



HAL
open science

A super-convergent quasi-discontinuous Galerkin method for long time behavior of transport-dominated problems in geophysical flows

Daniel Le Roux

► **To cite this version:**

Daniel Le Roux. A super-convergent quasi-discontinuous Galerkin method for long time behavior of transport-dominated problems in geophysical flows. 2024. hal-04447547

HAL Id: hal-04447547

<https://hal.science/hal-04447547>

Preprint submitted on 8 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A super-convergent quasi-discontinuous Galerkin method for long time behavior of transport-dominated problems in geophysical flows

Daniel Y. Le Roux^{a,*}

^a*Université Lyon 1, CNRS, Institut Camille Jordan, 43, blvd du 11 novembre 1918, 69622 Villeurbanne Cedex, France*

Abstract

The numerical approximation of the shallow-water equations, which support geophysical wave propagation, namely the fast external Poincaré and Kelvin waves and the slow large-scale planetary Rossby waves, is a delicate and difficult problem. Indeed, the coupling between the momentum and the continuity equations may lead to the presence of erratic or spurious solutions for these waves, *e.g.* the spurious pressure and inertial modes. In addition, the presence of spurious branches and the emergence of spectral gaps in the dispersion relation at specific wavenumbers may lead to anomalous dissipation/dispersion in the representation of both fast and slow waves. The aim of the present study is to propose a class of possible discretization schemes, via the discontinuous Galerkin method, that is not affected by the above-mentioned problems. A Fourier/stability analysis of the 2D shallow-water model is performed by analysing the finite volume method and the linear discontinuous (DG) and non conforming (NC) Galerkin discretizations. These are stabilized by means of a family of numerical fluxes via the Polynomial Viscosity Matrix approach. A long time stability result is proven for all schemes and fluxes examined here. Further, a super-convergent result is demonstrated for the discrete frequencies of the NC method compared to the DG one, except for the slow mode in the Roe flux case. Indeed, the Roe flux yields spurious frequencies and sub-optimal rates of convergence for the slow mode for both the DG and NC methods. Finally, numerical solutions of linear and non-linear test problems confirm the theoretical results and reveal the computational cost efficiency of the NC approach.

Keywords: Geophysical fluid dynamics, Shallow water equations, Discontinuous Galerkin methods, Fourier/stability analysis, Inertia-gravity waves, Rossby waves

1. Introduction

The shallow-water (SW) equations have been extensively employed in environmental studies to model hydrodynamics in estuaries, lakes, coastal regions, tides, atmospheric and oceanic motions. They are derived from the Navier-Stokes system by vertical integration under Boussinesq and hydrostatic pressure assumptions. Although the SW system is simplified it is often used as a prototype of the primitive equations and frequently employed as a benchmark for numerical schemes to be used in more complex oceanic or atmospheric models. Indeed, it retains much of the dynamical complexity of three-dimensional flows on the rotating Earth in representing waves that govern most of the atmospheric and oceanic circulations: the fast external Poincaré and Kelvin waves, mainly driven by gravity, and the large-scale planetary Rossby waves travelling much slower as they react on variations in the Coriolis force with the latitude. As for the primitive system, attempts to apply numerical methods to the solution of the SW equations were confronted with a number of difficulties arising out of inconsistencies in the discrete model relative to the continuous problem. This is especially true on meshes made up of triangular elements, which are considered in this paper.

A first difficulty has been the occurrence of computational modes in the pressure (surface-elevation) field, as is the case for the finite-difference (FD) A-grid. Indeed, the mixed Galerkin finite-element (FE) $P_1 - P_1$

*Corresponding Author: Daniel Y. Le Roux (dleroux@math.univ-lyon1.fr)

pair, where velocity components and surface-elevation are represented as piecewise-defined polynomials of degree 1, is plagued by spurious pressure oscillations preventing the uniqueness of the discrete pressure [3, 21, 25]. The appearance of such oscillations mainly results from an inappropriate placement of variables on the mesh and/or a bad choice of approximation function spaces. This difficulty with mixed methods is not specific to the SW problem alone but is also well known in the context of the Navier-Stokes system when the so-called discrete inf-sup or LBB condition is not satisfied [3]. The spurious pressure oscillations reflect the presence of erratic eigenvectors which correspond to physical eigenmodes of the system which have their phase speed reduced to zero by the numerical method. These appear as stationary internode oscillations and lead to an accumulation of energy in the smallest-resolvable scale [41]. Two main options have been proposed to remove the troublesome spurious pressure modes in the FE community: firstly, using mixed-order FE interpolation methods or equal-order elements with staggered variables and secondly, developing stabilization procedures. Unfortunately, both options entail a number of difficulties.

For discrete velocities having two components per node, the first option [17, 19] implies much more discrete vector momentum equations than discrete continuity equations, yielding the appearance of purely inertial modes [7, 24, 26], for which the discrete velocity rotates with frequency f (the Coriolis parameter). Such modes trigger noise in the discrete velocity field with suboptimal rates of convergence in the case of inviscid flows [5]. Unfortunately, the attempt of filtering the spurious inertial solutions unacceptably degrades the accuracy of the Rossby modes for long-term simulations [6, 26]. For FE pairs having normal velocity components on element edges belonging to $H(\text{div})$, *e.g.* the RT [33], BDM and BDFM [3] elements, spurious inertial modes are not present. This class of FE methods can be seen as generalization of the FD C-grid. However, as for the C-grid, the discrete Coriolis matrix is rank-deficient for these pairs due to the required spatial averaging of the Coriolis terms [36], yielding the so-called spurious f -modes [35, 42]. As for the C-grid, a sufficiently fine mesh with respect to the deformation radius or suitable strategies [11, 43] permit to attenuate the amplitude of the f -modes. The choice of the FE space for pressure is crucial when the discrete velocities belong to the RT, BDM and BDFM spaces. Indeed, *e.g.* the $RT_0 - P_1$ and $BDM_1 - P_1$ pairs are plagued by spurious pressure oscillations [36]. In fact, the discrete method should mimic fundamental properties of mathematical and physical systems including conservation laws, exact mathematical identities of the vector and tensor calculus, symmetry and positivity of solutions etc. Examples of mimetic properties for the SW equations was first proposed in [35] for classical FD grids and FE pairs, and then precised in [7, 8] for FE methods that satisfy the conditions of FE exterior calculus, summarized in [2], using the sequence of spaces and mappings, called the de Rham complex, for velocities belonging to $H(\text{div})$. Unfortunately, the desirable mimetic properties are not sufficient to ensure that compatible Galerkin methods have good wave dispersion properties. In particular, two types of problems in the discrete dispersion relation, relating the wavelength or wavenumber of a wave to its frequency, are known to arise: spurious (extra) branches and spectral gaps. For example, the $RT_n - P_n^{DG}$ and $BDM_n - P_n^{DG}$ pairs (n being a integer) have spurious inertia-gravity and Rossby branches, respectively, while the $BDFM_1 - P_1^{DG}$ pair having an exact 2 : 1 ratio of velocity degrees of freedom to pressure degrees of freedom is free of both inertia-gravity and Rossby spurious branches [7]. Spectral gaps correspond to wavenumbers for which the dispersion relation is multivalued, namely piecewise continuous, and the group velocity, which crucially controls energy propagation, drops to zero [13, 27, 30]. When a wave packet has significant energy at the wavenumbers approaching the spectral gaps, it will have anomalous dispersion and incorrect propagation. The exact cause of spectral gaps is still lacking. Regrettably, the aforementioned three FE pairs lead to the presence of spectral gaps in their 1-D representation [37], and since the discrete frequencies, solution of the dispersion relations, are real there is no intrinsic numerical diffusion to control the gaps.

The second option regroup the generalized wave continuity equation formulations (GWCE), the use of FE bubbles and stabilization techniques. The GWCE [22, 29] is formed by differentiating the continuity equation in time, substituting from the momentum equations, and rearranging terms. However, conservation properties, advection instabilities, lack of consistency with scalar transport schemes and poor accuracy using implicit schemes are a potential challenge for the method. Further, the GWCE is mainly used for coastal simulations since the method is not designed for the computation of slow large-scale planetary Rossby waves. Alternatively, the use of bubbles [1, 3] and stabilization techniques [20] aim to recover the information lost by projecting the discrete pressure gradient in inadequate velocity spaces. This may permit to recover the

expected convergence order for the discrete variables but at the price of introducing an undesirable amount of numerical viscosity, making questionable a sufficiently accurate computation of the slow Rossby modes.

Finally, the discontinuous Galerkin (DG) approach belongs to the group of stabilization techniques, in the sense that a Riemann solver compute the numerical fluxes at element interfaces by upwinding the characteristic variables. The DG method is thus based on a very different stabilization strategy from that adopted in [20]. Due to the potential of the DG approach in geophysical fluid dynamics, there are reasons to achieve a Fourier mode interpretation, via a Fourier analysis, in order to explore the properties of the method. These include the ability to derive the dispersion relation and group velocity, investigate the eventual presence of erratic higher frequency modes and study the long-time stability and accuracy of the discrete frequencies, which is the main objective of the present work. The SW equations are discretized by examining three schemes employing identical approximations for all variables: the low order finite-volume (FV) and the linear discontinuous P_1^{DG} and non conforming P_1^{NC} finite elements. Moreover, the Rusanov, Roe, FORCE and PVM – 4 fluxes are considered in the subsequent analyses via the Polynomial Viscosity Matrix approach in order to stabilize the numerical solutions. The P_1^{NC} FE was first used in a two-layer model to approximate the discrete velocities while a P_1 discretization of the surface-elevation was adopted [19]. It was however demonstrated in [24] that such a pair support spurious inertial frequencies. The linear SW model was later discretized by employing the $P_1^{NC} - P_1^{NC}$ pair with an approximate Riemann solver utilizing the Roe flux [5], but the numerical experiments were performed without analysis.

This work represents a starting point for the development of high order DG schemes for ocean modelling, and as a first step it focuses on initial smooth solutions to compute the waves of interest (inertia-gravity and Rossby). The continuous model is presented in Section 2 and the discretization of the model equations is performed in Section 3 for the FV, P_1^{DG} and P_1^{NC} methods. The Fourier/dispersion analyses are then conducted in Section 4, and the results are analysed. Numerical solutions of linear and non-linear test problems are simulated and discussed in Section 5. Some concluding remarks complete the study in Section 6.

2. Non-linear and linear model problems

Let Ω be a bounded open subset of \mathbb{R}^2 and consider an inviscid layer of constant- and uniform- density fluid. In this paper, attention is focused on the two-dimensional non-linear SW equations, which in Cartesian coordinates are expressed as

$$\frac{\partial \mathbf{q}}{\partial t} + \nabla \cdot \left(\mathbf{q} \otimes \frac{\mathbf{q}}{\xi} \right) + f \mathbf{k} \times \mathbf{q} + g \xi \nabla (\xi - H) = 0, \quad (1)$$

$$\frac{\partial \xi}{\partial t} + \nabla \cdot \mathbf{q} = 0, \quad (2)$$

where $\mathbf{u}(\mathbf{x}, t) = {}^t(u, v)$ is the velocity field with $\mathbf{x} = {}^t(x, y)$, $\mathbf{q} = \xi \mathbf{u} = {}^t(q^x, q^y)$, is the discharge, $\xi(\mathbf{x}, t)$ and $H(\mathbf{x})$ are the thickness of the fluid layer and the depth at rest, respectively, with $\xi = \eta + H$, where $\eta(\mathbf{x}, t)$ is the surface elevation with respect to the reference level $z = 0$, as shown in Fig. 1.

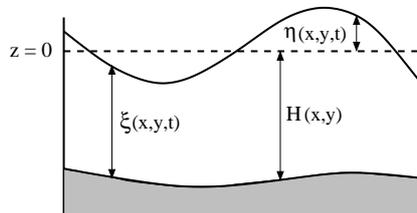


Figure 1: Schematic illustration of the free surface $\eta(\mathbf{x}, t)$ with respect to the reference level $z = 0$, the total depth $\xi(\mathbf{x}, t)$ and the depth at rest or bathymetry $H(\mathbf{x})$.

Further, g is the gravitational acceleration, $\mathbf{k} = {}^t(0, 0, 1)$ is a unit vector in the vertical direction leading to $\mathbf{k} \times \mathbf{q} = {}^t(-q^y, q^x)$. The coordinates x and y are oriented eastward and northward, respectively, and the

latter is measured from a reference latitude φ_0 . The Coriolis parameter f is allowed to vary with latitude, and the β -plane approximation $f = f_0 + \beta y$ is used [32], with $f_0 = 2\Omega_e \sin \varphi_0$ and $\beta = 2(\Omega_e/R_e) \cos \varphi_0$, where Ω_e is the angular frequency of the Earth's rotation and R_e is the Earth's radius.

Let the flux tensor $\mathbf{F}(\mathbf{w})$, the source function $\mathbf{S}(\mathbf{w})$ and the vector of conserved variables \mathbf{w} be

$$\mathbf{F}(\mathbf{w}) = \begin{pmatrix} \frac{(q^x)^2}{\xi} + \frac{1}{2}g\xi^2 & \frac{q^x q^y}{\xi} \\ \frac{q^x q^y}{\xi} & \frac{(q^y)^2}{\xi} + \frac{1}{2}g\xi^2 \end{pmatrix}, \quad \mathbf{S}(\mathbf{w}) = {}^t(-f\mathbf{k} \times \mathbf{q} + g\xi \nabla H, 0), \quad \mathbf{w} = {}^t(\mathbf{q}, \xi), \quad (3)$$

respectively, where $\mathbf{F}(\mathbf{w}) = (\mathbf{F}^x \ \mathbf{F}^y)$, \mathbf{F}^x and \mathbf{F}^y being the first and second columns of \mathbf{F} , respectively, and $\mathbf{S}(\mathbf{w}) = {}^t(S^x, S^y, 0)$. Equations (1) and (2) are then rewritten in the so-called divergence form

$$\frac{\partial \mathbf{w}}{\partial t}(\mathbf{x}, t) + \nabla \cdot \mathbf{F}(\mathbf{w}(\mathbf{x}, t)) = S(\mathbf{w}(\mathbf{x}, t)), \quad (4)$$

and the conservative form (4) describes a first order hyperbolic system. Boundary and initial conditions $\mathbf{w}(\mathbf{x}, t = 0)$ defined for all $\mathbf{x} \in \Omega$, complete the mathematical statement of the problem.

The discretization of (4) is mainly used in Section 5.2 to simulate the propagation of fast inertia-gravity waves and slow Rossby modes in the non-linear case. However, for the purpose of performing the Fourier/dispersion analysis in Section 4 and validate the theoretical results by conducting the numerical simulations of Section 5.1, Eqs. (1) and (2) are linearized about a state of rest by assuming constant Coriolis parameter f and bathymetry H , and we obtain [23]

$$\frac{\partial \mathbf{u}}{\partial t} + f\mathbf{k} \times \mathbf{u} + g\nabla \eta = \mathbf{0}, \quad (5)$$

$$\frac{\partial \eta}{\partial t} + H \nabla \cdot \mathbf{u} = 0. \quad (6)$$

Periodic boundary conditions are employed in Section 3 to derive the discretization of (5) and (6) and in Section 4 to perform the subsequent Fourier analyses, while no-normal flow boundary conditions are used in Section 5 to perform the numerical simulations. For the convenience of the subsequent analysis, we let $\bar{\mathbf{w}} = (u, v, \eta)$, and (5) and (6) are rewritten on the divergence form

$$\frac{\partial \bar{\mathbf{w}}}{\partial t}(\mathbf{x}, t) + \nabla \cdot \bar{\mathbf{F}}(\bar{\mathbf{w}}(\mathbf{x}, t)) = \bar{\mathbf{S}}(\bar{\mathbf{w}}(\mathbf{x}, t)), \quad (7)$$

where $\bar{\mathbf{F}}(\bar{\mathbf{w}})$ and $\bar{\mathbf{S}}(\bar{\mathbf{w}})$ are the flux matrix and the vector containing the source term, respectively, with

$$\bar{\mathbf{F}}(\bar{\mathbf{w}}) = \begin{pmatrix} \bar{\mathbf{F}}^x & \bar{\mathbf{F}}^y \end{pmatrix} = \begin{pmatrix} g\eta & 0 \\ 0 & g\eta \\ Hu & Hv \end{pmatrix} \quad \text{and} \quad \bar{\mathbf{S}}(\bar{\mathbf{w}}) = \begin{pmatrix} fv \\ -fu \\ 0 \end{pmatrix}. \quad (8)$$

The free modes of (5) and (6) are examined by perturbing about the zero basic state. Because the governing equations are linear with constant coefficients and periodic boundary conditions are employed, the solution may be examined by considering the behavior of one Fourier mode. We then seek solutions of (5) and (6) of the form $(u, v, \eta) = (\hat{u}, \hat{v}, \hat{\eta}) e^{i(kx+ly-\omega t)}$, where k and l are the wave numbers in the x - and y -directions, respectively, and ω is the harmonic frequency. Substitution into (5) and (6) leads to the 3×3 matrix system for the amplitudes. A non trivial solution exists if the determinant of the matrix equals zero, and this yields a relationship between the wave numbers k and l and the frequency ω , the so-called dispersion relation

$$\omega (\omega^2 - f^2 - gH(k^2 + l^2)) = 0. \quad (9)$$

Solving (9) leads to three continuous (C) solutions for the frequency

$$\omega_{1,2}^C = \pm \sqrt{f^2 + gH(k^2 + l^2)} \quad \text{and} \quad \omega_3^C = 0. \quad (10)$$

The two first solutions, named $\omega_{1,2}^C$, correspond to the free-surface gravitational modes with rotational correction, namely inertia-gravity or Poincaré waves, while ω_3^C is the geostrophic mode, and it would correspond to the slow Rossby mode on a β -plane. We point out that all modes are neutrally stable and neither amplify nor decay since ω is purely real in (10).

3. Discretization of the non-linear and linear model equations

The DG and NC weak formulations and the FV scheme are first introduced. The general Polynomial Viscosity Matrix method is then presented in order to derive the numerical flux or trace. Finally, the discrete equations are computed at selected nodes for the FV, DG and NC schemes, for the purpose of performing the Fourier analyses in Section 4.

3.1. The DG and NC weak formulations and the FV scheme

In order to describe the weak formulation, a few notations are first defined. Let $\{\tau_h\}_{h>0}$ denote a partition of the domain Ω into a finite number N_{el} of disjoint open non degenerate elements or triangles K_{el} , $el = 1, 2, \dots, N_{el}$, where h is the maximal element diameter. The closure of Ω satisfies $\bar{\Omega} = \cup_{el=1}^{N_{el}} \bar{K}_{el}$ and for $K_{el^+}, K_{el^-} \in \tau_h$ we have $\overset{\circ}{K}_{el^+} \cap \overset{\circ}{K}_{el^-} = \emptyset$ if $K_{el^+} \neq K_{el^-}$, where $\overset{\circ}{K}_{el}$ is the interior of K_{el} . In the case $K_{el^+} \cap K_{el^-} \neq \emptyset$, then $K_{el^+} \cap K_{el^-}$ is either a common face or a common vertex of K_{el^+} and K_{el^-} . Each element K_{el} is supposed to have a Lipschitz boundary ∂K_{el} . Let Γ be the finite ensemble of M_{ib} interelement boundaries $\Gamma_{ed} = \bar{\partial K}_{el^+} \cap \bar{\partial K}_{el^-}$ with $el^+ > el^-$ inside the domain, with all possible combinations

$$\bar{\Gamma} = \cup_{ed=1}^{M_{ib}} \bar{\Gamma}_{ed} \quad \text{and} \quad \Gamma_{ed^+} \cap \Gamma_{ed^-} = \emptyset \quad \text{if} \quad ed^+ \neq ed^-.$$

Each interelement boundary $\Gamma_{ed} \in \Gamma$ is associated with a *unique* outward pointing unit normal vector denoted by $\mathbf{n}_{ed} = (n_{ed}^x, n_{ed}^y)$. Further, for any function $\chi \in V_{K_{el}} := H^1(K_{el})$, where χ represents the variables of the model problem, and for each element K_{el} , the trace of χ on Γ_{ed} is represented by χ^\pm , namely χ_L (L for left) and χ_R (R for right). For $\mathbf{x} \in \Gamma_{ed}$ we have

$$\chi_L(\mathbf{x}) := \chi^-(\mathbf{x}) = \lim_{\epsilon \rightarrow 0^-} \chi(\mathbf{x} + \epsilon \mathbf{n}_{ed}), \quad \chi_R(\mathbf{x}) := \chi^+(\mathbf{x}) = \lim_{\epsilon \rightarrow 0^+} \chi(\mathbf{x} + \epsilon \mathbf{n}_{ed}). \quad (11)$$

The variational or weak formulation is obtained by multiplying (4) by an arbitrary test function $\psi(\mathbf{x}) \in V_{K_{el}}$ and integrating the flux term by parts, using Green's theorem, over each element K_{el} . A boundary term appears which will allow the coupling between two neighbouring elements, yielding

$$\int_{K_{el}} \frac{\partial \mathbf{w}}{\partial t} \psi \, d\mathbf{x} - \int_{K_{el}} \mathbf{F}(\mathbf{w}) \nabla \psi \, d\mathbf{x} + \int_{\partial K_{el}} \mathbf{F}(\mathbf{w}) \mathbf{n}_{ed} \psi \, ds = \int_{K_{el}} \mathbf{S}(\mathbf{w}) \psi \, d\mathbf{x}, \quad \forall \psi \in V_{K_{el}}. \quad (12)$$

Let $\mathcal{P}^r(K_{el})$ be the space of polynomials of degree at most $r \geq 0$, defined on K_{el} and consider the polynomial space $V_{h,K_{el}}$ of piecewise smooth functions that are differentiable over an element, but which allows discontinuities between elements. The space $V_{h,K_{el}}$ will be defined more precisely in Sections 3.3.2 and 3.3.3 for the DG and NC schemes, respectively. The components of \mathbf{w} are approximated in $V_{h,K_{el}}$ and these are denoted by \mathbf{w}_h . Further, we require $\psi \in V_{h,K_{el}}$. The approximate solution $\mathbf{w}_h \in V_{h,K_{el}}$, given by the spatial discretization, is then obtained from (12) as the solution of

$$\sum_{el=1}^N \int_{K_{el}} \frac{\partial \mathbf{w}_h}{\partial t} \psi \, d\mathbf{x} - \sum_{el=1}^N \int_{K_{el}} \mathbf{F}(\mathbf{w}_h) \nabla \psi \, d\mathbf{x} + \sum_{el=1}^N \int_{\partial K_{el}} \mathbf{F}(\mathbf{w}^*) \mathbf{n}_{ed} \psi \, ds = \sum_{el=1}^N \int_{K_{el}} \mathbf{S}(\mathbf{w}_h) \psi \, d\mathbf{x}, \quad (13)$$

where $\mathbf{w}^* = (\mathbf{q}^*, \xi^*)$ denotes the numerical trace of (\mathbf{q}, ξ) on the boundary element ∂K_{el} . Note that the FV scheme is obtained from (13) when $r = 0$, and hence the second term in the left hand side (LHS) of (13) is zero as the test function belongs to $\mathcal{P}^0(K_{el})$.

