



**HAL**  
open science

# Numerical methods around shallow water flows : dispersive effects, Coriolis force

Virgile Dubos

► **To cite this version:**

Virgile Dubos. Numerical methods around shallow water flows: dispersive effects, Coriolis force. Numerical Analysis [cs.NA]. Sorbonne Université, 2021. English. NNT: 2021SORUS260. tel-03557797v2

**HAL Id: tel-03557797**

**<https://theses.hal.science/tel-03557797v2>**

Submitted on 4 Feb 2022

**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.

SORBONNE UNIVERSITÉ  
LABORATOIRE JACQUES-LOUIS LIONS

INRIA  
ÉQUIPE-PROJET ANGE

ÉCOLE DOCTORALE DE SCIENCES MATHÉMATIQUES DE PARIS CENTRE

# THÈSE DE DOCTORAT

Discipline : Sciences Mathématiques

présentée par  
**Virgile DUBOS**

---

## Méthodes numériques autour d'écoulements en eaux peu profondes : effets dispersifs, force de Coriolis

---

sous la direction de Cindy GUICHARD, Yohan PENEL et Jacques SAINTE-MARIE

soutenue le 14 décembre 2021 devant un jury composé de :

Mme	Claire CHAINAIS-HILLAIRET	Professeure	Rapportrice
Mme	Cindy GUICHARD	Maîtresse de conférences	Encadrante
Mme	Raphaèle HERBIN	Professeure	Examinatrice
M.	Pierre-Yves LAGRÉE	Directeur de recherche	Examineur
M.	Pascal OMNES	Directeur de recherche	Président du jury
M.	Yohan PENEL	Chargé de recherche	Encadrant
M.	Jacques SAINTE-MARIE	Directeur de recherche	Directeur de thèse

après avis favorable des rapporteurs :

Mme	Claire CHAINAIS-HILLAIRET	Professeure
M.	Rainer HELMIG	Professeur



# Résumé

Nous étudions dans cette thèse des méthodes numériques pour les écoulements en eaux peu profondes à surface libre. D'une part, nous nous intéressons aux modèles non-hydrostatique de Saint-Venant prenant en compte des effets dispersifs. D'autre part, nous étudions le modèle de Saint-Venant avec terme source de Coriolis et ses équilibres géostrophiques.

Dans un premier temps nous considérons des méthodes numériques pour une famille de modèles d'Euler moyennés sur la hauteur d'eau provenant de la littérature. Il s'agit de modèles de Saint-Venant avec une pression non-hydrostatique. Les méthodes numériques considérées sont basées sur une méthode de prédiction-correction consistant à décomposer le problème en deux étapes à chaque itération en temps. L'étape de prédiction mène à résoudre un système de Saint-Venant qui est habituellement résolu par une méthode de volumes finis, alors que l'étape de correction mène à résoudre un problème elliptique régissant la pression non-hydrostatique. Notre but premier est d'analyser la convergence d'un schéma mixte avec condensation de masse introduit dans la littérature en utilisant la Méthode de Discrétisation du Gradient (GDM) qui forme un cadre englobant des schémas de discrétisation usuels et récents pour des problèmes de diffusion. Cette méthode nous permet d'obtenir la convergence du schéma proposé pour le problème elliptique. Par la suite, nous proposons une nouvelle formulation conforme du problème en pression et donnons les estimateurs d'erreur correspondant via la GDM. Enfin, nous donnons un exemple d'application avec la méthode des éléments finis  $\mathbb{P}_1$ .

Dans un second temps nous souhaitons avoir des schémas explicites colocalisés volumes finis pour les équations de Saint-Venant non-linéaires avec un terme source de Coriolis qui seraient précis aux alentours de l'équilibre géostrophique et stables dans le cadre non-linéaire. Nous construisons plusieurs schémas volumes finis et étudions les deux propriétés suivantes : la décroissance de l'énergie semi-discrète et la préservation de l'équilibre géostrophique par la version linéarisée. Nous proposons également une version conservative de l'un de nos schéma. Enfin, nous observons le comportement de ces schémas à travers plusieurs cas tests et obtenons de meilleurs résultats en comparaison d'un schéma volumes finis classique.

**Mots-clés** : écoulements à surface libre, équations de Saint-Venant, pression non-hydrostatique, méthode de discrétisation du gradient, éléments finis, force de Coriolis, équilibre géostrophique, schéma équilibre, diffusion numérique, volumes finis.



# Abstract

In this work we study some numerical methods for free surface shallow water flows. On one hand, we are interested in non-hydrostatic shallow water models which take into account dispersive effects. On the other hand, we study the shallow water model with Coriolis source term and the associated geostrophic equilibrium.

First, we consider numerical methods for a family of depth-averaged Euler models from the literature. These models are shallow water models with non-hydrostatic pressure. The considered numerical methods are based on a prediction-correction method which leads to decompose the problem into two parts at each time step. The prediction part consists in solving a shallow water system which is usually solved by finite volume methods, while the correction part consists in solving an elliptic problem governing the non-hydrostatic pressure. Our first goal is to analyse the convergence of a mixed mass-lumped scheme introduced in the literature through the Gradient Discretisation Method (GDM) which is a framework comprising classic and recent discretisation schemes for diffusion problems. This method allows us to obtain a convergence result on the proposed scheme for the elliptic problem. Then, we propose a new conforming formulation of the problem on the pressure and give the corresponding error estimate provided by the GDM. Finally, an application to the  $\mathbb{P}_1$  conforming finite element method is given.

Then, we aim at designing explicit collocated finite volume schemes for the nonlinear Shallow Water equations with Coriolis source term that are proved to be accurate around the geostrophic equilibrium and stable in the nonlinear framework. We define some finite volume schemes and study the two properties we are interested in : the decrease of the semi-discrete energy and the preservation of the geostrophic equilibrium in the linearized version. We also propose a conservative version of one of our scheme. Finally, we illustrate the behaviour of the schemes for some standard test cases and we exhibit a great improvement when compared to a classic finite volume scheme.

**Keywords** : free surface flows, shallow water equations, non-hydrostatic pressure, Gradient Discretisation Method, finite element, Coriolis force, geostrophic equilibrium, well-balanced scheme, numerical diffusion, finite volume.



# Table des matières

<b>1</b>	<b>Introduction</b>	<b>11</b>
1.1	Équations d'Euler et de Saint-Venant . . . . .	11
1.2	Une famille de modèles non-hydrostatiques . . . . .	13
1.3	Modèle de Saint-Venant avec force de Coriolis . . . . .	16
1.4	Organisation du manuscrit . . . . .	18
<b>2</b>	<b>Gradient Discretisation Methods to analyse numerical schemes for the elliptic part of a dispersive Shallow Water system.</b>	<b>21</b>
2.1	Introduction . . . . .	22
2.2	The Gradient Discretisation Method on <i>classic</i> operators . . . . .	24
2.2.1	A conforming formulation on the pressure . . . . .	25
2.2.2	Weak formulation . . . . .	25
2.2.3	Gradient Scheme . . . . .	27
2.2.4	Key properties of the considered Gradient Discretisation . . . . .	28
2.2.5	Error estimates . . . . .	29
2.2.6	The $\mathbb{P}_1$ conforming finite elements . . . . .	33
2.3	The Abstract Gradient Discretisation Method on <i>special</i> operators . . . . .	34
2.3.1	Weak formulation . . . . .	34
2.3.2	Abstract Gradient Scheme . . . . .	35
2.3.3	Key properties of the considered Abstract Gradient Discretisation . . . . .	36
2.3.4	Error estimates . . . . .	37
2.3.5	The $\mathbb{P}_1$ conforming finite elements . . . . .	39
2.4	Analyse of a mixed mass-lumped scheme . . . . .	40
2.4.1	A weak mixed formulation with lifting . . . . .	40
2.4.2	The associated Abstract Gradient Discretisation Method framework . . . . .	41
2.4.3	Mixed mass-lumping finite element methods & AGD . . . . .	43
2.4.4	The $\mathbb{P}_1/\mathbb{P}_1$ approximation . . . . .	47
2.5	Conclusion . . . . .	56
<b>3</b>	<b>Energy stable and linearly well-balanced numerical schemes for the non-linear Shallow Water equations with Coriolis force</b>	<b>57</b>
3.1	Introduction . . . . .	58
3.2	Shallow water equations and geostrophic equilibrium . . . . .	58
3.3	Discrete operators . . . . .	60
3.3.1	Definition of the mesh . . . . .	60
3.3.2	Discrete operators . . . . .	61
3.3.3	Mimetic properties . . . . .	63



3.4	Well-balanced and stable finite volume schemes . . . . .	64
3.4.1	Edge based 9 points schemes . . . . .	64
3.4.2	Cell-based 13 points scheme . . . . .	68
3.4.3	Vertex-based 25 points scheme . . . . .	70
3.4.4	Time discretisation . . . . .	75
3.5	Numerical results . . . . .	75
3.5.1	River test case . . . . .	75
3.5.2	Stationary vortex test case . . . . .	77
3.5.3	Translated vortex test case . . . . .	80
3.5.4	Water-column test case . . . . .	83
3.6	Conclusion . . . . .	85
<b>A</b>	<b>Numerical approximation of the Shallow Water equations with Coriolis source term</b>	<b>87</b>
A.1	Design of the numerical schemes . . . . .	89
A.1.1	Collocated scheme . . . . .	89
A.1.2	Staggered scheme . . . . .	95
A.1.3	Time discretisation . . . . .	98
A.2	Numerical assessments of the schemes . . . . .	98
A.2.1	Presentation of the test cases . . . . .	98
A.2.2	Numerical results . . . . .	100
A.3	Conclusion . . . . .	100
	<b>Bibliographie</b>	<b>101</b>

# Table des figures

1.1	Domaine et notations pour le système de Saint-Venant monodimensionnel . . . . .	13
1.2	Cas test de Dingemans . . . . .	14
1.3	Comparaison entre modèles SW hydrostatique et non-hydrostatique . . . . .	14
2.1	Notations for the one-dimensional Shallow Water system. Source : [1] . . . . .	23
3.1	Geometric notations . . . . .	60
3.2	Physical and phantom domains . . . . .	61
3.3	Initial water depth. . . . .	76
3.4	Cross section in $y = 0$ of solution at $t = 20s$ difference to initial state. . . . .	76
3.5	Energy of the system in function of time. . . . .	77
3.6	Initial state of the stationary vortex. . . . .	78
3.7	Water depth at $t = 200s$ . . . . .	78
3.8	Velocity $v$ at $t = 200s$ . . . . .	79
3.9	Cross sections $y = 0$ of solution at $t = 200s$ . . . . .	79
3.10	Water depth relative error to initial state. . . . .	79
3.11	Energy of the system in function of time for $\epsilon = 0.01$ . . . . .	80
3.12	Initial state of the translated vortex. . . . .	80
3.13	Water depth at $t = 20s$ . . . . .	81
3.14	Velocity $v$ at $t = 20s$ . . . . .	81
3.15	Water depth relative error to initial state in function of time . . . . .	81
3.16	Cross section of water depth at $t = 20s$ in $x = 0$ (a) and $y = 0$ (b). . . . .	82
3.17	Cross section of velocity $v$ at $t = 20s$ in $x = 0$ (a) and $y = 0$ (b). . . . .	82
3.18	Energy of the system in function of time for the different schemes. . . . .	82
3.19	Water depth at $t = 100s$ . . . . .	83
3.20	Velocity $v$ at $t = 100s$ . . . . .	83
3.21	Cross section in $y = 0$ of simulation results for the different schemes at different times. . . . .	84
3.22	Energy of the system in function of time for the different schemes. . . . .	85
A.1	Geometric settings . . . . .	90
A.2	Geometric settings : for the primal mesh with cells $K$ and $K_e$ , refer to Figs. A.1 – for the dual mesh with diamond cells $D$ and $D_f$ separated by the interface $f = \partial D \cap \partial D_f$ , $\mathbf{n}_{f,D}$ is the outward normal of $f$ pointing to $D_f$ . . . . .	95
A.3	Stationary vortex : cross-section of the water height at final time for different values of $\epsilon$ ; comparisons between the exact solution ( <i>green curves</i> ), the numerical solution obtained with the classic Godunov scheme ( <i>orange curve</i> ), the collocated scheme (A.12) (A-B-C) and the staggered scheme (A.25) (D-E-F) ( <i>blue curves</i> ) . . . . .	99
A.4	Error in $L^2$ norm between the numerical solution and the exact solution at final time . . . . .	99



# Chapitre 1

## Introduction

Une branche de la mécanique des fluides consiste à étudier les écoulements à surface libre, typiquement constitués d'eau plus ou moins pure, et considérés incompressibles. On entend par "surface libre", un fluide dont la surface serait contrainte uniquement par la pression atmosphérique. Historiquement, on sépare ces écoulements en deux grandes catégories : d'une part, les océans/mers qui sont considérés comme presque plats et pour lesquels la description verticale du fluide est la priorité ; d'autre part, les rivières pour lesquelles on s'intéresse le plus souvent à la hauteur d'eau et au débit. Cette distinction s'explique mathématiquement par le fait que dans le premier cas les solutions sont régulières sur de très grands domaines, alors que dans le second, il est possible d'observer des chocs *i.e.* des solutions discontinues sur des domaines réduits. On utilise donc très souvent des modèles différents en fonction de la situation rencontrée. En océanographie, on utilise des modèles de type Navier-Stokes (en pratique on néglige les effets visqueux ce qui permet d'utiliser une version simplifiée *i.e.* les équations d'Euler), alors que pour les études fluviales on préfère utiliser des modèles de type Saint-Venant. La modélisation, l'analyse et la simulation d'écoulements à surface libre sont des sujets complexes. En effet, les équations rencontrées sont à la fois difficiles à analyser et à résoudre avec de fortes non-linéarités (e.g. équations de Navier-Stokes), les échelles de temps et d'espace sont souvent très grandes et diverses (e.g. érosion des côtes), enfin, les phénomènes à modéliser sont multi-physiques et nécessitent de nombreux couplages (e.g. salinité, température).

### 1.1 Équations d'Euler et de Saint-Venant

Les équations d'Euler ont été introduites en 1757 par Leonhard Euler [28] dans le but de décrire des écoulements incompressibles sans viscosité et à densité uniforme égale à 1. Elles s'obtiennent à partir des équations de Navier-Stokes en négligeant les termes de friction et de viscosité. Le système d'équations a la forme suivante :

$$\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} = 0, \quad (1.1a)$$

$$\frac{\partial u}{\partial t} + \frac{\partial u^2}{\partial x} + \frac{\partial uv}{\partial y} + \frac{\partial uw}{\partial z} + \frac{\partial p}{\partial x} = 0, \quad (1.1b)$$

$$\frac{\partial v}{\partial t} + \frac{\partial uv}{\partial x} + \frac{\partial v^2}{\partial y} + \frac{\partial vw}{\partial z} + \frac{\partial p}{\partial y} = 0, \quad (1.1c)$$

$$\frac{\partial w}{\partial t} + \frac{\partial uw}{\partial x} + \frac{\partial vw}{\partial y} + \frac{\partial w^2}{\partial z} + \frac{\partial p}{\partial z} = g, \quad (1.1d)$$

où  $(u, v, w)$  est la vitesse,  $p$  est la pression et  $g$  est la constante de gravitation.

L'équation (1.1a) est l'équation de conservation de la masse, les équations (1.1b)-(1.1d) sont les équations de conservation de la quantité de mouvement. Numériquement ces équations sont très coûteuses à résoudre, en particulier dans le cadre de l'océanographie où les échelles de temps et d'espace sont très grandes et où les écoulements sont à surface libre.

Pour contourner cet obstacle, il est possible de réduire le système d'Euler à un modèle plus simple, celui de Saint-Venant. En effet, ce dernier est une très bonne approximation du système d'Euler pour de nombreux problèmes (rupture de barrage, inondations, charriage). Les équations dites "de Barré de Saint-Venant", publiées en 1871, décrivent les écoulements à surface libre en eaux peu profondes, d'où leur appellation anglaise *Shallow Water Equations* (SW) [46]. On les rencontre en géophysique par exemple pour décrire les écoulements de rivières, la circulation de courants océaniques, la propagation de tsunamis ou encore pour modéliser d'autres phénomènes comme les avalanches, le trafic routier ou la circulation sanguine. Du fait de sa validité expérimentale et de son efficacité numérique, le modèle SW est très utilisé pour la simulation de nombreux phénomènes d'actualité : production d'énergies renouvelables, prédiction de catastrophes naturelles, étude du climat. En effet, les deux points forts de ce modèle sont sa validité expérimentale, dépassant largement sa validité théorique, ainsi que son excellent rapport précision / coût numérique. Nous rappelons ici brièvement son obtention à partir des équations d'Euler (1.1).

Sous l'hypothèse d'eaux peu profondes et en négligeant la pression à la surface du fluide, les équations d'Euler adimensionnées mènent à une expression de la pression  $p = g(\eta - z)$  dite hydrostatique car ne dépendant que du poids du fluide à l'équilibre. Dans ce cas, on dit que la pression non-hydrostatique est négligée et le modèle SW est alors hydrostatique. Pour fermer le système, on fait de plus l'hypothèse que la vitesse est homogène selon la verticale. Ainsi en moyennant selon la verticale les équations (1.1) on obtient le modèle SW suivant :

$$\frac{\partial H}{\partial t} + \frac{\partial H\bar{u}}{\partial x} + \frac{\partial H\bar{v}}{\partial y} = 0, \quad (1.2a)$$

$$\frac{\partial H\bar{u}}{\partial t} + \frac{\partial H\bar{u}^2}{\partial x} + \frac{\partial H\bar{u}\bar{v}}{\partial y} + g\frac{\partial}{\partial x} \left( \frac{H^2}{2} \right) = -gH\frac{\partial z_b}{\partial x}, \quad (1.2b)$$

$$\frac{\partial H\bar{v}}{\partial t} + \frac{\partial H\bar{u}\bar{v}}{\partial x} + \frac{\partial H\bar{v}^2}{\partial y} + g\frac{\partial}{\partial y} \left( \frac{H^2}{2} \right) = -gH\frac{\partial z_b}{\partial y}, \quad (1.2c)$$

où  $(\bar{u}, \bar{v})$  est la vitesse moyennée sur la verticale,  $H$  est la hauteur d'eau et  $z_b$  est la cote du fond du domaine, appelée bathymétrie. En introduisant  $\eta$  la cote de la surface libre on peut définir la vitesse moyenne par :

$$\bar{u} = \frac{1}{H} \int_{z_b}^{\eta} u dz, \quad \bar{v} = \frac{1}{H} \int_{z_b}^{\eta} v dz, \quad (1.3)$$

et la hauteur d'eau par :

$$H = \eta - z_b. \quad (1.4)$$

Ces différentes quantités sont illustrées dans la Figure 1.1 suivante pour le modèle SW monodimensionnel.

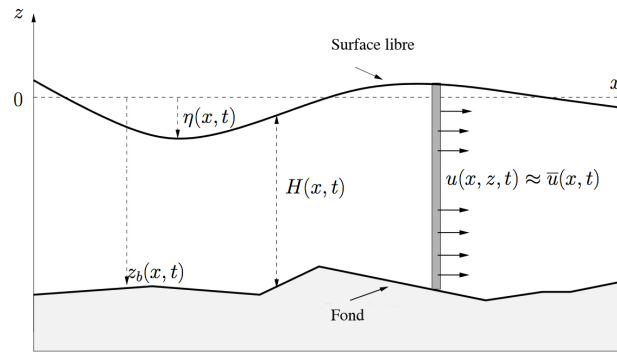


FIGURE 1.1 – Domaine et notations pour le système de Saint-Venant monodimensionnel

Par construction, la vitesse verticale  $w$  est négligeable dans les équations de Saint-Venant. Ainsi le système d'équations de Saint-Venant n'aura pas d'équation selon la direction  $z$ . La vitesse verticale  $w$  pourra néanmoins être reconstruite à partir de la solution du problème si besoin. Ce modèle permet donc d'aborder un problème physique tridimensionnel (resp. bidimensionnel) et instationnaire au travers de l'étude d'un système posé sur un domaine bidimensionnel (resp. monodimensionnel), réduisant drastiquement le coût numérique.

Cependant, tous les phénomènes physiques étudiés ne peuvent pas être modélisés de cette manière et ce modèle hydrostatique échoue à les capturer. Pour remédier à ce problème, il est possible de dériver des modèles SW non-hydrostatiques au prix d'un coût numérique et d'une complexité mathématique plus importants.

## 1.2 Une famille de modèles non-hydrostatiques

Le modèle SW se déduit donc des équations d'Euler via deux hypothèses majeures : une pression hydrostatique et une vitesse constante selon la verticale. On peut cependant se demander quel est le domaine de validité d'une telle approche ? Dans quel cas le modèle n'est pas suffisamment précis et doit être remplacé ? Dans la suite, on s'intéresse en particulier au cas où la pression ne peut pas être considérée hydrostatique.

Dans le cas où la pression non-hydrostatique ne peut plus être négligée, le modèle SW échoue à simuler précisément les écoulements (cf. Figure 1.3). Cette limitation a mené à l'étude de modèles dits non-hydrostatiques dans lesquels l'hypothèse de pression hydrostatique n'est pas faite. On illustre ci-dessous la différence entre un modèle SW hydrostatique et un autre non-hydrostatique. Le cas test choisi (cf. Figure 1.2) est issu du livre [21], il s'agit d'un canal, dans lequel un obstacle est présent. On génère une vague de petite amplitude sur le bord gauche du domaine et on mesure expérimentalement la hauteur d'eau en aval de l'obstacle où une condition de type libre est imposée. Puis on compare ces données aux simulations numériques pour les deux modèles. On observe le modèle hydrostatique échouer à capturer précisément les oscillations, contrairement au modèle non-hydrostatique. En effet, la présence de l'obstacle augmente la part hydrodynamique de la pression dans le fluide, causant ainsi l'apparition d'effets dispersifs non modélisés par le modèle classique SW.

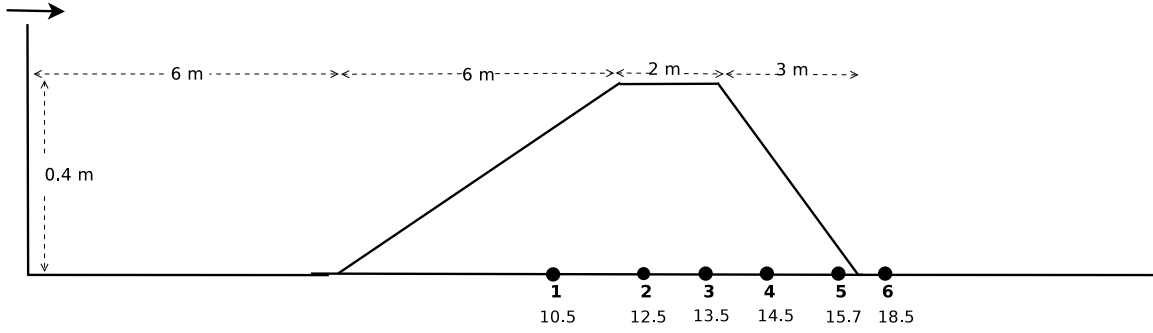


FIGURE 1.2 – Cas test de Dingemans

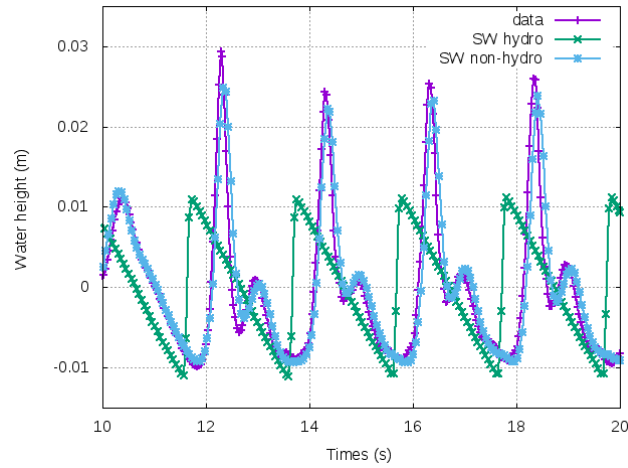


FIGURE 1.3 – Comparaison entre modèles SW hydrostatique et non-hydrostatique

Il existe de nombreux modèles non-hydrostatiques capables de rendre compte de ces effets dispersifs, chacun d'entre eux admet son propre domaine de validité. Nous nous intéressons particulièrement au modèle *depth-averaged Euler* (DAE) introduit dans [15] et étendu dans [2] en une famille de modèles dispersifs. Cette famille est issue du système d'Euler incompressible à surface libre par une approximation de type eaux peu profondes mais sans l'hypothèse hydrostatique. Le modèle est décrit par le système suivant :

$$\frac{\partial H}{\partial t} + \frac{\partial Hu}{\partial x} + \frac{\partial Hv}{\partial y} = 0, \quad (1.5a)$$

$$\frac{\partial Hu}{\partial t} + \frac{\partial}{\partial x} \left( Hu^2 + \frac{g}{2} H^2 + Hp \right) + \frac{\partial Huv}{\partial y} = - \left( gH + \frac{\alpha^2}{2} p \right) \frac{\partial z_b}{\partial x}, \quad (1.5b)$$

$$\frac{\partial Hv}{\partial t} + \frac{\partial Huv}{\partial x} + \frac{\partial}{\partial y} \left( Hv^2 + \frac{g}{2} H^2 + Hp \right) = - \left( gH + \frac{\alpha^2}{2} p \right) \frac{\partial z_b}{\partial y}, \quad (1.5c)$$

$$\frac{\partial Hw}{\partial t} + \frac{\partial Huw}{\partial x} + \frac{\partial Hvw}{\partial y} = \alpha p, \quad (1.5d)$$

$$\alpha w = -H \frac{\partial u}{\partial x} + \frac{\alpha^2}{2} u \frac{\partial z_b}{\partial x} - H \frac{\partial v}{\partial y} + \frac{\alpha^2}{2} v \frac{\partial z_b}{\partial y}, \quad (1.5e)$$

où  $(u, v, w)$  est le vecteur vitesse,  $p$  est la pression non-hydrostatique (la pression totale est donnée par

$p + gH/2$ ),  $g$  est la constante de gravitation et  $\alpha$  un paramètre réel.

Ce système apparaît comme une famille de modèles SW complétés par des termes dispersifs, dépendant du paramètre  $\alpha$ . Il peut être réécrit sous la forme condensée suivante :

$$\frac{\partial H}{\partial t} + \nabla_0 \cdot (H\mathbf{u}) = 0, \quad (1.6a)$$

$$\frac{\partial H\mathbf{u}}{\partial t} + \nabla_0 \cdot (H\mathbf{u} \otimes \mathbf{u}) + \nabla_0 \left( \frac{g}{2} H^2 \right) + \nabla_s^\alpha p = -gH\nabla_0 z_b, \quad (1.6b)$$

$$\operatorname{div}_s^\alpha \mathbf{u} = 0, \quad (1.6c)$$

où  $\nabla_0 = \left( \frac{\partial}{\partial x}, \frac{\partial}{\partial y}, 0 \right)^T$  et où les opérateurs  $\nabla_s^\alpha$  et  $\operatorname{div}_s^\alpha$  sont définis par :

$$\nabla_s^\alpha \varphi = (H\nabla\varphi + \varphi\nabla\zeta, -\alpha\varphi)^T, \quad (1.7a)$$

$$\operatorname{div}_s^\alpha \boldsymbol{\omega} = \operatorname{div}(H\mathbf{v}) - \mathbf{v} \cdot \nabla\zeta + \alpha w, \quad (1.7b)$$

pour  $\boldsymbol{\omega} = (\mathbf{v}, w)^T$  et  $\zeta = H + \frac{\alpha^2}{2} z_b$ .

En changeant la valeur du paramètre  $\alpha$  on retrouve deux modèles dispersifs présent dans la littérature. En effet,  $\alpha = 2$  permet d'obtenir le modèle DAE présenté dans [15], alors que  $\alpha = \sqrt{3}$  permet d'écrire le modèle de Green-Naghdi [35] à un petit terme d'erreur près pour la propagation des ondes (voir [2]).

Dans [2] les auteurs proposent une stratégie numérique pour le modèle (1.6) basée sur la décomposition de Chorin-Teman appliquée au système d'Euler incompressible [44]. Par un *splitting* en temps, les auteurs obtiennent un schéma de prédiction-correction qui revient à résoudre consécutivement les deux schémas semi-discrets en temps suivants. Pour la discrétisation en temps, on introduit  $t^{n+1} = t^n + \Delta t^n$ , où le pas de temps  $\Delta t^n$  est déterminé par une condition de type CFL. L'état  $X^k$  est l'approximation de  $X(t^k)$  pour  $X \in \{H, \mathbf{u}, p\}$  et pour  $k \in \{n, n + \frac{1}{2}, n + 1\}$ , où  $t^{n+\frac{1}{2}}$  est un temps intermédiaire entre  $t^n$  et  $t^{n+1}$ .

On résout l'étape de prédiction :

$$\frac{H^{n+\frac{1}{2}} - H^n}{\Delta t^n} + \nabla_0 \cdot (H\mathbf{u})^n = 0, \quad (1.8a)$$

$$\frac{(H\mathbf{u})^{n+\frac{1}{2}} - (H\mathbf{u})^n}{\Delta t^n} + \nabla_0 \cdot (H\mathbf{u} \otimes \mathbf{u})^n + \nabla_0 \left( \frac{g}{2} (H^n)^2 \right) = -gH^n \nabla_0 z_b. \quad (1.8b)$$

Puis, l'étape de correction :

$$H^{n+1} - H^{n+\frac{1}{2}} = 0, \quad (1.9a)$$

$$\frac{(H\mathbf{u})^{n+1} - (H\mathbf{u})^{n+\frac{1}{2}}}{\Delta t^n} + \nabla_s^\alpha p^{n+1} = 0, \quad (1.9b)$$

$$\operatorname{div}_s^\alpha \mathbf{u}^{n+1} = 0. \quad (1.9c)$$



Le problème (1.8) est donc le système SW (1.2) pouvant être résolu par les solveurs usuels (typiquement un schéma volumes finis). Le problème (1.9) est quant à lui un problème elliptique que les auteurs de [2] choisissent de résoudre par un schéma éléments finis mixtes. Ce schéma sera détaillé et étudié dans le Chapitre 2 de cette thèse.

### 1.3 Modèle de Saint-Venant avec force de Coriolis

Outre cette question sur la représentation de la pression et de la vitesse verticale à l'intérieur du fluide, le modèle SW permet la modélisation d'une grande variété de phénomènes physiques par l'ajout de termes sources correspondants. Le vent, la pluie, la topographie du fond sont autant de phénomènes pouvant être pris en compte dans le modèle SW par ce processus. Parmi ces phénomènes, on s'intéresse ici à la force de Coriolis qui résulte de la rotation de la Terre. Physiquement, son influence est importante lorsque de très grandes étendues de fluides sont considérées, typiquement de grands lacs, des mers ou des océans dans le cadre étudié ici. En océanographie, par exemple, il est essentiel de prendre en compte cette force afin d'étudier précisément la circulation des courants marins ayant un impact sur le climat. La force de Coriolis est une pseudo-force agissant perpendiculairement à la direction du mouvement d'un corps en déplacement dans un référentiel lui-même en rotation uniforme. Il ne s'agit pas d'une "force" au sens strict, *i.e.* l'action d'un corps sur un autre, mais plutôt d'une force fictive résultant du mouvement non linéaire du référentiel lui-même.

Classiquement, la force de Coriolis est représentée, elle aussi, sous la forme d'un terme source à l'équation de conservation de la quantité de mouvement du modèle SW (1.2) donné par l'expression suivante :

$$\omega H \mathbf{u}^\perp, \quad (1.10)$$

où  $\omega$  est la vitesse angulaire de rotation du référentiel (ici la Terre) supposée constante,  $\mathbf{u} = (u_x, u_y)$  est la vitesse horizontale et où  $\mathbf{u}^\perp = (-u_y, u_x)$  est le vecteur orthogonal correspondant.

Considérons dans la suite un fond plat dans un souci de simplicité, le terme de topographie disparaît et le système de Saint-Venant avec force de Coriolis (SWC) s'écrit alors de manière condensée :

$$\partial_t H + \operatorname{div}(H \mathbf{u}) = 0, \quad (1.11a)$$

$$\partial_t(H \mathbf{u}) + \operatorname{div}(H \mathbf{u} \otimes \mathbf{u}) + H(g \nabla H + \omega \mathbf{u}^\perp) = 0, \quad (1.11b)$$

où  $g$  est la constante de gravitation et où  $H$  est la hauteur d'eau.

On s'intéresse maintenant aux méthodes numériques pour résoudre ce système, notre priorité est d'en assurer la stabilité et la précision.

**États stationnaires :** Lorsque l'on considère le système SW (1.2), *i.e.* sans le terme de Coriolis, il existe des états stationnaires non triviaux, pour lesquels les inconnues ne sont pas constantes sur le domaine. Du fait de la présence du terme source topographique, le système SW présente en effet la particularité de posséder des états stationnaires complexes. En dimension un, et pour des solutions régulières, ils sont caractérisés par les relations suivantes :

$$\frac{\partial h u}{\partial x} = 0, \quad \text{et} \quad \frac{\partial}{\partial x} \left( \frac{u^2}{2} + g(h + z_b) \right) = 0.$$

Un équilibre très important à capturer dans ce cas est "le lac au repos", défini par  $u = 0$  et  $h + z_b = cst.$ , car de nombreuses observations expérimentales sont en réalité de petites perturbations

autour de cet équilibre. Ainsi la précision de la méthode numérique choisie est directement liée à sa capacité à préserver cet équilibre. Il est devenu usuel de construire des schémas dits équilibrés (ou *well-balanced*) qui préservent le lac au repos. On citera par exemple, la reconstruction hydrostatique introduite dans [4], basée sur une discrétisation du terme source de topographie à partir d'une reconstruction de la hauteur d'eau aux interfaces des cellules. Les flux aux interfaces sont alors construits à partir de cette hauteur reconstruite. Cette technique a l'avantage d'être indépendante du solveur homogène choisi, mais adapte ce dernier au cas non-homogène en assurant la préservation de certains états stationnaires du système continu, en particulier de l'équilibre associé au lac au repos, tout en conservant intactes les propriétés de stabilité obtenues au préalable dans le cas homogène.

Cependant la présence du terme de Coriolis dans le système (1.11) amène à considérer le système linéarisé autour du "lac au repos"  $(\tilde{H}, \tilde{\mathbf{u}}) = (H_0, 0)$  avec  $H_0$  constant :

$$\partial_t H = -H_0 \operatorname{div} \mathbf{u} , \quad (1.12a)$$

$$\partial_t \mathbf{u} = -(g \nabla H + \omega \mathbf{u}^\perp) . \quad (1.12b)$$

Ce système linéarisé a pour état stationnaire, l'équilibre géostrophique suivant :

$$g \nabla H + \omega \mathbf{u}^\perp = 0. \quad (1.13)$$

En effet, cet équilibre implique la contrainte de divergence libre :

$$\operatorname{div} \mathbf{u} = 0. \quad (1.14)$$

On remarque que l'équilibre géostrophique (1.13) n'est pas un équilibre du système SWC (1.11), mais en pratique il est nécessaire que nos schémas linéarisés puissent le préserver dans un souci de stabilité. Dans [5], les auteurs ont montré que le schéma de Godunov volumes finis colocalisé classique ne préserve pas cet équilibre à cause des termes de diffusion numériques. En effet, il est nécessaire de construire des schémas dédiés dont les termes de diffusion numériques seront nuls aux alentours d'un équilibre géostrophique discret. On peut alors citer les auteurs de [14] qui ont adapté la reconstruction hydrostatique au cas avec terme source de Coriolis en introduisant artificiellement une topographie, menant alors à la méthode dite de topographie apparente (*Apparent Topography*).

**Stabilité entropique :** Concernant la stabilité, on remarque que le système (1.11) admet une énergie  $E$  qui vérifie (pour les solutions régulières) l'équation de conservation suivante :

$$\partial_t E + \operatorname{div} \left[ \left( gh + \frac{1}{2} \|\mathbf{u}\|^2 \right) h \mathbf{u} \right] = 0, \quad (1.15)$$

où  $E = \frac{1}{2} gh^2 + \frac{1}{2} h \|\mathbf{u}\|^2$ .

Ainsi, à l'échelle du domaine, cette énergie reste constante au cours du temps. Pour les solutions discontinues, on obtient seulement un résultat de non croissance de l'énergie. Dans tous les cas,  $E$  correspond à une entropie mathématique, en construisant une approximation numérique qui vérifie une version discrète de (1.15) on obtient alors une stabilité entropique. Les méthodes de viscosité introduites dans le cadre linéaire [5] et non-linéaire [20] fournissent un outil puissant pour assurer la stabilité des schémas au sens de l'énergie. Dans le contexte des équations de Saint-Venant non-linéaires, très peu de schémas type Volumes-Finis sont capables de fournir des inégalités d'entropie discrètes en temps et en espace, et ce même en l'absence de topographie. Dans des travaux récents [10], ces termes

de stabilisation ont été utilisés pour garantir des inégalités d'entropie discrètes pour le schéma de reconstruction hydrostatique [4], introduit pour assurer les équilibres du type lac au repos pour les équations de Saint-Venant (schémas *well-balanced*).

Dans le Chapitre 3 de cette thèse, nous proposons plusieurs schémas vérifiant ces propriétés de stabilité entropique et/ou de préservation d'équilibre géostrophique.

## 1.4 Organisation du manuscrit

Cette thèse se compose de deux chapitres indépendants rédigés en anglais dans le cadre de publications. Dans le Chapitre 2, réalisé en collaboration avec Cindy Guichard et Yohan Penel, nous nous intéressons au problème issu de la partie correction (1.9) du schéma de *splitting* pour le modèle DAE [2].

Soit  $\Omega \subset \mathbb{R}^d$ , où  $d = 1$  ou  $d = 2$ . Le problème revient à trouver  $p^\dagger : \Omega \rightarrow \mathbb{R}$  la pression non-hydrostatique et  $\mathbf{u}^\dagger = (\mathbf{v}^\dagger, w^\dagger)^T : \Omega \rightarrow \mathbb{R}^{d+1}$  le champ de vitesse tels que :

$$H\mathbf{u}^\dagger + \nabla_s^\alpha p^\dagger = \mathbf{g} \text{ sur } \Omega, \quad (1.16a)$$

$$\operatorname{div}_s^\alpha \mathbf{u}^\dagger = f \text{ sur } \Omega, \quad (1.16b)$$

où  $H$  et  $\mathbf{g} = (\mathbf{g}_1, g_2)$  sont donnés par l'étape de prédiction (1.8) et où  $f$  est la fonction nulle sur  $\Omega$  dans le modèle DAE mais qui sera considérée quelconque pour l'analyse du problème ici.

En introduisant  $\Gamma_n$  et  $\Gamma_d$  tels que  $\partial\Omega = \Gamma_n \cup \Gamma_d$ , nous complétons le problème par les conditions de bord suivantes :

$$H\mathbf{u}^\dagger \cdot \mathbf{n}_s = H\mathbf{v}^\dagger \cdot \mathbf{n}_\Gamma = \phi \text{ sur } \Gamma_n, \quad (1.17a)$$

$$p^\dagger = 0 \text{ sur } \Gamma_d, \quad (1.17b)$$

où  $\mathbf{n}_\Gamma$  est le vecteur normal unitaire sortant de  $\Omega$ ,  $\phi$  un flux imposé et où nous définissons  $\mathbf{n}_s := (\mathbf{n}_\Gamma, 0)^T$ .

Nous étudions deux stratégies pour résoudre le problème (1.16)-(1.17). La première consiste à exprimer  $\mathbf{u}^\dagger$  en fonction de  $p^\dagger$  pour éliminer une inconnue en utilisant l'hypothèse théorique d'une hauteur d'eau  $H$  strictement positive. Nous pouvons alors écrire une formulation conforme en pression munie des conditions de bord correspondantes :

Trouver  $p^\dagger : \Omega \rightarrow \mathbb{R}$  tel que,

$$-\operatorname{div}_s^\alpha \left( \frac{1}{H} \nabla_s^\alpha p^\dagger \right) = f - \operatorname{div}_s^\alpha \left( \frac{\mathbf{g}}{H} \right), \quad (1.18a)$$

$$(H\nabla p^\dagger + p^\dagger \nabla \zeta) \cdot \mathbf{n}_\Gamma = \mathbf{g}_1 \cdot \mathbf{n}_\Gamma - \phi \text{ sur } \Gamma_n, \quad (1.18b)$$

$$p^\dagger = 0 \text{ sur } \Gamma_d. \quad (1.18c)$$

La seconde stratégie consiste à garder la formulation mixte (1.16) munie des conditions aux limites (1.17). Les auteurs de [2] ont employé cette stratégie et ont introduit une formulation faible mixte avec un relèvement sur la vitesse, puis au niveau discret, les auteurs ont utilisé une méthode de condensation de masse (ou *mass-lumping*) pour réduire le coût numérique du schéma.

Ces deux stratégies sont étudiées dans le Chapitre 2 dans le cadre des Méthodes de Discrétisation du Gradient (GDM) [22]. Ces dernières permettent d'étudier la convergence d'une grande variété de méthodes numériques en identifiant quelques propriétés clés. Par exemple, les éléments finis conformes et non-conformes sont des GDM. Cet outil consiste à choisir des opérateurs et des espaces discrets pour la formulation faible du problème considéré, menant ainsi à un schéma numérique appelé Schéma Gradient. Dans ce chapitre, nous prouvons que le Schéma Gradient associé au problème (1.18) converge pour toute méthode appartenant aux GDM et un exemple d'application est donné pour la méthode des éléments finis conformes. Par la suite, nous montrons que le schéma éléments finis mixtes proposé dans la Section 4.1 du papier [2] est une Méthode de Discrétisation du Gradient et nous donnons les estimateurs d'erreur associés.

Dans nos travaux futurs, nous voulons implémenter le schéma éléments finis pour la formulation conforme (1.18) et le comparer au schéma proposé dans [2] en termes de précision et de temps de calcul.

Dans le Chapitre 3, réalisé en collaboration avec Emmanuel Audusse, Noémie Gaveau et Yohan Penel, nous rappelons que les schémas classiques de type Godunov pour les équations de Saint-Venant avec terme de Coriolis (1.11) manquent de précision autour de l'équilibre géostrophique (1.13). Ceci est dû au fait que les termes de diffusion qui apparaissent dans ces schémas ne s'annulent pas autour de l'équilibre géostrophique. Il est donc nécessaire de construire spécialement des schémas qui préserveront linéairement cet équilibre. Nous proposons donc des schémas impliquant une version discrétisée de l'équilibre géostrophique dans leurs termes de diffusion afin qu'ils s'annulent autour de l'équilibre. De plus, nous cherchons à avoir une stabilité entropique pour les versions semi-discrète de nos schémas lorsque c'est possible.

Nous présentons quatre schémas semi-discrets volumes finies sur maillages cartésiens qui seront comparés à un schéma HLLC de référence. Chacun de ces schémas ayant ses avantages et inconvénients, nous les présentons succinctement ici :

- le schéma *edge-based entropic* (3.7) a été présenté pour la première fois dans [25] (cf. Section A.11 de l'Annexe A) et a été reformulé dans ce chapitre dans le but d'être comparé aux schémas suivants.

