implicitly while the non-stiff part is handled explicitly. In both works they begin with nonlinear handling of the stiff implicit part but in [3], the authors linearize the problem relying on a specific numerical flux while those in [4] maintain the nonlinear structure and use Newton's method to linearize the problem and make use of the Jacobian-free Newton-Krylov method [5]. The work that we describe here continues the work first proposed in [6,7] where we rely on the linearization of the implicit nonlinear problem using a reference solution, typical of semi-implicit methods [8,9] used in the atmospheric modeling community. This approach has been recently analyzed in the so-called RS-IMEX method (the RS denotes the reference solution linearization) in [10-14]. Our goal is to complete the work that we began in [7] where we sought a Schur complement to the construction of the DG IMEX method but were only able to do so for a specific class of boundary conditions (namely, periodic in all directions). In what follows, we generalize this approach to reflection boundary conditions, although other types of boundary conditions are also possible. Furthermore, we show that the eigenvalues of the resulting Helmholtz-like operator resulting from the Schur complement are all real and positive, indicating that the resulting operator is symmetric positive-definite, similar to its CG counterpart (see, e.g., [6,15]); to our knowledge this is a new result and one which has repercussions on the class of preconditioners that are now applicable to this formulation. We then construct Schur complements for the full 3D-IMEX problem (when the grid is isotropic and the stiffness is due to all directions) as well as for the 1D-IMEX problem (when the stiffness results from one direction, as occurs along the vertical direction in geophysical fluid dynamics applications of atmospheric models).

The remainder of the paper is organized as follows: in Sec. 2 we discuss the two equation sets used which, using the naming convention from [6,16], we refer to as sets 2C and 3C. In Sec. 3 we describe the implicit-explicit time-integration strategy which is the main point of this paper, and describe both the No-Schur and Schur forms, where the former is the approach used by all previous works on IMEX DG formulations and the latter is the main contribution of this paper. Section 4 discusses the discontinuous Galerkin spatial discretization and the role that its numerical fluxes play in the construction of the No-Schur and Schur forms of the IMEX formulation. In Secs. 5 and 6 we compare the accuracy, stability, and efficiency of both the No-Schur and Schur forms for four test cases: the first two tests are standard test cases used to verify nonhydrostatic atmospheric models at the Large-Eddy-Simulation (LES) scale ($\mathcal{O}(100 \text{ m})$ grid resolutions), the third is used to verify models at the cloud-resolving scale ($\mathcal{O}(1 \text{ km})$ grid resolutions) and the fourth has been used previously for verification at global scales ($\mathcal{O}(100 \text{ km})$ grid resolutions). We summarize the paper in Sec. 7 and offer directions for future work.

2. Governing equations

This work presents semi-implicit formulations² for two different equation sets often used in non-hydrostatic atmospheric modeling. We limit our discussions and formulations to equations sets in conservation form. The nomenclature that we use below was first defined in [16] where various forms of the Euler equations were defined, discretized, and compared. We only choose to discuss the equations in conservation form since these are the only practical options for using the discontinuous Galerkin method. We consider balance laws written in the form

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \mathbf{F} = \mathbf{S} \tag{2.1}$$

where \mathbf{q} , \mathbf{F} , and \mathbf{S} are vectors of prognostic, conserved variables and their total fluxes and sources, respectively. These vectors are defined in the subsequent sections.

2.1. Equation Set2C

This equation set is used in its conservation form in [17,18] and in non-conservation form in [15]. This equation set features the standard formulations for the continuity and momentum equations but utilizes the thermodynamic equation in terms of potential temperature. This can be advantageous in atmospheric modeling as this quantity is unaffected by the rising and sinking motion of the fluid over obstacles (such as mountains) or in large-scale atmospheric turbulence; the potential temperature is also a useful measure of the static stability of the atmosphere [19, Ch. 8]. One drawback of this equation set is that its derivation (from the first law of thermodynamics) assumes a constant composition atmosphere (i.e., gas constants remain constant); therefore, this equation set is not a good choice for atmospheric simulations reaching into the thermosphere. Another drawback (shown in [16]) is that the equation for computing pressure is expensive due to the required exponentials. This equation set is written as Eq. (2.1) with

² We use the term <u>semi-implicit</u> to the refer to the construction of IMEX methods whereby we construct an explicit solution followed by an implicit correction.

$$\mathbf{q}_{2C} = \begin{pmatrix} \rho \\ \mathbf{U} \\ \Theta \end{pmatrix}, \quad \mathbf{F}_{2C} = \begin{pmatrix} \mathbf{U} \\ \mathbf{U} \otimes \mathbf{U} \\ \rho \\ \left(\frac{\Theta}{\rho} \right) \mathbf{U} \end{pmatrix}, \quad \mathbf{S}_{2C} = \begin{pmatrix} 0 \\ -(\rho \nabla \phi + 2\Omega \times \mathbf{U}) \\ 0 \end{pmatrix}, \tag{2.2}$$

where the conserved, prognostic variables are $(\rho, \mathbf{U}, \Theta)$, $\mathbf{U} = \rho \mathbf{u} = (\rho u, \rho v, \rho w)^T$, $\Theta = \rho \theta$ is the potential temperature density, $\phi = gr$ is the geopotential height with g representing the gravitational constant and r representing the height coordinate (in flow-in-a-box, it is the Cartesian z-component and on the sphere it represents the radial coordinate), \mathbf{I} is a rank-3 identity tensor, Θ is the tensor-product operator, and \mathcal{T} is the transpose operator. The pressure P is given as

$$P = P_A \left(\frac{R\Theta}{P_A}\right)^{\gamma} \tag{2.3}$$

where P_A is the reference pressure at the surface $(P_A=1\times 10^5 \text{ Pa}),\ R=c_p-c_v$ is the specific gas constant, c_p and c_v are the specific heats at constant pressure and volume, respectively, and $\gamma=c_p/c_v$ is the ratio of specific heats. Introducing splitting about some reference field for density, $\rho(\mathbf{x},t)=\rho_0(r)+\rho'(\mathbf{x},t)$, pressure, $P(\mathbf{x},t)=P_0(r)+P'(\mathbf{x},t)$, and potential temperature, $\Theta(\mathbf{x},t)=\Theta_0(r)+\Theta'(\mathbf{x},t)$ where the reference values (ρ_0,P_0,Θ_0) are in hydrostatic balance $(\nabla P_0=-\rho_0\nabla\phi)$, Eq. (2.2) can be written in perturbation form as Eq. (2.1) with $\mathbf{q}:=\mathbf{q}'_{2C}$, $\mathbf{F}:=\mathbf{F}'_{2C}$ and, $\mathbf{S}:=\mathbf{S}'_{2C}$ defined as

$$\mathbf{q}_{2C}' = \begin{pmatrix} \rho' \\ \mathbf{U} \\ \Theta' \end{pmatrix}, \quad \mathbf{F}_{2C}' = \begin{pmatrix} \mathbf{U} \\ \frac{\mathbf{U} \otimes \mathbf{U}}{\rho} + P' \mathbf{I} \\ \left(\frac{\Theta}{\rho}\right) \mathbf{U} \end{pmatrix}, \quad \mathbf{S}_{2C}' = \begin{pmatrix} 0 \\ -\left(\rho' \nabla \phi + 2\Omega \times \mathbf{U}\right) \\ 0 \end{pmatrix}. \tag{2.4}$$

2.2. Equation Set3C

This equation set is the most common form of the Euler equations in CFD [20] and has been used in two recent dynamical cores related to the CliMa model [21,22]. This equation set shares the continuity and momentum equations with Set2C but differs in the thermodynamic equation and the equation of state. One of the advantages of this equation set is that it allows for conservation of not only mass and momentum, but also total energy. Moreover, since it is the equation of choice in CFD, it has been used in many recent works that rigorously prove the stability of this equation set using DG methods even at high-order (see [22,23]). Yet another advantage of this equation set is that its derivation does not assume anything about the composition of the atmosphere and therefore applicable to high-altitude simulations. A disadvantage of this equation set, at least in the context of atmospheric modeling, is that much of the physical parameterizations require either potential temperature or temperature as inputs and so one would need to extract this variable from total energy. One other likely disadvantage (based on our experience) is that this equation set requires more care to stabilize since it is completely conservative (i.e., non-dissipative) and so the numerical method has to be designed carefully to maintain stability. Set3C is written as Eq. (2.1) with

$$\mathbf{q}_{3C} = \begin{pmatrix} \rho \\ \mathbf{U} \\ E \end{pmatrix}, \quad \mathbf{F}_{3C} = \begin{pmatrix} \mathbf{U} \\ \frac{\mathbf{U} \otimes \mathbf{U}}{\rho} + P\mathbf{I} \\ \frac{(E+P)}{\rho} \mathbf{U} \end{pmatrix}, \quad \mathbf{S}_{3C} = \begin{pmatrix} 0 \\ -(\rho \nabla \phi + 2\Omega \times \mathbf{U}) \\ 0 \end{pmatrix}$$
(2.5)

where the conserved, prognostic variables are (ρ, \mathbf{U}, E) , E is the total energy given by $E = \rho c_V T + \frac{\mathbf{U} \cdot \mathbf{U}}{2\rho} + \rho \phi$ which denote the internal, kinetic, and potential energies, respectively. The pressure P is given as

$$P = \rho RT \equiv (\gamma - 1) \left(E - \frac{\mathbf{U} \cdot \mathbf{U}}{2\rho} - \rho \phi \right). \tag{2.6}$$

Note that Eq. (2.6) contains no exponentials, as does Set2C, and is therefore less expensive to compute. Since this equation uses the total energy as the thermodynamic variable and is written in conservation form, we can conserve energy; this is not guaranteed for Set2C.

Splitting the density and pressure as done for Set2C and total energy as $E(\mathbf{x},t) = E_0(r) + E'(\mathbf{x},t)$, and assuming the reference field is in hydrostatic balance, Set3C can be written in perturbation form as Eq. (2.1) with $\mathbf{q} := \mathbf{q}'_{3C}$, $\mathbf{F} := \mathbf{F}'_{3C}$ and, $\mathbf{S} := \mathbf{S}'_{3C}$ defined as

$$\mathbf{q}'_{3C} = \begin{pmatrix} \rho' \\ \mathbf{U} \\ E' \end{pmatrix}, \quad \mathbf{F}'_{3C} = \begin{pmatrix} \mathbf{U} \\ \frac{\mathbf{U} \otimes \mathbf{U}}{\rho} + P' \mathbf{I} \\ \frac{(E+P)}{\rho} \mathbf{U} \end{pmatrix}, \quad \mathbf{S}'_{3C} = \begin{pmatrix} 0 \\ -(\rho' \nabla \phi + 2\Omega \times \mathbf{U}) \\ 0 \end{pmatrix}. \tag{2.7}$$

3. IMplicit-EXplicit formulations

3.1. Semi-implicit time integration

The governing equations for both Set2C and Set3C can be written in vector form as

$$\frac{\partial \mathbf{q}}{\partial t} = S(\mathbf{q}) = \{E(\mathbf{q})\} + [I(\mathbf{q})]
= \{S(\mathbf{q}) - L(\mathbf{q})\} + [L(\mathbf{q})]$$
(3.1)

where for Set2C and Set3C, $\mathbf{q}_{2C} := (\rho, \mathbf{U}, \Theta)$ and $\mathbf{q}_{3C} := (\rho, \mathbf{U}, E)$, respectively, $S(\mathbf{q})$ contains all spatial terms, $E(\mathbf{q})$ contains the operators to be treated explicitly, $I(\mathbf{q})$ contains the operators to be treated implicitly, and $L(\mathbf{q})$ is a linear operator to be defined in the following sections. The semi-implicit formulation allows the full problem $(S(\mathbf{q}))$ to be decomposed into subsets that can be treated differently. In Eq. (3.1), the terms in the curly brackets $\{\cdot\}$ contain the slow moving advection terms and are treated explicitly, whereas the terms in the square brackets $[\cdot]$ contain the fast moving acoustic and gravity terms and are treated implicitly. This implicit treatment of the restrictive dynamics removes the CFL constraint due to the fast propagating acoustic and gravity waves and allows for larger time steps that are dictated only by the slow dynamics. The implicit treatment of linear operators requires the solution of a linear system whereas the treatment of nonlinear operators requires a nonlinear system to be solved.

Consider a single-step multistage method with s stages written as

$$\mathbf{Q}^{(i)} = \mathbf{q}^{n} + \Delta t \left[\sum_{j=1}^{i-1} a_{ij} \left(S\left(\mathbf{Q}^{(j)}\right) - L\left(\mathbf{Q}^{(j)}\right) \right) \right] + \Delta t \left[\sum_{j=1}^{i-1} \widetilde{a}_{ij} \left(L\left(\mathbf{Q}^{(j)}\right) \right) \right] + \Delta t \widetilde{a}_{ii} L\left(\mathbf{Q}^{(i)}\right)$$
(3.2)

where $\mathbf{Q}^{(i)}$ is the i^{th} stage approximation, Δt is the time-step, a_{ij} , \widetilde{a}_{ij} , b_i , and \widetilde{b}_i are the coefficients in the Butcher tableaux for the explicit and implicit components, respectively. Following [15] we define

$$\mathbf{q}_{tt} = \mathbf{Q}^{(i)} + \sum_{i=1}^{i-1} \frac{\widetilde{a}_{ij} - a_{ij}}{\widetilde{a}_{ii}} \mathbf{Q}^{(j)}$$
(3.3)

and $\hat{\mathbf{q}}$ as

$$\widehat{\mathbf{q}} = \mathbf{q}^n + \sum_{i=1}^{i-1} \frac{\widetilde{a}_{ij} - a_{ij}}{\widetilde{a}_{ii}} \mathbf{Q}^{(j)} + \Delta t \sum_{j=1}^{i-1} a_{ij} S\left(\mathbf{Q}^{(j)}\right), \tag{3.4}$$

where \mathbf{q}^n is the solution at time $n\Delta t$. Eq. (3.2) can be written compactly as

$$\mathbf{q}_{tt} = \widehat{\mathbf{q}} + \alpha L(\mathbf{q}_{tt}) \tag{3.5}$$

where $\alpha = \Delta t \, \widetilde{a}_{ii}$. The stage value $\mathbf{Q}^{(i)}$ can be obtained from \mathbf{q}_{tt} using Eq. (3.3) and the solution update at t_{n+1} can be constructed as

$$\mathbf{q}^{n+1} = \mathbf{q}^n + \Delta t \sum_{i=1}^s b_i \left[S\left(\mathbf{Q}^{(i)}\right) - L\left(\mathbf{Q}^{(i)}\right) \right] + \Delta t \sum_{i=1}^s \widetilde{b_i} L\left(\mathbf{Q}^{(i)}\right). \tag{3.6}$$

Assuming that the coefficients in the Butcher tableaux satisfy $b_i = \tilde{b_i}$ needed for conservation (see [15]), we obtain

$$\mathbf{q}^{n+1} = \mathbf{q}^n + \Delta t \sum_{i=1}^s b_i S\left(\mathbf{Q}^{(i)}\right). \tag{3.7}$$

In what follows, we derive the *No-Schur* and Schur forms for equation sets 2C and 3C, however, to simplify the exposition, we perform the derivation using continuous differential operators. In the appendices, we substitute the continuous differential operators with their discrete forms to complete the discussion.