Following the same approach and letting $\bar{\mathbf{w}}^* = (\mathbf{u}^*, \eta^*)$ be the numerical trace of (\mathbf{u}, η) on ∂K_{el} with $\bar{\mathbf{w}}_h = (u_h, v_h, \eta_h)$, the variational formulation for the linearized system (7) reads

$$\sum_{el=1}^N \int_{K_{el}} \frac{\partial \bar{\mathbf{w}}_h}{\partial t} \psi \, d\mathbf{x} - \sum_{el=1}^N \int_{K_{el}} \bar{\mathbf{F}}(\bar{\mathbf{w}}_h) \nabla \psi \, d\mathbf{x} + \sum_{el=1}^N \int_{\partial K_{el}} \bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed} \psi \, ds = \sum_{el=1}^N \int_{K_{el}} \bar{\mathbf{S}}(\bar{\mathbf{w}}_h) \psi \, d\mathbf{x}. \quad (14)$$

To complete the definition of \mathbf{w}_h and $\bar{\mathbf{w}}_h$ in (13) and (14), it remains to choose unique numerical fluxes $\mathbf{F}(\mathbf{w}^*) \mathbf{n}_{ed}$ and $\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed}$ on interelement boundaries Γ_{ed} as to render the method consistent and stable, which is the subject of the next section. The numerical flux is a single valued function defined on Γ_{ed} and it usually depends on the numerical values of the numerical solution from both sides of the interface.

3.2. Derivation of the numerical trace for the non-linear and linear model problems

Let us first derive the Roe flux in the context of the non-linear problem (4). Assuming the depth at rest H is constant and the Coriolis term is disregarded, Eq. (4) is rewritten in the quasi-linear form

$$\frac{\partial \mathbf{w}}{\partial t} + \frac{\partial \mathbf{F}^x}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial x} + \frac{\partial \mathbf{F}^y}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial y} = 0. \quad (15)$$

The Jacobian matrix $J_{\mathbf{F}}(\mathbf{w})$ of the normal flux $\mathbf{F}(\mathbf{w}) \mathbf{n}_{ed}$ then reads

$$J_{\mathbf{F}}(\mathbf{w}) = \frac{\partial (\mathbf{F}(\mathbf{w}) \mathbf{n}_{ed})}{\partial \mathbf{w}} = \frac{\partial \mathbf{F}^x}{\partial \mathbf{w}} n_{ed}^x + \frac{\partial \mathbf{F}^y}{\partial \mathbf{w}} n_{ed}^y = \begin{pmatrix} un_{ed}^x + \mathbf{u} \cdot \mathbf{n}_{ed} & un_{ed}^y & g\xi n_{ed}^x - u\mathbf{u} \cdot \mathbf{n}_{ed} \\ vn_{ed}^x & vn_{ed}^y + \mathbf{u} \cdot \mathbf{n}_{ed} & g\xi n_{ed}^y - v\mathbf{u} \cdot \mathbf{n}_{ed} \\ n_{ed}^x & n_{ed}^y & 0 \end{pmatrix}, \quad (16)$$

and we deduce from (16) that $J_{\mathbf{F}}(\mathbf{w})$ has three distinct real eigenvalues named λ_j , $j = 1, 2, 3$, with

$$\lambda_1 = \mathbf{u} \cdot \mathbf{n}_{ed} - \sqrt{g\xi}, \quad \lambda_2 = \mathbf{u} \cdot \mathbf{n}_{ed}, \quad \lambda_3 = \mathbf{u} \cdot \mathbf{n}_{ed} + \sqrt{g\xi}. \quad (17)$$

Assume that \mathbf{w}_R and \mathbf{w}_L are the values of \mathbf{w} on the right and left sides of a given interelement boundary Γ_{ed} , respectively. A Roe matrix with constant coefficients depending on \mathbf{w}_R and \mathbf{w}_L , named $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$, is then searched, in order to approximate the Jacobian matrix $J_{\mathbf{F}}(\mathbf{w})$. The linearized matrix $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$ is constructed such that the following properties are satisfied

- (i) As $\mathbf{w}_R \rightarrow \mathbf{w}$ and $\mathbf{w}_L \rightarrow \mathbf{w}$, then $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L) \rightarrow J_{\mathbf{F}}(\mathbf{w})$ (consistency).
- (ii) For any $\mathbf{w}_R, \mathbf{w}_L$, we have $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L) (\mathbf{w}_R - \mathbf{w}_L) = (\mathbf{F}(\mathbf{w}_R) - \mathbf{F}(\mathbf{w}_L)) \mathbf{n}_{ed}$ (conservation across discontinuities).
- (iii) The matrix $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$ is diagonalizable with real eigenvalues.

The Roe scheme then requires the resolution of a Riemann problem at each edge or interelement boundary $\Gamma_{ed} \in \Gamma$. The Roe approximate Riemann solver was first devised for the Euler system [34], and later on applied to the shallow-water model [15, 38]. Classical linearized Roe matrices are known for the most usual systems [28] and $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$ is obtained from $J_{\mathbf{F}}(\mathbf{w})$ by replacing (u, v, ξ) in (16) by (u^*, v^*, ξ^*) defined as

$$\mathbf{u}^* := (u^*, v^*) := \left(\frac{\sqrt{\xi_R} u_R + \sqrt{\xi_L} u_L}{\sqrt{\xi_R} + \sqrt{\xi_L}}, \frac{\sqrt{\xi_R} v_R + \sqrt{\xi_L} v_L}{\sqrt{\xi_R} + \sqrt{\xi_L}} \right), \quad \xi^* := \frac{\xi_R + \xi_L}{2}. \quad (18)$$

By using the properties of the Roe matrix, the Roe flux is then defined as

$$\mathbf{F}(\mathbf{w}^*) \mathbf{n}_{ed} = \frac{1}{2} (\mathbf{F}(\mathbf{w}_R) + \mathbf{F}(\mathbf{w}_L)) \mathbf{n}_{ed} - \frac{1}{2} |J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)| (\mathbf{w}_R - \mathbf{w}_L), \quad (19)$$

where the viscosity matrix $|J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)|$ is written as follows

$$|J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)| = R(\mathbf{w}_R, \mathbf{w}_L) |\Lambda(\mathbf{w}_R, \mathbf{w}_L)| R^{-1}(\mathbf{w}_R, \mathbf{w}_L). \quad (20)$$

In (20), $|\Lambda(\mathbf{w}_R, \mathbf{w}_L)|$ is the diagonal matrix whose coefficients are the absolute value of the eigenvalues of $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$, namely $\lambda_j^*(\mathbf{w}_R, \mathbf{w}_L)$, $j = 1, 2, 3$, obtained from (17) by replacing (u, v, ξ) by (u^*, v^*, ξ^*) , with

$$\lambda_1^* = \mathbf{u}^* \cdot \mathbf{n}_{ed} - \sqrt{g\xi^*}, \quad \lambda_2^* = \mathbf{u}^* \cdot \mathbf{n}_{ed}, \quad \lambda_3^* = \mathbf{u}^* \cdot \mathbf{n}_{ed} + \sqrt{g\xi^*}, \quad (21)$$

and $R(\mathbf{w}_R, \mathbf{w}_L)$ is the matrix whose the j -th column is the eigenvector corresponding to the j -th eigenvalue of $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$, $j = 1, 2, 3$. In order to avoid the computation of the eigenvectors of $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$ in (20) at each side $\Gamma_{ed} \in \Gamma$, a class of methods, named PVM (Polynomial Viscosity Matrix), has been developed in [4] following an idea proposed in [12], leading to a reduced computational cost compared to the evaluation of (20). It is shown in [4] that a few well known schemes such as the Roe, Rusanov, Lax-Friedrichs, FORCE, GFORCE and HLL solvers, can be rewritten on the form of a PVM method.

For $m \in \mathbb{N}$, the PVM method corresponding to the Roe flux in (19) defines the viscosity matrix through a suitable polynomial evaluation of the Roe linearized matrix $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$ on the form

$$|J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)| = \sum_{m=0}^{\max(m)} \alpha_m (J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L))^m, \quad (22)$$

with $\max(m) = 2$. Equation (22) can be interpreted as an approximate Riemann solver based on a polynomial that interpolates the absolute value function at the eigenvalues λ_j^* , $j = 1, 2, 3$, of $J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L)$, yielding

$$\begin{pmatrix} 1 & \lambda_1^* & (\lambda_1^*)^2 \\ 1 & \lambda_2^* & (\lambda_2^*)^2 \\ 1 & \lambda_3^* & (\lambda_3^*)^2 \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{pmatrix} = \begin{pmatrix} |\lambda_1^*| \\ |\lambda_2^*| \\ |\lambda_3^*| \end{pmatrix}. \quad (23)$$

Since the eigenvalues λ_j^* , $j = 1, 2, 3$, are all distinct, Eq. (23) has a unique solution and we obtain

$$\begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{pmatrix} = \frac{1}{2g\xi^*} \begin{pmatrix} |\lambda_1^*|\lambda_2^*\lambda_3^* - 2\lambda_1^*|\lambda_2^*|\lambda_3^* + \lambda_1^*\lambda_2^*|\lambda_3^*| \\ -(\lambda_2^* + \lambda_3^*)|\lambda_1^*| + 2(\lambda_1^* + \lambda_3^*)|\lambda_2^*| - (\lambda_1^* + \lambda_2^*)|\lambda_3^*| \\ |\lambda_1^*| - 2|\lambda_2^*| + |\lambda_3^*| \end{pmatrix}. \quad (24)$$

The Roe flux employed to perform the numerical simulations of Section 5.2 in the non-linear case then yields

$$\mathbf{F}(\mathbf{w}^*) \mathbf{n}_{ed} = \frac{1}{2} (\mathbf{F}(\mathbf{w}_R) + \mathbf{F}(\mathbf{w}_L)) \mathbf{n}_{ed} - \frac{1}{2} \sum_{m=0}^{\max(m)} \alpha_m (J_{\mathbf{F}}(\mathbf{w}_R, \mathbf{w}_L))^m (\mathbf{w}_R - \mathbf{w}_L), \quad (25)$$

with $\max(m) = 2$. The following Rusanov flux is also used in Section 5.2 for comparison purposes with

$$\mathbf{F}(\mathbf{w}^*) \mathbf{n}_{ed} = \frac{1}{2} (\mathbf{F}(\mathbf{w}_R) + \mathbf{F}(\mathbf{w}_L)) \mathbf{n}_{ed} - \frac{1}{2} \max_{j=1,2,3} |\lambda_j^*| I_{3,3} (\mathbf{w}_R - \mathbf{w}_L), \quad (26)$$

where $\max_{j=1,2,3} |\lambda_j^*| = |\mathbf{u}^* \cdot \mathbf{n}_{ed}| + \sqrt{g\xi^*}$ and $I_{3,3}$ is the 3×3 identity matrix.

In order to perform the Fourier analysis in Section 4, the flux $\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed}$ in (14) needs to be derived for the linear SW model. Let $\{\bar{\mathbf{w}}\}$ and $\llbracket \bar{\mathbf{w}} \rrbracket$ denote the mean and the jump of the components of $\bar{\mathbf{w}}$, respectively,

$$\begin{aligned} \{\bar{\mathbf{w}}\} &:= (\{\mathbf{u}\}, \{\eta\}) := (\{u\}, \{v\}, \{\eta\}) := \frac{1}{2} (u_R + u_L, v_R + v_L, \eta_R + \eta_L), \\ \llbracket \bar{\mathbf{w}} \rrbracket &:= (\llbracket \mathbf{u} \rrbracket, \llbracket \eta \rrbracket) := (\llbracket u \rrbracket, \llbracket v \rrbracket, \llbracket \eta \rrbracket) := (u_R - u_L, v_R - v_L, \eta_R - \eta_L). \end{aligned}$$

By using the notation defined in (11) for $\bar{\mathbf{w}}$, the Jacobian matrix $\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L)$ of the normal flux $\bar{\mathbf{F}}(\bar{\mathbf{w}}) \mathbf{n}_{ed}$ is computed from $\bar{\mathbf{F}}(\bar{\mathbf{w}})$ in (8) and the eigenvalues $\bar{\lambda}_j$, $j = 1, 2, 3$, of $\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L)$ are then obtained, leading to

$$\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L) = \begin{pmatrix} 0 & 0 & gn_{ed}^x \\ 0 & 0 & gn_{ed}^y \\ Hn_{ed}^x & Hn_{ed}^y & 0 \end{pmatrix}, \quad \text{and} \quad \begin{pmatrix} \bar{\lambda}_1 \\ \bar{\lambda}_2 \\ \bar{\lambda}_3 \end{pmatrix} = \begin{pmatrix} -\sqrt{gH} \\ 0 \\ \sqrt{gH} \end{pmatrix}. \quad (27)$$

By replacing λ_j^* by $\bar{\lambda}_j$, $j = 1, 2, 3$, in (24) we obtain $\alpha_0 = \alpha_1 = 0$ and $\alpha_2 = 1/\sqrt{gH}$ for the Roe flux, with $\xi^* = H$. In the context of the PVM approach, the Roe flux is then rewritten from (25) in the linear case as

$$\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed} = \frac{1}{2} (\bar{\mathbf{F}}(\bar{\mathbf{w}}_R) + \bar{\mathbf{F}}(\bar{\mathbf{w}}_L)) \mathbf{n}_{ed} - \frac{1}{2\sqrt{gH}} (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^2 \llbracket \bar{\mathbf{w}} \rrbracket, \quad (28)$$

and the Rusanov flux yields

$$\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed} = \frac{1}{2} (\bar{\mathbf{F}}(\bar{\mathbf{w}}_R) + \bar{\mathbf{F}}(\bar{\mathbf{w}}_L)) \mathbf{n}_{ed} - \frac{1}{2} \sqrt{gH} I_{3,3} \llbracket \bar{\mathbf{w}} \rrbracket. \quad (29)$$

Thanks to the particular form of $\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L)$ in (27), the following properties hold

$$(\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^{2m+1} = (gH)^m \bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L), \quad (30)$$

$$(\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^{2m} = (gH)^{m-1} (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^2, \quad (31)$$

for $m = 0, 1, 2, \dots, \max(m)$. By letting

$$\tilde{\alpha}_0 = \alpha_0, \quad \tilde{\alpha}_1 = \sum_{m=1}^{\lfloor \frac{\max(m)+1}{2} \rfloor} \alpha_{2m-1} (gH)^{m-1}, \quad \tilde{\alpha}_2 = \sum_{m=1}^{\lfloor \frac{\max(m)}{2} \rfloor} \alpha_{2m} (gH)^{m-1}, \quad (32)$$

a general PVM method is then rewritten using (25) for the purpose of the Fourier analysis in Section 4 as

$$\sum_{m=0}^{\max(m)} \alpha_m (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^m = \tilde{\alpha}_0 I_{3,3} + \tilde{\alpha}_1 \bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L) + \tilde{\alpha}_2 (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^2. \quad (33)$$

Equation (33) then permits the compact writing of a whole family of PVM fluxes on Γ_{ed} , including the Roe and Rusanov fluxes on the general form

$$\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed} = \frac{1}{2} (\bar{\mathbf{F}}(\bar{\mathbf{w}}_R) + \bar{\mathbf{F}}(\bar{\mathbf{w}}_L)) \mathbf{n}_{ed} - \frac{1}{2} (\tilde{\alpha}_0 I_{3,3} \llbracket \bar{\mathbf{w}} \rrbracket + \tilde{\alpha}_1 \bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L) \llbracket \bar{\mathbf{w}} \rrbracket + \tilde{\alpha}_2 (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^2 \llbracket \bar{\mathbf{w}} \rrbracket). \quad (34)$$

In the remaining part of the paper we let $\tilde{\alpha}_1 = 0$. Indeed, the eigenvalues of $\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L)$ are symmetric, and hence, the polynomials which defines the matrix of numerical dissipation should also be symmetric. Further, we let $\tilde{\alpha}_0 = p\sqrt{gH}$ and $\tilde{\alpha}_2 = q/\sqrt{gH}$, where p and q are positive parameters, and the general PVM flux on Γ_{ed} is finally rewritten from (34) in the following form that is used in Sections 3.3, 4 and 5.1

$$\bar{\mathbf{F}}(\bar{\mathbf{w}}^*) \mathbf{n}_{ed} = \underbrace{\begin{pmatrix} g\{\eta\} \mathbf{n}_{ed} \\ H\{\mathbf{u}\} \cdot \mathbf{n}_{ed} \end{pmatrix}}_{\text{centered contribution}} - \frac{\sqrt{gH}}{2} \underbrace{\begin{pmatrix} p \llbracket \mathbf{u} \rrbracket + q (\llbracket \mathbf{u} \rrbracket \cdot \mathbf{n}_{ed}) \mathbf{n}_{ed} \\ (p+q) \llbracket \eta \rrbracket \end{pmatrix}}_{\text{stabilization part}}. \quad (35)$$

In (35) the PVM flux is naturally split into the sum of a centered part and a stabilizing or diffusive part. The eigenvalues of the stabilization matrix appearing in the right hand side of (34), namely

$$\sum_{m=0}^2 \tilde{\alpha}_m (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^m = p\sqrt{gH} I_{3,3} + \frac{q}{\sqrt{gH}} (\bar{J}_{\mathbf{F}}(\bar{\mathbf{w}}_R, \bar{\mathbf{w}}_L))^2 = \sqrt{gH} \begin{pmatrix} p+q(n_{ed}^x)^2 & qn_{ed}^x n_{ed}^y & 0 \\ qn_{ed}^x n_{ed}^y & p+q(n_{ed}^y)^2 & 0 \\ 0 & 0 & p+q \end{pmatrix}, \quad (36)$$

are $p+q$ (double) and p . Consequently, since p and q are real numbers, the choice $p+q > 0$ and $p > 0$ makes the stabilization matrix definite positive and the jumps $\llbracket \bar{\mathbf{w}} \rrbracket$ are expected to render the scheme stable by acting as a source of numerical dissipation. The choices for (p, q) employed for the linear model are summarized in Table 1. The Rusanov and Roe fluxes coincide with a PVM method with $(p, q) = (1, 0)$ and $(p, q) = (0, 1)$, respectively. The FORCE flux occurs for $(p, q) = (1/2, 1/2)$, namely the PVM – 2 flux, according to the nomenclature defined in [4], while the PVM – 4 flux, introduced in [4] by using a fourth degree polynomial in (22), namely with $\max(m) = 4$, is obtained for $(p, q) = (3/8, 5/8)$.

Table 1: The choices of (p,q) employed in the linear SW model leading to the centered, Rusanov, Roe, PVM-2 and PVM-4 fluxes.