Il s'agit d'un schéma à 9 points dont les termes de diffusion sont liés à une discrétisation de l'équilibre géostrophique sur les arêtes du maillage. En revanche, l'équilibre géostrophique qui apparaît dans l'équation de conservation de quantité de mouvement (1.11b) est discrétisé aux cellules afin de s'assurer de la stabilité entropique du schéma semi-discret. Ainsi, dans ce schéma l'équilibre géostrophique est discrétisé de deux façons différentes ce qui empêche le schéma linéarisé de préserver l'une ou l'autre. Le schéma perd donc le caractère linéairement équilibre au profit de la stabilité entropique du schéma linéarisé. Cet inconvénient vient avec un avantage, contrairement aux schémas suivants, il est très difficile d'avoir une solution parasite dans le noyau du schéma linéarisé (cf. Remarque 3.4.6), en pratique nous n'en observons pas dans nos cas tests numériques en Section 3.5.

- le schéma *edge-based well-balanced* (3.14) est un schéma à 9 points dont les termes de diffusion sont liés à une discrétisation de l'équilibre géostrophique sur les arêtes du maillage. De même, l'équilibre géostrophique qui apparaît dans l'équation de conservation de quantité de mouvement (1.11b) est discrétisé sur les arêtes ce qui permet d'assurer le caractère linéairement équilibre du schéma. En revanche, la stabilité entropique du schéma semi-discret n'est plus assurée avec ce choix. En effet, un terme de création d'énergie non contrôlé apparaît. Ainsi, ce schéma est le pendant du schéma (3.7) précédent, dans lequel la stabilité entropique est sacrifiée au profit

du caractère linéairement équilibre. Autre inconvénient, des modes parasites sur la vitesse sont présents dans le noyau des schémas non-linéaire et linéaire (cf. Remarque 3.4.6).

- le schéma *cell-based* (3.16) est un schéma à 13 points dont les termes de diffusion sont liés à une discrétisation de l'équilibre géostrophique sur les cellules du maillage. De même, l'équilibre géostrophique qui apparaît dans l'équation de conservation de quantité de mouvement (1.11b) est discrétisé sur les cellules ce qui permet d'assurer le caractère linéairement équilibre du schéma. La stabilité entropique du schéma semi-discret est également assurée avec ce choix. En revanche, des modes parasites sur la hauteur d'eau sont présents dans le noyau des schémas non-linéaire et linéaire (cf. Remarque 3.4.9).
- le schéma *vertex-based* (3.20) est un schéma à 25 points dont les termes de diffusion sont liés à une discrétisation de l'équilibre géostrophique sur les nœuds du maillage. De même, l'équilibre géostrophique qui apparaît dans l'équation de conservation de quantité de mouvement (1.11b) est discrétisé sur les nœuds ce qui permet d'assurer le caractère linéairement équilibre du schéma. La stabilité entropique du schéma semi-discret est également assurée avec ce choix. En revanche, des modes parasites sur la hauteur d'eau et sur la vitesse sont présents dans le noyau des schémas non-linéaire et linéaire (cf. Remarque 3.4.15). Ce schéma admet une version conservative (3.28) qui possède les mêmes propriétés.

Concernant la discrétisation en temps, les dérivées temporelles sont discrétisées explicitement et le terme de Coriolis est discrétisé de manière semi-implicite pour tous les schémas. Ainsi, nous considérons uniquement des systèmes résolus explicitement. Le choix de la discrétisation du terme de Coriolis est justifié par les problèmes de stabilité rencontrés par le schéma HLLC dans le cas d'une discrétisation explicite (cf. Section 3.4.4). Numériquement, les autres schémas ne semblent pas nécessiter une telle discrétisation, nous l'imposons uniquement afin de comparer les schémas entre eux.

Dans la Section 3.5, nous présentons des cas tests se concentrant sur des équilibres géostrophiques afin de montrer la précision autour de ces équilibres de nos schéma par rapport à un schéma de type Godunov classique. Nous présentons également une rupture de barrage pour montrer le comportement des schémas lorsque la condition initiale est très éloignée d'un équilibre géostrophique.

Nos futurs travaux porteront sur l'étude des modes parasites présents dans les noyaux des différents schémas, le but étant de les réduire au maximum, voire de les faire disparaître. Par la suite, il sera intéressant d'étudier la stabilité des schémas entièrement discrétisés. En parallèle, la question du maillage sera explorée afin d'étendre nos travaux aux maillages triangulaires non-structurés.

Dans l'Annexe A est présenté un *proceeding* issu d'un travail collaboratif réalisé lors du CEMRACS 2019 - Geophysical Fluids, Gravity Flows. Il s'agit du travail préliminaire dans lequel le schéma *edge-based entropic* (3.7) a été introduit sous une forme différente puisque les opérateurs ont été depuis remaniés pour être cohérents avec ceux nécessaires dans les autres schémas du Chapitre 3. Une démonstration de la décroissance de l'énergie semi-discrète ainsi que le résultat de non-préservation de l'équilibre géostrophique discret peuvent être trouvés en Section A.1.1. Un schéma semi-discret sur grille décalée (A.25) est présenté en Section A.1.2. On y montre la stabilité entropique ainsi que le caractère linéairement équilibre du schéma. Enfin, en Section A.2 les deux schémas sont comparés à un schéma HLLC sur le cas test d'un vortex stationnaire et se sont révélés plus précis que ce dernier.

# Chapitre 2

## Gradient Discretisation Methods to analyse numerical schemes for the elliptic part of a dispersive Shallow Water system.

*This work has been done in collaboration with Cindy Guichard and Yohan Penel. To be submitted*

### Contents

---

<b>2.1</b>	<b>Introduction</b>	<b>22</b>
<b>2.2</b>	<b>The Gradient Discretisation Method on <i>classic</i> operators</b>	<b>24</b>
2.2.1	A conforming formulation on the pressure	25
2.2.2	Weak formulation	25
2.2.3	Gradient Scheme	27
2.2.4	Key properties of the considered Gradient Discretisation	28
2.2.5	Error estimates	29
2.2.6	The $\mathbb{P}_1$ conforming finite elements	33
<b>2.3</b>	<b>The Abstract Gradient Discretisation Method on <i>special</i> operators</b>	<b>34</b>
2.3.1	Weak formulation	34
2.3.2	Abstract Gradient Scheme	35
2.3.3	Key properties of the considered Abstract Gradient Discretisation	36
2.3.4	Error estimates	37
2.3.5	The $\mathbb{P}_1$ conforming finite elements	39
<b>2.4</b>	<b>Analyse of a mixed mass-lumped scheme</b>	<b>40</b>
2.4.1	A weak mixed formulation with lifting	40
2.4.2	The associated Abstract Gradient Discretisation Method framework	41
2.4.3	Mixed mass-lumping finite element methods & AGD	43
2.4.4	The $\mathbb{P}_1/\mathbb{P}_1$ approximation	47
<b>2.5</b>	<b>Conclusion</b>	<b>56</b>

---

## 2.1 Introduction

Modelisation of free surface flows plays an important role in many engineering applications such as ocean circulation, coastal exploitation, man-made structures in rivers or lakes. Usually, this kind of flow is described by the three dimensional Navier-Stokes Equations (NSE) [32], assuming the fluid to be Newtonian, viscous and incompressible. Obviously, the real phenomena is always much more complicated, but in practice, the NSE provide a relevant mathematical model to approach the reality of oceans and river flows. At the numerical level, difficulties arise when solving these equations for real applications, especially for domains with large dimensions (e.g. oceans, lakes, rivers). Computationally, the complete resolution of the NSE for a free surface flow is known to be dramatically onerous. In order to reduce the numerical cost, we are interested in models where the dimensions are reduced. From the mathematical point of view, geophysical free surface models can be derived from the Navier-Stokes system under some assumptions.

For these reasons, and when the fluid domain can be regarded as a thin layer of fluid, it is usual to consider the nonlinear shallow water equations. This assumption implies that the model is hydrostatic and states that the velocity is constant along the water height. More generally, Non-linear Shallow Water equations (NLSW) model the dynamics of a shallow, rotating layer of homogeneous incompressible and inviscid fluid. NLSW equations are well-suited for the study and numerical simulations of a large class of geophysical phenomena, such as river flows, coastal domains, ocean circulation, or even run-off or avalanches when modified with adapted source terms. The usual conservative form of NLSW equations introduced in [46] is written as a first order hyperbolic system with various source terms such as bed slope or bottom friction terms.

However these equations rely on the hydrostatic assumption and hence when the vertical acceleration of the fluid can no longer be neglected, the shallow water system fails to represent dispersive effects e.g. in the context of wave propagation [8, 33]. It is a challenging issue to reduce a model from NSE with good properties and for which we can develop a practicable numerical method. We stress that there is not a single satisfactory model which is valid for many regimes with a good mathematical structure, each of them has a limited range of validity.

In this work, we are interested in a family of dispersive models, introduced in [2], depending on a parameter that characterises the type of model. For a given value of this parameter the authors obtain the dispersive model proposed in [1, 15] and for another value of this parameter the studied model corresponds to the Green-Naghdi model [12, 17, 35] up to some small error terms involving the bathymetry gradient. In the model considered in [2], only first order derivatives appear unlike most of existing model. The authors use a time splitting based on a Chorin-Temam projection-correction scheme [19, 47], which leads to decompose the problem into two steps at each time step. In the first one, a hyperbolic problem is solved with a finite volume scheme yielding water height  $H$  and discharge  $Hu$  whereas, in the second one, an elliptic problem governing the non-hydrostatic pressure  $p$  is solved with a finite element scheme (computationally, this is the most onerous part). From a triangular primal mesh a dual mesh centered on the vertices is used for a finite volume scheme. The unknowns of the hyperbolic problem in the predictive step are constant on the dual cells. Then, for the finite element scheme, the unknowns are approximated at the vertices of the triangles. In the correction step, we use the constant values computed on the dual cell at the prediction step. In this chapter, we focus on the elliptic problem and its approximation.

Let  $\Omega$  be a subset of  $\mathbb{R}^d$ , where  $d = 1$  or  $d = 2$ . We introduce  $\eta : \Omega \rightarrow \mathbb{R}$  the water elevation and

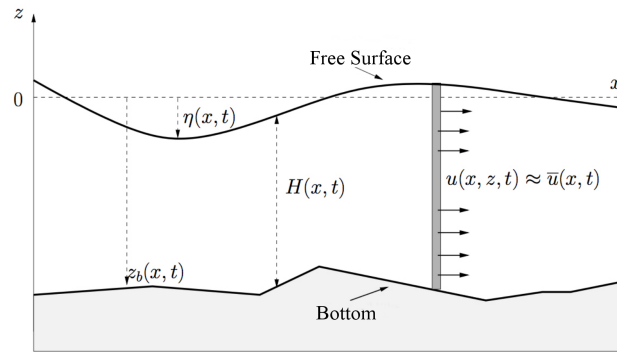


FIGURE 2.1 – Notations for the one-dimensional Shallow Water system. Source : [1]

$z_b : \Omega \rightarrow \mathbb{R}$  the bathymetry. The water height is given by  $H = \eta - z_b$ . Let us introduce the notation  $\zeta = H + \frac{\alpha^2}{2} z_b$ .

As in [2, section 2.4.1], in order to write the model under a convenient form, we introduce specific first order differential operators that we call in the following *special* gradient and *special* divergence, respectively denoted by  $\nabla_s^\alpha$  and  $\text{div}_s^\alpha$ , which are defined for  $\varphi : \Omega \rightarrow \mathbb{R}$  by,

$$\nabla_s^\alpha \varphi = (H \nabla \varphi + \varphi \nabla \zeta, -\alpha \varphi)^T \quad (2.1a)$$

and for  $\boldsymbol{\omega} = (\mathbf{v}, w)^T$  with  $\mathbf{v} : \Omega \rightarrow \mathbb{R}^d$  and  $w : \Omega \rightarrow \mathbb{R}$ ,

$$\text{div}_s^\alpha \boldsymbol{\omega} = \text{div}(H \mathbf{v}) - \mathbf{v} \cdot \nabla \zeta + \alpha w, \quad (2.1b)$$

where, as specified above,  $\alpha$  is the parameter which allows us to deal with several models such as the Green-Naghdi model, up to small error terms if  $\alpha = \sqrt{3}$  (exact if  $z_b = \text{cst.}$ ) [35] and the Depth-Averaged Euler model [1, 15] if  $\alpha = 2$ .

Then the problem reads :

Find  $p^\dagger : \Omega \rightarrow \mathbb{R}$  the non-hydrostatic pressure and  $\mathbf{u}^\dagger = (\mathbf{v}^\dagger, w^\dagger)^T : \Omega \rightarrow \mathbb{R}^{d+1}$  the velocity field such that,

$$H \mathbf{u}^\dagger + \nabla_s^\alpha p^\dagger = \mathbf{g} \text{ in } \Omega, \quad (2.2a)$$

$$\text{div}_s^\alpha \mathbf{u}^\dagger = f \text{ in } \Omega, \quad (2.2b)$$

with the following boundary conditions,

$$H \mathbf{u}^\dagger \cdot \mathbf{n}_s = H \mathbf{v}^\dagger \cdot \mathbf{n}_\Gamma = \phi \text{ on } \Gamma_n, \quad (2.3)$$

and

$$p^\dagger = 0 \text{ on } \Gamma_d, \quad (2.4)$$

where  $\mathbf{n}_\Gamma$  is the unit normal vector outward  $\Omega$  and where we define  $\mathbf{n}_s := (\mathbf{n}_\Gamma, 0)^T$ . These boundary conditions are in the spirit of the ones used in the litterature, the reader can refer to [2, section 3.4.1] for more details about their compatibility with the full model.

We consider the following hypotheses :



- **Domain**

$\Omega$  is an open connected polygonal bounded domain with Lipschitz boundary, denoted by  $\Gamma$ , such that  $\Gamma_d, \Gamma_n$  are two disjoint relatively open subsets of  $\Gamma$  with  $|\Gamma \setminus (\Gamma_d \cup \Gamma_n)| = 0$ , (2.5a)

( $|\cdot|$  denotes the  $(d-1)$ -dimensional measure);

- **Water height**

$$H \in \mathcal{C}^2(\overline{\Omega}) \text{ and there exist } \underline{H}, \overline{H} > 0 \text{ such that, for } \mathbf{x} \in \Omega, \underline{H} \leq H(\mathbf{x}) \leq \overline{H}, \quad (2.5b)$$

- **Other data**

$$\zeta \in \mathcal{C}^2(\overline{\Omega}), \mathbf{g} = (g_1, g_2)^T \in L^2(\Omega)^d \times L^2(\Omega), \operatorname{div}(\mathbf{g}) \in L^2(\Omega), f \in L^2(\Omega), \phi \in L^2(\Gamma_n), \alpha \in \mathbb{R}^*. \quad (2.5c)$$

Our first goal is to analyse the convergence of a scheme introduced in [2] through the Gradient Discretisation Method (GDM). The GDM, detailed in [22], is a framework which contains classic and recent discretisation schemes for diffusion problems of different kinds : linear or non-linear, steady-state or time-dependent. The schemes may be conforming or non-conforming, low or high order, and may be built on very general meshes. A GDM is based on the choice of a set of discrete spaces and operators, referred to as a ‘‘Gradient Discretisation’’ (GD). Replacing, in the weak formulation of a diffusion problem, the continuous space and operators by the discrete elements provided by a particular GD yields a numerical scheme called a ‘‘Gradient Scheme’’ (GS). The concept of gradient scheme has been introduced in [31] for diffusive flows in heterogeneous anisotropic porous media with the aim of deriving a scheme which leads to symmetric positive-definite matrices. Then key properties are shown to be sufficient to prove the convergence of gradient schemes for linear and nonlinear elliptic and parabolic problems in [23]. It is shown that a number of methods are GDMs such as the conforming finite elements, the non-conforming finite elements, one discontinuous Galerkin method, the hybrid mimetic mixed or nodal mimetic finite differences methods, some discrete duality finite volume schemes, and some multi-point flux approximation schemes. The GDM has been recently extended in [24] with the Abstract Gradient Discretisation Method (AGDM), which has been developed in order to provide a unified convergence analysis that simultaneously covers all usual boundary conditions for the approximation in Banach spaces of operators acting in duality. In this work, this method allows us to study the considered problem with abstract operators and obtain a convergence result for a scheme proposed in [2] for the elliptic problem.

The chapter is organised as follows. In the next section, we propose a conforming formulation of (2.2) on the pressure and the corresponding GS. Then an application to the conforming finite element method is given. Section 2.3 is devoted to the AGDM applied to the same formulation as in Section 2.2 and to a weak mixed formulation of the strong problem (2.2) where a lifting has been operated. Finally in Section 2.4, we prove the convergence of a scheme proposed in [2] through the AGDM.

## 2.2 The Gradient Discretisation Method on *classic* operators

To our knowledge, Problem (2.2) has only been studied in the literature thanks to a mixed approach. We propose in this section a conforming formulation only on the pressure  $p^\dagger$  which has the advantage to be easier to analyse and discretise. Then the velocity  $\mathbf{u}^\dagger$  can be reconstructed from the pressure  $p^\dagger$ . We write a gradient scheme resulting from the weak problem associated to the pressure. Then, error estimates on the pressure are obtained. Finally, we choose a finite-element method as an example.

### 2.2.1 A conforming formulation on the pressure

From Problem (2.2), we deduce a conforming formulation by expressing the velocity  $\mathbf{u}^\dagger$  as a function of the pressure  $p^\dagger$ . Indeed, from Equation (2.2a) and thanks to Assumption (2.5b) (strict positivity of  $H$ ), we can express  $\mathbf{u}^\dagger$  as follows :

$$\mathbf{u}^\dagger = \frac{1}{H} \left( \mathbf{g} - \nabla_s^\alpha p^\dagger \right) \stackrel{(2.1a)}{=} \frac{1}{H} \begin{pmatrix} \mathbf{g}_1 - H \nabla p^\dagger - p^\dagger \nabla \zeta \\ g_2 + \alpha p^\dagger \end{pmatrix} = \begin{pmatrix} \mathbf{v}^\dagger \\ w^\dagger \end{pmatrix}. \quad (2.6)$$

#### Remark 2.2.1

Using the classic Gradient Discretisation Method here, we deal with the classic gradient and divergence operators. In Section 2.3 we use the Abstract Gradient Discretisation framework in order to deal with the special operators (2.1a) and (2.1b).

From now on, the velocity  $\mathbf{u}^\dagger$  is a secondary unknown since it can be expressed according to the main unknown  $p^\dagger$ . By injecting (2.6) into Equation (2.2b) we get,

$$-\operatorname{div}_s^\alpha \left( \frac{1}{H} \nabla_s^\alpha p^\dagger \right) = f - \operatorname{div}_s^\alpha \left( \frac{\mathbf{g}}{H} \right).$$

By using the expression of the *special* divergence (2.1b), finally the problem (2.2)-(2.4) becomes, with the *classic* operators :

Find  $p^\dagger : \Omega \rightarrow \mathbb{R}$  such that,

$$-\operatorname{div}(H \nabla p^\dagger + p^\dagger \nabla \zeta) + \nabla p^\dagger \cdot \nabla \zeta + \frac{p^\dagger}{H} (|\nabla \zeta|^2 + \alpha^2) = f - \operatorname{div}(\mathbf{g}_1) + \frac{1}{H} (\mathbf{g}_1 \cdot \nabla \zeta - \alpha g_2) \text{ in } \Omega, \quad (2.7a)$$

$$(H \nabla p^\dagger + p^\dagger \nabla \zeta) \cdot \mathbf{n}_\Gamma = \mathbf{g}_1 \cdot \mathbf{n}_\Gamma - \phi \text{ on } \Gamma_n, \quad (2.7b)$$

$$p^\dagger = 0 \text{ on } \Gamma_d, \quad (2.7c)$$

and with  $\mathbf{u}^\dagger$  defined by (2.6).

### 2.2.2 Weak formulation

We now derive a weak formulation of Problem (2.7) in order to introduce the associated gradient scheme. Hence, a function  $p^\dagger$  is said to be a weak solution of Problem (2.7) if the following holds :

Find  $p^\dagger \in H_d^1(\Omega)$ , such that  $\forall q^\dagger \in H_d^1(\Omega)$ ,

$$\begin{aligned} & \int_\Omega \frac{(H \nabla p^\dagger + p^\dagger \nabla \zeta) \cdot (H \nabla q^\dagger + q^\dagger \nabla \zeta)}{H} + \alpha^2 \frac{p^\dagger q^\dagger}{H} \, d\mathbf{x} \\ & = \int_\Omega \left( f + \frac{\mathbf{g}_1 \cdot \nabla \zeta - \alpha g_2}{H} \right) q^\dagger + \mathbf{g}_1 \cdot \nabla q^\dagger \, d\mathbf{x} - \int_{\Gamma_n} \phi \gamma(q^\dagger) \, d\sigma(\mathbf{x}), \end{aligned} \quad (2.8)$$

where  $\gamma(\cdot)$  is the trace operator on  $\Gamma$  and the Hilbert space,

$$H_d^1(\Omega) = \{\varphi \in H^1(\Omega) \mid \gamma(\varphi) = 0 \text{ on } \Gamma_d\}, \quad (2.9)$$

with  $\|\varphi\|_{H_d^1(\Omega)}^2 = \|\varphi\|_{L^2(\Omega)}^2 + \|\nabla\varphi\|_{L^2(\Omega)^d}^2$ .

A solution of Problem (2.8) exists and is unique as stated in the following lemma.

**Lemma 2.2.2** (*Existence and uniqueness for the conforming formulation*)

Under Assumptions (2.5), there exists one and only one solution to Problem (2.8).

**Proof:** In order to apply the Lax-Milgram theorem, we introduce the bilinear form  $a(\cdot, \cdot)$  and the linear form  $l(\cdot)$  :

$$a(p, q) = \int_{\Omega} \frac{(H\nabla p + p\nabla\zeta) \cdot (H\nabla q + q\nabla\zeta)}{H} + \alpha^2 \frac{pq}{H} \, d\mathbf{x},$$

$$l(q) = \int_{\Omega} \left( f + \frac{\mathbf{g}_1 \cdot \nabla\zeta - \alpha g_2}{H} \right) q + \mathbf{g}_1 \cdot \nabla q \, d\mathbf{x} - \int_{\Gamma_n} \phi \gamma(q) \, d\sigma(\mathbf{x}),$$

defined for all  $p, q \in H_d^1(\Omega)$ . We can underline that the bilinear form  $a(\cdot, \cdot)$  is symmetric thus this proof could also be viewed as an application of the Riesz representation theorem. The continuity of these two forms are proved thanks to Assumptions (2.5) and Cauchy-Schwarz inequality ; the trace continuity is also needed for  $l(\cdot)$ . The main point is the *coercivity* – also said *ellipticity* – of  $a(\cdot, \cdot)$ . First of all, notice that for all  $\nu$  and  $\omega$  in  $L^2(\Omega)^d$ , we can prove that,

$$\int_{\Omega} \frac{|\nu|^2}{H} \, d\mathbf{x} \leq 2 \int_{\Omega} \frac{|\nu + \omega|^2}{H} \, d\mathbf{x} + 2 \int_{\Omega} \frac{|\omega|^2}{H} \, d\mathbf{x},$$

by taking  $\nu = H\nabla p^\dagger$  and  $\omega = p\nabla\zeta$ , we can write, also by using (2.5),

$$\underbrace{\int_{\Omega} H|\nabla p|^2 \, d\mathbf{x}}_{\geq \underline{H}\|\nabla p^\dagger\|_{L^2(\Omega)^d}^2} \leq 2 \underbrace{\int_{\Omega} \frac{|H\nabla p + p\nabla\zeta|^2}{H} \, d\mathbf{x}}_{\leq a(p, p)} + 2 \int_{\Omega} \frac{p^2|\nabla\zeta|^2}{H} \, d\mathbf{x}.$$

Since, still using (2.5),

$$\int_{\Omega} \frac{p^2|\nabla\zeta|^2}{H} \, d\mathbf{x} \leq \frac{\|\nabla\zeta\|_{L^\infty(\Omega)^d}^2}{\alpha^2} \underbrace{\int_{\Omega} \alpha^2 \frac{p^2}{H} \, d\mathbf{x}}_{\leq a(p, p)}.$$

So we have,

$$a(p, p) \geq \frac{\alpha^2 \underline{H}}{2\alpha^2 + 2\|\nabla\zeta\|_{L^\infty(\Omega)^d}^2} \|\nabla p\|_{L^2(\Omega)^d}^2.$$

Moreover we can write, thanks to (2.5),

$$\|p\|_{L^2(\Omega)}^2 \leq \frac{\overline{H}}{\alpha^2} \underbrace{\int_{\Omega} \frac{\alpha^2 p^2}{H} \, d\mathbf{x}}_{\leq a(p, p)}.$$

Finally, we conclude the proof by writing,

$$a(p, p) \geq \frac{\alpha^2}{2} \min \left( \frac{\underline{H}}{2\alpha^2 + 2\|\nabla\zeta\|_{L^\infty(\Omega)^d}^2}, \frac{1}{\overline{H}} \right) \|p\|_{H_d^1(\Omega)}^2. \quad (2.10) \quad \blacksquare$$

**Remark 2.2.3** (*The Dirichlet boundary condition is not mandatory*)

Let us mention that the assumption  $|\Gamma_d| > 0$  is not used in order to prove Lemma 2.2.2.

This is due to the term  $\alpha^2 \int_{\Omega} \frac{p^2}{H} d\mathbf{x}$  in  $a(p, p)$  (where  $a(\cdot, \cdot)$  is defined in the proof) which provides a bound with respect to  $\|p\|_{H_d^1(\Omega)}^2$  without using a Poincaré inequality.

**2.2.3 Gradient Scheme**

In order to introduce a gradient scheme, four discrete objects  $\mathcal{D} = (X_{\mathcal{D}} = X_{\mathcal{D}, \Omega, \Gamma_n} \oplus X_{\mathcal{D}, \Gamma_d}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, \nabla_{\mathcal{D}})$ , forming a Gradient Discretisation (GD), are introduced.

**Definition 2.2.4** (*GD, homogeneous Dirichlet and non-homogeneous Neumann BCs*)

Under the assumption (2.5a), a gradient discretisation  $\mathcal{D}$  for mixed boundary conditions is defined by  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, \nabla_{\mathcal{D}})$  where :

1. the set of discrete unknowns  $X_{\mathcal{D}} = X_{\mathcal{D}, \Omega, \Gamma_n} \oplus X_{\mathcal{D}, \Gamma_d}$  is the direct sum of two finite-dimensional vector spaces on  $\mathbb{R}$ , corresponding respectively to the unknowns in  $\Omega$  and on  $\Gamma_n$ , and with  $X_{\mathcal{D}, \Gamma_d} = \{0\}$  due to homogeneous Dirichlet condition on  $\Gamma_d$ ,
2. the function reconstruction  $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  is linear,
3. the trace reconstruction  $\mathbb{T}_{\mathcal{D}, \Gamma_n} : X_{\mathcal{D}} \rightarrow L^2(\Gamma_n)$  is linear,
4. the gradient reconstruction  $\nabla_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^d$  is linear,
5. the operators are such that  $\|\cdot\|_{\mathcal{D}}$ , with, for  $v \in X_{\mathcal{D}, \Omega, \Gamma_n}$ ,  $\|v\|_{\mathcal{D}}^2 := \|\Pi_{\mathcal{D}}v\|_{L^2(\Omega)}^2 + \|\nabla_{\mathcal{D}}v\|_{L^2(\Omega)^d}^2$ , is a norm on  $X_{\mathcal{D}, \Omega, \Gamma_n}$ .

**Remark 2.2.5**

The definition of a GD depends on the boundary conditions (see Part I. of [22]) of the studied problem. Here, boundary conditions of the strong problem (2.7) include a homogeneous Dirichlet condition (2.7c) and a Robin/Fourier condition (2.7b). However in the weak formulation (2.8), the term on  $\Gamma_n$  is closer to a non-homogeneous Neumann condition since there is no term under the form  $\int_{\Gamma_n} \gamma(p^\dagger) \gamma(q^\dagger) d\sigma(\mathbf{x})$ . This is the reason why Definition 2.2.4 of a GD has been established for homogeneous Dirichlet and non-homogeneous Neumann conditions.

We now introduce the Gradient Scheme (GS) for the approximation of Problem (2.8) which is designed by *replacing* the continuous spaces and operators by discrete counterparts.

**Definition 2.2.6 (GS, Conforming formulation)**

Let  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, \nabla_{\mathcal{D}})$  a GD in the sense of Definition 2.2.4. The related gradient scheme for Problem (2.8) is defined by :

Find  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  such that for any  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$ ,

$$\int_{\Omega} \frac{(H \nabla_{\mathcal{D}} p + \Pi_{\mathcal{D}} p \nabla \zeta) \cdot (H \nabla_{\mathcal{D}} q + \Pi_{\mathcal{D}} q \nabla \zeta)}{H} + \alpha^2 \frac{\Pi_{\mathcal{D}} p \Pi_{\mathcal{D}} q}{H} \, dx$$

$$= \int_{\Omega} \left( f + \frac{\mathbf{g}_1 \cdot \nabla \zeta - \alpha g_2}{H} \right) \Pi_{\mathcal{D}} q + \mathbf{g}_1 \cdot \nabla_{\mathcal{D}} q \, dx - \int_{\Gamma_n} \phi \mathbb{T}_{\mathcal{D}, \Gamma_n} q \, d\sigma(\mathbf{x}). \quad (2.11)$$

**2.2.4 Key properties of the considered Gradient Discretisation**

In the following, we will provide general properties on the GD, given by Definition 2.2.4, that ensure the convergence of the corresponding GS (2.11). These properties are named GD-Coercivity, GD-Consistency, GD-Trace-Consistency and GD-Limit-Conformity.

**Definition 2.2.7 (GD-Coercivity)**

Under the assumption (2.5a), if  $\mathcal{D}$  is a gradient discretisation in the sense of Definition 2.2.4, we define,

$$C_{\mathcal{D}} = \max_{v \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \left( \frac{\|\mathbb{T}_{\mathcal{D}, \Gamma_n} v\|_{L^2(\Gamma_n)}}{\|v\|_{\mathcal{D}}} \right). \quad (2.12)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of gradient discretisations in the sense of Definition 2.2.4 is **coercive** if there exists  $C_P \in \mathbb{R}_+$  such that  $C_{\mathcal{D}_m} \leq C_P$  for all  $m \in \mathbb{N}$ .

**Remark 2.2.8 (Norm on the discrete space)**

In [22] the norm  $\|\cdot\|_{\mathcal{D}}$  is mainly defined by  $\|\nabla_{\mathcal{D}} \cdot\|_{L^2(\Omega)^d}$  and the GD-Coercivity introduces a kind of discrete Poincaré inequality. However, it is not the case for the problem studied here, (2.12) defines a discrete trace inequality, it is the consequence, at the discrete level, of Remark 2.2.3.

**Definition 2.2.9 (GD-Consistency)**

Under the assumption (2.5a) and using the definition (2.9) for  $H_d^1(\Omega)$ , if  $\mathcal{D}$  is a gradient discretisation in the sense of Definition 2.2.4, we define  $S_{\mathcal{D}} : H_d^1(\Omega) \rightarrow [0, +\infty)$  by,

$$\forall \varphi \in H_d^1(\Omega), \quad S_{\mathcal{D}}(\varphi) = \min_{v \in X_{\mathcal{D}, \Omega, \Gamma_n}} \left\{ \|\Pi_{\mathcal{D}} v - \varphi\|_{L^2(\Omega)} + \|\nabla_{\mathcal{D}} v - \nabla \varphi\|_{L^2(\Omega)^d} \right\}. \quad (2.13)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of gradient discretisations in the sense of Definition 2.2.4 is **consistent** if,

$$\forall \varphi \in H_d^1(\Omega), \quad \lim_{m \rightarrow \infty} S_{\mathcal{D}_m}(\varphi) = 0. \quad (2.14)$$

**Definition 2.2.10 (GD-Trace-Consistency)**

Under the assumption (2.5a) and using the definition (2.9) for  $H_d^1(\Omega)$ , if  $\mathcal{D}$  is a gradient discretisation in the sense of Definition 2.2.4, we define  $\bar{S}_{\mathcal{D}} : H_d^1(\Omega) \rightarrow [0, +\infty)$  by,

$$\forall \varphi \in H_d^1(\Omega), \bar{S}_{\mathcal{D}}(\varphi) = \min_{v \in X_{\mathcal{D}, \Omega, \Gamma_n}} \left\{ \|\Pi_{\mathcal{D}} v - \varphi\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D}, \Gamma_n} v - \gamma(\varphi)\|_{L^2(\Gamma_n)} + \|\nabla_{\mathcal{D}} v - \nabla \varphi\|_{L^2(\Omega)^d} \right\}. \quad (2.15)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of gradient discretisations in the sense of Definition 2.2.4 is **trace-consistent** if,

$$\forall \varphi \in H_d^1(\Omega), \lim_{m \rightarrow \infty} \bar{S}_{\mathcal{D}_m}(\varphi) = 0. \quad (2.16)$$

The Definitions of the GD-Consistency and GD-Trace-Consistency have to be understood in the sense of an interpolation error. Moreover, we can underline that the quantity  $\bar{S}_{\mathcal{D}}$  in (2.15) is needed in order to control the discrete trace  $\mathbb{T}_{\mathcal{D}, \Gamma_n}$ .

Finally, since a GDM could be a nonconforming method, we need that the dual of the discrete gradient be “close to” a discrete divergence.

**Definition 2.2.11 (GD-Limit-Conformity)**

Under the assumption (2.5a), let  $\mathcal{D}$  is a gradient discretisation in the sense of Definition 2.2.4. We introduce,

$$H_{\text{div}, \Gamma_n}(\Omega) = \{\varphi \in L^2(\Omega)^d : \text{div} \varphi \in L^2(\Omega), \gamma_{\mathbf{n}}(\varphi) \in L^2(\Gamma_n)\}, \quad (2.17)$$

where  $\gamma_{\mathbf{n}}(\varphi)$  is the normal trace of  $\varphi$  on  $\Gamma_n$ . We now define  $W_{\mathcal{D}} : H_{\text{div}, \Gamma_n}(\Omega) \rightarrow [0, +\infty)$  by,

$$\forall \varphi \in H_{\text{div}, \Gamma_n}(\Omega),$$

$$W_{\mathcal{D}}(\varphi) = \max_{v \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \frac{1}{\|v\|_{\mathcal{D}}} \left| \int_{\Omega} (\nabla_{\mathcal{D}} v \cdot \varphi + \Pi_{\mathcal{D}} v \text{div} \varphi) \, d\mathbf{x} - \int_{\Gamma_n} \mathbb{T}_{\mathcal{D}, \Gamma_n} v \gamma_{\mathbf{n}}(\varphi) \, d\sigma(\mathbf{x}) \right|. \quad (2.18)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of gradient discretisations in the sense of Definition 2.2.4 is **limit-conforming** if,

$$\forall \varphi \in H_{\text{div}, \Gamma_n}(\Omega), \lim_{m \rightarrow \infty} W_{\mathcal{D}_m}(\varphi) = 0. \quad (2.19)$$

**2.2.5 Error estimates**

The error estimates of the GS (2.11) are stated in the next two theorems. We then state a corollary on the convergence result that stems from these error estimates.

**Theorem 2.2.12 (Error estimate of a GS, Conforming formulation)**

Under Assumptions (2.5), let  $p^\dagger \in H_d^1(\Omega)$  be the solution of Problem (2.8). Let  $\mathcal{D}$  be a GD in the sense of Definition 2.2.4. Then there exists one and only one  $p \in X_{\mathcal{D},\Omega,\Gamma_n}$  solution to the GS (2.11); this solution satisfies the following inequality :

$$\|p^\dagger - \Pi_{\mathcal{D}}p\|_{L^2(\Omega)} + \|\nabla p^\dagger - \nabla_{\mathcal{D}}p\|_{L^2(\Omega)^d} \leq S_{\mathcal{D}}(p^\dagger) + \frac{W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) + S_{\mathcal{D}}(p^\dagger)C_S}{C_{LM}}, \quad (2.20)$$

where  $S_{\mathcal{D}}$  and  $W_{\mathcal{D}}$  are respectively defined by the GD-Consistency (2.13) and the GD-Limit-Conformity (2.18), moreover,  $C_S$  and  $C_{LM}$  only depend on the data of the problem (see (2.5)).

**Proof:** Let us first prove that, if (2.20) holds for any solution  $p \in X_{\mathcal{D},\Omega,\Gamma_n}$  to Scheme (2.11), then the solution to Scheme (2.8) exists and is unique. By introducing a basis of the space  $X_{\mathcal{D},\Omega,\Gamma_n}$ , it is easy to show that the Scheme (2.11) can be rewritten under the form of a square linear system denoted  $A_{\mathcal{D}}P_{\mathcal{D}} = F_{\mathcal{D}}$  where the unknown vector  $P_{\mathcal{D}}$  is made of the coordinates of  $p$  in the mentioned basis. By choosing  $f = 0$ ,  $\mathbf{g} = 0$  and  $\phi = 0$ , the right hand side of the continuous problem (2.8) vanishes and we have also  $F_{\mathcal{D}} = 0$ . In this case the unique solution of (2.8) is  $p^\dagger = 0$ . Then the right hand side of the inequality of (2.20) is equal to zero, thus we get that any solution  $p$  to the scheme satisfies  $\|p\|_{\mathcal{D}} = 0$ . Since  $\|\cdot\|_{\mathcal{D}}$  is a norm on  $X_{\mathcal{D},\Omega,\Gamma_n}$  (see Definition 2.2.4 of GD), this leads to  $p = 0$ . Therefore the kernel of the matrix  $A_{\mathcal{D}}$  is  $\{0\}$  which is equivalent to say that the linear system has a unique solution.

Let us now prove that any solution  $p \in X_{\mathcal{D},\Omega,\Gamma_n}$  to Scheme (2.11) satisfies (2.20). Firstly we can notice that  $p^\dagger \in H_d^1(\Omega)$  implies, in the distribution sense, that  $-\operatorname{div}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \in L^2(\Omega)^d$  and  $(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \cdot \mathbf{n}_\Gamma = -\phi \in L^2(\Gamma_n)$  therefore, by the definition (2.17),  $(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \in H_{\operatorname{div},\Gamma_n}(\Omega)$ . Thus we can take  $\boldsymbol{\varphi} = H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1$  in  $W_{\mathcal{D}}$  (2.18). We then obtain, for a given  $q \in X_{\mathcal{D},\Omega,\Gamma_n}$ ,

$$\left| \int_{\Omega} \nabla_{\mathcal{D}}q \cdot (H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) + \Pi_{\mathcal{D}}q \operatorname{div}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \, d\mathbf{x} - \int_{\Gamma_n} \mathbb{T}_{\mathcal{D},\Gamma_n}q \, \gamma_{\mathbf{n}}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \, d\sigma(\mathbf{x}) \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1).$$

Since, by (2.7a),  $-\operatorname{div}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) = f + \frac{\mathbf{g}_1 \cdot \nabla\zeta - \alpha g_2}{H} - \nabla p^\dagger \cdot \nabla\zeta - \frac{p^\dagger(|\nabla\zeta|^2 + \alpha^2)}{H}$  a.e. in  $\Omega$  and, by (2.7b),  $(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) \cdot \mathbf{n}_\Gamma = -\phi$  a.e. on  $\Gamma_n$ , this leads to,

$$\left| \int_{\Omega} \nabla_{\mathcal{D}}q \cdot (H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1) - \Pi_{\mathcal{D}}q \left( f + \frac{\mathbf{g}_1 \cdot \nabla\zeta - \alpha g_2}{H} - \nabla p^\dagger \cdot \nabla\zeta - \frac{p^\dagger(|\nabla\zeta|^2 + \alpha^2)}{H} \right) \, d\mathbf{x} + \int_{\Gamma_n} \mathbb{T}_{\mathcal{D},\Gamma_n}q \, \phi \, d\sigma(\mathbf{x}) \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1).$$

Now, since  $p$  is a solution of (2.11), we get,

$$\left| \int_{\Omega} \frac{(H\nabla_{\mathcal{D}}q + \Pi_{\mathcal{D}}q\nabla\zeta) \cdot (H(\nabla p^\dagger - \nabla_{\mathcal{D}}p) + (p^\dagger - \Pi_{\mathcal{D}}p)\nabla\zeta)}{H} + \alpha^2 \frac{\Pi_{\mathcal{D}}q(p^\dagger - \Pi_{\mathcal{D}}p)}{H} \, d\mathbf{x} \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger\nabla\zeta - \mathbf{g}_1). \quad (2.21)$$

We now define,

$$I_{\mathcal{D}}p^\dagger = \operatorname{argmin}_{s \in X_{\mathcal{D},\Omega,\Gamma_n}} (\|\Pi_{\mathcal{D}}s - p^\dagger\|_{L^2(\Omega)} + \|\nabla_{\mathcal{D}}s - \nabla p^\dagger\|_{L^2(\Omega)^d}), \quad (2.22)$$

and notice that, by definition (2.13) of  $S_{\mathcal{D}}$ ,

$$\|\Pi_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - p^\dagger\|_{L^2(\Omega)} + \|\nabla_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - \nabla p^\dagger\|_{L^2(\Omega)^d} = S_{\mathcal{D}}(p^\dagger). \quad (2.23)$$

We also introduce the bilinear form on  $X_{\mathcal{D},\Omega,\Gamma_n} \times X_{\mathcal{D},\Omega,\Gamma_n}$ ,

$$a_{\mathcal{D}}(q, s) = \int_{\Omega} \frac{(H\nabla_{\mathcal{D}}q + \Pi_{\mathcal{D}}q \nabla\zeta) \cdot (H\nabla_{\mathcal{D}}s + \Pi_{\mathcal{D}}s \nabla\zeta)}{H} + \alpha^2 \frac{\Pi_{\mathcal{D}}q \Pi_{\mathcal{D}}s}{H} \, d\mathbf{x}.$$

By triangular inequality and by (2.21), we get,

$$\begin{aligned} & \left| a_{\mathcal{D}}(q, I_{\mathcal{D}}p^\dagger - p) \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla\zeta - \mathbf{g}_1) \\ & + \left| \int_{\Omega} \frac{(H\nabla_{\mathcal{D}}q + \Pi_{\mathcal{D}}q \nabla\zeta) \cdot (H(\nabla_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - \nabla p^\dagger) + (\Pi_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - p^\dagger) \nabla\zeta)}{H} + \alpha^2 \frac{\Pi_{\mathcal{D}}q (\Pi_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - p^\dagger)}{H} \, d\mathbf{x} \right|. \end{aligned}$$

We now have to establish an upper bound of the last term of the previous inequality. Recalling the definition of  $\|\cdot\|_{\mathcal{D}}$  in Definition 2.2.4, by using the Cauchy-Schwarz inequality, the assumption (2.5b) and (2.23), the right hand side is bounded by

$$\begin{aligned} & \|q\|_{\mathcal{D}} \left[ W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla\zeta - \mathbf{g}_1) + S_{\mathcal{D}}(p^\dagger) \underbrace{\left( \bar{H} + 2\|\nabla\zeta\|_{L^\infty(\Omega)^d} + \frac{\|\nabla\zeta\|_{L^\infty(\Omega)^d}^2 + \alpha^2}{\underline{H}} \right)}_{:= C_S} \right] \\ & \geq \left| a_{\mathcal{D}}(q, I_{\mathcal{D}}p^\dagger - p) \right|. \end{aligned}$$

Choosing  $q = I_{\mathcal{D}}p^\dagger - p \in X_{\mathcal{D},\Omega,\Gamma_n}$  yields,

$$0 \leq a_{\mathcal{D}}(I_{\mathcal{D}}p^\dagger - p, I_{\mathcal{D}}p^\dagger - p) \leq \|I_{\mathcal{D}}p^\dagger - p\|_{\mathcal{D}} \left( W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla\zeta - \mathbf{g}_1) + S_{\mathcal{D}}(p^\dagger) C_S \right). \quad (2.24)$$

We now have to express a lower bound of  $a_{\mathcal{D}}(I_{\mathcal{D}}p^\dagger - p, I_{\mathcal{D}}p^\dagger - p)$  with respect to  $\|I_{\mathcal{D}}p^\dagger - p\|_{\mathcal{D}}^2$ . To do so, it is possible to use the same tricks as in the proof of Lemma 2.2.2. And thus by mimicking at the discrete level (2.10), we obtain,

$$a_{\mathcal{D}}(I_{\mathcal{D}}p^\dagger - p, I_{\mathcal{D}}p^\dagger - p) \geq \underbrace{\frac{\alpha^2}{2} \min \left( \frac{\underline{H}}{2\alpha^2 + 2\|\nabla\zeta\|_{L^\infty(\Omega)^d}^2}, \frac{1}{\bar{H}} \right)}_{:= C_{LM} > 0} \|I_{\mathcal{D}}p^\dagger - p\|_{\mathcal{D}}^2.$$

Finally thanks to (2.24), we obtain,

$$\|I_{\mathcal{D}}p^\dagger - p\|_{\mathcal{D}} \leq \frac{W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla\zeta - \mathbf{g}_1) + S_{\mathcal{D}}(p^\dagger) C_S}{C_{LM}}. \quad (2.25)$$



Since there holds,

$$\|p^\dagger - \Pi_{\mathcal{D}} p\|_{L^2(\Omega)} + \|\nabla p^\dagger - \nabla_{\mathcal{D}} p\|_{L^2(\Omega)^d} \leq \underbrace{\|p^\dagger - \Pi_{\mathcal{D}} I_{\mathcal{D}} p^\dagger\|_{L^2(\Omega)} + \|\nabla p^\dagger - \nabla_{\mathcal{D}} I_{\mathcal{D}} p^\dagger\|_{L^2(\Omega)^d}}_{= S_{\mathcal{D}}(p^\dagger)} + \|I_{\mathcal{D}} p^\dagger - p\|_{\mathcal{D}},$$

we obtain the estimate (2.20). ■

In [22], the obtained estimates depend on the quantity  $C_{\mathcal{D}}$ , introduced in (2.12), and related to the GD-Coercivity given by Definition 2.2.7. It is not the case here, this is due to the definition of the discrete norm  $\|\cdot\|_{\mathcal{D}}$  which introduces a straightforward bound of  $\|\Pi_{\mathcal{D}} \cdot\|_{L^2(\Omega)}$  with respect to  $\|\cdot\|_{\mathcal{D}}$ . The two last GD properties, i.e. GD-Coercivity and the GD-Trace-Consistency, are required to establish an estimate on the trace operator.

