



2011-07

# High-order discontinuous element-based schemes for the inviscid shallow water equations: spectral multidomain penalty and discontinuous Galerkin methods



Calhoun is a project of the Dudley Knox Library at NPS, furthering the precepts and goals of open government and government transparency. All information contained herein has been approved for release by the NPS Public Affairs Officer.

**Dudley Knox Library / Naval Postgraduate School  
411 Dyer Road / 1 University Circle  
Monterey, California USA 93943**

# High-order discontinuous element-based schemes for the inviscid shallow water equations: spectral multidomain penalty and discontinuous Galerkin methods

J. A. Escobar-Vargas<sup>a,\*</sup>, P. J. Diamessis<sup>a</sup>, F. X. Giraldo<sup>b</sup>

<sup>a</sup>*School of Civil and Environmental Engineering, Cornell University, Ithaca, NY 14853*

<sup>b</sup>*Department of Applied Mathematics, Naval Postgraduate School, Monterey, CA 93943*

---

## Abstract

Two commonly used types of high-order-accuracy element-based schemes, collocation-based spectral multidomain penalty methods (SMPM) and nodal discontinuous Galerkin methods (DGM), are compared in the framework of the inviscid shallow water equations. Differences and similarities in formulation are identified, with the primary difference being the dissipative term in the Rusanov form of the numerical flux for the DGM that provides additional numerical stability; however, it should be emphasized that to arrive at this equivalence between SMPM and DGM requires making specific choices in the construction of both methods; these choices are addressed. In general, both methods offer a multitude of choices in the penalty terms used to introduce boundary conditions and stabilize the numerical solution. The resulting specialized class of SMPM and DGM are then applied to a suite of six commonly considered geophysical flow test cases, three linear and three non-linear; we also include results for a classical continuous Galerkin (i.e., spectral element) method for comparison. Both the analysis and numerical experiments show that the SMPM and DGM are essentially identical; both methods can be shown to be equivalent for very special choices of quadrature rules and Riemann solvers in the DGM along with special choices in the type of penalty term in the SMPM. Although we only focus our studies on the inviscid shallow water equations the results presented should be applicable to other systems of nonlinear hyperbolic equations (such as the compressible Euler equations) and extendable to the compressible and incompressible Navier-Stokes equations, where viscous terms are included.

*Keywords:* Shallow Water Equations, Spectral multidomain methods, Discontinuous Galerkin methods, interfacial treatment

---

\*Jorge A. Escobar-Vargas, School of Civil and Environmental Engineering, Cornell University, Ithaca, NY 14853, U.S.A. . email: jae47@cornell.edu

## 1. Introduction

Geophysical flows exhibit a complex structure and dynamics over a broad range of scales that render their numerical simulation a formidable task for state-of-the-art computational methods and resources. Through a complex interplay between the earth's rotation, ambient stratification and the constraining effects of lateral and vertical boundaries, flow processes in geophysical fluids commonly exhibit a characteristic horizontal lengthscale that can be a few orders of magnitude larger than its vertical counterpart [1]. Hydrostatic wave motions occur from the basin/planetary scale roughly down to the mesoscale. As the wave scales decrease, non-linear effects become significant in the form of internal/surface bores [2, 3]. At wavelengths of  $O(1km)$ , the waves also become strongly non-hydrostatic [4]. The fully non-linear and non-hydrostatic waves propagate nearly non-dissipatively and non-dispersively over long distances. Turbulent events, driven by wave breaking, current-topography interactions and other mechanisms, can be highly localized in space and time and span a broad range of scales within their region of occurrence. Finally, the dissipative effect of molecular viscosity is only felt at the smallest,  $O(1mm)$ , scales of the flow field.

As a result, the numerical methods used in the investigation of geophysical flows need to exhibit a number of preferred features. These include: a) front/wave propagation that is effectively non-dissipative and non-dispersive, b) minimum artificial dissipation at the smallest resolved scales to enable as broad a scale separation as possible, c) efficient resolution of localized flow features and complex geometries and d) optimal use of computational resources. High-order accurate element-based schemes [5, 6] are particularly appealing in addressing such needs. These schemes combine the exponential convergence and weak artificial dissipation and dispersion of standard single-domain spectral methods [7] with the spatial adaptivity of classical finite element/volume techniques [8, 9]. Furthermore, the domain decomposition philosophy inherent in these techniques renders them highly amenable for efficient parallelization [10].

On account of the inevitable impossibility of capturing the full range of scales intrinsic to a highly nonlinear, and steep, front/wave or any resulting localized turbulent event, geophysical flow simulations are inherently under-resolved. Under-resolved high-order simulations are prone towards, often catastrophic, numerical instability as Gibbs oscillations are compounded by aliasing driven by the nonlinear terms in the governing equations [11]. In high-order element-based simulations, these numerical instabilities are most pronounced at the element interfaces when strong continuity of the solution is enforced across neighboring elements [12] as is typically done in continuous Galerkin methods.

In discontinuous high-order element-based methods, neighboring subdomains carry separate values of the solution at a fixed spatial location thereby relaxing the constraint of strong continuity of the solution and significantly mitigating the above concerns of numerical instability. The two prevalent categories of such methods are spectral multidomain methods (with and without a

penalty scheme) [13, 14, 15, 16, 17, 18, 19] and discontinuous Galerkin methods (DGM) [20, 21, 22, 23, 6, 24, 25, 26]. Spectral multidomain methods, first introduced by Orzag [27] for elliptic problems were originally formulated with a strong enforcement of continuity of the solution and its derivative at the subdomain interfaces. Subsequently, Kopriva [13] extended this approach to hyperbolic problems, where the interfacial patching was implemented with an upwind scheme based on a modified method of characteristics. This approach was further refined through introducing a correction method based combination of characteristic information at the interfacial points [14]. In the framework of the compressible Navier-Stokes equations, Kopriva [15] introduced a penalty formulation to patch subdomains when higher (i.e., 2nd) derivatives were present. The formulation used in our work follows the Spectral Multidomain Penalty Method (SMPM) presented by Hesthaven [16] and expanded upon by Don [19] but implemented, to our knowledge for the first time, to the shallow water equations. In the SMPM, the strong interfacial patching conditions are replaced with a linear combination of the governing equation and the patching condition, the latter multiplied by an appropriately chosen penalty coefficient. On the other hand, DGM are based on a Galerkin weighted residual formulation where the integration is performed at the level of an individual element. Since adjacent elements are not continuously coupled, as is the case with finite and spectral elements, interfacial flux integrals do not vanish and are represented in the form of an appropriately chosen numerical flux that preserves consistency and numerical stability.

SMPM have been successfully applied to high  $Re$  incompressible stratified flow process studies in vertically non-periodic domains such as internal solitary wave-induced bottom boundary layers, turbulent wakes and propagating internal wave packets [28, 29, 30]. DGM have been effectively used in the simulation of the shallow water equations (SWE) both on the sphere and on planar but fully unstructured domains [20, 21, 23, 25, 24] and for compressible atmospheric models [22, 26].

However, the literature exploring the similarities and differences of the SMPM and DGM is limited to the recent work by Gottlieb and Jung [31] who considered the modal form of SMPM and DGM, both in Galerkin (integral) formulation. Focusing on one-dimensional conservation laws, that particular study established the equivalence between the two techniques for a specific value of the penalty coefficient and emphasized the additional flexibility of the penalty scheme in varying the value of this coefficient in space and time and splitting the advective flux at the subdomain interfaces, which provided for greater stability in regions of strong inhomogeneity of subdomain thickness. The trade-offs of accuracy vs. stability as a function of the penalty coefficient value were also examined as was the potential of the coefficient truncation method [32] in suppressing rapid error growth when using high-order polynomials in the penalty method. Finally, the impact of inconsistent evaluation of integrals (exact versus numerical quadrature) in the left and right-hand sides of the modal Galerkin formulation of the penalty method was also considered in the framework of linear and nonlinear problems. Note that both the coefficient truncation method

and the issues with integral evaluation are restricted to the modal Galerkin form of the SMPM.

No investigations are known so far that compare the collocation-based SMPM and the nodal Galerkin formulation of the DGM, the most commonly used formulations of the two methods which this paper focuses on. Furthermore, we are unaware of any comparison of the two methods in the framework of a system of multi-dimensional equations, particularly in a geophysical context. Such a comparison is the objective of the present paper. The platform for this comparison are the SWEs for a variety of reasons: a) the relative facility of their spatial and temporal discretization with respect to more complex partial differential equations, such as the Navier-Stokes equations, b) their capability for non-dissipative propagation of highly non-linear waves, which renders them an ideal experimentation tool for testing numerical schemes for nonlinear advection, the primary source of the aliasing-driven instabilities mentioned above and c) their role as a predictive tool of ocean wave phenomena for the purpose of coastal engineering applications [33] and tsunami propagation [34]. We specifically aim to compare the two methods in terms of formulation (with a focus on subdomain communication), accuracy, conservation properties, numerical stability and computational cost in the framework of specific linear and non-linear test-cases.

The remainder of the paper is organized as follows. The inviscid SWE are introduced in §2 along with their representation in linear and quasi-linear form. The formulation of SMPM and DGM is presented in §3 along with an overview of the accompanying temporal discretization. SMPM and DGM are applied to six basic test cases in §4 followed by a comparative discussion of the two methods in §5. Conclusions are offered in §6.

## 2. Inviscid Shallow Water Equations

The inviscid shallow water equations (SWE) govern the behavior of a fluid with a horizontal extent much larger than its depth, and are derived by applying the hydrostatic approximation to the incompressible Navier-Stokes equations [35]. The primitive variable formulation of the SWE is given by

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} - Z(u, v) = -g \frac{\partial h}{\partial x} \quad (1)$$

$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} + Z(u, v) = -g \frac{\partial h}{\partial y} \quad (2)$$

$$\frac{\partial h}{\partial t} + \frac{\partial}{\partial x} [(H + h)u] + \frac{\partial}{\partial y} [(H + h)v] = 0 \quad (3)$$

where  $u, v$  are the horizontal velocities,  $H$  is the mean depth,  $h$  is the displacement of the free surface,  $Z(u, v)$  is the external forcing and  $g$  is the gravitational constant.

### 2.1. Conservative form of the SWE

The inviscid shallow water equations (equations (1),(2) and (3)) can also be written in conservative form:

$$\frac{\partial \mathbf{q}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{q})}{\partial x} + \frac{\partial \mathbf{G}(\mathbf{q})}{\partial y} = \mathbf{S}(\mathbf{q}), \quad (4)$$

where the conservative variables  $\mathbf{q}$  are

$$\mathbf{q} = \begin{bmatrix} \phi \\ \phi u \\ \phi v \end{bmatrix} = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix} \quad (5)$$

the horizontal and vertical fluxes  $\mathbf{F}(\mathbf{q})$  and  $\mathbf{G}(\mathbf{q})$  are defined as

$$\mathbf{F}(\mathbf{q}) = \begin{bmatrix} \phi u \\ \phi u^2 + \frac{1}{2}\phi^2 \\ \phi uv \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \\ F_3 \end{bmatrix}, \quad \mathbf{G}(\mathbf{q}) = \begin{bmatrix} \phi v \\ \phi uv \\ \phi v^2 + \frac{1}{2}\phi^2 \end{bmatrix} = \begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix} \quad (6)$$

and the source terms  $\mathbf{S}(\mathbf{q})$  are

$$\mathbf{S}(\mathbf{q}) = \begin{bmatrix} 0 \\ f\phi v + \frac{\tau_x}{\rho} - \gamma\phi u \\ -f\phi u + \frac{\tau_y}{\rho} - \gamma\phi v \end{bmatrix}. \quad (7)$$

In Eqs. (5), (6) and (7),  $\phi = gh$  is the geopotential height,  $f = f_0 + \beta(y - y_m)$  is the Coriolis force,  $\tau_x, \tau_y$  are the components of the wind stress,  $\rho$  is the fluid density, and  $\gamma$  is a bottom friction constant.

### 2.2. Linearized SWE

Assuming a mean depth much larger than the free surface elevation ( $H \gg h$ ), and neglecting the nonlinear terms in (4), a linearized version of the conservative SWE is obtained. The modified set of conservation variables is defined as

$$\mathbf{q} = \begin{bmatrix} \phi \\ \Phi u \\ \Phi v \end{bmatrix} = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}, \quad \mathbf{F}(\mathbf{q}) = \begin{bmatrix} \Phi u \\ \Phi\phi \\ 0 \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \\ F_3 \end{bmatrix}, \quad \mathbf{G}(\mathbf{q}) = \begin{bmatrix} \Phi v \\ 0 \\ \Phi\phi \end{bmatrix} = \begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix} \quad (8)$$

where  $\Phi = gH$  is the mean depth geopotential height.