	Centered	Rusanov	Roe	PVM – 2	PVM – 4
p	0	1	0	1/2	3/8
q	0	0	1	1/2	5/8

3.3. The discrete equations for the linear model problem

The discretized weak formulation (14) is now computed using the FV, DG and NC methods. For the purpose of the subsequent Fourier analyses in Section 4, a uniform mesh made up of biased right triangles is considered, denoted as Mesh 1, and the representative meshlength parameter h is thus taken as a constant in the x - and y -directions. Further, only a few set of discrete equations have to be computed for the FV, DG and NC methods at typical nodal locations, due to symmetry reasons and the needs of the Fourier analysis. For the FV, DG and NC schemes, the discretized weak formulation may be expressed on the compact form

$$\frac{\partial}{\partial t} \begin{pmatrix} M(u) \\ M(v) \end{pmatrix} + f \begin{pmatrix} -M(v) \\ M(u) \end{pmatrix} + \frac{g}{h} \begin{pmatrix} G_1(\eta) \\ G_2(\eta) \end{pmatrix} - \frac{\sqrt{gH}}{h} \begin{pmatrix} (p+q)D_1(u) + pD_2(u) + qD_3(v) \\ pD_1(v) + (p+q)D_2(v) + qD_3(u) \end{pmatrix} = 0, \quad (37)$$

$$\frac{\partial}{\partial t} M(\eta) + \frac{H}{h} (G_1(u) + G_2(v)) - \frac{\sqrt{gH}}{h} (p+q) (D_1(\eta) + D_2(\eta)) = 0, \quad (38)$$

where (37) and (38) are computed over triangles K_1 and K_2 on Fig. 2 (middle panel) for the FV method, at nodes 1, 2, 3, 4, 5, 6, on Fig. 3 (middle panel) for the DG scheme, and finally at nodes 1, 2, 3, on Fig. 4 (middle panel) for the NC method. A selected number of scalar discrete equations are hence obtained and written on the form of (37) and (38): 6 for the FV scheme, and 18 and 9, respectively, for the DG and NC schemes. In (37) and (38) G_1 and G_2 denote the gradient operators in the $-x$ and $-y$ directions, respectively, D_1 , D_2 and D_3 are dissipation operators originating from the jumps $[[\mathbf{u}]]$ and $[[\eta]]$ in (35) and M is the mass operator. These operators are now computed in Sections 3.3.1 to 3.3.3 for the FV, DG and NC methods.

3.3.1. The FV discretization

The space of polynomials $\mathcal{P}^0(K_{el})$, with $r = 0$, is first considered. The components (u_h, v_h, η_h) of the approximate solution $\bar{\mathbf{w}}_h$ are constant over each triangle K_{el} of τ_h , and we let $\bar{\mathbf{w}}_h = (u_{K_{el}}, v_{K_{el}}, \eta_{K_{el}})$, while the basis function ψ is nothing else than the characteristic function over K_{el} which is defined to have the value 1 on K_{el} and zero elsewhere, as shown in Fig. 2 (right panel).

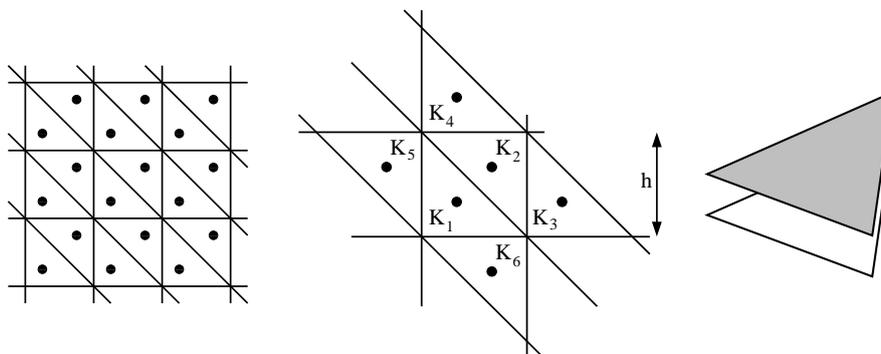


Figure 2: Nodes distribution displayed by the symbol \bullet (left and middle panels) and the constant test function ψ belonging to $\mathcal{P}^0(K_{el})$ (right panel) over a triangle K_{el} of τ_h , for the FV scheme.

As an example, the equation for u in (14) is computed over triangle K_1 in Fig. 2 (middle panel), yielding

$$\frac{\partial u_{K_1}}{\partial t} - f v_{K_1} + \frac{2}{h^2} \int_{\partial K_1} \left(g n_{ed}^x \{\eta_h\} - \frac{\sqrt{gH}}{2} (p \llbracket u_h \rrbracket + q (\llbracket \mathbf{u}_h \rrbracket \cdot \mathbf{n}_{ed}) n_{ed}^x) \right) ds = 0, \quad (39)$$

since the area of K_1 is $h^2/2$ and $\psi_{K_1} = 1$. The computation of the boundary integral along ∂K_1 requires its evaluation along the three boundaries of K_1 , involving the neighbouring triangles K_2, K_5 and K_6 due to the computation of the mean of η_h and the jumps of u_h and v_h . Simple calculations lead to

$$\begin{aligned} \frac{\partial u_{K_1}}{\partial t} - f v_{K_1} + \frac{g}{h} (\eta_{K_2} - \eta_{K_5}) - \frac{\sqrt{gH}}{h} \left(p \left(\sqrt{2}(u_{K_2} - u_{K_1}) - 2u_{K_1} + u_{K_5} + u_{K_6} \right) \right. \\ \left. + q \left(\frac{1}{\sqrt{2}}(-u_{K_1} + u_{K_2} - v_{K_1} + v_{K_2}) - u_{K_1} + u_{K_5} \right) \right) = 0. \end{aligned} \quad (40)$$

In order to perform the Fourier analysis, Eq. (14) is computed similarly for the two remaining equations (v and η) on K_1 and for the three equations (u, v and η) on K_2 . Since χ stands for u, v or η , we obtain

$$G_1(\chi) = \begin{pmatrix} \chi_{K_2} - \chi_{K_5} \\ -\chi_{K_1} + \chi_{K_3} \end{pmatrix}, \quad D_3(\chi) = \frac{1}{\sqrt{2}} \begin{pmatrix} -\chi_{K_1} + \chi_{K_2} \\ \chi_{K_1} - \chi_{K_2} \end{pmatrix}, \quad D_1(\chi) = D_3(\chi) - \begin{pmatrix} \chi_{K_1} - \chi_{K_5} \\ \chi_{K_2} - \chi_{K_3} \end{pmatrix}, \quad (41)$$

$$G_2(\chi) = \begin{pmatrix} \chi_{K_2} - \chi_{K_6} \\ -\chi_{K_1} + \chi_{K_4} \end{pmatrix}, \quad M(\chi) = \begin{pmatrix} \chi_{K_1} \\ \chi_{K_2} \end{pmatrix}, \quad D_2(\chi) = D_3(\chi) - \begin{pmatrix} \chi_{K_1} - \chi_{K_6} \\ \chi_{K_2} - \chi_{K_4} \end{pmatrix}. \quad (42)$$

3.3.2. The P_1^{DG} discretization

The space $V_{h,K_{el}} \subset V_{K_{el}}$, defined as

$$V_{h,K_{el}} := V_{h,K_{el}}^{DG} = \{\chi : \chi|_{K_{el}} \in \mathcal{P}^1(K_{el}) \text{ for all } K_{el} \in \tau_h\},$$

is now considered. The basis function ψ and the approximate solution $\bar{\mathbf{w}}_h$ belong to the space of polynomials $\mathcal{P}^1(K_{el})$, and hence the components u_h, v_h, η_h of $\bar{\mathbf{w}}_h$ are approximated linearly *inside* each triangle K_{el} of τ_h . For example, $\bar{\mathbf{w}}_h = \sum_{j=1}^3 \bar{\mathbf{w}}_j(t) \psi_j(x, y)$ over triangle K_1 , shown in Fig. 3 (middle panel). The nodal values $\bar{\mathbf{w}}_j(t)$ exclusively belong to a given element K_{el} and hence such an approximation is completely discontinuous between the triangles of τ_h . The degrees of freedom $\bar{\mathbf{w}}_j(t)$ are values of the numerical solution at the vertices of K_{el} and the linear basis function $\psi_j(x, y)$ takes the value 1 at a given vertex j and the value 0 at the remaining two vertices of K_{el} , as shown in Fig. 3 (right panel).

The equation for u_h in (14) is computed at node 1 over K_1 in Fig. 3 (middle panel) and we obtain

$$\int_{K_1} \frac{\partial u_h}{\partial t} \psi_1 d\mathbf{x} - f \int_{K_1} v_h \psi_1 d\mathbf{x} - g \int_{K_1} \eta_h \frac{\partial \psi_1}{\partial x} d\mathbf{x} + \int_{\partial K_1} g \eta^* n_{ed}^x \psi_1 ds = 0. \quad (43)$$

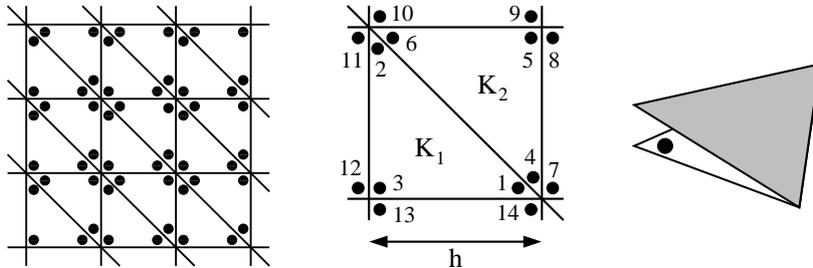


Figure 3: As for Fig. 2 but for the DG scheme with ψ belonging to $V_{h,K_{el}}^{DG}$.

The integrals in (43) are then computed by expanding u_h, v_h and η_h *inside* triangle K_1 in terms of the linear basis functions ψ_1, ψ_2 and ψ_3 . The evaluation of the three first terms in the LHS of (43) yields

$$\int_{K_1} \frac{\partial u_h}{\partial t} \psi_1 d\mathbf{x} = \frac{h^2}{24} \frac{\partial}{\partial t} (2u_1 + u_2 + u_3), \quad -f \int_{K_1} v_h \psi_1 d\mathbf{x} = -f \frac{h^2}{24} (2v_1 + v_2 + v_3), \quad (44)$$

$$-g \int_{K_1} \eta_h \frac{\partial \psi_1}{\partial x} d\mathbf{x} = -g \frac{h}{6} (\eta_1 + \eta_2 + \eta_3). \quad (45)$$

For the computation of the integral along ∂K_1 in (43), let $\Gamma_{i,j}$ and $\mathbf{n}_{i,j}$ be the edge connecting node i to node j and the associated normal, respectively, such that $\partial K_1 = \Gamma_{1,2} \cup \Gamma_{2,3} \cup \Gamma_{3,1}$. The integral along $\Gamma_{2,3}$ cancels as ψ_1 is zero on $\Gamma_{2,3}$ by construction. For the integral along $\Gamma_{1,2}$ we have $\mathbf{n}_{1,2} = \frac{1}{\sqrt{2}}(1, 1)$ and $\psi_1(s) = 1 - s$, with $0 \leq s \leq 1$, leading to

$$\begin{aligned} & \int_{\Gamma_{1,2}} \left(gn_{ed}^x \{\eta_h\} - \frac{\sqrt{gH}}{2} (p \llbracket u_h \rrbracket + q (\llbracket \mathbf{u}_h \rrbracket \cdot \mathbf{n}_{ed}) n_{ed}^x) \right) \psi_1 ds \\ &= \frac{h}{2} \int_0^1 g(\eta_{\Gamma_{4,6}} + \eta_{\Gamma_{1,2}}) (1-s) ds - \frac{h}{2} \sqrt{\frac{gH}{2}} \int_0^1 \left((2p+q)(u_{\Gamma_{4,6}} - u_{\Gamma_{1,2}}) + q(v_{\Gamma_{4,6}} - v_{\Gamma_{1,2}}) \right) (1-s) ds \\ &= \frac{h}{12} g(2\eta_1 + \eta_2 + 2\eta_4 + \eta_6) + \frac{h}{12} \sqrt{\frac{gH}{2}} \left((2p+q)(2u_1 + u_2 - 2u_4 - u_6) + q(2v_1 + v_2 - 2v_4 - v_6) \right), \end{aligned} \quad (46)$$

since $\bar{\mathbf{w}}|_{\Gamma_{1,2}} = \bar{\mathbf{w}}_1 + (\bar{\mathbf{w}}_2 - \bar{\mathbf{w}}_1)s$, and $\bar{\mathbf{w}}|_{\Gamma_{4,6}} = \bar{\mathbf{w}}_4 + (\bar{\mathbf{w}}_6 - \bar{\mathbf{w}}_4)s$.

Similarly, along $\Gamma_{3,1}$ we have $\mathbf{n}_{3,1} = (0, -1)$ and $\psi_1(s) = s$, with $0 \leq s \leq 1$, and

$$\begin{aligned} & \int_{\Gamma_{3,1}} \left(gn_{ed}^x \{\eta_h\} - \frac{\sqrt{gH}}{2} (p \llbracket u_h \rrbracket + q (\llbracket \mathbf{u}_h \rrbracket \cdot \mathbf{n}_{ed}) n_{ed}^x) \right) \psi_1 ds \\ &= -p\sqrt{gH} \frac{h}{2} \int_0^1 (u_{\Gamma_{13,14}} - u_{\Gamma_{3,1}}) s ds = p\sqrt{gH} \frac{h}{12} (2u_1 + u_3 - u_{13} - 2u_{14}), \end{aligned} \quad (47)$$

using $\bar{\mathbf{w}}|_{\Gamma_{3,1}} = \bar{\mathbf{w}}_3 + (\bar{\mathbf{w}}_1 - \bar{\mathbf{w}}_3)s$ and $\bar{\mathbf{w}}|_{\Gamma_{13,14}} = \bar{\mathbf{w}}_{13} + (\bar{\mathbf{w}}_{14} - \bar{\mathbf{w}}_{13})s$. This completes the computation of the equation for u_h in (14) at node 1 and we finally obtain

$$\begin{aligned} & \frac{\partial}{\partial t} (2u_1 + u_2 + u_3) - f(2v_1 + v_2 + v_3) + \frac{2g}{h} (-\eta_2 - 2\eta_3 + 2\eta_4 + \eta_6) \\ &+ \frac{\sqrt{2gH}}{h} \left((2p+q)(2u_1 + u_2 - 2u_4 - u_6) + q(2v_1 + v_2 - 2v_4 - v_6) + p\sqrt{2}(2u_1 + u_3 - u_{13} - 2u_{14}) \right) = 0. \end{aligned} \quad (48)$$

For the purpose of the Fourier analysis, the equation for u_h is computed similarly in (14) at nodes 2, 3, \dots , 6, of Fig. 3, as well as the equations for v_h and η_h in (14) at nodes 1, 2, 3, \dots , 6, leading to a set of 18 discrete equations written on the compact form (37) and (38) with

$$G_1(\chi) = \begin{pmatrix} -\chi_2 - 2\chi_3 + 2\chi_4 + \chi_6 \\ \chi_1 - \chi_3 + \chi_4 + 2\chi_6 - 2\chi_{11} - \chi_{12} \\ 2\chi_1 + \chi_2 - \chi_{11} - 2\chi_{12} \\ -2\chi_1 - \chi_2 + \chi_5 - \chi_6 + 2\chi_7 + \chi_8 \\ -\chi_4 - 2\chi_6 + \chi_7 + 2\chi_8 \\ -\chi_1 - 2\chi_2 + \chi_4 + 2\chi_5 \end{pmatrix}, \quad D_3(\chi) = \frac{1}{\sqrt{2}} \begin{pmatrix} -2\chi_1 - \chi_2 + 2\chi_4 + \chi_6 \\ -\chi_1 - 2\chi_2 + \chi_4 + 2\chi_6 \\ 0 \\ 2\chi_1 + \chi_2 - 2\chi_4 - \chi_6 \\ 0 \\ \chi_1 + 2\chi_2 - \chi_4 - 2\chi_6 \end{pmatrix}, \quad (49)$$

$$G_2(\chi) = \begin{pmatrix} \chi_2 - \chi_3 + 2\chi_4 + \chi_6 - \chi_{13} - 2\chi_{14} \\ -\chi_1 - 2\chi_3 + \chi_4 + 2\chi_6 \\ \chi_1 + 2\chi_2 - 2\chi_{13} - \chi_{14} \\ -2\chi_1 - \chi_2 + 2\chi_5 + \chi_6 \\ -2\chi_4 - \chi_6 + 2\chi_9 + \chi_{10} \\ -\chi_1 - 2\chi_2 - \chi_4 + \chi_5 + \chi_9 + 2\chi_{10} \end{pmatrix}, \quad D_1(\chi) = D_3(\chi) - \begin{pmatrix} 0 \\ 2\chi_2 + \chi_3 - 2\chi_{11} - \chi_{12} \\ \chi_2 + 2\chi_3 - \chi_{11} - 2\chi_{12} \\ 2\chi_4 + \chi_5 - 2\chi_7 - \chi_8 \\ \chi_4 + 2\chi_5 - \chi_7 - 2\chi_8 \\ 0 \end{pmatrix}, \quad (50)$$

$$M(\chi) = \frac{1}{2} \begin{pmatrix} 2\chi_1 + \chi_2 + \chi_3 \\ \chi_1 + 2\chi_2 + \chi_3 \\ \chi_1 + \chi_2 + 2\chi_3 \\ 2\chi_4 + \chi_5 + \chi_6 \\ \chi_4 + 2\chi_5 + \chi_6 \\ \chi_4 + \chi_5 + 2\chi_6 \end{pmatrix}, \quad D_2(\chi) = D_3(\chi) - \begin{pmatrix} 2\chi_1 + \chi_3 - \chi_{13} - 2\chi_{14} \\ 0 \\ \chi_1 + 2\chi_3 - 2\chi_{13} - \chi_{14} \\ 0 \\ 2\chi_5 + \chi_6 - 2\chi_9 - \chi_{10} \\ \chi_5 + 2\chi_6 - \chi_9 - 2\chi_{10} \end{pmatrix}. \quad (51)$$

3.3.3. The P_1^{NC} discretization

The polynomial space $V_{h,K_{el}}$ is finally considered with

$$V_{h,K_{el}} := V_{h,K_{el}}^{NC} = \{\chi : \chi|_{K_{el}} \in \mathcal{P}^1(K_{el}) \text{ for all } K_{el} \in \tau_h, \chi \text{ is continuous at the midpoints of } \Gamma_{ed}\},$$

As for the P_1^{DG} case, the basis function ψ and the discrete solution $\bar{\mathbf{w}}_h$ belong to the space of polynomials $\mathcal{P}^1(K_{el})$, and hence the components of $\bar{\mathbf{w}}_h$ are approximated linearly over triangles K_{el} of τ_h . However, the nodal values $\bar{\mathbf{w}}_j$ of the P_1^{NC} element are located at midpoints of edges of K_{el} [9], as shown in Fig. 4 (left and middle panels), and two adjacent triangles share the *same* nodal value at the midpoint of their common edge. The shape of the P_1^{NC} linear basis function ψ that is used to approximate the components of $\bar{\mathbf{w}}_h$ on the element's two-triangle support is represented in grey in Fig. 4 (right panel). The basis function has the value 1 at the midpoint and along the common edge of the two adjacent triangles and 0 at the midpoints of the remaining four edges, represented by the symbol \bullet in Fig. 4 (right panel). Since such a representation is only continuous across triangle boundaries at midpoint of the edges, and discontinuous everywhere else around a triangle boundary, this element is termed nonconforming (*NC*) in the FE literature. A useful property of the P_1^{NC} element is that the mass matrix is diagonal since the P_1^{NC} basis functions are orthogonal [39].

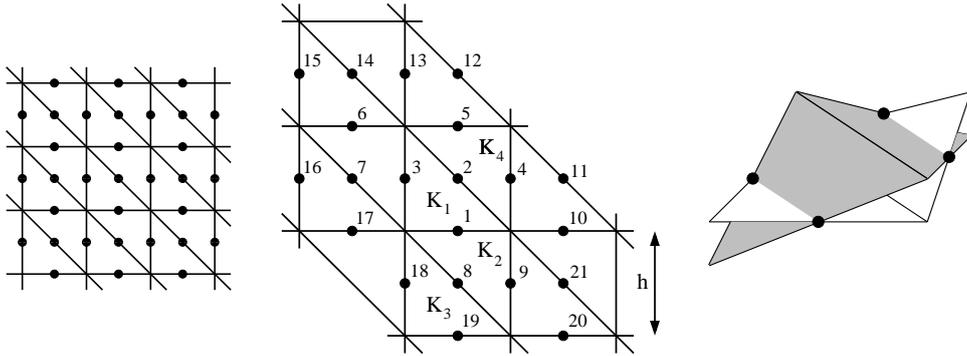


Figure 4: As for Fig. 2 but for the NC scheme with ψ belonging to $V_{h,K_{el}}^{NC}$.