**Proposition 2.2.13** (*Error estimate on the trace operator, Conforming formulation*)

Under Assumptions (2.5), let  $p^\dagger \in H_d^1(\Omega)$  be the solution of Problem (2.8). Let  $\mathcal{D}$  be a GD in the sense of Definition 2.2.4 and  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  be the solution of the GS (2.11) given by Theorem 2.2.12; this solution satisfies the following inequality :

$$\|\gamma(p^\dagger) - \mathbb{T}_{\mathcal{D}, \Gamma_n} p\|_{L^2(\Gamma_n)} \leq C_{\mathcal{D}} \frac{W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla \zeta - \mathbf{g}_1) + \bar{S}_{\mathcal{D}}(p^\dagger) C_S}{C_{LM}} + \bar{S}_{\mathcal{D}}(p^\dagger), \quad (2.26)$$

where  $C_{\mathcal{D}}$ ,  $\bar{S}_{\mathcal{D}}$  and  $W_{\mathcal{D}}$  are respectively defined by the GD-Coercivity (2.12), the GD-Trace-Consistency (2.15) and the GD-Limit-Conformity (2.18), moreover,  $C_S$  and  $C_{LM}$  only depend on the data of the problem (see (2.5)).

**Proof:** We define,

$$\bar{I}_{\mathcal{D}} p^\dagger = \operatorname{argmin}_{s \in X_{\mathcal{D}, \Omega, \Gamma_n}} (\|\Pi_{\mathcal{D}} s - p^\dagger\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D}, \Gamma_n} s - \gamma(p^\dagger)\|_{L^2(\Gamma_n)} + \|\nabla_{\mathcal{D}} s - \nabla p^\dagger\|_{L^2(\Omega)^d}).$$

We can remark, by using (2.15), that

$$\|\Pi_{\mathcal{D}} \bar{I}_{\mathcal{D}} p^\dagger - p^\dagger\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D}, \Gamma_n} \bar{I}_{\mathcal{D}} p^\dagger - \gamma(p^\dagger)\|_{L^2(\Gamma_n)} + \|\nabla_{\mathcal{D}} \bar{I}_{\mathcal{D}} p^\dagger - \nabla p^\dagger\|_{L^2(\Omega)^d} = \bar{S}_{\mathcal{D}}(p^\dagger).$$

Moreover, by using (2.12),

$$\begin{aligned} \|\mathbb{T}_{\mathcal{D}, \Gamma_n} p - \gamma(p^\dagger)\|_{L^2(\Gamma_n)} &\leq \underbrace{\|\mathbb{T}_{\mathcal{D}, \Gamma_n} (p - \bar{I}_{\mathcal{D}} p^\dagger)\|_{L^2(\Gamma_n)}}_{\leq C_{\mathcal{D}} \|p - \bar{I}_{\mathcal{D}} p^\dagger\|_{\mathcal{D}}} + \underbrace{\|\mathbb{T}_{\mathcal{D}, \Gamma_n} \bar{I}_{\mathcal{D}} p^\dagger - \gamma(p^\dagger)\|_{L^2(\Gamma_n)}}_{\leq \bar{S}_{\mathcal{D}}(p^\dagger)}. \end{aligned}$$

Finally, by mimicking (2.25), we can prove that,

$$\|p - \bar{I}_{\mathcal{D}} p^\dagger\|_{\mathcal{D}} \leq \frac{W_{\mathcal{D}}(H\nabla p^\dagger + p^\dagger \nabla \zeta - \mathbf{g}_1) + \bar{S}_{\mathcal{D}}(p^\dagger) C_S}{C_{LM}},$$

which enables us to conclude. ■

**Corollary 2.2.14 (Convergence, Conforming formulation)**

Under Assumptions (2.5), let  $(D_m)_{m \in \mathbb{N}}$  be a sequence of GDs in the sense of Definition 2.2.4, which is coercive, consistent, trace-consistent and limit-conforming in the sense of Definitions 2.2.7, 2.2.9, 2.2.10 and 2.2.11.

Then, for any  $m \in \mathbb{N}$ , there exists a unique solution  $p_m \in X_{\mathcal{D}_m, \Omega, \Gamma_n}$  to the gradient scheme (2.11) and, if  $p^\dagger$  is the solution of the weak problem (2.8) then, as  $m \rightarrow \infty$ ,  $\Pi_{\mathcal{D}_m} p_m$  converges to  $p^\dagger$  in  $L^2(\Omega)$ . Similarly,  $\nabla_{\mathcal{D}_m} p_m$  converges to  $\nabla p^\dagger$  in  $L^2(\Omega)^d$  and  $\mathbb{T}_{\mathcal{D}_m, \Gamma_n} p_m$  converges to  $\gamma(p^\dagger)$  in  $L^2(\Gamma_n)$ .

**Remark 2.2.15 (Velocity reconstruction, Conforming formulation)**

Under the same assumptions and notation of the previous corollary, we can mimic the relation (2.6) in order to define a discrete velocity  $\mathbf{u}_m = \begin{pmatrix} \mathbf{g}_1 - H \nabla_{\mathcal{D}_m} p_m - \Pi_{\mathcal{D}_m} p_m \nabla \zeta \\ \mathbf{g}_2 + \alpha \Pi_{\mathcal{D}_m} p_m \end{pmatrix}$  which converges to  $\mathbf{u}^\dagger$  in  $L^2(\Omega)^{d+1}$  as  $m \rightarrow \infty$ .

**2.2.6 The  $\mathbb{P}_1$  conforming finite elements**

The  $\mathbb{P}_1$  finite elements method is a particular conforming Galerkin method and thus is encompassed in the GDM framework (see section 8.1 in [22]). We briefly discuss here why the  $\mathbb{P}_1$  conforming finite elements method can be seen as a GDM in the sense of Definition 2.2.4. We consider  $d = 2$ . Let  $\mathfrak{T} = (\mathcal{M}, \mathcal{F}, \mathcal{C}, \mathcal{V})$  be a conforming simplicial mesh of  $\Omega$  in the sense of Definition 7.4 in [22] where  $\mathcal{M}$  is the set of *cells*,  $\mathcal{F}$  of *faces*,  $\mathcal{C}$  of the *cell centers* and  $\mathcal{V}$  of the *vertices*. A conforming simplicial mesh of  $\Omega$  is a polytopal mesh in the sense of Definition 7.2 in [22] (see also Definition 1.50 in [27]) such that for each  $K \in \mathcal{M}$  we have  $\text{Card}(\mathcal{F}_K) = \text{Card}(\mathcal{V}_K) = 3$ .

We introduce  $N = \text{card}(\mathcal{V})$  and the family  $(\phi_i)_{i=1}^N$  of  $\mathbb{P}_1$  finite element shape functions associated to the vertices. The  $\mathbb{P}_1$  conforming finite elements method is a GDM in the sense of Definition 2.2.4 by using :

- $X_{\mathcal{D}} = X_{\mathcal{D}, \Omega, \Gamma_n} \oplus X_{\mathcal{D}, \Gamma_d}$  with  $X_{\mathcal{D}} = \mathbb{R}^N$  and  $X_{\mathcal{D}, \Gamma_d} = \mathbb{R}^{N_d} = \{0\}$  where  $N_d$  is the number of vertices belonging to  $\Gamma_d$ ,
- $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  is defined, for all  $\alpha := (\alpha_i)_{i=1}^N \in X_{\mathcal{D}}$ , by  $\Pi_{\mathcal{D}}(\alpha) = \sum_{i=1}^N \alpha_i \phi_i$  with  $\alpha_i = 0$  if it is associated to a node on  $\Gamma_d$ ; we can remark that  $\Pi_{\mathcal{D}}(\alpha) \in H_d^1(\Omega)$  which allows us to define the last two operators as follows,
- $\mathbb{T}_{\mathcal{D}, \Gamma_n} : X_{\mathcal{D}} \rightarrow L^2(\Gamma_n)$  is defined by  $\mathbb{T}_{\mathcal{D}, \Gamma_n}(\cdot) = \gamma(\Pi_{\mathcal{D}}(\cdot))|_{\Gamma_n}$ ,
- $\nabla_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^d$  is defined by  $\nabla_{\mathcal{D}}(\cdot) = \nabla(\Pi_{\mathcal{D}}(\cdot))$ .

Using these definitions, it is straightforward to show that  $\|\cdot\|_{\mathcal{D}}$  is a norm on  $X_{\mathcal{D}, \Omega, \Gamma_n}$  : if  $\|v\|_{\mathcal{D}} = 0$  then  $\Pi_{\mathcal{D}} v = 0$ . Then, for all  $s \in \mathcal{V}$ ,  $v_s = \Pi_{\mathcal{D}} v(s) = 0$ , which shows that  $v = 0$ .

We now list the arguments that allow us to say that, with the above definitions, the GD properties are satisfied :

- the **GD-Coercivity** is a direct application of the continuity of  $\gamma$ , the continuous trace operator; this gives a uniform estimate on  $C_{\mathcal{D}}$  (defined by (2.2.7)) depending only on  $\Omega$  and the chosen  $L^p$  space, here  $L^2$ .

- the **GD-Consistency** and **GD-Trace-Consistency** are deduced from results on the interpolation error, see Corollary 3.23 in [27], and the continuity of  $\gamma$ . We can underline that an assumption on the mesh regularity is here needed as mentioned in Definition 1.107 of [27],
- the choice of discrete trace operator leads the method to be a conforming one, thus  $W_{\mathcal{D}}(\cdot) = 0$  (2.18), so we have the **GD-Limit-Conformity**.

We denote by  $h_{\mathcal{D}}$  the *mesh size*, i.e. the maximal diameter of the *cells*. Then, if we assume that the solution  $p^\dagger$  is more regular and belongs to  $H_d^2(\Omega)$ , we can express the quantities  $S_{\mathcal{D}}(\cdot)$  and  $\bar{S}_{\mathcal{D}}(\cdot)$  as  $O(h_{\mathcal{D}})$ , see [27, Theorem 3.16]. This result on  $\bar{S}_{\mathcal{D}}(\cdot)$  is more complex to obtain for non-conforming methods.

## 2.3 The Abstract Gradient Discretisation Method on *special* operators

In the previous section, we studied the strong problem (2.2)-(2.4) through the classic GDM framework based on the weak formulation of the problem with the *classic* operators  $\nabla$  and  $\text{div}$ . The Abstract Gradient Discretisation Method (AGDM) enables to mimic, at the discrete level, the *special* operators  $\nabla_s^\alpha$  and  $\text{div}_s^\alpha$  and obtain suitable error estimates.

### 2.3.1 Weak formulation

First, we need to ensure that the *special* operators  $\nabla_s^\alpha$  and  $\text{div}_s^\alpha$  defined by (2.1) fulfil a duality relation.

#### Lemma 2.3.1

Operators  $\nabla_s^\alpha$  and  $\text{div}_s^\alpha$  given by (2.1) satisfy :  $\forall \varphi \in H^1(\Omega)$  and  $\forall \boldsymbol{\omega} \in H_{\text{div}_s^\alpha}(\Omega)$ ,

$$\int_{\Omega} \nabla_s^\alpha \varphi \cdot \boldsymbol{\omega} \, d\mathbf{x} = - \int_{\Omega} \varphi \, \text{div}_s^\alpha \boldsymbol{\omega} \, d\mathbf{x} + \int_{\Gamma} H \boldsymbol{\omega} \cdot \mathbf{n}_s \, \gamma(\varphi) \, d\sigma(\mathbf{x}), \quad (2.27)$$

with  $H_{\text{div}_s^\alpha}(\Omega) = \{ \boldsymbol{\omega} \in L^2(\Omega)^{d+1} \mid \text{div}_s^\alpha(\boldsymbol{\omega}) \in L^2(\Omega), \boldsymbol{\omega} \cdot \mathbf{n}_s \in L^2(\Gamma) \}$ .

**Proof:** Using (2.1) and Green formula we get :  $\forall \varphi \in H^1(\Omega)$  and  $\forall \boldsymbol{\omega} = (\mathbf{v}, w)^T \in H_{\text{div}_s^\alpha}(\Omega)$ ,

$$\begin{aligned} \int_{\Omega} \nabla_s^\alpha \varphi \cdot \boldsymbol{\omega} \, d\mathbf{x} &= \int_{\Omega} (H \nabla \varphi + \varphi \nabla \zeta) \cdot \mathbf{v} - \alpha \varphi w \, d\mathbf{x} \\ &= - \int_{\Omega} \varphi \, \text{div}(H \mathbf{v}) - \varphi (\nabla \zeta \cdot \mathbf{v} - \alpha w) \, d\mathbf{x} + \int_{\Gamma} H \mathbf{v} \cdot \mathbf{n}_{\Gamma} \, \gamma(\varphi) \, d\sigma(\mathbf{x}) \\ &= - \int_{\Omega} \varphi \, \text{div}_s^\alpha(H \mathbf{v}) \, d\mathbf{x} + \int_{\Gamma} H \boldsymbol{\omega} \cdot \mathbf{n}_s \, \gamma(\varphi) \, d\sigma(\mathbf{x}). \end{aligned}$$

■

Contrary to subsection 2.2.1, we keep the *special* operators in the conforming formulation of Problem (2.2)-(2.4). Thus, the strong problem can be read as :

Find  $p^\dagger : \Omega \rightarrow \mathbb{R}$  the non hydrostatic pressure such that,

$$-\operatorname{div}_s^\alpha \left( \frac{1}{H} \nabla_s^\alpha p^\dagger \right) = f - \operatorname{div}_s^\alpha \left( \frac{1}{H} \mathbf{g} \right) \text{ in } \Omega, \quad (2.28a)$$

$$\nabla_s^\alpha p^\dagger \cdot \mathbf{n}_s = \mathbf{g} \cdot \mathbf{n}_s - \phi \text{ on } \Gamma_n, \quad (2.28b)$$

$$p^\dagger = 0 \text{ on } \Gamma_d, \quad (2.28c)$$

where we still have  $\mathbf{u}^\dagger = \frac{\mathbf{g} - \nabla_s^\alpha p^\dagger}{H}$  as in (2.6).

We now write a weak formulation of Problem (2.28) thanks to relation (2.27) in order to introduce the associated Abstract Gradient Scheme (AGS). A function  $p^\dagger$  is said to be a weak solution of Problem (2.28) if the following holds :

Find  $p^\dagger \in H_d^1(\Omega)$  such that  $\forall q^\dagger \in H_d^1(\Omega)$ ,

$$\int_\Omega \frac{\nabla_s^\alpha p^\dagger \cdot \nabla_s^\alpha q^\dagger}{H} \, d\mathbf{x} = \int_\Omega f q^\dagger \, d\mathbf{x} + \int_{\Gamma_n} -\phi \gamma(q^\dagger) \, d\sigma(\mathbf{x}) + \int_\Omega \frac{\mathbf{g} \cdot \nabla_s^\alpha q^\dagger}{H} \, d\mathbf{x}. \quad (2.29)$$

Equation (2.29) is Equation (2.8) using the operator  $\nabla_s^\alpha$  (2.1a). Thus existence and uniqueness of a solution are a direct consequence of Lemma 2.2.2.

### 2.3.2 Abstract Gradient Scheme

The approximation of Problem (2.29) relies on four discrete objects  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, \mathbf{G}_{\mathcal{D}})$  forming altogether an Abstract Gradient Discretisation.

#### Definition 2.3.2 (*Abstract Gradient Discretisation*)

Under the assumption (2.5a), an abstract gradient discretisation is defined by  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, \mathbf{G}_{\mathcal{D}})$ , where :

1. the set of discrete unknowns  $X_{\mathcal{D}} = X_{\mathcal{D}, \Omega, \Gamma_n} \oplus X_{\mathcal{D}, \Gamma_d}$  is the direct sum of two finite dimensional vector spaces on  $\mathbb{R}$ , corresponding respectively to the unknowns in  $\Omega$  and on  $\Gamma_n$ , and to the unknowns on  $\Gamma_d$  ( $X_{\mathcal{D}, \Gamma_d} = \{0\}$  due to (2.4)),
2. the function reconstruction  $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  is a linear mapping that reconstructs, from an element of  $X_{\mathcal{D}}$ , an element in  $L^2(\Omega)$ ,
3. the trace reconstruction  $\mathbb{T}_{\mathcal{D}, \Gamma_n} : X_{\mathcal{D}} \rightarrow L^2(\Gamma_n)$  is a linear mapping that reconstructs, from an element of  $X_{\mathcal{D}}$ , an element in  $\Gamma_n$ ,
4. the special gradient reconstruction  $\mathbf{G}_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^{d+1}$  is a linear mapping that reconstructs, from an element of  $X_{\mathcal{D}}$ , an element in  $L^2(\Omega)^{d+1}$ ,
5. the mapping  $\mathbf{G}_{\mathcal{D}}$  is such that the following quantity is a norm on  $X_{\mathcal{D}, \Omega, \Gamma_n} : \|q\|_{\mathcal{D}} := \|\mathbf{G}_{\mathcal{D}} q\|_{L^2(\Omega)^{d+1}}$ .

Now, we define the Abstract Gradient Scheme associated to the weak problem by replacing the continuous spaces and the *special* operators by discrete counterparts.

**Definition 2.3.3 (AGS, Conforming formulation)**

Let  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, \mathbb{T}_{\mathcal{D}, \Gamma_n}, G_{\mathcal{D}})$  be an AGD in the sense of Definition 2.3.2. The related abstract gradient scheme for Problem (2.29) is defined by :

Find  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  such that for any  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$ ,

$$\int_{\Omega} \frac{G_{\mathcal{D}} p \cdot G_{\mathcal{D}} q}{H} \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}} q \, d\mathbf{x} + \int_{\Gamma_n} -\phi \mathbb{T}_{\mathcal{D}, \Gamma_n} q \, d\sigma(\mathbf{x}) + \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}} q}{H} \, d\mathbf{x}. \quad (2.30)$$

**2.3.3 Key properties of the considered Abstract Gradient Discretisation**

Once again, after Definition 2.3.2 of an AGD, we introduce its key properties that ensure the convergence of the corresponding AGS (2.30). These properties are defined in the following and named here AGD-Coercivity, AGD-Consistency and AGD-Limit-Conformity.

**Definition 2.3.4 (AGD-Coercivity)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.3.2, let  $C_{\mathcal{D}}$  be,

$$C_{\mathcal{D}} = \max_{q \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \frac{\|\Pi_{\mathcal{D}} q\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D}, \Gamma_n} q\|_{L^2(\Gamma_n)}}{\|q\|_{\mathcal{D}}}. \quad (2.31)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **coercive** if there exists  $C_P \in \mathbb{R}_+$  such that  $C_{\mathcal{D}_m} \leq C_P$  for all  $m \in \mathbb{N}$ .

**Remark 2.3.5 (About the definition of the coercivity)**

In comparison with the previous section, the norm  $\|\cdot\|_{\mathcal{D}}$  includes a discrete version of the *special* gradient and thus includes the term in  $\alpha$ . This is why, we keep here the standard definition of the coercivity regarding the literature on the GDM.

**Definition 2.3.6 (AGD-Consistency)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.3.2, let  $S_{\mathcal{D}} : H_d^1(\Omega) \rightarrow [0, +\infty)$  be given by,

$$\forall \varphi \in H_d^1(\Omega), \quad S_{\mathcal{D}}(\varphi) = \min_{q \in X_{\mathcal{D}, \Omega, \Gamma_n}} \left( \|\Pi_{\mathcal{D}} q - \varphi\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D}, \Gamma_n} q - \gamma(\varphi)\|_{L^2(\Gamma_n)} + \|G_{\mathcal{D}} q - \nabla_s^\alpha \varphi\|_{L^2(\Omega)^{d+1}} \right). \quad (2.32)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **consistent** if,

$$\forall \varphi \in H_d^1(\Omega), \quad \lim_{m \rightarrow \infty} S_{\mathcal{D}_m}(\varphi) = 0. \quad (2.33)$$

The AGD-consistency can still be understood as interpolation errors but we can underline that the last term is related to the *special* gradient operator.

**Definition 2.3.7 (AGD-Limit-Conformity)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.3.2, let  $W_{\mathcal{D}} : H_{\text{div}_s^\alpha}(\Omega) \rightarrow [0, +\infty)$  be given by,

$$\forall \varphi \in H_{\text{div}_s^\alpha}(\Omega), W_{\mathcal{D}}(\varphi) = \sup_{q \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \frac{\left| \int_{\Omega} G_{\mathcal{D}} q \cdot \varphi + \Pi_{\mathcal{D}} q \operatorname{div}_s^\alpha \varphi \, d\mathbf{x} - \int_{\Gamma_n} H \varphi \cdot \mathbf{n}_s \, \mathbb{T}_{\mathcal{D}, \Gamma_n} q \, d\sigma(\mathbf{x}) \right|}{\|q\|_{\mathcal{D}}}. \quad (2.34)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **limit-conforming** if,

$$\forall \varphi \in H_{\text{div}_s^\alpha}(\Omega), \lim_{m \rightarrow \infty} W_{\mathcal{D}_m}(\varphi) = 0.$$

The AGD-Limit-Conformity is needed in order to deal with non-conforming methods. The aim of this last property is to approach the duality relation given by Lemma (2.27).

**2.3.4 Error estimates**

In order to establish error estimates, we can apply Theorem 5.2 in [24] on the studied problem. The proof is close to the one done for Theorem 2.2.12 and Proposition 2.2.13 of the previous section.

**Theorem 2.3.8 (Error estimates of an AGS, Conforming formulation)**

Under Assumptions (2.5), let  $p^\dagger \in H_d^1(\Omega)$  be the solution of Problem (2.29). Let  $\mathcal{D}$  be an AGD in the sense of Definition 2.3.2. Then there exists one and only one  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  solution to the AGS given by (2.30). This solution satisfies the following inequalities :

$$\|\nabla_s^\alpha p^\dagger - G_{\mathcal{D}} p\|_{L^2(\Omega)^{d+1}} \leq \overline{H} \left[ W_{\mathcal{D}} \left( \frac{\nabla_s^\alpha p^\dagger - \mathbf{g}}{H} \right) + \left( \frac{1}{\underline{H}} + \frac{1}{\overline{H}} \right) S_{\mathcal{D}}(p^\dagger) \right], \quad (2.35)$$

$$\|p^\dagger - \Pi_{\mathcal{D}} p\|_{L^2(\Omega)} + \|\gamma(p^\dagger) - \mathbb{T}_{\mathcal{D}, \Gamma_n} p\|_{L^2(\Gamma_n)} \leq \overline{H} \left[ C_{\mathcal{D}} W_{\mathcal{D}} \left( \frac{\nabla_s^\alpha p^\dagger - \mathbf{g}}{H} \right) + \left( \frac{C_{\mathcal{D}}}{\underline{H}} + \frac{1}{\overline{H}} \right) S_{\mathcal{D}}(p^\dagger) \right], \quad (2.36)$$

where  $C_{\mathcal{D}}$ ,  $S_{\mathcal{D}}$  and  $W_{\mathcal{D}}$  are respectively related to the reconstruction operators, the consistency defect and the conformity defect, defined by (2.31), (2.32) and (2.34).

**Proof:**

Let us first prove that, if (2.35) holds for any solution  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  to Scheme (2.30), then the solution exists and is unique. By introducing a basis of the space  $X_{\mathcal{D}, \Omega, \Gamma_n}$ , it is easy to show that the Scheme (2.30) can be rewritten under the form of a square linear system denoted  $A_{\mathcal{D}} P_{\mathcal{D}} = F_{\mathcal{D}}$  where the unknown vector  $P_{\mathcal{D}}$  is made of the coordinates of  $p$  in the mentioned basis. By choosing  $f = 0$ ,  $\mathbf{g} = 0$  and  $\phi = 0$  the right hand side of the continuous problem (2.29) vanishes and we also have  $F_{\mathcal{D}} = 0$ . In this case the unique solution of (2.29) is  $p^\dagger = 0$ . Then the right hand side of the inequality of (2.35) is equal to zero, thus we get that any solution  $p$  to the scheme satisfies  $\|p\|_{\mathcal{D}} = 0$ , which leads to  $p = 0$ . Therefore the kernel of the matrix  $A_{\mathcal{D}}$  is  $\{0\}$  which is equivalent to say that the linear system has a unique solution.

By the definition of  $W_{\mathcal{D}}$  (2.34), we have  $\forall q \in X_{\mathcal{D},\Omega,\Gamma_n}$  and  $\forall \varphi \in H_{\text{div}_s^\alpha}(\Omega)$ ,

$$\left| \int_{\Omega} G_{\mathcal{D}}q \cdot \varphi + \Pi_{\mathcal{D}}q \operatorname{div}_s^\alpha \varphi \, d\mathbf{x} - \int_{\Gamma_n} H\varphi \cdot \mathbf{n}_s \mathbb{T}_{\mathcal{D},\Gamma_n}q \, d\sigma(\mathbf{x}) \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(\varphi).$$

We take  $\varphi = \frac{1}{H} (\nabla_s^\alpha p^\dagger - \mathbf{g})$  where  $p^\dagger$  is the solution of (2.29). We can easily prove that  $\varphi \in H_{\text{div}_s^\alpha}(\Omega)$  and, in the distribution sense,  $-\operatorname{div}_s^\alpha \varphi = f$  on  $\Omega$  and  $H\varphi \cdot \mathbf{n}_s = -\phi$  on  $\Gamma_n$ . Thus we have,

$$\left| \int_{\Omega} \frac{1}{H} G_{\mathcal{D}}q \cdot (\mathbf{g} - \nabla_s^\alpha p^\dagger) + \Pi_{\mathcal{D}}q f \, d\mathbf{x} - \int_{\Gamma_n} \phi \mathbb{T}_{\mathcal{D},\Gamma_n}q \, d\sigma(\mathbf{x}) \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(\varphi).$$

Moreover, by using (2.30), we obtain,

$$\left| \int_{\Omega} \frac{1}{H} G_{\mathcal{D}}q \cdot (G_{\mathcal{D}}p - \nabla_s^\alpha p^\dagger) \, d\mathbf{x} \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(\varphi).$$

We now define,

$$I_{\mathcal{D}}p^\dagger = \operatorname{argmin}_{w \in X_{\mathcal{D},\Omega,\Gamma_n}} \left( \|\Pi_{\mathcal{D}}w - p^\dagger\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D},\Gamma_n}w - \gamma(\varphi)\|_{L^2(\Gamma_n)} + \|G_{\mathcal{D}}w - \nabla_s^\alpha p^\dagger\|_{L^2(\Omega)^{d+1}} \right).$$

Thus we can write that,

$$\left| \int_{\Omega} \frac{1}{H} G_{\mathcal{D}}q \cdot (G_{\mathcal{D}}p - G_{\mathcal{D}}I_{\mathcal{D}}p^\dagger) \, d\mathbf{x} \right| \leq \|q\|_{\mathcal{D}} W_{\mathcal{D}}(\varphi) + \left| \int_{\Omega} \frac{1}{H} G_{\mathcal{D}}q \cdot (\nabla_s^\alpha p^\dagger - G_{\mathcal{D}}I_{\mathcal{D}}p^\dagger) \, d\mathbf{x} \right|,$$

which implies, by the definition of  $S_{\mathcal{D}}$  (2.32), and the assumption (2.5b),

$$\left| \int_{\Omega} \frac{1}{H} G_{\mathcal{D}}q \cdot (G_{\mathcal{D}}p - G_{\mathcal{D}}I_{\mathcal{D}}p^\dagger) \, d\mathbf{x} \right| \leq \|q\|_{\mathcal{D}} \left( W_{\mathcal{D}}(\varphi) + \frac{S_{\mathcal{D}}(p^\dagger)}{H} \right).$$

We take  $q = p - I_{\mathcal{D}}p^\dagger \in X_{\mathcal{D},\Omega,\Gamma_n}$ , then by using the definition of  $\|\cdot\|_{\mathcal{D}}$ , the Cauchy-Schwarz inequality and assumption (2.5b), we can write,

$$\frac{1}{H} \|G_{\mathcal{D}}(p - I_{\mathcal{D}}p^\dagger)\|_{L^2(\Omega)^{d+1}}^2 \leq \|G_{\mathcal{D}}(p - I_{\mathcal{D}}p^\dagger)\|_{L^2(\Omega)^{d+1}} \left( W_{\mathcal{D}}(\varphi) + \frac{S_{\mathcal{D}}(p^\dagger)}{H} \right).$$

By using the triangular inequality and the definition of  $S_{\mathcal{D}}$  (2.32), we have,

$$\|G_{\mathcal{D}}p - \nabla_s^\alpha p^\dagger\|_{L^2(\Omega)^{d+1}} \leq \overline{H} \left( W_{\mathcal{D}}(\varphi) + \frac{S_{\mathcal{D}}(p^\dagger)}{H} \right) + S_{\mathcal{D}}(p^\dagger)$$

In order to write an estimate on the solution, we can also use Definition 2.3.4 of the AGD-Coercivity,

$$\begin{aligned} \|\Pi_{\mathcal{D}}p - p^\dagger\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D},\Gamma_n}p - \gamma(p^\dagger)\|_{L^2(\Gamma_n)} &\leq \overbrace{\|\Pi_{\mathcal{D}}(p - I_{\mathcal{D}}p^\dagger)\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D},\Gamma_n}(p - I_{\mathcal{D}}p^\dagger)\|_{L^2(\Gamma_n)}}^{\leq C_{\mathcal{D}}\|G_{\mathcal{D}}(p - I_{\mathcal{D}}p^\dagger)\|_{L^2(\Omega)^{d+1}}} \\ &\quad + \underbrace{\|\Pi_{\mathcal{D}}I_{\mathcal{D}}p^\dagger - p^\dagger\|_{L^2(\Omega)} + \|\mathbb{T}_{\mathcal{D},\Gamma_n}I_{\mathcal{D}}p^\dagger - \gamma(p^\dagger)\|_{L^2(\Gamma_n)}}_{\leq S_{\mathcal{D}}(p^\dagger)}. \end{aligned} \quad \blacksquare$$

### Remark 2.3.9

We have convergence results similar to Corollary 2.2.14 and Remark 2.2.15.

### 2.3.5 The $\mathbb{P}_1$ conforming finite elements

We discuss here why the  $\mathbb{P}_1$  conforming finite elements method can also be seen as an Abstract Gradient Discretisation method in the sense of Definition 2.3.2. To do so we consider the same notations for the triangular mesh as in section 2.2.6. The vector of unknowns  $X_{\mathcal{D}}$ , the function reconstruction  $\Pi_{\mathcal{D}}$  and the trace reconstruction  $\mathbb{T}_{\mathcal{D},\Gamma_n}$  are also the same as in section 2.2.6. The main point is the definition of the discrete gradient.

Thus the  $\mathbb{P}_1$  conforming finite elements method is a AGDM in the sense of Definition 2.3.2 by using :

- $X_{\mathcal{D}} = X_{\mathcal{D},\Omega,\Gamma_n} \oplus X_{\mathcal{D},\Gamma_d}$  with  $X_{\mathcal{D}} = \mathbb{R}^N$  and  $X_{\mathcal{D},\Gamma_d} = \mathbb{R}^{N_d} = \{0\}$  where  $N_d$  is the number of vertices belonging to  $\Gamma_d$ ,
- $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  is defined, for all  $\alpha := (\alpha_i)_{i=1}^N \in X_{\mathcal{D}}$ , by  $\Pi_{\mathcal{D}}(\alpha) = \sum_{i=1}^N \alpha_i \phi_i$  with  $\alpha_i = 0$  if it is associated to a node on  $\Gamma_d$  ( $\phi_i$  still denote the  $\mathbb{P}_1$  finite element shape function at vertex  $i$ ); we still have  $\Pi_{\mathcal{D}}(\alpha) \in H_d^1(\Omega)$  which allows us to define the last two operators as follows,
- $\mathbb{T}_{\mathcal{D},\Gamma_n} : X_{\mathcal{D}} \rightarrow L^2(\Gamma_n)$  is defined by  $\mathbb{T}_{\mathcal{D},\Gamma_n}(\cdot) = \gamma(\Pi_{\mathcal{D}}(\cdot))|_{\Gamma_n}$ ,
- $G_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^{d+1}$  is defined by  $G_{\mathcal{D}}(\cdot) = \nabla_s^\alpha(\Pi_{\mathcal{D}}(\cdot)) = (H\nabla(\Pi_{\mathcal{D}}(\cdot)) + \Pi_{\mathcal{D}}(\cdot)\nabla\zeta, -\alpha\Pi_{\mathcal{D}}(\cdot))^T$ .

Using these definitions, we have for any  $v \in X_{\mathcal{D},\Omega,\Gamma_n}$ ,

$$\|v\|_{\mathcal{D}}^2 = \|H\nabla(\Pi_{\mathcal{D}}v) + \Pi_{\mathcal{D}}v\nabla\zeta\|_{L^2(\Omega)^d}^2 + \alpha^2\|\Pi_{\mathcal{D}}v\|_{L^2(\Omega)}^2, \quad (2.37)$$

thus it is easy to show that  $\|\cdot\|_{\mathcal{D}}$  is a norm on  $X_{\mathcal{D},\Omega,\Gamma_n}$  : if  $\|v\|_{\mathcal{D}} = 0$ , since  $\alpha \neq 0$ , then  $\Pi_{\mathcal{D}}v = 0$ , which implies that  $v = 0$ .

We now study the arguments that allow us to say that, with the above definitions, the AGD properties are satisfied :

- **AGD-Coercivity.** Thanks to (2.37) we clearly have  $\|\Pi_{\mathcal{D}}v\|_{L^2(\Omega)} \leq \frac{1}{\alpha}\|v\|_{\mathcal{D}}$ . Moreover we have,

$$\|H\nabla(\Pi_{\mathcal{D}}v)\|_{L^2(\Omega)^d} \leq \|v\|_{\mathcal{D}} + \|\nabla\zeta\|_{L^\infty(\Omega)^d}\|\Pi_{\mathcal{D}}v\|_{L^2(\Omega)}.$$

Thus, using assumption (2.5b), there exists  $C > 0$  independent of the mesh, such that

$\|\nabla(\Pi_{\mathcal{D}}v)\|_{L^2(\Omega)^d} \leq C\|v\|_{\mathcal{D}}$ . Finally using the continuity of the trace  $\gamma$ , we obtain the wanted property.

- **AGD-Consistency.** This property is deduced from results on the interpolation error, see Corollary 3.23 in [27], the continuity of  $\gamma$ , and the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c). An assumption on the mesh regularity is still needed here.
- **AGD-Limit-Conformity.** Since we deals with a conforming method, we can used the duality relation on the *special* operators, given by (2.27), to obtain  $W_{\mathcal{D}}(\cdot) = 0$  (2.18).

#### Remark 2.3.10 (*Non-conforming construction of $G_{\mathcal{D}}$* )

*It is possible to make a non-conforming choice to build  $G_{\mathcal{D}}$ . Indeed, if we have an other function reconstruction  $\tilde{\Pi}_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  (for instance the  $\mathbb{P}_1$  mass-lumping) we can define  $G_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^{d+1}$  by  $G_{\mathcal{D}}(\cdot) = \nabla_s^\alpha(\tilde{\Pi}_{\mathcal{D}}(\cdot)) = (H\nabla(\tilde{\Pi}_{\mathcal{D}}(\cdot)) + \tilde{\Pi}_{\mathcal{D}}(\cdot)\nabla\zeta, -\alpha\tilde{\Pi}_{\mathcal{D}}(\cdot))^T$ . If the operator  $\tilde{\Pi}_{\mathcal{D}}$  satisfies  $\|\tilde{\Pi}_{\mathcal{D}}(\cdot) - \Pi_{\mathcal{D}}(\cdot)\|_{L^2(\Omega)} \leq h^p C \|G_{\mathcal{D}}(\cdot)\|_{L^2(\Omega)^{d+1}}$ , with  $C > 0$  independent of the mesh and  $p > 0$ , the AGD-properties are satisfied.*



## 2.4 Analyse of a mixed mass-lumped scheme

As mentioned in the introduction, one goal of this work is to analyse the mixed mass-lumped schemes proposed in [2]. In this cited work, the authors introduce a mixed weak formulation with a lifting on the velocity. At the discrete level, the authors use a lumping strategy in order to reduce the numerical cost of the schemes. This section plugs these schemes into the AGD framework.

### 2.4.1 A weak mixed formulation with lifting

From problem (2.2)-(2.4), we write a weak mixed formulation where the functions  $p^\dagger$  and  $\mathbf{u}^\dagger$  are said to be a weak solution if the following holds :

Find  $(\mathbf{u}^\dagger, p^\dagger) \in H_{\text{div}_s^\alpha}^\phi(\Omega) \times L^2(\Omega)$  such that,

$$\forall \boldsymbol{\omega}^\dagger \in H_{\text{div}_s^\alpha}^0(\Omega), \int_{\Omega} H\mathbf{u}^\dagger \cdot \boldsymbol{\omega}^\dagger - p^\dagger \text{div}_s^\alpha \boldsymbol{\omega}^\dagger \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \boldsymbol{\omega}^\dagger \, d\mathbf{x}, \quad (2.38a)$$

$$\forall q^\dagger \in L^2(\Omega), \int_{\Omega} q^\dagger \text{div}_s^\alpha \mathbf{u}^\dagger \, d\mathbf{x} = \int_{\Omega} f q^\dagger \, d\mathbf{x}, \quad (2.38b)$$

with  $H_{\text{div}_s^\alpha}^\phi(\Omega) = \{\boldsymbol{\omega} \in L^2(\Omega)^{d+1} \mid \text{div}_s^\alpha(\boldsymbol{\omega}) \in L^2(\Omega), H\boldsymbol{\omega} \cdot \mathbf{n}_s = \phi \text{ on } \Gamma_n\}$ .

Now, as in [2], we introduce a lifting on the velocity  $\mathbf{u}^\dagger$ . To do so, we assume that there exist  $\tilde{\boldsymbol{\omega}} \in H_{\text{div}_s^\alpha}^\phi(\Omega)$ . Then we define  $\tilde{\mathbf{u}} := \mathbf{u}^\dagger - \tilde{\boldsymbol{\omega}}$  and, by construction, we have  $\tilde{\mathbf{u}} \in H_{\text{div}_s^\alpha}^0(\Omega)$ . Then the problem reads :

Find  $(\tilde{\mathbf{u}}, p^\dagger) \in H_{\text{div}_s^\alpha}^0(\Omega) \times L^2(\Omega)$  such that,

$$\forall \boldsymbol{\omega}^\dagger \in H_{\text{div}_s^\alpha}^0(\Omega), \int_{\Omega} H\tilde{\mathbf{u}} \cdot \boldsymbol{\omega}^\dagger - p^\dagger \text{div}_s^\alpha \boldsymbol{\omega}^\dagger \, d\mathbf{x} = \int_{\Omega} (\mathbf{g} - H\tilde{\boldsymbol{\omega}}) \cdot \boldsymbol{\omega}^\dagger \, d\mathbf{x}, \quad (2.39a)$$

$$\forall q^\dagger \in L^2(\Omega), \int_{\Omega} q^\dagger \text{div}_s^\alpha \tilde{\mathbf{u}} \, d\mathbf{x} = \int_{\Omega} q^\dagger (f - \text{div}_s^\alpha \tilde{\boldsymbol{\omega}}) \, d\mathbf{x}. \quad (2.39b)$$

So if we take  $\tilde{\mathbf{g}} := \mathbf{g} - H\tilde{\boldsymbol{\omega}}$  and  $\tilde{f} := f - \text{div}_s^\alpha \tilde{\boldsymbol{\omega}}$ , we recover (2.38) by dropping the tildes on  $f$  and  $\mathbf{g}$ , and by replacing  $\tilde{\mathbf{u}}$  with  $\mathbf{u}^\dagger$ . Thus, in the remaining of this chapter, we replace  $\tilde{f}$ , respectively  $\tilde{\mathbf{g}}$ , by  $f$ , respectively by  $\mathbf{g}$ .