3.2. IMplicit-EXplicit formulation: Set2C

3.2.1. 3D-IMEX formulation

Let us now describe the IMEX formulation for Set2C. Following [24], the linear operator for Set2C is taken as

$$L_{2C}(\mathbf{q}) = -\begin{pmatrix} \nabla \cdot \mathbf{U} \\ \nabla P_L' + \rho' \nabla \phi \\ \nabla \cdot \left(\frac{\Theta_0}{\rho_0} \mathbf{U} \right) \end{pmatrix}$$
(3.8)

with the linearized pressure as

$$P_L' = \frac{\gamma P_0}{\Theta_0} \Theta' \tag{3.9}$$

which is derived from a first order Taylor series approximation. It can be seen that the linear operator contains the fast moving acoustic (P'_L) and gravity wave (ϕ) terms. Applying the semi-implicit formulation (Eq. (3.5)) to Set2C and letting $h_{2C} = \frac{\Theta_0}{\Omega_0}$, yields the *No-Schur* form of the time-discretized equations

$$\rho_{tt} = \widehat{\rho} - \alpha \nabla \cdot \mathbf{U}_{tt} \tag{3.10a}$$

$$\mathbf{U}_{tt} = \widehat{\mathbf{U}} - \alpha \left(\nabla P_{tt} + \rho_{tt} \nabla \phi \right) \tag{3.10b}$$

$$\Theta_{tt} = \widehat{\Theta} - \alpha \nabla \cdot (h_{2C} \mathbf{U}_{tt}) \tag{3.10c}$$

$$P_{tt} = \frac{\gamma P_0}{\Theta_0} \Theta_{tt}. \tag{3.10d}$$

Since the linearization is performed about time-independent reference fields ($\rho_0 \neq \rho_0(t)$, $\Theta_0 \neq \Theta_0(t)$, $E_0 \neq E_0(t)$), the Jacobian is constant in time and only needs to be constructed once. To derive the Schur complement of Eqs. (3.10), we begin by substituting Eq. (3.10c) into Eq. (3.10d) to obtain

$$P_{tt} = \frac{\gamma P_0}{\Theta_0} \left(\widehat{\Theta} - \alpha \nabla \cdot [h_{2C} \mathbf{U}_{tt}] \right). \tag{3.11}$$

Multiplying Eq. (3.10a) by h_{2C} and subtracting from Eq. (3.10c) yields

$$\Theta_{tt} - h_{2C} \rho_{tt} = \widehat{\Theta} - h_{2C} \widehat{\rho} - \alpha \left(\nabla \cdot (h_{2C} \mathbf{U}_{tt}) - h_{2C} \nabla \cdot \mathbf{U}_{tt} \right) \\
\Theta_{tt} - h_{2C} \rho_{tt} = \widehat{\Theta} - h_{2C} \widehat{\rho} - \alpha \left(\mathbf{U}_{tt} \cdot \nabla h_{2C} \right)$$
(3.12)

where we have made the assumption that the product rule is satisfied (although for element-based Galerkin methods this is not the case when using inexact integration - see, e.g., [1]); we will see in the results section that this assumption is reasonable. Substituting Eq. (3.10d) into Eq. (3.12) yields

$$\rho_{tt} = \widehat{\rho} - \frac{\widehat{\Theta}}{h_{2C}} + \frac{\Theta_0 P_{tt}}{h_{2C} \gamma P_0} + \frac{\alpha}{h_{2C}} \left(\mathbf{U}_{tt} \cdot \nabla h_{2C} \right). \tag{3.13}$$

Letting $G_{2C} = \frac{\nabla \phi}{h_{2C}}$, substituting Eq. (3.13) into Eq. (3.10b), and solving for \mathbf{U}_{tt} yields

$$\mathbf{U}_{tt} = \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha \left[\nabla P_{tt} + \left(\widehat{\rho} \nabla \phi - G_{2C} \widehat{\Theta} + G_{2C} \frac{\Theta_0 P_{tt}}{\gamma P_0} \right) \right] \right)$$
(3.14)

where $\mathbf{A} = \mathbf{I} + \alpha^2 G_{2C} \otimes \nabla h_{2C}$. Substituting Eq. (3.14) into Eq. (3.11) and rearranging yields the Schur complement for Set2C

$$P_{tt} - \frac{\gamma P_0 \alpha^2}{\Theta_0} \nabla \cdot (h_{2C} \mathcal{L}_{2C}) = \frac{\gamma P_0}{\Theta_0} \left(\widehat{\Theta} - \alpha \nabla \cdot [h_{2C} \mathcal{R}_{2C}] \right)$$
(3.15)

where
$$\mathcal{L}_{2C} = \mathbf{A}^{-1} \left(\nabla P_{tt} + G_{2C} \frac{\Theta_0 P_{tt}}{\gamma P_0} \right)$$
 and $\mathcal{R}_{2C} = \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha \left[\widehat{\rho} \nabla \phi - G_{2C} \widehat{\Theta} \right] \right)$.

This reduces the system of equations from five prognostic variables to a single (Helmholtz-like) equation for pressure. Once the pressure is computed, the velocities can be extracted using Eq. (3.14), followed by the density Eq. (3.13), and then the potential temperature Eq. (3.12).

3.2.2. 1D-IMEX formulation

The formulation in Eq. (3.15), which treats the acoustics terms in all directions, can be further reduced to implicitly treat only the vertical/radial direction. This uni-directional formulation will be referred to as 1D-IMEX and is obtained by replacing all three-dimensional (3D) differential operators $\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z}\right)$ with radial derivatives $\nabla_r = \frac{\partial}{\partial r}$ that are one-dimensional (1D). Letting $A_r = 1 + \alpha^2 G_{2C} \frac{\partial h_{2C}}{\partial r}$, the Schur complement of the 1D-IMEX formulation is given as

$$P_{tt} - \frac{\gamma P_0 \alpha^2}{\Theta_0} \nabla_r \cdot \left(h_{2C} \mathcal{L}_{2C}^r \right) = \frac{\gamma P_0}{\Theta_0} \left(\widehat{\Theta} - \alpha \nabla_r \cdot \left[h_{2C} \mathcal{R}_{2C}^r \right] \right)$$
(3.16)

where $\mathcal{L}_{2C}^r = A_r^{-1} \left(\frac{\partial P_{tt}}{\partial r} + G_{2C} \frac{\Theta_0 P_{tt}}{\gamma P_0} \right)$, and $\mathcal{R}_{2C}^r = A_r^{-1} \left(\widehat{U_r} - \alpha \left[\widehat{\rho} \nabla \phi - G_{2C} \widehat{\Theta} \right] \right)$. It should be noted that the Schur complement is valid when **A** is nonsingular, a condition that is satisfied for a stable stratified reference atmosphere used in this work (see [24]).

3.3. IMplicit-EXplicit formulation: Set3C

3.3.1. 3D-IMEX formulation

We now present the IMEX formulation for Set3C. Following [24], the linear operator for Set3C is taken as

$$L_{3C}(\mathbf{q}) = -\begin{pmatrix} \nabla \cdot \mathbf{U} \\ \nabla P_L' + \rho' \nabla \phi \\ \nabla \cdot \left(\frac{(E_0 + P_0)}{\rho_0} \mathbf{U} \right) \end{pmatrix}$$
(3.17)

with linearized pressure

$$P'_{I} = (\gamma - 1) \left(E' - \rho' \phi \right). \tag{3.18}$$

Then applying the semi-implicit formulation (Eq. (3.5)) to Set3C and letting $h_{3C} = \frac{(E_0 + P_0)}{\rho_0}$ be the reference enthalpy, yields the *No-Schur* form of the time-discretized equations

$$\rho_{tt} = \widehat{\rho} - \alpha \nabla \cdot \mathbf{U}_{tt}$$
 (3.19a)

$$\mathbf{U}_{tt} = \widehat{\mathbf{U}} - \alpha \left(\nabla P_{tt} + \rho_{tt} \nabla \phi \right) \tag{3.19b}$$

$$E_{tt} = \widehat{E} - \alpha \nabla \cdot (h_{3C} \mathbf{U}_{tt}) \tag{3.19c}$$

$$P_{tt} = (\gamma - 1) (E_{tt} - \rho_{tt} \phi).$$
 (3.19d)

To derive the Schur complement of Eqs. (3.19), we begin by substituting Eq. (3.19a) and Eq. (3.19c) into Eq. (3.19d), and letting $F_{3C} = h_{3C} - \phi$ to obtain

$$P_{tt} = (\gamma - 1)(\widehat{E} - \phi\widehat{\rho}) - \alpha(\gamma - 1)(F_{3C}\nabla \cdot \mathbf{U}_{tt} + \nabla h_{3C} \cdot \mathbf{U}_{tt}). \tag{3.20}$$

Multiplying Eq. (3.19a) by h_{3C} and subtracting from Eq. (3.19c) yields

$$E_{tt} - h_{3C}\rho_{tt} = \widehat{E} - h_{3C}\widehat{\rho} - \alpha \nabla h_{3C} \cdot \mathbf{U}_{tt}. \tag{3.21}$$

Substituting Eq. (3.21) into Eq. (3.19d) yields

$$P_{tt} = (\gamma - 1)\left(\widehat{E} + h_{3C}\rho_{tt} - h_{3C}\widehat{\rho} - \alpha\nabla h_{3C} \cdot \mathbf{U}_{tt} - \rho_{tt}\phi\right)$$
(3.22)

which can be solved for density ρ_{tt} as

$$\rho_{tt} = \frac{1}{F_{3C}} \left(\frac{1}{\gamma - 1} P_{tt} + \alpha \nabla h_{3C} \cdot \mathbf{U}_{tt} + h_{3C} \widehat{\rho} - \widehat{\mathbf{E}} \right). \tag{3.23}$$

Letting $G_{3C} = \frac{\nabla \phi}{F_{3C}}$, substituting Eq. (3.23) into Eq. (3.19b), and solving for \mathbf{U}_{tt} yields

$$\mathbf{U}_{tt} = \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha G_{3C} \left(h_{3C} \widehat{\rho} - \widehat{E} \right) - \alpha \nabla P_{tt} - \frac{\alpha G_{3C}}{(\gamma - 1)} P_{tt} \right)$$
(3.24)

where $\mathbf{A} = \mathbf{I} + \alpha^2 G_{3C} \otimes \nabla h_{3C}$. Then, substituting Eq. (3.24) into Eq. (3.20) and rearranging yields the Schur complement for Set3C

$$P_{tt} - \alpha^{2} (\gamma - 1) F_{3C} \nabla \cdot \mathcal{L}_{3C} - \alpha^{2} (\gamma - 1) \nabla h_{3C} \cdot \mathcal{L}_{3C}$$

$$= (\gamma - 1) (\widehat{E} - \phi \widehat{\rho}) - \alpha (\gamma - 1) F_{3C} \nabla \cdot \mathcal{R}_{3C} - \alpha (\gamma - 1) \nabla h_{3C} \cdot \mathcal{R}_{3C}$$
(3.25)

where $\mathcal{L}_{3C} = \mathbf{A}^{-1} \left(\nabla P_{tt} + \frac{G_{3C}}{(\gamma - 1)} P_{tt} \right)$ and $\mathcal{R}_{3C} = \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha G_{3C} \left(h_{3C} \widehat{\rho} - \widehat{E} \right) \right)$. The Schur complement can be written more compactly as follows

$$P_{tt} - \alpha^{2} (\gamma - 1) \left(\nabla \cdot (h_{3C} \mathcal{L}_{3C}) - \phi \nabla \cdot \mathcal{L}_{3C} \right)$$

$$= (\gamma - 1) \left(\widehat{E} - \phi \widehat{\rho} - \alpha \left(\nabla \cdot (h_{3C} \mathcal{R}_{3C}) - \phi \nabla \cdot \mathcal{R}_{3C} \right) \right)$$
(3.26)

and is the form that will be used for deriving the discrete forms.

3.3.2. 1D-IMEX formulation

As was done for Set2C (Eq. (3.16)), the 1D-IMEX formulation for Set3C can be obtained by substituting $\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z}\right)$ in Eq. (3.26) with radial derivatives $\nabla_r = \frac{\partial}{\partial r}$. Letting $A_r = 1 + \alpha^2 G_{3C} \frac{\partial h_{3C}}{\partial r}$, the Schur complement of the 1D-IMEX formulation for Set3C is written as

$$P_{tt} - \alpha^{2} (\gamma - 1) \left(\nabla_{r} \cdot \left(h_{3C} \mathcal{L}_{3C}^{r} \right) - \phi \nabla_{r} \cdot \mathcal{L}_{3C}^{r} \right)$$

$$= (\gamma - 1) \left(\widehat{E} - \phi \widehat{\rho} - \alpha \left(\nabla_{r} \cdot \left(h_{3C} \mathcal{R}_{3C}^{r} \right) - \phi \nabla_{r} \cdot \mathcal{R}_{3C}^{r} \right) \right)$$
where $\mathcal{L}_{3C}^{r} = A_{r}^{-1} \left(\frac{\partial P_{tt}}{\partial r} + \frac{G_{3C}}{(\gamma - 1)} P_{tt} \right)$, and $\mathcal{R}_{3C}^{r} = A_{r}^{-1} \left(\widehat{U}_{r} - \alpha G_{3C} \left(h_{3C} \widehat{\rho} - \widehat{E} \right) \right)$. (3.27)

3.4. Complexity comparison between Schur and no-Schur forms

For both sets 2C and 3C, the Schur form, compared to the No-Schur form, reduces the size of the resulting matrix problem by a factor of $(d+2)^2$ where d denotes the number of spatial dimensions; d represents the number of momentum equations and the 2 refers to the two scalar equations for the density and thermodynamic variable. For example, in three-dimensions the savings for the 3D-IMEX formulation is a factor of 5^2 since we go from 5 equations (ρ, U, V, W, E) to one equation for pressure. For the 1D-IMEX formulation, we let d=1. There are some additional costs incurred by the Schur form in order to extract the d+2 prognostic variables from the pressure variable but the cost is equivalent to one iteration of the No-Schur form; as long as more than one elliptic solver iteration is required, the Schur form will be faster. For the 1D-IMEX formulation, the cost comparison is more straightforward since we use direct (LU) solvers. For a standard LU decomposition, the LU decomposition costs $\mathcal{O}(N^3)$ for N degrees of freedom while the back substitution costs $\mathcal{O}(N^2)$. Therefore, if we use a constant reference solution for the linearization of the IMEX method, then the LU decomposition only occurs once at the beginning of the simulation whereas the back substitution occurs at each of the s-stages of each time-step. This results in savings of a factor of $(d+2)^2=9$ by using the Schur form over its No-Schur counterpart. For the 3D-IMEX, there is additional savings for the Schur form due to the use of elliptic solvers. Since the Schur form results in a better conditioned matrix (shown in Sec. 6), this results in a faster time-to-solution compared to its No-Schur counterpart (since the cost of, e.g., GMRES, increases quadratically with iteration count).

4. Numerical implementation

4.1. Spatial discretization

We now turn to the description of the discontinuous Galerkin discretization ([1]) of the Schur forms for Set2C and Set3C. Multiplying Eq. (3.15) and Eq. (3.26) by a test function ψ and integrating over an element Ω_e

$$\int_{\Omega_{e}} \psi P_{tt} d\Omega_{e} - \alpha^{2} \int_{\Omega_{e}} \psi \frac{\gamma P_{0}}{\Theta_{0}} \nabla \cdot (h_{2C} \mathcal{L}_{2C}) d\Omega_{e} = \int_{\Omega_{e}} \psi \frac{\gamma P_{0}}{\Theta_{0}} \widehat{\Theta} d\Omega_{e} - \alpha \int_{\Omega_{e}} \psi \frac{\gamma P_{0}}{\Theta_{0}} \nabla \cdot (h_{2C} \mathcal{R}_{2C}) d\Omega_{e}$$

$$\tag{4.1}$$

$$\int_{\Omega_{e}} \psi P_{tt} d\Omega_{e} - \alpha^{2} (\gamma - 1) \int_{\Omega_{e}} \psi (\nabla \cdot (h_{3C} \mathcal{L}_{3C}) - \phi \nabla \cdot \mathcal{L}_{3C}) d\Omega_{e} = (\gamma - 1) \int_{\Omega_{e}} \psi (\widehat{E} - \phi \widehat{\rho}) d\Omega_{e}$$

$$- \alpha (\gamma - 1) \int_{\Omega_{e}} \psi (\nabla \cdot (h_{3C} \mathcal{R}_{3C}) - \phi \nabla \cdot \mathcal{R}_{3C}) d\Omega_{e},$$

$$(4.2)$$

we obtain the variational forms of the Schur complements for Set2C and Set3C, respectively. Without loss of generality, the $\mathcal L$ and $\mathcal R$ operators for both equation sets can be written as

$$\mathcal{L} = \mathbf{A}^{-1} \left(\nabla P_{tt} + a P_{tt} \right)$$

$$\mathcal{R} = \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha \left(\delta \widehat{\rho} - c \widehat{T} \right) \right)$$
(4.3)

where $a = G_{2C} \frac{\Theta_0}{\gamma P_0}$, $b = \nabla \phi$, $c = G_{2C}$ and the thermodynamic variable $\widehat{T} = \widehat{\Theta}$ for Set2C, and $a = \frac{G_{3C}}{\gamma - 1}$, $b = G_{3C}h_{3C}$, c = 1 and $\widehat{T} = \widehat{E}$ for Set3C.