The equation for u_h in (14) is computed at node 1 of Fig. 4 (middle panel), and the compact support of the NC basis function ψ_1 at node 1 is now made up of triangles K_1 and K_2 , namely $K_1 \cup K_2$, leading to

$$\int_{K_1 \cup K_2} \frac{\partial u_h}{\partial t} \psi_1 \, d\mathbf{x} - f \int_{K_1 \cup K_2} v_h \psi_1 \, d\mathbf{x} - g \int_{K_1 \cup K_2} \eta_h \frac{\partial \psi_1}{\partial x} \, d\mathbf{x} + \int_{\partial K_1 \cup \partial K_2} g\eta^* n_{ed}^x \psi_1 \, ds = 0. \quad (52)$$

The integrals in (52) are computed by expanding u_h, v_h and η_h over triangles K_1 and K_2 in terms of the NC basis functions $\psi_1, \psi_2, \psi_3, \psi_8$ and ψ_9 . Due to the orthogonality property of the NC basis functions and because $\frac{\partial \psi_1}{\partial x}|_{K_1} = \frac{\partial \psi_1}{\partial x}|_{K_2} = 0$, we obtain

$$\int_{K_1 \cup K_2} \frac{\partial u_h}{\partial t} \psi_1 d\mathbf{x} = \frac{h^2}{3} \frac{\partial u_1}{\partial t}, \quad -f \int_{K_1 \cup K_2} v_h \psi_1 d\mathbf{x} = -f \frac{h^2}{3} v_1, \quad -g \int_{K_1 \cup K_2} \eta_h \frac{\partial \psi_1}{\partial x} d\mathbf{x} = 0. \quad (53)$$

Let Γ_j and \mathbf{n}_j be the edge corresponding to midpoint j and the associated normal, respectively, and let $\bar{\mathbf{w}}_{\Gamma_i}|_{K_j}$ denote the evaluation of $\bar{\mathbf{w}}$ along Γ_i on K_j . We have $\partial K_1 = \Gamma_1 \cup \Gamma_2 \cup \Gamma_3$ and $\partial K_2 = \Gamma_1 \cup \Gamma_8 \cup \Gamma_9$. The contribution of K_1 and K_2 in the last integral in the LHS of (52) cancels along Γ_1 due to the *unique* normal vector \mathbf{n}_1 on Γ_1 . Along Γ_8 we have $\mathbf{n}_8 = \frac{-1}{\sqrt{2}}(1, 1)$ and $\psi_1(s) = 1 - 2s$, with $0 \leq s \leq 1$, yielding

$$\begin{aligned} & \int_{\Gamma_8} \left(g n_{ed}^x \{\eta_h\} - \frac{\sqrt{gH}}{2} (p \llbracket u_h \rrbracket + q (\llbracket \mathbf{u}_h \rrbracket \cdot \mathbf{n}_{ed}) n_{ed}^x) \right) \psi_1 ds \\ &= \frac{h}{2} \int_0^1 g (\eta_{\Gamma_8}|_{K_3} + \eta_{\Gamma_8}|_{K_2}) (2s-1) ds - \frac{h}{2} \sqrt{\frac{gH}{2}} \int_0^1 \left((2p+q)(u_{\Gamma_8}|_{K_3} - u_{\Gamma_8}|_{K_2}) + q(v_{\Gamma_8}|_{K_3} - v_{\Gamma_8}|_{K_2}) \right) (1-2s) ds \\ &= \frac{h}{6} g (-\eta_1 + \eta_9 - \eta_{18} + \eta_{19}) + \frac{h}{6} \sqrt{\frac{gH}{2}} \left((2p+q)(u_1 - u_9 - u_{18} + u_{19}) + q(v_1 - v_9 - v_{18} + v_{19}) \right), \end{aligned} \quad (54)$$

since $\bar{\mathbf{w}}|_{\Gamma_8}$ on $K_2 = \bar{\mathbf{w}}_8 - \bar{\mathbf{w}}_9 + \bar{\mathbf{w}}_1 + 2s(\bar{\mathbf{w}}_9 - \bar{\mathbf{w}}_1)$ and $\bar{\mathbf{w}}|_{\Gamma_8}$ on $K_3 = \bar{\mathbf{w}}_8 - \bar{\mathbf{w}}_{19} + \bar{\mathbf{w}}_{18} + 2s(\bar{\mathbf{w}}_{19} - \bar{\mathbf{w}}_{18})$. The computations of the last integral in the LHS of (52) are conducted similarly on Γ_9, Γ_2 and Γ_3 , and (52) finally reads as

$$\begin{aligned} & \frac{\partial u_1}{\partial t} - f v_1 + \frac{g}{2h} (\eta_2 - \eta_3 + \eta_4 - \eta_5 + \eta_6 - \eta_7 - \eta_8 + \eta_9 - \eta_{18} + \eta_{19} - \eta_{20} + \eta_{21}) \\ & + \frac{\sqrt{gH}}{2\sqrt{2}h} \left((2p+q)(2u_1 - u_3 - u_4 + u_5 - u_9 - u_{18} + u_{19}) + q(2v_1 - v_3 - v_4 + v_5 \right. \\ & \quad \left. - v_9 - v_{18} + v_{19}) + \sqrt{2}(p+q)(2u_1 - u_2 + u_6 - u_7 - u_8 + u_{20} - u_{21}) \right) = 0. \end{aligned} \quad (55)$$

In order to perform the Fourier analysis, the equation for u_h in (14) is computed similarly at nodes 2 and 3 of Fig. 4, as well as the equations for v_h and η_h in (14) at nodes 1, 2 and 3, leading to a set of 9 discrete equations written on the compact form (37) and (38) with

$$G_1(\chi) = \begin{pmatrix} \chi_2 - \chi_3 + \chi_4 - \chi_5 + \chi_6 - \chi_7 - \chi_8 + \chi_9 - \chi_{18} + \chi_{19} - \chi_{20} + \chi_{21} \\ -\chi_1 - 2\chi_3 + 2\chi_4 + \chi_5 - \chi_6 + \chi_7 + \chi_{10} - \chi_{11} \\ \chi_1 + 2\chi_2 - \chi_4 + \chi_5 - \chi_6 - 2\chi_7 + \chi_{16} - \chi_{17} \end{pmatrix}, \quad (56)$$

$$G_2(\chi) = \begin{pmatrix} 2\chi_2 + \chi_3 + \chi_4 - \chi_5 - 2\chi_8 - \chi_9 - \chi_{18} + \chi_{19} \\ -2\chi_1 - \chi_3 + \chi_4 + 2\chi_5 + \chi_8 - \chi_9 - \chi_{12} + \chi_{13} \\ -\chi_1 + \chi_2 - \chi_4 + \chi_5 + \chi_6 - \chi_7 - \chi_8 + \chi_9 + \chi_{14} - \chi_{15} + \chi_{16} - \chi_{17} \end{pmatrix}, \quad (57)$$

$$D_3(\chi) = \frac{1}{\sqrt{2}} \begin{pmatrix} -2\chi_1 + \chi_3 + \chi_4 - \chi_5 + \chi_9 + \chi_{18} - \chi_{19} \\ 0 \\ \chi_1 - 2\chi_3 - \chi_4 + \chi_5 + \chi_6 - \chi_{16} + \chi_{17} \end{pmatrix}, \quad (58)$$

$$D_1(\chi) = D_3(\chi) - \begin{pmatrix} 2\chi_1 - \chi_2 + \chi_6 - \chi_7 - \chi_8 + \chi_{20} - \chi_{21} \\ -\chi_1 + 2\chi_2 - \chi_5 - \chi_6 + \chi_7 - \chi_{10} + \chi_{11} \\ 0 \end{pmatrix}, \quad (59)$$

$$D_2(\chi) = D_3(\chi) - \begin{pmatrix} 0 \\ 2\chi_2 - \chi_3 - \chi_4 + \chi_8 - \chi_9 + \chi_{12} - \chi_{13} \\ -\chi_2 + 2\chi_3 - \chi_7 - \chi_8 + \chi_9 - \chi_{14} + \chi_{15} \end{pmatrix}, \quad M(\chi) = 2 \begin{pmatrix} \chi_1 \\ \chi_2 \\ \chi_3 \end{pmatrix}. \quad (60)$$

4. Fourier/Dispersion analyses of the linear model problem

4.1. The generalized eigenvalue system

The Fourier / dispersion analysis has proven practical and beneficial to detect dispersion and dissipation problems as well as the presence of spurious modes in a few discretization schemes [6, 24, 25, 26]. It is adopted in this study to analyse the stability and accuracy of the FV, DG and NC discrete schemes examined earlier in Section 3. Note that the meshlength parameter h remains constant in the present section.

When continuous linear finite elements are employed, the Fourier analysis is straightforward since all the discrete unknowns are located at only one type of mesh nodes, namely vertices. It is hence sufficient to consider one discrete equation for u_h , v_h and η_h at such a typical node to compute the dispersion relation, due to symmetry reasons. As for the continuous case, a polynomial of degree three for the frequency is then obtained, but spurious pressure oscillations prevent the uniqueness of the solution, as mentioned in Section 1.

In this study, such a simple case does not occur since the nodal unknowns appearing in (37) and (38) are located at the barycenter of triangles for the FV method, local vertices per triangle for the DG scheme and midpoints of edges for the NC method. Further, due to symmetry reasons, two possible types of triangles need to be considered for the FV and DG schemes: the lower left and upper right triangles, namely K_1 and K_2 in Figs. 2 and 3. Consequently, we need to compute the discrete equations for u_h , v_h and η_h at the:

- two types of barycenters corresponding to upward and downward pointing triangles for the FV scheme,
- six types of vertices for the DG scheme, namely three vertices per triangle for the upward and downward pointing triangles,
- three types of faces: namely, horizontal, vertical and diagonal faces for the NC scheme.

This explains why we have computed (37) and (38) over triangles K_1 and K_2 on Fig. 2 for the FV method, at nodes 1, 2, 3, 4, 5, 6, on Fig. 3 for the DG scheme, and finally at nodes 1, 2, 3, on Fig. 4 for the NC method. Note that the nodal unknowns belonging to the same set are distributed on a regular grid of size h , as shown in Fig. 5 for the FV, DG and NC schemes. The numbering of the typical nodal unknowns in Fig. 5 corresponds to that employed in Figs. 2 to 4 for the three schemes.

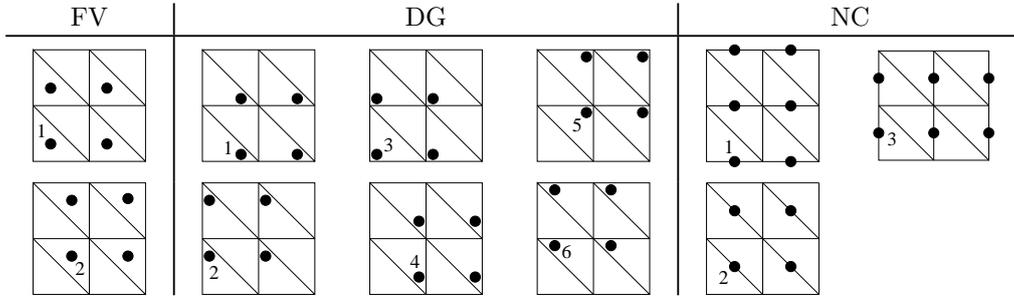


Figure 5: The sets of nodal unknowns for the FV, DG and NC schemes. The nodal unknowns belonging to the same set are distributed on a regular grid of size h .

In order to obtain the dispersion relations at the discrete level and to find the discrete geostrophic and inertia-gravity frequencies, we follow the same procedure as in Section 2 for the continuous case. The discrete solutions corresponding to $(u_j, v_j, \eta_j) = (\hat{u}_{j_0}, \hat{v}_{j_0}, \hat{\eta}_{j_0}) e^{i(kx_j + ly_j - \omega t)}$ are sought, where (u_j, v_j, η_j) are the nodal unknowns that appear in the selected discrete equations (37) and (38) and $(\hat{u}_{j_0}, \hat{v}_{j_0}, \hat{\eta}_{j_0})$ are the Fourier amplitudes (for the FV scheme j is written as K_j). The (x_j, y_j) coordinates are expressed in terms of a distance to a reference node. By letting $\mathcal{I}_n = \{1, 2, \dots, n\}$, where n is a positive integer, j_0 belongs to the set of nodal unknowns for the discrete schemes shown in Fig. 5, with $j_0 \in \mathcal{I}_2$ for the FV method, $j_0 \in \mathcal{I}_6$

for the DG scheme and $j_0 \in \mathcal{I}_3$ for the NC discretization. Indeed, j_0 is intimately associated with each type of node for which (37) and (38) have been computed, as illustrated in Sections 4.2 to 4.4. Let

$$\tilde{\omega} = \frac{h\omega}{\sqrt{gH}}, \quad R_d = \frac{\sqrt{gH}}{f}, \quad \lambda = \frac{R_d}{h}, \quad (61)$$

be the normalized frequency, the Rossby radius of deformation and the λ parameter, a dimensionless quantity measuring resolution, with e.g. $\lambda = 2$ and $\lambda = 1/10$ corresponding to fine and coarse meshes, respectively.

Substitution of (u_j, v_j, η_j) in the selected discrete equations (37) and (38) leads to a matrix system for the amplitudes on the form of the following generalized eigenvalue problem

$$\hat{A}_{3n,3n} X = i\tilde{\omega} \hat{B}_{3n,3n} X, \quad (62)$$

with $X = \left(\{\hat{u}_j\}_{j \in \mathcal{I}_n}, \{\hat{v}_j\}_{j \in \mathcal{I}_n}, \{\hat{\eta}_j\}_{j \in \mathcal{I}_n} \right)$ and

$$\hat{A}_{3n,3n} = \begin{pmatrix} (p+q)\hat{D}_1 + p\hat{D}_2 & -\frac{1}{\lambda}\hat{M} + q\hat{D}_3 & \sqrt{\frac{g}{H}}\hat{G}_1 \\ \frac{1}{\lambda}\hat{M} + q\hat{D}_3 & p\hat{D}_1 + (p+q)\hat{D}_2 & \sqrt{\frac{g}{H}}\hat{G}_2 \\ \sqrt{\frac{H}{g}}\hat{G}_1 & \sqrt{\frac{H}{g}}\hat{G}_2 & (p+q)(\hat{D}_1 + \hat{D}_2) \end{pmatrix}, \quad \hat{B}_{3n,3n} = \begin{pmatrix} \hat{M} & O & O \\ O & \hat{M} & O \\ O & O & \hat{M} \end{pmatrix}, \quad (63)$$

where O is the $n \times n$ null matrix. The matrices $\hat{M}, \hat{G}_1, \hat{G}_2, \hat{D}_1, \hat{D}_2$ and \hat{D}_3 of size $n \times n$ in (63) are defined in Sections 4.2 to 4.4 for the FV ($n = 2$), DG ($n = 6$) and NC ($n = 3$) schemes, respectively.

4.2. The FV scheme

Inserting $(u_{K_j}, v_{K_j}, \eta_{K_j}) = (\hat{u}_{j_0}, \hat{v}_{j_0}, \hat{\eta}_{j_0}) e^{i(kx_j + ly_j - \omega t)}$ in the three first terms in the LHS of Eq. (40), with $j_0 \in \mathcal{I}_2$, yields

$$\begin{aligned} & \frac{\partial u_{K_1}}{\partial t} - f v_{K_1} + \frac{g}{h} (\eta_{K_2} - \eta_{K_5}) \\ &= -i\omega \hat{u}_1 e^{i(kx_{K_1} + ly_{K_1})} - f \hat{v}_1 e^{i(kx_{K_1} + ly_{K_1})} + \frac{g}{h} (\hat{\eta}_2 e^{i(kx_{K_2} + ly_{K_2})} - \hat{\eta}_5 e^{i(kx_{K_5} + ly_{K_5})}) \\ &= e^{i(kx_{K_1} + ly_{K_1})} \left(-i\omega \hat{u}_1 - f \hat{v}_1 + \frac{g}{h} \hat{\eta}_2 (e^{i(k+l)\frac{h}{3}} - e^{i(-2k+l)\frac{h}{3}}) \right), \end{aligned} \quad (64)$$

since node K_5 in Fig. 2 (middle) belongs to the set of amplitudes of type 2 in Fig. 5. By treating the stabilizing term in a similar manner, Eq. (40) finally leads to

$$\begin{aligned} & -i\omega \hat{u}_1 - f \hat{v}_1 + \frac{g}{h} \hat{\eta}_2 e^{il\frac{h}{3}} \left(e^{ik\frac{h}{3}} - e^{-2ik\frac{h}{3}} \right) + \frac{\sqrt{gH}}{h} \left(\frac{1}{\sqrt{2}} (2p+q) (\hat{u}_1 - \hat{u}_2 e^{i(k+l)\frac{h}{3}}) \right. \\ & \left. + \frac{1}{\sqrt{2}} q (\hat{v}_1 - \hat{v}_2 e^{i(k+l)\frac{h}{3}}) + (p+q) (\hat{u}_1 - \hat{u}_2 e^{i(l-2k)\frac{h}{3}}) + p (\hat{u}_1 - \hat{u}_2 e^{i(k-2l)\frac{h}{3}}) \right) = 0. \end{aligned} \quad (65)$$

By using the normalized frequency $\tilde{\omega}$ and the parameter λ , letting $a_1 = e^{\frac{ikh}{3}}$ and $a_2 = e^{\frac{i lh}{3}}$, and for $X = (\hat{u}_1, \hat{u}_2, \hat{v}_1, \hat{v}_2, \hat{\eta}_1, \hat{\eta}_2)$ in (62), we obtain from (41) and (42)

$$\hat{G}_1 = (a_1^3 - 1) \begin{pmatrix} 0 & a_2 \\ 1 & a_1^2 \\ a_1 a_2 & 0 \end{pmatrix}, \quad \hat{D}_1 = \begin{pmatrix} 1 + \frac{1}{\sqrt{2}} & \frac{-a_1 a_2}{\sqrt{2}} - \frac{a_2}{a_1^2} \\ -1 & 1 + \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2} a_1 a_2} - \frac{a_2}{a_1} & 1 + \frac{1}{\sqrt{2}} \end{pmatrix}, \quad \hat{D}_3 = \frac{-1}{\sqrt{2}} \begin{pmatrix} -1 & a_1 a_2 \\ 1 & -1 \\ a_1 a_2 & -1 \end{pmatrix}, \quad (66)$$

$$\hat{G}_2 = (a_2^3 - 1) \begin{pmatrix} 0 & a_1 \\ 1 & a_2^2 \\ a_1 a_2 & 0 \end{pmatrix}, \quad \hat{D}_2 = \begin{pmatrix} 1 + \frac{1}{\sqrt{2}} & \frac{-a_1 a_2}{\sqrt{2}} - \frac{a_1}{a_2^2} \\ -1 & 1 + \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2} a_1 a_2} - \frac{a_2}{a_1} & 1 + \frac{1}{\sqrt{2}} \end{pmatrix}, \quad \hat{M} = I_{2,2}. \quad (67)$$

4.3. The P_1^{DG} scheme

The Fourier expansion $(u_j, v_j, \eta_j) = (\hat{u}_{j_0}, \hat{v}_{j_0}, \hat{\eta}_{j_0}) e^{i(kx_j + ly_j - \omega t)}$ is now inserted in (48). By following the same procedure than in Section 4.2, but employing the set of six possible amplitudes of Fig. 5, with $j_0 \in \mathcal{I}_6$, we obtain at node 1

$$\begin{aligned} & -i\omega(2\hat{u}_1 + \hat{u}_2 e^{i(-k+l)h} + \hat{u}_3 e^{-ikh}) - f(2\hat{v}_1 + \hat{v}_2 e^{i(-k+l)h} + \hat{v}_3 e^{-ikh}) \\ & + \frac{2g}{h} (-\hat{\eta}_2 e^{i(-k+l)h} - 2\hat{\eta}_3 e^{-ikh} + 2\hat{\eta}_4 + \hat{\eta}_6 e^{i(-k+l)h}) + \frac{\sqrt{2gH}}{h} \left((2p+q)(2\hat{u}_1 + \hat{u}_2 e^{i(-k+l)h} - 2\hat{u}_4 - \hat{u}_6 e^{i(-k+l)h}) \right. \\ & \left. + q(2\hat{v}_1 + \hat{v}_2 e^{i(-k+l)h} - 2\hat{v}_4 - \hat{v}_6 e^{i(-k+l)h}) + p\sqrt{2}(2\hat{u}_1 + \hat{u}_3 e^{-ikh} - 2\hat{u}_5 - \hat{u}_6 e^{-ikh}) \right) = 0. \end{aligned} \quad (68)$$