We introduce as in [2], for  $\boldsymbol{\omega} \in H_{\text{div}_s^\alpha}^0(\Omega)$ ,

$$\|\boldsymbol{\omega}\|_{H_{\text{div}_s^\alpha}^0(\Omega)}^2 = \sum_{i=1}^{d+1} \|\boldsymbol{\omega}_i\|_{L^2(\Omega)}^2 + \|\text{div}_s^\alpha \boldsymbol{\omega}\|_{L^2(\Omega)}^2. \quad (2.40)$$

Due to the definition (2.1b), we rewrite, with  $\boldsymbol{\omega} = (\mathbf{v}, w)^T$ ,

$$\|\boldsymbol{\omega}\|_{H_{\text{div}_s^\alpha}^0(\Omega)}^2 = \sum_{i=1}^d \|\mathbf{v}_i\|_{L^2(\Omega)}^2 + \|w\|_{L^2(\Omega)}^2 + \|\text{div}(H\mathbf{v}) - \mathbf{v} \cdot \nabla \zeta + \alpha w\|_{L^2(\Omega)}^2.$$

The lifting on the velocity means that the boundary condition (2.3) can be considered homogeneous on  $\Gamma_n$ . Thus, with  $\mathbf{u}^\dagger \in H_{\text{div}_s^\alpha}^0(\Omega)$ , the weak problem (2.29) involving only the pressure unknown  $p^\dagger$  can be rewritten :

Find  $p^\dagger \in H_d^1(\Omega)$  such that  $\forall q^\dagger \in H_d^1(\Omega)$ ,

$$\int_{\Omega} \frac{\nabla_s^\alpha p^\dagger \cdot \nabla_s^\alpha q^\dagger}{H} \, d\mathbf{x} = \int_{\Omega} f q^\dagger \, d\mathbf{x} + \int_{\Omega} \frac{\mathbf{g} \cdot \nabla_s^\alpha q^\dagger}{H} \, d\mathbf{x}. \quad (2.41)$$

The equivalence between (2.39) and (2.41) can be deduced from the *inf-sup* condition proved in [2, section 3.4.4]. The above formulation is needed in order to introduce the associated AGDM framework.

### 2.4.2 The associated Abstract Gradient Discretisation Method framework

In order to take into account the lifting on the velocity, we need to adapt the AGDM framework previously developed in the sections 2.3.2, 2.3.3 and 2.3.4. The adaption mainly consists in removing the terms related to  $\Gamma_n$ , including the trace reconstruction.

In this way, an Abstract Gradient Discretisation associated to problem (2.41) is defined as :

#### Definition 2.4.1 (AGD, Lifted problem)

Under the assumption (2.5a), an abstract gradient discretisation is defined by  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, G_{\mathcal{D}})$ , where :

1. the set of discrete unknowns  $X_{\mathcal{D}} = X_{\mathcal{D},\Omega,\Gamma_n} \oplus X_{\mathcal{D},\Gamma_d}$  is the direct sum of two finite dimensional vector spaces on  $\mathbb{R}$ , corresponding respectively to the unknowns in  $\Omega$  and on  $\Gamma_n$ , and to the unknowns on  $\Gamma_d$  ( $X_{\mathcal{D},\Gamma_d} = \{0\}$  due to (2.4)),
2. the function reconstruction  $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  is a linear mapping,
3. the special gradient reconstruction  $G_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^{d+1}$  is a linear mapping,
4. the mapping  $G_{\mathcal{D}}$  is such that the following quantity is a norm on  $X_{\mathcal{D},\Omega,\Gamma_n}$  :  $\|q\|_{\mathcal{D}} := \|G_{\mathcal{D}}q\|_{L^2(\Omega)^{d+1}}$ .

Now we can define the Abstract Gradient Scheme associated to the weak problem (2.41).

#### Definition 2.4.2 (AGS, Lifted problem)

Let  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, G_{\mathcal{D}})$  an AGD in the sense of Definition 2.4.1. The related abstract gradient scheme for Problem (2.41) is defined by :

Find  $p \in X_{\mathcal{D},\Omega,\Gamma_n}$  such that for any  $q \in X_{\mathcal{D},\Omega,\Gamma_n}$ ,

$$\int_{\Omega} \frac{G_{\mathcal{D}}p \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}}q \, d\mathbf{x} + \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x}. \quad (2.42)$$

In the following, we provide the key properties of an AGD, in the sense of Definition 2.4.1, that ensure the convergence of the corresponding AGS (2.42).

**Definition 2.4.3 (AGD-Coercivity, Lifted problem)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.4.1, let  $C_{\mathcal{D}}$  be,

$$C_{\mathcal{D}} = \max_{q \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \frac{\|\Pi_{\mathcal{D}} q\|_{L^2(\Omega)}}{\|q\|_{\mathcal{D}}}. \quad (2.43)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **coercive** if there exists  $C_P \in \mathbb{R}_+$  such that  $C_{\mathcal{D}_m} \leq C_P$  for all  $m \in \mathbb{N}$ .

**Definition 2.4.4 (AGD-Consistency, Lifted problem)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.4.1, let  $S_{\mathcal{D}} : H_d^1(\Omega) \rightarrow [0, +\infty)$  be given by,

$$\forall \varphi \in H_d^1(\Omega), \quad S_{\mathcal{D}}(\varphi) = \min_{q \in X_{\mathcal{D}, \Omega, \Gamma_n}} \left( \|\Pi_{\mathcal{D}} q - \varphi\|_{L^2(\Omega)} + \|G_{\mathcal{D}} q - \nabla_s^\alpha \varphi\|_{L^2(\Omega)^{d+1}} \right). \quad (2.44)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **consistent** if,

$$\forall \varphi \in H_d^1(\Omega), \quad \lim_{m \rightarrow \infty} S_{\mathcal{D}_m}(\varphi) = 0. \quad (2.45)$$

**Definition 2.4.5 (AGD-Limit-conformity, Lifted problem)**

Under the assumption (2.5a), if  $\mathcal{D}$  is an abstract gradient discretisation in the sense of Definition 2.4.1, let  $W_{\mathcal{D}} : H_{\text{div}_s^\alpha}^0(\Omega) \rightarrow [0, +\infty)$  be given by,

$$\forall \varphi \in H_{\text{div}_s^\alpha}^0(\Omega), \quad W_{\mathcal{D}}(\varphi) = \sup_{q \in X_{\mathcal{D}, \Omega, \Gamma_n} \setminus \{0\}} \frac{\left| \int_{\Omega} G_{\mathcal{D}} q \cdot \varphi + \Pi_{\mathcal{D}} q \operatorname{div}_s^\alpha \varphi \, dx \right|}{\|q\|_{\mathcal{D}}}. \quad (2.46)$$

A sequence  $(\mathcal{D}_m)_{m \in \mathbb{N}}$  of abstract gradient discretisations is **limit-conforming** if,

$$\forall \varphi \in H_{\text{div}_s^\alpha}^0(\Omega), \quad \lim_{m \rightarrow \infty} W_{\mathcal{D}_m}(\varphi) = 0.$$

We now can plug Theorem 2.3.8 into the framework of Problem (2.41).

**Theorem 2.4.6 (Error estimates of an AGS, Lifted problem)**

Under Assumptions (2.5), let  $p^\dagger \in H_d^1(\Omega)$  be the solution of Problem (2.41). Let  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, G_{\mathcal{D}})$  be an AGD in the sense of Definition 2.4.1. Then there exists one and only one  $p \in X_{\mathcal{D}, \Omega, \Gamma_n}$  solution to the AGS (2.42); this solution satisfies the following inequality :

$$\|\Pi_{\mathcal{D}} p - p^\dagger\|_{L^2(\Omega)} \leq C_{\mathcal{D}} \bar{H} \left( W_{\mathcal{D}} \left( \frac{1}{H} \left( \nabla_s^\alpha p^\dagger - \mathbf{g} \right) \right) + \frac{S_{\mathcal{D}}(p^\dagger)}{\bar{H}} \right) + S_{\mathcal{D}}(p^\dagger), \quad (2.47a)$$

$$\|G_{\mathcal{D}} p - \nabla_s^\alpha p^\dagger\|_{L^2(\Omega)^{d+1}} \leq \bar{H} \left( W_{\mathcal{D}} \left( \frac{1}{H} \left( \nabla_s^\alpha p^\dagger - \mathbf{g} \right) \right) + \frac{S_{\mathcal{D}}(p^\dagger)}{\bar{H}} \right) + S_{\mathcal{D}}(p^\dagger), \quad (2.47b)$$

where  $S_{\mathcal{D}}$  and  $W_{\mathcal{D}}$  are respectively defined by the AGD-Consistency (2.44) and the AGD-Limit-Conformity (2.46), and  $C_p$  is related to the AGD-Coercivity (2.43).

### 2.4.3 Mixed mass-lumping finite element methods & AGD

In [2, section 4.] the authors applied a finite element approximation to solve the mixed problem (2.39). Thus they needed two discrete spaces, one for the velocity and one for the pressure, denoted respectively  $W_h^0$  and  $Q_h$  in the following. The choice of these spaces is in particular done in order to ensure a discrete *inf-sup* condition. The authors used the  $\mathbb{P}_1/\mathbb{P}_1$  spaces and then the  $\mathbb{P}_1 - \text{iso}\mathbb{P}_2/\mathbb{P}_1$  spaces.

By the way, in both cases, the discrete approximation can be seen as choosing finite element spaces such that  $W_h^0 \subset H_{\text{div}_s^\alpha}^0(\Omega)$  and  $Q_h \subset L^2(\Omega)$ . Then the discrete formulation of Problem (2.39) reads :

Find  $(\mathbf{u}_h, p_h) \in W_h^0 \times Q_h$  such that,

$$\forall \boldsymbol{\omega}_h \in W_h^0, \int_{\Omega} H \mathbf{u}_h \cdot \boldsymbol{\omega}_h - p_h \text{div}_s^\alpha \boldsymbol{\omega}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \boldsymbol{\omega}_h \, d\mathbf{x} , \quad (2.48a)$$

$$\forall q_h \in Q_h, \int_{\Omega} q_h \text{div}_s^\alpha \mathbf{u}_h \, d\mathbf{x} = \int_{\Omega} f q_h \, d\mathbf{x} . \quad (2.48b)$$

Then the authors of [2] introduce a mass-lumping strategy in order to algebraically eliminate the velocity unknowns which leads to reduce the linear system to the discrete pressure unknowns and thus the numerical cost. After all, the authors have to resolve the following problem :

Find  $(\mathbf{u}_h, p_h) \in W_h^0 \times Q_h$  such that,

$$\forall \boldsymbol{\omega}_h \in W_h^0, \int_{\Omega} H \widehat{\mathbf{u}}_h \cdot \widehat{\boldsymbol{\omega}}_h - p_h \text{div}_s^\alpha \boldsymbol{\omega}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} , \quad (2.49a)$$

$$\forall q_h \in Q_h, \int_{\Omega} q_h \text{div}_s^\alpha \mathbf{u}_h \, d\mathbf{x} = \int_{\Omega} f q_h \, d\mathbf{x} , \quad (2.49b)$$

where  $\widehat{\boldsymbol{\omega}}_h$  is the constant piece-wise approximations on the dual cells of the finite element method.

#### Remark 2.4.7 (About the mass-lumping)

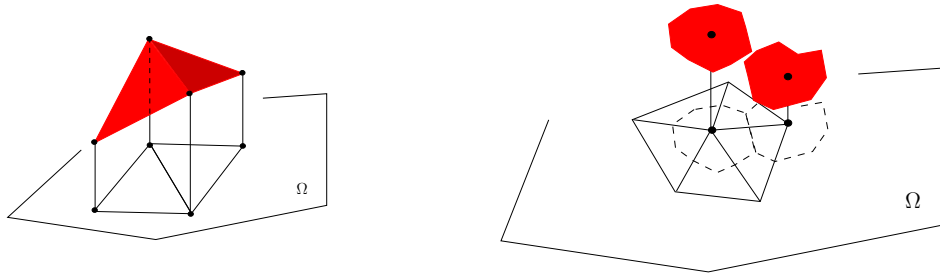
In order to better understand the notation regarding the mass-lumping, we can introduce, for  $d = 2$ ,

$$W_h^0 = \{k \in \{1, 2, 3\}, u^{(k)}(\mathbf{x}) = \sum_{i=1}^N u_i^{(k)} \phi_i(\mathbf{x}), (u_i^{(1)}, u_i^{(2)})^T \cdot \mathbf{n}_\Gamma = 0 \text{ for } \mathbf{x}_i \in \bar{\Gamma}_n\} ,$$

where  $N$  is the number of nodes of the conforming simplicial mesh as defined in section 2.2.6, and for the  $i$ -th node,  $\mathbf{x}_i$  is its coordinates and  $\phi_i$  the  $\mathbb{P}_1$  finite element shape function of this vertex. Thus, if  $\boldsymbol{\omega}_h \in W_h^0$ , we can write,

$$\omega_h(\mathbf{x}) = \begin{pmatrix} \omega_h^{(1)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(1)} \phi_i(\mathbf{x}) \\ \omega_h^{(2)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(2)} \phi_i(\mathbf{x}) \\ \omega_h^{(3)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(3)} \phi_i(\mathbf{x}) \end{pmatrix}, \quad \widehat{\omega}_h(\mathbf{x}) = \begin{pmatrix} \widehat{\omega}_h^{(1)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(1)} \widehat{\phi}_i(\mathbf{x}) \\ \widehat{\omega}_h^{(2)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(2)} \widehat{\phi}_i(\mathbf{x}) \\ \widehat{\omega}_h^{(3)}(\mathbf{x}) = \sum_{i=1}^N \omega_i^{(3)} \widehat{\phi}_i(\mathbf{x}) \end{pmatrix},$$

where  $\widehat{\phi}_i$  is the piece-wise constant function equal to 1 on the dual cell  $K_i$  and 0 otherwise, where  $K_i$  corresponds to the Donald cell around the node  $i$  build using the medians of the triangles as shown at right side of the figure below.



On the figure above : at the left an example of a function  $\omega_h^{(k)}$  and at the right an example of a function  $\widehat{\omega}_h^{(k)}$ , for  $k = 1, 2$  or  $3$ .

#### Remark 2.4.8 (Discrete linear system)

Since in [2, section 4] the mass-lumping is used at the algebraical level, the rigorous equivalence between Problem (2.49) and the schemes used in [2] should be proved through the linear system resulting of the schemes.

The aim of the end of this section is to show how the scheme (2.49) is an AGS in the sense of Definition 2.4.2. To do so, the first step is to define the appropriated AGD. We re-use the notations used in section 2.2.6 regarding the definition of a conforming simplicial mesh of  $\Omega$ .

#### Definition 2.4.9 (Mixed mass-lumping AGD)

Let  $\mathcal{D} = (X_{\mathcal{D}}, \Pi_{\mathcal{D}}, G_{\mathcal{D}})$  be an AGD, in the sense of Definition 2.4.1, defined by :

- $X_{\mathcal{D}} = X_{\mathcal{D}, \Omega, \Gamma_n} \oplus X_{\mathcal{D}, \Gamma_d}$  with  $X_{\mathcal{D}} = \mathbb{R}^N$  and  $X_{\mathcal{D}, \Gamma_d} = \mathbb{R}^{N_d} = \{0\}$  where  $N_d$  is the number of vertices belonging to  $\Gamma_d$ ,
- $\Pi_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)$  ;  $(p_i)_i \mapsto \sum_{i=1}^N p_i \varphi_i$  with  $Q_h := \text{span}(\varphi_i)_{1 \leq i \leq N}$ ,
- $G_{\mathcal{D}} : X_{\mathcal{D}} \rightarrow L^2(\Omega)^{d+1}$  ;  $p \mapsto G_{\mathcal{D}} p$  such that :
  - (i)  $\forall \omega_h \in W_h^0, \int_{\Omega} G_{\mathcal{D}} p \cdot \widehat{\omega}_h + \Pi_{\mathcal{D}} p \operatorname{div}_s^\alpha \omega_h \, d\mathbf{x} = 0$ ,
  - (ii)  $\exists G_h \in W_h^0$  such that  $G_{\mathcal{D}} p = H \widehat{G}_h$ .

The definition of the discrete *special* gradient  $G_{\mathcal{D}}$  must be specified through the following remark.

**Remark 2.4.10 (Construction of  $G_{\mathcal{D}}$ )**

In the item (ii) we look for  $G_h \in W_h^0 := \text{span}(\psi_j)_{j \in J}$  (where  $J$  is related to the d.o.f. of  $W_h^0$ ), so we have to find  $(\beta_j)_{j \in J}$  such that  $\widehat{G}_h = \sum_{j \in J} \beta_j \widehat{\psi}_j$ . Moreover, through the item (i), we have for any

$k \in J$  :

$$\int_{\Omega} G_{\mathcal{D}} p_{\mathcal{D}} \cdot \widehat{\psi}_k + \Pi_{\mathcal{D}} p_{\mathcal{D}} \operatorname{div}_s^{\alpha} \psi_k \, d\mathbf{x} = 0.$$

Then, by using  $\text{supp}(\widehat{\psi}_j) \cap \text{supp}(\widehat{\psi}_k) = \emptyset$  for  $j \neq k$  and by denoting  $\text{supp}(\widehat{\psi}_k) = K_k$ , the dual cell centered around the  $k$ -nth vertex, we obtain,

$$\beta_k \int_{K_k} H |\psi_k|^2 \, d\mathbf{x} = - \int_{\Omega} \Pi_{\mathcal{D}} p_{\mathcal{D}} \operatorname{div}_s^{\alpha} \psi_k \, d\mathbf{x}. \quad (2.50)$$

Thus (2.50) gives a way to compute the family  $(\beta_j)_{j \in J}$  which ensures the computation of  $G_h$  and at last of  $G_{\mathcal{D}}$ .

Now we will state, and prove, the main result of this section.

**Theorem 2.4.11 (Equivalence between mixed mass lumping and abstract gradient scheme)**

- On one hand, let  $(\mathbf{u}_h, p_h) \in W_h^0 \times Q_h$  be the solution of the mixed mass lumping finite element scheme (2.49) with  $W_h^0 := \text{span}(\psi_j)_{j \in J}$  and  $Q_h := \text{span}(\varphi_i)_{i \in I}$ .
- On the other hand, let  $(X_{\mathcal{D}}, \Pi_{\mathcal{D}}, G_{\mathcal{D}})$  be the AGD of Definition 2.4.9 and let  $p_{\mathcal{D}} := (\lambda_i)_{i \in I}$  be the solution of the associated AGS (2.42).
- We also defined  $\mathbf{g}_h \in W_h^0$  such that,

$$\int_{\Omega} H \widehat{\mathbf{g}}_h \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x}, \quad \forall \boldsymbol{\omega}_h \in W_h^0. \quad (2.51)$$

We have the following result :  $(\mathbf{u}_h = \mathbf{g}_h - G_h, p_h = \Pi_{\mathcal{D}} p_{\mathcal{D}})$  is the solution of the mixed scheme (2.49) if and only if  $p_{\mathcal{D}}$  is the solution of the AGS (2.42) with  $\lambda_i = \int_{K_i} p_h \varphi_i \, d\mathbf{x}$  for any vertex numbering  $i$ .

**Proof:**

( $\Rightarrow$ ). Let  $(\mathbf{u}_h, p_h) \in W_h^0 \times Q_h$  the solution of scheme (2.49). We introduce  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$  and  $q_h := \Pi_{\mathcal{D}} q \in Q_h$ . We inject  $q_h$  in (2.49b). Thus we have,

$$\int_{\Omega} \operatorname{div}_s^{\alpha} \mathbf{u}_h \Pi_{\mathcal{D}} q \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}} q \, d\mathbf{x}.$$

By using the item (i) in the construction of the operator  $G_{\mathcal{D}}$ , we have,

$$\int_{\Omega} G_{\mathcal{D}} q \cdot \widehat{\mathbf{u}}_h \, d\mathbf{x} = - \int_{\Omega} \operatorname{div}_s^{\alpha} \mathbf{u}_h \Pi_{\mathcal{D}} q \, d\mathbf{x}.$$

Thus we obtain,

$$- \int_{\Omega} G_{\mathcal{D}} q \cdot \widehat{\mathbf{u}}_h \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}} q \, d\mathbf{x}.$$

Now by using the item (ii) in the construction of  $G_{\mathcal{D}}$ , there exists  $F_h \in W_h^0$  such that  $G_{\mathcal{D}}q = H\widehat{F}_h$  with  $F_h \in W_h^0$ . Then by construction of  $\mathbf{g}_h$ , we can write,

$$\int_{\Omega} \widehat{\mathbf{g}}_h \cdot G_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x} .$$

So we can write that,

$$\int_{\Omega} (\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h) \cdot G_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} f \, \Pi_{\mathcal{D}}q \, d\mathbf{x} + \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x} . \quad (2.52)$$

Let  $p_{\mathcal{D}} := (\lambda_i)_{i \in I} \in X_{\mathcal{D}, \Omega, \Gamma_n}$  such that  $\lambda_i = \int_{K_i} p_h \varphi_i \, d\mathbf{x}$ , this choice implies that  $p_h = \Pi_{\mathcal{D}}p_{\mathcal{D}}$ . Thanks to the item (i) in the construction of the operator  $G_{\mathcal{D}}$ , we can write,

$$\int_{\Omega} (G_{\mathcal{D}}p_{\mathcal{D}} \cdot \widehat{\boldsymbol{\omega}}_h + \Pi_{\mathcal{D}}p_{\mathcal{D}} \operatorname{div}_s^{\alpha} \boldsymbol{\omega}_h) = 0 , \quad \forall \boldsymbol{\omega}_h \in W_h^0 .$$

By combining with (2.49a), we obtain,

$$\int_{\Omega} G_{\mathcal{D}}p_{\mathcal{D}} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} + \int_{\Omega} H\widehat{\mathbf{u}}_h \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} , \quad \forall \boldsymbol{\omega}_h \in W_h^0 .$$

Thanks to the construction of  $\boldsymbol{\omega}_h$ , we can write,

$$\int_{\Omega} \mathbf{g} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{\Omega} H\widehat{\mathbf{g}}_h \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} , \quad \forall \boldsymbol{\omega}_h \in W_h^0 .$$

At this stage, we have,

$$\int_{\Omega} G_{\mathcal{D}}p_{\mathcal{D}} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{\Omega} H(\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h) \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} , \quad \forall \boldsymbol{\omega}_h \in W_h^0 ,$$

which allows us to check, using the item (ii) of  $G_{\mathcal{D}}$ , that  $G_h = \mathbf{g}_h - \mathbf{u}_h$ . To conclude, by taking  $\boldsymbol{\omega}_h = \frac{G_{\mathcal{D}}p_{\mathcal{D}}}{H}$  and by using (2.52), we recover that  $p_{\mathcal{D}}$  is the solution of the AGS (2.42).

( $\Leftarrow$ ). Let  $p_{\mathcal{D}} \in X_{\mathcal{D}, \Omega, \Gamma_n}$  the solution of the AGS (2.42). The item (ii) in the construction of  $G_{\mathcal{D}}$  implies that there exists  $G_h \in W_h^0$  such that  $G_{\mathcal{D}}p_{\mathcal{D}} = H\widehat{G}_h$ . We have defined  $\mathbf{u}_h = \mathbf{g}_h - G_h$ , thus  $\mathbf{u}_h \in W_h^0$  and  $G_{\mathcal{D}}p_{\mathcal{D}} = H(\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h)$ . Then the item (i) allows us to write,

$$\int_{\Omega} H(\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h) \cdot \widehat{\boldsymbol{\omega}}_h + \Pi_{\mathcal{D}}p_{\mathcal{D}} \operatorname{div}_s^{\alpha} \boldsymbol{\omega}_h \, d\mathbf{x} , \quad \forall \boldsymbol{\omega}_h \in W_h^0 . \quad (2.53)$$

Moreover, by construction of  $\mathbf{g}_h$ , we can write,

$$\int_{\Omega} H\widehat{\mathbf{g}}_h \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} . \quad (2.54)$$

Since we have also defined  $p_h = \Pi_{\mathcal{D}}p_{\mathcal{D}}$ , we have  $p_h \in Q_h$  and by injecting (2.54) in (2.53), we recover the first equation (2.49a) of the mixed mass lumping finite element scheme. We now have to recover the second equation of this scheme. To do so, we inject  $G_{\mathcal{D}}p_{\mathcal{D}} = H(\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h)$  in the AGS (2.42), thus we obtain,

$$\int_{\Omega} (\widehat{\mathbf{g}}_h - \widehat{\mathbf{u}}_h) \cdot G_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} f \, \Pi_{\mathcal{D}}q \, d\mathbf{x} + \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x} , \quad \forall q \in X_{\mathcal{D}, \Omega, \Gamma_n} .$$

Thanks to the item (ii) of  $G_{\mathcal{D}}$ , there exists  $F_h \in W_h^0$  such that  $G_{\mathcal{D}}q = H\widehat{F}_h$  which involves, also by using the definition of  $\mathbf{g}_h$ , that,

$$\int_{\Omega} \widehat{\mathbf{g}}_h \cdot G_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} H\widehat{\mathbf{g}}_h \cdot \widehat{F}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{g} \cdot \widehat{F}_h \, d\mathbf{x} = \int_{\Omega} \frac{\mathbf{g} \cdot G_{\mathcal{D}}q}{H} \, d\mathbf{x} .$$

Thus, at this stage, we have,

$$- \int_{\Omega} \widehat{\mathbf{u}}_h \cdot G_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}}q \, d\mathbf{x} , \quad \forall q \in X_{\mathcal{D},\Omega,\Gamma_n} .$$

Thanks to the item (i) of  $G_{\mathcal{D}}$ , and since  $\mathbf{u}_h \in W_h^0$ , we obtain,

$$- \int_{\Omega} G_{\mathcal{D}}q \cdot \widehat{\mathbf{u}}_h \, d\mathbf{x} = \int_{\Omega} \Pi_{\mathcal{D}}q \operatorname{div}_s^{\alpha} \mathbf{u}_h \, d\mathbf{x} , \quad \forall q \in X_{\mathcal{D},\Omega,\Gamma_n} .$$

So, by combining the two last equations, we have,

$$\int_{\Omega} \operatorname{div}_s^{\alpha} \mathbf{u}_h \Pi_{\mathcal{D}}q \, d\mathbf{x} = \int_{\Omega} f \Pi_{\mathcal{D}}q \, d\mathbf{x} , \quad \forall q \in X_{\mathcal{D},\Omega,\Gamma_n} .$$

Finally, by taking  $q_h := \Pi_{\mathcal{D}}q$ , we recover (2.49b). So we have proved that  $(\mathbf{u}_h, p_h)$  is the solution of (2.49). ■

#### 2.4.4 The $\mathbb{P}_1/\mathbb{P}_1$ approximation

We study here the scheme (2.49) in the case of the  $\mathbb{P}_1/\mathbb{P}_1$  finite elements approximation. This method has been used in [2, subsection 4.1]. The aim of this last section is to check that, in this case, the associated AGD, see Definition 2.4.9, satisfies the mentioned AGD-properties.

As in section 2.2.6, for  $d = 2$ , we denote by  $\mathfrak{T} = (\mathcal{M}, \mathcal{F}, \mathcal{C}, \mathcal{V})$  a conforming simplicial mesh of  $\Omega$  in the sense of Definition 7.4 in [22] where  $\mathcal{M}$  is the set of *cells* (or *triangles*),  $\mathcal{F}$  of *faces* (or *edges*),  $\mathcal{C}$  of the *cell centers* and  $\mathcal{V}$  of the *vertices* (or *nodes*). A conforming simplicial mesh of  $\Omega$  is a polytopal mesh in the sense of Definition 7.2 in [22] (see also Definition 1.50 in [27]) such that for each  $K \in \mathcal{M}$  we have  $\operatorname{Card}(\mathcal{F}_K) = \operatorname{Card}(\mathcal{V}_K) = 3$ . We then introduce  $N = \operatorname{card}(\mathcal{V})$  and the family  $(\phi_i)_{i=1}^N$  of  $\mathbb{P}_1$  finite element shape functions associated to the vertices. We can now introduce the following space,

$$V_h = \{v_h \in \mathcal{C}^0(\overline{\Omega}) , v_h|_K \in \mathbb{P}_1(\mathbb{R}^d) , \forall K \in \mathcal{M}\} .$$

So we take, in the AGD-Defintion 2.4.9,  $Q_h = V_h$  and  $W_h^0 = \{\mathbf{u}_h \in V_h^{d+1} , \mathbf{u}_h = (\mathbf{v}_h, w_h)^T , \mathbf{v}_h \cdot \mathbf{n}_{\Gamma} = 0 \text{ on } \Gamma_n\}$ . We denote  $I$  the set of vertices of the mesh, where  $\mathbf{x}_i \in \mathbb{R}^d$  is the vertex numbered  $i$ . Then we write that  $V_h = \operatorname{span}(\phi_i)_{i \in I}$  and also,

$$W_h^0 = \{k \in \{1, \dots, d+1\} = \{1, 2, 3\}, u_h^{(k)}(\mathbf{x}) = \sum_{i \in I} u_i^{(k)} \phi_i(\mathbf{x}) , (u_i^{(1)}, u_i^{(2)})^T \cdot \mathbf{n}_{\Gamma} = 0 \text{ for } \mathbf{x}_i \in \overline{\Gamma}_n\} .$$



**Remark 2.4.12** (*The boundary condition on  $\Gamma_n$* )

We recall that, thanks to the lifting, the boundary condition for  $\mathbf{u}^\dagger$  on  $\Gamma_n$  reads  $H\mathbf{u}^\dagger \cdot \mathbf{n}_s = 0$  which is equivalent to  $\mathbf{u}^\dagger \cdot \mathbf{n}_s = 0$  on  $\Gamma_n$  thanks to hypothesis (2.5b).

At the discrete level, regarding the grid, we have to include the nodes at the end of  $\Gamma_n$  in order to fulfill the Neumann condition all along  $\Gamma_n$ . Indeed, for  $\mathbf{x} \in \bar{\Gamma}_n$  there exists  $(i, j) \in I^2$  and  $\beta \in (0, 1)$  such that  $\mathbf{x} = \beta\mathbf{x}_i + (1 - \beta)\mathbf{x}_j$ . Then for any  $\mathbf{u}_h = (\mathbf{v}_h, w_h)^T \in W_h^0$ , we have,

$$\mathbf{v}_h(\mathbf{x}) \cdot \mathbf{n}_\Gamma(\mathbf{x}) = \beta\mathbf{v}_h(\mathbf{x}_i) \cdot \mathbf{n}_\Gamma(\mathbf{x}_i) + (1 - \beta)\mathbf{v}_h(\mathbf{x}_j) \cdot \mathbf{n}_\Gamma(\mathbf{x}_j).$$

So it is equivalent to say  $\mathbf{v}_h \cdot \mathbf{n}_\Gamma = 0$  on  $\Gamma_n$  as imposed  $\mathbf{v}_h(\mathbf{x}_i) \cdot \mathbf{n}_\Gamma(\mathbf{x}_i) = 0$  for any  $i \in I_{\Gamma_n}$ , with  $\mathbf{v}_h(\mathbf{x}_i) = 0$  if  $\mathbf{x}_i$  corresponds to a boundary angle and where  $I_{\Gamma_n} = \{i \in I \mid \mathbf{x}_i \in \bar{\Gamma}_n\}$ .

In the following  $h$  denotes the *mesh size* which can be seen here as the length of the longest edge of the mesh.

**AGD-Coercivity**

Following the Definition 2.4.3, we are looking for a constant  $C_{\mathcal{D}}$ , independent of  $h$  such that, for any  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$ ,  $\|\Pi_{\mathcal{D}}q\|_{L^2(\Omega)} \leq C_{\mathcal{D}}\|G_{\mathcal{D}}q\|_{L^2(\Omega)^{d+1}}$ .

Let  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$ , we take  $\boldsymbol{\omega}_h = (0, 0, \Pi_{\mathcal{D}}q)^T \in W_h^0$ . Thus we have, using (2.1b),  $\operatorname{div}_s^\alpha(\boldsymbol{\omega}_h) = \alpha \Pi_{\mathcal{D}}q$ . And by using the item (i) in the construction of  $G_{\mathcal{D}}$ , we obtain,

$$\int_{\Omega} G_{\mathcal{D}}q \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} + \alpha \int_{\Omega} (\Pi_{\mathcal{D}}q)^2 \, d\mathbf{x} = 0.$$

Then using the Cauchy-Schwarz inequality, we can write,

$$\alpha \|\Pi_{\mathcal{D}}q\|_{L^2(\Omega)}^2 \leq \|G_{\mathcal{D}}q\|_{L^2(\Omega)^{d+1}} \|\widehat{\boldsymbol{\omega}}_h\|_{L^2(\Omega)^{d+1}}.$$

We can remark that  $\|\widehat{\boldsymbol{\omega}}_h\|_{L^2(\Omega)^{d+1}} = \|\widehat{\Pi_{\mathcal{D}}q}\|_{L^2(\Omega)}$  and also that,

$$\|\Pi_{\mathcal{D}}q\|_{L^2(\Omega)}^2 = \sum_{K \in \mathcal{M}} \int_K \left( \sum_{j \in I_K} \beta_j \phi_j \right)^2 \, d\mathbf{x}, \quad \|\widehat{\Pi_{\mathcal{D}}q}\|_{L^2(\Omega)}^2 = \sum_{K \in \mathcal{M}} \int_K \left( \sum_{j \in I_K} \beta_j \widehat{\phi}_j \right)^2 \, d\mathbf{x},$$

where  $q = (\beta_1, \dots, \beta_N)^T$  and  $I_K = \{i \in I \mid \mathbf{x}_i \in K\}$ . But, due to the mass lumping, we can write that, for any  $K \in \mathcal{M}$ ,

$$\int_K \left( \sum_{j \in I_K} \beta_j \widehat{\phi}_j \right)^2 \, d\mathbf{x} = \int_K \left( \beta_{j_1} \widehat{\phi}_{j_1} + \beta_{j_2} \widehat{\phi}_{j_2} + \beta_{j_3} \widehat{\phi}_{j_3} \right)^2 \, d\mathbf{x} = \frac{|K|}{3} (\beta_{j_1}^2 + \beta_{j_2}^2 + \beta_{j_3}^2),$$

where  $I_K := \{j_1, j_2, j_3\}$  is the set of indices of the vertices of the triangular cell  $K$ . On the other hand we have, with  $\boldsymbol{\beta}_K := (\beta_{j_1}, \beta_{j_2}, \beta_{j_3})^T$ ,

$$\int_K \left( \sum_{j \in I_K} \beta_j \phi_j \right)^2 \, d\mathbf{x} = \boldsymbol{\beta}_K^T M_K \boldsymbol{\beta}_K,$$

where  $M_K = \frac{|K|}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$  is the mass matrix which is positive definite, thus satisfies,

$$\boldsymbol{\beta}_K^T M_K \boldsymbol{\beta}_K \geq \min\{\lambda, \lambda \in \text{Spec}(M_K)\} \boldsymbol{\beta}^T \boldsymbol{\beta}.$$

So we obtain,

$$\int_K \left( \sum_{j \in I_K} \beta_j \phi_j \right)^2 d\mathbf{x} \geq \frac{|K|}{12} (\beta_{j_1}^2 + \beta_{j_2}^2 + \beta_{j_3}^2).$$

Finally we have proved that,

$$\|\widehat{\Pi_{\mathcal{D}}} q\|_{L^2(\Omega)}^2 \leq 4 \|\Pi_{\mathcal{D}} q\|_{L^2(\Omega)}^2.$$

Thus, we can write that,

$$\|\Pi_{\mathcal{D}} q\|_{L^2(\Omega)} \leq \frac{2}{\alpha} \|\mathbf{G}_{\mathcal{D}} q\|_{L^2(\Omega)^{d+1}},$$

which involves  $C_{\mathcal{D}} = \frac{2}{\alpha}$ .

**Remark 2.4.13 (Norm on  $X_{\mathcal{D}}$ )**

When the AGD-Coercivity holds,  $\|\mathbf{G}_{\mathcal{D}}(\cdot)\|_{L^2(\Omega)^{d+1}}$  is a norm on  $X_{\mathcal{D},\Omega,\Gamma_n}$ . Indeed, let  $q \in X_{\mathcal{D},\Omega,\Gamma_n}$  such that  $\|\mathbf{G}_{\mathcal{D}} q\|_{L^2(\Omega)^{d+1}} = 0$ , thus  $\|\Pi_{\mathcal{D}} q\|_{L^2(\Omega)} = 0$ . Since the family  $(\phi_i)_{i \in I}$  is linearly independent, we infer that  $q = 0$ .

**AGD-Limit-conformity**

As mentioned in [22, Lemma A.9], it is sufficient to prove the limit-conformity in a dense subset of  $H_{\text{div}_s^\alpha}^0(\Omega)$ . Thus we take,  $\boldsymbol{\omega} \in H_{\text{div}_s^\alpha}^0(\Omega) \cap (\mathcal{C}^2(\overline{\Omega}))^{d+1}$ . Then we denote  $\boldsymbol{\omega} := (\omega^{(1)}, \omega^{(2)}, \omega^{(3)})^T$  and  $\boldsymbol{\omega}_h = (\omega_h^{(1)}, \omega_h^{(2)}, \omega_h^{(3)})^T$  where, for  $k \in \{1, 2, 3\}$ ,

$$\omega_h^{(k)} := \sum_{i \in I} \omega^{(k)}(\mathbf{x}_i) \phi_i.$$

We have  $\boldsymbol{\omega}_h \in W_h^0$  since  $\boldsymbol{\omega}_h \cdot \mathbf{n}_s = 0$  (due to the fact that  $\boldsymbol{\omega} \cdot \mathbf{n}_s = 0$ ).

Then we can use the item (i) in the construction of  $\mathbf{G}_{\mathcal{D}}$  to write,

$$\int_{\Omega} \mathbf{G}_{\mathcal{D}} q \cdot \boldsymbol{\omega} + \Pi_{\mathcal{D}} q \text{div}_s^\alpha \boldsymbol{\omega} d\mathbf{x} = \int_{\Omega} \mathbf{G}_{\mathcal{D}} q \cdot (\boldsymbol{\omega} - \widehat{\boldsymbol{\omega}}_h) + \Pi_{\mathcal{D}} q \text{div}_s^\alpha (\boldsymbol{\omega} - \boldsymbol{\omega}_h) d\mathbf{x}, \quad \forall q \in X_{\mathcal{D},\Omega,\Gamma_n}.$$

Using the Cauchy-Schwarz inequality and the AGD-Coercivity, we obtain,

$$\left| \int_{\Omega} \mathbf{G}_{\mathcal{D}} q \cdot \boldsymbol{\omega} + \Pi_{\mathcal{D}} q \text{div}_s^\alpha \boldsymbol{\omega} d\mathbf{x} \right| \leq ( \|\widehat{\boldsymbol{\omega}}_h - \boldsymbol{\omega}_h\|_{L^2(\Omega)^{d+1}} + \|\boldsymbol{\omega} - \boldsymbol{\omega}_h\|_{H_{\text{div}_s^\alpha}^0(\Omega)} (1 + C_{\mathcal{D}}) ) \|q\|_{\mathcal{D}}, \quad \forall q \in X_{\mathcal{D},\Omega,\Gamma_n}.$$

Component by component, for  $k \in \{1, 2, 3\}$ , we get (see for instance equation (8.18) in [22]) that,

$$\|\widehat{\boldsymbol{\omega}}_h^{(k)} - \boldsymbol{\omega}_h^{(k)}\|_{L^2(\Omega)} \leq h \|\nabla \boldsymbol{\omega}_h^{(k)}\|_{L^2(\Omega)^d},$$

where  $\nabla \boldsymbol{\omega}_h$  remains bounded since  $\|\nabla \widehat{\boldsymbol{\omega}}_h^{(k)} - \nabla \boldsymbol{\omega}_h^{(k)}\|_{L^2(\Omega)^d} \leq h \|\boldsymbol{\omega}^{(k)}\|_{2,\infty}$ , with  $\|\boldsymbol{\omega}^{(k)}\|_{2,\infty} = \sup_{\Omega} |\partial^2 \boldsymbol{\omega}^{(k)}|$ . We can also write,

$$\|\boldsymbol{\omega}^{(k)} - \boldsymbol{\omega}_h^{(k)}\|_{H^1(\Omega)} \leq h \|\boldsymbol{\omega}^{(k)}\|_{2,\infty}, \quad \forall k \in \{1, 2, 3\}.$$

Moreover, by definition of the norm on  $H_{\text{div}_s^\alpha}^0(\Omega)$  given by (2.40), we have  $C > 0$  independent of the mesh such that,

$$\|\boldsymbol{\omega} - \boldsymbol{\omega}_h\|_{H_{\text{div}_s^\alpha}^0(\Omega)} \leq C \|\boldsymbol{\omega} - \boldsymbol{\omega}_h\|_{H^1(\Omega)^d},$$

which allows us to conclude that our AGD is limit-conforming in the sens of Definition 2.4.5.

### AGD-Consistency

Following the Definition 2.4.4, we have to prove that,  $\forall \psi \in H_d^1(\Omega)$ ,

$$\lim_{h \rightarrow 0} S_{\mathcal{D}_h}(\psi) = \lim_{h \rightarrow 0} \min_{q \in X_{\mathcal{D}, \Omega, \Gamma_n}} \left( \|\Pi_{\mathcal{D}_h} q - \psi\|_{L^2(\Omega)} + \|\mathbf{G}_{\mathcal{D}_h} q - \nabla_s^\alpha \psi\|_{L^2(\Omega)^{d+1}} \right) = 0.$$

As mentioned in [22, Lemma A.13], it is sufficient to prove this result in a dense subset of  $H_d^1(\Omega)$ . Thus we will take,  $\psi \in H_d^1(\Omega) \cap \mathcal{C}^2(\overline{\Omega})$ .

We have, for all  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$ ,

$$0 \leq S_{\mathcal{D}}(\psi) \leq \|\Pi_{\mathcal{D}_h} q - \psi\|_{L^2(\Omega)} + \|\mathbf{G}_{\mathcal{D}_h} q - \nabla_s^\alpha \psi\|_{L^2(\Omega)^{d+1}}.$$

We introduce  $q = (\psi(\mathbf{x}_i))_{i \in I}$ , thus  $q \in X_{\mathcal{D}, \Omega, \Gamma_n}$  and  $\Pi_{\mathcal{D}} q = \sum_{i \in I} \psi(\mathbf{x}_i) \phi_i$ .

By definition of  $Q_h$ ,  $\Pi_{\mathcal{D}} q$  is the  $\mathbb{P}_1$  finite element approximation of  $\psi$ , thus we have,

$$\lim_{h \rightarrow 0} \|\Pi_{\mathcal{D}} q - \psi\|_{L^2(\Omega)} = 0.$$

Now we have to study the term  $\|\mathbf{G}_{\mathcal{D}} q - \nabla_s^\alpha \psi\|_{L^2(\Omega)^{d+1}}$ .

Thanks to the item (ii) in the construction of  $\mathbf{G}_{\mathcal{D}}$ , we know that there exists  $G_h \in W_h^0$  such that  $\mathbf{G}_{\mathcal{D}} q = H \widehat{G}_h$ . For any  $i \in I$ , and almost everywhere in  $K_i$  (the dual cell around the vertex  $i$ ), the function  $\widehat{G}_h$  is constant and, in the following, we denote this quantity  $\widehat{G}_{h,i}$ .

We are looking for a bound of the difference  $\mathbf{G}_{\mathcal{D}} q - \nabla_s^\alpha \psi$ . For all  $\mathbf{x} \in \overline{\Omega}$ , there exists  $i \in I$  such that  $\mathbf{x} \in K_i$  and we can write,

$$\mathbf{G}_{\mathcal{D}} q(\mathbf{x}) - \nabla_s^\alpha \psi(\mathbf{x}) = (H(\mathbf{x}_i) \widehat{G}_{h,i} - \nabla_s^\alpha \psi(\mathbf{x}_i)) + (H(\mathbf{x}) - H(\mathbf{x}_i)) \widehat{G}_{h,i} - (\nabla_s^\alpha \psi(\mathbf{x}) - \nabla_s^\alpha \psi(\mathbf{x}_i)).$$

Since for all  $(a, b, c) \in \mathbb{R}^3$ , we have  $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$ , we can write,

$$\begin{aligned} \|\mathbf{G}_{\mathcal{D}} q - \nabla_s^\alpha \psi\|_{L^2(\Omega)^{d+1}}^2 &\leq 3 \sum_{i \in I} \left( \|H(\mathbf{x}_i) \widehat{G}_{h,i} - \nabla_s^\alpha \psi(\mathbf{x}_i)\|_{L^2(K_i)^{d+1}}^2 + \right. \\ &\quad \left. \|(H - H(\mathbf{x}_i)) \widehat{G}_{h,i}\|_{L^2(K_i)^{d+1}}^2 + \|\nabla_s^\alpha \psi - \nabla_s^\alpha \psi(\mathbf{x}_i)\|_{L^2(K_i)^{d+1}}^2 \right). \end{aligned} \quad (2.55)$$

Now we have to study the three terms inside the sum. To do so we have to compute  $\widehat{G}_{h,i}^{(k)}$  for any vertex  $i$  and any coordinate  $k \in \{1, 2, 3\}$ . To do so, we need to distinguish the value of  $k$ , but also if the vertex  $i$  is, or not, on  $\Gamma_n$ . This is why we distinguish if  $\mathbf{x}_i$  belongs, or not, to  $\Gamma_n$ .