It can be seen from Eq. (4.1) and Eq. (4.2) that the differential operators are similar, therefore, we will show the discretization for Set3C with $\mathcal{F} \in \{\mathcal{L}_{3C}, \mathcal{R}_{3C}\}$, $h_C = h_{3C}$ and $\xi = \phi$ but discretizations can be obtained for Set2C with $\mathcal{F} \in \{\mathcal{L}_{2C}, \mathcal{R}_{2C}\}$, $h_C = h_{2C}$, and $\xi = 0$. Similarly, the 1D-IMEX discretization can be obtained by retaining only the radial differential operators. It can be seen that $\nabla \cdot \mathcal{L}$ results in both first and second order operators while $\nabla \cdot \mathcal{R}$ results in only first order operators. The second order operators are discretized using the local discontinuous Galerkin (LDG) method [25].

We begin by first discretizing the first-order operators in the Schur form

$$\int_{\Omega_{e}} \psi \left(\omega \nabla \cdot (h_{C} \mathcal{F}) - \xi \nabla \cdot \mathcal{F} \right) d\Omega_{e} = \int_{\Omega_{e}} \psi \omega \nabla \cdot (h_{C} \mathcal{F}) d\Omega_{e} - \int_{\Omega_{e}} \psi \xi \nabla \cdot \mathcal{F} d\Omega_{e}$$

$$\tag{4.4}$$

where $\omega = \frac{\gamma P_0}{\Theta_0}$ for Set2C, $\omega = 1$ for Set3C and ψ are the multi-dimensional Lagrange polynomials defined as a tensor product of 1D Lagrange polynomials supported at Legendre-Gauss-Lobotto (LGL) points. The integrals are approximated by means of inexact integration using LGL quadrature (for more details see, e.g., [1]). Applying integration by parts to the first term on the right-hand side yields

$$\int_{\Omega_{e}} \psi \omega \nabla \cdot (h_{C} \mathcal{F}) d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (\psi \omega h_{C} \mathcal{F}) d\Omega_{e} - \int_{\Omega_{e}} h_{C} \mathcal{F} \cdot \nabla (\psi \omega) d\Omega_{e}$$

$$= \int_{\Omega_{e}} \widetilde{\psi \omega h_{C} \mathcal{F}} \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \omega h_{C} \nabla \psi \cdot \mathcal{F} d\Omega_{e} - \int_{\Omega_{e}} \psi h_{C} \nabla \omega \cdot \mathcal{F} d\Omega_{e}$$
(4.5)

where $\psi \omega h_C \mathcal{F}$ is some numerical flux to be defined. Due to the use of conforming elements and symmetry of the LGL nodes, this simplifies to $\psi \omega h_C \mathcal{F} = \psi \omega h_C \mathcal{F}$, i.e., the basis function satisfies $\psi = \widetilde{\psi}$. The second term on the right-hand side in Eq. (4.4) can be written as

$$\int_{\Omega_{e}} \psi \xi \nabla \cdot \mathcal{F} \, d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (\psi \xi \mathcal{F}) \, d\Omega_{e} - \int_{\Omega_{e}} \mathcal{F} \cdot \nabla (\psi \xi) \, d\Omega_{e}$$

$$= \int_{\Gamma_{e}} \psi \widetilde{\xi} \widetilde{\mathcal{F}} \cdot \widehat{\mathbf{n}} \, d\Gamma_{e} - \int_{\Omega_{e}} \xi \nabla \psi \cdot \mathcal{F} \, d\Omega_{e} - \int_{\Omega_{e}} \psi \nabla \xi \cdot \mathcal{F} \, d\Omega_{e}.$$
(4.6)

In flow-in-a-box mode, since for Set3C, $\xi = \phi = gz$ then $\nabla \phi = g\hat{k}$ where \hat{k} is the unit vector along the z direction. Substituting Eq. (4.5) and Eq. (4.6) into Eq. (4.4), and rearranging yields the weak form of the differential operator

$$\int_{\Omega_{e}} \psi \left(\omega \nabla \cdot (h_{C} \mathcal{F}) - \xi \nabla \cdot \mathcal{F} \right) d\Omega_{e} = \int_{\Gamma_{e}} \psi \left(\widetilde{\omega h_{C} \mathcal{F}} - \widetilde{\xi} \widetilde{\mathcal{F}} \right) \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \left(\omega h_{C} - \xi \right) \nabla \psi \cdot \mathcal{F} d\Omega_{e}
- \int_{\Omega_{e}} \psi \left(h_{C} \nabla \omega - \nabla \xi \right) \cdot \mathcal{F} d\Omega_{e}.$$
(4.7)

Applying integration by parts to the second term on the right-hand side of Eq. (4.7) yields

$$\int_{\Omega_{e}} (\omega h_{C} - \xi) \nabla \psi \cdot \mathcal{F} \, d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (\psi \, (\omega h_{C} - \xi) \, \mathcal{F}) \, d\Omega_{e} - \int_{\Omega_{e}} \psi \nabla \cdot ((\omega h_{C} - \xi) \, \mathcal{F}) \, d\Omega_{e}$$

$$= \int_{\Gamma_{e}} \psi \, (\omega^{e} h_{C}^{e} - \xi^{e}) \, \mathcal{F}^{e} \cdot \hat{\mathbf{n}} \, d\Gamma_{e} - \int_{\Omega_{e}} \psi \, ((\omega h_{C} - \xi) \, \nabla \cdot \mathcal{F} + \nabla (\omega h_{C} - \xi) \cdot \mathcal{F}) \, d\Omega_{e}, \tag{4.8}$$

where we have added the superscript e to the surface integral only (first term on the right-hand side) to denote a local element-wise quantity. Substituting Eq. (4.8) into Eq. (4.7), we obtain the strong form

$$\int_{\Omega_{e}} \psi \left(\nabla \cdot (\omega h_{C} \mathcal{F}) - \xi \nabla \cdot \mathcal{F} \right) d\Omega_{e} = \int_{\Gamma_{e}} \psi \left(\left(\widetilde{\omega h_{C} \mathcal{F}} - \widetilde{\xi} \widetilde{\mathcal{F}} \right) - \left(\omega^{e} h_{C}^{e} - \phi^{e} \right) \mathcal{F}^{e} \right) \cdot \widehat{\mathbf{n}} d\Gamma_{e}
+ \int_{\Omega_{e}} \psi \left(\omega h_{C} - \xi \right) \nabla \cdot \mathcal{F} d\Omega_{e} + \int_{\Omega_{e}} \psi \omega \nabla h_{C} \cdot \mathcal{F} d\Omega_{e} \tag{4.9}$$

where it should be understood that terms with no designation (either a superscript e or tilde) strictly represent local element-wise quantities. If we assume the reference fields to be continuous then $(\omega h_C \mathcal{F} - \widetilde{\xi} \widetilde{\mathcal{F}}) = (\omega h_C - \xi) \widetilde{\mathcal{F}}$. Furthermore, if $\mathcal{F} = \mathcal{R}$ then

$$\widetilde{\mathcal{R}} = \mathbf{A}^{-1} \left(\widetilde{\widetilde{\mathbf{U}}} - \alpha \left(\mathbf{b} \ \widetilde{\widehat{\rho}} - c \ \widetilde{\widehat{T}} \right) \right)$$

since \mathbf{A} , a, b and c are also continuous (as they are dependent on the reference fields). The case where $\mathcal{F} = \mathcal{L}$ results in both first and second order operators. Expanding $\mathcal{L} = \mathbf{A}^{-1} \nabla P_{tt} + \mathbf{A}^{-1} a P_{tt}$, the second term on the right-hand side can be treated by letting $\mathcal{F} = \mathbf{A}^{-1} a P_{tt}$ in Eq. (4.7) and Eq. (4.9) (and setting the coefficients appropriately), but the first term requires computation of the variational form of the gradient consistent with the DG scheme. This is done by multiplying ∇P_{tt} by a tensor test function $\Psi = \psi \mathbf{I}$ and invoking the divergence theorem resulting in the variational weak form of the gradient

$$\int_{\Omega_{e}} \Psi \cdot \nabla P_{tt} \, d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (\Psi P_{tt}) \, d\Omega_{e} - \int_{\Omega_{e}} P_{tt} \nabla \cdot \Psi \, d\Omega_{e}$$

$$= \int_{\Gamma_{e}} \widetilde{P_{tt} \Psi} \cdot \widehat{\mathbf{n}} \, d\Gamma_{e} - \int_{\Omega_{e}} P_{tt} \nabla \cdot \Psi \, d\Omega_{e}$$
(4.10)

where $\widetilde{P_{tt}\Psi} = \widetilde{P_{tt}\Psi}$, since $\widetilde{\psi} = \psi$. Applying integration by parts to the second term on the right hand side of Eq. (4.10) yields

$$\int_{\Omega_{e}} P_{tt} \nabla \cdot \Psi \, d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (P_{tt} \Psi) \, d\Omega_{e} - \int_{\Omega_{e}} \nabla P_{tt} \cdot \Psi \, d\Omega_{e}$$

$$= \int_{\Gamma_{e}} P_{tt}^{e} \Psi \cdot \hat{\mathbf{n}} \, d\Gamma_{e} - \int_{\Omega_{e}} \nabla P_{tt} \cdot \Psi \, d\Omega_{e}.$$
(4.11)

Substituting Eq. (4.11) into Eq. (4.10) yields the variational strong form for the gradient

$$\int_{\Omega_{e}} \Psi \cdot \nabla P_{tt} \, d\Omega_{e} = \int_{\Omega_{e}} \nabla \cdot (\Psi P_{tt}) \, d\Omega_{e} - \int_{\Omega_{e}} P_{tt} \nabla \cdot \Psi \, d\Omega_{e}$$

$$= \int_{\Gamma_{e}} \left(\widetilde{P}_{tt} - P_{tt}^{e} \right) \Psi \cdot \widehat{\mathbf{n}} \, d\Gamma_{e} + \int_{\Omega_{e}} \nabla P_{tt} \cdot \Psi \, d\Omega_{e}.$$
(4.12)

Then the weak and strong forms for \mathcal{L} are

$$\mathcal{L}_{weak} = \mathbf{A}^{-1} \left[\int_{\Gamma_e} \widetilde{P}_{tt} \Psi \cdot \widehat{\mathbf{n}} \, d\Gamma_e - \int_{\Omega_e} P_{tt} \nabla \cdot \Psi \, d\Omega_e + a P_{tt} \right]$$

$$= \mathbf{A}^{-1} \left[\mathcal{C}^{(e,k)} - \widetilde{\mathbf{D}}^{(e)} + a \mathbf{I} \right] P_{tt}$$
(4.13)

and

$$\mathcal{L}_{strong} = \mathbf{A}^{-1} \left[\int_{\Gamma_e} \left(\widetilde{P}_{tt} - P_{tt}^e \right) \Psi \cdot \widehat{\mathbf{n}} \, d\Gamma_e + \int_{\Omega_e} \nabla P_{tt} \cdot \Psi \, d\Omega_e + a P_{tt} \right]$$

$$= \mathbf{A}^{-1} \left[\mathbf{\Delta}^{(e,k)} + \mathbf{D}^{(e)} + a \mathbf{I} \right] P_{tt}$$
(4.14)

respectively, where **I** is the identity matrix, $C^{(e,k)}$ and $\Delta^{(e,k)}$ are the centered and difference matrices, and $D^{(e)}$ and $D^{(e)}$ are the strong and weak form differentiation matrices, which are defined in Appendix C with more details found in [[1], Ch. 21]. The matrices with superscript (e,k) are global matrices that require communication between the element e and its neighbor e, whereas the matrices with superscript (e) are local matrices that do not require communication with neighboring elements. Since the fully discrete representation of the Schur forms are unwieldy, they are presented in Appendix C. Unless explicitly stated the use of the strong form, the results and analysis of various test cases are obtained using the weak form.

4.2. Conditions on numerical flux for consistent IMEX splitting

Here we show the conditions required for the IMEX splitting to be consistent on the continuous and discrete levels. By <u>consistent</u>, we mean that the discrete form of the implicitly and explicitly treated operators appropriately represent their <u>continuous</u> forms (e.g., the wavespeeds in the discrete form are those for the continuous operator). We only show this for Set3C, but the same conditions apply for Set2C. Consider a balance law (Set3C) written in vector form defined by Eq. (2.1) and Eq. (2.5), where, for clarity, we denote the total flux as $\mathbf{F}_T = \mathbf{F}$ with components $(F^\rho, \mathbf{F}^U, F^E)$ that represent the mass, momentum and total energy flux. The coloring and subscripts T, L, and NL denote total, linear (blue), and nonlinear (red) components, respectively.

Here onwards, for convenience, the source terms **S** due to gravity and Coriolis are neglected; thereby effectively replacing the original balance law with a conservation law with no loss in generality (because the stiffness in the equations is derived from the terms in the flux terms). Casting Eq. (2.1) into the variational weak form (and invoking the divergence theorem) yields

$$\int_{\Omega_e} \psi \frac{\partial \mathbf{q}}{\partial t} d\Omega_e + \int_{\Gamma_e} \widetilde{\psi} \widetilde{\mathbf{F}}_T \cdot \widehat{n} d\Gamma_e - \int_{\Omega_e} \nabla \psi \cdot \mathbf{F}_T d\Omega_e = 0$$
(4.15)

where \widetilde{Q} represents the numerical flux for the quantity Q. Due to the symmetry of the Legendre-Gauss-Lobotto (LGL) points, tensor product representation of the Lagrange polynomials, and conforming elements, $\widetilde{\psi}\widetilde{\mathbf{F}_T} = \psi\widetilde{\mathbf{F}_T}$. Taking the numerical flux \widetilde{Q} to be the Rusanov flux for some quantity Q with flux \mathbf{F}_O as

$$\widetilde{Q} = \left\{ \left\{ \mathbf{F}_{Q} \right\} \right\} - \frac{\lambda}{2} \mathbb{I} Q \, \mathbb{I} \widehat{n}, \tag{4.16}$$

Eq. (4.15) then becomes

$$\int_{\Omega_{e}} \psi \frac{\partial \mathbf{q}}{\partial t} d\Omega_{e} + \int_{\Gamma_{e}} \psi \left(\{ \{ \mathbf{F}_{T} \} \} - \frac{\lambda_{T}}{2} [\mathbf{q}] \widehat{n} \right) \cdot \widehat{n} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{F}_{T} d\Omega_{e} = 0$$

$$(4.17)$$

where $\{\{\cdot\}\}\$ and $[\![\cdot]\!]$ represent the average and jump operators, $\lambda_T = (u+a)^\star$ is the maximum eigenvalue (wave speed) of the total flux Jacobian $\left(\frac{\partial F_T}{\partial \mathbf{q}}\right)$, \star denotes the maximum of the left and right states (i.e., $A^\star = \max(|A^L|, |A^R|)$), and $a = \sqrt{\frac{\gamma P}{\rho}}$ is the speed of sound.

The flux vector $\mathbf{F}_T := \mathbf{F}$ in Eq. (2.5) on the continuous level can be decomposed into its linear \mathbf{F}_L and nonlinear $\mathbf{F}_{NL} = \mathbf{F}_T - \mathbf{F}_L$ components as

$$\mathbf{F}_{T} = \begin{pmatrix} \mathbf{0} \\ \frac{\mathbf{U} \otimes \mathbf{U}}{\rho} + P\mathbf{I} \\ \left(\frac{E+P}{\rho}\right)\mathbf{U} \end{pmatrix}, \quad \mathbf{F}_{L} = \begin{pmatrix} \mathbf{U} \\ P_{L}\mathbf{I} \\ h_{0}\mathbf{U} \end{pmatrix}, \quad \mathbf{F}_{NL} = \begin{pmatrix} \mathbf{0} \\ \frac{\mathbf{U} \otimes \mathbf{U}}{\rho} + (P-P_{L})\mathbf{I} \\ \left(\frac{E+P}{\rho} - h_{0}\right)\mathbf{U} \end{pmatrix}, \tag{4.18}$$

where $h_0 = \frac{E_0 + P_0}{\rho_0}$ is the reference enthalpy and $P_L = (\gamma - 1) (E - \rho \phi)$ where \mathbf{F}_L is constructed via a first order Taylor series approximation of \mathbf{F}_T . Using such a decomposition and rearranging terms, Eq. (4.17) can also be decomposed into its linear (blue) and nonlinear components (red) on the discrete level as

$$\int_{\Omega_{e}} \psi \frac{\partial \mathbf{q}}{\partial t} d\Omega_{e} + \int_{\Gamma_{e}} \psi \left(\{ \{ \mathbf{F}_{NL} \} \} - \frac{\lambda_{T}}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{F}_{NL} d\Omega_{e}
+ \int_{\Gamma_{e}} \psi \left\{ \{ \mathbf{F}_{L} \} \} \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{F}_{L} d\Omega_{e} = 0 \right\}
\int_{\Omega_{e}} \psi \nabla \cdot \mathbf{F}_{L} d\Omega_{e} = 0$$
(4.19)

where the linear (blue) term represents the variational weak form of the linear flux terms using centered fluxes and the nonlinear (red) term represents the variational weak form of the nonlinear flux terms using the Rusanov flux.