### 2.3. Quasilinear form of the SWE

Using the chain rule, Eq. (4) can be rewritten in the quasi-linear form [9, 36]

$$\begin{aligned} \frac{\partial \mathbf{q}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{q})}{\partial \mathbf{q}} \frac{\partial \mathbf{q}}{\partial x} + \frac{\partial \mathbf{G}(\mathbf{q})}{\partial \mathbf{q}} \frac{\partial \mathbf{q}}{\partial y} &= \mathbf{S}(\mathbf{q}) \\ \frac{\partial \mathbf{q}}{\partial t} + \mathbf{A} \frac{\partial \mathbf{q}}{\partial x} + \mathbf{B} \frac{\partial \mathbf{q}}{\partial y} &= \mathbf{S}(\mathbf{q}) \end{aligned} \quad (9)$$

where  $\mathbf{A}$  and  $\mathbf{B}$  are the flux Jacobian matrices, that can be decomposed (via an eigendecomposition or characteristic decomposition) as

$$\mathbf{A} = \mathbf{S}_A \mathbf{\Lambda}_A \mathbf{S}_A^{-1} \quad (10)$$

$$\mathbf{B} = \mathbf{S}_B \mathbf{\Lambda}_B \mathbf{S}_B^{-1} \quad (11)$$

where  $\mathbf{\Lambda}_A$  and  $\mathbf{\Lambda}_B$  are diagonal matrices containing the eigenvalues of  $\mathbf{A}$  and  $\mathbf{B}$ , and  $\mathbf{S}_A$ ,  $\mathbf{S}_B$  are orthogonal matrices whose columns are the respective eigenvectors.

The positive and negative flux vectors ( $\mathbf{F}^+$ ,  $\mathbf{F}^-$ ,  $\mathbf{G}^+$ ,  $\mathbf{G}^-$ ) are defined by

$$\mathbf{F}^+ = \int \mathbf{S}_A \mathbf{\Lambda}_A^+ \mathbf{S}_A^{-1} dq \quad (12)$$

$$\mathbf{F}^- = \int \mathbf{S}_A \mathbf{\Lambda}_A^- \mathbf{S}_A^{-1} dq \quad (13)$$

$$\mathbf{G}^+ = \int \mathbf{S}_B \mathbf{\Lambda}_B^+ \mathbf{S}_B^{-1} dq \quad (14)$$

$$\mathbf{G}^- = \int \mathbf{S}_B \mathbf{\Lambda}_B^- \mathbf{S}_B^{-1} dq \quad (15)$$

where  $\mathbf{\Lambda}_A^\pm$  and  $\mathbf{\Lambda}_B^\pm$  are the diagonal matrices composed of positive and negative eigenvalues of  $\mathbf{A}$  and  $\mathbf{B}$ , respectively. Based on the above decomposition, the flux vectors have the properties

$$\mathbf{\Lambda}_A = \mathbf{\Lambda}_A^+ + \mathbf{\Lambda}_A^- \quad \rightarrow \quad \mathbf{F} = \mathbf{F}^+ + \mathbf{F}^- \quad (16)$$

$$\mathbf{\Lambda}_B = \mathbf{\Lambda}_B^+ + \mathbf{\Lambda}_B^- \quad \rightarrow \quad \mathbf{G} = \mathbf{G}^+ + \mathbf{G}^-. \quad (17)$$

The eigenvalue matrices and flux vectors are the building blocks for the penalty formulation of the SWE via SMPM, and for the definition of the numerical flux of the DGM used in this work [20].

### 3. Numerical Methods

#### 3.1. Spectral Multidomain Penalty Method (SMPM)

The SMPM implemented in this work is based on the formulation first introduced by Hesthaven [18] and further refined by Don et.al. [19]. Specifically, this SMPM consists of a multidomain collocation approach based on discontinuous non-overlapping rectangular subdomains that are connected by a penalty term that ensures stability of the solution by imposing weak continuity at the subdomain interfaces. On account of the intrinsic discontinuity of the method and the critical role of interfacial patching, the formulation of the SMPM will be presented in two parts: the subdomain interior and the treatment of interfaces.

### 3.1.1. Subdomain Interior

The SMPM is based on a collocation approach in 2D quadrilateral discontinuous subdomains, where, within each subdomain, any function  $q(x, y, t)$  can be approximated by using  $N$ -th order Lagrange interpolating polynomials on a Gauss-Lobatto-Legendre (GLL) grid [18] as

$$q(x, y, t) = \sum_{i=0}^N \sum_{j=0}^N q(x_i, y_j, t) l_i(x) l_j(y) \quad (18)$$

where  $q(x_i, y_j, t)$  is the value of the function at the discrete point  $(x_i, y_j)$ , and  $l_i(x), l_j(y)$  are the  $i$ -th and  $j$ -th Lagrange interpolating polynomials based on the GLL nodes in the  $x$  and  $y$  directions, respectively. The spatial derivatives in the  $x$ -direction in the global coordinate system are approximated as

$$\frac{\partial q(x_i, y_j, t)}{\partial x} = \frac{\partial q(x_i, y_j, t)}{\partial \xi} \frac{\partial \xi}{\partial x} = \frac{\partial \xi}{\partial x} \sum_{k=0}^N D_{ik} q(x_k, y_j, t) \quad (19)$$

where, here, we assume that  $x = x(\xi)$  and  $\xi = \xi(x)$  with  $\eta \neq \eta(x)$ . In Eq. (19)  $d\xi/dx$ , represents the mapping from the local coordinate system  $\xi \in [-1, 1]$ , given by the GLL points, to the global coordinate system  $x \in \mathbf{R}$ , and  $D_{ij}$  is the Legendre spectral differentiation matrix, that is computed following Costa and Don [37]. The  $y$ -derivative is approximated in a similar manner.

### 3.1.2. Interfacial Treatment and Boundary Conditions

The penalized form of the SWE at a collocation point located along the boundaries of a subdomain requires that (see reference [19] for a similar formulation of the compressible Navier Stokes equations for chemically reacting flow)

$$\begin{aligned} \frac{\partial \mathbf{q}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{q})}{\partial x} + \frac{\partial \mathbf{G}(\mathbf{q})}{\partial y} &= \mathbf{S}(\mathbf{q}) \\ &+ \tau_1 \mathcal{Q}(\mathbf{x}) [\mathbf{F}^+(\mathbf{q}) - \mathbf{F}^+(\mathbf{q}^*)] \\ &+ \tau_2 \mathcal{Q}(\mathbf{x}) [\mathbf{F}^-(\mathbf{q}) - \mathbf{F}^-(\mathbf{q}^*)] \\ &+ \tau_3 \mathcal{Q}(\mathbf{x}) [\mathbf{G}^+(\mathbf{q}) - \mathbf{G}^+(\mathbf{q}^*)] \\ &+ \tau_4 \mathcal{Q}(\mathbf{x}) [\mathbf{G}^-(\mathbf{q}) - \mathbf{G}^-(\mathbf{q}^*)]. \end{aligned} \quad (20)$$

In (20),  $\tau_i$  ( $i = 1, \dots, 4$ ) are the penalty coefficients,  $\mathcal{Q}(\mathbf{x})$  are effectively Dirac delta functions that are non-zero only at the interfaces of the subdomain, where the penalty terms are active, and  $\mathbf{F}^\pm(\mathbf{q})$ ,  $\mathbf{G}^\pm(\mathbf{q})$ ,  $\mathbf{F}^\pm(\mathbf{q}^*)$ , and  $\mathbf{G}^\pm(\mathbf{q}^*)$  represent the positive and negative fluxes at the grid points on the particular interfaces of the subdomain (with  $*$  indicating the corresponding point on the neighboring interface) on the subdomain under consideration. In a general sense, the penalty coefficients can be viewed as weighting factors for the positive and negative fluxes across the interfaces.



In what follows, the penalized form of the SWE will be presented for the case of structured quadrilateral grids with rectangular subdomains, where the treatment for vertical interfaces is determined by the horizontal fluxes  $\partial \mathbf{F} / \partial x$ , and for the horizontal interfaces by the vertical fluxes  $\partial \mathbf{G} / \partial y$ . Embedded in the penalty coefficients  $\tau_i$  ( $i = 1, \dots, 4$ ) are mapping factors to enable consistency in units between the different terms in Eq. (20).

*Vertical interfaces.* Figure 1 presents a schematic of the vertical interface between subdomains  $I$  and  $II$ , where  $L$  or  $R$  represent any collocation point at the left and right edges of the interface.

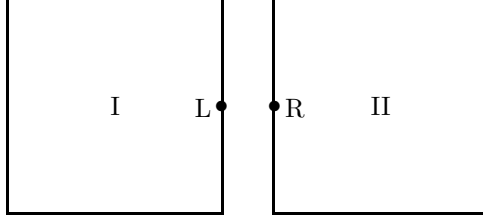


Figure 1: Vertical interface

Based on (20), the penalized form of the SWE for a point located at the left edge of the interface is

$$\begin{aligned} \frac{\partial \mathbf{q}^L}{\partial t} + \frac{\partial \mathbf{F}^L}{\partial x} + \frac{\partial \mathbf{G}^L}{\partial y} &= \mathbf{S}(\mathbf{q})^L \\ &+ \tau_1 \mathcal{Q}_L[(\mathbf{F}^+)^L - (\mathbf{F}^+)^R] \\ &+ \tau_2 \mathcal{Q}_L[(\mathbf{F}^-)^L - (\mathbf{F}^-)^R]. \end{aligned} \quad (21)$$

Similarly, for a point along the right edge of the interface the penalized form is

$$\begin{aligned} \frac{\partial \mathbf{q}^R}{\partial t} + \frac{\partial \mathbf{F}^R}{\partial x} + \frac{\partial \mathbf{G}^R}{\partial y} &= \mathbf{S}(\mathbf{q})^R \\ &+ \tau_5 \mathcal{Q}_R[(\mathbf{F}^+)^R - (\mathbf{F}^+)^L] \\ &+ \tau_6 \mathcal{Q}_R[(\mathbf{F}^-)^R - (\mathbf{F}^-)^L]. \end{aligned} \quad (22)$$

In Eq. (22)  $\tau_5, \tau_6$  are the corresponding penalty coefficients for the right edge of the interface.

*Horizontal interfaces.* Figure 2 presents a schematic of a horizontal interface between subdomains  $I$  and  $III$ . In this case,  $B$  and  $T$  represent the collocation points along the bottom and top edges of the interface. The penalized equations for a point located at the bottom edge of the horizontal interface are

$$\begin{aligned} \frac{\partial \mathbf{q}^B}{\partial t} + \frac{\partial \mathbf{F}^B}{\partial x} + \frac{\partial \mathbf{G}^B}{\partial y} &= \mathbf{S}(\mathbf{q})^B \\ &+ \tau_3 \mathcal{Q}_B[(\mathbf{G}^+)^B - (\mathbf{G}^+)^T] \\ &+ \tau_4 \mathcal{Q}_B[(\mathbf{G}^-)^B - (\mathbf{G}^-)^T] \end{aligned} \quad (23)$$

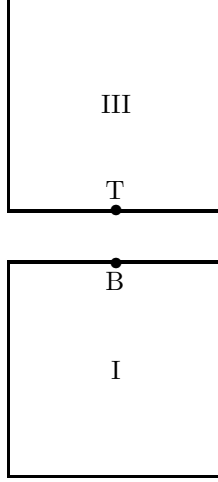


Figure 2: Horizontal interface

whereas for a point located on the top side are

$$\begin{aligned}
\frac{\partial \mathbf{q}^T}{\partial t} + \frac{\partial \mathbf{F}^T}{\partial x} + \frac{\partial \mathbf{G}^T}{\partial y} &= \mathbf{S}(\mathbf{q})^T \\
&+ \tau_7 \mathcal{Q}_T[(\mathbf{G}^+)^T - (\mathbf{G}^+)^B] \\
&+ \tau_8 \mathcal{Q}_T[(\mathbf{G}^-)^T - (\mathbf{G}^-)^B].
\end{aligned} \tag{24}$$

In Eq. (24)  $\tau_7, \tau_8$  are the corresponding penalty coefficients for the top edge of the interface.

The approach of Don et al. [19, 38] for a one-dimensional conservation law can be readily extended to the penalized equations (21)-(24) to show that the penalty scheme formally conserves mass. Moreover, the energy of the system can be shown to be bounded by its initial value [19, 38] if

$$\begin{aligned}
2\omega^L \tau_1 &\leq 1, & 2\omega^L \tau_2 &\geq 1 \\
2\omega^B \tau_3 &\leq 1, & 2\omega^B \tau_4 &\geq 1 \\
2\omega^R \tau_5 &\leq -1, & 2\omega^R \tau_6 &\geq -1 \\
2\omega^T \tau_7 &\leq -1, & 2\omega^T \tau_8 &\geq -1 \\
\omega^L \tau_1 - \omega^R \tau_5 &= 1, & \omega^L \tau_2 - \omega^R \tau_6 &= 1 \\
\omega^B \tau_3 - \omega^T \tau_7 &= 1, & \omega^B \tau_4 - \omega^T \tau_8 &= 1
\end{aligned}$$

where  $\omega^L, \omega^B, \omega^R$  and  $\omega^T$  are the GLL quadrature weights assigned to points along the left, bottom, right and top interfaces, respectively. For a uniform order of polynomial approximation,  $N$ , in each subdomain a single value of  $\omega = 2/N(N+1)$  can be used instead.