- If  $\mathbf{x}_i \in \overline{\Omega} \setminus \overline{\Gamma}_n$ .

**[I]** We take  $\boldsymbol{\omega}_h = (\phi_i, 0, 0)^T$ . By localisation of the vertex  $i$ , we have  $\boldsymbol{\omega}_h = 0$  on  $\Gamma_n$ . Thus  $\boldsymbol{\omega}_h \in W_h^0$ . Moreover  $\int_{\Omega} G_{\mathcal{D}} q \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{K_i} G_{\mathcal{D}}^{(1)} q \, d\mathbf{x}$  where  $K_i$  still denotes the dual cell around the vertex  $i$ . We can write,

$$\int_{\Omega} G_{\mathcal{D}} q \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = H_i |K_i| \widehat{G}_{h,i}^{(1)},$$

with  $H_i = \frac{1}{|K_i|} \int_{K_i} H \, d\mathbf{x}$ . Then, using the item (i) in the construction of  $G_{\mathcal{D}}$  and the duality relation (2.27) (since  $\Pi_{\mathcal{D}} q \in H_d^1(\Omega)$  and  $\boldsymbol{\omega}_h \in W_h^0 \subset H_{\text{div}_s^\alpha}(\Omega)$ ), we can write,

$$H_i |K_i| \widehat{G}_{h,i}^{(1)} = \int_{\Omega} \boldsymbol{\omega}_h \cdot \nabla_s^\alpha \Pi_{\mathcal{D}} q \, d\mathbf{x},$$

which implies by definitions of  $\boldsymbol{\omega}_h$  and  $\nabla_s^\alpha$  (2.1a),

$$H_i |K_i| \widehat{G}_{h,i}^{(1)} = \sum_{T \in \mathcal{M}_i} \left( \int_T \phi_i (H (\nabla \Pi_{\mathcal{D}} q)^{(1)} + \Pi_{\mathcal{D}} q (\nabla \zeta)^{(1)}) \, d\mathbf{x} \right),$$

where  $\mathcal{M}_i$  is the set of triangles of the mesh connected to the vertex  $i$ . We can find in [27] that, for every  $T \in \mathcal{M}_i$ ,

$$\|(\nabla \Pi_{\mathcal{D}} q)^{(1)} - (\nabla \psi)^{(1)}\|_{L^\infty(T)} \leq h \|\psi\|_{2,\infty},$$

and

$$\|\Pi_{\mathcal{D}} q - \psi\|_{L^\infty(T)} \leq h \|\psi\|_{2,\infty}.$$

Thus, we can introduce  $R_1$  and  $R_2$ , two bounded functions defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$(\nabla \Pi_{\mathcal{D}} q)^{(1)}(\mathbf{x}) = (\nabla \psi)^{(1)}(\mathbf{x}_i) + h R_1(\mathbf{x}),$$

and

$$\Pi_{\mathcal{D}} q(\mathbf{x}) = \psi(\mathbf{x}_i) + h R_2(\mathbf{x}),$$

thus with  $|R_l(\mathbf{x})| \leq \|\psi\|_{2,\infty}$  for  $l = 1, 2$  and  $\mathbf{x} \in \bigcup_{T \in \mathcal{M}_i} T$ . This implies,

$$\sum_{T \in \mathcal{M}_i} \int_T \phi_i (H (\nabla \Pi_{\mathcal{D}} q)^{(1)} + \Pi_{\mathcal{D}} q (\nabla \zeta)^{(1)}) \, d\mathbf{x} = \sum_{T \in \mathcal{M}_i} \left[ \int_T \phi_i (H (\nabla \psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i) (\nabla \zeta)^{(1)}) \, d\mathbf{x} + h \mathcal{J}_T \right].$$

where  $\mathcal{J}_T = \int_T \phi_i H R_1 + (\nabla \zeta)^{(1)} R_2 \, d\mathbf{x}$ . Since  $0 \leq \phi_i \leq 1$  and the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c), all the terms inside the integral  $\mathcal{J}_T$  are bounded on  $T$  and thus  $|\mathcal{J}_T| \leq |T|C$  with  $C$  independent of the mesh.

So we can introduce  $R_3$  a bounded function defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$H_i |K_i| \widehat{G}_{h,i}^{(1)} = \sum_{T \in \mathcal{M}_i} \left[ \int_T \phi_i (H (\nabla \psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i) (\nabla \zeta)^{(1)}) \, d\mathbf{x} \right] + h |K_i| R_3.$$

Then we write that, for  $T \in \mathcal{M}_i$  and  $\mathbf{x} \in T$ ,

$$\begin{aligned} H(\mathbf{x}) (\nabla \psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i) (\nabla \zeta)^{(1)}(\mathbf{x}) &= H(\mathbf{x}_i) (\nabla \psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i) (\nabla \zeta)^{(1)}(\mathbf{x}_i) \\ &\quad + (H(\mathbf{x}) - H(\mathbf{x}_i)) (\nabla \psi)^{(1)}(\mathbf{x}_i) + ((\nabla \zeta)^{(1)}(\mathbf{x}) - (\nabla \zeta)^{(1)}(\mathbf{x}_i)) \psi(\mathbf{x}_i). \end{aligned} \quad (2.56)$$

Due to the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c), we write that there exists a bounded function  $R_4$  defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$H(\mathbf{x})(\nabla\psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla\zeta)^{(1)}(\mathbf{x}) = H(\mathbf{x}_i)(\nabla\psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla\zeta)^{(1)}(\mathbf{x}_i) + h R_4(\mathbf{x}).$$

Since  $\sum_{T \in \mathcal{M}_i} \int_T \phi_i(x) \, d\mathbf{x} = |K_i|$ , we obtain,

$$\begin{aligned} H_i |K_i| \widehat{G}_{h,i}^{(1)} &= |K_i| \left( \underbrace{H(\mathbf{x}_i)(\nabla\psi)^{(1)}(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla\zeta)^{(1)}(\mathbf{x}_i)}_{=(\nabla_s^\alpha \psi)^{(1)}(\mathbf{x}_i)} + h(R_3 + R_4) \right). \end{aligned}$$

Since  $|H_i - H(\mathbf{x}_i)| \leq C h$  for a constant  $C > 0$ , we can prove that,

$$|H(\mathbf{x}_i) \widehat{G}_{h,i}^{(1)} - \nabla_s^\alpha \psi(\mathbf{x}_i)^{(1)}| \leq h C_*,$$

and thus, there exists  $R_* > 0$  independent of the mesh such that,

$$\|(\mathbb{G}_{\mathcal{D}q})^{(1)} - (\nabla_s^\alpha \psi)^{(1)}\|_{L^\infty(\Omega_{int})} \leq h R_*,$$

where  $\Omega_{int} = \bigcup_{T \in \mathcal{M}_i | i \in I_{int}} T$  with  $I_{int} = \{i \in I \mid \mathbf{x}_i \notin \bar{\Gamma}_n\}$ . We can write,

$$\|(\mathbb{G}_{\mathcal{D}q})^{(1)} - (\nabla_s^\alpha \psi)^{(1)}\|_{L^2(\Omega_{int})}^2 \leq |\Omega_{int}| C h^2. \quad (2.57)$$

[II] If we take  $\boldsymbol{\omega}_h = (0, \phi_i, 0)^T$  a similar demonstration proves that,

$$\|(\mathbb{G}_{\mathcal{D}q})^{(2)} - (\nabla_s^\alpha \psi)^{(2)}\|_{L^\infty(\Omega_{int})} \leq h R_*,$$

and thus,

$$\|(\mathbb{G}_{\mathcal{D}q})^{(2)} - (\nabla_s^\alpha \psi)^{(2)}\|_{L^2(\Omega_{int})}^2 \leq |\Omega_{int}| C h^2 \quad (2.58)$$

[III] If we take  $\boldsymbol{\omega}_h = (0, 0, \phi_i)^T$ , by following the same strategy, we obtain,

$$\begin{aligned} H_i |K_i| (G_h)_i^{(3)} &= \int_{\Omega} \boldsymbol{\omega}_h \cdot \nabla_s^\alpha \Pi_{\mathcal{D}q} \, d\mathbf{x} = \int_{\Omega} \phi_i (\nabla_s^\alpha \Pi_{\mathcal{D}q})^{(3)} \, d\mathbf{x} = \int_{\Omega} \phi_i (-\alpha \Pi_{\mathcal{D}q}) \, d\mathbf{x} \\ &= -\alpha \int_{\Omega} \phi_i \left( \sum_{j \in I} \psi(\mathbf{x}_j) \phi_j \right) \, d\mathbf{x} = -\alpha \sum_{T \in \mathcal{M}_i} \left( \int_T \phi_i \sum_{j \in I_T} \psi(\mathbf{x}_j) \phi_j \, d\mathbf{x} \right) \\ &= -\alpha \sum_{T \in \mathcal{M}_i} \int_T \phi_i \left( \sum_{j \in I_T} \psi(\mathbf{x}_i) \phi_j + \sum_{j \in I_T} (\psi(\mathbf{x}_j) - \psi(\mathbf{x}_i)) \phi_j \right) \, d\mathbf{x}, \end{aligned}$$

where  $I_T = \{i \in I \mid \mathbf{x}_i \in T\}$ . We have,

$$\sum_{j \in I_T} \psi(\mathbf{x}_i) \phi_j = \psi(\mathbf{x}_i) \underbrace{\sum_{j \in I_T} \phi_j}_{=1} = \psi(\mathbf{x}_i).$$

We have also  $0 \leq \phi_j \leq 1$  and, due to the regularity of  $\psi$  there exists  $R_T$ , a bounded function defined on  $T$ , such that  $(\psi(\mathbf{x}_j) - \psi(\mathbf{x}_i)) \phi_j = h R_T$ . Since  $(\nabla_s^\alpha \psi)^{(3)}(\mathbf{x}_i) = -\alpha \psi(\mathbf{x}_i)$  and  $\sum_{T \in \mathcal{M}_i} \int_T \phi_i \, d\mathbf{x} = |K_i|$ , we can write,

$$|H_i (G_h)_i^{(3)} - (\nabla_s^\alpha \psi)^{(3)}| \leq C h.$$

And, by using the same trick between  $H_i$  and  $H(\mathbf{x}_i)$ , we can write,

$$|H(\mathbf{x}_i)\widehat{G}_{h,i}^{(3)} - (\nabla_s^\alpha \psi(\mathbf{x}_i))^{(3)}| \leq h C_*,$$

and thus,

$$\|(G_{\mathcal{D}q})^{(3)} - (\nabla_s^\alpha \psi)^{(3)}\|_{L^\infty(\Omega_{int})} \leq h C_*.$$

We remark that these computations are independent of the assumption  $\mathbf{x}_i \in \overline{\Omega} \setminus \overline{\Gamma}_n$ . Thus we can obtain the same estimation if  $\mathbf{x}_i \in \overline{\Gamma}_n$  and we can write,

$$\|(G_{\mathcal{D}q})^{(3)} - (\nabla_s^\alpha \psi)^{(3)}\|_{L^2(\Omega)}^2 \leq |\Omega| C h^2. \quad (2.59)$$

- If  $\mathbf{x}_i \in \overline{\Gamma}_n$ .

We take  $\boldsymbol{\omega}_h = (\beta\phi_i, \lambda\phi_i, 0)^T$  (implicitly  $\beta$  and  $\lambda$  depend of the vertex  $i$ ). Thus  $\boldsymbol{\omega}_h(\mathbf{x}) \cdot \mathbf{n}_{\Gamma_n}(\mathbf{x}) = \beta \mathbf{n}_{\Gamma_n}^{(1)}(\mathbf{x}) + \lambda \mathbf{n}_{\Gamma_n}^{(2)}(\mathbf{x})$  for  $\mathbf{x} \in \overline{\Gamma}_n$ . So in order to obtain  $\boldsymbol{\omega}_h \in W_h^0$ , we need to have (see Remark 2.4.12)  $\beta \mathbf{n}_{\Gamma_n}^{(1)}(\mathbf{x}_i) + \lambda \mathbf{n}_{\Gamma_n}^{(2)}(\mathbf{x}_i) = 0$ . For instance, we can chose  $\beta = \mathbf{n}_{\Gamma_n}^{(2)}(\mathbf{x}_i) := \mathbf{n}_{\Gamma_n,i}^{(2)}$  and  $\lambda = -\mathbf{n}_{\Gamma_n}^{(1)}(\mathbf{x}_i) := -\mathbf{n}_{\Gamma_n,i}^{(1)}$ .

As recalled, we have  $G_{\mathcal{D}q} = H\widehat{G}_h$  with  $G_h \in W_h^0$ , thus  $G_h \cdot \mathbf{n}_{\Gamma_n} = 0$ . Since  $(\widehat{G}_h)_i = G_h(\mathbf{x}_i)$  we can write  $(\widehat{G}_h)_i^{(1)} \mathbf{n}_{\Gamma_n,i}^{(1)} + (\widehat{G}_h)_i^{(2)} \mathbf{n}_{\Gamma_n,i}^{(2)} = 0$ .

Moreover,

$$\int_{\Omega} G_{\mathcal{D}q} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = \int_{K_i} G_{\mathcal{D}q} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} \, d\mathbf{x},$$

where  $K_i$  still denotes the dual cell around the vertex  $i$ . We can write,

$$\int_{\Omega} G_{\mathcal{D}q} \cdot \widehat{\boldsymbol{\omega}}_h \, d\mathbf{x} = H_i |K_i| \widehat{G}_{h,i} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} \, d\mathbf{x},$$

with  $H_i = \frac{1}{|K_i|} \int_{K_i} H \, d\mathbf{x}$ . Then using the item (i) in the construction of  $G_{\mathcal{D}}$  and the duality relation (2.27) (since  $\Pi_{\mathcal{D}q} \in H_d^1(\Omega)$  and  $\boldsymbol{\omega}_h \in W_h^0 \subset H_{\text{div}_s^\alpha}(\Omega)$ ), we can write,

$$H_i |K_i| \widehat{G}_{h,i} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} = \int_{\Omega} \boldsymbol{\omega}_h \cdot \nabla_s^\alpha \Pi_{\mathcal{D}q} \, d\mathbf{x},$$

which implies, by definitions of  $\boldsymbol{\omega}_h$  and  $\nabla_s^\alpha$  (2.1a),

$$H_i |K_i| \widehat{G}_{h,i} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} = \sum_{T \in \mathcal{M}_i} \int_T \phi_i (H(\nabla \Pi_{\mathcal{D}q}) + \Pi_{\mathcal{D}q}(\nabla \zeta)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} \, d\mathbf{x},$$

where  $\mathcal{M}_i$  is the set of triangles of the mesh connected to the vertex  $i$ . We can reused that, for  $k = 1, 2$ ,

$$\|(\nabla \Pi_{\mathcal{D}q})^{(k)} - (\nabla \psi)^{(k)}\|_{L^\infty(T)} \leq h \|\psi\|_{2,\infty},$$

and

$$\|\Pi_{\mathcal{D}}q - \psi\|_{L^\infty(T)} \leq h\|\psi\|_{2,\infty}.$$

Thus, we can introduce  $\mathbf{R}_1$  and  $R_2$  two bounded functions defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$(\nabla \Pi_{\mathcal{D}}q)(\mathbf{x}) = (\nabla \psi)(\mathbf{x}_i) + h \mathbf{R}_1(\mathbf{x}),$$

and

$$\Pi_{\mathcal{D}}q(\mathbf{x}) = \psi(\mathbf{x}_i) + h R_2(\mathbf{x}),$$

thus with  $\max_{l=1,2} \|\mathbf{R}_1^{(l)}\|_{L^\infty(\bigcup_{T \in \mathcal{M}_i} T)} \leq \|\psi\|_{2,\infty}$  and  $|R_2(\mathbf{x})| \leq \|\psi\|_{2,\infty}$  for  $\mathbf{x} \in \bigcup_{T \in \mathcal{M}_i} T$ . This involves,

$$\begin{aligned} \sum_{T \in \mathcal{M}_i} \int_T \phi_i (H(\nabla \Pi_{\mathcal{D}}q) + \Pi_{\mathcal{D}}q(\nabla \zeta)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} d\mathbf{x} = \\ \sum_{T \in \mathcal{M}_i} \left[ \int_T \phi_i (H(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} d\mathbf{x} + h \mathcal{J}_T \right], \end{aligned}$$

with,

$$\mathcal{J}_T = \int_T \phi_i (H \mathbf{R}_1 + R_2(\nabla \zeta)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} d\mathbf{x}.$$

Since  $0 \leq \phi_i \leq 1$  and the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c), all the terms inside the integral  $\mathcal{J}_T$  are bounded on  $T$  and thus  $\|\mathcal{J}_T\|_{L^\infty(T)} \leq |T| C$  with  $C$  independent of the mesh.

So we can introduce  $R_3$  a bounded function defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$H_i |K_i| \widehat{G}_{h,i} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} = \sum_{T \in \mathcal{M}_i} \left[ \int_T \phi_i (H(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} d\mathbf{x} \right] + h |K_i| R_3.$$

Then we write that, for  $T \in \mathcal{M}_i$  and  $\mathbf{x} \in T$ ,

$$\begin{aligned} (H(\mathbf{x})(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)(\mathbf{x})) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} = \\ (H(\mathbf{x}_i)(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)(\mathbf{x}_i)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} + (H(\mathbf{x}) - H(\mathbf{x}_i))(\nabla \psi)(\mathbf{x}_i) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} \\ + \psi(\mathbf{x}_i)((\nabla \zeta)(\mathbf{x}) - (\nabla \zeta)(\mathbf{x}_i)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix}. \quad (2.60) \end{aligned}$$

Due to the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c), there exists a bounded function  $R_4$  defined on  $\bigcup_{T \in \mathcal{M}_i} T$  such that,

$$(H(\mathbf{x})(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)(\mathbf{x})) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} = (H(\mathbf{x}_i)(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)(\mathbf{x}_i)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix} + h R_4(\mathbf{x}).$$

Since  $\sum_{T \in \mathcal{M}_i} \int_T \phi_i(x) d\mathbf{x} = |K_i|$ , we obtain,

$$\begin{aligned} H_i |K_i| \widehat{G}_{h,i} \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} = |K_i| \left[ \underbrace{(H(\mathbf{x}_i)(\nabla \psi)(\mathbf{x}_i) + \psi(\mathbf{x}_i)(\nabla \zeta)(\mathbf{x}_i)) \cdot \begin{pmatrix} \beta \\ \lambda \end{pmatrix}}_{((\nabla_s^\alpha \psi)^{(1)}(\mathbf{x}_i), (\nabla_s^\alpha \psi)^{(2)}(\mathbf{x}_i))^T} + h(R_3 + R_4) \right]. \end{aligned}$$

We can introduce  $R_{\star,i} > 0$  such that,

$$\left( H_i \widehat{G}_{h,i} - \nabla_s^\alpha \psi(\mathbf{x}_i) \right) \cdot \begin{pmatrix} \beta \\ \lambda \\ 0 \end{pmatrix} = h R_{\star,i}.$$

As mentioned previously, we can chose  $\beta = \mathbf{n}_{\Gamma_n}^{(2)}(\mathbf{x}_i) := \mathbf{n}_{\Gamma_n,i}^{(2)}$  and  $\lambda = -\mathbf{n}_{\Gamma_n}^{(1)}(\mathbf{x}_i) := -\mathbf{n}_{\Gamma_n,i}^{(1)}$  with  $\beta^2 + \lambda^2 = 1$  since  $\mathbf{n}_{\Gamma_n}$  is a unit vector. So we have the system,

$$\begin{cases} \beta^2 + \lambda^2 &= 1, \\ \beta \widehat{G}_{h,i}^{(1)} + \lambda \widehat{G}_{h,i}^{(2)} &= \frac{1}{H_i} (\beta (\nabla_s^\alpha \psi)^{(1)}(\mathbf{x}_i) + \lambda (\nabla_s^\alpha \psi)^{(2)}(\mathbf{x}_i) + h R_{\star,i}), \\ -\lambda \widehat{G}_{h,i}^{(1)} + \beta \widehat{G}_{h,i}^{(2)} &= 0. \end{cases}$$

Due to the assumptions on  $H$  (2.5b) and  $\zeta$  (2.5c), the regularity of  $\psi$ , we can deduce that there exists  $C_0$  independent of the mesh such that, for  $k = 1, 2$ ,

$$|\widehat{G}_{h,i}^{(k)}| \leq C_0. \quad (2.61)$$

We close the discussion on the case  $\mathbf{x}_i \in \bar{\Gamma}_n$  and come back to establish an upper bound of (2.55).

- Due to the regularity of  $H$ , we have, for  $\mathbf{x} \in K_i$ ,  $|H(\mathbf{x}) - H(\mathbf{x}_i)| \leq h C$  with  $C > 0$  independent of the mesh, and thus  $\forall i \in I$ ,

$$\|(H - H(\mathbf{x}_i)) \widehat{G}_{h,i}\|_{L^2(K_i)^{d+1}} \leq \sqrt{|K_i|} \max_i \|\widehat{G}_{h,i}\|_2 C h.$$

The previous computations allow us to say that  $|\widehat{G}_{h,i}^{(k)}|$  is bounded for every node  $i \in I$  and for  $k \in \{1, 2, 3\}$ . So we can write,

$$\sum_{i \in I} \|(H - H(\mathbf{x}_i)) \widehat{G}_{h,i}\|_{L^2(K_i)^{d+1}}^2 \leq |\Omega| C_1 h^2.$$

- Now we use the definition of  $\nabla_s^\alpha$  (2.1a), and we have, for all  $i \in I$  and  $\mathbf{x} \in K_i$ ,

$$\nabla_s^\alpha \psi(\mathbf{x}) - \nabla_s^\alpha \psi(\mathbf{x}_i) = (H(\mathbf{x}) \nabla \psi(\mathbf{x}) - H(\mathbf{x}_i) \nabla \psi(\mathbf{x}_i) + \psi(\mathbf{x}) \nabla \zeta(\mathbf{x}) - \psi(\mathbf{x}_i) \nabla \zeta(\mathbf{x}_i), -\alpha(\psi(\mathbf{x}) - \psi(\mathbf{x}_i))^T).$$

Due to the regularity of  $H, \psi, \zeta$ , we can write,

$$\sum_{i \in I} \|\nabla_s^\alpha \psi - \nabla_s^\alpha \psi(\mathbf{x}_i)\|_{L^2(K_i)^{d+1}}^2 \leq |\Omega| C_2 h^2.$$

- Finally, we have for any  $i \in I$ ,

$$\|(H(\mathbf{x}_i) \widehat{G}_{h,i} - \nabla_s^\alpha \psi(\mathbf{x}_i))\|_{L^2(K_i)^{d+1}}^2 = \sum_{k=1,2,3} \|(H(\mathbf{x}_i) \widehat{G}_{h,i}^{(k)} - (\nabla_s^\alpha \psi)^{(k)}(\mathbf{x}_i))\|_{L^2(K_i)}^2.$$

By (2.59), we can write,

$$\sum_{i \in I} \|(H(\mathbf{x}_i) \widehat{G}_{h,i}^{(3)} - (\nabla_s^\alpha \psi)^{(3)}(\mathbf{x}_i))\|_{L^2(K_i)}^2 \leq |\Omega| C_3 h^2.$$



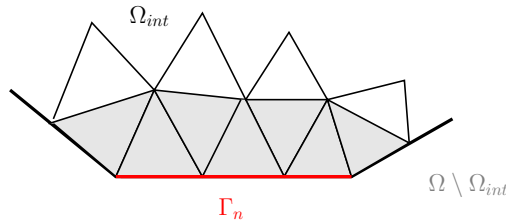
And, thanks to (2.57) and (2.58), we can write,

$$\sum_{i \in I \mid \mathbf{x}_i \notin \bar{\Gamma}_n} \sum_{k=1,2} \|(H(\mathbf{x}_i) \widehat{G}_{h,i}^{(k)} - (\nabla_s^\alpha \psi)^{(k)}(\mathbf{x}_i))\|_{L^2(K_i)}^2 \leq |\Omega_{int}| C_4 h^2.$$

And by using (2.61), and the assumptions on data, we can write,

$$\sum_{i \in I \mid \mathbf{x}_i \in \bar{\Gamma}_n} \sum_{k=1,2} \|(H(\mathbf{x}_i) \widehat{G}_{h,i}^{(k)} - (\nabla_s^\alpha \psi)^{(k)}(\mathbf{x}_i))\|_{L^2(K_i)}^2 \leq \underbrace{|\Omega \setminus \Omega_{int}|}_{\leq (|\Gamma_n| + m) h} C_4,$$

where  $m \geq 0$  is associated to the geometry of  $\Gamma_n$ . The figure below can help to understand the bound of  $|\Omega \setminus \Omega_{int}|$ .



## 2.5 Conclusion

In this chapter, some novel results and open questions have been exhibited regarding the analysis of the elliptic problem (2.2) and the GDM framework.

- In Section 2.2 a new conforming formulation the only unknown is the pressure has been exhibited, the velocity being post-processed. An error estimate has been proved thanks to the GDM framework. The  $\mathbb{P}_1$  finite element method is encompassed in this framework. A future work could be to implement this strategy in order to compare it with the mixed strategy. The comparison could be done in terms of accuracy and numerical cost by paying attention to the cost of the linear solver and, if applicable, to the chosen preconditioner. The application of this novel strategy to simulate the full model (including the hyperbolic part) implies to solve the problem of handling the boundary conditions in the full model and the localisation of the unknowns (some reconstructions have to be done due to the splitting technique).
- In Section 2.3, as far as we know, it is the first time that the Abstract GDM framework is applied to analyse *special* operators. Moreover we have seen that this framework allows us (contrary to section 2.2) to obtain error estimates even if the function reconstruction used in the *special* gradient is not the same as in the GD-Definition.
- In Section 2.4 we have seen how the two schemes used in [2] can be interpreted as Gradient Schemes. We focused, at the end of this chapter, on the proof that the  $\mathbb{P}_1/\mathbb{P}_1$  approximation satisfies the AGD-properties. A future work could be to prove that the  $\mathbb{P}_1 - \text{iso}\mathbb{P}_2/\mathbb{P}_1$  approximation also satisfies the AGD-properties.
- Finally, regarding the literature on the GDM, in [30], the authors firstly studied how a mixed mixed element method could be encompassed in the GDM framework. By the way, in this work, we have managed to ensure the  $H_{\text{div}}$  conformity of our gradient discretisation without problem dependency (contrary to [30, Remark 3.1]). Additionally we have been able here to deal with a source term in  $H^{-1}(\Omega)$  and not only in  $L^2(\Omega)$ . Thus we have increased the cases where the mixed method can be seen as a GDM.

# Chapitre 3

## Energy stable and linearly well-balanced numerical schemes for the non-linear Shallow Water equations with Coriolis force

*This work has been done in collaboration with Emmanuel Audusse, Noémie Gaveau and Yohan Penel.  
To be submitted*

### Contents

---

<b>3.1</b>	<b>Introduction</b>	<b>58</b>
<b>3.2</b>	<b>Shallow water equations and geostrophic equilibrium</b>	<b>58</b>
<b>3.3</b>	<b>Discrete operators</b>	<b>60</b>
3.3.1	Definition of the mesh	60
3.3.2	Discrete operators	61
3.3.3	Mimetic properties	63
<b>3.4</b>	<b>Well-balanced and stable finite volume schemes</b>	<b>64</b>
3.4.1	Edge based 9 points schemes	64
3.4.2	Cell-based 13 points scheme	68
3.4.3	Vertex-based 25 points scheme	70
3.4.4	Time discretisation	75
<b>3.5</b>	<b>Numerical results</b>	<b>75</b>
3.5.1	River test case	75
3.5.2	Stationary vortex test case	77
3.5.3	Translated vortex test case	80
3.5.4	Water-column test case	83
<b>3.6</b>	<b>Conclusion</b>	<b>85</b>

---

### 3.1 Introduction

The question of the accuracy of numerical schemes for hyperbolic systems with source terms around stationary solutions and/or in asymptotic regimes has been a subject of great interest over the last two decades, see the seminal works [9, 36, 38] in late nineties and the reference books [13, 34] ten years later. In the context of geophysical flows and for colocated finite-volume methods applied to shallow water equations, a lot of works have been devoted to the accuracy around the so-called lake-at-rest equilibrium and more recently extended to nonzero velocity one dimensional stationary states, see [11] and references therein. But for large scales atmospheric or oceanographic flows, the relevant stationary state is the geostrophic equilibrium, see [53] for a general introduction to geophysical rotating fluid dynamics. The accuracy of colocated finite volume numerical schemes around such an equilibrium was less investigated, To our knowledge, the first work in this field is due to Bouchut, Le Sommer and Zeitlin [14], see also [16] and [18], but was fully accurate only for one-dimensional flows, as exhibited in [5]. Recently two independent works [42, 52] proposed IMplicit-EXplicit type schemes for fully nonlinear equations which are proved to be accurate near the geostrophic equilibrium but, due to their implicit part, need the solution of a global Laplacian at each time step, see also [51] for a study about the time discretization of the Coriolis term. Note that there exists also a lot of works devoted to the approximation of the Coriolis term in staggered schemes, see for example [48] for a linear analysis and [45] for the fully nonlinear case. Let us also mention the work [50] where the authors compare the dispersion relations of some mixed Finite Volume / Finite Difference and Finite Volume / Finite Element methods in a large scale oceanographic context including the Coriolis force.

In this work, we aim at designing explicit colocated finite volume schemes that are proved to be accurate around the geostrophic equilibrium and stable in the nonlinear framework. Our work is based on the ideas developed in [5] where accurate and stable Godunov type schemes were designed for the linear two-dimensional rotating wave equation but we will see in the sequel that further developments are needed to take charge of the non linear case in a conservative way. All the numerical finite volume scheme we consider in this paper belong to the AUSM family where the flux is divided in an advective part and a pressure part, see the seminal works [41] and [40] and the recent review [26]. More precisely, in Section 3.2, we first introduce the system of equations and we characterize the geostrophic equilibrium. In Section 3.3, we define some discrete operators and we prove some of their properties. Equipped with these definitions, in Section 3.4, we can define some finite volume schemes and study the two properties we are interested in : the decrease of the semi-discrete energy and the preservation of the geostrophic equilibrium in the linearized version. Moreover, in the last part of this section, we propose a conservative version of one of our scheme. Finally, in Section 3.5, we illustrate the behaviour of the schemes for some standard test cases and we exhibit a great improvement when compared to a classic finite volume scheme.

### 3.2 Shallow water equations and geostrophic equilibrium

Let  $\Omega$  be an open bounded domain of  $\mathbb{R}^2$  and let  $T > 0$ . The nonlinear Shallow Water equations with Coriolis force formulated on  $\Omega \times (0, T)$  read :

$$\begin{cases} \partial_t h + \operatorname{div}(h\mathbf{u}) = 0, \\ \partial_t(h\mathbf{u}) + \operatorname{div}(h\mathbf{u} \otimes \mathbf{u}) + h(\nabla\phi + \omega\mathbf{u}^\perp) = 0, \end{cases} \quad (3.1)$$

where  $h$  is the water height and  $\mathbf{u} = (u_x, u_y)$  the horizontal velocity,  $\mathbf{u}^\perp = (-u_y, u_x)$  denoting its orthogonal vector in the  $(x, y)$  plane. The Coriolis force is accounted for in the momentum equations

through the angular speed  $\omega$ . Following [20, 43], the pressure forces appear under a non conservative form through the scalar potential  $\phi = gh$ , where  $g$  is the standard gravity constant. For the sake of simplicity, a flat topography is considered in the present work, but the proposed approaches naturally extended to varying bottoms by taking  $\phi = g(h + b)$  where  $b$  denotes the topography. It is also easy to extend the method to a varying Coriolis parameter in the  $\beta$ -plane approximation since all the modifications we introduce in this paper are purely local.

It is well-known that the total energy associated to the system (3.1) decomposes as  $E = \mathcal{P} + \mathcal{K}$  where

$$\mathcal{P} = \frac{1}{2}gh^2 \quad \text{and} \quad \mathcal{K} = \frac{1}{2}h\|\mathbf{u}\|^2$$

stand respectively for potential and kinetic energies. We recall that the energy  $E$  plays the role of a mathematical entropy associated to the hyperbolic system (3.1) and regular solutions satisfy the following conservation law

$$\partial_t E + \operatorname{div} \left[ \left( \phi + \frac{1}{2}\|\mathbf{u}\|^2 \right) h\mathbf{u} \right] = 0. \quad (3.2)$$

whereas for discontinuous solution, the total energy is only non-increasing in time.

When developing numerical methods, main objectives are accuracy and stability. To get stability, a crucial objective is to build numerical approximations satisfying a discrete counterpart of (3.2) that ensures that the discrete energy is nonincreasing. To achieve this, a general strategy is to consider a sufficient amount of numerical diffusion in the scheme. But in some physical contexts such as low Froude number regimes or near specific stationary states, these diffusive terms may considerably degrade the accuracy of the approximations and specific schemes are needed. Here we are interested in flows around the geostrophic balance

$$\nabla\phi + \omega\mathbf{u}^\perp = 0, \quad \operatorname{div}\mathbf{u} = 0 \quad (3.3)$$

To address such an issue, based on the study for the linear case [5], we propose a numerical approach involving discrete versions of these equilibria in the numerical fluxes. As a preliminary step, the strategy can be understood at the continuous level by investigating how the model (3.1) behaves with respect to some generic perturbations  $(\mathbf{q}, \pi)$  :

$$\begin{cases} \partial_t h + \operatorname{div}(h\mathbf{u} - \mathbf{q}) = 0, \\ \partial_t(h\mathbf{u}) + \operatorname{div}(\mathbf{u} \otimes (h\mathbf{u} - \mathbf{q})) + (h\nabla\phi - \nabla\pi) + \omega(h\mathbf{u} - \mathbf{q})^\perp = 0. \end{cases} \quad (3.4)$$

or, in conservative form

$$\begin{cases} \partial_t h + \operatorname{div}(h\mathbf{u} - \mathbf{q}) = 0 \\ \partial_t(h\mathbf{u}) + \operatorname{div}(\mathbf{u} \otimes (h\mathbf{u} - \mathbf{q})) + \nabla \left( \frac{1}{2}gh^2 - \pi \right) + \omega(h\mathbf{u} - \mathbf{q})^\perp = 0. \end{cases} \quad (3.5)$$

where  $\mathbf{q}$  and  $\pi$  can be respectively seen as (small) perturbations with respect to the flow rate and to the hydrostatic pressure.

The solutions to the modified equations (3.4) satisfy the following energy balance :

$$\partial_t E + \operatorname{div} \left[ \left( \phi + \frac{1}{2}\|\mathbf{u}\|^2 \right) (h\mathbf{u} - \mathbf{q}) - \pi\mathbf{u} \right] = -\mathbf{q} \cdot (\nabla\phi + \omega\mathbf{u}^\perp) - \pi \operatorname{div}\mathbf{u}, \quad (3.6)$$

which motivates a choice for  $\mathbf{q}$  and  $\pi$  involving resp. the quantities  $\nabla\phi + \omega\mathbf{u}^\perp$  and  $\operatorname{div}\mathbf{u}$ . Let us remark that these quantities govern the geostrophic equilibrium (3.3) associated to System (3.1) linearized around the steady state  $(\tilde{h}, \tilde{\mathbf{u}}) = (h_0, 0)$  for a constant  $h_0$  :

$$\begin{cases} \partial_t h = -h_0 \operatorname{div}\mathbf{u}, \\ \partial_t \mathbf{u} = -(\nabla\phi + \omega\mathbf{u}^\perp). \end{cases}$$

From a numerical point of view, diffusion terms are thus expected to have regularizing effects in the sense that they allow to recover a discrete counterpart of (3.6). Moreover, such terms are intended to vanish close to the geostrophic equilibrium, which must improve the quality of the approximations in this regime. The rest of the article is then devoted to the presentation of different ways to implement this idea in a discrete setting.

### 3.3 Discrete operators

#### 3.3.1 Definition of the mesh

Let us first introduce some generic notations related to the discretization of the equations. We consider a tessellation  $\mathbb{T}$  of the computational domain  $\Omega \subset \mathbb{R}^2$  made of non-overlapping rectangular cells of sizes  $\Delta x$ ,  $\Delta y$ . The set of all edges of the mesh is denoted by  $\mathcal{E}$  and the set of vertices by  $\mathcal{V}$ .

- A generic cell of  $\mathbb{T}$  is denoted by  $K$ , its area by  $m_K$  and its boundary by  $\partial K$ . A given quantity  $\Phi$  located on  $K$  is numbered  $\Phi_{i,j}$ .
- A generic edge of  $\mathcal{E}$  is denoted by  $e$  and its length by  $m_e$ . A given quantity  $\Phi$  located on  $e$  is numbered  $\Phi_{i+1/2,j}$  (respectively  $\Phi_{i,j+1/2}$ ) for y-axis (respectively x-axis) edge.
- Given a cell  $K$  and an edge  $e \in \partial K$ ,  $K_e$  is the neighbouring cell to  $e$  (other than  $K$ ) and  $\mathbf{n}_{e,K}$  the outward normal pointing to  $K_e$ .
- A generic vertex of  $\mathcal{V}$  is denoted by  $v$ . A given quantity  $\Phi$  located on  $v$  is numbered  $\Phi_{i+1/2,j+1/2}$ .

Notations are pictured on Figs. 3.1.