Without loss of generality, we can write $\lambda_T = \lambda_L + \lambda_{NL}$ to obtain a form where each component can be treated using the Rusanov flux using the appropriate wave speeds,

$$\int_{\Omega_{e}} \psi \frac{\partial \mathbf{q}}{\partial t} d\Omega_{e} + \int_{\Gamma_{e}} \psi \left(\{ \{ \mathbf{F}_{NL} \} \} - \frac{\lambda_{NL}}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{F}_{NL} d\Omega_{e} \\
+ \int_{\Gamma_{e}} \psi \left(\{ \{ \mathbf{F}_{L} \} \} - \frac{\lambda_{L}}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{F}_{L} d\Omega_{e} = 0$$

$$\int_{\Omega_{e}} \psi \nabla \cdot \mathbf{F}_{L} d\Omega_{e} \tag{4.20}$$

where λ_L and λ_{NL} are now the maximum eigenvalues of the Jacobians of the linear flux \mathbf{F}_L and nonlinear flux \mathbf{F}_{NL} , respectively.

In one dimension, the linear and total flux terms are given as follows

$$\mathbf{F}_{L} = \begin{pmatrix} U \\ P_{L} \\ h_{0}U \end{pmatrix} \equiv \begin{pmatrix} U \\ P_{0} + (\gamma - 1)\left(E' - \rho'\phi\right) \\ \left(\frac{\gamma E_{0} - (\gamma - 1)\rho_{0}\phi}{\rho_{0}}\right)U \end{pmatrix}, \quad \mathbf{F}_{T} = \begin{pmatrix} U \\ \frac{U^{2}}{\rho} + P \\ \left(\frac{E + P}{\rho}\right)U \end{pmatrix}$$

$$(4.21)$$

with the resulting nonlinear flux

$$\mathbf{F}_{NL} \equiv \mathbf{F}_{T} - \mathbf{F}_{L} = \begin{pmatrix} 0 \\ \frac{U^{2}}{\rho} + P - P_{L} \\ \left(\frac{E+P}{\rho}\right)U - h_{0}U \end{pmatrix} = \begin{pmatrix} 0 \\ \left(\frac{3-\gamma}{2}\right)\frac{U^{2}}{\rho} \\ \gamma U\left(\frac{E}{\rho} - \frac{E_{0}}{\rho_{0}}\right) - (\gamma - 1)\frac{U^{3}}{2\rho^{2}} \end{pmatrix}$$
(4.22)

where the buoyancy term in the momentum equation is omitted because it is considered a source. Computing the Jacobian of the linear flux yields

$$J_{L} = \frac{\partial \mathbf{F}_{L}}{\partial \mathbf{q}} = \begin{pmatrix} 0 & 1 & 0 \\ -(\gamma - 1)\phi & 0 & \gamma - 1 \\ 0 & h_{0} & 0 \end{pmatrix}$$
(4.23)

with eigenvalues $\lambda_L = (-a_0, 0, a_0)$ where $a_0 = \sqrt{\frac{\gamma P_0}{\rho_0}}$. Computing the Jacobian of the total flux \mathbf{F}_T yields

$$J_{T} = \frac{\partial \mathbf{F}_{T}}{\partial \mathbf{q}} = \begin{pmatrix} 0 & 1 & 0 \\ (\gamma - 3)\frac{U^{2}}{2\rho^{2}} - (\gamma - 1)\phi & (3 - \gamma)\frac{U}{\rho} & \gamma - 1 \\ -\frac{U}{\rho^{2}}\left(\gamma E - (\gamma - 1)\frac{U^{2}}{\rho}\right) & \frac{\gamma E}{\rho} - \frac{3(\gamma - 1)U^{2}}{2\rho^{2}} - (\gamma - 1)\phi & \frac{\gamma U}{\rho} \end{pmatrix}$$
(4.24)

with eigenvalues $\lambda_T = (u - a, u, u + a)$.

Table 1Flux combinations and wave speeds (magnitude of the jump term) in the Rusanov flux for the IMEX formulations.

Name	Abbrv.	λ_L	λ_{NL}	λ_T
Acoustic-Total Rusanov	AT	a_0	u + a'	u + a
Acoustic-Linearized Rusanov	AL	a_0	и	$u + a_0$
Centered-Total Rusanov	CT	0	u + a	u + a
Centered-Linearized Rusanov	CL	0	$u + a_0$	$u + a_0$
Centered-Advective Rusanov	CA	0	и	и

The Jacobian $(\Delta J = \frac{\partial \mathbf{F}_{NL}}{\partial \mathbf{q}})$ of the nonlinear flux $\mathbf{F}_{NL} = \mathbf{F}_T - \mathbf{F}_L$ is the following lower triangular matrix

$$\Delta J \equiv \frac{\partial \mathbf{F}_{NL}}{\partial \mathbf{q}} = \begin{pmatrix} 0 & 0 & 0 \\ (\gamma - 3)\frac{U^2}{2\rho^2} & (3 - \gamma)\frac{U}{\rho} & 0 \\ -\frac{U}{\rho^2} \left(\gamma E - (\gamma - 1)\frac{U^2}{\rho}\right) & \gamma \left(\frac{E}{\rho} - \frac{E_0}{\rho_0}\right) - \frac{3(\gamma - 1)U^2}{2\rho^2} & \gamma \frac{U}{\rho} \end{pmatrix}$$
(4.25)

with eigenvalues $\Delta \lambda = \lambda_{NL} = (0, \gamma u, (3 - \gamma)u)$. The key point is that even though the operator F_L has a linearized pressure term, the eigenvalues of ΔJ only have velocity components and no acoustic ones.³

This one-dimensional analysis shows that for the splitting of the linear and nonlinear terms to be consistent on both the continuous and discrete levels, the wave speed for \mathbf{F}_T should be taken as $\lambda_T = \lambda_L + \lambda_{NL}$, where $\lambda_L = 0$ when the implicit term is treated with a centered flux and $\lambda_L = a_0$ when the implicit term is treated with a linearized Rusanov flux; implications of various choices for the wavespeeds in the implicit and explicit components are presented in Appendix A. To keep the IMEX formulation for Runge-Kutta methods that we presented in [15] unchanged, we construct the nonlinear operator as $\mathbf{F}_{NL} = \mathbf{F}_T - \mathbf{F}_L$. Since the nonzero eigenvalues of \mathbf{F}_{NL} are all of relatively the same magnitude (i.e., $\mathcal{O}(u) = \mathcal{O}(\gamma u) = \mathcal{O}((3 - \gamma)u)$) we take $\lambda_{NL} = u$.

Table 1 presents the various flux combinations that can be used for the implicit and explicit dynamics. In the following sections, for all IMEX methods we call the AT flux that which uses the \underline{A} coustic Rusanov flux ($\lambda_L = a_0$) for the linear implicit part (IM) combined with \underline{T} otal Rusanov flux with $\lambda_T = u + a$ (and $\lambda_{NL} = u + a'$) accounting for both the advective and acoustic wave speeds for the nonlinear explicit (EX) part; in contrast, we call the CA flux that which uses the \underline{C} entered flux for the linear implicit (IM) part (i.e., $\lambda_L = 0$) combined with only the \underline{A} dvective Rusanov flux ($\lambda_T = \lambda_{NL} = u$) for the nonlinear explicit (EX) part (see **Case 1** in Appendix A). Insignificant differences are seen between the total Rusanov (AT Flux) and linearized Rusanov (AL Flux) when constructing the total flux \mathbf{F}_T , the reasons for which are presented in **Case 3** in Appendix A.

The *No-Schur* system can be solved using either the *AT* flux or the *CA* flux formulation, whereas the Schur system will only be solved using the *CA* flux approach. The Schur form can also be solved using the centered flux for the linear term $(\lambda_L = 0)$ and a Rusanov flux for the total/nonlinear terms (e.g., $\lambda_T = u + a_0$ as in the *CL* flux or $\lambda_T = u + a$ as in the *CT* flux); such a choice will subject the explicit time stepping to the CFL restriction due to the acoustic dynamics and negate the benefits of the IMEX methods (see **Case 2** in Appendix A). In Appendix B, we show why it is challenging to construct a Schur form when the linear term is treated with the Rusanov flux or any flux containing a jump term $(\lambda_L \neq 0)$.

4.3. Numerical stabilization

It should be noted that the complete system discretized using the CA flux is less dissipative than the system discretized using the AT flux, particularly for cases where $u \ll a$ (see **Case 1** in Appendix A).

This lack of dissipation can lead to spurious oscillations that can drive the system unstable, leading to divergence (i.e., blow-up). It was previously shown that artificial viscosity and sub-grid-scale (SGS) models can be used to stabilize the continuous Galerkin (which is also inherently non-dissipative) [26] and discontinuous Galerkin formulations [27]. For this reason, for some test problems we introduce an artificial diffusion operator (Laplacian) with uniform viscosity or SGS viscosity in the governing equations. The SGS viscosity used in this work is obtained using the Vreman model. It was shown in [26] that the Vreman model is able to stabilize the non-dissipative system while limiting the amount of artificial dissipation introduced. For brevity, the reader is referred to [28] for the complete description of the Vreman model. The second order diffusion/Laplacian operator ($\nabla \cdot (\mu \nabla \mathbf{U})$) is discretized using the weak-form local discontinuous Galerkin (LDG) method [25]. In the future, it will be interesting to explore entropy-stable methods in this context that do not require additional dissipation mechanisms (see, e.g., [22]).

4.4. Linear solvers

Before the linear solvers used to solve the resulting linear systems are discussed, the domain decomposition for multicore processing is briefly presented. As is often done in NWP or climate modeling (e.g., see [29]), the spherical manifold

³ Although we only show the results for one dimension, the two- and three-dimensional results bring us to the same conclusion.

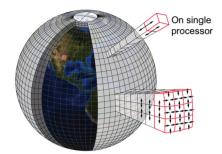


Fig. 1. The decomposition of the computational domain for multi-core processing. [35, Ch. 2].

is partitioned using either a cubed-sphere or icosahedral mesh and extruded along the radial direction (e.g., see [30]) such that each column of elements lies completely on a single processor (Fig. 1). This approach allows for the construction of all radial operators without requiring intra-processor communication, thereby, increasing scaling efficiency [30–33]. All graph partitioning and parallel mesh decomposition is performed using the *p4est* (for 3D-IMEX) and *p6est* (for 1D-IMEX) library [34].

The choice of the formulation (No-Schur or Schur) can directly influence the efficiency of the solver for the linear system. Since the vertical/radial degrees-of-freedom (DOF) N_z (i.e., DOF per mesh column) is much lower than the total DOF N_{tot} , and since the linear system in the 1D-IMEX formulation is independent of the adjacent columns, it is beneficial to use direct solvers. As stated previously, the Jacobian is constant in time and, therefore, can be factored once and stored for all subsequent linear solves. The linear system resulting from the 3D-IMEX discretization is often too large to store, too sparse to solve using direct methods, and requires intra-processor communication. For such systems, iterative Krylov methods such as GMRES, conjugate gradient (CG), or Bi-Conjugate Gradient-Stabilized (BiCGStab) are preferred. In this study we solve all 3D-IMEX systems using GMRES with a relative tolerance of 10^{-4} , and all 1D-IMEX systems, which have a banded form, using banded LU factorization from the LAPACK library. We show in the following sections that the Schur form of the IMEX formulations results in all real, positive eigenvalues, indicating a symmetric positive-definite (SPD) system. Therefore, the Schur systems can also be solved using conjugate gradient or Cholesky factorization. It should be mentioned that for the 1D-IMEX systems, Cholesky decomposition offers negligible gains over LU factorization since the matrix is constructed and decomposed (either LU or Cholesky) only once and reused for each subsequent linear solve (i.e. the forward reduction is done once while the backward substitution is performed at each s-stage of each time-step).

5. Test problems

This section briefly presents the four test problems used to analyze the semi-implicit discretizations. These problems are common benchmarks for verifying atmospheric models. For a more in-depth description of the test problems, readers are referred to their corresponding literature.

5.1. Rising thermal bubble

The rising thermal bubble (RTB) problem features the evolution of a warm bubble in a constant temperature ambient environment. The rising bubble deforms due to the shear stress caused by the velocity gradients. The test problem conditions [16] are similar to those of Robert [36] and features an ambient environment initially at rest and in hydrostatic balance. A warm air bubble with a potential temperature perturbation θ' is placed in the ambient atmosphere with uniform potential temperature θ_0 . The potential temperature perturbation distribution is defined by

$$\theta' = \begin{cases} 0, & \text{for } r > r_c \\ \frac{\theta_c}{2} \left(1 + \cos\left(\frac{\pi r}{r_c}\right) \right), & \text{for } r \leqslant r_c \end{cases}$$
 (5.1)

where $\theta_c=0.5$ K, $r=\sqrt{(x-x_c)^2+(z-z_c)^2}$ with the following constants: $\theta_0=300$ K, $r_c=250$ m, and $(x,z)\in[0,1000]^2$ m with $t\in[0,650]$ s and $(x_c,z_c)=(500,350)$ m. The no-flux boundary condition is imposed on all four walls of the boundary.

5.2. Density current

This two-dimensional test case, introduced in [37], is a standard benchmark for the development and verification of atmospheric models. It consists of a bubble of cold air descending to the ground in a neutrally stratified atmosphere. Upon hitting the lower boundary (all four boundaries are set as no-flux), the bubble develops Kelvin-Helmholtz shear instability rotors as it spreads laterally. This case is often solved by assuming a constant and uniform diffusion coefficient $\mu = 75 \text{ m}^2 \text{s}^{-1}$. The initial distribution of the potential temperature is given as

$$\theta' = \begin{cases} 0, & \text{for } r > r_c \\ \frac{\theta_c}{2} \left(1 + \cos \left(\frac{\pi r}{r_c} \right) \right), & \text{for } r \leqslant r_c \end{cases}$$
 (5.2)

where $\theta_c = -15$ K, $r = \sqrt{\left(\frac{x - x_c}{x_r}\right)^2 + \left(\frac{z - z_c}{z_r}\right)^2}$, and $r_c = 1$. The domain is defined as $(x, z) \in [0, 25600] \times [0, 6400]$ m with $t \in [0, 900]$ s and the center of the bubble as $(x_c, z_c) = (0, 3000)$ m with the size of the bubble taken to be $(x_r, z_r) = (4000, 2000)$ m.

5.3. Inertia gravity wave

The inertia gravity wave problem features the symmetric propagation of potential temperature perturbation in a periodic channel with a constant mean flow [38, Section 3]. The initial state of the atmosphere is taken to have a constant mean flow of $\bar{u} = 20$ m/s with the reference field defined as

$$T_0 = 250 \text{ K}, \quad p_0(z) = p_s e^{-\delta z}, \quad \rho_0(z) = \rho_s e^{-\delta z}$$

where $p_s = 10^5$ Pa, $\rho_s = \frac{p_s}{RT_0}$, and $\delta = \frac{g}{RT_0}$. The initial perturbation fields are given by

$$T'(x,z) = e^{\frac{\delta z}{2}} T_b(x,z), \quad p'(x,z) = 0, \quad \rho'(x,z) = e^{-\frac{\delta z}{2}} \rho_b(x,z)$$

with

$$T_b(x, z) = \Delta T \sin\left(\pi \frac{z}{h_c}\right) e^{-\left(\frac{x-x_c}{a_c}\right)^2}, \quad \rho_b(x, z) = -\rho_s \frac{T_b(x, z)}{T_0},$$

where $\Delta T = 0.001$ K, $h_c = 10$ km, $a_c = 5$ km, and $x_c = 100$ km. The domain is defined as $(x, z) \in [0, 300] \times [0, 10]$ km with $t \in [0, 2500]$ s. The left and right sides of the domain are taken to be periodic, while the top and bottom boundaries are treated as no-flux boundaries. This test problem (modified from [39]) admits an analytical solution (see [38, Section 2]).