*Implementation Issues.* In this work, the averaging method [19, 39] is implemented such that the penalty coefficients for positive and negative fluxes (Eqs. (21)-(24)) at the sides of the interfaces are taken to be equal. This leads to

$$\begin{aligned}\tau_L = \tau_1 = \tau_2 &= \frac{1}{2\omega} \frac{\partial \xi}{\partial x} \\ &= \frac{1}{\omega \Delta x}\end{aligned}\quad (25)$$

$$\begin{aligned}\tau_B = \tau_3 = \tau_4 &= \frac{1}{2\omega} \frac{\partial \eta}{\partial y} \\ &= \frac{1}{\omega \Delta y}\end{aligned}\quad (26)$$

$$\begin{aligned}\tau_R = \tau_5 = \tau_6 &= -\frac{1}{2\omega} \frac{\partial \xi}{\partial x} \\ &= -\frac{1}{\omega \Delta x}\end{aligned}\quad (27)$$

$$\begin{aligned}\tau_T = \tau_7 = \tau_8 &= -\frac{1}{2\omega} \frac{\partial \eta}{\partial y} \\ &= -\frac{1}{\omega \Delta y}\end{aligned}\quad (28)$$

where  $\partial \xi / \partial x, \partial \eta / \partial y$  are the mapping factors for the penalty terms acting on vertical and horizontal interfaces respectively (see Eqs. (48) and (49)). This approach ensures stability of the penalty scheme. Moreover, the positive and negative fluxes of Eqs. (16) and (17), have been lumped into a single total flux in the penalty term.

The penalized SWE ( eqs. (21)- (24) ) may now be recast accordingly for each possible orientation of subdomain interfaces:

- Vertical interfaces

- Left edge of the interface

$$\frac{\partial \mathbf{q}^L}{\partial t} + \frac{\partial \mathbf{F}^L}{\partial x} + \frac{\partial \mathbf{G}^L}{\partial y} = \mathbf{S}(\mathbf{q})^L + \tau_L \mathcal{Q}_L[\mathbf{F}^L - \mathbf{F}^R] \quad (29)$$

- Right edge of the interface

$$\frac{\partial \mathbf{q}^R}{\partial t} + \frac{\partial \mathbf{F}^R}{\partial x} + \frac{\partial \mathbf{G}^R}{\partial y} = \mathbf{S}(\mathbf{q})^R + \tau_R \mathcal{Q}_R[\mathbf{F}^R - \mathbf{F}^L] \quad (30)$$

- Horizontal interfaces

- Bottom edge of the interface

$$\frac{\partial \mathbf{q}^B}{\partial t} + \frac{\partial \mathbf{F}^B}{\partial x} + \frac{\partial \mathbf{G}^B}{\partial y} = \mathbf{S}(\mathbf{q})^B + \tau_B \mathcal{Q}_B[\mathbf{G}^B - \mathbf{G}^T] \quad (31)$$

– Top edge of the interface

$$\frac{\partial \mathbf{q}^T}{\partial t} + \frac{\partial \mathbf{F}^T}{\partial x} + \frac{\partial \mathbf{G}^T}{\partial y} = \mathbf{S}(\mathbf{q})^T + \tau_T \mathcal{Q}_T[\mathbf{G}^T - \mathbf{G}^B] \quad (32)$$

Note that, in this scheme, unlike Hesthaven [18] no special formulation is used for the corners, which are simply treated as points that belong to two edges of the same subdomain orthogonal to each other. This simplified approach is found to be more stable than the theoretically derived one. In addition, the formulation of the penalty term is the same form used by Hesthaven [17, 18], Don et al. [19] and the 2nd author of this paper [12]. Variations of this formulation are possible and a particular one, involving the incorporation of dissipative Rusanov flux-like term, is examined in more detail in §5.3.

*Compact Representation of the SMPM.* A compact form of representing Eqs. (29) - (32) is

$$\frac{\partial \mathbf{q}^e}{\partial t} + \frac{\partial \mathbf{F}^e}{\partial x} + \frac{\partial \mathbf{G}^e}{\partial y} = \mathbf{S}(\mathbf{q})^e + \sum_{l=1}^4 \hat{\tau}_e \mathcal{Q}_e \mathbf{n}^{(e,l)} \cdot [\mathbf{F}^e - \mathbf{F}^l] \quad (33)$$

where  $\mathbf{n}^{(e,l)}$  is the outward pointing unit vector in the direction from control volume  $e$  to  $l$ ,

$$\hat{\tau} \equiv |\tau| = \frac{1}{\omega \Delta s}$$

with  $\Delta s = (\Delta x, \Delta y)$  depending on the orientation of the subdomain interfaces.

### 3.2. Discontinuous Galerkin Method (DGM)

The discontinuous Galerkin (DG) discretization of SWE (4) is as follows: we begin with the governing equations in continuous flux-form

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \mathbf{F}(\mathbf{q}) = S(\mathbf{q}). \quad (34)$$

Next we introduce a basis function expansion

$$\mathbf{q}_N(\mathbf{x}) = \sum_{i=1}^{(N+1)^2} \psi_i(\mathbf{x}) \mathbf{q}_i \quad (35)$$

where  $\psi$  represents the basis functions of order  $N$  and  $\mathbf{q}_i$  are the solution variables at specially chosen interpolation points; in this work they are chosen to be the Gauss-Legendre-Lobatto (GLL) points in order to make the comparison with the SMPM more relevant and because we have used these points in previous DG formulations (e.g., [20, 22]). Using Eq. (35) we can now construct approximations for the remainder of the spatial terms in Eq. (34). For example, we can now represent the flux tensor as

$$\mathbf{F}_N = \mathbf{F}(\mathbf{q}_N) \quad (36)$$

and the source function as

$$\mathbf{S}_N = S(\mathbf{q}_N). \quad (37)$$

Upon defining these expansions, we can then substitute them into Eq. (34), multiply the equations by a test function, and integrate to obtain the element-wise integral problem: find  $\mathbf{q}_N \in \mathcal{S}(\Omega_e) \forall \psi \in \mathcal{S}(\Omega_e)$  on each element  $\Omega_e$  such that

$$\int_{\Omega_e} \psi_i \left( \frac{\partial \mathbf{q}_N}{\partial t} + \nabla \cdot \mathbf{F}_N \right) d\Omega_e = \int_{\Omega_e} \psi_i S_N d\Omega_e \quad (38)$$

where  $\mathcal{S}$  is the finite-dimensional space

$$\mathcal{S} = \{ \psi \in \mathcal{L}^2(\Omega) : \psi|_{\Omega_e} \in P_N(\Omega_e) \forall \Omega_e \},$$

$P_N$  is the polynomial space of order  $N$  defined on  $\Omega_e$  and the union of these  $N_e$  elements defines the global domain, i.e.,  $\Omega = \bigcup_{e=1}^{N_e} \Omega_e$ . Next, we integrate the divergence term by parts to get

$$\begin{aligned} \int_{\Omega_e} \psi_i \frac{\partial \mathbf{q}_N^{(e)}}{\partial t} d\Omega_e + \sum_{l=1}^4 \int_{\Gamma_e} \psi_i \mathbf{n}^{(e,l)} \cdot \mathbf{F}_N^{(e)} d\Gamma_e & - \int_{\Omega_e} \nabla \psi_i \cdot \mathbf{F}_N^{(e)} d\Omega_e \\ & = \int_{\Omega_e} \psi_i S_N^{(e)} d\Omega_e \end{aligned} \quad (39)$$

where  $\mathbf{n}^{(e,l)}$  is the outward normal vector going from element  $e$  to element  $l$  that defines a specific edge of the (in this specific case) quadrilateral control volume. Now, since the solutions are discontinuous across element boundaries then it becomes critical (in order to construct a consistent and stable numerical approximation to the governing continuous equations) to choose the flux tensor carefully. To resolve this inconsistency, a numerical flux is introduced that we denote by  $F^{(*,l)}$ . The simplest choice is the mean value between the two elements claiming the same interface

$$\mathbf{F}_N^{(*,l)} = \frac{1}{2} \left[ \mathbf{F}_N^{(e)} + \mathbf{F}_N^{(l)} \right]$$

where the superscripts  $e$  and  $l$  represent the element under consideration and the side (interface) neighbor; unfortunately this numerical flux is not the best choice. Another easy but better choice is the local Lax-Friedrichs (or Rusanov) flux defined as

$$\mathbf{F}_N^{(*,l)} = \frac{1}{2} \left[ \mathbf{F}_N^{(e)} + \mathbf{F}_N^{(l)} - \delta_{diss} |\lambda_{max}| \mathbf{n}^{(e,l)} \left( \mathbf{q}_N^{(l)} - \mathbf{q}_N^{(e)} \right) \right] \quad (40)$$

where  $\lambda_{max}$  is the maximum wave speed of the shallow water equations (the maximum eigenvalue of the Jacobian matrix at the edge  $l$ ) and we have included the switch  $\delta_{diss}$  that controls whether the dissipation term is included. With a

specific numerical flux defined, the DG formulation becomes

$$\begin{aligned} \int_{\Omega_e} \psi_i \frac{\partial \mathbf{q}_N^{(e)}}{\partial t} d\Omega_e + \sum_{l=1}^4 \int_{\Gamma_e} \psi_i \mathbf{n}^{(e,l)} \cdot \mathbf{F}_N^{(*,l)} d\Gamma_e - \int_{\Omega_e} \nabla \psi_i \cdot \mathbf{F}_N^{(e)} d\Omega_e \\ = \int_{\Omega_e} \psi_i S_N^{(e)} d\Omega_e \end{aligned} \quad (41)$$

that is in fact the *weak* form DGM. Integrating by parts one more time yields the following mathematically equivalent system

$$\begin{aligned} \int_{\Omega_e} \psi_i \frac{\partial \mathbf{q}_N^{(e)}}{\partial t} d\Omega_e + \sum_{l=1}^4 \int_{\Gamma_e} \psi_i \mathbf{n}^{(e,l)} \cdot \left( \mathbf{F}_N^{(*,l)} - \mathbf{F}_N^{(e)} \right) d\Gamma_e + \int_{\Omega_e} \psi_i \nabla \cdot \mathbf{F}_N^{(e)} d\Omega_e \\ = \int_{\Omega_e} \psi_i S_N^{(e)} d\Omega_e \end{aligned} \quad (42)$$

which is the *strong* form DGM and is the form that we shall use to compare and contrast with the SMPM described in §3.1. Next, let us expand the terms  $\mathbf{q}_N$  and  $\mathbf{S}_N$  in order to rewrite Eq. (42) in matrix-vector form. Expanding these terms in Eq. (42) gives

$$M_{ij}^{(e)} \frac{dq_j^{(e)}}{dt} + \left( \mathbf{D}_{ij}^{(e)} \right)^T \mathbf{F}_j^{(e)} + \sum_{l=1}^4 \left( \mathbf{M}_{ij}^{(l)} \right)^T \left( \mathbf{F}_j^{(*,l)} - \mathbf{F}_j^{(e)} \right) = M_{ij}^{(e)} S_j^{(e)} \quad (43)$$

where the elemental matrices are defined as follows:

$$M_{ij}^{(e)} = \int_{\Omega_e} \psi_i \psi_j d\Omega_e, \quad \mathbf{D}_{i,j}^{(e)} = \int_{\Omega_e} \psi_i \nabla \psi_j d\Omega_e, \quad \mathbf{M}_{ij}^{(l)} = \int_{\Gamma_e} \psi_i \psi_j \mathbf{n}^{(e,l)} d\Gamma_e \quad (44)$$

where  $\mathcal{T}$  denotes the transpose operator. At this point in the DG formulation, we have to introduce numerical quadrature in order to evaluate the integrals defined in Eq.(44) in the following way

$$\begin{aligned} M_{ij}^{(e)} &= \sum_{k=1}^{(Q+1)^2} \omega_k^{(e)} |J_k^{(e)}| \psi_i(\mathbf{x}_k) \psi_j(\mathbf{x}_k), \\ \mathbf{D}_{ij}^{(e)} &= \sum_{k=1}^{(Q+1)^2} \omega_k^{(e)} |J_k^{(e)}| \psi_i(\mathbf{x}_k) \nabla \psi_j(\mathbf{x}_k), \\ \mathbf{M}_{ij}^{(l)} &= \sum_{k=1}^{(Q+1)} \omega_k^{(l)} |J_k^{(l)}| \psi_i(\mathbf{x}_k) \psi_j(\mathbf{x}_k) \end{aligned} \quad (45)$$

where  $Q$  is the number of quadrature points along each direction of the quadrilateral element, and  $\omega$  and  $J$  are quadrature weights and Jacobians, respectively.

Using GLL points for both interpolation and integration we obtain the following element matrices

$$\begin{aligned} M_{ij}^{(e)} &= \omega_i^{(e)} |J_i^{(e)}| \delta_{ij}, \\ D_{ij}^{(e)} &= \omega_i^{(e)} |J_i^{(e)}| \nabla \psi_j(\mathbf{x}_i), \\ M_{ij}^{(l)} &= \omega_i^{(l)} |J_i^{(l)}| \delta_{ij} \end{aligned} \quad (46)$$

where  $\delta$  denotes the usual Kronecker delta function. Using Eq. (46) in Eq. (43) and dividing by the mass matrix yields:

$$\frac{dq_i^{(e)}}{dt} + (\nabla \psi_j(\mathbf{x}_i))^T \mathbf{F}_j^{(e)} = S_i^{(e)} + \sum_{l=1}^4 \tau_i^{(l)} \mathcal{Q}_i^{(l)} \mathbf{n}_i^{(e,l)} \cdot (\mathbf{F}_i^{(e)} - \mathbf{F}_i^{(*,l)}) \quad (47)$$

where

$$\mathcal{Q}_i^{(l)} = \begin{cases} 1 & \text{if } i \text{ is on the edge } l \\ 0 & \text{otherwise} \end{cases}$$

and

$$\tau_i^{(l)} = \frac{\omega_i^{(l)} |J_i^{(l)}|}{\omega_i^{(e)} |J_i^{(e)}|};$$

note that Eq. (47) is quite similar to Eq. (33) for the SMPM.