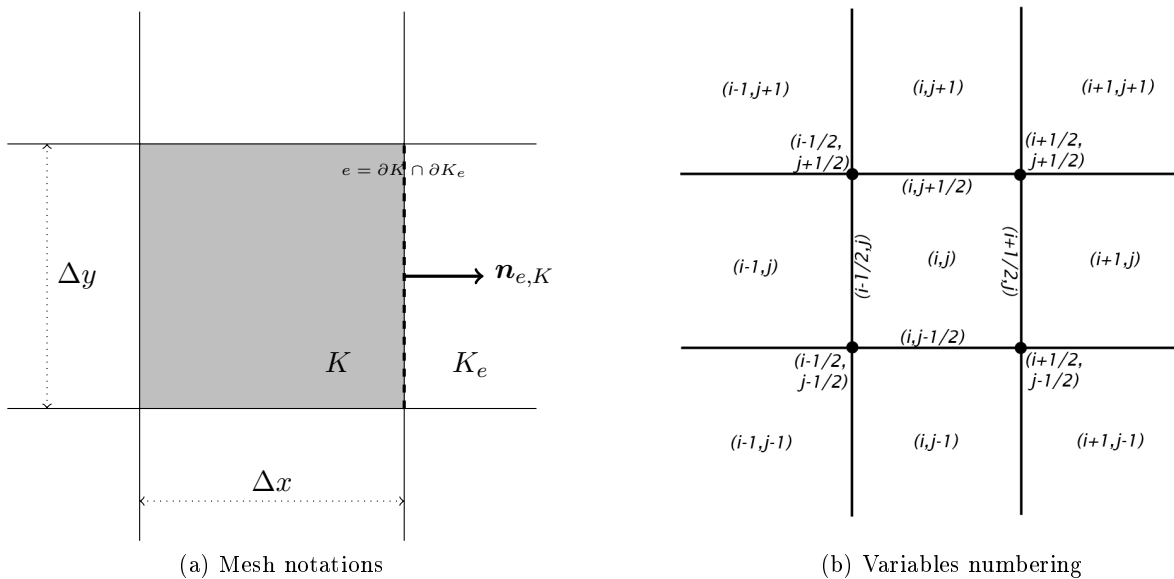


FIGURE 3.1 – Geometric notations

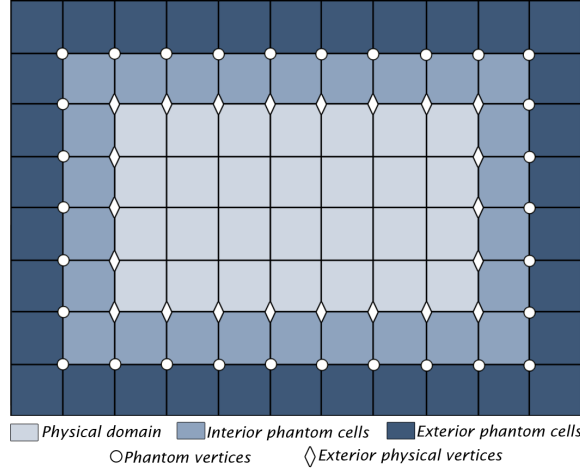


FIGURE 3.2 – Physical and phantom domains

### 3.3.2 Discrete operators

Equipped with these geometrical settings, we can now introduce discrete operators that will be needed to construct numerical schemes. Since we only consider collocated finite volume schemes, all the unknowns are defined on the cells  $K \in \mathbb{T}$ . But we will see in the next sections that some other quantities (including the numerical diffusion terms) need to be computed on the edges  $e \in \mathcal{E}$  or at the vertices  $v \in \mathcal{V}$ . Then we need to define discrete operators from cells to edges (and vice-versa) and from cells to vertices (and vice-versa) – see Figure 3.1. In the following definitions, the notations  $X_i^j(\phi_j)$  means that the operator  $X$  is applied to a quantity  $\phi$  defined at the location  $j$  and computes a quantity that is defined at the location  $i$ . For example, the first operator below  $\nabla_e^K \varphi_K$  is a discrete gradient operator that is defined for quantities that are defined on a cell  $K$  and that allows to construct a consistent gradient on an edge  $e$ . Discrete gradient and divergence are denoted with classic notations. The notation  $f$  always denotes an algebraic reconstruction operator. Let us begin with the operators from cells to edges (and vice-versa).

$$\begin{aligned}
 \nabla_e^K \varphi_K &= \frac{m_e}{m_K} (\varphi_{Ke} - \varphi_K) \mathbf{n}_{e,K} & \nabla_K^e \varphi_e &= \frac{1}{m_K} \sum_{e \in \partial K} m_e \varphi_e \mathbf{n}_{e,K} \\
 \operatorname{div}_e^K \varphi_K &= \frac{m_e}{m_K} (\varphi_{Ke} - \varphi_K) \cdot \mathbf{n}_{e,K} & \operatorname{div}_K^e \varphi_e &= \frac{1}{m_K} \sum_{e \in \partial K} m_e \varphi_e \cdot \mathbf{n}_{e,K} \\
 f_e^K(\varphi_K) &= \frac{1}{2} (\varphi_{Ke} + \varphi_K) & f_K^e(\varphi_e) &= \frac{1}{2} \sum_{e \in \partial K} \varphi_e \cdot \mathbf{n}_{e,K} \mathbf{n}_{e,K} \\
 f_e^K(\varphi_K) &= \frac{1}{2} (\varphi_{Ke} + \varphi_K) \cdot \mathbf{n}_{e,K} \mathbf{n}_{e,K} & &
 \end{aligned}$$

Then we define operators from cells to vertices, where we use the notation  $\varphi = (\varphi, \psi)$

$$\begin{aligned}
 [\nabla_v^K \varphi_K]_{i+1/2, j+1/2} &= \frac{1}{2} \left( \frac{\varphi_{i+1, j+1} - \varphi_{i, j+1}}{\Delta x} + \frac{\varphi_{i+1, j} - \varphi_{i, j}}{\Delta x} \right. \\
 &\quad \left. \frac{\varphi_{i+1, j+1} - \varphi_{i+1, j}}{\Delta y} + \frac{\varphi_{i, j+1} - \varphi_{i, j}}{\Delta y} \right) \\
 [\operatorname{div}_v^K \varphi_K]_{i+1/2, j+1/2} &= \frac{1}{2} \left[ \frac{\varphi_{i+1, j+1} - \varphi_{i, j+1}}{\Delta x} + \frac{\varphi_{i+1, j} - \varphi_{i, j}}{\Delta x} \right] + \frac{1}{2} \left[ \frac{\psi_{i+1, j+1} - \psi_{i+1, j}}{\Delta y} + \frac{\psi_{i, j+1} - \psi_{i, j}}{\Delta y} \right] \\
 [f_v^K(\varphi_K)]_{i+1/2, j+1/2} &= \frac{\varphi_{i+1, j+1} + \varphi_{i, j+1} + \varphi_{i+1, j} + \varphi_{i, j}}{4}
 \end{aligned}$$

and from vertices to cells

$$\begin{aligned} [\nabla_K^v \varphi_v]_{i,j} &= \frac{1}{2} \left( \frac{\varphi_{i+1/2,j+1/2} - \varphi_{i-1/2,j+1/2}}{\Delta x} + \frac{\varphi_{i+1/2,j-1/2} - \varphi_{i-1/2,j-1/2}}{\Delta x} \right. \\ &\quad \left. \frac{\varphi_{i+1/2,j+1/2} - \varphi_{i+1/2,j-1/2}}{\Delta y} + \frac{\varphi_{i-1/2,j+1/2} - \varphi_{i-1/2,j-1/2}}{\Delta y} \right) \\ [\operatorname{div}_K^v \varphi_v]_{i,j} &= \frac{1}{2} \left[ \frac{\varphi_{i+1/2,j+1/2} - \varphi_{i-1/2,j+1/2}}{\Delta x} + \frac{\varphi_{i+1/2,j-1/2} - \varphi_{i-1/2,j-1/2}}{\Delta x} \right] \\ &\quad + \frac{1}{2} \left[ \frac{\psi_{i+1/2,j+1/2} - \psi_{i+1/2,j-1/2}}{\Delta y} + \frac{\psi_{i-1/2,j+1/2} - \psi_{i-1/2,j-1/2}}{\Delta y} \right] \\ [f_K^v(\varphi_v)]_{i,j} &= \frac{\varphi_{i+1/2,j+1/2} + \varphi_{i-1/2,j+1/2} + \varphi_{i+1/2,j-1/2} + \varphi_{i-1/2,j-1/2}}{4} \end{aligned}$$

We will also need a divergence operator from edges to vertices

$$[\operatorname{div}_v^e \Phi_e]_{i+1/2,j+1/2} = \frac{1}{\Delta x} (\Phi_{i+1,j+1/2} - \Phi_{i,j+1/2}) \cdot \mathbf{e}_x + \frac{1}{\Delta y} (\Phi_{i+1/2,j+1} - \Phi_{i+1/2,j}) \cdot \mathbf{e}_y.$$

Finally we need to define upwind divergence operators that will be used to discretize the transport part of the equations in order to ensure the stability of the numerical schemes. In the following definitions, the quantity

$$\varphi^\pm = \frac{1}{2}(\varphi \pm |\varphi|)$$

will refer to the positive and negative parts of any scalar function  $\varphi$ . From edges to cells, the operator reads

$$\operatorname{div}_K^{e,up}(\psi_K \otimes \varphi_e) = \frac{1}{m_K} \sum_{e \in \partial K} m_e (\psi_K(\varphi_e \cdot \mathbf{n}_{e,K})^+ + \psi_{K_e}(\varphi_e \cdot \mathbf{n}_{e,K})^-)$$

and from vertices to cells, it reads

$$\begin{aligned} \operatorname{div}_K^{v,up}(\psi_K \otimes \varphi_v)_{i,j} &= \frac{1}{2\Delta x} \left( \psi_{i,j} ((\varphi_{i+1/2,j+1/2} + \varphi_{i+1/2,j-1/2}) \cdot \mathbf{e}_x)^+ + \psi_{i+1,j} ((\varphi_{i+1/2,j+1/2} + \varphi_{i+1/2,j-1/2}) \cdot \mathbf{e}_x)^- \right. \\ &\quad \left. + \psi_{i,j} ((\varphi_{i-1/2,j+1/2} + \varphi_{i-1/2,j-1/2}) \cdot (-\mathbf{e}_x))^+ + \psi_{i-1,j} ((\varphi_{i-1/2,j+1/2} + \varphi_{i-1/2,j-1/2}) \cdot (-\mathbf{e}_x))^- \right) \\ &\quad + \frac{1}{2\Delta y} \left( \psi_{i,j} ((\varphi_{i+1/2,j+1/2} + \varphi_{i-1/2,j+1/2}) \cdot \mathbf{e}_y)^+ + \psi_{i,j+1} ((\varphi_{i+1/2,j+1/2} + \varphi_{i-1/2,j+1/2}) \cdot \mathbf{e}_y)^- \right. \\ &\quad \left. + \psi_{i,j} ((\varphi_{i+1/2,j-1/2} + \varphi_{i-1/2,j-1/2}) \cdot (-\mathbf{e}_y))^+ + \psi_{i,j-1} ((\varphi_{i+1/2,j-1/2} + \varphi_{i-1/2,j-1/2}) \cdot (-\mathbf{e}_y))^- \right). \end{aligned}$$

### 3.3.3 Mimetic properties

These discrete operators satisfy some important properties that will be used to prove some results for the numerical schemes we propose in the next section.

#### Lemma 3.3.1

We first mention some local properties about the permutation of the derivative, reconstruction and orthogonal operators. Computations are obvious.

- i)  $f_v^K(\mathbf{u}_K^\perp) = (f_v^K \mathbf{u}_K)^\perp$  and  $f_K^v(\mathbf{u}_K^\perp) = (f_K^v \mathbf{u}_K)^\perp$
- ii)  $-\operatorname{div}_v^e((f_e^K(\mathbf{u}_K^\perp))^\perp) = \operatorname{div}_v^K \mathbf{u}_K$
- iii)  $\operatorname{div}_K^v(f_v^K \mathbf{u}_K) = f_K^v(\operatorname{div}_v^K \mathbf{u}_K)$

We now define the three scalar products

$$\langle \Phi_K, \Psi_K \rangle = \sum_{K \in \mathbb{T}} \Phi_K \cdot \Psi_K, \quad \langle \Phi_v, \Psi_v \rangle = \sum_{v \in \mathcal{V}} \Phi_v \cdot \Psi_v$$

$$\text{and } \langle \varphi_e, \Psi_e \rangle = \sum_{i,j} [(\varphi_{i+1/2,j} \cdot \mathbf{e}_x)(\Psi_{i+1/2,j} \cdot \mathbf{e}_x) + (\varphi_{i,j+1/2} \cdot \mathbf{e}_y)(\Psi_{i,j+1/2} \cdot \mathbf{e}_y)].$$

#### Lemma 3.3.2

We have the following properties for the reconstruction and orthogonal operators

- i)  $\langle \mathbf{u}_K, (f_K^e \mathbf{q}_e)^\perp \rangle = -\langle f_e^K(\mathbf{u}_K^\perp), \mathbf{q}_e \rangle$
- ii)  $\langle \mathbf{u}_K, f_K^v(\mathbf{q}_v^\perp) \rangle = -\langle f_v^K(\mathbf{u}_K^\perp), \mathbf{q}_v \rangle$
- iii)  $\langle f_K^v(f_v^K(\mathbf{u}_K^\perp)), \mathbf{u}_K \rangle = 0$
- iv)  $\langle f_K^v(f_v^K(h_k) f_v^K(\mathbf{u}_K^\perp)), \mathbf{u}_K \rangle = 0$

**Proof:** Property i) is obtained after a rearrangement of sum under the condition  $\mathbf{q}_e = 0$  on the edges on the boundary of the physical domain (see Figure 3.2) to eliminate the boundary terms. Property ii) is obtained after a rearrangement of sum under the periodic boundary conditions on  $h_K$  and  $\mathbf{u}_K$  and  $\mathbf{q}_v = 0$  on the vertices on the boundary of the interior domain (see Figure 3.2) to eliminate the boundary terms. Property iii) and lem :reconstruction :vii are obtained after a rearrangement of sum under the periodic boundary conditions on  $h_K$  and  $\mathbf{u}_K$  to eliminate the boundary terms. ■



**Lemma 3.3.3**

We have the following mimetic properties for the discrete gradient and divergence operators, including for some of them the reconstruction operators

- i)  $\langle \phi_K, \operatorname{div}_K^e \mathbf{q}_e \rangle = -\langle \nabla_e^K \phi_K, \mathbf{q}_e \rangle$
- ii)  $\langle \pi_v, \operatorname{div}_v^K \mathbf{u}_K \rangle = -\langle \nabla_v^v \pi_v, \mathbf{u}_K \rangle$
- iii)  $\langle \phi_K, \operatorname{div}_K^v \mathbf{q}_v \rangle = -\langle \nabla_v^K \phi_K, \mathbf{q}_v \rangle$
- iv)  $\langle f_K^e(\nabla_e^K \phi_K), h_K \mathbf{u}_K \rangle = -\langle \phi_K, \operatorname{div}_K^e(f_K^e(h_K \mathbf{u}_K)) \rangle$
- v)  $\langle f_K^v(\nabla_v^K \phi_K), \mathbf{u}_K \rangle = -\langle \phi_K, f_K^v(\operatorname{div}_v^K \mathbf{u}_K) \rangle$
- vi)  $\langle f_K^v(f_v^K(h_K) \nabla_v^K \phi_K), \mathbf{u}_K \rangle = -\langle \phi_K, \operatorname{div}_K^v(f_v^K(h_K) f_v^K(\mathbf{u}_K)) \rangle$

**Proof:** Property i) (respectively ii)) is obtained after a rearrangement of sum under the condition  $\mathbf{q}_e = 0$  on the edges (respectively  $\pi_v = 0$  on the vertices) on the boundary of the physical domain (see Figure 3.2) to eliminate the boundary terms.

Property iii) is obtained after a rearrangement of sum under the periodic boundary conditions on  $h_K$  and  $\mathbf{u}_K$  and  $\mathbf{q}_v = 0$  on the vertices on the boundary of the interior domain (see Figure 3.2) to eliminate the boundary terms.

Property iv), v) and vi) are obtained after a rearrangement of sum under the periodic boundary conditions on  $h_K$  and  $\mathbf{u}_K$  to eliminate the boundary terms. ■

## 3.4 Well-balanced and stable finite volume schemes

When considering hyperbolic equations, classic finite volume schemes in collocated two-dimensional cartesian framework are often referred as five points schemes. Indeed the update of the quantities of interest in a cell  $K$  of the tessellation  $\mathbb{T}$  needs the computation of the fluxes through the four edges of its boundary  $\partial K$ , see Figure 3.1. For a first order scheme, the numerical flux through an edge  $e \in \mathcal{E}$  is generally computed from the values of the quantities in the two neighbouring cells. Hence, the update of the quantities in a cell  $K$  involves five cells of the tessellation  $\mathbb{T}$ . We refer to [39, 49] for more details about classic first order finite volume schemes for hyperbolic problems. In the last section of the paper, devoted to the numerical tests, such a five point scheme, typically the HLLC scheme, see [49], will be considered as a standard scheme to which we will compare the well-balanced schemes we designed in the next sections. Because of the essentially two-dimensional character of the geostrophic equilibrium, these well-balanced schemes have to involve a larger stencil. Note that to consider enlarged stencils is also a common way to design high order schemes in a finite volume framework through the MUSCL strategy, see [39, 49], and is also commonly used in diffusion problems since the computation of the numerical flux needs the reconstruction of a complete two-dimensional gradient, see [29].

### 3.4.1 Edge based 9 points schemes

In this section, we present two nine points schemes for which the geostrophic equilibrium (3.3) and the perturbation  $q$ , see (3.4), are defined at the edges  $e$  whereas the perturbation  $\pi$  is defined on the vertices  $v$ . None of them succeeds to ensure both the preservation of the linearized stationary states

and the decreasing of the nonlinear energy. nevertheless, the numerical results presented in the last section of the paper exhibit a reasonably good behaviour for all the test cases we have performed.

### Energy decreasing semi-discrete scheme

In [6], we proposed a first numerical scheme for which we prove a semi-discrete counterpart of (3.6), i.e. the semi-discrete energy is non-increasing. Note that all along the paper the term semi-discrete will refer to quantities that are discrete in space but continuous in time. We also exhibited in that work that the linearized version of the scheme failed to preserve the geostrophic balance. Using the discrete operators introduced in Section 3.3, the scheme now reads

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K^e(\mathcal{F}_e) = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \mathbf{div}_K^{e,up}(\mathbf{u}_K \otimes \mathcal{F}_e) + h_K f_K^e (\nabla_e^K \phi_K) - \nabla_K^v \pi_v = -\omega (h_K \mathbf{u}_K - f_K^e \mathbf{q}_e)^\perp, \end{cases} \quad (3.7)$$

where the mass fluxes are defined at the level of the edge  $e$  as :

$$\mathcal{F}_e = f_e^K (h_K \mathbf{u}_K) - \mathbf{q}_e, \quad (3.8)$$

with the numerical diffusion term on the flow rate

$$\mathbf{q}_e = \gamma \frac{\Lambda}{g} \max\{\Delta x, \Delta y\} (\omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K), \quad (3.9)$$

while the numerical diffusion term on the hydrostatic pressure is defined at the vertices

$$\pi_v = \nu \Lambda \max\{\Delta x, \Delta y\} f_v^K (h_K) \operatorname{div}_v^K (\mathbf{u}_K), \quad (3.10)$$

where  $\gamma$  and  $\nu$  are positive dimensionless constants, and  $\Lambda$  is a positive characteristic velocity, typically, we take

$$\Lambda = \max_{K \in \mathbb{T}} \left\{ \|\mathbf{u}_K\| + \sqrt{gh_K} \right\}.$$

**Semi-discrete energy :** We first show the scheme (3.7) ensures a discrete counterpart of (3.6) through semi-discrete mechanic energy estimates. To do so, we first need the two following lemmas, describing the evolution of potential and kinetic energies.

#### Lemma 3.4.1 (*Semi-discrete potential energy*)

We set  $\mathcal{P}_K = \frac{1}{2}g (h_K)^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt} \mathcal{P}_K + \phi_K \operatorname{div}_K^e (f_e^K (h_K \mathbf{u}_K)) - \phi_K \operatorname{div}_K^e \mathbf{q}_e = 0. \quad (3.11)$$

#### Lemma 3.4.2 (*Semi-discrete kinetic energy*)

We set  $\mathcal{K}_K = \frac{1}{2}h_K \|\mathbf{u}_K\|^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt} \mathcal{K}_K + \frac{1}{2} \mathbf{div}_K^{e,up} (\|\mathbf{u}_K\|^2 \mathcal{F}_e) + h_K \mathbf{u}_K \cdot f_K^e (\nabla_e^K \phi_K) \leq \omega \mathbf{u}_K \cdot (f_K^e \mathbf{q}_e)^\perp + \mathbf{u}_K \cdot \nabla_K^v \pi_v. \quad (3.12)$$

**Proof:** We write :

$$\begin{aligned} \frac{d}{dt} \left( \frac{1}{2} h_K \|\mathbf{u}_K\|^2 \right) &= \mathbf{u}_K \cdot \frac{d}{dt} (h_K \mathbf{u}_K) - \frac{1}{2} \|\mathbf{u}_K\|^2 \frac{d}{dt} h_K \\ &= -h_K \mathbf{u}_K \cdot f_K^e (\nabla_e^K \phi_K) - \mathbf{u}_K \cdot \mathbf{div}_K^{e,up} (\mathbf{u}_K \otimes \mathcal{F}_e) + \omega \mathbf{u}_K \cdot (f_K^e \mathbf{q}_e)^\perp \\ &\quad + \frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K^e \mathcal{F}_e + \mathbf{u}_K \cdot \nabla_K^v \pi_v. \end{aligned}$$

After some basic computations, we get the relation

$$\begin{aligned} \frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K^e \mathcal{F}_e - \mathbf{u}_K \cdot \mathbf{div}_K^{e,up} (\mathbf{u}_K \otimes \mathcal{F}_e) &= \\ - \frac{1}{2} \mathbf{div}_K^{e,up} (\|\mathbf{u}_K\|^2 \mathcal{F}_e) + \frac{1}{2m_K} \sum_{e \subset \partial K} \|\mathbf{u}_K - \mathbf{u}_{K_e}\|^2 (\mathcal{F}_e \cdot \mathbf{n}_{e,k})^- . \end{aligned} \quad (3.13)$$

The second term of the right hand side being non-positive, we get the announced result .  $\blacksquare$

**Proposition 3.4.3 (Decreasing of the semi-discrete energy)**

We define the total energy  $E_K = \mathcal{P}_K + \mathcal{K}_K$ . Then we obtain a discrete counterpart of (3.6)

$$\begin{aligned} \frac{d}{dt} \left( \sum_{K \in \mathbb{T}} m_K E_K \right) \\ \leq - \max\{\Delta x, \Delta y\} \sum_{K \in \mathbb{T}} m_K \left[ \gamma \frac{\Lambda}{g} \|\omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K\|^2 + \nu \Lambda f_v^K (h_K) (\operatorname{div}_v^K \mathbf{u}_K)^2 \right]. \end{aligned}$$

**Proof:** Gathering relations (3.11) and (3.12), we obtain the following estimate for the total energy

$$\begin{aligned} \frac{d}{dt} E_K + \frac{1}{2} \mathbf{div}_K^{e,up} (\|\mathbf{u}_K\|^2 \mathcal{F}_e) + \phi_K \operatorname{div}_K^e (f_e^K (h_K \mathbf{u}_K)) + h_K \mathbf{u}_K \cdot f_K^e (\nabla_e^K \phi_K) \\ \leq -\mathbf{u}_K \cdot (f_K^e \mathbf{q}_e)^\perp + \phi_K \operatorname{div}_K^e \mathbf{q}_e + \mathbf{u}_K \cdot \nabla_K^v \pi_v. \end{aligned}$$

Hence we get the following semi-discrete inequality using Lemma 3.3.2 and 3.3.3

$$\frac{d}{dt} \left( \sum_{K \in \mathbb{T}} m_K E_K \right) \leq \sum_{K \in \mathbb{T}} m_K \left[ -\mathbf{q}_e \cdot (\omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K) - \pi_v \operatorname{div}_v^K \mathbf{u}_K \right].$$

The results follows from the choices of the numerical diffusion terms (3.9) and (3.10).  $\blacksquare$

**Linearized well-balanced property :** We refer to the next paragraph for the details of the proof. Here it is obvious that the linearized version of the scheme can not exactly preserve the geostrophic equilibrium. Indeed the numerical diffusion term  $\mathbf{q}_e \propto \omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K$  (3.9) is defined on the edges whereas the geostrophic term  $\omega \mathbf{u}_K^\perp + f_K^e (\nabla_e^K \phi_K)$  that appears in the update of the momentum in (3.7) is defined on the cells. It follows that if we define the discrete geostrophic equilibrium on the edges, the numerical diffusion term will vanish but not the other one, and vice versa. This is due to the chosen operators which ensure  $f_K^e (f_e^K (\mathbf{u}_K^\perp)) \neq \mathbf{u}_K^\perp$ .

### Well-balanced semi-discrete scheme

For completeness, we propose here a slightly modified scheme whose linearized version is meant to preserve the geostrophic equilibrium but for which it is not possible to ensure that the semi-discrete energy is non-increasing. This semi-discrete scheme reads

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K^e \mathcal{F}_e = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \operatorname{div}_K^{e,up} (\mathbf{u}_K \otimes \mathcal{F}_e) + h_K f_K^e (\nabla_e^K \phi_K) - \nabla_K^v \pi_v = -\omega h_K f_K^e (f_e^K (\mathbf{u}_K^\perp)) + \omega (f_K^e \mathbf{q}_e)^\perp, \end{cases} \quad (3.14)$$

with definitions (3.8), (3.9) and (3.10) for the mass flux, the numerical correction on the flow rate and the numerical correction on the pressure term. We see that only the definition of the Coriolis term differs from the previous version.

**Linearized well-balanced property :** The linearized version of the scheme (3.14) reads

$$\begin{cases} \frac{d}{dt} h_K + h_0 \operatorname{div}_K^e (f_e^K (\mathbf{u}_K)) - h_0 \operatorname{div}_K^e \tilde{\mathbf{q}}_e = 0, \\ \frac{d}{dt} \mathbf{u}_K + f_K^e (\nabla_e^K \phi_K) - \nabla_K^v \tilde{\pi}_v = -\omega f_K^e (f_e^K (\mathbf{u}_K^\perp)) + \omega (f_K^e \tilde{\mathbf{q}}_e)^\perp, \end{cases} \quad (3.15)$$

with the numerical diffusion terms

$$\tilde{\mathbf{q}}_e = \gamma \frac{\Lambda}{gh_0} \max\{\Delta x, \Delta y\} \left( \omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K \right) \quad \text{and} \quad \tilde{\pi}_v = \nu \Lambda \max\{\Delta x, \Delta y\} \operatorname{div}_v^K \mathbf{u}_K.$$

#### Proposition 3.4.4

Without considering the term  $\operatorname{div}_K^e (f_e^K \mathbf{u}_K)$ , the linearized scheme (3.15) preserves the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_e = 0$ .

**Proof:** When the geostrophic balance expressed on the edges holds, i.e when

$\omega f_e^K (\mathbf{u}_K^\perp) + \nabla_e^K \phi_K = 0$ , the geostrophic term that appears in the update of the momentum obviously vanishes

$$f_K^e (\nabla_e^K \phi_K) + \omega f_K^e (f_e^K (\mathbf{u}_K^\perp)) = f_K^e \left( \nabla_e^K \phi_K + \omega f_e^K (\mathbf{u}_K^\perp) \right) = 0.$$

Moreover, thanks to Lemma 3.3.1, we have

$$\operatorname{div}_v^K \mathbf{u}_K = -\operatorname{div}_v^e ((f_e^K (\mathbf{u}_K^\perp))^\perp) = \operatorname{div}_v^e \left( \frac{(\nabla_e^K \phi_K)^\perp}{\omega} \right)$$

and obvious computations show that  $\operatorname{div}_v^e ((\nabla_e^K \phi_K)^\perp) = 0$ . ■

#### Remark 3.4.5

The term  $\operatorname{div}_K^e (f_e^K \mathbf{u}_K)$  is not strictly equal to zero when the geostrophic equilibrium holds on the edges of the considered cell. However, by enforcing the geostrophic equilibrium on enough edges of the domain it is possible to write this term on the cell  $(i, j)$  using solely velocities on cells  $(i \pm 2$  and/or  $j \pm 2)$ . Due to the complexity of the computations, we failed to prove by recurrence that this term on the cell  $(i, j)$  can be written using solely velocities on cells  $(i \pm k$  and/or  $j \pm k)$ , with  $k$  arbitrarily large. Nevertheless, the test cases presented in Section 3.5 tend to show that this term is not major and that scheme (3.14) has similar results to other well-balanced scheme. Hence, we assume in the following that scheme (3.14) can be considered as a well-balanced scheme and will be referred to as such.

**Remark 3.4.6**

As it contains a reconstruction operator on the velocity, the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_e = 0$  also contains spurious solution

$$\forall(i, j), \quad \phi_{i,j} = \text{cst} \quad \text{and} \quad \mathbf{u}_{i,j} = ((-1)^i a_j, (-1)^j b_i)^T.$$

Due to the nonlinear term on the velocity, this spurious solution does not belong to the kernel of the nonlinear scheme (3.14). However, when  $\forall j, a_j = a$  and  $\forall i, b_i = b$ , this solution does appear in the kernel of the nonlinear scheme. Numerically, these solutions do not appear in the test cases presented in Section 3.5. One can note that in the edge-based entropic scheme (3.7), there is no reconstruction operator on the velocity and it is difficult to formulate a spurious solution due to the fact that such solution has to verify :

$$\begin{cases} \omega f_e^K(\mathbf{u}_K^\perp) + \nabla_e^K \phi_K = 0 \\ \omega \mathbf{u}_K^\perp + f_K^e(\nabla_e^K \phi_K) = 0. \end{cases}$$

Numerically, simulations in Section 3.5 do not show such solutions.

**Semi-discrete energy :** The computations are similar to the previous scheme but the Coriolis term is no more the same and some computations show that

$$\langle f_K^e(f_e^K(\mathbf{u}_K^\perp)), \mathbf{u}_K \rangle = \frac{1}{2} \sum_{i,j} \left( v_{i,j} \frac{u_{i,j+1} + u_{i,j-1}}{2} - u_{i,j} \frac{v_{i+1,j} + v_{i-1,j}}{2} \right).$$

Hence, it remains a term with no *a priori* sign on the right hand side of the semi-discrete energy inequality. Note that numerical simulations presented in Section 3.5 tend to show that this term is usually negligible with respect to the non positive correction terms.

**3.4.2 Cell-based 13 points scheme**

In this section, we slightly enlarged the stencil in order to propose a scheme for which we are able to prove that the semi-discrete energy is non-increasing and that the linearized version does preserve the geostrophic equilibrium. To do so, we define the geostrophic equilibrium and the perturbation  $q$  on the cells  $K$  and the perturbation  $\pi$  on the vertices  $v$  and finally obtain a scheme with a thirteen points stencil.

The semi-discrete scheme reads

$$\begin{cases} \frac{d}{dt} h_K + \text{div}_K^e(f_e^K \mathcal{F}_K) = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \text{div}_K^{e,up}(\mathbf{u}_K \otimes f_e^K \mathcal{F}_K) + h_K \nabla_K^e(f_e^K \phi_K) - \nabla_K^v \pi_v = -\omega (h_K \mathbf{u}_K - \mathbf{q}_K)^\perp, \end{cases} \quad (3.16)$$

where the mass flux  $f_e^K \mathcal{F}_K$  is defined on the edge  $e$  thanks to a reconstruction operator applied to a modified momentum term initially defined on the cells  $K$  by :

$$\mathcal{F}_K = h_K \mathbf{u}_K - \mathbf{q}_K, \quad (3.17)$$

with the numerical diffusion term

$$\mathbf{q}_K = \gamma \frac{\Lambda}{g} \max\{\Delta x, \Delta y\} \left( \omega \mathbf{u}_K^\perp + \nabla_K^e(f_e^K \phi_K) \right). \quad (3.18)$$

The numerical diffusion term based on the velocity divergence is the same as for the nine points scheme, see (3.10).

**Semi-discrete energy :** As for the first 9 points scheme, we are able to establish a discrete counterpart of (3.6) for the 13 points scheme (3.16). We do not detail the computations since they are very similar to the ones for the 9 points scheme in Proposition 3.4.3 and even simpler since the Coriolis term is defined on the cells and then obviously vanishes when multiplied by  $\mathbf{u}_K$ .

**Proposition 3.4.7**

*The 13- points scheme (3.16) satisfies the semi-discrete energy inequality*

$$\frac{d}{dt} \left( \sum_{K \in \mathbb{T}} m_K E_K \right) \leq \sum_{K \in \mathbb{T}} m_K \left[ -\mathbf{q}_K \cdot (\omega \mathbf{u}_K^\perp + \nabla_K^e (f_e^K \phi_K)) - \pi_v \operatorname{div}_v^K \mathbf{u}_K \right].$$

**Linearized well-balanced property :** The linearized version of the scheme (3.16) reads

$$\begin{cases} \frac{d}{dt} h_K + h_0 \operatorname{div}_K^e (f_e^K \mathbf{u}_K) - h_0 \operatorname{div}_K^e (f_e^K \tilde{\mathbf{q}}_K) = 0, \\ \frac{d}{dt} \mathbf{u}_K + \nabla_K^e (f_e^K \phi_K) - \nabla_K^v \tilde{\pi}_v = -\omega (\mathbf{u}_K - \tilde{\mathbf{q}}_K)^\perp, \end{cases} \quad (3.19)$$

with the numerical corrections

$$\tilde{\mathbf{q}}_K = \gamma \frac{\Lambda}{gh_0} \max\{\Delta x, \Delta y\} \left( \omega \mathbf{u}_K^\perp + \nabla_K^e (f_e^K \phi_K) \right) \quad \text{and} \quad \tilde{\pi}_v = \nu \Lambda \max\{\Delta x, \Delta y\} \operatorname{div}_v^K \mathbf{u}_K.$$

**Proposition 3.4.8**

*The linearized scheme (3.19) preserves the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_K = 0$ .*

**Proof :** When the geostrophic balance expressed on the cells holds, i.e when

$$\omega \mathbf{u}_K^\perp + \nabla_K^e (f_e^K (\phi_K)) = 0, \text{ we get}$$

$$\operatorname{div}_K^e (f_e^K \mathbf{u}_K) = \frac{1}{\omega} \operatorname{div}_K^e \left( f_e^K \left( (\nabla_K^e (f_e^K \phi_K))^\perp \right) \right)$$

and

$$\operatorname{div}_v^K \mathbf{u}_K = \frac{1}{\omega} \operatorname{div}_v^K \left( (\nabla_K^e (f_e^K \phi_K))^\perp \right),$$

obvious computations show that these two quantities vanish. Hence the geostrophic balance expressed on the cells is included in the discrete kernel of the system. ■

**Remark 3.4.9**

As it contains a reconstruction operator on the water depth, the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_K = 0$  also contains spurious checkerboard solution

$$\forall(i, j), \quad \mathbf{u}_{i,j} = 0 \quad \text{and} \quad \begin{cases} \phi_{2i,2j} = a, \\ \phi_{2i,2j+1} = b, \\ \phi_{2i+1,2j} = c, \\ \phi_{2i+1,2j+1} = d. \end{cases}$$

Since it is constructed on a zero velocity mode, this spurious solution does belong to the kernel of the nonlinear scheme (3.16) but remains bounded in time. Numerically, spurious solutions appear in the test cases 3.5.1, 3.5.3 and 3.5.4.

**3.4.3 Vertex-based 25 points scheme**

The previous 13 points scheme (3.16) ensures the preservation of a discrete geostrophic equilibrium and the decreasing of the semi-discrete energy. But its kernel contains spurious solutions and we do not succeed to write a conservative version of the discrete pressure term. To address this latter issue, here, we define the discrete geostrophic equilibrium and the perturbation  $q$  on the vertices  $v$  and then the perturbation  $\pi$  on the cells  $K$ . It leads to a scheme with a 25 points stencil.

**Non-conservative semi-discrete scheme**

With the discrete operators introduced in Section 3.3, the semi-discrete scheme reads :

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K^v \mathcal{F}_v = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \operatorname{div}_K^{v,up}(\mathbf{u}_K \otimes \mathcal{F}_v) + h_K f_K^v(\nabla_v^K \phi_K) - f_K^v(\nabla_v^K \pi_K) = -\omega (h_K f_K^v(f_v^K \mathbf{u}_K) - f_K^v \mathbf{q}_v)^\perp, \end{cases} \quad (3.20)$$

where the interface fluxes are defined at the level of the vertices  $v$  as :

$$\mathcal{F}_v = f_v^K (h_K \mathbf{u}_K) - \mathbf{q}_v, \quad (3.21)$$

with the numerical diffusion term on the flow rate

$$\mathbf{q}_v = \gamma \frac{\Lambda}{g} \max\{\Delta x, \Delta y\} (\omega f_v^K (\mathbf{u}_K^\perp) + \nabla_v^K \phi_K), \quad (3.22)$$

while the numerical diffusion term on the hydrostatic pressure is defined at the vertices

$$\pi_K = \nu \Lambda \max\{\Delta x, \Delta y\} h_K f_K^v (\operatorname{div}_v^K \mathbf{u}_K). \quad (3.23)$$

**Remark 3.4.10**

The term  $f_K^v(\nabla_v^K \pi_K)$  enforces the 25 points stencil. All the other operators are limited to a 9 points stencil.

**Semi-discrete energy :** We first show the scheme (3.20) ensures a discrete counterpart of (3.6) through semi-discrete mechanic energy estimates. To do so, and as for the 9 points scheme, we first need the two following lemmas, describing the evolution of potential and kinetic energies.

**Lemma 3.4.11** (*Semi-discrete scheme for the potential energy*)

We set  $\mathcal{P}_K = \frac{1}{2}g(h_K)^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt}\mathcal{P}_K + \phi_K \operatorname{div}_K^v(f_v^K(h_K \mathbf{u}_K)) - \phi_K \operatorname{div}_K^v q_v = 0. \quad (3.24)$$

**Lemma 3.4.12** (*Semi-discrete scheme for the kinetic energy*)

We set  $\mathcal{K}_K = \frac{1}{2}h_K \|\mathbf{u}_K\|^2$  for  $K \in \mathbb{T}$ . Then :

$$\begin{aligned} \frac{d}{dt}\mathcal{K}_K + \frac{1}{2} \operatorname{div}_K^{v,up}(\|\mathbf{u}_K\|^2 \otimes \mathcal{F}_v) + h_K \mathbf{u}_K \cdot f_v^K(\nabla_v^K \phi_K) \\ \leq -\omega(h_K \mathbf{u}_K \cdot f_v^K(f_v^K(\mathbf{u}_K^\perp)) - \mathbf{u}_K \cdot f_v^K(\mathbf{q}_v^\perp)) + \mathbf{u}_K \cdot f_v^K(\nabla_v^K \pi_K). \end{aligned} \quad (3.25)$$

**Proof:** Using i) of Lemma 3.3.1 and the first equation of (3.20) we write :

$$\begin{aligned} \frac{d}{dt} \left( \frac{1}{2} h_K \|\mathbf{u}_K\|^2 \right) &= \mathbf{u}_K \cdot \frac{d}{dt}(h_K \mathbf{u}_K) - \frac{1}{2} \|\mathbf{u}_K\|^2 \frac{d}{dt} h_K \\ &= -h_K \mathbf{u}_K \cdot f_v^K(\nabla_v^K \phi_K) - \mathbf{u}_K \cdot \operatorname{div}_K^{v,up}(\mathbf{u}_K \otimes \mathcal{F}_v) + \omega \mathbf{u}_K \cdot (f_v^K \mathbf{q}_v)^\perp \\ &\quad + \frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K^v \mathcal{F}_v + \mathbf{u}_K \cdot f_v^K(\nabla_v^K \pi_K) - \omega h_K \mathbf{u}_K \cdot f_v^K(f_v^K(\mathbf{u}_K^\perp)). \end{aligned}$$

After some basic computations, we get the relation

$$\begin{aligned} \frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K^v \mathcal{F}_v - \mathbf{u}_K \cdot \operatorname{div}_K^{v,up}(\mathbf{u}_K \otimes \mathcal{F}_v) \\ = -\frac{1}{2} \operatorname{div}_K^{v,up}(\|\mathbf{u}_K\|^2 \otimes \mathcal{F}_v) + \frac{1}{2m_K} \sum_{e \in \partial K} \|\mathbf{u}_K - \mathbf{u}_{Ke}\|^2 ((f_e^v \mathcal{F}_v) \cdot \mathbf{n}_{e,k})^-. \end{aligned}$$

The second term of the right hand side being non-positive, we get the announced result.  $\blacksquare$

**Proposition 3.4.13** (*Decreasing of the semi-discrete energy*)

We define the total energy  $E_K = \mathcal{P}_K + \mathcal{K}_K$ . Then we obtain a discrete counterpart of (3.6)

$$\begin{aligned} \frac{d}{dt} \left( \sum_{K \in \mathbb{T}} E_K \right) \\ \leq -\max\{\Delta x, \Delta y\} \left[ \sum_{K \in \mathbb{T}} \left[ \nu \Lambda h_K (f_v^K(\operatorname{div}_v^K \mathbf{u}_K))^2 \right] + \sum_{v \in \mathbb{V}} \left[ \gamma \frac{\Lambda}{g} \|\omega f_v^K(\mathbf{u}_K^\perp) + \nabla_v^K \phi_K\|^2 \right] \right]. \end{aligned} \quad (3.26)$$

**Proof:** Gathering relations (3.24) and (3.25), we obtain the following estimate for the total energy  $E_K = \mathcal{P}_K + \mathcal{K}_K$  :

$$\begin{aligned} \frac{d}{dt} E_K + \frac{1}{2} \operatorname{div}_K^{v,up}(\|\mathbf{u}_K\|^2 \otimes \mathcal{F}_v) + \phi_K \operatorname{div}_K^v(f_v^K(h_K \mathbf{u}_K)) + h_K \mathbf{u}_K \cdot f_v^K(\nabla_v^K \phi_K) \\ \leq -\omega(h_K \mathbf{u}_K \cdot f_v^K(f_v^K(\mathbf{u}_K^\perp)) - \mathbf{u}_K \cdot f_v^K(\mathbf{q}_v^\perp)) + \phi_K \operatorname{div}_K^v q_v + \mathbf{u}_K \cdot f_v^K(\nabla_v^K \pi_K). \end{aligned}$$



By telescoping and using periodic boundary condition we get :

$$\sum_{K \in \mathbb{T}} (\mathbf{div}_K^{v,up} (\|\mathbf{u}_K\|^2 \otimes \mathcal{F}_v)) = 0.$$

Using iii) of Lemma 3.3.1 and v) of Lemma 3.3.3, we get :

$$\sum_{K \in \mathbb{T}} [\phi_K \operatorname{div}_K^v (f_v^K (h_K \mathbf{u}_K)) + h_K \mathbf{u}_K \cdot f_v^K (\nabla_v^K \phi_K)] = 0.$$

Thanks to iii) of Lemma 3.3.2 we get :

$$\sum_{K \in \mathbb{T}} [h_K \mathbf{u}_K \cdot f_v^K (f_v^K (\mathbf{u}_K^\perp))] = 0.$$

Finally, ii) and iii) of Lemma 3.3.3 and ii) of Lemma 3.3.2 leads us to the following semi-discrete inequality :

$$\Delta x \Delta y \frac{d}{dt} \left( \sum_{K \in \mathbb{T}} E_K \right) \leq -\Delta x \Delta y \left[ \sum_{K \in \mathbb{T}} [\pi_K f_v^K (\operatorname{div}_v^K \mathbf{u}_K)] + \sum_{v \in \mathbb{V}} [\mathbf{q}_v \cdot (\omega f_v^K (\mathbf{u}_K^\perp) + \nabla_v^K \phi_K)] \right].$$

With the choices (3.22) and (3.23), we finally obtain the result.  $\blacksquare$

**Linearized well-balanced property :** The linearized version of the scheme (3.20) is the following :

$$\begin{cases} \frac{d}{dt} h_K + h_0 \operatorname{div}_K^v (f_v^K \mathbf{u}_K) - h_0 \operatorname{div}_K^v \tilde{\mathbf{q}}_v = 0, \\ \frac{d}{dt} \mathbf{u}_K + f_v^K (\nabla_v^K \phi_K) - f_v^K (\nabla_v^K \tilde{\pi}_K) = -\omega (f_v^K (f_v^K \mathbf{u}_K) - f_v^K \tilde{\mathbf{q}}_v)^\perp, \end{cases} \quad (3.27)$$

where

$$\tilde{\mathbf{q}}_v = \gamma \frac{\Lambda}{gh_0} \max\{\Delta x, \Delta y\} \left( \omega f_v^K (\mathbf{u}_K^\perp) + \nabla_v^K \phi_K \right) \quad \text{and} \quad \tilde{\pi}_K = \nu \Lambda \max\{\Delta x, \Delta y\} f_v^K (\operatorname{div}_v^K \mathbf{u}_K).$$

**Proposition 3.4.14**

*The linearized scheme (3.27) preserves the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_v = 0$ .*

**Proof:** When the geostrophic equilibrium holds on vertices, i.e. when  $\omega f_v^K (\mathbf{u}_K^\perp) + \nabla_v^K \phi_K = 0$ , we get

$$f_v^K (\omega f_v^K (\mathbf{u}_K^\perp) + \nabla_v^K \phi_K) = 0 \quad \text{and} \quad \operatorname{div}_K^v (f_v^K (\mathbf{u}_K)) = \operatorname{div}_K^v \left( \frac{(\nabla_v^K \phi_K)^\perp}{\omega} \right).$$

Basic computations then yield  $\operatorname{div}_K^v (f_v^K (\mathbf{u}_K)) = 0$ . Finally using property iii) of Lemma (3.3.3), we obtain that  $f_v^K (\operatorname{div}_v^K \mathbf{u}_K) = 0$ . Hence the corrections  $\tilde{\mathbf{q}}_v$  and  $\tilde{\pi}_K$  equal zero when the geostrophic equilibrium holds on the vertices and the scheme (3.27) preserves this equilibrium. Note that the natural divergence operator  $\operatorname{div}_v^K$  fails to ensure  $\operatorname{div}_v^K \mathbf{u}_K = 0$  when the geostrophic equilibrium holds at the vertices. It explains the choice of the expression of  $\pi_K$  in (3.23).  $\blacksquare$

**Remark 3.4.15**

As it contains a reconstruction operator on both the velocity and the water depth, the discrete geostrophic equilibrium  $\tilde{\mathbf{q}}_v = 0$  also contains spurious solutions

$$\forall(i, j), \quad \begin{cases} \mathbf{u}_{2i,2j} = \mathbf{a} \\ \mathbf{u}_{2i+1,2j+1} = \mathbf{b}, \\ \mathbf{u}_{2i,2j+1} = \mathbf{c} \\ \mathbf{u}_{2i+1,2j} = -\mathbf{a} - \mathbf{b} - \mathbf{c}, \end{cases} \quad \text{and} \quad \begin{cases} \phi_{2i,2j} = \phi_{2i+1,2j+1} = a, \\ \phi_{2i,2j+1} = \phi_{2i+1,2j} = b, \end{cases}$$

Due to the nonlinear term on the velocity, the spurious mode on the velocity does not belong to the kernel of the nonlinear scheme (3.14) but the checkerboard mode on the water depth does. However, the following spurious mode on the velocity does belong to the kernel of the nonlinear scheme :

$$\forall(i, j), \quad \begin{cases} \mathbf{u}_{2i+1,j} = \mathbf{u}_0, \\ \mathbf{u}_{2i,j} = -\mathbf{u}_0, \end{cases}$$

Numerically, small oscillations appear in the test case 3.5.4.

**Conservative semi-discrete scheme**

Here we propose a conservative form of the scheme (3.20), i.e. a version of the scheme with a conservative discrete pressure term as it should be in a finite volume framework. With the discrete operators introduced in Section 3.3, the conservative semi-discrete scheme reads :

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K^v \mathcal{F}_v = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \operatorname{div}_K^{e,up} (\mathbf{u}_K \otimes f_e^v \mathcal{F}_v) + f_K^v (f_v^K (h_K) \nabla_v^K (\phi_K)) - f_K^v (\nabla_v^K \pi_K) \\ \quad = -\omega (f_K^v (f_v^K (h_K) f_v^K (\mathbf{u}_K))) - f_K^v \mathbf{q}_v^\perp, \end{cases} \quad (3.28)$$

where the interface fluxes are defined at the level of the vertices  $v$  as :

$$\mathcal{F}_v = f_v^K (h_K) f_v^K (\mathbf{u}_K) - \mathbf{q}_v, \quad (3.29)$$

and the numerical diffusion terms  $\mathbf{q}_v$  and  $\pi_K$  are defined by (3.22) and (3.23). Note the change. on the pressure term has for consequence a change on the mass flux and on the Coriolis term in order to preserve the energy balance.