For $X = (\hat{u}_1, \hat{u}_2, \hat{u}_3, \hat{u}_6, \hat{u}_4, \hat{u}_5, \hat{v}_1, \hat{v}_2, \hat{v}_3, \hat{v}_6, \hat{v}_4, \hat{v}_5, \hat{\eta}_1, \hat{\eta}_2, \hat{\eta}_3, \hat{\eta}_6, \hat{\eta}_4, \hat{\eta}_5)$ in (62), we let

$$b_1 = e^{ikh}, \quad b_2 = e^{ilh}, \quad b_3 = b_1 - 1, \quad b_4 = b_2 - 1, \quad b_5 = \frac{b_1}{\sqrt{2}b_2} + \frac{1}{b_2}, \quad b_6 = \frac{b_2}{\sqrt{2}b_1} + \frac{1}{b_1},$$

and after long et tedious algebra, Eqs. (49) to (51) yields

$$\hat{D}_j = \begin{pmatrix} \hat{D}_{j1}(b_1, b_2) & \hat{D}_{j2}(b_1, b_2) \\ \hat{D}_{j2}(\frac{1}{b_1}, \frac{1}{b_2}) & \hat{D}_{j1}(\frac{1}{b_1}, \frac{1}{b_2}) \end{pmatrix}, \quad \hat{G}_j = \begin{pmatrix} \hat{G}_{j1}(b_1, b_2) & \hat{G}_{j2}(b_1, b_2) \\ -\hat{G}_{j2}(\frac{1}{b_1}, \frac{1}{b_2}) & -\hat{G}_{j1}(\frac{1}{b_1}, \frac{1}{b_2}) \end{pmatrix}, \quad \hat{M} = \begin{pmatrix} \hat{M}_1(b_1, b_2) & O_{3,3} \\ O_{3,3} & \hat{M}_1(\frac{1}{b_1}, \frac{1}{b_2}) \end{pmatrix},$$

with $j = 1, 2, 3$, for \hat{D}_j , and $j = 1, 2$, for \hat{G}_j , where the submatrices read as

$$\hat{D}_{11}(b_1, b_2) = \begin{pmatrix} \sqrt{2} & \frac{b_2}{\sqrt{2}b_1} & 0 \\ \frac{b_1}{\sqrt{2}b_2} & 2 + \sqrt{2} & \frac{1}{b_2} \\ 0 & b_2 & 2 \end{pmatrix}, \quad \hat{D}_{21}(b_1, b_2) = \begin{pmatrix} 2 + \sqrt{2} & \frac{b_2}{\sqrt{2}b_1} & \frac{1}{b_1} \\ \frac{b_1}{\sqrt{2}b_2} & \sqrt{2} & 0 \\ b_1 & 0 & 2 \end{pmatrix}, \quad \hat{D}_{31}(b_1, b_2) = \frac{1}{\sqrt{2}} \begin{pmatrix} 2 & \frac{b_2}{b_1} & 0 \\ \frac{b_1}{b_2} & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix},$$

$$\hat{D}_{12}(b_1, b_2) = - \begin{pmatrix} \frac{b_2}{\sqrt{2}b_1} & \sqrt{2} & 0 \\ \sqrt{2} & b_5 & 2 \\ 0 & 2 & b_2 \end{pmatrix}, \quad \hat{D}_{22}(b_1, b_2) = - \begin{pmatrix} b_6 & \sqrt{2} & 2 \\ \sqrt{2} & \frac{b_1}{\sqrt{2}b_2} & 0 \\ 2 & 0 & b_1 \end{pmatrix}, \quad \hat{D}_{32}(b_1, b_2) = -\frac{1}{\sqrt{2}} \begin{pmatrix} \frac{b_2}{b_1} & 2 & 0 \\ 2 & \frac{b_1}{b_2} & 0 \\ 0 & 0 & 0 \end{pmatrix},$$

$$\hat{G}_{11}(b_1, b_2) = \begin{pmatrix} 0 & \frac{-b_2}{b_1} & \frac{-2}{b_1} \\ \frac{b_1}{b_2} & 0 & \frac{-1}{b_2} \\ 2b_1 & b_2 & 0 \end{pmatrix}, \quad \hat{G}_{12}(b_1, b_2) = \begin{pmatrix} \frac{b_2}{b_1} & 2 & 0 \\ 2 & \frac{b_3}{b_2} & -2 \\ 0 & -2 & -b_2 \end{pmatrix},$$

$$\hat{G}_{21} = \begin{pmatrix} 0 & \frac{b_2}{b_1} & \frac{-1}{b_1} \\ \frac{-b_1}{b_2} & 0 & \frac{-2}{b_2} \\ b_1 & 2b_2 & 0 \end{pmatrix}, \quad \hat{G}_{22} = \begin{pmatrix} \frac{b_4}{b_1} & 2 & -2 \\ 2 & \frac{b_1}{b_2} & 0 \\ -2 & 0 & -b_1 \end{pmatrix}. \quad \hat{M}_1(b_1, b_2) = \frac{1}{2} \begin{pmatrix} 2 & \frac{b_2}{b_1} & \frac{1}{b_1} \\ \frac{b_1}{b_2} & 2 & \frac{1}{b_2} \\ b_1 & b_2 & 2 \end{pmatrix}.$$

4.4. The P_1^{NC} scheme

The expansion $(u_j, v_j, \eta_j) = (\hat{u}_{j_0}, \hat{v}_{j_0}, \hat{\eta}_{j_0}) e^{i(kx_j + ly_j - \omega t)}$ is finally inserted in (55), and using the set of three possible amplitudes of Fig. 5, with $j_0 \in \mathcal{I}_3$, leads to

$$\begin{aligned} & -i\omega \hat{u}_1 - f\hat{v}_1 + 2i \sin \frac{kh}{2} \frac{g}{h} \left(-\hat{\eta}_1 \cos \frac{(k-2l)h}{2} + \hat{\eta}_2 \cos \frac{(k-l)h}{2} + \hat{\eta}_3 \cos \frac{lh}{2} \right) \\ & + \frac{\sqrt{gH}}{h} \left(\sqrt{2}(2p+q) (\hat{u}_1 \cos^2 \frac{lh}{2} - \hat{u}_3 \cos \frac{kh}{2} \cos \frac{lh}{2}) + \sqrt{2}q (\hat{v}_1 \cos^2 \frac{lh}{2} - \hat{v}_3 \cos \frac{kh}{2} \cos \frac{lh}{2}) \right. \\ & \left. + 2(p+q) (\hat{u}_1 \cos^2 \frac{(k-l)h}{2} - \hat{u}_2 \cos \frac{kh}{2} \cos \frac{(k-l)h}{2}) \right) = 0. \end{aligned} \quad (69)$$

By letting

$$\begin{aligned} c_1 &= \cos \frac{kh}{2}, & c_3 &= \sin \frac{kh}{2}, & c_5 &= \cos \frac{(k-l)h}{2}, & c_7 &= \cos \frac{(k-2l)h}{2}, \\ c_2 &= \cos \frac{lh}{2}, & c_4 &= \sin \frac{lh}{2}, & c_6 &= \cos \frac{(2k-l)h}{2}. \end{aligned}$$

and for $X = (\hat{u}_1, \hat{u}_2, \hat{u}_3, \hat{v}_1, \hat{v}_2, \hat{v}_3, \hat{\eta}_1, \hat{\eta}_2, \hat{\eta}_3)$ in (62), we obtain from Eqs. (56) to (60)

$$\begin{aligned} \hat{D}_1 &= \sqrt{2} \begin{pmatrix} c_2^2 + \sqrt{2}c_5^2 & -\sqrt{2}c_1c_5 & -c_1c_2 \\ -\sqrt{2}c_1c_5 & \sqrt{2}c_1^2 & 0 \\ -c_1c_2 & 0 & c_1^2 \end{pmatrix}, & \hat{D}_2 &= \sqrt{2} \begin{pmatrix} c_2^2 & 0 & -c_1c_2 \\ 0 & \sqrt{2}c_2^2 & -\sqrt{2}c_2c_5 \\ -c_1c_2 & -\sqrt{2}c_2c_5 & c_1^2 + \sqrt{2}c_5^2 \end{pmatrix}, & \hat{M} &= I_{3,3} \\ \hat{D}_3 &= \sqrt{2} \begin{pmatrix} c_2^2 & 0 & -c_1c_2 \\ 0 & 0 & 0 \\ -c_1c_2 & 0 & c_1^2 \end{pmatrix}, & \hat{G}_1 &= 2ic_3 \begin{pmatrix} -c_7 & c_5 & c_2 \\ c_5 & -c_1 & 1 \\ c_2 & 1 & -c_1 \end{pmatrix}, & \hat{G}_2 &= 2ic_4 \begin{pmatrix} -c_2 & 1 & c_1 \\ 1 & -c_2 & c_5 \\ c_1 & c_5 & -c_6 \end{pmatrix}. \end{aligned}$$

4.5. The dispersion relations

The dispersion relation or characteristic polynomial $P_{3n}(\tilde{\omega})$ is computed from (62) by setting

$$P_{3n}(\tilde{\omega}) := \det \left(\hat{A}_{3n,3n} - i\tilde{\omega} \hat{B}_{3n,3n} \right) = 0. \quad (70)$$

A polynomial of degree $3n$ in $\tilde{\omega}$ is then obtained, namely a polynomial of degree 6, 18 and 9, for the FV ($n = 2$), DG ($n = 6$) and NC ($n = 3$) schemes, respectively, whose coefficients depend on the parameters p, q, k, l, h, g and H . On the other hand, the continuous frequencies are rewritten from (10) by using the definition of $\tilde{\omega}$ and the parameter λ in (61), and this leads to

$$\tilde{\omega}_{1,2}^C = \pm \sqrt{\lambda^{-2} + (kh)^2 + (lh)^2} \quad \text{and} \quad \tilde{\omega}_3^C = 0. \quad (71)$$

The purpose of the subsequent analyses is to compare the $3n$ discrete frequencies $\tilde{\omega}$ solution of (70) with their 3 continuous counterparts in (71).

Remark 1. *By letting $p = q = 0$ and multiplying the two first lines of (62) by $\sqrt{H/g}$ and the last line of (62) by $\sqrt{g/H}$, the matrix in the LHS of (62) is found to be a skew-Hermitian matrix. Indeed, the mass matrix \hat{M} is symmetric and the gradient matrices \hat{G}_1 and \hat{G}_2 are skew-Hermitian. Further, \hat{M} is positive definite and hence \hat{B} is also a symmetric and positive definite matrix. Consequently, the eigenvalues $\tilde{\omega}$ of (62) are real. However, the frequencies $\tilde{\omega}$ are complex numbers as soon as $p \neq 0$ or $q \neq 0$.*

4.5.1. The centered case with $(p, q) = (0, 0)$

In the case $(p, q) = (0, 0)$, the solutions of (70) yield the following $3n$ real roots in the absence of numerical stabilization: n roots equal to zero corresponding to $\tilde{\omega}_3^C$, and n positive and n negative roots corresponding to $\tilde{\omega}_{1,2}^C$. For the FV scheme, when $n = 2$, the inertia-gravity frequency is a double root. The normalized positive non zero roots $\tilde{\omega}$ are displayed as a function of kh in Fig. 6 in the case $f = 0$ and $l = 0$ for simplicity, and each root is given its own color. The dotted black line indicates the analytical frequency $\tilde{\omega}_{1,2}^C = kh$.

As expected for the three schemes, $\tilde{\omega}$ is folding back to zero at least once, at $kh = \pi$ or/and at the end of the spectrum, namely $kh = 2\pi$. Indeed, the unknown $\bar{\mathbf{w}}_h$ and the test function ψ both belong to $V_{h, K_{el}}$, and hence the LBB or inf-sup stability condition [3, 26] is no longer satisfied. For $kh \neq 0$, the eigenvectors associated with zero frequency do not propagate. These stationary modes correspond to physical eigenmodes of the system which have their phase speed reduced to zero because of the deficiency of the numerical method and appear as stationary internode oscillations. If such erratic eigenmodes are left undamped, without resorting to stabilization procedures, they usually cause an accumulation of energy in the smallest-resolvable scale, leading to noisy/unstable approximations of the discrete elevation field η_h .

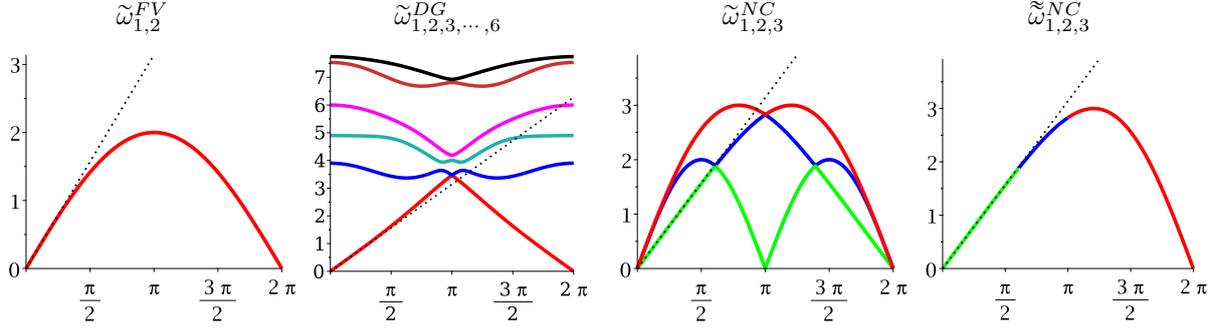


Figure 6: The positive normalized frequencies $\tilde{\omega}$ corresponding to the inertia-gravity waves for the FV, DG and NC schemes, in the case $p = q = 0$ (centered flux) and for $f = 0$ and $l = 0$. The dotted black line indicates the analytical frequency $\tilde{\omega}_{1,2}^C = kh$.

Apart from the presence of zero frequencies, the results of Fig. 6 are difficult to interpret for $\tilde{\omega}_{1,2,3,\dots,6}^{DG}$ and $\tilde{\omega}_{1,2,3}^{NC}$ due to the appearance of n positive solutions $\tilde{\omega}$ representing the n Fourier modes present in each eigenmode solution of the discrete equations. This fact has caused some confusion in the literature, with several works associating the single solution with n different solutions, one physical and $n - 1$ unphysical spurious modes [16, 18]. This interpretation would result in a multiple valued dispersion relation. Another approach, followed in [27, 31] in the case of the one-dimensional first-order advection equation for high-order FE and DG methods, considers each of the n solutions valid over only a limited wavenumber range, termed a branch. The physical solution is obtained by inspecting the spatial structure of the discrete solution for each branch, using the eigenvectors. The union of all branches then gives the complete dispersion relation named $\tilde{\omega}$. Each eigenmode is then associated with a single kh rather than treated as n different solutions with n different values of kh . In fact, such a procedure aim to eliminate the mathematical artifacts arising from symmetries in the Fourier analysis. It works well in the centered case for the NC scheme and $\tilde{\omega}_{1,2,3}^{NC}$ is shown in the last column of Fig. 6. However, the DG case is more complex to solve and the procedure described above does not seem appropriate.

4.5.2. The stabilized case with $(p, q) \neq (0, 0)$

When the stabilization procedure described in Section 3 is employed, namely with $(p, q) \neq (0, 0)$, the asymptotic expansion of the complex frequencies solutions of $P_{3n}(\tilde{\omega}) = 0$ in (70), for infinitesimal mesh spacing h are computed by using the MAPLE software for the FV, DG and NC schemes. After long and tedious algebraic manipulations, a super-convergent result is obtain for the discrete frequencies of the NC scheme compared to the DG one. When the Roe flux is employed with $q = 1$, however, the discrete frequencies of both schemes (ω_3^{DG} and ω_3^{NC}) exhibit sub-optimal rates of convergence for the slow mode.

Theorem 4.1. *Let \mathcal{E}_j , $j = 1, 2, 3, \dots, 12$, and \mathcal{F}_j and \mathcal{G}_j , $j = 1, 2, 3, 4$, be polynomial functions depending on k, l, p, q , except \mathcal{G}_3 and \mathcal{G}_4 , which only depend on k, l and not p, q , with*

$$\hat{\mathcal{E}}_j = \frac{\mathcal{E}_j + \frac{f^2}{gH} \mathcal{E}_{j+1}}{\frac{f^2}{gH} + k^2 + l^2}, \quad \hat{\mathcal{F}} = \frac{\sum_{j=1}^4 \left(\frac{f^2}{gH}\right)^{j-1} \mathcal{F}_j}{\left(\frac{f^2}{gH} + k^2 + l^2\right)^{5/2}}, \quad \hat{\mathcal{G}}_j = \frac{\mathcal{G}_j + \frac{f^2}{gH} \mathcal{G}_{j+1}}{\sqrt{\frac{f^2}{gH} + k^2 + l^2}}.$$

satisfying $\hat{\mathcal{E}}_{2j-1} > 0$, $j = 1, 2, 3, \dots, 6$, $\hat{\mathcal{F}} > 0$, and $\hat{\mathcal{G}}_{2j-1} > 0$, $j = 1, 2$, for all k, l, p and q .

In the limit as mesh spacing $h \rightarrow 0$ the following asymptotic results are obtained for the non normalized inertia-gravity ($\omega_{1,2}$) and geostrophic (ω_3) frequencies solutions of $P_{3n}(\tilde{\omega}) = 0$ in (70) in the case of:

- The FV scheme:

$$\begin{aligned} \omega_{1,2}^{FV} &= \omega_{1,2}^C - \hat{\mathcal{E}}_1 \mathbf{i}h \pm \hat{\mathcal{F}}h^2 + O(h^3), & \text{for } (p, q) \neq (0, 0), \\ \omega_3^{FV} &= \omega_3^C - \hat{\mathcal{E}}_3 \mathbf{i}h + O(h^3), & \text{for } (p, q) \neq (0, 0). \end{aligned}$$

- The DG scheme:

$$\omega_{1,2}^{DG} = \omega_{1,2}^C - \widehat{\mathcal{E}}_5 \mathbf{i} h^3 \pm \widehat{\mathcal{G}}_1 h^4 + O(h^5), \quad \text{for } (p, q) \neq (0, 0),$$

$$\omega_3^{DG} = \omega_3^C - \widehat{\mathcal{E}}_7 \mathbf{i} h^3 + O(h^5), \quad \text{for } p \neq 0,$$

$$\omega_3^{DG} = \omega_3^C \quad \text{and} \quad \omega_3^{DG} = \omega_3^C - \frac{(2\sqrt{2}-1)}{18} \frac{f^2}{\sqrt{gH}} \mathbf{i} h + O(h^2) \text{ triple root}, \quad \text{for } q = 1.$$

- The NC scheme:

$$\omega_{1,2}^{NC} = \omega_{1,2}^C \pm \widehat{\mathcal{G}}_3 h^4 - \widehat{\mathcal{E}}_9 \mathbf{i} h^5 + O(h^6), \quad \text{for } (p, q) \neq (0, 0),$$

$$\omega_3^{NC} = \omega_3^C - \widehat{\mathcal{E}}_{11} \mathbf{i} h^5 + O(h^7), \quad \text{for } p \neq 0,$$

$$\omega_3^{NC} = \omega_3^C - \sqrt{gH} \left(\frac{(9\sqrt{2}-1)}{84} (k-l)^2 + \frac{1}{6} kl + \frac{(\sqrt{2}+2)}{48} \frac{f^2}{gH} \right) \mathbf{i} h + O(h^2) \text{ double root}, \quad \text{for } q = 1.$$

There are $3n$ frequencies, solutions of $P_{3n}(\tilde{\omega}) = 0$ in (70), and hence at most $3(n-1)$ remaining frequencies need to be considered in addition to those mentioned in Theorem 4.1. This is the purpose of Theorem 4.2.

Theorem 4.2. *In the limit as mesh spacing $h \rightarrow 0$ with $f = 0$ (for simplicity) and $q = 1 - p$ (as in Table 1), the $3(n-1)$ remaining frequencies of $P_{3n}(\tilde{\omega}) = 0$ in (70) behave as high-frequency modes of the form*

$$\omega_j^{FV} = i \sqrt{gH} \frac{\alpha_j(p)}{h} + i O(h), \quad j = 4, 5, 6, \quad \text{with } -2(\sqrt{2} + 2) \leq \alpha_j \leq -2,$$

$$\omega_j^{DG} = \sqrt{gH} \frac{\beta_j(p) + i \gamma_j(p)}{h} + i O(h), \quad j = 4, 5, 6, \dots, 18, \quad \text{with } |\beta_j| \leq 7.617 \text{ and } -6\sqrt{2} \leq \gamma_j \leq 0,$$

$$\omega_j^{NC} = i \sqrt{gH} \frac{\delta_j(p)}{h} + i O(h), \quad j = 4, 5, 6, \dots, 9, \quad \text{with } -2(2\sqrt{2} + 1) \leq \delta_j \leq 0.$$

The values of $\alpha_j(p)$, $j = 4, 5, 6$, $\gamma_j(p)$, $j = 4, 5, 6, \dots, 18$, and $\delta_j(p)$, $j = 4, 5, 6, \dots, 9$, are graphed in Fig. 7 for the FV, DG and NC schemes, respectively, when $p \in [0, 1]$.