Next, we need to simplify the penalty-like term that we have called  $\tau$ . To do so requires explicitly stating the value of the Jacobians of both the element and edges. For the sake of simplicity, if we assume that  $\xi = \xi(x)$  and  $\eta = \eta(y)$ , that is, that the computational axes are aligned exactly with the physical axes, then we can write

$$\begin{aligned} \xi &= \frac{2(x - x_0)}{\Delta x} - 1 \\ \eta &= \frac{2(y - y_0)}{\Delta y} - 1 \end{aligned} \quad (48)$$

where  $x_0, y_0$  is the left-bottom most point on each element and  $\Delta x, \Delta y$  is the length of the element along the x and y directions, respectively.

This mapping yields the following metric terms

$$\begin{aligned} \frac{\partial \xi}{\partial x} &= \frac{2}{\Delta x} \\ \frac{\partial \eta}{\partial y} &= \frac{2}{\Delta y} \end{aligned} \quad (49)$$

with the following Jacobians

$$|J^{(e)}| \equiv \frac{\partial x}{\partial \xi} \frac{\partial y}{\partial \eta} - \frac{\partial x}{\partial \eta} \frac{\partial y}{\partial \xi} = \frac{\Delta x \Delta y}{4}$$

and

$$|J^{(l)}| = \begin{cases} \frac{\Delta y}{2} & \text{along a vertical interface (Left-Right edge)} \\ \frac{\Delta x}{2} & \text{along a horizontal interface (Top-Bottom edge)}. \end{cases}$$

From the definition of these metric terms we can see that the penalty-like term simplifies to

$$\tau_i^{(l)} = \begin{cases} \frac{2}{\omega \Delta x} & \text{along a vertical interface (Left-Right edge)} \\ \frac{2}{\omega \Delta y} & \text{along a horizontal interface (Top-Bottom edge)} \end{cases}$$

where  $\omega = \omega_0 = \omega_N$  is the value of the quadrature weight at the beginning or end point (they are equal by symmetry). Introducing the DGM numerical flux given in Eq. (40) into Eq. (47) yields

$$\begin{aligned} & \frac{dq_i^{(e)}}{dt} + (\nabla \psi_j(\mathbf{x}_i))^T \mathbf{F}_j^{(e)} \\ &= S_i^{(e)} + \sum_{l=1}^4 \hat{\tau}_i^{(l)} \mathcal{Q}_i^{(l)} \mathbf{n}_i^{(e,l)} \cdot \left[ \mathbf{F}_i^{(e)} - \mathbf{F}_i^{(l)} - \delta_{diss} |\lambda_{max}| \mathbf{n}_i^{(e,l)} \left( \mathbf{q}_i^{(l)} - \mathbf{q}_i^{(e)} \right) \right] \end{aligned} \quad (50)$$

where

$$\hat{\tau} \equiv \frac{\tau}{2} = \frac{1}{\omega \Delta s}$$

and  $\Delta s = (\Delta x, \Delta y)$  depending in which direction the interface is oriented. At this point, we have not made too many sacrifices or simplifications in deriving Eq. (50). This equation is in fact a valid DGM representation of the shallow water equations with only the very slight assumptions that:

1. The computational coordinates  $(\xi, \eta)$  are aligned with the physical coordinates  $(x, y)$ .
2. Co-located interpolation and integration points are used. The fact that we have chosen these points to be the GLL points results in inexact integration.
3. The numerical flux used is the simple Rusanov flux.

Taking the special case  $\delta_{diss} = 0$ , that is, no dissipation in the flux term, yields

$$\frac{dq_i^{(e)}}{dt} + (\nabla \psi_j(\mathbf{x}_i))^T \mathbf{F}_j^{(e)} = S_i^{(e)} + \sum_{l=1}^4 \hat{\tau}_i^{(l)} \mathcal{Q}_i^{(l)} \mathbf{n}_i^{(e,l)} \cdot \left[ \mathbf{F}_i^{(e)} - \mathbf{F}_i^{(l)} \right] \quad (51)$$

which is identical to the SMPM representation given in Eq. (33). Eq. (51) shows that another way of viewing the penalty term is as an extra differencing term (as is evident by the  $\frac{1}{\Delta s}$  term in  $\hat{\tau}$  and  $\Delta F$  in the numerator) that considers the information from the neighboring elements, which is in fact what we mean



s	$\alpha_{ik}$				$\beta_{ik}$			
	4	1	0	1	0	1/2	0	1/2
	0	2/3	0	1/3	0	0	0	1/6
	0	0	0	0	1	0	0	1/2

Table I: Coefficients for the third order - four stage SSP-RK(34) method

by the usual term *flux*. In §4 we use Eq. (50) with and without the dissipation term to compare the SMPM with the DGM. We now turn our discussion to the time-integrator we use to advance the SMPM and DGM solutions forward in time.

### 3.3. Temporal Discretization

To retain the high-order accuracy of the SMPM and the DGM, a high-order time advancement scheme is needed. The explicit strongly stability preserving Runge-Kutta (SSP-RK) method [40, 41] is implemented for both approaches. Consider the following initial value problem

$$\frac{dq}{dt} = R(q). \quad (52)$$

The prediction at the time  $n + 1$  is based on the existing solution at the time  $n$  and the forcing terms  $R(q)$ . The scheme can be written as [41]

$$q^{(0)} = q^n \quad (53)$$

$$q^{(i)} = \sum_{k=0}^{i-1} \left( \alpha_{ik} q^{(k)} + \Delta t \beta_{ik} R(q^{(k)}) \right), \quad i = 1, 2, \dots, s \quad (54)$$

$$q^{(n+1)} = q^{(s)} \quad (55)$$

where  $s$  are the number of stages of the SSP-RK approach,  $\alpha_{ik}$  and  $\beta_{ik}$  are constant coefficients given in Table I [41], and  $\Delta t$  is the size of the time step at a specific time.

## 4. Test cases: Description and Results

Six test cases are examined to compare the performance of the SMPM and DGM in terms of accuracy, dynamic stability, robustness and conservation properties: three linear (standing wave, Kelvin wave, and Stommel problem), where accuracy can be evaluated through the availability of analytic solutions, and three non-linear (nonlinear Stommel, equatorial Rossby wave, and Riemann problem) that provide a platform for assessing the dynamic stability and robustness of the methods. In addition, results obtained with the spectral element method (SEM) [42, 43, 44] are included to compare, for each case, the behavior of a continuous method with a discontinuous element-based approach. For the

linear cases an additional error analysis based on the normalized  $L_\infty$  and  $L_2$  norms of the error is performed.

The normalized  $L_\infty$  and  $L_2$  error norms are defined as

$$\|h\|_{L_\infty} = \frac{\max_{\mathbf{x} \in \Omega} (h_{exact} - h)}{\max_{\mathbf{x} \in \Omega} h_{exact}} \quad (56)$$

$$\|h\|_{L_2} = \sqrt{\frac{\int_{\Omega} (h_{exact} - h)^2 d\Omega}{\int_{\Omega} h_{exact}^2 d\Omega}}. \quad (57)$$

The mass ( $M$ ) and energy ( $E$ ) of the system are measured in the following way

$$M = \int_{\Omega} \phi d\Omega \quad (58)$$

$$E = \int_{\Omega} [\phi(u^2 + v^2) + \phi^2] d\Omega. \quad (59)$$

The metric for assessing mass and energy conservation is the respective relative error, defined with respect to the corresponding initial values of  $M$  and  $E$ . It is computed as

$$R_M = \left| \frac{M_t - M_0}{M_0} \right|, \quad R_E = \left| \frac{E_t - E_0}{E_0} \right| \quad (60)$$

where  $R_M$  and  $R_E$  are the relative errors in mass and energy, and  $M_0, E_0, M_t, E_t$  are the corresponding values for mass and energy at the initial and final times of the simulation, respectively. For each test case, it is specified explicitly if mass and energy are lost or generated by the end of simulation.

For all simulations no boundary conditions are applied to the continuity equation. For the momentum equation no-flux (i.e., reflecting) boundary conditions are applied along all four walls of the basins; for the SEM and SMPM methods this is accomplished via strong homogeneous Dirichlet boundary conditions whereas for the DGM they are satisfied in a weak sense.

To compute the Courant number a high-order cell technique is used, where the cells are defined based on the GLL points on each subdomain. A mean velocity and geopotential height is defined at the center of each cell [23]. With these considerations, the Courant number is defined as

$$\text{Courant Number} = \max \left( \frac{\Delta t (U + \sqrt{\phi})}{\Delta s} \right)$$

where  $\Delta t$  is the size of the time step,  $U$  is the mean velocity magnitude at the cell,  $\phi$  is the average geopotential height in the cell and  $\Delta s = \sqrt{\Delta x^2 + \Delta y^2}$  is the grid spacing. For SMPM and DGM, the maximum Courant number  $\Delta t$  that ensures stability of the numerical simulations 0.5 (Courant Number  $\leq 0.5$ ). The equivalent value for SEM is 1. As specified in the relevant sections, two

test cases (standing and Kelvin wave) are run with a significantly smaller time step to prevent the time-stepping error from dominating the error associated with the spatial discretization. Nonetheless, as the conservation properties of the SMPM are negatively impacted by a linearly growing loss of mass which is of order machine epsilon at each time step, all other test cases are run with a time step that is 80% the maximum time step associated with Courant number limits indicated above. We refer the reader to §5.1 and 5.2 for further discussion on time-stepping error and the impact of time step on the conservation properties of the spatial discretization methods under consideration. The degree of polynomial approximation is varied from  $N = 4$  to 20. The number of subdomains is also varied within a range dependent on the geometry of each case, and the SSP-RK34 method defined previously is used to advance in time the simulations.

#### 4.1. Linear Problems

In this section, we compare the three methods quantitatively using linear test cases that have analytic solutions.

##### 4.1.1. Linear Standing Wave

This case represents the evolution in time of a wave driven only by gravitational effects ( $S = 0$ ) through an initial perturbation of the free surface. From references [23, 45], the analytic solution for this case is given by

$$h(x, y, t) = \cos(\pi x) \cos(\pi y) \cos(\pi t\sqrt{2}) \quad (61)$$

$$u(x, y, t) = \frac{1}{\sqrt{2}} \sin(\pi x) \sin(\pi y) \sin(\pi t\sqrt{2}) \quad (62)$$

$$v(x, y, t) = \frac{1}{\sqrt{2}} \cos(\pi x) \sin(\pi y) \sin(\pi t\sqrt{2}) \quad (63)$$

with  $(x, y) \in [0, 1] \times [0, 1]$ .

The simulations are run for  $t \in [0, 0.5]$ . Figure 3 shows results for SMPM, DGM and SEM simulations for a fixed number of subdomains and variable order of polynomial approximation  $N$ . A time step which is 1/50th of that associated with a Courant number value of 0.4 is used, to make time-stepping errors sufficiently small. The results are indistinguishable if an even smaller time step is employed. Exponential convergence of the error norms for free surface elevation and horizontal velocity is attained for each method for polynomial degree less or equal than  $N = 8$ . At higher values of  $N$ , the convergence rate is finally reduced, reaching a plateau of the order of  $O(10^{-12})$ , the cause of which we are unable to determine. The Galerkin based methods (i.e. DGM, SEM) conserve mass up to machine precision. The SMPM mass cumulatively loses mass over time. All three methods show improved energy conservation with increasing  $N$  with the relative error reaching a value of  $O(10^{-12})$  at  $N = 8$ . An interpretation for the performance of the SMPM in terms of mass conservation is offered in §5.1.

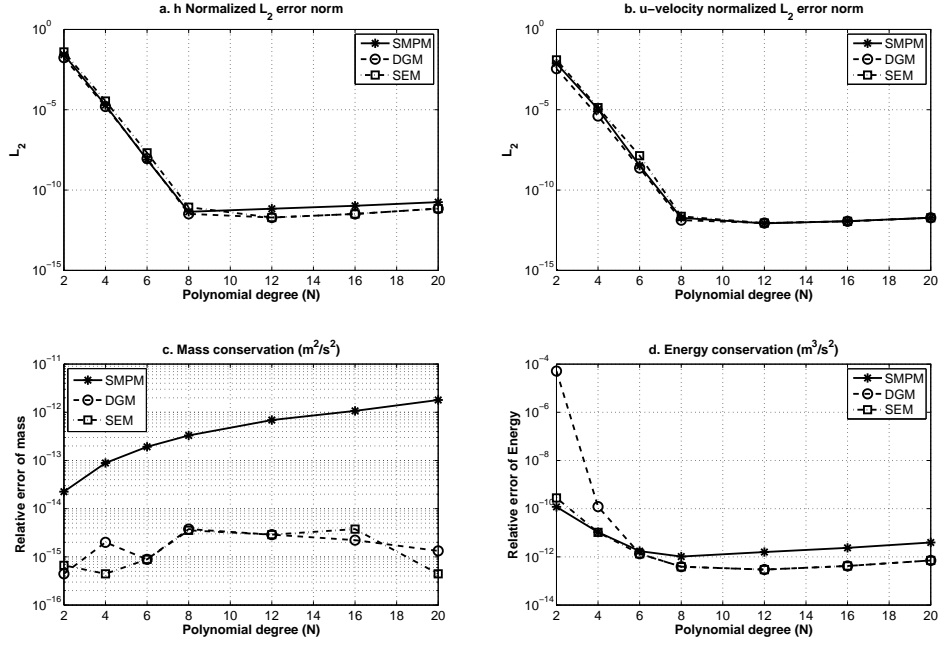


Figure 3: Analysis of the standing wave ( $5 \times 5$  subdomains) at  $t = 0.5$  seconds for a varying number of GLL points. a)  $L_2$  normalized relative error in the free surface elevation  $h$ . b)  $L_2$  normalized relative error in  $u$  velocity. c) Relative error in mass. d) Relative error in energy.