5.4. Acoustic wave on a sphere

The acoustic wave propagation problem features an acoustic wave caused by a pressure disturbance traveling around the globe. This problem, proposed by Tomita and Satoh [40], is defined with a hydrostatically balanced state and an isothermal background potential temperature of $\theta_0 = 300$ K. The pressure perturbation is defined as $P' = f(\lambda, \phi)g(r)$, where

$$f(\lambda, \phi) = \begin{cases} 0, & \text{for } r > r_c \\ \frac{\Delta P}{2} \left(1 + \cos \left(\frac{\pi r}{r_c} \right) \right), & \text{for } r \leqslant r_c \end{cases}$$

$$g(r) = \sin \left(\frac{n_v \pi r}{r_T} \right),$$

$$(5.3)$$

where $\Delta P = 100$ Pa, $n_v = 1$, $r_c = r_e/3$ is one third of the radius of the earth $r_e = 6371$ km, and a model top of $r_T = 10$ km. The geodesic distance r is calculated as

$$r = r_e \cos^{-1} (\sin \phi_0 \sin \phi + \cos \phi_0 \cos \phi \cos (\lambda - \lambda_0))$$

where (λ_0, ϕ_0) is the origin of the acoustic wave. The no-flux boundary condition is imposed at the top and bottom boundaries. The spatial discretization for this problem often features high aspect ratio cells where the characteristic length of the elements in the radial direction is much smaller than the length in the horizontal direction.

5.5. Metrics of performance

To evaluate the performance of various semi-implicit formulations, we consider metrics that evaluate discrete conservation, efficiency, and accuracy of the methods. The methods' conservation properties can be quantified using relative mass and energy loss over simulation time t. The relative mass loss $\Delta \mathcal{M}$ and energy loss $\Delta \mathcal{E}$ are defined as

$$\Delta \mathcal{M} = \left| \frac{\mathcal{M}(t) - \mathcal{M}(0)}{\mathcal{M}(0)} \right| \quad \text{and} \quad \Delta \mathcal{E} = \left| \frac{\mathcal{E}(t) - \mathcal{E}(0)}{\mathcal{E}(0)} \right|$$

where $\mathcal{M}(t) = \int_{\Omega} \rho(\mathbf{x},t) d\Omega$ and $\mathcal{E}(t) = \int_{\Omega} E(\mathbf{x},t) d\Omega$. The accuracy is quantified using the mean-absolute error (MAE) relative to an analytical or reference solution \mathbf{q}^{ref} as

$$MAE = \frac{1}{N} \left\| \mathbf{q}^{ref} - \mathbf{q} \right\|_{1}$$

where N is the number of degrees-of-freedom (i.e., size of the mesh), \mathbf{q} is the set of prognostic variables for the equation set. Due to the lack of known analytical solutions for realistic problems in atmospheric sciences, the accuracy of the method is quantified relative to the solution obtained using an explicit scheme with a very small time-step (the small time-step is used in order to isolate the spatial error). The reference solutions are generated with a third-order, five-stage Runge-Kutta (RK35) scheme [41] with a Courant number of $CN \approx 0.0002$. We define the Courant number as

$$CN = \max\left(\frac{c\Delta t}{\Delta s}\right) \tag{5.4}$$

where $c = |u_n + a|$ is the characteristic wave speed, $u_n = \mathbf{u} \cdot \widehat{\mathbf{n}}$ is the velocity in the direction $\widehat{\mathbf{n}}$, a is the speed of sound and $\Delta s = \sqrt{\Delta x^2 + \Delta y^2 + \Delta z^2}$ is the characteristic length. Furthermore, a separate Courant number can be specified for each of the dynamics (acoustic and advection) by considering only the appropriate wave speeds in the expression for CN. Finally, the efficiency of the formulation is studied by analyzing the time-to-solution and number of iterations needed to converge to the solution. For each of the different IMEX formulations, we present the condition number $\kappa(A)$ defined as

$$\kappa(A) = \frac{\sigma_{\text{max}}}{\sigma_{\text{min}}},$$

where σ are the singular values of an $N \times N$ matrix A.

6. Results

6.1. Case 1: rising thermal bubble

The rising thermal bubble problem is solved on a 30 \times 30 mesh of fourth-order elements ($\Delta x = \Delta z = 8.3$ m) and is evolved until t = 650 s. It should be mentioned that all spatial resolutions Δx , Δy and Δz are defined as the average distance between the LGL nodes. The smaller dissipation in the CA flux IMEX discretization drives the simulation unstable whereas the added dissipation of the AT flux IMEX formulation keeps the model stable for the duration of the simulation. For a consistent comparison between all formulations, the Vreman SGS model, with $C_s = 0.21$, is used to stabilize all simulations of the rising thermal bubble.

The distribution of the potential temperature at t = 650 s using the various formulations of the two equation sets is shown in Fig. 2. The two equation sets result in similar distributions of the potential temperature. The 3D-IMEX and 1D-IMEX formulations, using the same flux scheme, also result in similar solutions. There is, however, a difference when comparing results obtained using the AT and CA flux formulations. The smaller magnitude of dissipation in the CA flux formulation leads to a warping (Rayleigh-Taylor instabilities) at the top of the bubble, whereas the solution obtained using the AT flux formulation shows a bubble that is relatively intact. A similar bubble breakdown behavior is seen when using a mixed-finite element, discontinuous Galerkin formulation (see [42]).

This difference can be due to the *jump* term in the Rusanov flux that is absent in the *CA* flux. The results of the Schur formulation show a similar distribution of the potential temperature as the *No-Schur* formulation for both equation sets and IMEX formulations (3D and 1D), indicating a correct derivation and implementation of the two forms. Based on the results in Fig. 2, we can conjecture that the CA flux is more appropriate in an IMEX setting since it is less dissipative (resulting in the Rayleigh-Taylor instabilities) while avoiding the production of unphysical grid imprinting (as exhibited in the AT flux).

Fig. 3 shows the eigenspectra of the flux formulations and semi-implicit discretizations. The linear system for the *No-Schur* case with the *CA* flux formulation, for both the 3D and 1D IMEX methods, has its eigenvalues along the imaginary axis. Similar spectra are seen for the IMEX formulations discretized using the continuous Galerkin method [24,43]. This is so because the continuous Galerkin formulation is inherently a spatially non-dissipative scheme. The centered flux is also non-dissipative and, therefore, the eigenvalues of the IMEX formulation using continuous Galerkin and centered-flux discontinuous Galerkin methods have similar eigenspectra (e.g., see Chs. 6 and 7 in [1]).

In the case of the IMEX methods with the dissipative *AT* flux, the eigenvalues seep into the real part of the complex plane. A similar behavior was presented for discontinuous Galerkin formulations in [44] where the non-dissipative Bassi-Rebay scheme [45] resulted in only imaginary eigenvalues while the dissipative Baumann-Oden scheme [46] resulted in eigenvalues spread across the complex plane.

The eigenvalues for both the 3D and 1D Schur formulations lie entirely along the real axis - as it should for an elliptic differential operator. Although not shown, the eigenspectra of Set2C are similar to those of Set3C for all formulations. The condition number as a function of Courant number for Set2C and Set3C are shown in Figs. 3e and 3f, respectively. The *No-Schur* linear system for Set2C is much better conditioned than that of Set3C; this can be attributed to $h_{2C} \ll h_{3C}$ in Eq. (3.8) and Eq. (3.17), respectively.

Comparing the Schur and *No-Schur* formulations for Set2C, the Schur formulation results in a system with a smaller condition number for lower Courant numbers, after which the condition number grows larger than that of the *No-Schur* system, particularly for the 1D-IMEX formulation. Comparing the Schur and *No-Schur* formulations for Set3C, the linear system for

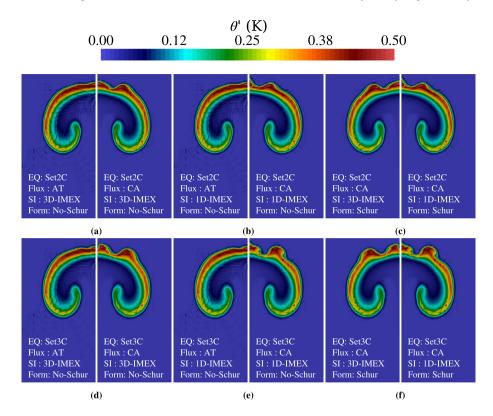


Fig. 2. Rising Thermal Bubble. The distribution of potential temperature perturbation at t = 650 s using Set2C with: (a) No-Schur 3D-IMEX, (b) No-Schur 1D-IMEX, and (c) Schur 3D vs. 1D IMEX; Set3C with: (d) No-Schur 3D-IMEX, (e) No-Schur 1D-IMEX, and (f) Schur 3D vs. 1D IMEX.

the Schur form is consistently better conditioned with condition numbers five orders of magnitude smaller than that of its *No-Schur* counterpart. This is expected to significantly decrease the number of iteration required to reach convergence, thereby reducing the time to solution. The condition number of the Schur forms for Set2C and Set3C are very similar; this is not surprising because, as we showed in [6] for the CG method, that regardless of which equation set is selected the resulting Schur form pressure equation should represent the same governing dynamics.

Table 2 shows the various performance metrics for the rising thermal bubble problem. It shows that the solution metrics (min., max., mass and energy loss) are all in good agreement for various formulations. The IMEX formulations using the AT flux are all in agreement with the solution obtained using the fully explicit third-order RK35 (also using the AT flux formulation). All equations sets and numerical formulations are able to conserve mass up to machine precision, whereas only Set3C is able to conserve energy up to machine precision. Set2C cannot conserve energy because the thermodynamic equation is written in terms of (density) potential temperature. Similar behaviors were observed when the two equation sets were analyzed using both the continuous Galerkin [24] and discontinuous Galerkin [16] methods. The minimum and maximum values of the thermodynamic variable are consistent across the different IMEX and flux formulations and the mean-absolute-error is low for all cases.

The results in Table 2 show that the 3D-IMEX *No-Schur CA* flux formulation is approximately 1.15 times faster than the *AT* flux formulation for both Set2C and Set3C. This is due to the better conditioning of the system from the *CA* flux formulation, thereby requiring fewer GMRES iteration to reach convergence. The Schur formulation is approximately 1.3 times faster than the *No-Schur AT* flux formulation due to the better conditioning of the system and a smaller overall size of the system, requiring fewer overall FLOPs (floating point operations); this conclusion holds for both Set2C and Set3C. Negligible differences are seen in the computing times between the two *No-Schur* flux formulations for 1D-IMEX. This is because the size of the Jacobian matrices is the same and, therefore, requires the same number of FLOPs for a direct solve (forward reduction and back substitution) of the linear system (standard LU decomposition requires $\mathcal{O}\left(\frac{2}{3}N^3\right)$ operations for *N* gridpoints, see, e.g., [47]). However, the linear system for the Schur formulation is much smaller (five times smaller) and therefore, the Schur formulation is approximately 4.8 times faster than the 1D-IMEX *No-Schur* discretization with *AT* flux. The Schur form for 1D-IMEX is slightly less than 5 times faster than the *No-Schur* form due to the additional FLOPs that are needed to extract the prognostic variables $(\rho, \mathbf{U}, \Theta, E)$ from the diagnostic variable (P). All timings are for simulations performed using 30 MPI ranks. Note that for this isotropic grid resolution, the 1D-IMEX method is at a disadvantage, which is evident by the large time-to-solution. This is due to the smaller time-steps required for stability of the explicitly treated horizontal acoustic dynamics. The time-step for the 1D-IMEX and fully-explicit schemes are the same; the additional

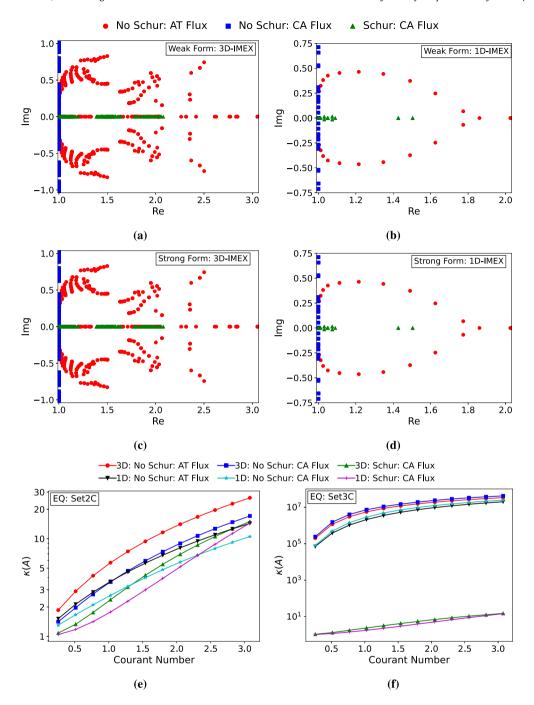


Fig. 3. Rising Thermal Bubble. The eigenspectra for Set3C for various flux formulations with Courant number ≈ 1 for variational weak form for: (a) 3D-IMEX and (b) 1D-IMEX; variational strong form for: (c) 3D-IMEX and (d) 1D-IMEX; condition number as a function of Courant number for (e) Set2C and (f) Set3C. (Similar eigenspectra are obtained with Set2C).

overhead of the 1D-IMEX over the fully explicit scheme comes from the added cost of solving a linear system at each timestep.

6.2. Case 2: density current

The density current test problem was solved on a 128×32 mesh of fourth-order elements ($\Delta x = \Delta z = 50$ m) and is evolved until t = 900 s. As previously mentioned, a uniform diffusion coefficient of $\mu = 75$ m² s⁻¹ is used to maintain stability.

Table 2Rising Thermal Bubble. Comparison of performance metrics for various formulations using Set2C and Set3C. The performance metrics include the mass $(\Delta \mathcal{M})$ and energy $(\Delta \mathcal{E})$ loss, mean absolute error (MAE) and time-to-solution (Sol. Time).

		Explicit (Rusanov)	3D-IMEX		1D-IMEX			
			No-Schur		Schur	No-Schur		Schur
			AT Flux	CA Flux	Schul	AT Flux	CA Flux	Schul
	Courant No.	0.1		1.5			0.1	
Set2C	Min θ' Max θ'	-1.45E-1 5.30E-1	-1.45E-1 5.30E-1	-2.94E-2 5.54E-1	-2.94E-2 5.54E-1	-1.45E-1 5.30E-1	-5.62E-2 5.43E-1	-5.62E-2 5.43E-1
	$\Delta \mathcal{M} \ \Delta \mathcal{E}$	0.0E0 2.50E-8	0.0E0 2.49E-8	0.0E0 2.47E-8	0.0E0 2.47E-8	2.09E-16 2.49E-8	0.0E0 2.26E-8	2.09E-16 2.26E-8
	MAE	-	1.63E-5	2.37E-05	2.08E-05	9.91E-7	2.51E-4	2.68E-4
	Sol. Time (s)	4384	319	275	250	15499	15580	3231
Set3C	Min $ heta'$ Max $ heta'$	-5.29E-2 5.20E-1	-5.29E-2 5.20E-1	-3.45E-2 5.19E-1	-3.45E-2 5.19E-1	-5.29E-2 5.20E-1	-2.74E-2 5.08E-1	-2.67E-2 5.08E-1
	$\Delta \mathcal{M} \ \Delta \mathcal{E}$	0.0E0 0.0E0	0.0E0 0.0E0	0.0E0 0.0E0	0.0E0 0.0E0	2.09E-16 0.0E0	2.09E-16 0.0E0	2.09E-16 0.0E0
	MAE	-	1.60E-5	2.18E-05	1.90E-05	6.00E-7	3.00E-4	3.20E-4
	Sol. Time (s)	4435	322	274	248	15484	15443	3187

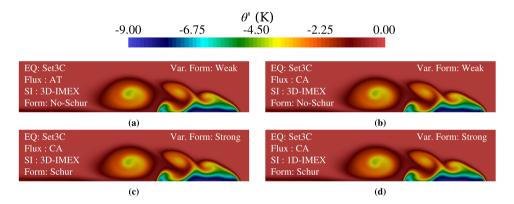


Fig. 4. Density Current. The distribution of potential temperature perturbation at t = 900 s using Set3C with variational weak *No-Schur* form and: (a) 3D-IMEX with *AT* flux, (b) 3D-IMEX with *CA* flux; and variational strong Schur form for: (c) 3D IMEX, (d) 1D-IMEX.