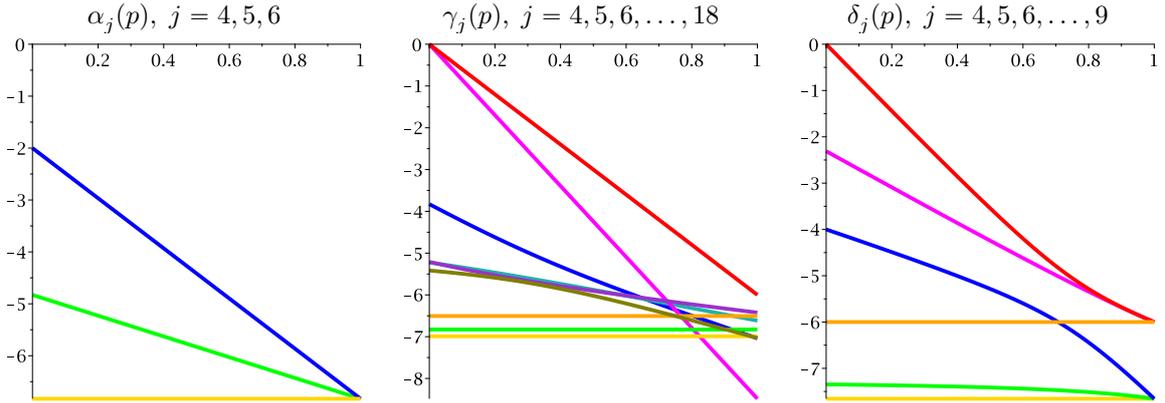


Figure 7: The values of $\alpha_j(p)$, $j = 4, 5, 6$, $\gamma_j(p)$, $j = 4, 5, 6, \dots, 18$, and $\delta_j(p)$, $j = 4, 5, 6, \dots, 9$, for the FV, DG and NC schemes, respectively, when $p \in [0, 1]$.

For the DG method, among the 15 values for $\gamma_j(p)$, $j = 4, 5, 6, \dots, 18$, we obtain $\gamma_{4,5}(p) = -6p$ (double root) and $\gamma_6(p) = -6p\sqrt{2}$, corresponding to the red (double) and magenta curves in Figure 7, respectively.

When $q = 1$ (in the case of the Roe flux) we have $\beta_{4,5,6}(p) = \gamma_{4,5,6}(p) = 0$ and hence $\omega_{4,5,6}^{DG} = iO(h)$. Such three frequencies coincide with the triple root $-\frac{(2\sqrt{2}-1)}{18} \frac{f^2}{\sqrt{gH}} ih + O(h^2)$ in Theorem 4.1 when $f \neq 0$.

For the NC scheme we obtain $\delta_4(p) = -2p(1 + \sqrt{2}) - 2 + 2\sqrt{(3 - \sqrt{2})p^2 - (1 + \sqrt{2})p + 1}$, corresponding to the red curve in Figure 7. The Roe flux then yields $\delta_4(p) = 0$ when $(p, q) = (0, 1)$, and the resulting frequency $\omega_4^{NC} = iO(h)$ coincides, along with the degenerated geostrophic frequency ω_3^{NC} , with the double root mentioned in Theorem 4.1.

As shown in Figure 7, $\alpha_j(p) \neq 0$, $j = 4, 5, 6$, for $p \in [0, 1]$, and that explains why the FV scheme gives rise to exactly two inertia-gravity and one geostrophic frequencies in Theorem 4.1, without being polluted by the presence of double or triple roots, contrary to the DG and NC methods.

Corollary 1. *The results of Theorems 4.1 to 4.2 and Fig. 7, demonstrate that the $3n$ discrete frequencies $\omega_{1,2,3,\dots,3n}$ solutions of $P_{3n}(\tilde{\omega}) = 0$ in (70) are stable at leading order for the FV, DG and NC methods, at least when $q = 1 - p$ with $p \in [0, 1]$. Further, it is observed graphically that the assumption $f = 0$ has no effect on the stability constraint of the high-frequency modes.*

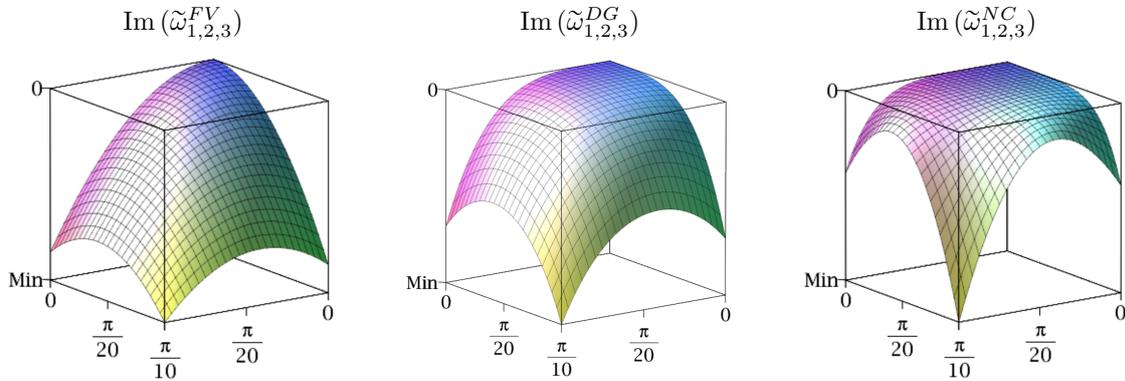


Figure 8: $\text{Im}(\tilde{\omega}_{1,2,3})$ corresponding to the inertia-gravity and geostrophic frequencies for the FV, DG and NC methods using the Rusanov flux when $0 \leq kh, lh \leq \pi/10$. The minimum (Min) values of the normalized frequencies are specified in Table 2.

To highlight the results of Theorem 4.1, the frequencies $\text{Im}(\tilde{\omega}_{1,2,3})$ corresponding to the FV, DG and NC methods, are shown in Fig. 8 for the Rusanov flux with $0 \leq kh, lh \leq \pi/10$. The minimum (Min) values of the frequencies are specified in Table 2 for $f = 0$ and $\lambda = 2$. For both values of the Coriolis parameter, we observe a better accuracy for the NC scheme compared with the DG method due to less numerical diffusion.

Table 2: Minimum values, corresponding to Min in Fig. 8, of $\text{Im}(\tilde{\omega}_{1,2,3})$ at $(kh, lh) = (\pi/10, \pi/10)$ for the FV, DG and NC schemes employing the Rusanov flux (Rus).

Rus	$\text{Im}(\tilde{\omega}_{1,2}^{FV})$	$\text{Im}(\tilde{\omega}_3^{FV})$	$\text{Im}(\tilde{\omega}_{1,2}^{DG})$	$\text{Im}(\tilde{\omega}_3^{DG})$	$\text{Im}(\tilde{\omega}_{1,2}^{NC})$	$\text{Im}(\tilde{\omega}_3^{NC})$
$f = 0$	-4.08×10^{-2}	-4.08×10^{-2}	-1.07×10^{-4}	0	-4.34×10^{-6}	0
$\lambda = 2$	-4.08×10^{-2}	-4.08×10^{-2}	-7.57×10^{-5}	-6.29×10^{-5}	-3.08×10^{-6}	-2.48×10^{-6}

We also computed $\text{Im}(\tilde{\omega}_{1,2,3})$ for the Roe flux and the three methods (not shown) and the minimum values of the frequencies are specified in Table 3. While the results for the inertia-gravity modes $\text{Im}(\tilde{\omega}_{1,2})$ compare well with those of the Rusanov flux, this is not the case for the geostrophic modes $\text{Im}(\tilde{\omega}_3^{DG})$ and $\text{Im}(\tilde{\omega}_3^{NC})$. Indeed, the values of $\text{Im}(\tilde{\omega}_3^{DG})$ and $\text{Im}(\tilde{\omega}_3^{NC})$ are close to those of $\text{Im}(\tilde{\omega}_3^{FV})$ and such a low

accuracy is reflected in Theorem 4.1 by the presence of the double and triple roots when $q = 1$, yielding a poor accuracy for $\text{Im}(\tilde{\omega}_3^{DG})$ and $\text{Im}(\tilde{\omega}_3^{NC})$ in the case of the Roe flux. The problematic values of $\text{Im}(\tilde{\omega}_3^{DG})$ and $\text{Im}(\tilde{\omega}_3^{NC})$ are framed in Table 3 to facilitate comparison with the corresponding values of Table 2.

Table 3: As for Table 2, but in the case of the Roe flux.

Roe	$\text{Im}(\tilde{\omega}_{1,2}^{FV})$	$\text{Im}(\tilde{\omega}_3^{FV})$	$\text{Im}(\tilde{\omega}_{1,2}^{DG})$	$\text{Im}(\tilde{\omega}_3^{DG})$	$\text{Im}(\tilde{\omega}_{1,2}^{NC})$	$\text{Im}(\tilde{\omega}_3^{NC})$
$f = 0$	-3.48×10^{-2}	0	-9.50×10^{-5}	0	-2.97×10^{-6}	-3.27×10^{-2}
$\lambda = 2$	-2.40×10^{-2}	-2.30×10^{-2}	-6.46×10^{-5}	-3.94×10^{-2}	-2.43×10^{-6}	-6.88×10^{-2}

Finally, in the case of the Rusanov flux, positive values of $\text{Re}(\tilde{\omega})$, corresponding to $\tilde{\omega}_1^C$, and $\text{Im}(\tilde{\omega})$ are displayed in Fig. 9 for the DG method (with $l = 0$) and the NC scheme (with $k = l$) when $f = 0$.

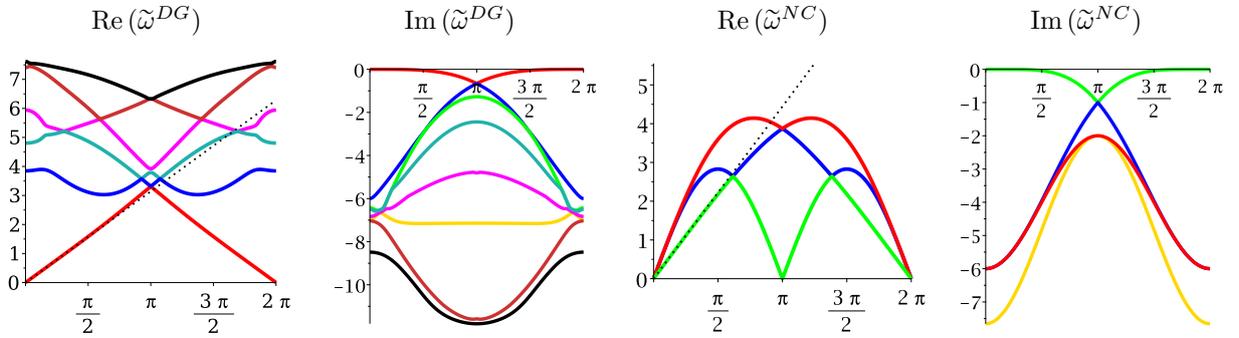


Figure 9: Positive values of $\text{Re}(\tilde{\omega})$, corresponding to $\tilde{\omega}_1^C$, and $\text{Im}(\tilde{\omega})$ for the DG (when $l = 0$) and NC (when $k = l$) schemes in the case of the Rusanov flux when $f = 0$. Several frequencies are double roots of $P_{3n}(\tilde{\omega})$ in (70). The dotted black line indicates the real analytical frequency $\tilde{\omega}_1^C$.

Each root is given its own color and several frequencies are double roots. The graphs of $\text{Re}(\tilde{\omega}^{DG})$ and $\text{Re}(\tilde{\omega}^{NC})$ are close to the results obtained in Fig. 6. However, the non zero imaginary parts $\text{Im}(\tilde{\omega}^{DG})$ and $\text{Im}(\tilde{\omega}^{NC})$ now prevent the presence of an erratic stationary mode both at the end of the spectrum and at $kh = \pi$ for the NC scheme. For $f \neq 0$ and other choices for k and l than those considered above, particularly in the 2D case, the graph of the dispersion relation is too complicated to be solved and analyzed, and as for the DG scheme in Section 4.5.1 (the five high-frequency modes in Fig. 6), the question of whether the high-frequency modes are mathematical artifacts or not remains, to our knowledge, open. Nevertheless, the high-frequency modes observed in Fig. 9 for $\text{Re}(\tilde{\omega}^{DG})$ are severely damped due to the presence of $\text{Im}(\tilde{\omega}^{DG})$.

Remark 2. We emphasize that the results obtained for the PVM-2 and PVM-4 methods in Section 4.5.2 give very similar results compared to those obtained with the Rusanov flux, and hence they are not shown.

To try to understand why the Roe flux leads to a loss of convergence for the geostrophic modes ω_3^{DG} and ω_3^{NC} , the stabilization term in (35), corresponding to the momentum equations, is rewritten in the case $(p, q) = (0, 1)$. This yields $([\mathbf{u}] \cdot \mathbf{n}_{ed}, \mathbf{n}_{ed})$, namely $([\mathbf{u}], 0)$ on vertical faces and $(0, [\mathbf{v}])$ on horizontal faces. The presence of such zero values for the flux is reflected, via the discrete equations obtained in Section 3.3, into the zero value entries of the stabilization matrices \hat{D}_j , $j = 1, 2, 3$, computed in Section 4 for the DG and NC schemes. For the DG scheme, \hat{D}_3 has two rows (3 and 6) and two columns (3 and 6) of zero elements (corresponding to nodes 3 and 6 of Fig. 5), and hence the Coriolis matrix in (63) is not stabilized for \hat{u}_3 , \hat{v}_3 , \hat{u}_6 and \hat{v}_6 . Similar conclusions can be drawn for \hat{D}_3 in the case of the NC method at node 2 of Fig. 5, while the FV method is clearly not concerned with such problems. These observations may reveal why the Roe flux leads to a loss of convergence for the geostrophic modes ω_3^{DG} and ω_3^{NC} .

5. Numerical simulations for the linear and non-linear model problems

In order to validate the results obtained in Theorem 4.1, three numerical tests are proposed in Section 5.1 using the linear variational formulation (14). The non-linear variational formulation (13) is then employed in Section 5.2. In both sections, the third-order Runge-Kutta time stepping scheme (RK3) is used with a Courant–Friedrichs–Lewy (CFL) stability criterion such that $\text{CFL} := \sqrt{gH}\Delta t/h = 0.1$. Let CFL^{NC} and CFL^{DG} denote the CFL criteria permitting the maximum allowable time step for the NC and DG schemes, respectively. For all the numerical simulations performed in Section 5 we observed that

$$\text{CFL}^{NC} \simeq 0.30 \quad \text{and} \quad \text{CFL}^{DG} \simeq 0.18, \quad \text{namely} \quad \text{CFL}^{NC} \simeq \frac{5}{3}\text{CFL}^{DG}, \quad (72)$$

regardless of the choice of the numerical flux. A fifth order strong-stability preserving Runge-Kutta scheme (SSPRK (5, 5)) [44] was also employed and it led to almost indistinguishable graphics compared to those obtained with the use of the RK3 scheme. Further, mass is conserved up to machine precision for all numerical experiments conducted in Section 5, and hence the results are not shown.

The domain extent is an idealized $L_1 \times L_2$ rectangular ocean basin discretized on Mesh 1 and Mesh 2, made up of biased right isocles and unstructured triangles with smoothing, respectively, as shown in Fig. 10. Remember that Mesh 1 is used in Section 4 for the computation of the dispersion relations. Further, the no-normal flow boundary condition $\mathbf{u} \cdot \mathbf{n} = 0$ is employed on the domain boundary $\partial\Omega$ in Section 5.

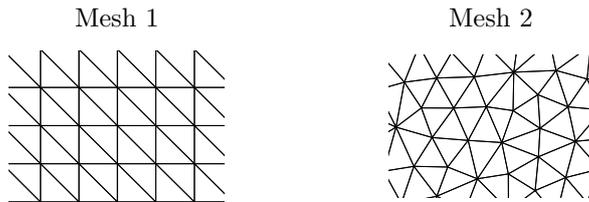


Figure 10: A window of Mesh 1 and Mesh 2, made up of biased right isocles and unstructured triangles with smoothing, respectively.

5.1. The linear case

5.1.1. The propagation of pure gravity waves

The first experiment examines the ability of the numerical schemes to reproduce the damping effect induces by $\hat{\mathcal{E}}_1$, $\hat{\mathcal{E}}_5$ and $\hat{\mathcal{E}}_9$ on the solution via $\omega_{1,2}$ in Theorem 4.1. To do so, the propagation and dispersion of pure gravity waves are simulated with $L_1 = L_2 = 10\,000$ km and the Coriolis parameter is set to zero. The discrete SW equations are solved with a Gaussian distribution of the surface elevation prescribed at initial time and initially centered in the domain. That is,

$$\eta(t = 0) := \eta_0 = A e^{-(x^2+y^2)/B^2}, \quad (73)$$

and the initial velocity field is set to zero. The Gaussian distribution parameters that define η_0 are set to $A = 1$ m and $B = 130$ km. A flat bottom is assumed and the mean depth is $H = 1\,000$ m, yielding a phase speed of the surface gravity waves $\sqrt{gH} = 100$ m s⁻¹ by assuming $g = 10$ m s⁻². The initial perturbation amplitude thus represents 0.1% of the total depth. Further, the e-folding radius of the initial Gaussian, namely the distance $\sqrt{x^2 + y^2} = B$ from the origin for which $\eta = A e^{-1}$, is resolved by at least six velocity nodes with a node spacing of approximately 20 km for the triangulation $\tau_{h/2}$.

In the absence of shocks, as is the case in this paper, the integral of total energy, named E_t in the linear case, with $E_t := \frac{1}{2} \int_{\Omega} (H(u^2 + v^2) + g\eta^2) d\Omega$ (up to an additive constant), should be preserved over time for smooth solutions, namely $\frac{dE_t}{dt} = 0$. However, due to the dissipative nature of the numerical fluxes employed in this study to stabilize the computed solutions, E_t is expected to decrease over time, particularly for simulations involving large time scales. The purpose of the present section is to examine this issue.

The damping effect induces by $\hat{\mathcal{E}}_1$, $\hat{\mathcal{E}}_5$ and $\hat{\mathcal{E}}_9$ on the solution via $\omega_{1,2}$ in Theorem 4.1 is first investigated. Let $\hat{\mathcal{E}}_5|_{Roe}$ (the coefficient of $-ih^3$) denote the function $\hat{\mathcal{E}}_5$ in the Roe flux case. The ratios graphed in Fig. 11 are then a measure of the damping introduced by the fluxes, compared to each other, for the DG scheme.

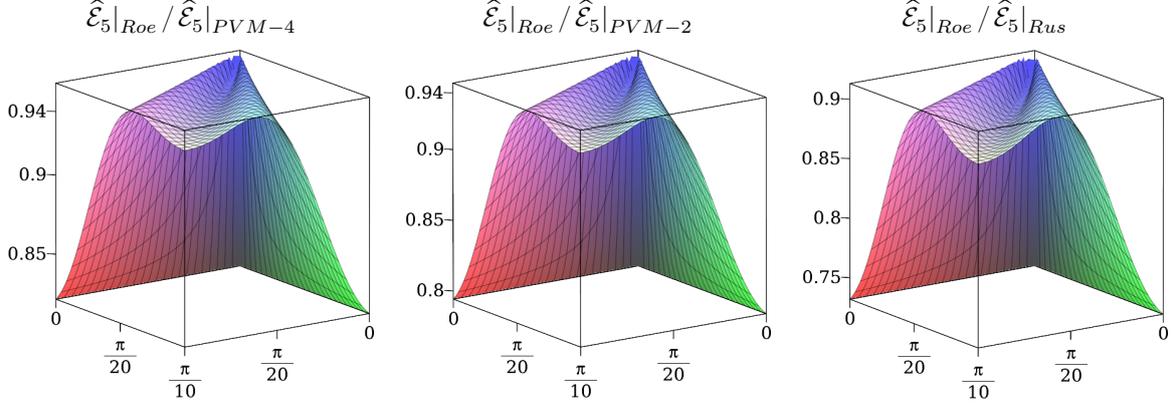


Figure 11: Three ratios (using the nomenclature defined in Theorem 4.1) measuring the damping introduced by the fluxes, compared to each other, for the DG scheme.

The ratios are graphed for $(kh, lh) \in [0, \pi/10]$ but identical curves and numerical values are observed if $(kh, lh) \in [0, \pi]$ for example. By employing $\hat{\mathcal{E}}_1$ (the coefficient of $-ih$) instead of $\hat{\mathcal{E}}_5$ to match the case of the FV scheme, curves of similar shape and values to that of Fig. 11 were obtained (not shown). For the NC scheme, $\hat{\mathcal{G}}_3$ (the coefficient of $\pm h^4$) only depends on k and l and not p and q , as mentioned in Theorem 4.1, and thus the damping term corresponding to $-\hat{\mathcal{E}}_9 ih^5$ is not expected to be discriminatory for the fluxes.

The numerical experiment described at the beginning of the section is now performed and E_l , normalized by its initial value at $t = 0$, is computed on Mesh 1 and Mesh 2. The results are graphed in Fig. 12 with $h = 40$ km and in Fig. 13 with $h = 20$ km for the DG and NC schemes and the four fluxes up to 13 hours and 20 mn, namely the time it takes for the solution to propagate over 4800 km, roughly 30 mn before reaching the domain boundary. The results for the FV method are not shown since the energy lost half and two thirds of its initial value, after only one and two hours of simulation, respectively.