#### 4.1.2. Linear Kelvin Wave

The equatorial Kelvin wave is a low amplitude non-dispersive wave trapped in the vicinity of the equator. It is driven by rotational and gravitational effects through an initial perturbation of the free surface. The analytic solution for this case [23, 46] is

$$h(x, y, t) = 1 + \exp\left(-\frac{y^2}{2}\right) \exp\left(-\frac{(x + 5 - t)^2}{2}\right) \quad (64)$$

$$u(x, y, t) = \exp\left(-\frac{y^2}{2}\right) \exp\left(-\frac{(x + 5 - t)^2}{2}\right) \quad (65)$$

$$v(x, y, t) = 0 \quad (66)$$

for  $f_0 = 0$ ,  $\beta = 1$  and  $(x, y) \in [-20, 20] \times [-10, 10]$ .

Simulations are run for  $t \in [0, 5]$ . Figure 4 shows results for this case for a domain discretized with  $20 \times 10$  elements and a varying value of  $N$ . As with the Standing wave, here the time step is 1/50th that associated with a Courant number value of 0.4. No further reduction in time step was required to make time-stepping errors sufficiently small. The behavior of the error norms is similar to that observed for the linear standing wave: exponential convergence

is observed for all the three methods. DGM and SEM conserve mass up to machine precision. On the contrary, SMPM again shows a loss of mass, which, in the end of simulations, is up to one order of magnitude larger than the value computed for DGM and SEM. The trend in relative error of total energy conserved is comparable to that observed for the linear standing wave in Fig. 3. Improved energy conservation occurs with increasing  $N$  with a relative error value of  $O(10^{-13})$  observed for  $N = 20$ .

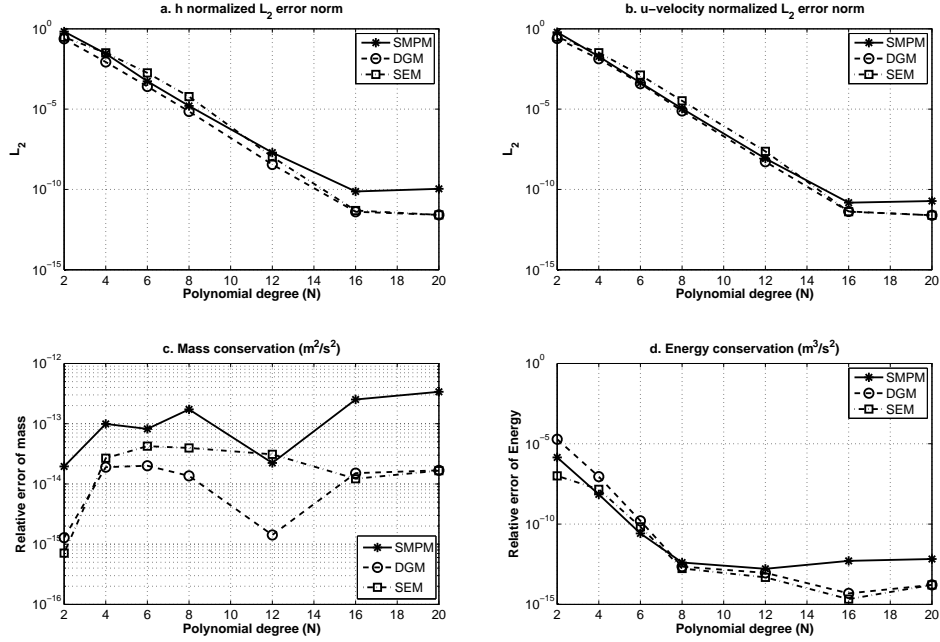


Figure 4: Kelvin wave results for  $20 \times 10$  subdomains at  $t = 5$ . Panels (a) through (d) show the same quantities with Fig. 3.

#### 4.1.3. Linear Stommel Problem

This problem [47] also known as westward intensification of wind-driven ocean currents, represents the steady balance between rotation, gravity, friction and wind stress in a square ocean basin. A sinusoidal wind stress forces an unperturbed free surface generating a small amplitude wave moving westward due to the Coriolis force that is compensated by bottom friction and gravitational effects and, eventually, reaches steady state. The analytic solution used for this

case is [23]

$$\begin{aligned}
h(x, y, t) &= \left( \frac{C_1}{\lambda_1} e^{\lambda_1 x} + \frac{C_2}{\lambda_2} e^{\lambda_2 x} \right) \frac{\gamma \pi}{l} \cos\left(\frac{\pi y}{l}\right) \\
&\quad + \frac{\tau \beta}{\gamma} \left(\frac{l}{\pi}\right)^2 \cos\left(\frac{\pi y}{l}\right) \\
&\quad + f \sin\left(\frac{\pi y}{l}\right) (C_1 e^{\lambda_1 x} + C_2 e^{\lambda_2 x} + C_3)
\end{aligned} \tag{67}$$

$$u(x, y, t) = -(C_1 e^{\lambda_1 x} + C_2 e^{\lambda_2 x} + C_3) \frac{\pi}{l} \cos\left(\frac{\pi y}{l}\right) \tag{68}$$

$$v(x, y, t) = (C_1 \lambda_1 e^{\lambda_1 x} + C_2 \lambda_2 e^{\lambda_2 x}) \sin\left(\frac{\pi y}{l}\right) \tag{69}$$

where

$$C_1 = C_3 \frac{1 - e^{\lambda_2 l}}{e^{\lambda_2 l} - e^{\lambda_1 l}} \tag{70}$$

$$C_2 = -C_3 \frac{1 - e^{\lambda_1 l}}{e^{\lambda_2 l} - e^{\lambda_1 l}} \tag{71}$$

$$C_3 = \frac{\tau l}{\pi \gamma} \tag{72}$$

For the case presented here,  $f_0 = 1 \times 10^{-4}$ ,  $\beta = 1 \times 10^{-11}$ ,  $\gamma = 1 \times 10^{-6}$ ,  $g = 10$ ,  $\rho = 1000$ ,  $\tau = 0.2$ ,  $H_0 = 1000$ , and  $(x, y) \in [0, 1 \times 10^6] \times [0, 1 \times 10^6]$ . Note that the solution is symmetric with respect to the  $y$  axis.

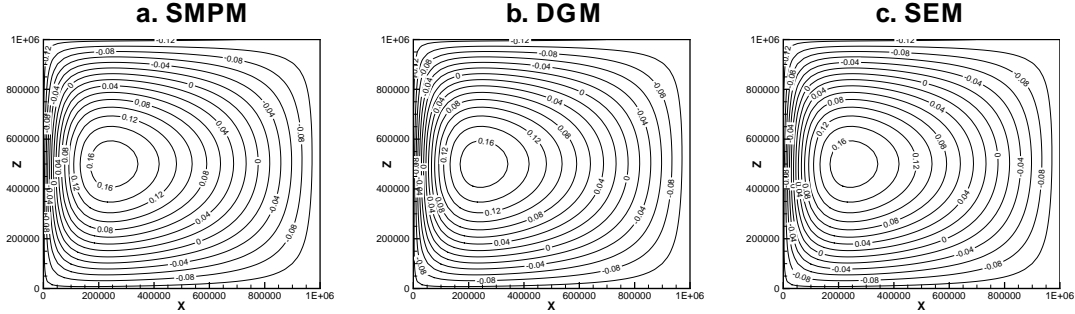


Figure 5: Free surface elevation computed by all three methods for the linear Stommel problem for  $5 \times 5$  subdomains and  $N = 12$  at  $t = 400$  days

Simulations are run until the solution is close to the steady state (i.e.  $t = 320$  days), and the structure of the steady state flow field, displaying the expected symmetry around the horizontal axis at  $z = 5 \times 10^5$ , is shown in Fig. 5 for all three methods. Figure 6 shows the error norm convergence curves for the case of a  $5 \times 5$  mesh for solutions obtained with different values of  $N$ . For all three methods, the error in the free surface displacement shows an exponential

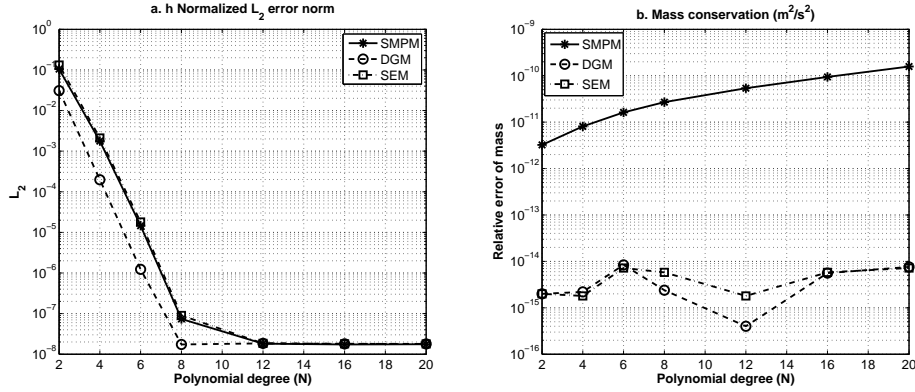


Figure 6: Linear Stommel problem results for  $5 \times 5$  subdomains at  $t = 320$  days. a) Convergence plot for the  $L_2$  normalized relative error in the free surface elevation  $h$ . b) Relative error in mass conservation.

convergence similar to the previous two linear cases for up to  $N = 8$ , beyond which the error norms level off to a constant value. This plateau is reached because an exact steady is almost never attained in practice, as simulations are dominated by slowly-decaying, weak-amplitude basin-scale modes, with the decay time of the gravest, longest-wavelength, mode reaching 60 years [48]. Mass is conserved up to machine precision by DGM and SEM, whereas SMPM shows a loss of total mass up to three orders of magnitude larger than DGM and SEM.

#### 4.2. Nonlinear Problems

In this section, we compare the three methods qualitatively using nonlinear test cases that, unfortunately, do not have analytic solutions. Instead, we use the conservation of mass and energy to compare the methods. All three models formally should conserve mass but are not guaranteed to conserve energy. It is possible to conserve energy (at least up to the time-truncation error) but this requires slight modifications to the discrete operators that we will not pursue in this work.

##### 4.2.1. Nonlinear Rossby Soliton

This case considers an equatorial non-linear Rossby wave of weak amplitude, driven by gravity and rotational forces. It is initialized by a Gaussian-like perturbation in the free surface elevation. An approximate asymptotic solution of the system of Korteweg-DeVries equations resulting from the SWE through application of the method of multiple scales is obtained for this problem in [49]. Although this first order solution does not provide a reference to assess the convergence rate of the numerically computed solution for the SEM, DGM, and SMPM, it is used to compare associated phase speed and solution structure

with the corresponding estimates computed by the three numerical methods. For this case  $(x, y) \in [-24, 24] \times [-8, 8]$ ,  $g = 1$ , and the Coriolis force  $f(y) = y$ .

Simulations are run for  $t \in [0, 40]$ . All three methods accurately reproduce the free surface/velocity structure of the soliton and its propagation at a constant phase speed equal to the analytically predicted value. The structure of the free surface elevation field at the end of the simulation, with its characteristic two-lobe structure, as computed by all three methods is shown in Fig. 7. Figure 8 shows results for mass and energy conservation for  $24 \times 8$  subdomains, and varying  $N$ , which are similar to their counterparts obtained for the linear cases. The SMPM is subject to a decrease in mass when the polynomial order increases. The DGM conserves mass up to machine precision, with the SEM offering comparable performance. The SEM and SMPM are the most and least energy conserving, respectively. As discussed in §5.1, the energy conservation properties of the DGM are highly dependent on the formulation of the numerical flux and the use of spectral filtering (see Fig. 12).

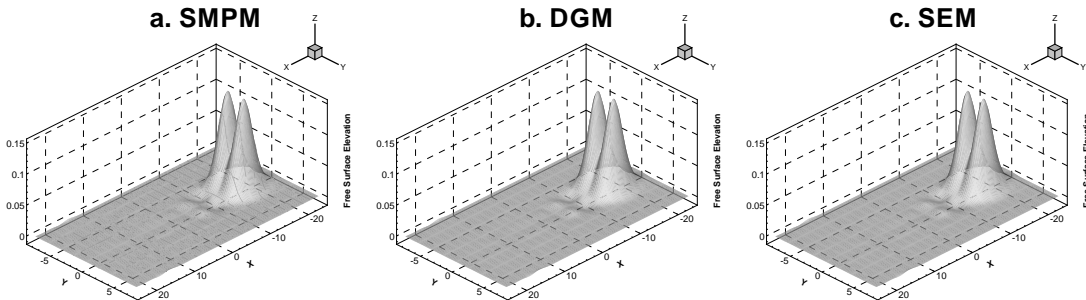


Figure 7: Qualitative comparison of the Non-linear Rossby wave results with  $24 \times 8$  subdomains,  $N = 12$ , and at time  $t = 40$ .

#### 4.2.2. Nonlinear Stommel Problem

The same configuration (forcing parameters, dimensions of the physical domain, and boundary conditions) is used as in the linear Stommel problem. However, the fully nonlinear set of Eqs. (4) are now solved. In this case, a shift of the gyre toward the northwest part of the basin is expected due to the effect of the nonlinear terms.