Fig. 4 shows the distribution of the potential temperature obtained using the two variational forms and different flux schemes for the 3D-IMEX and 1D-IMEX formulations. Comparing Fig. 4, it can be seen that the different flux formulations have little to no effect on the solution. This is because the dissipation introduced by the artificial viscous operator is significantly greater than the dissipation from the numerical fluxes. The potential temperature obtained using 3D-IMEX and 1D-IMEX in both weak and strong forms are all similar.

Implicit time integration (whether treating the linear or nonlinear terms) introduces greater dissipation into the system than explicit methods [48]. Therefore, the accuracy of the IMEX formulations must be investigated. The accuracy is quantified using the mean absolute error (MAE) of the difference between the solution resulting from the IMEX discretization and some reference solution. The reference solution is computed using a third-order, five-stage, explicit Runge-Kutta method (RK35) using a very small time step (CN \approx 0.0002) and the AT flux formulation, and is compared at t=100 s. The IMEX schemes are solved using a two-stage, second-order ARK (ARK2) [15], a three-stage, third-order ARK (ARK3) and a four-stage, fourth-order ARK (ARK4) [49] schemes. The different Courant numbers for the 3D-IMEX cases are obtained by modifying the time step appropriately. For the 1D-IMEX cases, the Courant numbers are modified by changing the resolution of the vertical discretization, so that the influence of only vertical/radial operators is affected. Since the grids for the 3D-IMEX and 1D-IMEX cases are different (even for the same Courant number), the errors between the two cannot be compared fairly.

Fig. 5 shows the accuracy of the IMEX formulations for the *No-Schur* form of Set3C. The IMEX methods using the *AT* flux formulation exhibit the expected order of convergence relative to the reference solution obtained using RK35 with *AT* flux. For smaller Courant numbers, the IMEX methods using the *CA* flux formulation differ significantly from the reference solution but are in good agreement for larger Courant numbers. This is due to the difference in solution caused by the difference in the amount of dissipation from the numerical flux. For larger Courant numbers, the dissipation from the time integration offsets this difference in dissipation from the numerical fluxes, thereby leading to similar results as the total *AT* flux. In

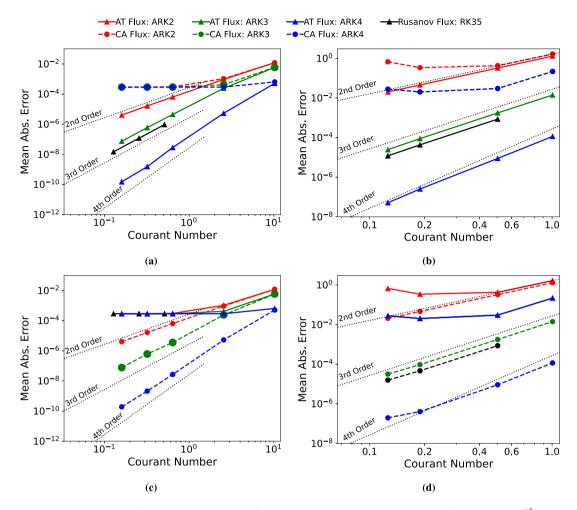


Fig. 5. Density Current. The accuracy of the IMEX discretizations and ARK time integrators for Set3C relative to a reference solution \mathbf{q}^{ref} computed using RK35 (with *AT* flux and Courant number ≈ 0.0002) for: (a) 3D-IMEX and (b) 1D-IMEX formulations; and relative to a reference solution \mathbf{q}^{ref} computed using RK35 (with *CA* flux and Courant number ≈ 0.0002) for: (c) 3D-IMEX and (d) 1D-IMEX formulations. (Note: the results for 1D-IMEX with *CA* flux are similar for both ARK3 and ARK4).

addition, the IMEX methods using the CA flux also exhibit the expected order of convergence when the reference solution from the RK35 is obtained using the CA flux formulation (Figs. 5c and 5d).

6.3. Case 3: inertia gravity waves

The inertia gravity wave problem was solved on a 96×32 mesh of fourth-order elements ($\Delta x = 780$ m, $\Delta z = 78$ m) and is evolved until t = 2500 s; at this final time, the initial perturbation will be exactly at the midpoint of the domain. The distribution of the potential temperature along the horizontal center line of the domain (at z = 5 km) obtained using various formulations of the strong form, 3D-IMEX discretization is shown in Fig. 6. No appreciable differences are seen between the distributions from the different discretizations and equation sets. To adhere to the analytical solution, the test case is solved without any artificial diffusion.

We now compare the efficiency of the various IMEX formulations. The average number of solver iterations per time step and the time-to-solution (wallclock time) for different Courant numbers are shown in Figs. 7a and 7b, respectively. To not be constrained by the CFL condition on the explicit treatment of the advective dynamics, the constant mean flow is set to $\bar{u}=0$ m/s (see Sec. 5.3). All analyses were performed on 40 cores with ARK2. Therefore, the number of iterations reported per time step are for two ARK stages. The efficiency results presented are shown for both Set2C and Set3C.

For the *No-Schur* system, the *CA* flux formulation requires fewer iterations to reach convergence. This is due to the better conditioning of the system and is consistent with the results for the rising thermal bubble test (see Sec. 6.1). As expected, a larger number of iterations are required with increasing Courant numbers. The Schur system requires fewer iterations than the *No-Schur* system for the same Courant number, due to the better conditioning of the linear system. Although the linear solvers are not the focus of this study, it was previously demonstrated (Fig. 3) that the Schur form results in a symmetric

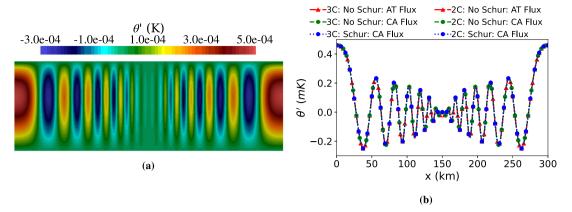


Fig. 6. Inertia Gravity Wave. The distribution of potential temperature perturbation at t = 2500 s using 3D-IMEX and: (a) Set2C in No-Schur form with *AT* flux and (b) distribution along the center line for each combination of equation set, numerical flux, and IMEX form (note the y-axis is in milli-Kelvin).

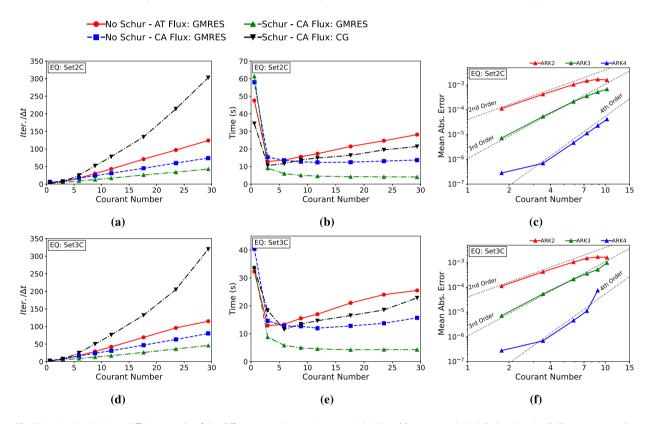


Fig. 7. Inertia Gravity Wave. Efficiency study of the different equation sets (top row is Set2C and bottom row is Set3C) showing: (a, d) The average number of iterations required to solve the linear system at each time step, (b, e) time required to compute 500 seconds and (c, f) convergence of the various IMEX time integrators with *No-Schur* form and *AT* flux formulation (reference solution was taken to be the analytical solution of Baldauf and Brdar [38]).

positive definite matrix and can be solved using conjugate gradient (CG). Fig. 7a shows that CG requires significantly more iterations than GMRES, particularly for larger Courant numbers. Although CG requires more iterations per time step, the number of FLOPs, storage requirements (no need to store Krylov vectors), and inter-processor communication per iteration is much lower than GMRES (when using modified Gram-Schmidt for the Arnoldi iteration). Hence, the CG solver is faster for smaller Courant numbers. Previous works have provided conditions for superlinear convergence of the CG [50] and GMRES [51] solvers based on the Ritz values, however, such analysis is beyond the scope of this study but of interest for future work. Both equation sets exhibit similar iteration counts.

Comparing the wallclock time required to reach the simulation time of 500 s, Figs. 7b and 7e show that there exists an optimal time step that yields the fastest time-to-solution, after which the efficiency degrades; this behavior can be delayed with preconditioners but we save this for future work (see Sec. 7). For smaller time steps where fewer iterations (one to

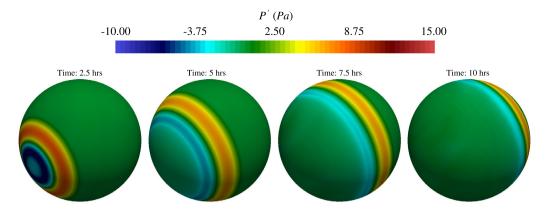


Fig. 8. Acoustic Wave on Sphere. Pressure perturbation at different times obtained using 1D-IMEX discretization of Set3C with the No-Schur AT flux formulation.

two) are needed to solve the linear system, the *No-Schur* formulations yield a faster time-to-solution; similar behavior was seen for continuous Galerkin IMEX formulations [24]. This is due to the computational cost of extracting the prognostic variables from the diagnostic variables being larger than the cost of performing few iterations of the linear solver. At the optimal time step ($CN \approx 3$), the Schur formulation begins to yield faster time-to-solution than the *No-Schur* formulations. For Courant numbers smaller than this optimal value, the Schur form with the conjugate gradient solver provides a faster time-to-solution than the *No-Schur* forms. For smaller Courant numbers (i.e., for more number of time steps), Set3C yields a much faster time-to-solution than Set2C. This can be attributed to the increased computational cost of the equation of state for Set2C over that of Set3C. For larger Courant numbers (i.e., fewer time steps), the speed-up of Set3C over Set2C is smaller (but Set3C is still faster) due to fewer evaluations of the equation of state.

Figs. 7c and 7f show results of the convergence study for the IMEX formulation with the AT flux. The reference solution was taken to be the analytical solution of Baldauf and Brdar [38] at t = 2500 s. All methods exhibit the expected order of convergence. A similar breakdown in the order of convergence was seen in [38] and [22]. Fig. 7f shows the breakdown of ARK4 at Set3C at higher Courant numbers; this behavior is described in [22,38].

6.4. Case 4: acoustic wave on a sphere

The acoustic wave on a sphere (AWS) problem is solved on a cubed-sphere grid with a resolution of $\Delta h = 200$ km and $\Delta v = 400$ m, in the horizontal and vertical directions, respectively. The simulation is evolved until t = 36000s with a time step of $\Delta t = 25$ s corresponding to a Courant number ≈ 37 . The high aspect ratio of such resolutions ($\Delta h/\Delta v = 500$) increases the stiffness due to the vertical resolution and dynamics. This stiffness leads to a drastic increase in the number of iterations required for convergence of the 3D-IMEX system. Therefore, the 1D-IMEX discretization with direct solver is ideal for such applications. Fig. 8 shows snapshots at different times for the acoustic wave traveling around the sphere. The solutions are obtained using the 1D-IMEX discretization of Set3C with the AT flux formulation and ARK2 time integration. No significant differences are seen between the AT and CA flux formulations.

The efficiency of the 1D-IMEX over the 3D-IMEX for global problems (i.e., problems on the sphere) is investigated. The number of iterations and time-to-solution for different Courant numbers is investigated for all formulations presented. Since the linear system in the 1D-IMEX formulation is solved using a direct method (banded LU decomposition from the LAPACK library), studying its performance by simply changing the time steps to obtain the required Courant number is not appropriate (since the matrix size and FLOPs will remain the same and, therefore, the performance will increase linearly with the increase in time step).

For this reason, the Courant number in the efficiency study of the 1D-IMEX formulations is updated by modifying the vertical resolution. That is, higher resolutions yield larger Courant numbers and larger system matrices. For the 3D-IMEX case the grid resolution is kept constant; instead the required Courant number is attained by adapting the time step. This keeps the size of the linear system constant between different Courant numbers but with different stiffness, thereby requiring different number of iterations to reach convergence. All analyses are performed with Set3C and the ARK2 scheme using 10 MPI ranks.

Fig. 9 shows the results of the convergence study for 1D-IMEX (with *AT* flux), average number of iterations (for the 3D-IMEX case) and the time required to evolve the system for 10 time steps. The Courant number presented is the Courant number of the acoustic term with vertical grid spacing (i.e., $a = \sqrt{\frac{\gamma P}{\rho}}$ and $\Delta s = \Delta v$ in Eq. (5.4)). Comparing Figs. 9a and 9b shows that the 1D-IMEX methods are more than an order of magnitude faster than the 3D-IMEX methods for such high aspect ratio problems. This is because the 3D-IMEX methods are restricted by the stiffness due to the vertical discretization and dynamics. The Schur forms for both 3D-IMEX and 1D-IMEX are significantly faster than their *No-Schur* forms. In the case of 1D-IMEX, the timings of the two *No-Schur* forms coincide as expected, since the two linear systems are of the same

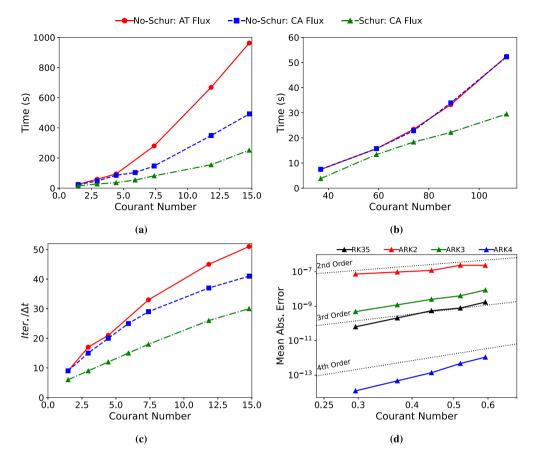


Fig. 9. Acoustic Wave on Sphere. The accuracy and efficiency of the IMEX discretizations with Set3C showing: the time required to compute 10 time steps as a function of Courant number using: (a) 3D-IMEX and (b) 1D-IMEX; (c) average number of iterations per time step for 3D-IMEX and (d) convergence of the 1D-IMEX scheme (reference solution was obtained using RK35 and $CN \approx 0.001$).

size and, because we are using direct solvers, the condition number of the matrices play no role with respect to time-to-solution. Although the 3D-IMEX formulations allow for a larger <u>stable</u> time step (since acoustic terms in all directions are handled implicitly), the larger computing/simulation times often render these methods unusable for such global models. Hence, the 1D-IMEX formulations are ideally suited for high aspect ratio meshes as they allow for much larger <u>usable</u> time steps while greatly reducing the time-to-solution.

Fig. 9c shows the average number of iterations per time step required to reach convergence. The Schur form of the 3D-IMEX discretization requires significantly fewer iterations to reach convergence; this can be attributed to the better conditioning (smaller condition number) as well as the location of the eigenvalues of the Schur system (Krylov methods prefer the eigenvalues near the real axis). The *No-Schur* form with AT and CA fluxes require similar number of iterations for smaller Courant numbers (CN < 5). For larger Courant numbers the CA flux formulation requires fewer iterations than the CA flux, thereby yielding a much faster time-to-solution. Due to the computational cost of the 3D-IMEX formulations for such highly anisotropic grids, the convergence analysis is only performed for 1D-IMEX formulations. These 1D-IMEX schemes (Fig. 9d) are shown to follow the expected order of convergence.

7. Conclusion

We present and analyze Schur complement IMplicit-EXplicit formulations for discontinuous Galerkin discretizations of different forms of the Euler equations. Two different IMEX methods that resolve the stiffness due to the governing dynamics and geometric discretization are shown. The 3D-IMEX method, which treats the acoustical stiffness in all directions, allows for significantly larger time steps and yields a faster time-to-solution for problems with mesh aspect ratios closer to one. In the case of global models, where the aspect ratio is much larger, the 1D-IMEX method is more efficient and more than an order of magnitude faster than 3D-IMEX. For each IMEX formulation and governing equation set, efficient Schur complements are derived to reduce the size of the linear system from a system of five variables to a single equation for pressure. The derivation and consequences of the different numerical fluxes used in these Schur complements are discussed. The Schur formulations of these IMEX methods result in a smaller, better conditioned system, and, therefore, require few iterations to reach convergence. Moreover, the resulting linear system is shown to have purely real eigenvalues confirming

the ellipticity of the resulting discrete form of the Helmholtz-like differential operator; this last point is important because it now opens the door to developing preconditioners that will further improve the performance of this method (e.g., multigrid methods). The ARK methods of different orders also yield the expected theoretical rates of convergence for the test cases considered. These different formulations of the IMEX methods are shown to be efficient and accurate for both mesoscale (box) and global (sphere) models.