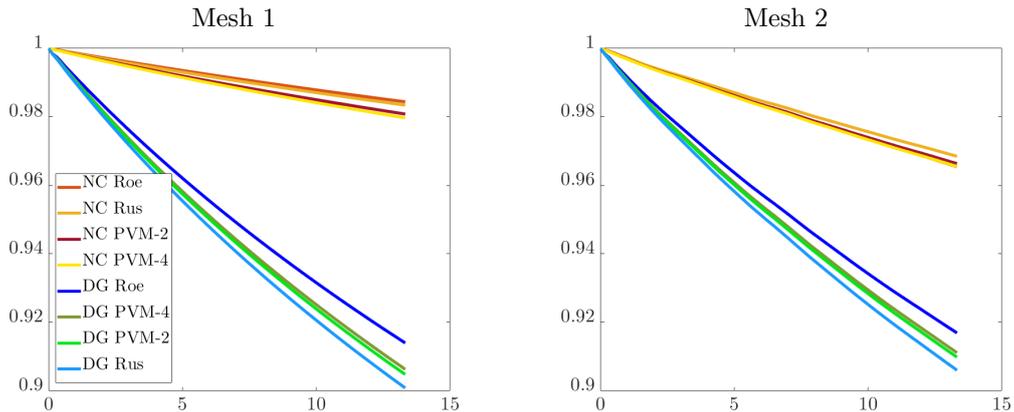


Figure 12: The decrease of the total energy E_l , normalized by its initial value at $t = 0$, on Mesh 1 and Mesh 2 with a resolution of $h = 40$ km, in the case of the DG and NC schemes by using the Rusanov, Roe, PVM - 2 and PVM - 4 fluxes examined in this paper during 13 hours and 20 mn, namely the time it takes for the solution to propagate over 4800 km.

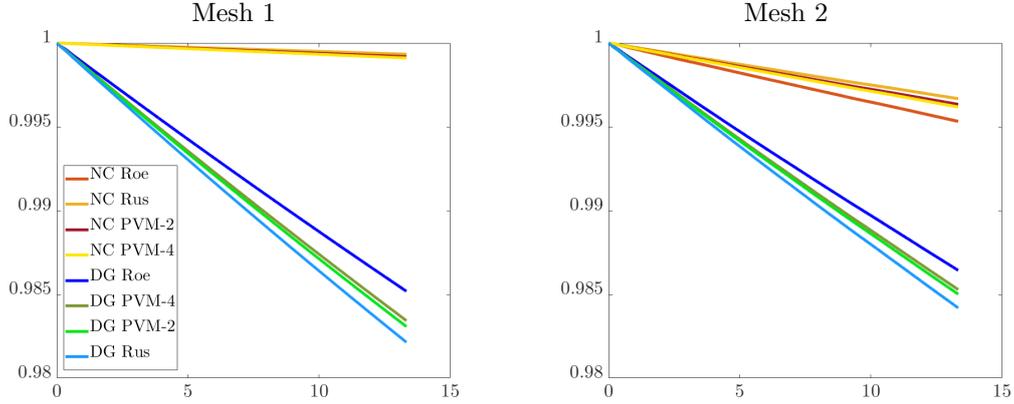


Figure 13: As for Fig. 12, but in the case $h = 20$ km.

The curves of Figs. 12 and 13 are in precise agreement with the results of Fig. 11. Indeed, we have $\hat{\mathcal{E}}_5|_{Roe} < \hat{\mathcal{E}}_5|_{PVM-4} < \hat{\mathcal{E}}_5|_{PVM-2} < \hat{\mathcal{E}}_5|_{Roe}$ for all $(kh, lh) \in [0, \pi/10]$ in Fig. 11, with

$$\begin{aligned} \hat{\mathcal{E}}_5|_{Roe} / \hat{\mathcal{E}}_5|_{PVM-4} &\in [0.822, 0.945], & \hat{\mathcal{E}}_5|_{Roe} / \hat{\mathcal{E}}_5|_{Rus} &\in [0.732, 0.886], \\ \hat{\mathcal{E}}_5|_{Roe} / \hat{\mathcal{E}}_5|_{PVM-2} &\in [0.794, 0.931], \end{aligned}$$

and this remains true for all $(kh, lh) \in [0, \pi]$. Consequently, as long as the DG method is concerned in the case of the propagation of pure gravity waves, the dissipation term being of the form $e^{-\hat{\mathcal{E}}_5 h^3 t}$ at leading order (see $\omega_{1,2}^{DG}$ in Theorem 4.1), we can expect the dissipation to be greater for the Rusanov flux than for the PVM-2 flux, itself greater than for the PVM-4 flux, itself greater than for the Roe flux. This is exactly what is observed in Figs. 12 and 13 on both Mesh 1 and Mesh 2.

We now examine the eventual numerical dispersion during the propagation of the initial Gaussian over 5000 km, at the time it is being reflected by the basin wall. Time sequences for η are displayed in Fig. 14 for the NC scheme and the Rusanov flux on Mesh 2 with $h = 40$ km. Each sequence is colored according to its own max (red) and min (blue). Similar results are obtained on Mesh 1 and also for the DG scheme. The radial symmetry of η and the wave speed (100 m s^{-1}) are very well respected, and no numerical dispersion is observed. Similar conclusions hold for u and v . For long-term simulations, even after the initial Gaussian is being reflected several times by the basin wall, the radial symmetry of the wave is very well preserved.

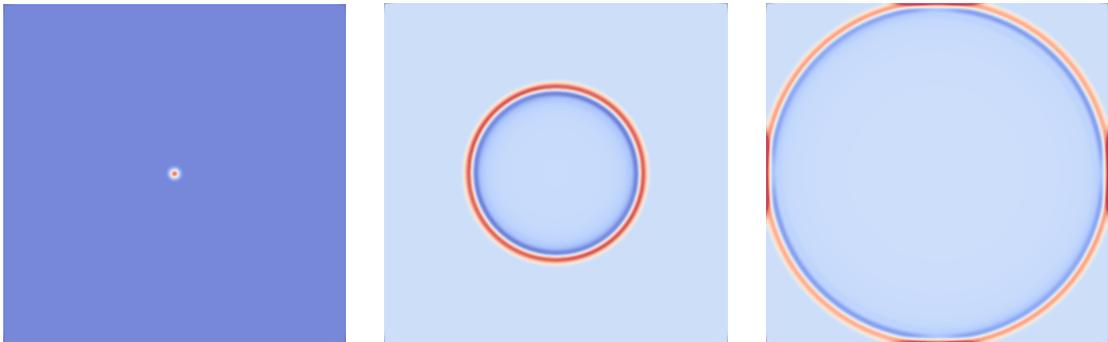


Figure 14: Time sequences for η at different stages of gravity wave propagation and dispersion: at initial time (left panel), after the wave has propagated over 2500 km (middle panel) and 5000 km (right panel) i.e. at the time it is being reflected by the basin wall. The computations are performed on Mesh 2 with $h = 40$ km by using the NC scheme and the Rusanov flux.

5.1.2. An almost-balanced flow

The second test examines the evolution of an almost-balanced flow in the case of the evolution of a vortex at midlatitudes with $L_1 = 2000$ km and $L_2 = 1200$ km. At initial time, an anticyclonic velocity field is in geostrophic equilibrium with the prescribed Gaussian distribution of the surface-elevation defined in (73), now centered at the point (1400 km, 600 km), namely

$$\eta_0 = A e^{-(x^2+y^2)/B^2}, \quad \mathbf{u}(\mathbf{x}, t = 0) = (g/f) \mathbf{k} \times \nabla \eta_0 = -(g/f) \mathbf{rot} \eta_0, \quad (74)$$

Both η and \mathbf{u} are evaluated pointwise at initial time by sampling the exact expressions in (73) and (74). Although η and \mathbf{u} are in exact analytic balance, such that their tendencies are identically zero, the initial conditions are not in exact *discrete* geostrophic balance and a non-zero time tendency of the order of the spatial truncation error is obtained. This leads to small-amplitude temporal oscillations being generated and we examine here how initial discrete imbalances propagate in time. Indeed, in most numerical models, initial conditions obtained from observations and interpolated to the model grid, are usually *out-of-balance*.

The coordinates x and y are oriented eastward and northward, respectively, and the latter is measured from a reference latitude $\varphi_0 = 28.6^\circ$ N. The Coriolis parameter f is evaluated at φ_0 and it is held constant with $f = f_0 = 2\Omega_e \sin \varphi_0$, where $\Omega_e = 7.2921 \times 10^{-5}$ rad s $^{-1}$ is the angular frequency of the Earth's rotation, leading to $f_0 = 6.9813 \times 10^{-5}$ s $^{-1}$. The choice $H = 100$ m results in a phase speed for gravity waves of 4 m s $^{-1}$ by using the reduced gravity $g' = 0.16$ m s $^{-2}$ instead of $g = 9.81$ m s $^{-2}$. The radius of deformation is then $R_d = \sqrt{g'H}/f_0 \simeq 57.3$ km. By setting $A = 58.25$ m and $B = 130$ km in (73), the initial maximum surface azimuthal velocity is 0.9 m s $^{-1}$. These values are representative of the oceanic eddy circulation at mid-latitude. The adopted configuration corresponds to that of a reduced-gravity model with a single active upper layer overlying a second passive layer, implicitly assumed, that is infinitely deep and at rest [10]. Let $\theta = \eta - \phi$, where $\phi(x, y, t)$ is the interface displacement. The thickness of the single moving layer above a motionless abyss is thus $H + \theta$. The resulting governing equations are identical to that of (5) and (6), and (14) is employed with the initial condition (74), provided that we replace g by g' and η by θ .

Since $\beta = 0$, the Coriolis parameter is held constant and this yields $\omega_3^C = 0$ in (10). Consequently, the solution does not propagate in space and the numerical solution should preserve the steady state. However, because the initial conditions are not in exact *discrete* geostrophic balance, the evolution of an almost-balanced flow in the presence of small-amplitude gravitational oscillations is examined.

At selected time steps n , a convergence analysis is performed by computing the ratio

$$R_\zeta(t^n = n\Delta t) = \frac{\|\zeta(t^0) - \zeta_h(t^n)\|_{L^2}}{\|\zeta(t^0) - \zeta_{h/2}(t^n)\|_{L^2}}, \quad n = 1, 2, 3, \dots$$

where ζ is equal to either θ or the flow-speed field $\sqrt{u^2 + v^2}$. The quantity $\ln R_\zeta / \ln 2$ is represented in Fig. 15 for the NC, DG and FV methods by using the Roe and Rusanov fluxes up to 15 weeks of simulation. The calculations are performed on Mesh 1, and two uniform meshes are used, a 50×30 grid with $h = 40$ km and a 100×60 grid with $h = 20$ km, to compute ζ_h and $\zeta_{h/2}$, respectively, namely for $h/2 = 20$ km.

The convergence rates observed in Fig. 15 for R_ζ , when $\zeta = \theta$ and $\zeta = \sqrt{u^2 + v^2}$, faithfully reproduce the results of Theorem 4.1 obtained on Mesh 1. Indeed, the results of Fig. 15 for the Rusanov flux on Mesh 1 clearly exhibit, after only a few days of simulation, the influence of the fourth-order accuracy obtained for the NC frequency in Theorem 4.1, while the DG scheme does not show such a super-convergence rate. However, when the Roe scheme is employed, the NC methods exhibits a second-order rate of accuracy for both θ and $\sqrt{u^2 + v^2}$, while the DG method yields only a first-order convergence rate for $\sqrt{u^2 + v^2}$. Such behaviors can likely be attributed to the presence of $iO(h)$ spurious solutions for the geostrophic modes ω_3^{DG} and ω_3^{NC} in Theorem 4.1. Consequently, the results obtained with the Roe scheme will not be considered in the following. We also point out the poor convergence rate of the FV method, as expected. Further, the graphics for $\zeta = u$ and $\zeta = v$ are almost indistinguishable from those of $\zeta = \sqrt{u^2 + v^2}$ and hence they are not shown. Likewise, the graphs corresponding to the PVM - 2 and PVM - 4 fluxes led to almost indistinguishable results from those of the Rusanov flux for both the DG and NC schemes.

The normalized total energy $E_i(t) = \frac{1}{2} \int_\Omega (H(u^2 + v^2) + g'\theta^2) d\Omega$ and maximum values of θ (i.e. $\|\theta\|_{L^\infty}$) have been computed on Mesh 1 and Mesh 2 for $h = 40$ km and $h = 20$ km. These are graphed in Fig. 16

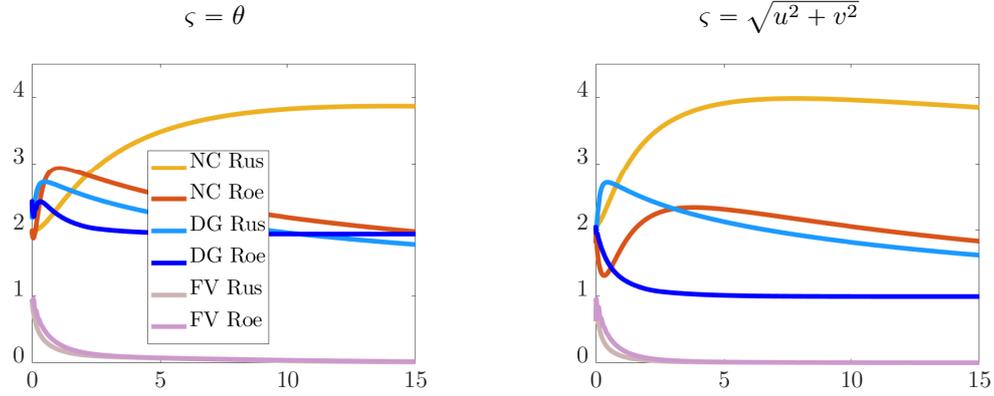


Figure 15: The quantity $\ln \mathcal{R}_c / \ln 2$ with $h/2 = 20$ km, for $\zeta = \theta$ and $\zeta = \sqrt{u^2 + v^2}$, in the case of the NC, DG and FV methods using the Roe and Rusanov fluxes up to 15 weeks of simulation on Mesh 1.

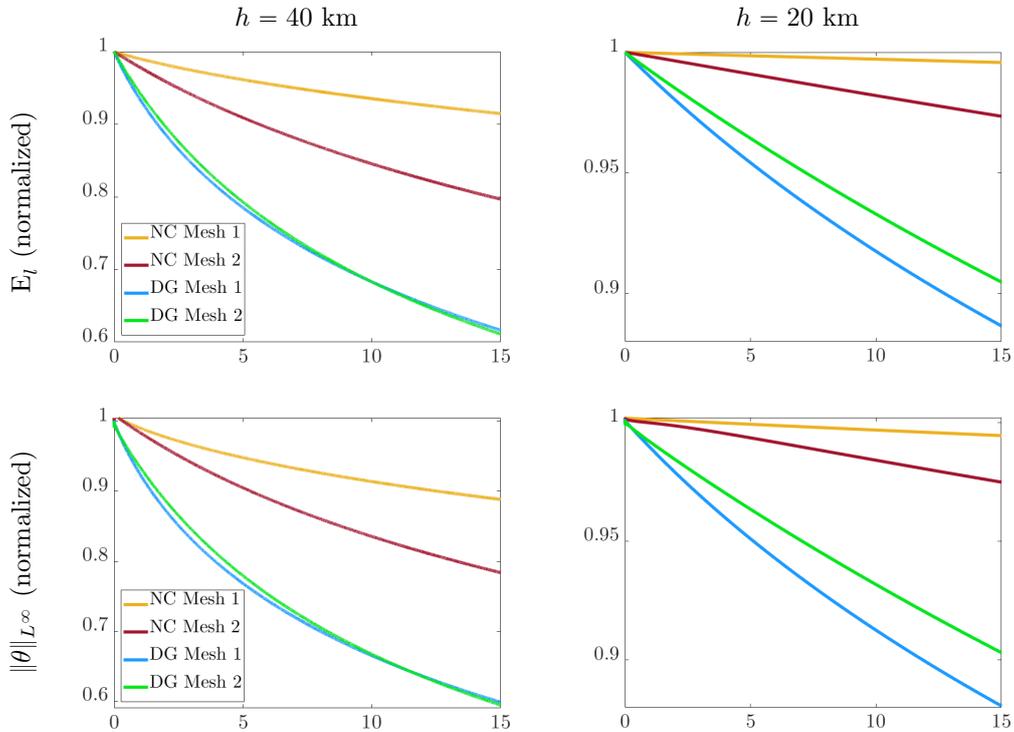


Figure 16: The normalized quantities E_l and $\|\theta\|_{L^\infty}$ in the case $f = f_0$ up to 15 weeks of simulation on Mesh 1 and Mesh 2, with $h = 40$ km and $h = 20$ km, for the DG and NC methods using the Rusanov flux.

for the DG and NC schemes in the case of the Rusanov flux during 15 weeks of simulation. The numerical results are supposed to reflect the steady state. Clearly, E_l and $\|\theta\|_{L^\infty}$ are much better conserved (less damped) with the NC method than with the DG one. Indeed, on Mesh 1 with $h = 20$ km, E_l and $\|\theta\|_{L^\infty}$ only lose 0.42% and 0.56% of their initial values, respectively, after 15 weeks of simulation. On Mesh 2 these values are 2.4% and 2.5%, respectively. The graphics of the normalized values of $\|\sqrt{u^2 + v^2}\|_{L^\infty}$ (not shown) are similar to those of $\|\theta\|_{L^\infty}$. Finally, the numerical solutions computed with the PVM - 2 and PVM - 4 fluxes are similar to those obtained using the Rusanov flux for both the DG and NC methods.

5.1.3. Simulation of an anticyclonic eddy at midlatitudes

The third experiment reproduces the numerical test of Section 5.1.2, with the same initial conditions, except that the Coriolis parameter f is no longer constant and instead it varies with latitude. The β -plane approximation $f = f_0 + \beta y$ is used [32], where $\beta = 2(\Omega_e/R) \cos \varphi_0$ is the so-called β -parameter and $R = 6371$ km is the Earth's radius, leading to $\beta = 2.0098 \times 10^{-11} \text{ m}^{-1} \text{ s}^{-1}$. Thanks to the slowing down of fast-moving surface gravity waves via the choice $\sqrt{gH} = 4 \text{ m s}^{-1}$, the slowly-propagating Rossby modes are simulated in the case of the evolution of a typical anticyclonic eddy at mid-latitude. During the first inertial period ($2\pi/f_0 \simeq 25$ h) the initial condition adjusts to the β -plane balance of the model. After this initial adjustment, the eddy evolves purely westward at an average translation speed of $\beta R_d^2 \simeq 5.7 \text{ km day}^{-1}$. Contrary to the test of Section 5.1.2, $\|\theta\|_{L^\infty}$ is now supposed to decrease as the Rossby mode propagates and evolves in time.

As in Section 5.1.2, the normalized quantities E_l and $\|\theta\|_{L^\infty}$ have been computed, for $h = 40$ km and $h = 20$ km, on Mesh 1 and Mesh 2, up to 15 weeks of simulation for the DG and NC schemes using the Rusanov flux. On both Mesh 1 and Mesh 2, the graphs of E_l (not displayed) are almost indistinguishable from those of Fig. 16. The results of $\|\theta\|_{L^\infty}$ are shown in Fig. 17. In addition, $\|\theta\|_{L^\infty}$ has been computed for $h = 10$ km (not displayed) and on such a high-resolution mesh the four curves representing $\|\theta\|_{L^\infty}$ appear superimposed on each other, suggesting that the numerical solutions have converged. Further, these are almost indistinguishable from the curve of $\|\theta\|_{L^\infty}$ obtained in the NC case when $h = 20$ km on Mesh 1, itself close to the case $h = 40$ km. From the results of Fig. 17, we deduce that the NC method has already almost converged on Mesh 1 when $h = 40$ km, reflecting the high order of convergence of the NC method. Again, the graphs obtained with the PVM – 2 and PVM – 4 fluxes are very similar to those of Fig. 17. Finally, no numerical dispersion was detected in the graphs of θ , u and v during the 15-week simulation period.

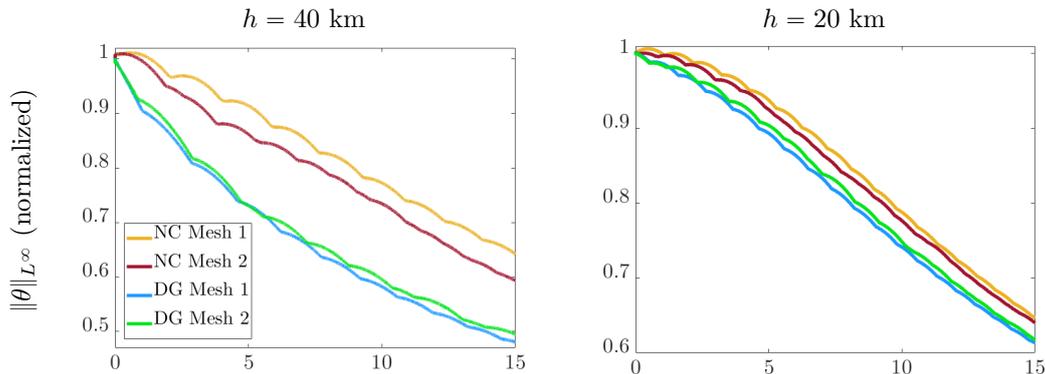


Figure 17: As for Fig. 16 for $\|\theta\|_{L^\infty}$, but in the case $f = f_0 + \beta y$, when the eddy propagates.