Figure 9 shows the steady state results, for a domain with  $5 \times 5$  subdomains. Similar trends are observed for all three methods. Note that in this particular case, the differences in subdomain interface treatment between SMPM and DGM give rise to challenges of numerical stability for the former, when values of polynomial degree  $N$  are used. In the SMPM, when  $5 \times 5$  subdomains are used and  $N \geq 12$ , weak spurious oscillations develop in the top left corner of the domain and intensify, as time advances, eventually forcing a catastrophic blow-up of the solution. As a counter-measure, a 16-th order Boyd-Vandeven filter [50] is used, which attenuates only the very highest modes of the solution,



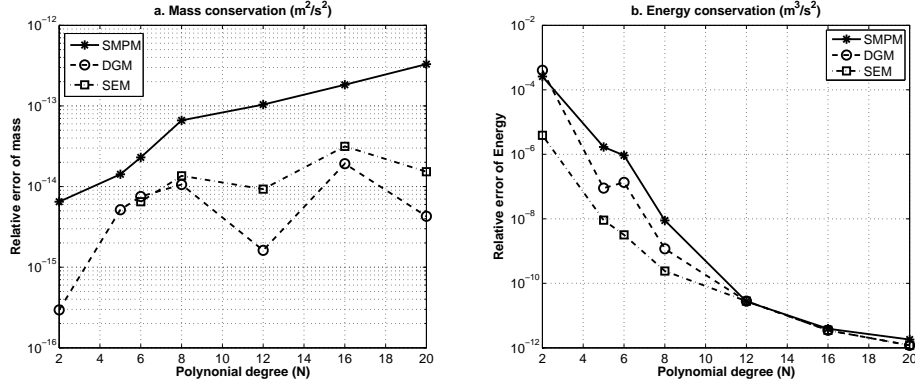


Figure 8: Non-linear Rossby wave results for  $24 \times 8$  subdomains at  $t = 40$ . a) Relative error in mass. b) Relative error in energy.

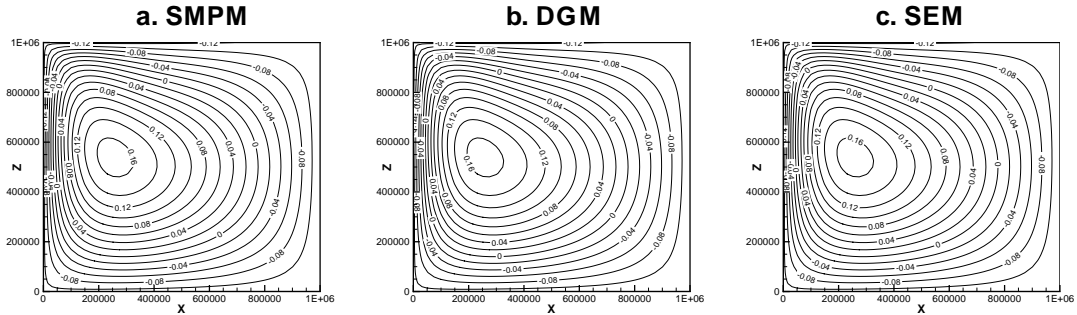


Figure 9: Nonlinear Stommel problem for  $5 \times 5$  subdomains and  $N = 12$  at  $t = 400$  days

to suppress these oscillations. This problem does not occur for the DGM, as the spurious oscillations are damped by the dissipative term  $\delta_{diss} = 1$  in the numerical flux. The sensitivity of the DGM and SMPM to the presence of dissipative terms is examined in greater detail in §5.3. Figure 10 shows the behavior of the relative error in mass as a function of  $N$ , which is similar to what is observed for the corresponding linear problem (Fig. 6). Results are restricted to  $N \leq 8$ , as high-order polynomial approximations require the use of a spectral filter to preserve stability.

#### 4.2.3. Nonlinear Riemann Problem

This modification of the circular dam break problem [51] is considered as a platform to assess the performance of the three methods in simulating strongly nonlinear flows, i.e. flow fields with distinct sharp spatial gradients. The initial condition, a Gaussian bump (used instead of a cylindrical step function), is characterized by such a sharp gradient and has free surface and velocity fields

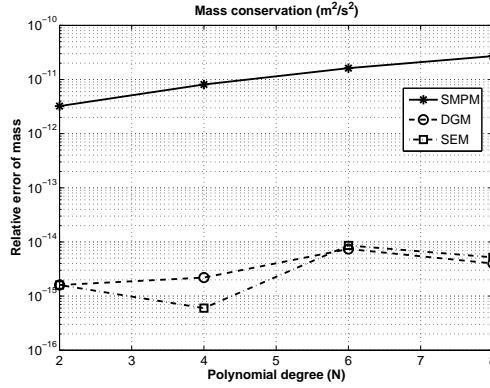


Figure 10: Relative error in mass as a function of polynomial order for the Nonlinear Stommel problem.  $5 \times 5$  subdomains at  $t = 360$  days.

given by:

$$\begin{aligned}
 h(x, y, t_0) &= H + A \exp\left(-\frac{(x - x_0)^2 + (y - y_0)^2}{2\sigma^2}\right) & (73) \\
 u(x, y, t_0) &= 0 \\
 v(x, y, t_0) &= 0 & (74)
 \end{aligned}$$

where  $(x, y) \in [0, 1] \times [0, 1]$ ,  $g = 9.8$ ,  $H = 1$ ,  $A = 0.2$ ,  $x_0 = y_0 = 0.5$ , and  $\sigma = 0.05$ . The flow is driven by gravity as in the standing wave problem. Simulations are run for  $t \in [0, 0.2]$ , i.e., up to a short time after the first reflection of the initial wave from the domain boundaries where reflecting boundary conditions are applied.

Figure 11 shows results for conservation properties in the case of a  $5 \times 5$  subdomains. In terms of mass conservation, it is difficult to discern which method offers superior performance. The energy conservation properties of each method improve with increasing  $N$ . At a given value of  $N$ , the DGM is found to produce a slightly larger relative error in terms of the total final energy. Note that for the time for which the simulations were run, no filtering was needed to preserve numerical stability at all values of  $N$  and subdomain thicknesses considered. Nevertheless, the smoothness of the solution is damaged at later times, as weak spurious wiggles emerge. As in the case of the non-linear Stommel problem, in the DGM, the dissipation term in the Rusanov flux stabilizes the solution while keeping it free of spurious oscillation, although somewhat adversely impacting the energy conservation properties of the method. The role of spectral filtering and dissipative terms on the conservation properties for the DGM is further discussed in §5.1.

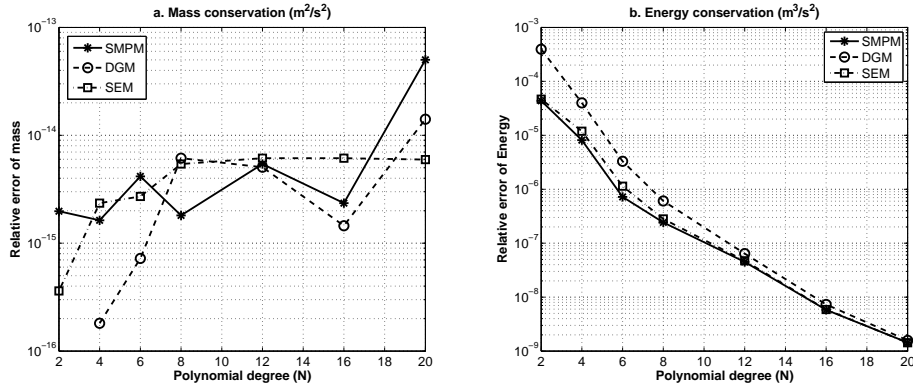


Figure 11: Nonlinear Riemann problem for  $5 \times 5$  subdomains at  $t = 0.2$ . Panels (a) and (b) are the same as Fig. 8

## 5. Discussion

### 5.1. Mass and Energy Conservation

All three methods are found to have very good conservation properties, a direct result of their formulation, see e.g. [19] for SMPM, [22] for DGM, and [52] for SEM. The DGM conserves mass up to machine precision. The SMPM is found to lose mass over long model times with the corresponding relative error as much as four orders of magnitude larger than that for the DGM. This error increases with number of time steps. Such observations might initially seem perplexing, given the analytical demonstration of Don et al. [19] that the averaging method-based penalty scheme is conservative. For all SMPM-driven test-cases we have found that the mass loss (not shown here) is a linear function of time, with a decay rate that is of the order of machine epsilon. The linear Stommel problem has a total mass loss that reaches values of  $10^{-10}$  at higher  $N$ , a value even higher than that observed for the standing and Kelvin wave test-cases where 1/50th the maximum time step is used. This difference is simply because  $10^6$  time steps are required for the linear Stommel problem to reach steady-state. Consequently, we attribute the observed loss of mass to an accumulation of round-off error.

The energy conservation properties of all three methods improve with increasing  $N$ , although both SMPM and DGM are found to be inferior in this regard to the SEM. Note that in simulations where no energy sink terms (such as bottom friction in the Stommel problems) are present, the performance of the discontinuous techniques in terms of energy conservation can be strongly influenced by spectral filtering and the structure of the numerical flux terms, such as the dissipative term used within the Rusanov flux. Figure 12 shows the differences in conservation of mass and energy in the DGM, for the Riemann problem, when spectral filtering, through a 10th-order Boyd-Vandeven

filter [50], is added to the simulation or the dissipation term is neglected in the numerical flux. The absence of both the dissipative term in the numerical flux and spectral filtering provides for the best energy conservation properties, although such behavior does not necessarily guarantee a smooth and stable solution for such a strongly nonlinear problem.

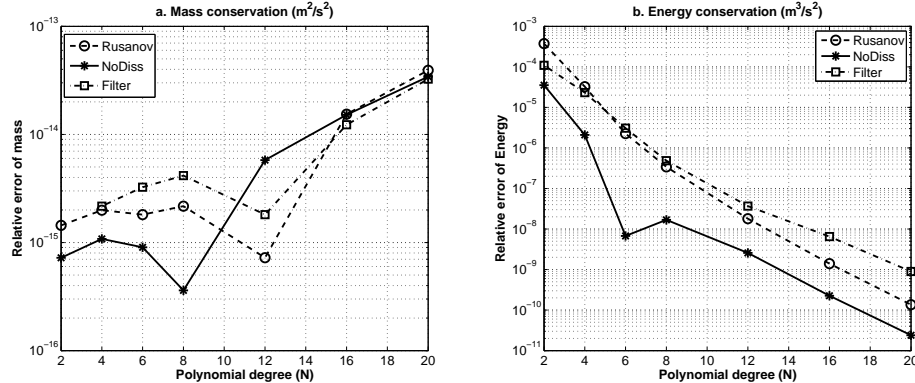


Figure 12: Comparison of conservation properties of the DGM for the Riemann problem. Results for  $5 \times 5$  at  $t = 0.2$ . (a) Mass conservation. (b) Energy conservation.

## 5.2. Effect of time step on convergence and conservation properties

For the purpose of demonstrating that the temporal discretization error does not dominate over the spatial error, we now perform an analysis of the effect of time step,  $\Delta t$ , size, on the convergence and conservation properties of each of the three methods. The base time step corresponds to that associated with a simulation with Courant Number of 0.4.  $\Delta t$  is then progressively decreased by a factor of 2, 10 and 50 (denoted by  $D2$ ,  $D10$ ,  $D50$  respectively). In figure 13 the convergence plots for the free surface elevation  $h$  of the Standing wave test case are presented for all three methods. For a given  $N$ , the increase in accuracy of all three methods is visible as  $\Delta t$  is decreased. Once a factor of 50 reduction is reached exponential convergence is obtained until  $N = 8$ .

The same exercise has been performed to assess the role of time-step on mass and energy conservation in all three methods. The results show (see figures 14 and 15) that the SMPM mass loss increases with decreasing  $\Delta t$ . This observation is consistent with the loss, at a linear decay rate of order machine epsilon, in the SMPM discussed in §5.1. In contrast, the DGM and SEM conserve mass to the order of machine epsilon regardless of the value of  $\Delta t$ . On the other hand, conservation of energy is improved by the three methods once the polynomial degree increases or the size of  $\Delta t$  decreases.

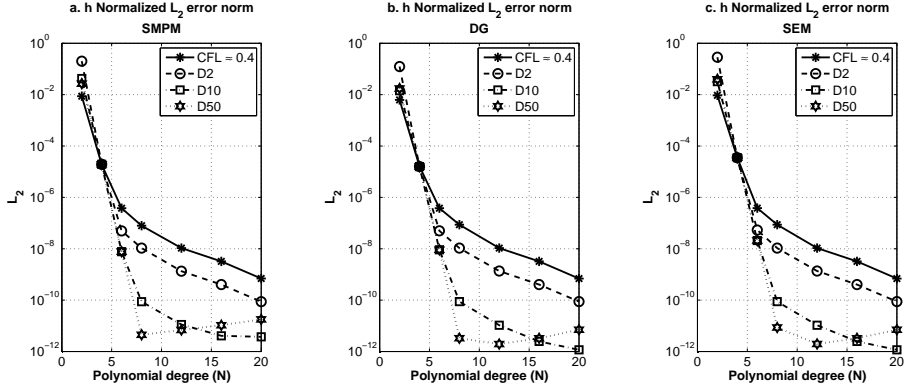


Figure 13: Convergence plots for the Standing wave problem when different  $\Delta t$  sizes are used.