CRediT authorship contribution statement

Sohail Reddy: Conceptualization, Formal analysis, Methodology, Software, Writing – original draft, Writing – review & editing. **Maciej Waruszewski:** Conceptualization, Methodology. **Felipe A.V. de Braganca Alves:** Software. **Francis X. Giraldo:** Conceptualization, Formal analysis, Methodology, Project administration, Software, Writing – original draft, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Data will be made available on request.

Acknowledgements

The authors gratefully acknowledge the support by the generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program. We are also grateful for funding provided by the Office of Naval Research under grant # N0001419WX00721 and National Science Foundation under grant AGS-1835881. This work was performed in part when S. Reddy, M. Waruszewski, and F. A.V. B. Alves held National Academy of Sciences' National Research Council (NRC) Fellowships at the Naval Postgraduate School. FXG would like to thank Prof. Maria Lukacova-Medvidova for our conversation on this topic at the Moist Processes in the Atmosphere meeting in Oberwolfach which formed the initial impetus for this work.

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

Appendix A. Effects of wave speeds on flux formulations for IMEX discretizations

Here, we present conditions on the wave speed of the total flux, given the treatment of the linear flux so that we obtain the appropriate eigenvalue/wave speed for the nonlinear flux. Computing the nonlinear flux in Eq. (4.20) as $\mathbf{F}_{NL} = \mathbf{F}_T - \mathbf{F}_L$ and writing $\lambda_{NL} = \lambda_T - \lambda_L$, let us consider the following cases:

Case 1 Treating the implicit system with a centered flux ($\lambda_L = 0$), letting $\lambda_{NL} = u$, the nonlinear (explicit) operator in Eq. (4.20) can be written as

$$\int_{\Omega} \psi \nabla \cdot \mathbf{F}_{NL} d\Omega = \int_{\Gamma} \psi \left(\{ \{ \mathbf{F}_{T} \} \} - \frac{u}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma - \int_{\Omega} \nabla \psi \cdot \mathbf{F}_{T} d\Omega
- \int_{\Gamma} \psi \left\{ \{ \mathbf{F}_{L} \} \right\} \cdot \widehat{\mathbf{n}} d\Gamma + \int_{\Omega} \nabla \psi \cdot \mathbf{F}_{L} d\Omega$$
(A.1)

where we obtain a form that treats the total flux (and the nonlinear flux) using Rusanov with the correct (advective) wave speed of the nonlinear flux Jacobian. This shows that when the implicit term is treated with centered fluxes, the total flux \mathbf{F}_T must be treated with a Rusanov flux with wave speed u. This analysis shows why the splitting used in [52,53] works in a semi-Lagrangian setting since the nonlinear (explicit) operator is treated in a semi-Lagrangian approach whereby only the advective speed is used in the method of characteristics and the remainder of the terms used centered fluxes, which is what we have analyzed here. Note that the centered flux ($\{\{\cdot\}\}$) is non-dissipative. Hence, it is clear that this approach results in a complete system that is less dissipative when $u \ll a$.

Case 2 Treating the total flux with a linearized Rusanov flux $(\lambda_T = u + a_0)$ with $a_0 = \sqrt{\frac{\gamma P_0}{\rho_0}}$, we get

$$\int_{\Omega} \psi \nabla \cdot \mathbf{F}_{NL} d\Omega = \int_{\Gamma} \psi \left(\{ \{ \mathbf{F}_{T} \} \} - \left(\frac{u + a_{0}}{2} \right) [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma - \int_{\Omega} \nabla \psi \cdot \mathbf{F}_{T} d\Omega
- \int_{\Gamma} \psi \left(\{ \{ \mathbf{F}_{L} \} \} - \frac{\lambda_{L}}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma + \int_{\Omega} \nabla \psi \cdot \mathbf{F}_{L} d\Omega$$
(A.2)

where $\lambda_{NL}=u$ only when $\lambda_L=a_0$. This shows that when the implicit term is treated with a linearized Rusanov flux with wave speed λ_L , the total flux \mathbf{F}_T should be treated with a linearized Rusanov flux with $\lambda_T=u+\lambda_L$. When $\lambda_L=0$ (as it is in the Schur form), then $\lambda_{NL}=u+a_0$ is the wave speed of the terms being treated explicitly. Hence, the explicit time integration is still subject to the CFL restrictions due to the acoustic dynamics which requires using a smaller time steps.

Case 3 Treating the total flux with the total nonlinear Rusanov flux ($\lambda_T = u + a$) we get

$$\int_{\Omega} \psi \nabla \cdot \mathbf{F}_{NL} d\Omega = \int_{\Gamma} \psi \left(\{ \{ \mathbf{F} \} \} - \frac{u+a}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma - \int_{\Omega} \nabla \psi \cdot \mathbf{F} d\Omega
- \int_{\Gamma} \psi \left(\{ \{ \mathbf{F}_{L} \} \} - \frac{\lambda_{L}}{2} [\mathbf{q}] \widehat{\mathbf{n}} \right) \cdot \widehat{\mathbf{n}} d\Gamma + \int_{\Omega} \nabla \psi \cdot \mathbf{F}_{L} d\Omega.$$
(A.3)

This appears valid but the pressure P appears nonlinearly in a. When $\lambda_L=0$ (i.e., implicit term is treated with a centered flux), then the wave speed of the nonlinear term is $\lambda_{NL}=u+a\gg u$ which is not an eigenvalue of the nonlinear flux Jacobian and can lead to a large portion of the eigenspectra to lie outside the region of stability of the explicit time integrator thereby requiring a smaller time step that is dictated by the acoustic dynamics. When $\lambda_L=a_0$ (i.e., implicit term is treated with a linearized Rusanov flux), then $\lambda_{NL}=u+a-a_0=u+a'$ is still not an eigenvalue of the nonlinear operator (as shown above) but yields eigenvalues nearer to it. Linearizing a about P_0 and ρ_0 , we obtain an expression for a'.

$$a \approx \underbrace{\sqrt{\frac{\gamma P_0}{\rho_0}}}_{a_0} + \underbrace{\frac{1}{2} \left(\sqrt{\frac{\gamma}{P_0 \rho_0}} P' - \sqrt{\frac{\gamma P_0}{\rho_0^3}} \rho' \right)}_{a'} \right)$$
$$\approx a_0 + \frac{a_0}{2} \left(\frac{P'}{P_0} - \frac{\rho'}{\rho_0} \right)$$

where $P' \ll P_0$ and $\rho' \ll \rho_0$. We can then perform a scale analysis with $\mathcal{O}(\gamma) = \mathcal{O}(\rho_0) = \mathcal{O}(1)$, $\mathcal{O}(P_0) = \mathcal{O}(10^5)$. Fig. 8 shows that $\mathcal{O}(P') = \mathcal{O}(10)$, and, numerically we can show that, $\mathcal{O}(\rho') < \mathcal{O}(10^{-4})$. Therefore $\mathcal{O}(a') < \mathcal{O}(10^{-1}) < \mathcal{O}(u)$, hence such a choice of the numerical flux is not problematic for low Mach number applications.

Appendix B. The necessity for special choices on the flux to construct the Schur complement

To understand the challenge of constructing the Schur complement for the DG method, let us consider the linearized shallow water equations

$$\frac{\partial \varphi}{\partial t} + \nabla \cdot \mathbf{U} = 0$$

$$\frac{\partial \mathbf{U}}{\partial t} + \nabla \varphi = 0$$
(B.1)

where (φ, \mathbf{U}) represent the geopotential height and momentum, respectively. Applying the element-based Galerkin method to Eq. (B.1) (multiplying by a test function and integrating) yields

$$\int_{\Omega_{e}} \psi \frac{\partial \varphi_{N}}{\partial t} d\Omega_{e} + \int_{\Omega_{e}} \psi \nabla \cdot \mathbf{U}_{N} d\Omega_{e} = 0$$

$$\int_{\Omega_{e}} \psi \frac{\partial \mathbf{U}_{N}}{\partial t} d\Omega_{e} + \int_{\Omega_{e}} \psi \nabla \varphi_{N} d\Omega_{e} = 0.$$
(B.2)

Next, note that the second terms for each equation can be written as follows

$$\int_{\Omega_{e}} \psi \nabla \cdot \mathbf{U}_{N} d\Omega_{e} = \int_{\Gamma_{e}} \psi \widetilde{\mathbf{U}}_{N} \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot \mathbf{U}_{N}^{(e)} d\Omega_{e}$$

$$\int_{\Omega_{e}} \psi \nabla \varphi_{N} d\Omega_{e} = \int_{\Gamma_{e}} \psi \widetilde{(\varphi \mathbf{I}_{d})}_{N} \cdot \widehat{\mathbf{n}} d\Gamma_{e} - \int_{\Omega_{e}} \nabla \psi \cdot (\varphi \mathbf{I}_{d})_{N}^{(e)} d\Omega_{e}$$
(B.3)

where \mathbf{I}_d is the rank-d identity matrix, where d denotes the spatial dimension of the problem. With little loss in generality let us write the numerical flux or the first terms on the right-hand side of Eq. (B.3) in the following form

$$\widetilde{\mathbf{U}}_{N} = \{\{\mathbf{U}_{N}\}\} - \frac{\lambda}{2} \llbracket \varphi_{N} \rrbracket \widehat{\mathbf{n}}
\widetilde{(\varphi \mathbf{I}_{d})}_{N} = \{\{\varphi_{N} \mathbf{I}_{d}\}\} - \frac{\lambda}{2} \llbracket \mathbf{U}_{N} \rrbracket \widehat{\mathbf{n}}
(B.4)$$

where $\{\{\cdot\}\}\$ and $[\cdot]$ represent the average and jump operators, $\widehat{\mathbf{n}}$ is the outward pointing unit normal vector of the face Γ_e , and λ is the eigenvalue of the flux Jacobian. With this in mind, we can rewrite Eq. (B.3) in the following matrix form (see, e.g., [1])

$$\int_{\Omega_{e}} \psi \nabla \cdot \mathbf{U}_{N} d\Omega_{e} = \left(\mathcal{C}^{(e,k)} \right)^{T} \mathbf{U}_{N} - \mathcal{J}^{(e,k)} \varphi_{N} - \left(\widetilde{\mathbf{D}}^{(e)} \right)^{T} \mathbf{U}_{N}
\int_{\Omega_{e}} \psi \nabla \varphi_{N} d\Omega_{e} = \mathcal{C}^{(e,k)} \varphi_{N} - \mathcal{J}^{(e,k)} \mathbf{U}_{N} - \widetilde{\mathbf{D}}^{(e)} \varphi_{N}$$
(B.5)

where $C^{(e,k)}$, $\mathcal{J}^{(e,k)}$, and $\widetilde{\mathbf{D}}^{(e)}$ represent the centered, jump, and weak form differentiation matrices. Note that $C^{(e,k)}$ and $\widetilde{\mathbf{D}}^{(e)}$ are vector matrices whereas $\mathcal{J}^{(e,k)}$ is a scalar matrix. Further note that the matrix superscript (e,k) denotes that the matrices of the element e also need the contribution of its neighbor k. Matrices with superscript (e) are purely local element-wise quantities.

Taking all of this into account and assuming a backward Euler time-integration, we can now write Eq. (B.2) as follows

$$M^{(e)}\varphi + \Delta t \left[\left(\mathcal{C}^{(e,k)} \right)^{\mathcal{T}} \mathbf{U} - \mathcal{J}^{(e,k)}\varphi - \left(\widetilde{\boldsymbol{D}}^{(e)} \right)^{\mathcal{T}} \mathbf{U} \right] = R_{\varphi}$$

$$M^{(e)}\mathbf{U} + \Delta t \left[\mathcal{C}^{(e,k)}\varphi - \mathcal{J}^{(e,k)}\mathbf{U} - \widetilde{\boldsymbol{D}}^{(e)}\varphi \right] = R_{\mathbf{U}}$$
(B.6)

where terms on the left-hand side are defined at the new time level n + 1 and terms on the right are at the current time level n; we also have dropped the subscript N for convenience.

Letting

$$\widehat{M} = M^{(e)} - \Delta t \mathcal{J}^{(e,k)}$$

$$\widehat{\mathbf{D}} = \mathcal{C}^{(e,k)} - \widetilde{\mathbf{D}}^{(e)}$$
(B.7)

allows us to write Eq. (B.6) as follows

$$\widehat{M}\varphi + \Delta t \widehat{\mathbf{D}}^{T} \mathbf{U} = R_{\varphi}$$

$$\widehat{M} \mathbf{U} + \Delta t \widehat{\mathbf{D}}\varphi = R_{\mathbf{U}}.$$
(B.8)

Finally, applying a block LU decomposition yields the Schur complement

$$\widehat{M}\varphi - \Delta t^2 \widehat{\boldsymbol{D}}^T \widehat{M}^{-1} \widehat{\boldsymbol{D}}\varphi = R_\varphi - \Delta t \widehat{\boldsymbol{D}} \widehat{M}^{-1} R_{\mathbf{U}}. \tag{B.9}$$

Looking at Eq. (B.7) and Eq. (B.9) we note that if $\mathcal{J}^{(e,k)}$ is not empty, then \widehat{M}^{-1} is not block diagonal and thereby requires a global solution. However, if $\mathcal{J}^{(e,k)}$ is empty, which it will be for the continuous Galerkin method and for the DG method using centered fluxes, then \widehat{M}^{-1} is block diagonal and trivial to invert - in fact, if we use inexact integration (as we do throughout this paper) it will be diagonal. Therefore, to make the Schur complement feasible requires using centered fluxes for the implicit linear operators. We say feasible because inverting \widehat{M} will increase the cost of the Schur complement by a factor dependent on how often it appears in the Schur form. To give you a sense of what we mean, in Eq. (B.9), the inverse of \widehat{M} appears once on the left and also on the right. Therefore, for non-block-diagonal \widehat{M} , Eq. (B.9) would require three global linear solves.