**Proposition 3.4.16**

The pressure term  $f_K^v(f_v^K(h_K) \nabla_v^K \phi_K)$  is in conservative form

$$[f_K^v(f_v^K(h_K) \nabla_v^K \phi_K)]_{i,j} = \frac{g}{2} \left( \frac{\left( \frac{h_{i+1,j+1} + h_{i+1,j}}{2} \right)^2 + \left( \frac{h_{i+1,j} + h_{i+1,j-1}}{2} \right)^2}{2} - \frac{\left( \frac{h_{i-1,j+1} + h_{i-1,j}}{2} \right)^2 + \left( \frac{h_{i-1,j} + h_{i-1,j-1}}{2} \right)^2}{2} \right) \frac{2\Delta x}{2\Delta y} - \frac{\left( \frac{h_{i+1,j-1} + h_{i,j-1}}{2} \right)^2 + \left( \frac{h_{i,j-1} + h_{i-1,j-1}}{2} \right)^2}{2}.$$

**Semi-discrete energy :** We now show the scheme (3.28) ensures a discrete counterpart of (3.6). We do not detail the computations since they are similar to the previous ones.

**Lemma 3.4.17 (Semi-discrete scheme for the potential energy)**

We set  $\mathcal{P}_K = \frac{1}{2}g(h_K)^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt} \mathcal{P}_K + \phi_K \operatorname{div}_K^v(f_v^K(h_K) f_v^K(\mathbf{u}_K)) - \phi_K \operatorname{div}_K^v q_v = 0. \quad (3.30)$$

**Lemma 3.4.18 (Semi-discrete scheme for the kinetic energy)**

We set  $\mathcal{K}_K = \frac{1}{2}h_K \|\mathbf{u}_K\|^2$  for  $K \in \mathbb{T}$ . Then :

$$\begin{aligned} \frac{d}{dt} \mathcal{K}_K + \frac{1}{2} \operatorname{div}_K^{e,up} (\|\mathbf{u}\|^2 f_e^v(\mathcal{F}_v)) + \mathbf{u}_K \cdot f_K^v(f_v^K(h_K) \nabla_K^v \phi_K) \\ \leq -\omega(\mathbf{u}_K \cdot f_K^v(f_v^K(h_K) f_v^K(\mathbf{u}_K^\perp)) - \mathbf{u}_K \cdot f_K^v(\mathbf{q}_v^\perp)) + \mathbf{u}_K \cdot f_K^v(\nabla_v^K \pi_K), \end{aligned} \quad (3.31)$$

**Proposition 3.4.19 (Decreasing of the semi-discrete energy)**

We define the total energy  $E_K = \mathcal{P}_K + \mathcal{K}_K$ . It satisfies the discrete energy inequality (3.26).

**Proof:** Gathering relations (3.30) and (3.31), we obtain the following estimate for the total energy  $E_K = \mathcal{P}_K + \mathcal{K}_K$  :

$$\begin{aligned} \frac{d}{dt} E_K + \frac{1}{2} \operatorname{div}_K^{e,up} (\|\mathbf{u}\|^2 f_e^v(\mathcal{F})) + \phi_K \operatorname{div}_K^e(f_e^K(f_v^K(h_K) f_v^K(\mathbf{u}_K))) + \mathbf{u}_K \cdot f_K^v(f_v^K(h_K) \nabla_v^K \phi_K) \\ \leq -\omega(\mathbf{u}_K \cdot f_K^v(f_v^K(h_K) f_v^K(\mathbf{u}_K^\perp)) - \mathbf{u}_K \cdot f_K^v(\mathbf{q}_v^\perp)) + \phi_K \operatorname{div}_K^v q_v + \mathbf{u}_K \cdot f_K^v(\nabla_v^K \pi_K). \end{aligned}$$

We get the result using properties mentioned in Lemmas 3.3.1, 3.3.2 and 3.3.3 and the particular form of the numerical diffusion terms. ■

**Linearized well-balanced property :** The linearized version of the conservative scheme (3.28) is the same as the non-conservative scheme (3.20). Hence, we obtain the same results : when the discrete geostrophic equilibrium holds at the vertices, the corrections  $\tilde{\mathbf{q}}_v$  and  $\tilde{\pi}_K$  equal zero and the geostrophic balance expressed on vertices is included in the kernel of the numerical scheme.

### 3.4.4 Time discretisation

For the discretisation in time, fluxes are taken explicit. Nevertheless it is well known that a fully explicit discretisation of the Coriolis term leads in that case to unstable schemes, see [16]. Then we consider the following discretisation of the Coriolis term for all the presented schemes :

$$\frac{u_x^{n+1} - u_x^n}{\Delta t} = \omega(\theta_u u_y^n + (1 - \theta_u)u_y^{n+1}), \quad \frac{u_y^{n+1} - u_y^n}{\Delta t} = -\omega(\theta_v u_x^n + (1 - \theta_v)u_x^{n+1}),$$

with  $\theta_u + \theta_v \leq 1$ . Here we choose  $\theta_u = 1$  and  $\theta_v = 0$  so that the system is solved explicitly.

The time step is chosen following [5] such that

$$\Delta t^n \leq \min \left\{ \frac{2}{\omega}, \frac{\min(\Delta x, \Delta y)}{\max(\|\mathbf{u}^n\| + \sqrt{gh^n})} \right\}.$$

## 3.5 Numerical results

In the following, we present different test cases to highlight the comportment of the edge-based entropic scheme (3.7), the edge-based well-balanced scheme (3.14), cell-based scheme (3.16) and vertex-based scheme (3.20), compared to a Godunov-type scheme HLLC and the edge-based entropic scheme (3.7). The conservative vertex-based scheme (3.28) give identical results as the non-conservative version (3.20), thus this scheme is omitted of this section. We study the water depth  $h$  and the velocity vector field  $\mathbf{u} = (u, v)$ , as well as the energy of the schemes.

The numerical experiments are performed with the gravitational constant  $g = 1$  and the angular velocity  $\omega = 1$ . The mesh is defined by a  $[101 \times 101]$  Cartesian grid and the numerical diffusion coefficients  $\gamma = \nu = 0.5$ .

### 3.5.1 River test case

This test case is the counterpart of the lake-at-rest for classic shallow water equations. It consist of the simulation of a flow in a stationary state through a channel, with wall-type boundary condition for  $x = -0.5$  and  $x = 0.5$ , and periodic boundary condition for  $y = -0.5$  and  $y = 0.5$ . The initial condition (see Figure 3.3) is as follow :

$$\mathbf{u} = \begin{pmatrix} 0 \\ \epsilon \\ \frac{\epsilon}{\omega} \end{pmatrix} \quad \text{and} \quad \nabla h = \begin{pmatrix} \frac{\epsilon}{g} \\ 0 \end{pmatrix},$$

with  $\epsilon = 0.01$ .

Hence, this initial condition is both a stationary solution of the nonlinear Shallow Water equations with Coriolis (3.1) and a geostrophic equilibrium  $g\nabla h + \omega\mathbf{u}^\perp = 0$ . We aim at numerical schemes being able to preserve this solution.

In Figures 3.4a and 3.4b, the vertex-based and edge-based well-balanced schemes preserve the initial condition. Due to the wall-type boundary condition the cell-based scheme (3.16) and the edge-based entropic scheme (3.7) present some anomalies on both the water depth and the velocities.

For the cell-based scheme, interferences on  $h$  and  $\mathbf{u}$  appear. One may note that these interferences also occur naturally through the water-column test case (see Section 3.5.4) and are then preserved.

These interference are in the scheme’s kernel (see Remark 3.4.9) and thus appear as a limitation of this scheme. Yet, one can remark that the mean values of  $h$  and  $\mathbf{u}$  equal the stationary solution.

For the edge-based entropic scheme, defaults on  $h$  and  $\mathbf{u}$  also appear on the wall boundaries but are not spread to the rest of the domain. Apart from these defaults, the scheme preserve the stationary solution.

Finally, the HLLC scheme do not preserve the initial condition and will tend to the lake-at-rest state.

In Figure 3.5 the energy fully discrete is roughly preserved by our schemes (small variations except for the vertex-based), whereas it is decreasing for HLLC scheme.

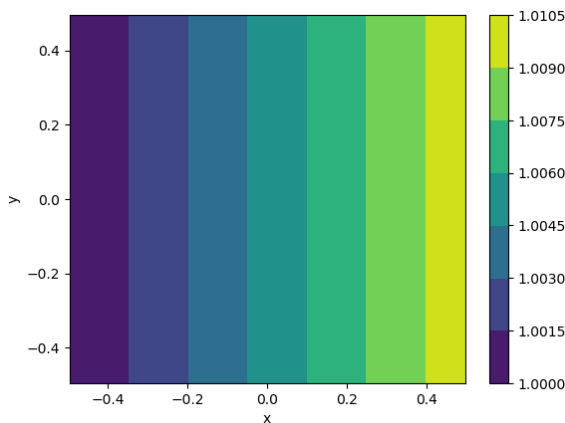


FIGURE 3.3 – Initial water depth.

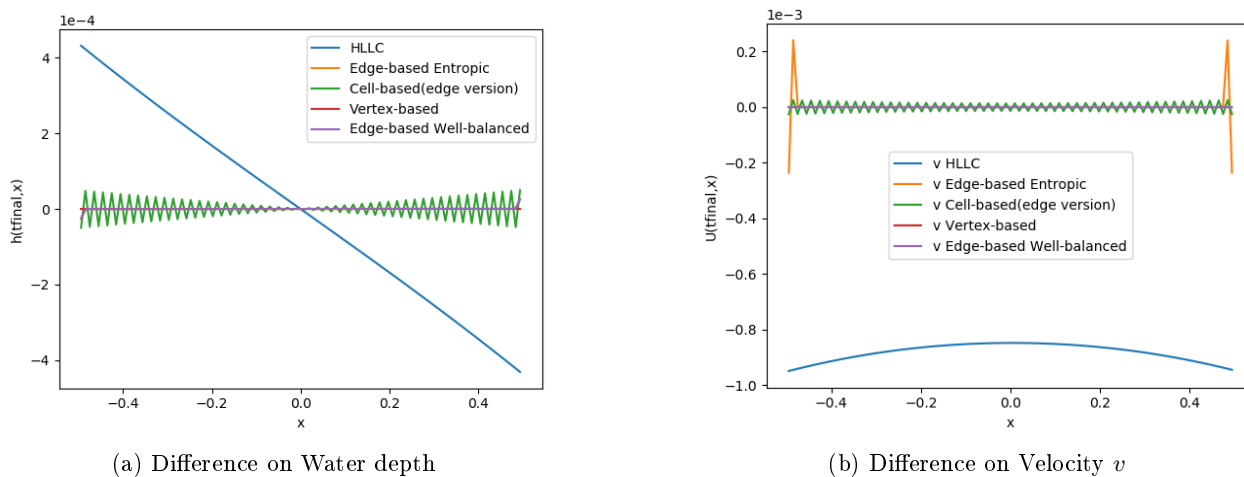


FIGURE 3.4 – Cross section in  $y = 0$  of solution at  $t = 20s$  difference to initial state.

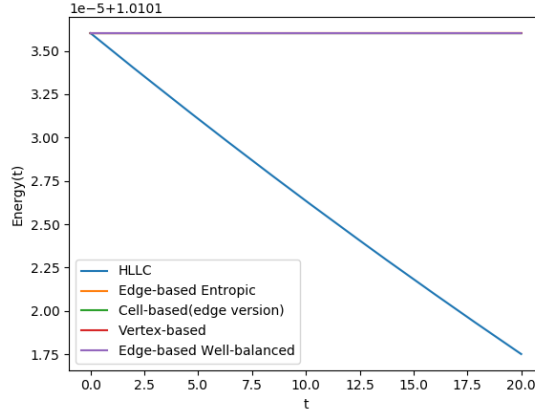


FIGURE 3.5 – Energy of the system in function of time.

### 3.5.2 Stationary vortex test case

We consider the test case introduced in [7] and defined by :

$$h(r) = \begin{cases} 1 + \frac{5\omega\epsilon}{2g}r^2 & \text{if } r \leq 0.2 \\ 1 + \frac{\omega\epsilon}{10g} - \frac{\omega\epsilon}{2g}(0.3 - 2r + 2.5r^2) + \frac{\epsilon^2}{g}(3.5 - 20r + 12.5r^2 + 4 \ln(5r)) & \text{if } 0.2 < r \leq 0.4 \\ 1 + \frac{\omega\epsilon}{5g} + \frac{\epsilon^2}{g}(4 \ln(2) - 2.5) & \text{if } r > 0.4 \end{cases}$$

$$\text{and } \mathbf{u}(r, \theta) = \begin{cases} -5\epsilon r^t(\sin(\theta), \cos(\theta)) & \text{if } r \leq 0.2 \\ -(2 - 5r)\epsilon^t(\sin(\theta), \cos(\theta)) & \text{if } 0.2 < r \leq 0.4 \\ 0 & \text{if } r > 0.4 \end{cases}$$

where  $r$  and  $\theta$  are the polar coordinates of a point of the domain. The parameter  $\epsilon$  influences the initial water velocity of the vortex, which is linked to the gradient of the water depth through the geostrophic equilibrium. This create a stationary vortex that is a stationary solution of the nonlinear Shallow Water equations with Coriolis (3.1). Let us recall that this stationary solution is different from the geostrophic equilibrium (3.3) since the definition of the water depth contains a second term that is related to the nonlinear advective term in the equations. Nevertheless, when the coefficient  $\epsilon$  introduced in the definition of the velocity field is small, the nonlinear term in the definition of the water depth is much smaller than the linear term and the stationary solution is very close to a geostrophic equilibrium. This initial condition can be seen on Figure 3.6. Finally, the domain has periodic boundary conditions. We take  $\epsilon = 0.01$  and a final simulation time of 200s. For this test case, the edge-based entropic, edge-based well-balanced, cell-based and vertex-based schemes all give similar results except the edge-based entropic one which slightly stand out. We thus choose to showcase only one scheme on the 2D graphs, the vertex-based one.

After 200s, on Figures 3.7, 3.8 and 3.9 we can notice that the vertex-based scheme indeed mainly preserves the vortex. In that amount of time, the velocities and the depth of the vortex slightly reduced, but the overall shape of the vortex is preserved. Meanwhile, we can see that the HLLC scheme failed to preserve the said shape : the velocities and the depth reduced greatly, qualitatively, stratified velocity can be seen on Figure 3.8 whereas a distorted vortex appears on Figure 3.7 hence the scheme seems

to try to reach a lake at rest state.

In order to quantify the preservation of the stationary vortex, we introduce the following relative error :

$$\frac{|\min(h_{final}) - \min(h_{init})|}{\max(h_{init}) - \min(h_{init})},$$

which track the bottom of the vortex. On Figure 3.10, we remark that the HLLC scheme quickly move away from the vortex, we can also note that the final water depth relative error made by the vertex-based scheme decreases with  $\epsilon$ , while this parameter has no influence on the error for the HLLC scheme. Indeed, for any value of  $\epsilon$ , the HLLC scheme has the same behaviour : the stationary vortex is not seen as an equilibrium and the difference of initial state thus has no effect while the scheme tries to reach the closest equilibrium in its kernel, the lake at rest. Meanwhile, the vertex-based scheme does consider the stationary vortex as an geostrophic equilibrium even if it is not a stationary solution for the scheme. Thus, when  $\epsilon$  decrease the non-linear term in the scheme become negligible and the vortex appears as a quasi-equilibrium for the scheme.

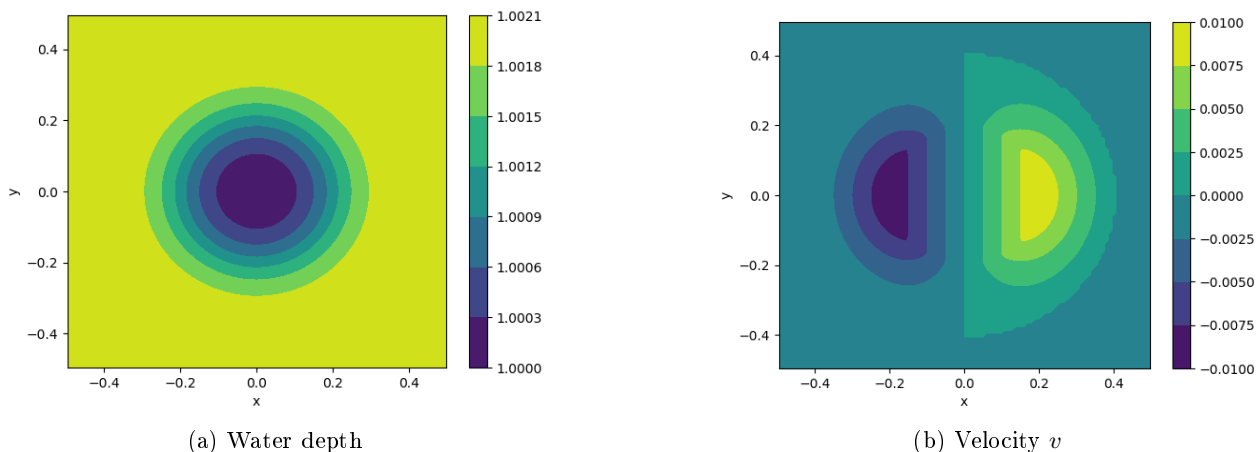


FIGURE 3.6 – Initial state of the stationary vortex.

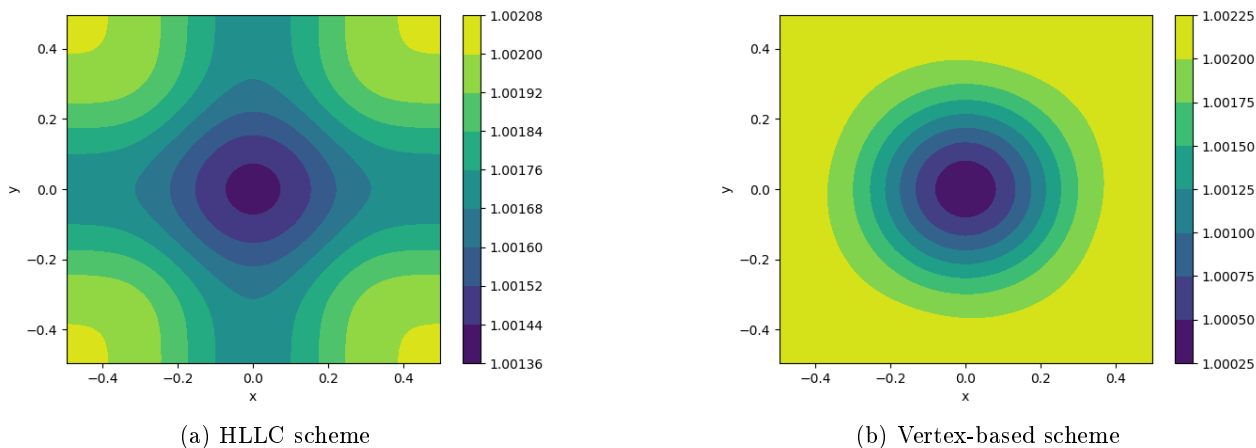


FIGURE 3.7 – Water depth at  $t = 200s$ .

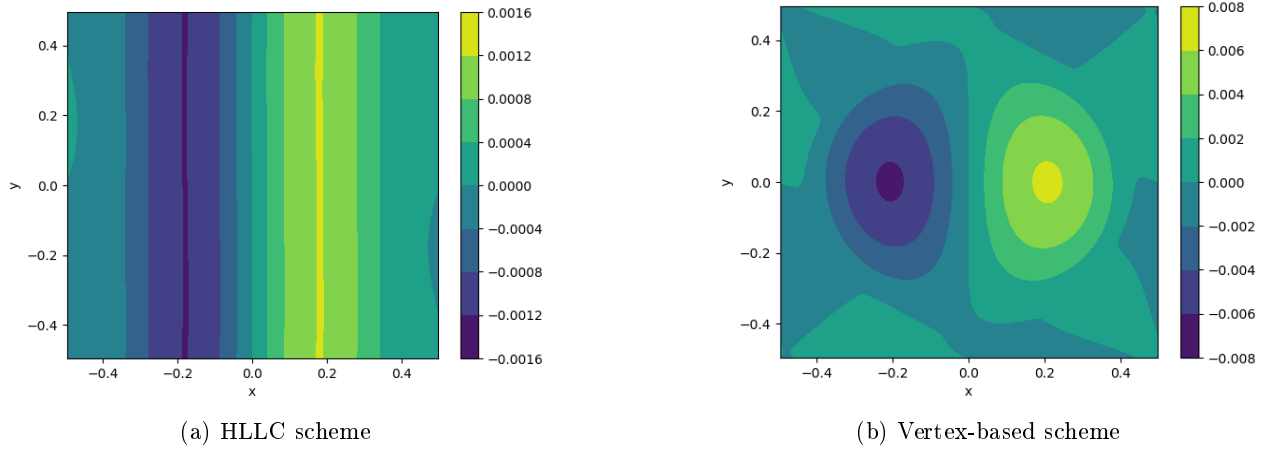


FIGURE 3.8 – Velocity  $v$  at  $t = 200s$ .

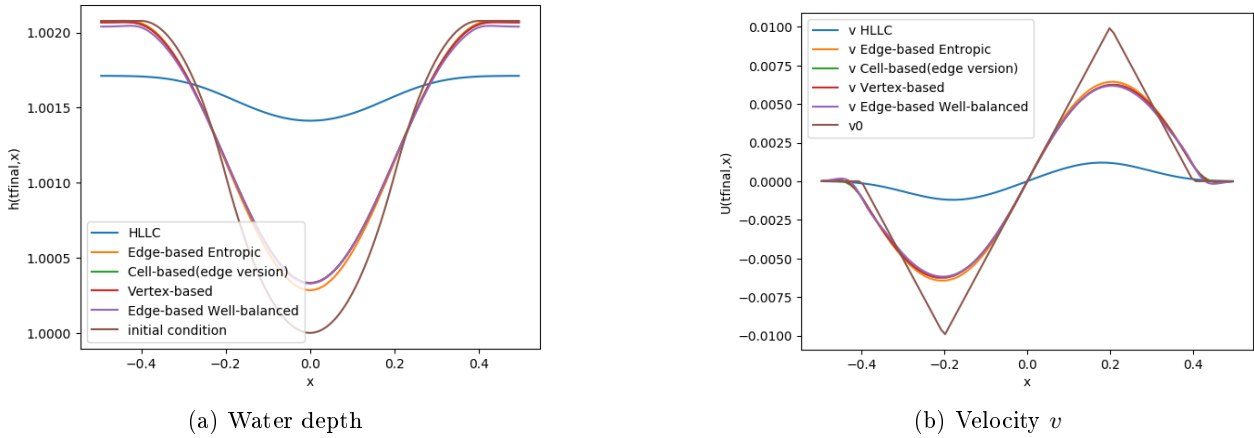


FIGURE 3.9 – Cross sections  $y = 0$  of solution at  $t = 200s$ .

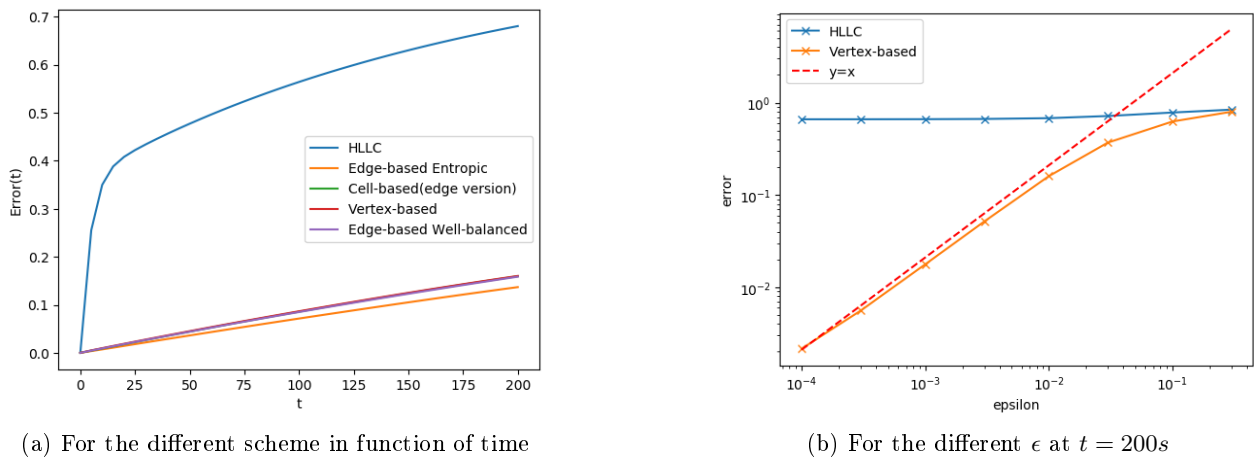
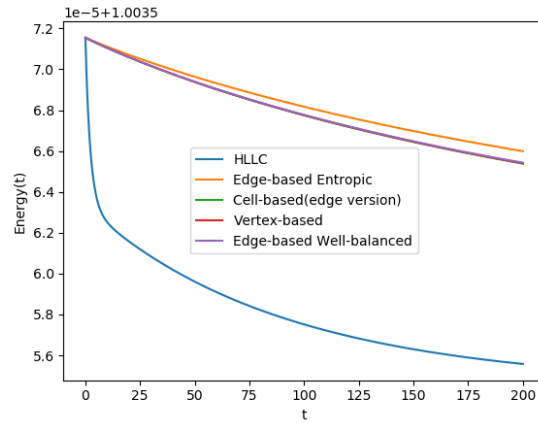


FIGURE 3.10 – Water depth relative error to initial state.



FIGURE 3.11 – Energy of the system in function of time for  $\epsilon = 0.01$ .

### 3.5.3 Translated vortex test case

This test case is the sum of two other ones : the river and the linear stationary vortex. Both are stationary solutions of the linear Shallow Water equations with Coriolis source term, and we are interested in the behaviour of our schemes when given that initial state (see Figure 3.12). As for the river test case, we enforce wall-type boundary condition for  $x = -0.5$  and  $x = 0.5$ , and periodic boundary condition for  $y = -0.5$  and  $y = 0.5$ . The slope of the river is  $\epsilon = 0.01$ , same as the vortex depth parameter.

After 20 seconds, we find in Figures 3.13 , 3.14 and 3.16 that the HLLC-scheme fails to preserve the initial state : the vortex is flattening out and the slope is decreasing. After 500s, it is brought back to the lake-at-rest. The other schemes also have a behaviour on par with the one they had on the two previous test cases and their graphs are overlapped. The vortex is preserved, except for a small depth reduction, and the slope of the river is constant, while interference appear for the cell-based scheme on Figure 3.16a. In addition, the river induces a translation of the vortex along the  $y$  axis.

On Figure 3.18, we can note that the energy is also what one could expect from this superposition of solution : it is almost constant for the vertex-, cell- and edge-based entropic/well-balanced schemes, while it drops abruptly for the HLLC one.

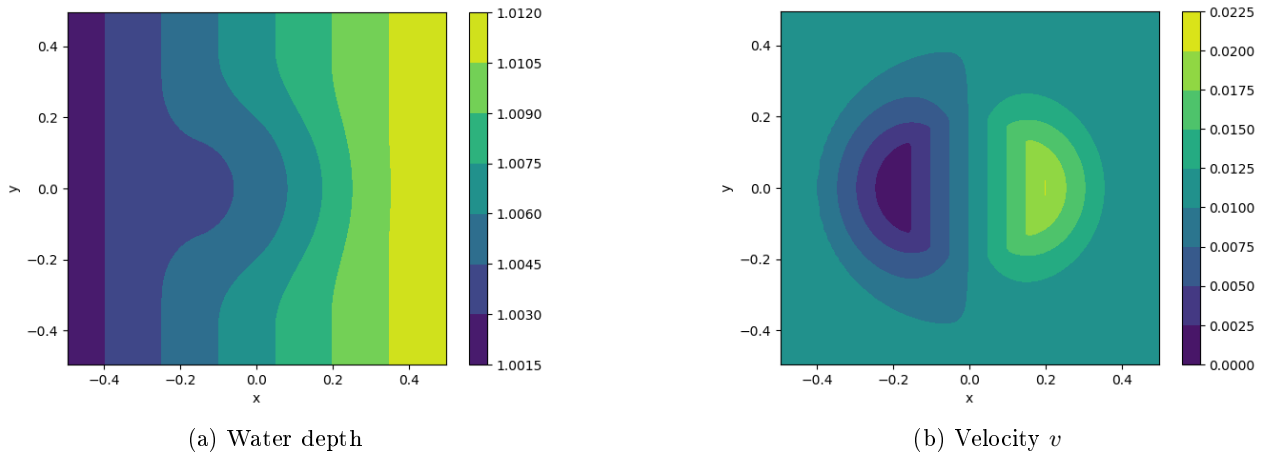


FIGURE 3.12 – Initial state of the translated vortex.

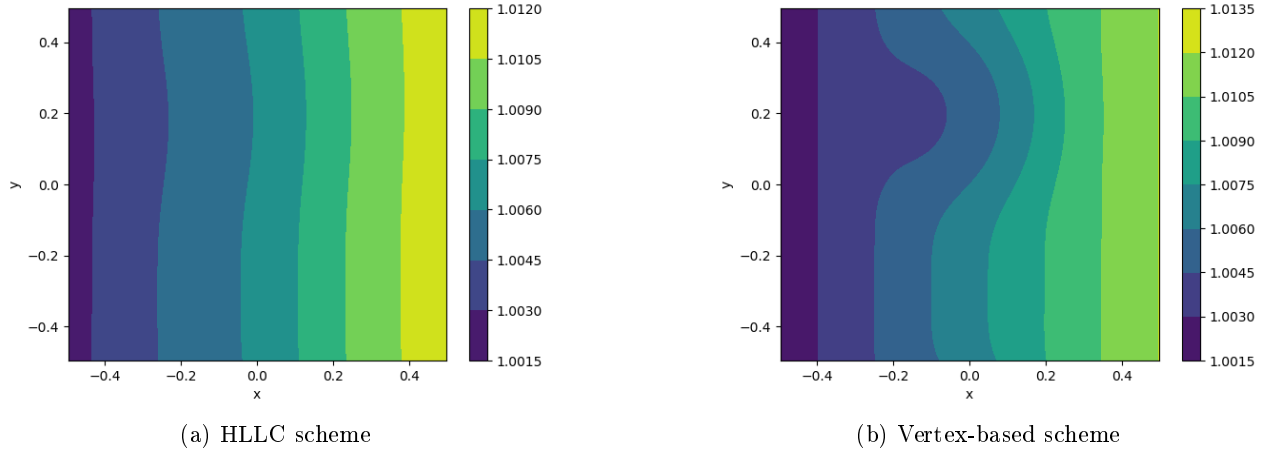
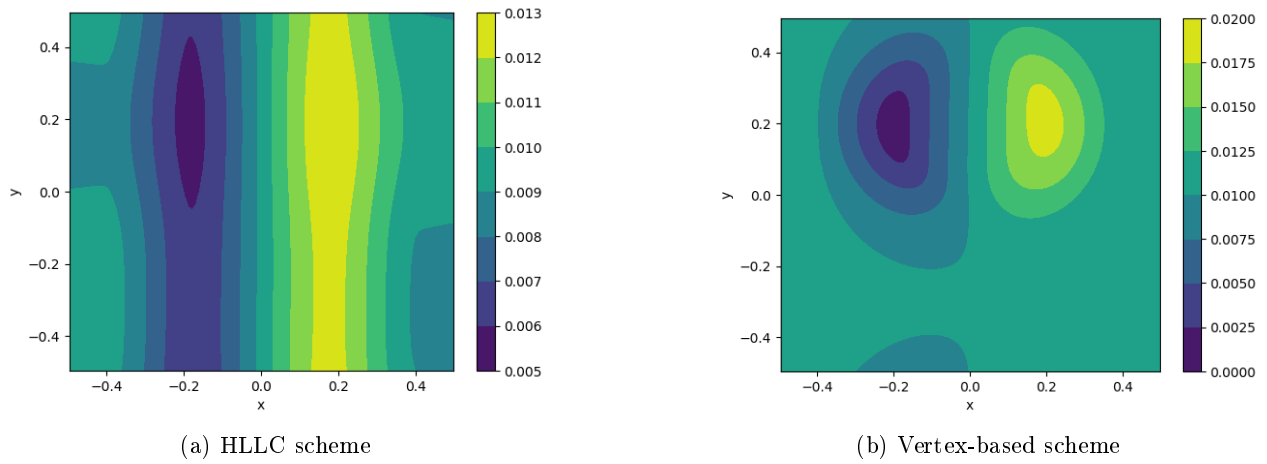
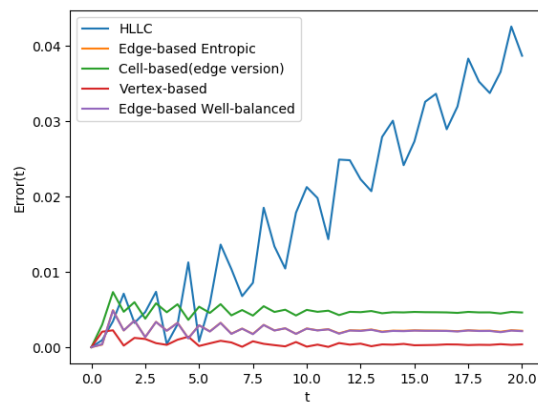
FIGURE 3.13 – Water depth at  $t = 20s$ .FIGURE 3.14 – Velocity  $v$  at  $t = 20s$ .

FIGURE 3.15 – Water depth relative error to initial state in function of time .

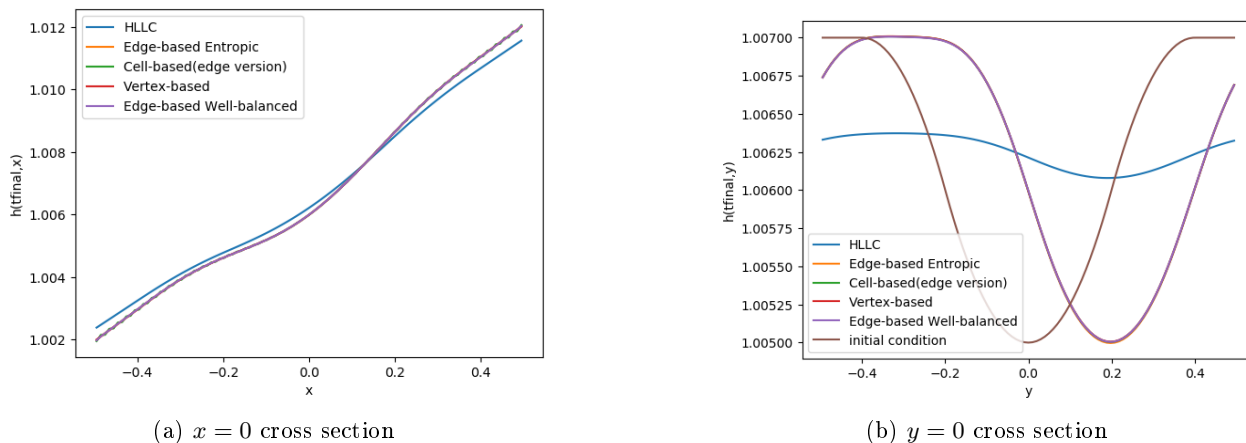


FIGURE 3.16 – Cross section of water depth at  $t = 20s$  in  $x = 0$  (a) and  $y = 0$  (b).

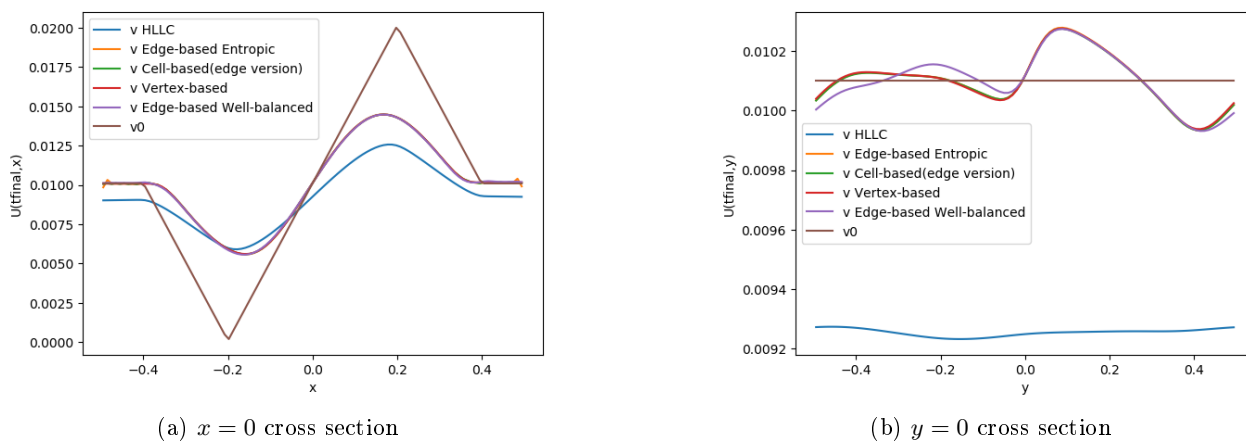


FIGURE 3.17 – Cross section of velocity  $v$  at  $t = 20s$  in  $x = 0$  (a) and  $y = 0$  (b).

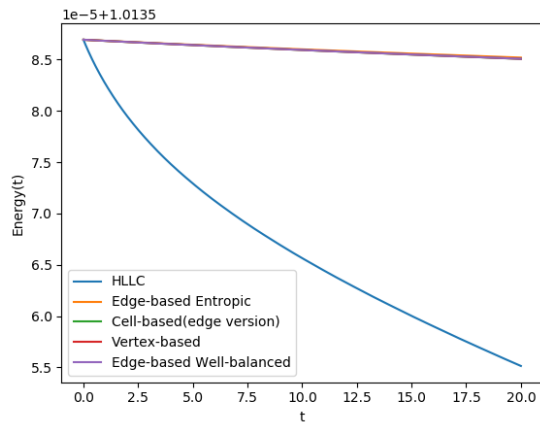


FIGURE 3.18 – Energy of the system in function of time for the different schemes.

### 3.5.4 Water-column test case

We consider a circular dam break with a radius of 1 and the domain has periodic boundary conditions. At  $t = 0$ , the velocity is zero throughout the domain and the water height is 1 except for a water column in the center of height 2.

This initial condition is very far from a geostrophic equilibrium, thus on Figure 3.22 we can see that the energy drop is sharper than previously. The edge-based well-balanced, cell-based and vertex-based schemes have similar energy, their graphs are overlapped.

On Figure 3.21, all schemes have the same overall behaviour at the beginning. However, at time goes by, the HLLC-scheme gets to a lake at rest equilibrium, while the vertex-, cell- and edge-based entropic/well-balanced schemes stabilize around another geostrophic equilibrium. The graphs of the vertex-, cell- and edge-based well-balanced schemes are overlapped. At  $t = 1s$ , we can also note that interferences appear once again for the cell-based scheme, as well as for the vertex-based one. The edge-based entropic/well-balanced schemes seem immune to it for this test case.

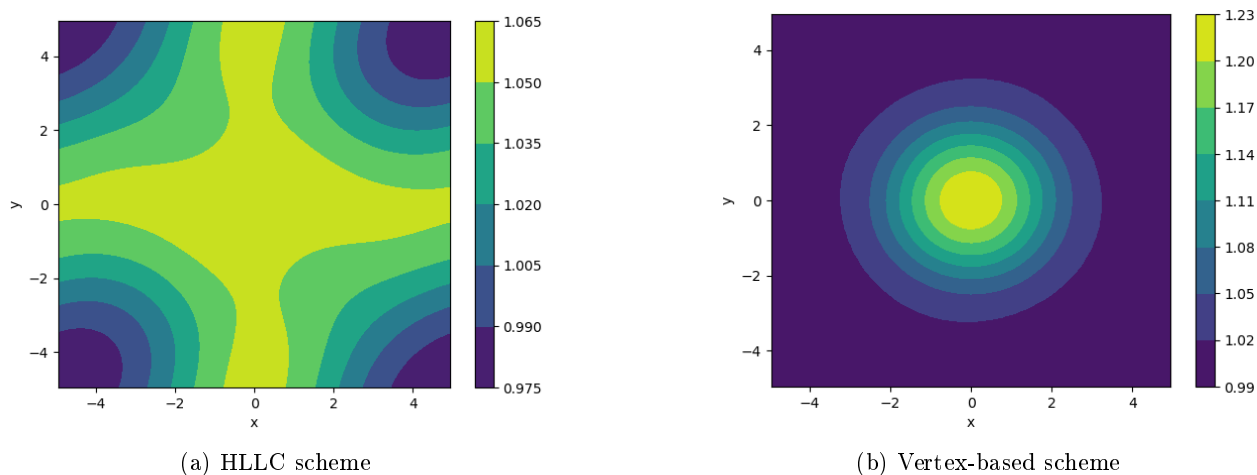


FIGURE 3.19 – Water depth at  $t = 100s$ .

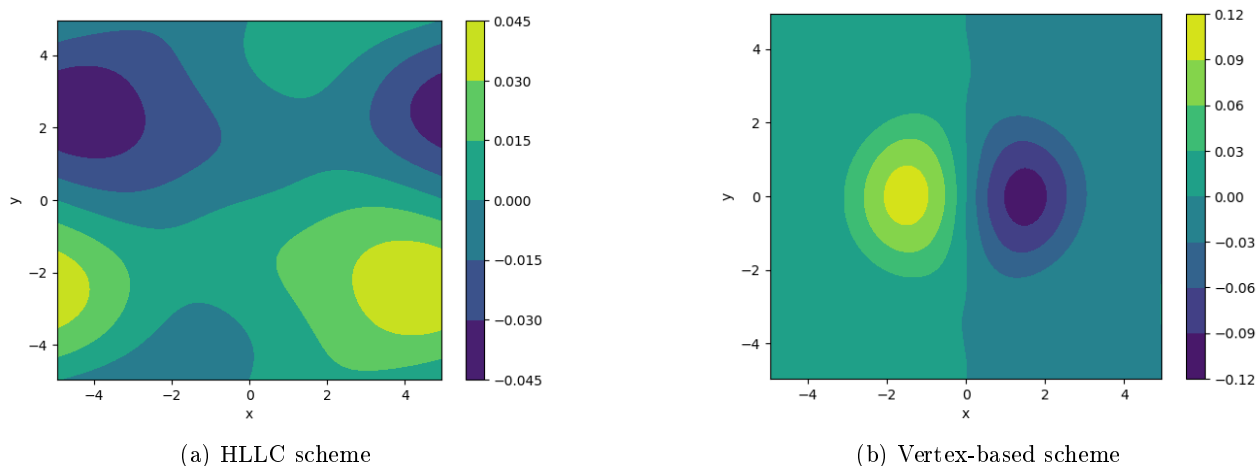
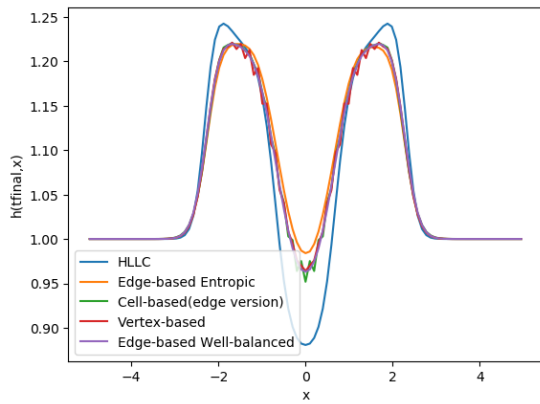
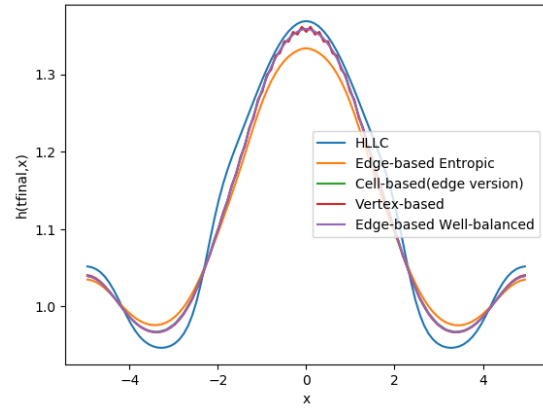


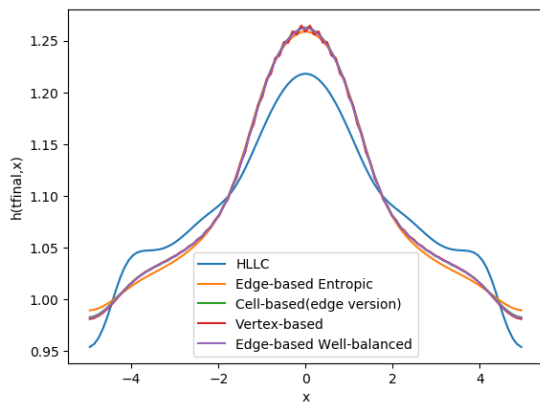
FIGURE 3.20 – Velocity  $v$  at  $t = 100s$ .