5.2. The non-linear case: Simulation of an anticyclonic eddy at midlatitudes

The test of Section 5.1.3 is now reproduced with the same initial conditions, but in the non-linear case. The governing equations are identical to that of (1) and (2) and the variational formulation (13) is employed, provided that g is replaced by g' and η by θ , as for the linear formulation in Sections 5.1.2 and 5.1.3. The total energy, named E_{nl} in the non-linear case, is defined as $E_{nl} := \frac{1}{2} \int_{\Omega} ((H + \theta)(u^2 + v^2) + g'\theta^2) d\Omega$ (up to an additive constant). As in Section 5.1, the total energy is examined to quantify the dissipative nature of the numerical fluxes. The quantities E_{nl} , $\|\theta\|_{L^\infty}$ are computed for two different resolutions, $h = 20$ km and $h = 40$ km, on Mesh 1 and Mesh 2 using the Rusanov flux (26) suitable for the non-linear case. The normalized results are graphed in Fig. 18 after 15 weeks of simulation for the DG and NC schemes. Clearly, E_l and $\|\theta\|_{L^\infty}$ are much better conserved (less damped) with the NC method than with the DG one.

The curves of E_{nl} are quite similar to that of E_l in Fig. 16 (the linear case) when $f = f_0$. On Mesh 1 with $h = 20$ km, E_{nl} loses 0.64% and 12.5% of its initial values, for the NC and DG methods, respectively, after 15 weeks of simulation. On Mesh 2 these values are 2.7% and 9.5%, respectively.

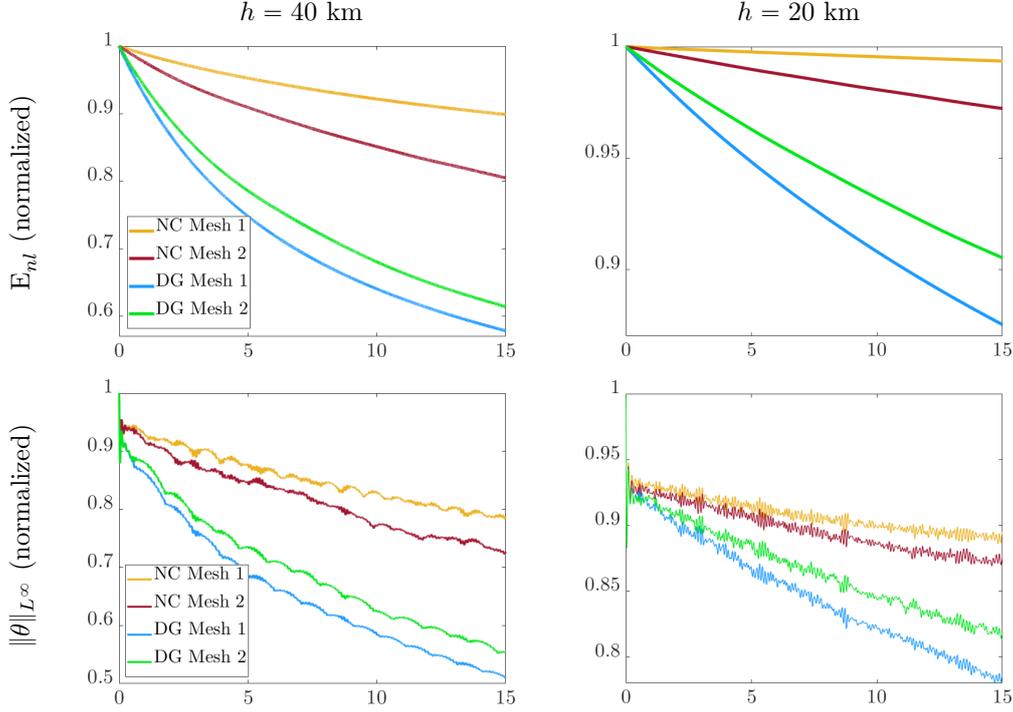


Figure 18: The normalized quantities E_{nl} and $\|\theta\|_{L^\infty}$ in the non-linear case with $f = f_0 + \beta y$ up to 15 weeks of simulation on Mesh 1 and Mesh 2, with $h = 40$ km and $h = 20$ km, for the DG and NC methods using the Rusanov flux.

Since u and v are initially in geostrophic balance, shortly after initialization, there will be a readjustment of the flow toward a gradient wind balance on the β -plane. During this time (the inertial period $2\pi/f_0 \simeq 25$ h), θ loses approximately 8% of its initial amplitude. However, the slight imbalance in the initial conditions does not affect the long-term evolution of the eddy. This initial 8% reduction in θ is followed by a further reduction in its amplitude, partly reflecting the dissipative nature of the numerical scheme. On Mesh 1 with $h = 20$ km, after the inertial period and up to 15 weeks of simulation, $\|\theta\|_{L^\infty}$ loses 3.1% and 13.5% of its value, for the NC and DG methods, respectively. On Mesh 2 these values are 4.9% and 10.3%, respectively.

The graphics for $\sqrt{u^2 + v^2}$ are first displayed in Fig. 19 after 15 weeks of propagation on Mesh 1, with $h = 40$ km and $h = 20$ km, for the NC and DG methods using the Rusanov flux. The legend at $t = 0$ is kept unchanged after 15 weeks of simulation. As predicted by Rossby wave dynamics, the eddy migrates westward and its westerly course exhibits a southwesterly drift that is due to non-linear effects. The average translation is in good agreement with that predicted by theory ($\beta R_d^2 \simeq 5.7$ km day $^{-1}$) and the cyclonic eddy formed by the wake progressively intensifies. The damping is quite severe for the DG method. Indeed, when $h = 20$ km, we obtain $\|\sqrt{u^2 + v^2}\|_{L^\infty} = 0.95$ m s $^{-1}$ and 0.75 m s $^{-1}$ for the NC and DG methods, respectively, after 15 weeks of propagation on Mesh 1. These values are 0.92 m s $^{-1}$ and 0.80 m s $^{-1}$, respectively, on Mesh 2. The test continued up to 42 weeks, after the eddy interacts with the western boundary and coastally trapped waves propagated around the entire basin. The graphs of $\sqrt{u^2 + v^2}$ on Mesh 1 with $h = 20$ km are displayed in Fig. 20 for the DG and NC methods. Again the DG solution exhibits a severe damping. At this stage of the propagation, a finer mesh would be needed close to the boundary, to capture the recirculation regions.

Finally, the non-linear Roe flux (25) is employed and $\sqrt{u^2 + v^2}$ is shown in Fig. 20 for the NC method (right panel) after 15 weeks of propagation on Mesh 1, with $h = 40$ km. We obtain $\|\sqrt{u^2 + v^2}\|_{L^\infty} = 0.55$ m s $^{-1}$, instead of 0.71 m s $^{-1}$ for the Rusanov scheme in Fig. 19 (first line, middle panel). Such values reflect the asymptotic result obtained for ω_3^{NC} (and also ω_3^{DG}) in Theorem 4.1 when $q = 1$ (the Roe flux), suggesting, at least, a pronounced damping of the geostrophic frequency in the linear case with $f = f_0$.

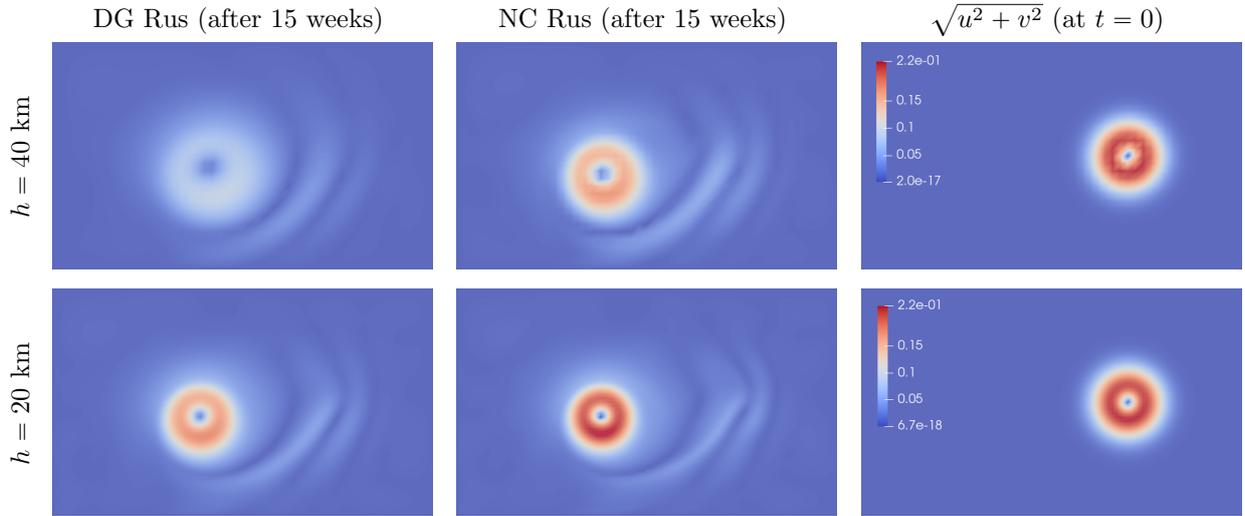


Figure 19: The flow-speed field $\sqrt{u^2 + v^2}$ at initial time (right panel) and after 15 weeks of propagation on Mesh 1, with $h = 40$ km and $h = 20$ km, for the NC (middle panel) and DG (left panel) methods using the Rusanov flux. The legend at $t = 0$ is kept unchanged after 15 weeks of simulation. The dimensionalized values of $\sqrt{u^2 + v^2}$ at $t = 0$ are recovered by multiplying those of the legend by $\sqrt{gH} = 4 \text{ m s}^{-1}$ leading to $\|\sqrt{u^2 + v^2}\|_{L^\infty} = 0.91 \text{ m s}^{-1}$.

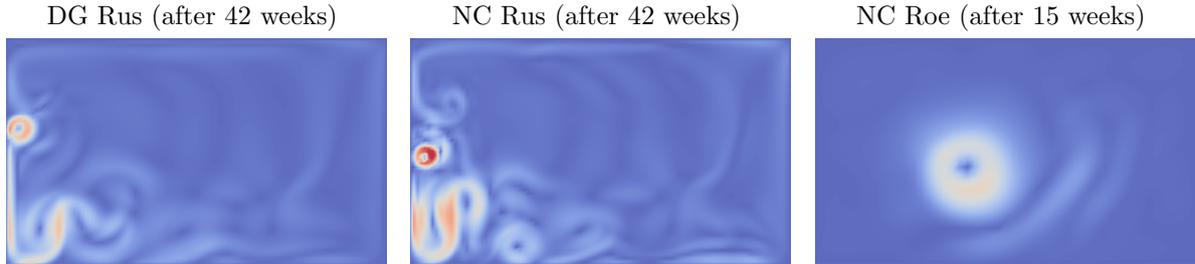


Figure 20: The flow-speed field $\sqrt{u^2 + v^2}$ after 42 weeks of propagation on Mesh 1, with $h = 20$ km, for the NC (middle panel) and DG (left panel) methods using the Rusanov flux. The flow-speed field $\sqrt{u^2 + v^2}$ after 15 weeks of propagation on Mesh 1, with $h = 40$ km, for the NC (right panel) method using the Roe flux. The legend of Fig. 19 is kept unchanged.

6. Conclusions

This paper appears to be the first study of the Fourier/stability analysis of the 2D shallow-water model for the linear discontinuous (DG) and non conforming (NC) Galerkin discretizations. The discrete formulations are stabilized by means of a family of numerical fluxes and the dispersion relations of the DG and NC schemes are constructed. A long time stability result is proven for all schemes and fluxes examined here by inspecting the discrete frequencies, solutions of the dispersion relations. A super-convergent result is also demonstrated for the discrete frequencies of the NC method compared to the DG ones, except for the slow mode in the Roe flux case. Indeed, the Roe flux yields spurious frequencies in $O(h)$, and hence sub-optimal rates of convergence are expected for the slow mode for both the DG and NC methods. The numerical simulation of pure gravity waves and an anticyclonic eddy at midlatitudes, using linear and non-linear model problems, confirm the theoretical results. Further, the computational cost efficiency of the NC approach is enhanced due to the fact that the CFL number of the NC method is $5/3$ larger than that of the DG method. These encouraging results suggest undertaking further experiments with a realistic bathymetry and wind forcing as a further step toward the possible construction of an ocean model based on the NC method described herein.

References

- [1] D.N. Arnold, F. Brezzi, M. Fortin, A stable finite element for the Stokes equations, *Calcolo* 21 (1984) 337–344.
- [2] D.N. Arnold, R.S. Falk, R. Winther, Finite element exterior calculus, homological techniques, and applications, *Acta Numerica* 15 (2006) 1–155.
- [3] D. Boffi, F. Brezzi, M. Fortin, *Mixed Finite Element Methods and Applications*, Springer Series in Computational Mathematics 44, Springer-Verlag Berlin-Heidelberg, 2013.
- [4] M.J. Castro-Díaz, E.D. Fernández-Nieto, A class of computationally fast first order finite volume solvers: PVM methods, *SIAM J. Sci. Comput.* 34 (2012) A2173–A2196.
- [5] R. Comblen, J. Lambrechts, J.F. Remacle, V. Legat, Practical evaluation of five partly discontinuous finite element pairs for the non-conservative shallow water equations, *Int. J. Numer. Meth. Fluids* 63 (2010) 701–724.
- [6] C.J. Cotter, D.A. Ham, Numerical wave propagation for the triangular P1DG-P2 finite element pair, *J. Comput. Phys.* 230 (2011) 2806–2820.
- [7] C.J. Cotter, J. Shipton, Mixed finite elements for numerical weather prediction, *J. Comput. Phys.* 231 (2012) 7076–7091.
- [8] C.J. Cotter, J. Thuburn, A finite element exterior calculus framework for the rotating shallow-water equations, *J. Comput. Phys.* 257 (2014) 1506–1526.
- [9] M. Crouzeix, P.A. Raviart, Conforming and non-conforming finite-element methods for solving the stationary Stokes equations, *RAIRO Anal. Numer.* 7 (1973) 33–76.
- [10] B. Cushman-Roisin, J.-M. Beckers, *Introduction to geophysical fluid dynamics: physical and numerical aspects*, Academic Press, 2011.
- [11] S. Danilov, On utility of triangular C-grid type discretization for numerical modeling of large-scale ocean flows, *Ocean Dynamics* 60 (2010) 1361–1369.
- [12] P. Degond, P.F. Peyrard, G. Russo, Ph. Villedieu, Polynomial upwind schemes for hyperbolic systems, *C. R. Acad. Sci. Paris* 328 (1999) 479–483.
- [13] C. Eldred, D.Y. Le Roux, Dispersion analysis of compatible Galerkin schemes for the 1D shallow water model, *J. Comput. Phys.* 371 (2018) 779–800.
- [14] C. Eldred, D.Y. Le Roux, Dispersion analysis of compatible Galerkin schemes on quadrilaterals for shallow water model, *J. Comput. Phys.* 387 (2019) 539–568.
- [15] P. Glaister, Approximate Riemann solutions of the shallow water equations, *J. Hydraul. Res.* 26 (1988) 293–306.
- [16] J.S. Hesthaven, T. Warburton, *Nodal Discontinuous Galerkin Methods: Algorithms, Analysis, and Applications*, Texts in Applied Mathematics 54, Springer, New York, 2008.
- [17] P. Hood, C. Taylor, Navier-Stokes equations using mixed interpolation, *Finite Elements in Flow Problems*, J. T. Oden et al., Ed., The University of Alabama in Huntsville (UAH) Press, Huntsville, AL 35899, U.S.A., 1974, pp. 121–132.
- [18] F.Q. Hu, M.Y. Hussaini, P. Rasetarinera, An analysis of the discontinuous Galerkin method for wave propagation problems, *J. Comput. Phys.* 151 (1999) 921–946.
- [19] B.L. Hua, F. Thomasset, A noise-free finite-element scheme for the two-layer shallow-water equations, *Tellus* 36 A (1984) 157–165.
- [20] T.J. Hughes, L.P. Franca, M. Balestra, A new finite element formulation for computational fluid dynamics: V. Circumventing the Babuska-Brezzi condition: A stable Petrov-Galerkin formulation of the Stokes problem accommodating equal-order interpolations, *Comput. Methods Appl. Mech. Eng.* 59 (1986) 85–99.
- [21] M. Iskandarani, D. Haidvogel, J. Boyd, A staggered spectral element model with application to the oceanic shallow water equations, *Int. J. Numer. Meth. Fluid.* 20 (1995) 393–414.
- [22] I.P.E. Kinnmark, *The shallow-water wave equations: Formulation, analysis and application*, Lecture Notes in Engineering 15, Springer-Verlag Berlin, Heidelberg, 1986.
- [23] P.H. LeBlond, L.A. Mysak, *Waves in the Ocean*, Elsevier, Amsterdam, 1978.
- [24] D.Y. Le Roux, Dispersion relation analysis of the $P_1^{NC} - P_1$ finite-element pair in shallow-water models, *SIAM J. Sci. Comput.* 27 (2005) 394–414.
- [25] D.Y. Le Roux, V. Rostand, B. Pouliot, Analysis of numerically induced oscillations in 2D finite-element shallow-water models, Part I: Inertia-gravity waves, *SIAM J. Sci. Comput.* 29 (2007) 331–360.
- [26] D.Y. Le Roux, Spurious inertial oscillations in shallow-water models, *J. Comput. Phys.* 231 (2012) 7959–7987.
- [27] D.Y. Le Roux, C. Eldred, M.A. Taylor, Fourier analyses of high order continuous and discontinuous Galerkin methods, *SIAM J. Numer. Anal.* 58 (2020) 1845–1866.
- [28] R.J. LeVeque, *Finite Volume Methods for Hyperbolic Problems*, Cambridge Texts in Applied Mathematics 31, 2002.
- [29] D.R. Lynch, W.G. Gray, A wave equation model for finite element tidal computations, *Comput. Fluids* 7 (1979) 207–228.
- [30] T. Melvin, A. Staniforth, J. Thuburn, Dispersion analysis of the spectral element method, *Q. J.R. Meteorol. Soc.* 138 (2012) 1934–1947.
- [31] R.C. Moura, S.J. Shervin, J. Peiró, Linear dispersion-diffusion analysis and its application to under-resolved turbulence simulations using discontinuous Galerkin spectral/hp methods, *J. Comput. Phys.* 298 (2015) 695–710.
- [32] J. Pedlosky, *Geophysical Fluid Dynamics*, Springer-Verlag, New York, 1987.
- [33] P. Raviart, J. Thomas, A mixed finite element method for 2nd order elliptic problems, in: Galligani, I., Magenes, E. (eds) *Mathematical Aspects of Finite Element Methods*, Lecture Notes in Mathematics 606, Springer, Berlin, Heidelberg, 1977.
- [34] P.L. Roe, Approximate Riemann solvers, Parameter Vectors, and Difference Schemes, *J. Comput. Phys.* 43 (1981) 357–372.
- [35] V. Rostand, D.Y. Le Roux, G. Carey, Kernel analysis of the discretized finite difference and finite element-shallow water models, *SIAM J. Sci. Comput.* 31 (2008) 531–556.

- [36] V. Rostand, D.Y. Le Roux, Raviart-Thomas and Brezzi-Douglas-Marini finite-element approximations of the shallow-water equations, *Int. J. Numer. Meth. Fluid.* 57 (2008) 951–976.
- [37] A. Staniforth, T. Melvin, C. Cotter, Analysis of a mixed finite-element pair proposed for an atmospheric dynamical core, *Q. J. R. Meteorol. Soc.* 139 (2013) 1239–1254.
- [38] E.F. Toro, Riemann problems and the WAF method for solving two-dimensional shallow-water equations, *Phil. Trans. Roy. Soc. London* 338 (1992) 43–68.
- [39] F. Thomasset, *Implementation of Finite Element Methods for Navier-Stokes Equations*, Springer-Verlag, Berlin, 1981.
- [40] P.A. Ullrich, D.R. Reynolds, J.E. Guerra, M.A. Taylor, Impact and importance of hyperdiffusion on the spectral element method: A linear dispersion analysis, *J. Comput. Phys.* 375 (2018) 427–446.
- [41] R.A. Walters, G.F. Carey, Analysis of spurious oscillation modes for the shallow-water and Navier-Stokes equations, *Comput. Fluids* 11 (1983) 51–68.
- [42] R.A. Walters, G.F. Carey, Numerical noise in ocean and estuarine models, *Adv. Water Res.* 7 (1984) 15–20.
- [43] P.J. Wolfram, O.B. Fringer, Mitigating horizontal divergence checker-board oscillations on unstructured triangular C-grids for nonlinear hydrostatic and nonhydrostatic flows, *Ocean Modelling*, 69 (2013) 64–78.
- [44] X. Zhong, C.-W. Shu, Numerical resolution of discontinuous Galerkin methods for time dependent wave equations, *Comput. Methods Appl. Mech. Eng.*, 200 (2011) 2814–2827.