Appendix C. Discrete formulations of the Schur complement for the IMEX discretizations

We now replace the continuous spatial derivatives of the Schur form with their discrete representations. Defining the weak form differentiation matrix $\tilde{\boldsymbol{D}}^{(e)}$ as

$$\widetilde{\boldsymbol{D}}_{ij}^{(e)} = \int\limits_{\Omega} \nabla \psi_i \psi_j \, d\Omega_e,$$

the strong form differentiation matrix as

$$\boldsymbol{D}_{ij}^{(e)} = \int_{\Omega_e} \psi_i \nabla \psi_j \, d\Omega_e,$$

and the centered flux matrix $\mathcal{C}^{(e,k)}$ as

$$C_{ij}^{(e,k)} = \int_{\Gamma_e} \psi_i \psi_j \widehat{\mathbf{n}}^{(e,k)} d\Gamma_e,$$

for element e with neighboring element k with an outward pointing normal vector $\widehat{\mathbf{n}}^{(e,k)}$ from e to k. Then the weak variational form for $\nabla \cdot \mathbf{Q}$ can be written as

$$\int_{\Omega_e} \psi \nabla \cdot \mathbf{Q} \, d\Omega_e = \int_{\Gamma_e} \psi \mathbf{Q} \cdot \hat{\mathbf{n}} \, d\Gamma_e - \int_{\Omega_e} \nabla \psi \cdot \mathbf{Q} \, d\Omega_e = (\mathcal{C}^{(e,k)})^T \mathbf{Q} - (\widetilde{\mathbf{D}}^{(e)})^T \mathbf{Q} = \widehat{\mathbf{D}}^T \mathbf{Q}$$

where $\widehat{\boldsymbol{D}} = (\boldsymbol{\mathcal{C}}^{(e,k)} - \widetilde{\boldsymbol{D}}^{(e)})$. Letting $\widecheck{\boldsymbol{D}} = \widehat{\boldsymbol{M}}^{-1} \widehat{\boldsymbol{D}}^{T}$, where

$$\widehat{M}_{ij}^e = M_{ij}^e - \alpha \mathcal{J}_{ij}^{(e,k)}$$

is the augmented mass matrix and $\mathcal{J}^{(e,k)}$ is the jump term (that is zero for centered flux and continuous fields), $\widehat{M}^{-1}\nabla \cdot \mathbf{Q} = \mathbf{D}^T\mathbf{Q}$ and $\widehat{M}^{-1}\nabla \mathbf{Q} = \mathbf{D}^T\mathbf{Q}$, the discrete weak Schur form for Set2C (Eq. (3.15)) becomes

$$P_{tt} - \alpha^{2} \frac{\gamma P_{0}}{\Theta_{0}} \mathbf{D}^{T} \left(\left[h_{2C} \left\{ \mathbf{A}^{-1} \left(\mathbf{D} P_{tt} + \frac{\mathbf{D} \phi}{h_{2C}} \frac{\Theta_{0} P_{tt}}{\gamma P_{0}} \right) \right\} \right] \right)$$

$$= \frac{\gamma P_{0}}{\Theta_{0}} \widehat{\Theta} - \alpha \frac{\gamma P_{0}}{\Theta_{0}} \mathbf{D}^{T} \left(\left[h_{2C} \left\{ \mathbf{A}^{-1} \left(\widehat{\mathbf{U}} - \alpha \left[\widehat{\rho} \mathbf{D} \phi - \frac{\mathbf{D} \phi}{h_{2C}} \widehat{\Theta} \right] \right) \right\} \right] \right),$$
(C.1)

where α is the IMEX coefficient defined in Sec. 3.1. The discrete strong Schur form can be obtained by expanding the divergence operators in Eq. (3.15) using the product rule and substituting \boldsymbol{D} and $\boldsymbol{D}^{\mathcal{T}}$ appropriately. Similarly, the discrete weak Schur form for Set3C (Eq. (3.26)) can be written as

$$P_{tt} - \alpha^{2} (\gamma - 1) \left(\mathbf{\tilde{D}}^{T} \left[h_{3C} \left\{ \mathbf{A}^{-1} \left(\mathbf{\tilde{D}} P_{tt} + \frac{\mathbf{\tilde{D}} \phi}{(\gamma - 1) (h_{3C} - \phi)} P_{tt} \right) \right\} \right] \right)$$

$$- \alpha^{2} (\gamma - 1) \phi \mathbf{\tilde{D}}^{T} \left\{ \mathbf{A}^{-1} \left(\mathbf{\tilde{D}} P_{tt} + \frac{\mathbf{\tilde{D}} \phi}{(\gamma - 1) (h_{3C} - \phi)} P_{tt} \right) \right\}$$

$$= (\gamma - 1) \left(\mathbf{\tilde{E}} - \phi \widehat{\rho} \right)$$

$$- \alpha (\gamma - 1) \left(\mathbf{\tilde{D}}^{T} \left[h_{3C} \left\{ \mathbf{A}^{-1} \left(\mathbf{\hat{U}} - \alpha \frac{\mathbf{\tilde{D}} \phi}{(h_{3C} - \phi)} \left(h_{3C} \widehat{\rho} - \widehat{E} \right) \right) \right\} \right] \right)$$

$$- \alpha (\gamma - 1) \phi \mathbf{\tilde{D}}^{T} \left\{ \mathbf{A}^{-1} \left(\mathbf{\hat{U}} - \alpha \frac{\mathbf{\tilde{D}} \phi}{(h_{3C} - \phi)} \left(h_{3C} \widehat{\rho} - \widehat{E} \right) \right) \right\}$$

$$(C.2)$$

and the discrete strong Schur form can be obtained by substituting \mathbf{D} and \mathbf{D}^{T} appropriately in Eq. (3.25). Note that the strong forms require two applications of integration by parts as described in [16,54].

References

- [1] Francis X. Giraldo, An Introduction to Element-Based Galerkin Methods on Tensor-Product Bases Analysis, Algorithms and Applications, Springer, 2020.
- [2] L. Yelash, A. Müller, M. Lukacova-Medvid'ova, F.X. Giraldo, V. Wirth, Adaptive discontinuous evolution Galerkin method for dry atmospheric flow, J. Comput. Phys. 268 (2014) 106.
- [3] V. Dolejší, M. Feistauer, A semi-implicit discontinuous Galerkin finite element method for the numerical solution of inviscid compressible flow, J. Comput. Phys. 198 (2) (Aug. 2004) 727–746.
- [4] Alex Kanevsky, Mark H. Carpenter, David Gottlieb, Jan S. Hesthaven, Application of implicit-explicit high order Runge-Kutta methods to discontinuous-Galerkin schemes, J. Comput. Phys. 225 (2) (Aug. 2007) 1753–1781.
- [5] D.A. Knoll, D.E. Keyes, Jacobian-free Newton-Krylov methods: a survey of approaches and applications, J. Comput. Phys. 193 (2) (JAN 20 2004) 357–397.
- [6] F.X. Giraldo, M. Restelli, High-order semi-implicit time-integrators for a triangular discontinuous Galerkin oceanic shallow water model, Int. J. Numer. Methods Fluids 63 (9) (Jul. 2010) 1077–1102.
- [7] M. Restelli, F.X. Giraldo, A conservative discontinuous Galerkin semi-implicit formulation for the Navier-Stokes equations in nonhydrostatic mesoscale modeling, SIAM J. Sci. Comput. (2009).
- [8] M. Kwizak, A.J. Robert, A semi-implicit scheme for grid point atmospheric models of the primitive equations, Mon. Weather Rev. 99 (1) (January 1971) 32–36.
- [9] A.J. Simmons, B.J. Hoskins, D.M. Burridge, Stability of the semi-implicit method of time integration, Mon. Weather Rev. 106 (3) (March 1978) 405–412.
- [10] Klaus Kaiser, Jochen Schutz, Ruth Schoebel, Sebastian Noelle, A new stable splitting for the isentropic Euler equations, J. Sci. Comput. 70 (3) (Mar 2017) 1390–1407.
- [11] Klaus Kaiser, Jochen Schutz, A high-order method for weakly compressible flows, Commun. Comput. Phys. 22 (4) (Oct. 2017) 1150-1174.
- [12] Georgij Bispen, Maria Lukacova-Medvid'ova, Leonid Yelash, Asymptotic preserving imex finite volume schemes for low Mach number Euler equations with gravitation, J. Comput. Phys. 335 (Apr. 2017) 222–248.
- [13] Jonas Zeifang, Klaus Kaiser, Andrea Beck, Jochen Schuetz, Claus-Dieter Munz, Efficient high-order discontinuous Galerkin computations of low Mach number flows, Commun. Appl. Math. Comput. Sci. 13 (2) (2018) 243–270.
- [14] Jonas Zeifang, Jochen Schuetz, Klaus Kaiser, Andrea Beck, Maria Lukacova-Medvid'ova, Sebastian Noelle, A novel full-Euler low Mach number imex splitting, Commun. Comput. Phys. 27 (1) (Jan 2020) 292–320.
- [15] Francis X. Giraldo, James F. Kelly, Emil M. Constantinescu, Implicit-explicit formulations of a three-dimensional nonhydrostatic unified model of the atmosphere (NUMA), SIAM J. Sci. Comput. 35 (5) (2013) B1162–B1194.
- [16] Francis X. Giraldo, Marco Restelli, A study of spectral element and discontinuous Galerkin methods for the Navier-Stokes equations in nonhydrostatic mesoscale atmospheric modeling: equation sets and test cases, J. Comput. Phys. 227 (2008) 3849–3877.
- [17] William C. Skamarock, Joseph B. Klemp, Jimy Dudhia, David O. Gill, Zhiquan Liu, Judith Berner, Wei Wang, Jordan G. Powers, Michael G. Duda, Dale Barker, Xiang Yu Huang, A description of the advanced research WRF model version 4.3, Technical report, NCAR, 2021.
- [18] Nash'at Ahmad, John Lindeman, Euler solutions using flux-based wave decomposition, Int. J. Numer. Methods Fluids 54 (2007) 41-72.
- [19] Robert V. Rohli, Chunyan Li, Meteorology for Coastal Scientists, Springer, 2021.
- [20] John.D. Anderson, Governing equations of fluid dynamics, in: John F. Wendt (Ed.), Computational Fluid Dynamics, Springer, Berlin, Heidelberg, 1992, pp. 15–51, Chapter 2.
- [21] Akshay Sridhar, Yassine Tissaoui, Simone Marras, Zhaoyi Shen, Simone Byrne, Charles Kawczynski, Kiran Pamnany, Maciej Waruzewski, Jeremy E. Kozdon, Thomas H. Gibson, Valentin Churavy, Lucas C. Wilcox, Francis X. Giraldo, Tapio Schneider, Large-eddy simulations with ClimateMachine v0.2.0: a new open-source code for atmospheric simulations on GPUs and CPUs, Geoscientific Model Development, 2021, Preprint: 1–41.
- [22] Maciej Waruszewski, Jeremy E. Kozdon, Lucas C. Wilcox, Thomas H. Gibson, Francis X. Giraldo, Entropy stable discontinuous Galerkin methods for balance laws in non-conservative form: applications to Euler with gravity, J. Comput. Phys. 468 (2022) 111507.
- [23] Gregor J. Gassner, Andrew R. Winters, David A. Kopriva, Split form nodal discontinuous Galerkin schemes with summation-by-parts property for the compressible Euler equations, J. Comput. Phys. 327 (2016) 39–66.
- [24] Francis X. Giraldo, Marco Restelli, M. Lauter, Semi-implicit formuations of the Navier-Stokes equations: application of nonhydrostatic atmospheric modeling, SIAM J. Sci. Comput. 32 (6) (2010) 3394–3425.
- [25] Bernardo Cockburn, Chi-Wang Shu, The local discontinuous Galerkin method for time-dependent convective-diffusion systems, SIAM J. Numer. Anal. 35 (6) (1998) 2440–2463.
- [26] Sohail Reddy, Yassine Tissaoui, Felipe A.V. de Branganca Alves, Simone Marras, Francis X. Giraldo, Comparison of sub-grid scale models for large-eddy simulation using a high-order spectral element approximation of the compressible Navier-Stokes equations at low Mach number, arXiv, 2021.
- [27] S. Marras, M. Nazarov, F.X. Giraldo, Stabilized high-order Galerkin methods based on a parameter-free dynamic SGS model for LES, J. Comput. Phys. 301 (2015) 77–101.
- [28] A.W. Vreman, An eddy-viscosity subgrid-scale model for turbulent shear flow: algebraic theory and applications, Phys. Fluids 16 (10) (2004) 3670-3681.
- [29] David L. Williamson, The evolution of dynamical cores for global atmospheric models, J. Meteorol. Soc. Jpn. 85B (2007) 241-269.
- [30] Andreas Muller, Michal A. Kopera, Simone Marras, Lucas C. Wilcox, Tobin Isaac, Francis X. Giraldo, Strong scaling for numerical weather prediction at petascale with the atmospheric model NUMA, Int. J. High Perform. Comput. Appl. 33 (2) (2019).
- [31] Daniel S. Abdi, Francis X. Giraldo, Efficient construction of unified continuous and discontinuous Galerkin formulations for the 3d Euler equations, J. Comput. Phys. 320 (2016) 46–68.
- [32] Daniel S. Abdi, Lucas C. Wilcox, Timothy C. Warburton, Francis X. Giraldo, A GPU-accelerated continuous and discontinuous Galerkin non-hydrostatic atmospheric model, Int. J. High Perform. Comput. Appl. 33 (1) (2017) 81–109.
- [33] Daniel S. Abdi, Francis X. Giraldo, Emil M. Constantinescu, Lester E. Carr, Lucas C. Wilcox, Timothy C. Warburton, Acceleration of the implicit-explicit non-hydrostatic unified model of the atmosphere (NUMA) on manycore processors, Int. J. High Perform. Comput. Appl. 33 (2) (2017) 242–267.
- [34] Carsten Burstedde, Lucas C. Wilcox, Omar Ghattas, p4est: scalable algorithms for parallel adaptive mesh refinement on forests of octrees, SIAM J. Sci. Comput. 33 (3) (2011) 1103–1133.
- [35] Rao Kotamarthi, Katharine Hayhoe, Linda O. Mearns, Donald Wuebbles, Jennifer Jacobs, Jennifer Jurado, Downscaling Techniques for High-Resolution Climate Projections: From Global Change to Local Impacts, Cambridge University Press, Cambridge, 2021.
- [36] Andre Robert, Bubble convection experiments with a semi-implicit formulation of the Euler equations, J. Atmos. Sci. 50 (13) (1993) 1865-1873.
- [37] J.M. Straka, Robert B. Wilhelmson, Louis J. Wicker, John R. Anderson, Kevin K. Droegemeier, Numerical solutions of a non-linear density current: a benchmark solution and comparisons, Int. J. Numer. Methods Fluids 17 (1) (1993) 1–22.
- [38] Michael Baldauf, Slavko Brdar, An analytic solution for linear gravity waves in a channel as a test for numerical models using the non-hydrostatic, compressible Euler equations, Q. J. R. Meteorol. Soc. 139 (677) (2013) 1977–1989.
- [39] William C. Skamarock, Joseph B. Klemp, Efficiency and accuracy of the Klemp-Wilhelmson time-splitting technique, Mon. Weather Rev. 122 (11) (1994) 2623–2630.
- [40] Hirofumi Tomita, Masaki Satoh, A new dynamical framework of nonhydrostatic global model using the icosahedral grid, Fluid Dyn. Res. 34 (6) (2004) 357–400.

- [41] Raymond J. Spiteri, Steven J. Ruuth, A new class of optimal high-order strong-stability-preserving time discretization methods, SIAM J. Numer. Anal. 40 (2) (2002).
- [42] Thomas M. Bendall, Thomas H. Gibson, Jemma Shipton, Colin J. Cotter, A compatible finite-element discretization for the moist compressible Euler equations, Q. J. R. Meteorol. Soc. 146 (2019) 3187–3205.
- [43] L.E. Carr, C.F. Borges, Francis X. Giraldo, Matrix-free polynomial-based nonlinear least squares optimized preconditioning and its application to discontinuous Galerkin discretizations of the Euler equations, J. Sci. Comput. 66 (2016) 917–940.
- [44] Robert M. Kirby, George Em Karniadakis, Selecting the numerical flux in discontinuous Galerkin methods for diffusion problems, J. Sci. Comput. 22–23 (2005).
- [45] F. Bassi, S. Rebay, A high-order accurate discontinuous finite element method for the numerical solution of the compressible Navier-Stokes equations, J. Comput. Phys. 131 (2) (1997) 267–279.
- [46] Carlos Erik Baumann, J. Tinsley Oden, A discontinuous hp finite element method for convection-diffusion problems, Comput. Methods Appl. Mech. Eng. 175 (3–4) (1999) 311–341.
- [47] L.N. Trefethen, D. Bau III, Numerical Linear Algebra, SIAM, Philadelphia, 1997.
- [48] Philipp Öffner, Jan Glaubitz, Hendrik Ranocha, Analysis of artificial dissipation of explicit and implicit time-integration methods, arXiv, 2016.
- [49] Christopher A. Kennedy, Mark H. Carpenter, Additive Runge-Kutta schemes for convection-diffusion-reaction equations, Appl. Numer. Math. 44 (1–2) (2003).
- [50] A. van der Sluis, H.A. van der Vorst, The rate of convergence of conjugate gradient, Numer. Math. 48 (1986) 543-560.
- [51] H.A. van der Vorst, C. Vuik, The superlinear convergence behaviour of gmres, J. Comput. Appl. Math. 48 (1993) 327-341.
- [52] G. Tumolo, L. Bonaventura, M. Restelli, A semi-implicit, semi-Lagrangian, p-adaptive discontinuous Galerkin method for the shallow water equations, J. Comput. Phys. 232 (1) (Jan 2013) 46–67.
- [53] Giovanni Tumolo, Luca Bonaventura, A semi-implicit, semi-Lagrangian discontinuous Galerkin framework for adaptive numerical weather prediction, Q. J. R. Meteorol. Soc. 141 (692A) (Oct 2015) 2582–2601.
- [54] F.X. Giraldo, J.S. Hesthaven, T. Warburton, Nodal high-order discontinuous Galerkin methods for the spherical shallow water equations, J. Comput. Phys. 181 (2) (September 2002) 499–525.