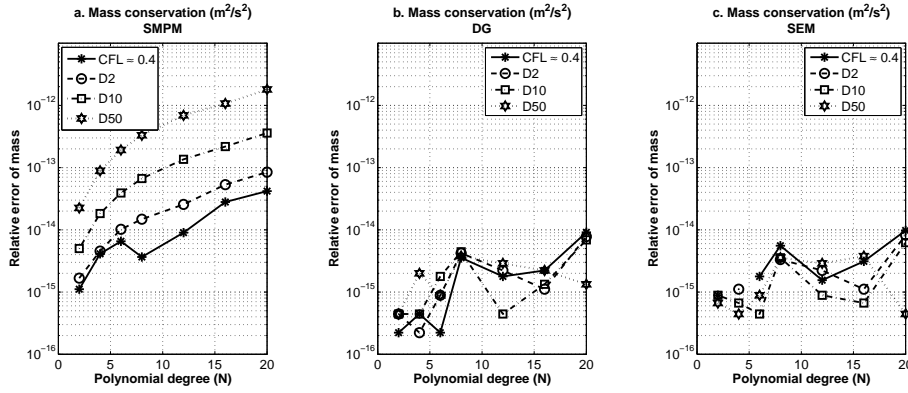


Figure 14: Conservation of mass for the Standing wave problem when different  $\Delta t$  sizes are used.

### 5.3. Effect of Filtering

In the interfacial treatment of the SMPM, there is no dissipative term that removes spurious high wavenumber oscillations that develop in highly nonlinear simulations. Thus, spectral filtering is needed when such simulations are run for long integration times, namely when sharp localized features emerge in the simulations (e.g., nonlinear Riemann problem) or even when the structure of the solution is apparently smooth and free of any localized features (e.g., nonlinear Stommel problem). In contrast, in the case of the DGM, the dissipation term introduces a dissipation mechanism that stabilizes the solution and renders it oscillation-free; for a very simple flow problem, this term reduces to a simple upwinding scheme. By neglecting it, the DGM-generated solution also becomes unstable. Without resorting to recasting the nonlinear terms in skew-symmetric

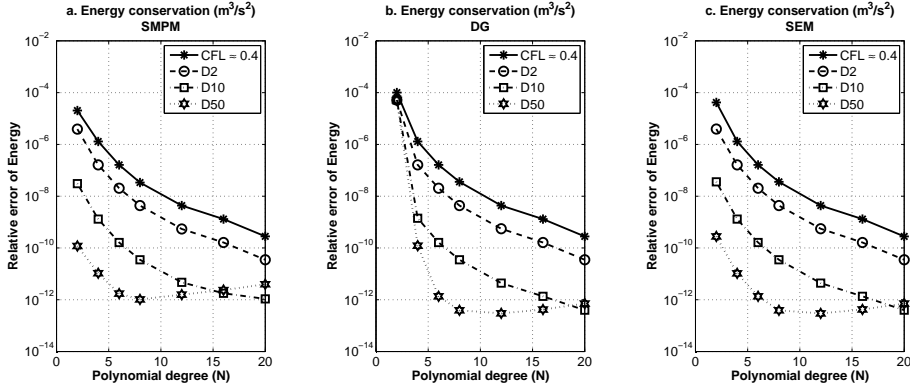


Figure 15: Conservation of energy for the Standing wave problem when different  $\Delta t$  sizes are used.

form [7] and in the absence of an over-integration-based de-aliasing strategy [53] (both which are out of the scope of this paper), spectral filtering is required to recover stability. In terms of mass and energy conservation, the performance of the DGM appears to be very similar when spectral filtering and no dissipative term is used or when only the dissipative term is used (Fig. 12).

The performance of both SMPM and DGM is further examined in problems where significantly sharp features are present. The dam-break problem [9] is simulated with a cylindrical step-function of the free surface elevation as an initial condition and with  $(x, y) \in [-20, 20] \times [-20, 20]$  and  $t \in [0, 0.1]$ . The effect of filtering (with a Boyd-Vandeven filter of  $p = 10$ ) and the dissipative term on the solution are shown in Fig. 16.

In the absence of a dissipative term in the DGM and any spectral filtering for both methods (panels a and d), spurious oscillations are localized in the vicinity of subdomain interfaces for the SMPM, whereas, in the DGM, these oscillations are more evenly distributed throughout the computational domain. When spectral filtering is applied to both methods (panels b and e), the oscillations are strongly damped in the subdomain interior where the effect of the filter is focused [11]. Nevertheless, some weaker oscillations remain at the subdomain interfaces [11]. If no spectral filtering is applied but an additional dissipative term is added to the penalty term in the SMPM (panel c), the solution has a near identical structure with the one computed by the DGM with the full Rusanov flux. For the purpose of comparison, Fig. 17 shows the filtered solution obtained from the SEM which is contrasted to its filtered counterparts (no Rusanov flux term present) computed from DGM and SMPM (Figs. 16b and e). The results for SEM with filtering show stronger spurious oscillations than SMPM or DGM with dissipation or spectral filtering.

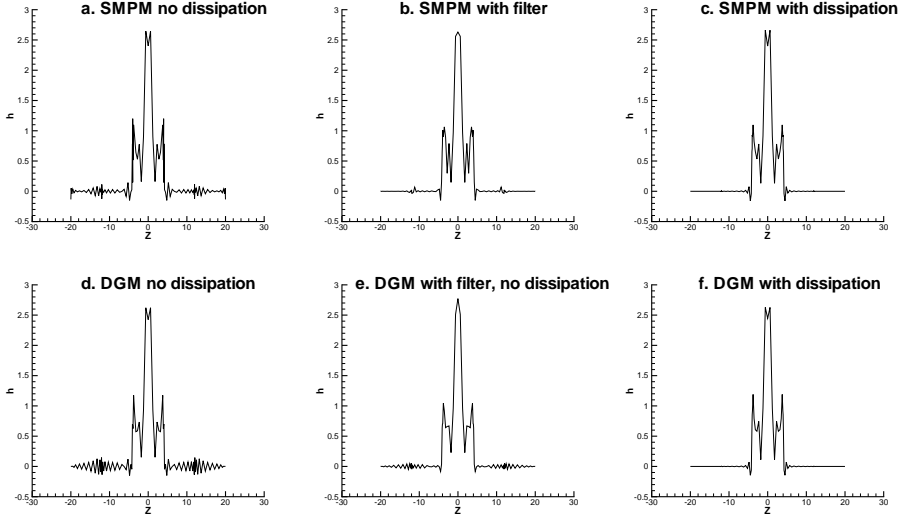


Figure 16: Cross section of the Dam-break problem for  $5 \times 5$  subdomains, and  $N = 20$  at  $t = 0.1$ . (a) SMPM without dissipation. (b) SMPM with filtering (Filter order  $p = 10$ ). (c) SMPM with dissipative term. (d) DGM without dissipative term ( $|\lambda|(\mathbf{q}^R - \mathbf{q}^L)$ ). (e) DGM with filtering (Filter order  $p = 10$ ), and without dissipation term. (f) DGM full Rusanov flux.

#### 5.4. Computational Efficiency and implementation

For all test cases, the order of magnitude of the CPU time per time step has been found to be comparable for both DGM and SMPM and increases when the number of degrees of freedom increases due to  $h$  or  $p$  refinement. Figure 18a shows the computational time for all three methods considered in this manuscript (SMPM, DGM and SEM) for different values of  $N$  for the Riemann problem with  $5 \times 5$  subdomains and the same time step value for each method, corresponding to Courant Number = 0.4.

Figure 18b shows the time needed to advance a simulation to the same final time as Fig. 18a, where the Courant Number is set to the empirically computed maximum value that enables a stable simulation for each method. SEM simulations are found to support double the maximum Courant Number value of DGM and SMPM and are thus twice as fast. DGM and SMPM simulations were also performed with a Courant Number value slightly above the empirically obtained stable limit value. In this case, DGM was found to destabilize faster than SMPM.

Theoretical justification for these observations is gained by examining the eigenvalue spectra of the discretized 1-D linear advection operator for each of the three discretization methods for a periodic domain with 5 subdomains and  $N = 4$  (Fig. 19). In the absence of the dissipation term in DGM, and as expected, all three methods have purely imaginary eigenvalues. The extreme eigenvalues of DGM are roughly 25% larger than their SMPM counterparts and

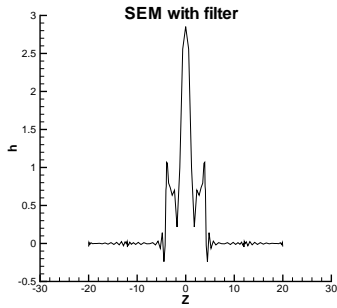


Figure 17: Cross section of the filtered Dam-break problem for  $5 \times 5$  subdomains, and  $N = 20$  at  $t = 0.1$  (Filter order  $p = 10$ ) for SEM.

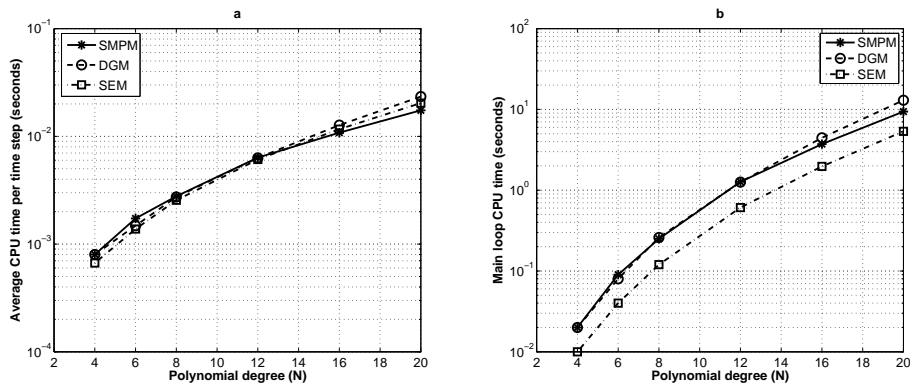


Figure 18: CPU time for the Riemann problem.  $5 \times 5$  subdomains with different polynomial orders at  $t = 0.2$ . (a) All methods with Courant Number = 0.4 and (b) DGM and SMPM with Courant Number = 0.4 and SEM with Courant Number = 0.8.

double the corresponding SEM eigenvalues. Incorporation of the numerical flux term in DGM gives rise to eigenvalues with a negative real part which equip the numerical solution with the necessary numerical dissipation. Moreover, the separation between the eigenvalues with the largest absolute imaginary values is reduced with respect to the case without dissipation but is still slightly larger than that in SMPM and almost double that of SEM. Taking into account the stability region of the SSP-RK34 scheme (which is stable along the imaginary axis) for Courant numbers below this eigenvalue separation can explain why SEM can attain double the Courant Number of DGM and SMPM and why DGM explodes a little faster than SMPM for a marginally unstable time step.

In terms of implementation, in the context of the SWE, both the SMPM and DGM can be written as a system of time-dependent ordinary differential equations where the vector of unknowns is the solution vector at the grid points [54]. In the matrix-vector product that appears on the right hand side of this system



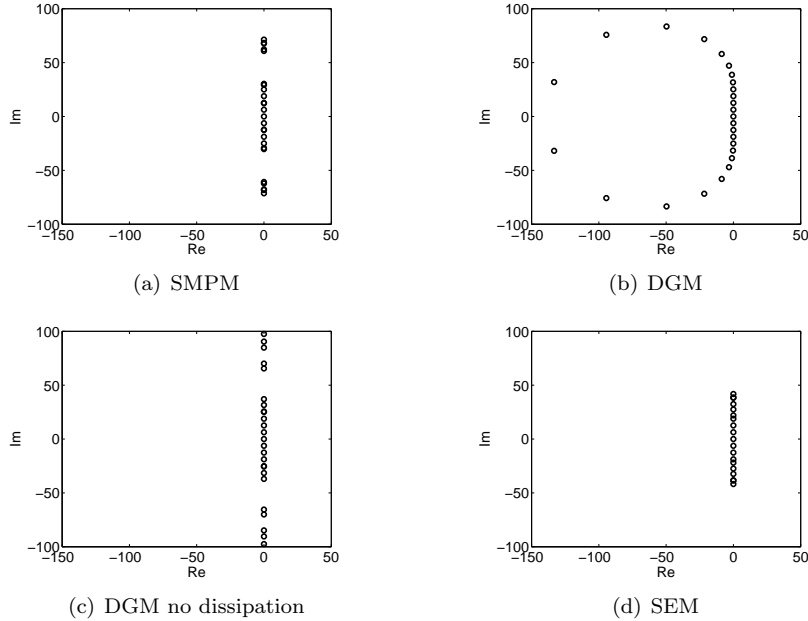


Figure 19: Eigenvalue distribution of the 1D discrete linear advection operator ( $u \frac{d}{dx}$ ) for all three methods, with an advective velocity of  $u = 1$ . In all cases  $x \in [-1, 1]$ , 5 subdomains,  $N = 4$ .

of equations, the associated matrix is simply a spectral differentiation matrix (Eq. (19)) for the SMPM due to its underlying collocation method framework with any modifications to this matrix incurred through communication with points on the edge of the neighboring subdomain. Similar modifications on account of the numerical flux term enter the construction of the corresponding right hand side matrix for the DGM, the core of which is built through additional numerical integration and, therefore, cost. This cost is, nevertheless, offset over the course of a long unsteady simulation. In summary, for hyperbolic systems of equations, the cost of SMPM and DGM are very similar. However, we expect the SMPM to have an advantage when elliptic operators are introduced since the addition of a Laplacian for the SMPM becomes simply a matter of introducing a Laplacian differentiation matrix whereas in DGM either local discontinuous Galerkin or interior penalty methods have to be introduced [55, 56, 57]. For SEM, the addition of Laplacian operators introduces only a slight cost.

## 6. Conclusions

The performance and properties of two commonly used high-order-accuracy element-based spatial discretization methods, spectral multidomain penalty (SMPM) and discontinuous Galerkin (DGM), are examined in the framework of the inviscid shallow water equations (SWE). Whereas a previous comparison study [31]

focused on one-dimensional conservation laws and considered a modally-based Galerkin formulation of SMPM and DGM, this paper applies both techniques to a system of nonlinear conservation equations and considers them in the more frequently used nodal form, in a collocation and Galerkin formulation, for SMPM and DGM, respectively. The two methods are applied to a suite of test cases that are of interest in oceanic shallow water flow: three linear (standing wave, Kelvin wave and linear Stommel problem) and three non-linear (Rossby soliton, nonlinear Stommel problem and Riemann problem). The analysis shows that the methods can be simplified to be the same method when specific choices of the penalty terms (for the SMPM) and numerical flux (for the DGM) and when the same collocation points are chosen for representing the discrete solution. The numerical solutions showed that the methods are extremely similar not only in achieving the same rate of convergence but also in their conservation of energy measures. The key difference between the SMPM and DGM is in their choice of penalty terms that enforce weak boundary conditions across element interfaces. The SMPM has much flexibility in selecting these terms whereas the DGM method is more rigid in its choices in the sense that a Riemann solver must be used; however, this idea offers much flexibility in handling a large variety of flows including those requiring wetting and drying algorithms, for example. Both methods can be used on fully-unstructured quadrilateral element grids but it is not clear how to extend the SMPM to unstructured triangular elements; in contrast, the formulation of the DGM is quite natural and can be extended to triangles rather straightforwardly, assuming that a good set of interpolation and integration points are known (see, e.g., [21, 23, 25, 6]). The SMPM proved to be slightly more efficient than the DGM, in terms of computational time, and we expect this trend to continue as Laplacian operators (as required by Navier-Stokes or even by more realistic shallow water ocean modeling simulations) are introduced.

To observe further differences between both methods, test cases with complex geometries, non-smooth solutions or additional forcing terms have to be executed with the methods. Additionally, parabolic and elliptic partial differential equations have to be assessed in the context of compressible and incompressible flows, where more challenging numerical difficulties appear for the implementation of both methods.

### Acknowledgements

We thank an anonymous reviewer for a number of insightful fundamental comments on the original version of the paper. The work of the first two authors was sponsored by National Science Foundation CAREER award grant OCE-0845558. F.X. Giraldo gratefully acknowledges the support of ONR through program element PE-0602435N.

## References

- [1] A. E. Gill, *Atmosphere-Ocean Dynamics*, Academic Press, San Diego, 1982.
- [2] P. A. Madsen, D. R. Fuhrman, H. A. Schaffer, On the solitary wave paradigm for tsunamis, *J. Geophys. Res.* 112 (2008) doi:10.1029/2008JC004932.
- [3] S. A. Thorpe, Fronts formed by obliquely reflecting internal waves at a sloping boundary, *J. Phys. Oceanogr.* 29 (1999) 2462–2467.
- [4] K. R. Helfrich, W. K. Melville, Long non-linear internal waves, *Ann. Rev. Fluid Mech.* 38 (2006) 395–425.
- [5] M. O. Deville, P. F. Fischer, E. H. Mund, *High Order Methods for Incompressible Fluid Flow*, Cambridge University Press, 2002.
- [6] J. S. Hesthaven, T. Warburton, *Nodal Discontinuous Galerkin Methods*, Springer-Verlag, 2008.
- [7] J. P. Boyd, *Chebyshev and Fourier Spectral Methods*, Dover, Mineola, New York, 2001.
- [8] L. H. Kantha, C. A. Clayson, *Numerical Models of Oceans and Oceanic Processes*, Academic Press, San Diego, 2000.
- [9] E. F. Toro, *Shock-capturing methods for Free-surface shallow flows*, John Wiley and Sons, 2001.
- [10] P. F. Fischer, D. Gottlieb, On the optimal number of subdomains for hyperbolic problems on parallel computers, *J. Supercomput. Appl.* 11 (1997) 65–76.
- [11] D. Gottlieb, J. S. Hesthaven, Spectral methods for hyperbolic problems, *Journal of Computational and Applied Mathematics* 128 (2001) 83–131.
- [12] P. J. Diamessis, J. A. Domaradzki, J. S. Hesthaven, A spectral multidomain penalty method model for the simulation of high Reynolds number localized stratified turbulence, *J. Comp. Phys.* 202 (2005) 298–322.
- [13] D. A. Kopriva, A spectral multidomain method for the solution of hyperbolic systems, *Applied Numerical Mathematics* 2 (1986) 221–241.
- [14] D. A. Kopriva, Computations of hyperbolic equations on complicated domains with patched and overset Chebyshev grids, *SIAM J. Sci. Stat. Comput* 10 (1989) 120–132.
- [15] D. A. Kopriva, Multidomain spectral solution of compressible viscous flows, *J. Comp. Phys.* 115 (1994) 184–199.

- [16] J. S. Hesthaven, D. Gottlieb, A stable penalty method for the compressible Navier-Stokes equations: I. Open boundary conditions, *SIAM J. Sci. Comput.* 17 (1996) 579–612.
- [17] J. S. Hesthaven, A stable penalty method for the compressible Navier-Stokes equations: II. One-dimensional domain decomposition schemes, *SIAM J. Sci. Comput.* 18 (1997) 658–685.
- [18] J. S. Hesthaven, A stable penalty method for the compressible Navier-Stokes equations: III. Multidimensional domain decomposition schemes, *SIAM J. Sci. Comput.* 20 (1998) 62–93.
- [19] W. S. Don, D. Gottlieb, J. H. Jung, A multidomain spectral method for supersonic reactive flows, *J. Comp. Phys.* 192 (2003) 325–354.
- [20] F. X. Giraldo, J. S. Hesthaven, T. Warburton, Nodal high-order discontinuous Galerkin methods for the spherical shallow water equations, *J. Comp. Phys.* 181 (2002) 499–525.
- [21] F. X. Giraldo, Higher-order triangle-based discontinuous Galerkin methods for hyperbolic equations on a sphere, *J. Comp. Phys.* In Press (2006).
- [22] F. X. Giraldo, M. Restelli, A study of spectral element and discontinuous Galerkin methods for mesoscale atmospheric modeling: Equation sets and test cases, *J. Comp. Phys.* 227 (2008) 3849–3877.
- [23] F. X. Giraldo, T. Warburton, A high-order triangular discontinuous Galerkin oceanic shallow water model, *Int. J. Numer. Meth. Fluids* 56 (2008) 899–925.
- [24] F. X. Giraldo, M. Restelli, High-order semi-implicit time-integrators for a triangular discontinuous galerkin oceanic shallow water model, *Int. J. Num. Meths. Fluids* 63 (2010) 1077–1102.
- [25] M. Läuter, F. X. Giraldo, D. Handorf, K. Dethloff, A discontinuous Galerkin method for the shallow water equations in spherical triangular coordinates, *Journal of Computational Physics* 227 (2008) 10226–10242.
- [26] M. Restelli, F. Giraldo, A conservative discontinuous galerkin semi-implicit formulation for the navier-stokes equations in nonhydrostatic mesoscale modeling, *SIAM J. Sci. Comp.* (2009).
- [27] S. A. Orszag, Spectral methods for problems in complex geometries, *J. Comp. Phys.* 37 (1980) 70–92.
- [28] P. J. Diamessis, L. G. Redekopp, Numerical investigation of solitary internal wave-induced global instability in shallow water benthic boundary layers, *J. Phys. Oceanogr.* 36 (2006) 784–812.

- [29] P. J. Diamessis, G. R. Spedding, J. A. Domaradzki, Similarity scaling and vorticity structure in high Reynolds number stably stratified turbulent wakes, *J. Fluid Mech.* 671 (2011) 52–95.
- [30] A. M. Abdilghanie, P. J. Diamessis, On the generation and evolution of numerically simulated large-amplitude internal gravity wave packets, *Theor. Comput. Fluid Dyn.* (2011) Available Online.
- [31] S. Gottlieb, J. H. Jung, Numerical issues in the implementation of high order polynomial multi-domain penalty spectral Galerkin methods for hyperbolic conservation laws, *Commun. Comput. Phys.* 5 (2009) 600–619.
- [32] J. H. Jung, B. D. Shizgal, On the numerical convergence with the inverse polynomial reconstruction method for the resolution of the gibbs phenomenon, *J. Comp. Phys.* 224 (2007) 477–488.
- [33] A. P. Engsig-Karup, J. S. Hesthaven, H. Bingham, T. Warburton, Dg-fem solution for nonlinear wave-structure interaction using boussinesq-type equations, *Coastal Engineering* 55 (2008) 197–208.
- [34] D. Alevras, Simulation of the Indian Ocean Tsunami with Real Bathymetry using a high-order triangular Discontinuous Galerkin shallow water model, Master’s thesis, Naval Postgraduate School, 2009.
- [35] P. K. Kundu, I. M. Cohen, *Fluid Mechanics*, Academic Press, San Diego, 2004.
- [36] J. D. Anderson, *Computational Fluid Dynamics*, McGraw Hill International Editions, 1995.
- [37] B. Costa, W. S. Don, On the computation of high order pseudospectral derivatives, *Appl. Num. Math.* 33 (2000) 151–159.
- [38] J. H. Jung, Multi-domain spectral penalty method for hyperbolic systems: Theory and applications, Ph.D. thesis, Brown University, 2002.
- [39] M. H. Carpenter, D. Gottlieb, C.-W. Shu, On the conservation and convergence to weak solutions of global schemes, *J. Sci. Comput.* 18 (2003) 111–132.
- [40] B. Cockburn, C.-W. Shu, Runge-kutta discontinuous galerkin methods for convection-dominated problems, *Journal of Scientific Computing* 16 (2001) 173–261.
- [41] R. J. Spiteri, S. J. Ruuth, A new class of optimal high-order strong-stability-preserving time discretization methods, *SIAM J. Numer. Anal.* 40 (2002) 469–491.
- [42] F. X. Giraldo, The Lagrange-Galerkin spectral element method on unstructured quadrilateral grids, *Journal of Computational Physics* 147 (1998) 114–146.

- [43] F. X. Giraldo, A spectral element shallow water model on spherical geodesic grids, *International Journal for Numerical Methods in Fluids* 35 (2001) 869–901.
- [44] F. X. Giraldo, T. E. Rosmond, A scalable spectral element eulerian atmospheric model (see-am) for nwp: Dynamical core tests, *Monthly Weather Review* 132 (2004) 133–153.
- [45] M. Iskandarani, D. B. Haidvogel, J. P. Boyd, A staggered spectral element model with application to the oceanic shallow water equations, *International Journal for Numerical Methods in Fluids* 20 (1995) 393–414.
- [46] C. Eskilsson, S. Sherwin, A triangular spectral/hp discontinuous Galerkin method for modelling 2d shallow water equations, *Int. J. Numer. Meth. Fluids* 45 (2004) 605–623.
- [47] H. Stommel, The westward intensification of wind-driven ocean currents, *Transactions of the American Geophysical Union* 29 (1948) 202–206.
- [48] D. B. Haidvogel, A. Beckmann, *Numerical Ocean Circulation Modelling*, Imperial College Press, Cambridge, 1999.
- [49] J. P. Boyd, Equatorial solitary waves. part 3: Westward-travelling modons, *J. Phys. Oceanogr.* 15 (1985) 46–54.
- [50] J. G. Levin, M. Iskandarani, D. B. Haidvogel, A spectral filtering procedure for eddy-resolving simulations with a spectral element ocean model, *J. Comp. Phys.* 137 (1997) 130–154.
- [51] C. Moler, *Experiments with MATLAB*, The Mathworks, 2009.
- [52] M. A. Taylor, A. Fournier, A compatible and conservative spectral element method on unstructured grids, *Journal of Computational Physics* 229 (2010) 5879–5895.
- [53] R. Kirby, G. Karniadakis, De-aliasing on non-uniform grids: algorithms and applications, *J. Comp. Phys.* 191 (2006) 249–264.
- [54] D. A. Kopriva, *Implementing Spectral Methods for Partial Differential Equations: Algorithms for Engineers and Scientists*, Springer Verlag, 2009.
- [55] P. Castillo, B. Cockburn, I. Perugia, D. Schötzau, An a priori error analysis of the local discontinuous Galerkin method for elliptic problems, *SIAM J. Num. Anal.* 38 (2000) 1676–1706.
- [56] D. N. Arnold, F. Brezzi, B. Cockburn, D. Marini, Unified analysis of discontinuous Galerkin methods for elliptic problems, *SIAM J. Num. Anal.* 39 (2002) 1749–1779.
- [57] B. Cockburn, B. Dong, J. Guzmán, A superconvergent LDG-hybridizable Galerkin method for second-order elliptic problems, *Math. Comp.* 77 (2008) 1887–1916.