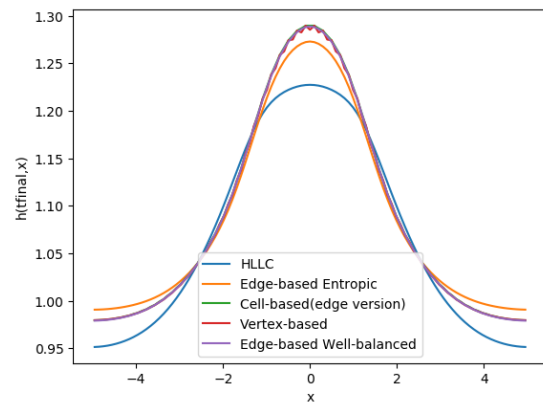
(a)  $t = 1s$



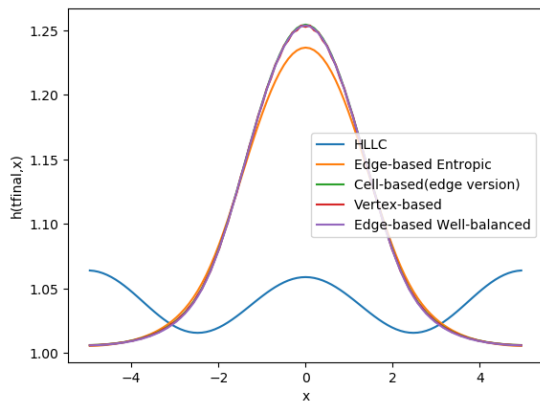
(b)  $t = 4s$



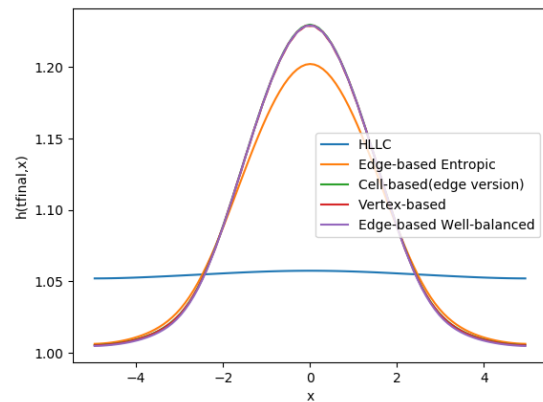
(c)  $t = 6s$



(d)  $t = 10s$



(e)  $t = 40s$



(f)  $t = 100s$

FIGURE 3.21 – Cross section in  $y = 0$  of simulation results for the different schemes at different times.

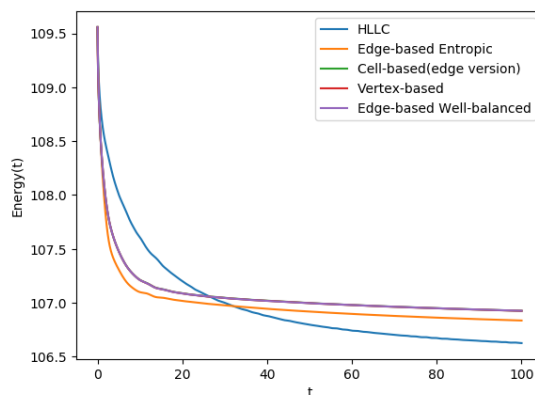


FIGURE 3.22 – Energy of the system in function of time for the different schemes.

### 3.6 Conclusion

In this work we have derived three semi-discrete schemes, namely edge- well-balanced, cell- and vertex-based for the shallow water equations with Coriolis force. For the edge-based well-balanced scheme we discretized the geostrophic equilibrium on the edges of the mesh and obtain a linearly well-balanced semi-discrete scheme. Although we failed to ensure the non-increasing energy, the numerical tests show a good short-term behaviour. For the cell-based scheme we discretized the geostrophic equilibrium on the cells of the mesh and obtain a linearly well-balanced semi-discrete scheme for which we have semi-discrete estimate. However, numerical results show that the kernel of this scheme contains spurious solutions and drive us to derived a final scheme. The vertex-based scheme has the same theoretical properties as the cell-based one and also has interference solutions but they appear less frequently.

Through different test cases, we saw that our schemes are more accurate around the geostrophic equilibrium than the HLLC one. However, spurious solutions appear for the cell- and vertex-based, further work has to be done in order to obtain a well-behaved scheme with restricted kernel. Also, the fully discrete energy of our schemes need to be studied.



# Annexe A

## Numerical approximation of the Shallow Water equations with Coriolis source term

E. AUDUSSE<sup>1</sup>, V. DUBOS<sup>2</sup>, A. DURAN<sup>3</sup>, N. GAVEAU<sup>4</sup>, Y. NASSERI<sup>5</sup>, Y. PENEL<sup>2</sup>

**Abstract** We investigate in this work a class of numerical schemes dedicated to the non-linear Shallow Water equations with topography and Coriolis force. The proposed algorithms rely on Finite Volume approximations formulated on collocated and staggered meshes, involving appropriate diffusion terms in the numerical fluxes, expressed as discrete versions of the linear geostrophic balance. It follows that, contrary to standard Finite-Volume approaches, the linear versions of the proposed schemes provide a relevant approximation of the geostrophic equilibrium. We also show that the resulting methods ensure semi-discrete energy estimates. Numerical experiments exhibit the efficiency of the approach in the presence of Coriolis force close to the geostrophic balance, especially at low Froude number regimes.

### Contents

---

<b>A.1</b>	<b>Design of the numerical schemes</b>	<b>89</b>
A.1.1	Collocated scheme	89
A.1.2	Staggered scheme	95
A.1.3	Time discretisation	98
<b>A.2</b>	<b>Numerical assessments of the schemes</b>	<b>98</b>
A.2.1	Presentation of the test cases	98
A.2.2	Numerical results	100
<b>A.3</b>	<b>Conclusion</b>	<b>100</b>

---

<sup>1</sup>LAGA, Institut Galilée, Université Sorbonne Paris Nord – 99 avenue Jean-Baptiste Clément, 93430 Villetaneuse.

<sup>2</sup>INRIA Paris - Sorbonne Université - CNRS (LJLL), team ANGE, 2 rue Simone Iff, CS 42112, 75589 Paris cedex 12.

<sup>3</sup>Institut Camille Jordan, Université Claude Bernard Lyon 1 – 43, boulevard du 11 novembre 1918, 69622 Villeurbanne.

<sup>4</sup>Institut Denis Poisson, Université D'Orléans, rue de Chartres, 45067 Orléans.

<sup>5</sup>Institut de Mathématiques de Marseille, Université d'Aix-Marseille – Technopôle Château Gombert, 13453, Marseille.



## Introduction

The question of the accuracy of numerical schemes around stationary solutions or/and in asymptotic regimes has been a subject of great interest over the last two decades, see the seminal works [9, 36, 38] in late nineties and the books [13, 34] ten years later. In the context of geophysical flows and more particularly for finite-volume methods applied to shallow water equations, a lot of works have been devoted to the accuracy around the *lake-at-rest* equilibrium and more recently extended to nonzero velocity one dimensional stationary states. The accuracy of numerical schemes around geostrophic equilibria which are of fundamental importance for large scales atmospheric or oceanographic flows was less investigated, see [53] for a general introduction to geophysical rotating fluid dynamics. To our knowledge, the first work in this field is due to Bouchut, Le Sommer and Zeitlin [14] but was fully accurate only for one-dimensional flows, as exhibited in [5] where accurate and stable Godunov type schemes were designed for linear two-dimensional shallow water equations. Recently two independent works [42, 52] proposed an IMEX type scheme for fully nonlinear equations.

Here we aim at designing explicit schemes that are proved to be accurate and stable in the nonlinear framework. Our work is mostly based on the ideas developed in [5] for the linear rotating case and in [20, 25] for the nonlinear non-rotating case. We investigate schemes applied to collocated and staggered schemes.

Let  $\Omega$  be an open bounded domain of  $\mathbb{R}^2$  and let  $T > 0$ . The nonlinear Shallow Water equations with Coriolis force formulated on  $\Omega \times (0, T)$  read :

$$\begin{cases} \partial_t h + \operatorname{div}(h\mathbf{u}) = 0, \\ \partial_t(h\mathbf{u}) + \operatorname{div}(h\mathbf{u} \otimes \mathbf{u}) + h(\nabla\phi + \omega\mathbf{u}^\perp) = 0, \end{cases} \quad (\text{A.1})$$

where  $h$  is the water height and  $\mathbf{u} = (u_x, u_y)$  the horizontal velocity<sup>1</sup>. The Coriolis force is accounted for in the momentum equations through the angular speed  $\omega$ . Following [20, 43], the pressure forces appear under a non conservative form through the scalar potential  $\phi = gh$ , where  $g$  is the standard gravity constant. For the sake of simplicity, a flat topography is considered in the present work, but the proposed approaches naturally extended to varying bottoms.<sup>2</sup>

It is well-known that the total energy of the system decomposes as  $E = \mathcal{E} + \mathcal{K}$  where

$$\mathcal{E} = \frac{1}{2}gh^2 \quad \text{and} \quad \mathcal{K} = \frac{1}{2}h\|\mathbf{u}\|^2$$

stand respectively for potential and kinetic energies. We recall that the energy  $E$  plays the role of a mathematical entropy associated to the hyperbolic system (A.1) and regular solutions satisfy the following conservation law :

$$\partial_t E + \operatorname{div} \left[ \left( \phi + \frac{1}{2}\|\mathbf{u}\|^2 \right) h\mathbf{u} \right] = 0. \quad (\text{A.2})$$

When developing numerical methods, main objectives are accuracy and stability. To get stability, a crucial objective is to build numerical approximations satisfying a discrete counterpart of (A.2) that ensures that the discrete energy is nonincreasing. To achieve this, a general strategy is to consider a sufficient amount of numerical diffusion in the scheme. But in some physical contexts such as low Froude number regimes or near specific stationary states, these diffusive terms may considerably degrade the accuracy of the approximations and specific methods are required. Here we are interested in flows around the geostrophic balance

$$\nabla\phi + \omega\mathbf{u}^\perp = 0, \quad \operatorname{div}\mathbf{u} = 0. \quad (\text{A.3})$$

<sup>1</sup> $\mathbf{u}^\perp = (-u_y, u_x)$ .

<sup>2</sup>In case of a non-flat topography  $z_b$ ,  $\phi = g(h + z_b)$ .

To address such an issue, based on the linear case [5], and on the non-linear case without Coriolis force [20, 25], we propose a numerical approach involving discrete versions of these equilibria in the numerical fluxes. As a preliminary step, the strategy can be understood at the continuous level by investigating how the model (A.1) behaves with respect to some generic perturbations  $(\mathbf{q}, \pi)$  :

$$\begin{cases} \partial_t h + \operatorname{div}(h\mathbf{u} - \mathbf{q}) = 0, \\ \partial_t(h\mathbf{u}) + \operatorname{div}(\mathbf{u} \otimes (h\mathbf{u} - \mathbf{q})) + (h\nabla\phi - \nabla\pi) + \omega(h\mathbf{u} - \mathbf{q})^\perp = 0. \end{cases} \quad (\text{A.4})$$

Hence  $\mathbf{q}$  and  $\pi$  must be respectively seen as (small) perturbations with respect to the flow rate and to the hydrostatic pressure. The solutions to the modified equations (A.4) satisfy the following energy balance :

$$\partial_t E + \operatorname{div} \left[ \left( \phi + \frac{1}{2} \|\mathbf{u}\|^2 \right) (h\mathbf{u} - \mathbf{q}) - \pi\mathbf{u} \right] = -\mathbf{q} \cdot (\nabla\phi + \omega\mathbf{u}^\perp) - \pi \operatorname{div}\mathbf{u}, \quad (\text{A.5})$$

which motivates a choice for  $\mathbf{q}$  and  $\pi$  involving resp. the quantities  $\nabla\phi + \omega\mathbf{u}^\perp$  and  $\operatorname{div}\mathbf{u}$ . Let us remark that these quantities govern the geostrophic equilibrium (A.3) associated to system (A.1) linearized around the steady state  $(\tilde{h}, \tilde{\mathbf{u}}) = (h_0, 0)$  for a constant  $h_0$  :

$$\begin{cases} \partial_t h = -h_0 \operatorname{div}\mathbf{u}, \\ \partial_t \mathbf{u} = -(\nabla\phi + \omega\mathbf{u}^\perp). \end{cases}$$

From a numerical point of view, diffusion terms are thus expected to have regularizing effects in the sense that they allow to recover a discrete counterpart of (A.5). Moreover, such terms are intended to vanish close to the geostrophic equilibrium, which must improve the quality of the approximations in this regime. Hence, the present approach can be seen as a non-linear extension of the linear case [5].

The paper is split into two parts : in the first section, we present the two geometric frameworks (collocated *vs.* staggered), we build some discrete operators with mimetic properties and we propose two well-balanced numerical schemes whose semi-discrete (in time) versions are proved to be stable. The second section is dedicated to the assessment of those schemes by means of the stationary vortex test case.

## A.1 Design of the numerical schemes

### A.1.1 Collocated scheme

#### Mesh and notations

To discretize System (A.1), we consider a tessellation  $\mathbb{T}$  of the computational domain  $\Omega \subset \mathbb{R}^2$  made of non-overlapping rectangular cells of sizes  $\Delta x$ ,  $\Delta y$ . The set of all edges of the mesh is denoted by  $\mathcal{E}$  and the set of vertices by  $\mathcal{V}$ .

- A generic cell of  $\mathbb{T}$  is denoted by  $K$ , its area by  $m_K$  and its boundary by  $\partial K$ .
- The length of any edge  $e$  is denoted by  $m_e$ .
- Given a cell  $K$  and an edge  $e \in \partial K$ ,  $K_e$  is the neighbouring cell to  $e$  (other than  $K$ ) and  $\mathbf{n}_{e,K}$  the outward normal pointing to  $K_e$ .
- Given a vertex  $\mathbf{x}$ , a dual cell is built by joining the centers of the neighbouring cells. Its area is denoted by  $m_{\mathbf{x}}$ .

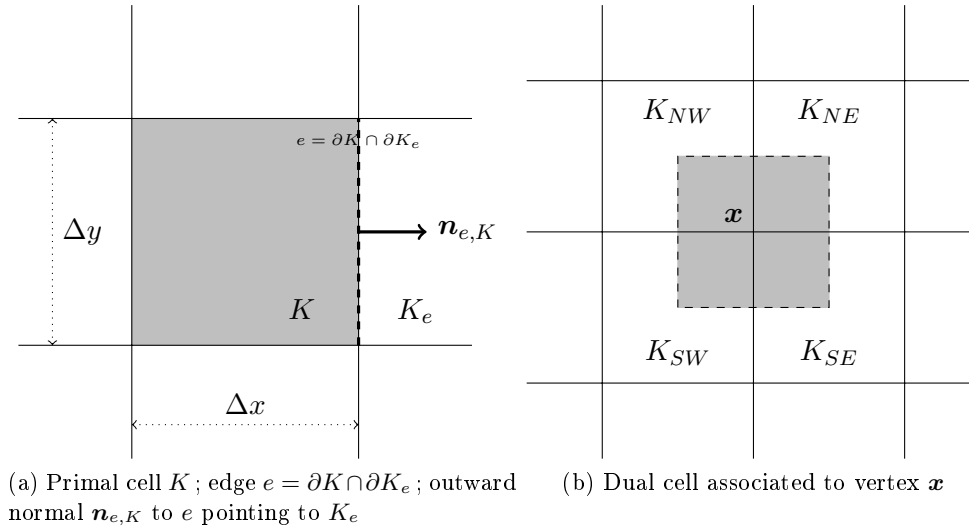


FIGURE A.1 – Geometric settings

Notations are pictured on Figs. A.1.

In particular, the discrete Green formula yields

$$\forall K \in \mathbb{T}, \sum_{e \subset \partial K} m_e \mathbf{n}_{e,K} = \mathbf{0}. \quad (\text{A.6})$$

In what follows, for any piecewise constant function  $\varphi$  we will denote

$$\bar{\varphi}_e = \frac{\varphi_{K_e} + \varphi_K}{2} \quad \text{and} \quad \bar{\varphi}_{\mathbf{x}} = \frac{1}{4} (\varphi_{K_{NW}} + \varphi_{K_{NE}} + \varphi_{K_{SE}} + \varphi_{K_{SW}})$$

and introduce the vector quantity

$$\delta \varphi_e = \frac{\varphi_{K_e} - \varphi_K}{2} \mathbf{n}_{e,K}.$$

In the same way, for a piecewise constant vector function  $\Phi$ , we will note

$$\delta \Phi_e = \frac{\Phi_{K_e} - \Phi_K}{2} \cdot \mathbf{n}_{e,K}.$$

Finally, the quantity

$$\varphi^\pm = \frac{1}{2} (\varphi \pm |\varphi|)$$

will refer to the the positive and negative parts of any scalar function  $\varphi$ . Before describing the schemes, we also need to define the following discrete divergence and gradient operators :

- Given a vector function  $(\Phi_e)_{e \in \mathcal{E}}$ , the discrete divergence operator  $\text{div}_K$  reads

$$\forall K \in \mathbb{T}, \text{div}_K \Phi = \frac{1}{m_K} \sum_{e \subset \partial K} m_e \Phi_e \cdot \mathbf{n}_{e,K}. \quad (\text{A.7})$$

If  $(\Phi_K)_{K \in \mathbb{T}}$  is rather considered, then  $\Phi_e$  is replaced by  $\bar{\Phi}_e$  in (A.7) and the corresponding operator is denoted by  $\text{div}_K^c \Phi$ ;

We can also define a divergence operator at the vertices when  $(\Phi_K)_{K \in \mathbb{T}}$  is considered by applying Formula (A.7) to the dual cell (see Fig.A.1(b)) :

$$\forall \mathbf{x} \in \mathcal{V}, \operatorname{div}_{\mathbf{x}}^c \Phi = \left( \frac{\Phi_{KSE} + \Phi_{KNE}}{2} - \frac{\Phi_{KSW} + \Phi_{KNW}}{2} \right) \cdot \frac{\mathbf{e}_1}{\Delta x} + \left( \frac{\Phi_{KNE} + \Phi_{KNW}}{2} - \frac{\Phi_{KSE} + \Phi_{KSW}}{2} \right) \cdot \frac{\mathbf{e}_2}{\Delta y}. \quad (\text{A.8})$$

- In the same way, considering  $(\Phi_e)_{e \in \mathcal{E}}$  and  $(\Psi_K)_{K \in \mathbb{T}}$ , the upwind divergence operator applied to the formal product  $\Psi \otimes \Phi$  is defined by

$$\forall K \in \mathbb{T}, \operatorname{div}_K^{up}(\Psi \otimes \Phi) = \frac{1}{m_K} \sum_{e \subset \partial K} m_e (\Psi_K (\Phi_e \cdot \mathbf{n}_{e,K})^+ + \Psi_{K_e} (\Phi_e \cdot \mathbf{n}_{e,K})^-). \quad (\text{A.9})$$

Likewise, if  $(\Phi_K)_{K \in \mathbb{T}}$  is rather considered, then  $\Phi_e$  is replaced by  $\bar{\Phi}_e$  in (A.9) and the corresponding operator is denoted by  $\operatorname{div}_K^{c,up}(\Psi \Phi)$ ;

- Given a scalar function  $(\phi_e)_{e \in \mathcal{E}}$ , the discrete gradient reads

$$\forall K \in \mathbb{T}, \nabla_K \phi = \frac{1}{m_K} \sum_{e \subset \partial K} m_e \phi_e \mathbf{n}_{e,K}. \quad (\text{A.10})$$

Likewise, if  $(\phi_K)_{K \in \mathbb{T}}$  is rather considered, then  $\phi_e$  is replaced by  $\bar{\phi}_e$  in (A.10) and the corresponding operator is denoted by  $\nabla_K^c \phi$ .

If  $(\phi_{\mathbf{x}})_{\mathbf{x} \in \mathcal{V}}$  is considered, then  $\phi_e$  is replaced by  $\frac{1}{2} \sum_{\mathbf{x} \in \partial e} \phi_{\mathbf{x}}$  and the corresponding operator is denoted by  $\nabla_K^x \phi$ .

In particular, we have the mimetic properties (see [5])

$$\operatorname{div}_K^c(\nabla_K^\perp \varphi) = 0 \quad \text{and} \quad \sum_{K \in \mathbb{T}} m_K \nabla_K^x \varphi \cdot \Phi_K = - \sum_{\mathbf{x} \in \mathcal{V}} m_{\mathbf{x}} \varphi_{\mathbf{x}} \operatorname{div}_{\mathbf{x}}^c \Phi. \quad (\text{A.11})$$

## Numerical scheme

In this collocated version, all unknowns are located at the center of the cells except the perturbations ( $\mathbf{q}$  located at the middle of the interfaces and  $\pi$  at the vertices). With the discrete operators (A.7-A.8-A.9-A.10) introduced above, the scheme reads, in the semi-discrete case :

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K(\mathcal{F}) = 0, \\ \frac{d}{dt} (h_K \mathbf{u}_K) + \operatorname{div}_K^{up}(\mathbf{u} \otimes \mathcal{F}) + h_K \nabla_K^c \phi - \nabla_K^x \pi = -\omega \left( h_K \mathbf{u}_K - \frac{1}{2} \sum_{e \subset \partial K} \mathbf{q}_e \right)^\perp, \end{cases} \quad (\text{A.12})$$

where the interface fluxes are defined at the level of the edge  $e$  as :

$$\mathcal{F}_e = \bar{h} \mathbf{u}_e - \mathbf{q}_e, \quad (\text{A.13})$$

with

$$\mathbf{q}_e = \gamma \Delta t \bar{h}_e \left( 2 \frac{m_e}{m_K} \delta \phi_e + \omega (\bar{\mathbf{u}}_e^\perp \cdot \mathbf{n}_{e,K}) \mathbf{n}_{e,K} \right), \quad (\text{A.14})$$

and

$$\pi_{\mathbf{x}} = \nu \Lambda \max\{\Delta x, \Delta y\} \bar{h}_{\mathbf{x}} \operatorname{div}_{\mathbf{x}}^c \mathbf{u}, \quad (\text{A.15})$$

where  $\gamma$  and  $\nu$  are positive dimensionless constants, and  $\Lambda$  is a positive characteristic velocity.<sup>3</sup>

<sup>3</sup>Typically, we take  $\Lambda = \max_K \left\{ \|\mathbf{u}_K^0\| + \sqrt{g h_K^0} \right\}$ .

**Remark A.1.1**

Some remarks must be done on the formulation of the scheme :

1. The origin of the expression (A.14) for the perturbation  $q$  will be explained in the proof below. We first mention that it is consistent with the geostrophic equilibrium (it is aimed to vanish at the equilibrium). Second, we only consider the normal component of  $\bar{\mathbf{u}}_e$  in order to have the correct number of constraints. Moreover, as the operator  $\delta\phi_e$  is along the normal vector,  $\mathbf{q}_e = 0$  would imply (without the normal vector) a too restrictive kernel (characterised by  $u_{K,x} = 0$  or  $u_{K,y} = 0$  depending on the edges).
2. As  $m_K = \Delta x \Delta y$  in this case,  $\mathbf{q}_e$  does not actually depend on  $K$ .
3. Let us now explain why there is a  $\frac{1}{2}$  coefficient the correction term in the right hand side of (A.12). The previous remark combined to the cartesian structure induces that there are only two terms by component in the mean  $\sum_e \mathbf{q}_e$ .<sup>4</sup>

**Semi-discrete energy**

We now show this scheme ensures a discrete counterpart of (A.5) through semi-discrete mechanic energy estimates. To obtain such a result, we first need the two following lemmas, describing the evolution of potential and kinetic energies.

**Lemma A.1.2 (Semi-discrete scheme for the potential energy)**

We set  $\mathcal{E}_K = \frac{1}{2}g(h_K)^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt}\mathcal{E}_K + \text{div}_K(\phi\mathcal{F}) - h_K\mathbf{u}_K \cdot \nabla_K^c\phi + \mathcal{A}_K^\mathcal{E} = -\frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \cdot \mathbf{q}_e, \quad (\text{A.16})$$

where the interface value for  $\phi\mathcal{F}$  is defined as :

$$(\phi\mathcal{F})_e = \bar{\phi}_e \mathcal{F}_e,$$

and the discrete conservative term :

$$\mathcal{A}_K^\mathcal{E} = \text{div}_K(\delta\phi \delta(h\mathbf{u})) = \frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \delta(h\mathbf{u})_e \cdot \mathbf{n}_{e,K}. \quad (\text{A.17})$$

**Proof:** Using the definition of  $\mathcal{E}_K$  and the mass conservation law (A.1), we get :

$$\frac{d}{dt}\mathcal{E}_K = \phi_K \frac{d}{dt}h_K = -\phi_K \text{div}_K(\mathcal{F}).$$

Then, invoking the decomposition  $\phi_K = \bar{\phi}_e + \frac{1}{2}(\phi_K - \phi_{K_e})$  and the expression of the numerical fluxes (A.13) :

$$\frac{d}{dt}\mathcal{E}_K = -\frac{1}{m_K} \sum_{e \subset \partial K} m_e \bar{\phi}_e \mathcal{F}_e \cdot \mathbf{n}_{e,K} + \frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \cdot \overline{h\mathbf{u}}_e - \frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \cdot \mathbf{q}_e.$$

Lastly, using the relation  $\overline{h\mathbf{u}}_e = h_K \mathbf{u}_K + \frac{1}{2}(h_{K_e} \mathbf{u}_{K_e} - h_K \mathbf{u}_K)$  and the discrete Green formula (A.6), basic computations yield :

$$\frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta \phi_e \cdot \overline{h\mathbf{u}}_e = h_K \mathbf{u}_K \cdot \nabla_K^c \phi + \frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta \phi_e \delta(h\mathbf{u})_e \cdot \mathbf{n}_{e,K},$$

which allows to conclude. ■

**Lemma A.1.3 (Semi-discrete scheme for the kinetic energy)**

We set  $\mathcal{K}_K = \frac{1}{2} h_K \|\mathbf{u}_K\|^2$  for  $K \in \mathbb{T}$ . Then :

$$\frac{d}{dt} \mathcal{K}_K + \frac{1}{2} \operatorname{div}_K^{up} (\|\mathbf{u}\|^2 \mathcal{F}) + h_K \mathbf{u}_K \cdot \nabla_K^c \phi - \mathcal{A}_K^K \leq -\frac{\omega}{2} \sum_{e \subset \partial K} \overline{\mathbf{u}}_e^\perp \cdot \mathbf{q}_e + \mathbf{u}_K \cdot \nabla_K^x \pi, \quad (\text{A.18})$$

with anti-symmetric contributions given by :

$$\mathcal{A}_K^K = \frac{\omega}{4} \sum_{e \subset \partial K} (\mathbf{u}_K - \mathbf{u}_{K_e}) \cdot \mathbf{q}_e^\perp. \quad (\text{A.19})$$

**Proof:** We write :

$$\begin{aligned} \frac{d}{dt} \left( \frac{1}{2} h_K \|\mathbf{u}_K\|^2 \right) &= \mathbf{u}_K \cdot \frac{d}{dt} (h_K \mathbf{u}_K) - \frac{1}{2} \|\mathbf{u}_K\|^2 \frac{d}{dt} h_K \\ &= -h_K \mathbf{u}_K \cdot \nabla_K^c \phi - \mathbf{u}_K \cdot \operatorname{div}_K^{up} (\mathbf{u} \otimes \mathcal{F}) + \omega \mathbf{u}_K \cdot \left( \frac{1}{2} \sum_{e \subset \partial K} \mathbf{q}_e^\perp \right) + \frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K (\mathcal{F}) + \mathbf{u}_K \cdot \nabla_K^x \pi. \end{aligned} \quad (\text{A.20})$$

First, we remark that :

$$\mathbf{u}_K \cdot \left( \frac{1}{2} \sum_{e \subset \partial K} \mathbf{q}_e^\perp \right) = \frac{1}{2} \sum_{e \subset \partial K} \overline{\mathbf{u}}_e \cdot \mathbf{q}_e^\perp + \frac{1}{4} \sum_{e \subset \partial K} (\mathbf{u}_K - \mathbf{u}_{K_e}) \cdot \mathbf{q}_e^\perp = -\frac{1}{2} \sum_{e \subset \partial K} \overline{\mathbf{u}}_e^\perp \cdot \mathbf{q}_e + \mathcal{A}_K^K. \quad (\text{A.21})$$

Then, after some basic computations, we get the relation

$$\frac{1}{2} \|\mathbf{u}_K\|^2 \operatorname{div}_K (\mathcal{F}) - \mathbf{u}_K \cdot \operatorname{div}_K^{up} (\mathbf{u} \otimes \mathcal{F}) = -\frac{1}{2} \operatorname{div}_K^{up} (\|\mathbf{u}\|^2 \mathcal{F}) + \frac{1}{2m_K} \sum_{e \subset \partial K} \|\mathbf{u}_K - \mathbf{u}_{K_e}\|^2 (\mathcal{F}_e \cdot \mathbf{n}_{e,k})^-. \quad (\text{A.22})$$

The second term of the right hand side being non-positive, we get the announced result inserting (A.21) and (A.22) into (A.20). ■

Gathering relations (A.16) and (A.18), we obtain the following estimate for the total energy  $E_K = \mathcal{E}_K + \mathcal{K}_K$  :

$$\frac{d}{dt} E_K + \operatorname{div}_K (\phi \mathcal{F}) + \frac{1}{2} \operatorname{div}_K^{up} (\|\mathbf{u}\|^2 \mathcal{F}) + \mathcal{A}_K^\mathcal{E} - \mathcal{A}_K^K \leq -\frac{1}{2} \sum_{e \subset \partial K} \left[ 2 \frac{m_e}{m_K} \delta \phi_e + \omega \overline{\mathbf{u}}_e^\perp \right] \cdot \mathbf{q}_e + \mathbf{u}_K \cdot \nabla_K^x \pi.$$

With the choices (A.14) and (A.15), we finally obtain a discrete counterpart of (A.5) thanks to (A.11).

**Remark A.1.4**

The contributions  $\mathcal{A}_K^\mathcal{E}$  (A.17) and  $\mathcal{A}_K^\mathcal{K}$  (A.19) involve anti-symmetric terms and therefore have no impact on the total energy budget :

$$\sum_{K \in \mathbb{T}} m_K (\mathcal{A}_K^\mathcal{E} - \mathcal{A}_K^\mathcal{K}) = 0,$$

up to boundary terms. In fact, these quantities can be interpreted as a bias to the conservative parts of potential and kinetic energies respectively.

Then :

**Theorem A.1.5 (Semi-discrete mechanic balance)**

The scheme (A.12, A.13, A.14, A.15) satisfies the following semi-discrete inequality :

$$\begin{aligned} & \frac{d}{dt} \left( \sum_{K \in \mathbb{T}} m_K E_K \right) \\ & \leq -\frac{\gamma \Delta t}{2} \sum_{K \in \mathbb{T}} m_K \sum_{e \subset \partial K} \bar{h}_e \left\| \frac{\phi_{K_e} - \phi_K}{m_K/m_e} + \omega \bar{\mathbf{u}}_e^\perp \cdot \mathbf{n}_{e,K} \right\|^2 - \nu \Lambda \max\{\Delta x, \Delta y\} \sum_{\mathbf{x} \in \mathcal{V}} m_{\mathbf{x}} \bar{h}_{\mathbf{x}} (\operatorname{div}_{\mathbf{x}}^c \mathbf{u})^2. \end{aligned}$$

**Steady states**

Let us first assess the ability of the scheme to capture the lake at rest, characterized by  $\phi = cte$  and  $\mathbf{u} = 0$ . One observes that the stabilization terms  $\mathbf{q}_e$  and  $\pi_{\mathbf{x}}$  vanish in such a configuration, leading directly to the exact preservation of the steady state in (A.12).

Then focus on the geostrophic equilibrium. In this case the nonlinear study is particularly complex and we choose to investigate the linear case. Although it leaves aside nonlinear interactions, this strategy may not seem irrelevant since most of low Froude number regimes encountered in practical contexts are governed by linear dynamics. Linearizing the scheme (A.12) around  $\tilde{\mathbf{u}} = 0$  and  $\tilde{h} = h_0$  for some constant  $h_0 > 0$ , we get :

$$\begin{cases} \frac{d}{dt} h_K + h_0 \operatorname{div}_K^c \mathbf{u} - \operatorname{div}_K \mathbf{q} = 0, \\ \frac{d}{dt} \mathbf{u}_K + \nabla_K^c \phi - \nabla_K^x \tilde{\pi} = -\omega \left( \mathbf{u}_K - \frac{1}{2} \sum_{e \subset \partial K} \mathbf{q}_e \right)^\perp, \end{cases}$$

where

$$\mathbf{q}_e = \gamma \Delta t \left( 2 \frac{m_e}{m_K} \delta \phi_e + \omega (\bar{\mathbf{u}}_e^\perp \cdot \mathbf{n}_{e,K}) \mathbf{n}_{e,K} \right) \text{ and } \tilde{\pi}_{\mathbf{x}} = -\nu \Lambda \max\{\Delta x, \Delta y\} \operatorname{div}_{\mathbf{x}}^c \mathbf{u}.$$

One should observe here that the discrete geostrophic equilibrium, characterized by :

$$\nabla_K^c \phi + \omega \mathbf{u}_K^\perp = 0 \stackrel{(A.11)}{\implies} \operatorname{div}_K^c \mathbf{u}_K = 0$$

does not imply  $\mathbf{q}_e = 0$  in the general case. Conversely,  $\mathbf{q}_e = 0 \implies \pi_{\mathbf{x}} = 0$  according to [5] but does not imply  $\nabla_K^c \phi + \omega \mathbf{u}_K^\perp = 0$ . As exhibited in the numerical validations, although it does not exactly preserve the geostrophic balance, the present scheme behaves much better than standard Godunov-type approaches. However, as detailed in the next section, it is possible to correct such a failure by means of a staggered mesh.

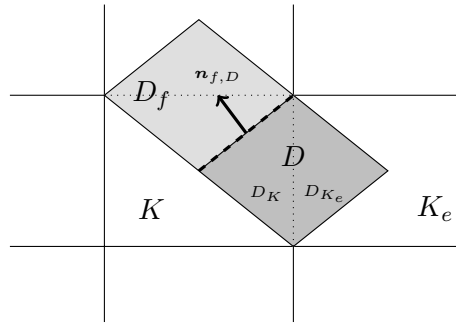


FIGURE A.2 – Geometric settings : for the primal mesh with cells  $K$  and  $K_e$ , refer to Figs. A.1 – for the dual mesh with diamond cells  $D$  and  $D_f$  separated by the interface  $f = \partial D \cap \partial D_f$ ,  $\mathbf{n}_{f,D}$  is the outward normal of  $f$  pointing to  $D_f$ .

### A.1.2 Staggered scheme

#### Mesh and notations

On the basis of the primal mesh  $\mathbb{T}$ , a dual grid  $\mathbb{T}^*$  can be built with respect to the Rannacher-Turek (RT) or Crouzeix-Raviart finite-element spaces, according to the formalism described in [3]. To set ideas, we focus here on the RT case considering a primal mesh  $\mathbb{T}$  made of regular quadrilateral elements like in the previous section, but the following framework also applies to simplicial meshes or mixed triangular/quadrilateral elements.

- A dual cell  $D \in \mathbb{T}^*$  is associated to each internal primal edge  $e$  as the quadrilateral admitting  $e$  as diagonal and the mass centers of the two neighbouring primal cells  $K$  and  $K_e$  as vertices (see Fig. A.2).<sup>5</sup>
- This dual element  $D$  is the union of the two triangles  $D_K, D_{K_e}$  separated by the edge  $e$ , whose areas will be denoted  $m_{D_K}$  and  $m_{D_{K_e}}$ , following the notations of the primal mesh.
- Given a dual cell  $D$  and an edge  $e \subset \partial K$ ,  $D_f$  is the neighbouring cell to  $f$  (other than  $D$ ) and  $\mathbf{n}_{f,D}$  the outward normal pointing to  $D_f$ .

In this context, the water height evolves on the primal mesh  $\mathbb{T}$  as for the previous scheme while the velocity field is located on the dual elements  $D \in \mathbb{T}^*$ . By analogy with the collocated frame, the upwind discrete divergence is defined on a dual element as :

$$\mathbf{div}_D^{up}(\Psi \otimes \Phi) = \frac{1}{m_D} \sum_{f \in \partial D} m_f (\Psi_D(\Phi_f \cdot \mathbf{n}_{f,D})^+ + \Psi_{D_f}(\Phi_f \cdot \mathbf{n}_{f,D})^-),$$

while the dual gradient is defined through the following formula :

$$\nabla_D \varphi = \frac{m_e}{m_D} (\varphi_{K_e} - \varphi_K) \mathbf{n}_{e,K}. \quad (\text{A.23})$$

Note that this definition ensures the grad/div duality with respect to the  $L^2$  inner product in the present staggered environment (for instance, see [37]) :

$$\sum_{K \in \mathbb{T}} m_K \varphi_K \mathbf{div}_K \Psi = - \sum_{D \in \mathbb{T}^*} m_D \nabla_D \varphi \cdot \Psi_D. \quad (\text{A.24})$$

<sup>5</sup>The edge  $e$  and the corresponding diamond cell  $D$  will be used indifferently.



### Numerical scheme

The scheme we consider is the following :

$$\begin{cases} \frac{d}{dt} h_K + \operatorname{div}_K(\mathcal{F}) = 0, \\ \frac{d}{dt} (h_D \mathbf{u}_D) + \operatorname{div}_D^{up}(\mathbf{u} \otimes \mathcal{F}) + h_D \nabla_D \phi - \nabla_D \pi = -\omega (h_D \mathbf{u}_D - \mathbf{q}_D)^\perp, \end{cases} \quad (\text{A.25})$$

where the primal fluxes are defined by interface by :

$$\mathcal{F}_e = h_D \mathbf{u}_D - \mathbf{q}_D, \quad (\text{A.26})$$

with

$$\mathbf{q}_D = \gamma \Delta t h_D \left( \nabla_D \phi + \omega \mathbf{u}_D^\perp \right), \quad (\text{A.27})$$

and

$$\pi_K = \nu \Lambda \max\{\Delta x, \Delta y\} h_K \operatorname{div}_K \mathbf{u}. \quad (\text{A.28})$$

As in the collocated frame,  $\gamma$  and  $\nu$  are positive constants and  $\Lambda$  is consistent with a velocity. Considering the momentum equations (A.25) and the diffusion term (A.27), it is necessary to define a dual water height  $h_D$ . Following [3], we compute

$$m_D h_D = m_{D_K} h_K + m_{D_{K_e}} h_{K_e}.$$

It is then possible to choose dual fluxes  $\mathcal{F}_f$  (involved in the divergence term of the momentum equations (A.25)) so that a discrete law of mass conservation also holds on diamond cells :

$$\frac{d}{dt} h_D + \frac{1}{m_D} \sum_{f \in \partial D} m_f \mathcal{F}_f \cdot \mathbf{n}_{f,D} = 0,$$

which can be rewritten in the more compact form :

$$\frac{d}{dt} h_D + \operatorname{div}_D(\mathcal{F}) = 0. \quad (\text{A.29})$$

As detailed in [3], these dual fluxes are expressed as linear combinations of the primal ones. These relations are purely geometrical and formulated locally, depending on the nature of the staggered mesh. Such a relation is mandatory to derive a kinetic energy budget as shown in the next paragraph.

#### Remark A.1.6

*Let us notice that in (A.27), both components are taken into account in the definition of  $\mathbf{q}_D$  unlike (A.14) where only the normal component was involved. In the present case,  $\mathbf{q}_D = 0$  implies that the normal component of the velocity field vanishes (since the gradient is normal to the interface) but the value of interest is the reconstructed mean velocity*

$$\mathbf{u}_K = \frac{1}{4} \sum_{e \subset \partial K} \mathbf{u}_D.$$

*Here  $\mathbf{u}_D$  is an unknown in itself (unlike in the collocated case where  $\bar{u}_e$  is reconstructed from the unknowns  $\mathbf{u}_K$ ).*

### Semi-discrete energy

Indeed, defining the dual kinetic energy by  $\mathcal{K}_D = \frac{1}{2}h_D\|\mathbf{u}_D\|^2$  and using the relation

$$\frac{d}{dt}\mathcal{K}_D = \mathbf{u}_D \cdot \frac{d}{dt}(h_D\mathbf{u}_D) - \frac{1}{2}\|\mathbf{u}\|^2 \frac{d}{dt}h_D$$

we have, with (A.29) and applying the relation (A.22) on the dual mesh :

$$\frac{d}{dt}\mathcal{K}_D + \operatorname{div}_D^{up}(\|\mathbf{u}_D\|^2\mathcal{F}) + h_D\mathbf{u}_D \cdot \nabla_D\phi \leq -\omega\mathbf{u}_D^\perp \cdot \mathbf{q}_D + \mathbf{u}_D \cdot \nabla_D\pi. \quad (\text{A.30})$$

Then, looking at the potential energy, basic computations give :

$$\frac{d}{dt}\mathcal{E}_K + \operatorname{div}_K(\phi\mathcal{F}) - \frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \cdot h_D\mathbf{u}_D = -\frac{1}{m_K} \sum_{e \subset \partial K} m_e \delta\phi_e \cdot \mathbf{q}_D. \quad (\text{A.31})$$

We are now left with the estimates (A.30) and (A.31), which are not expressed on the same mesh. A convenient way to understand the global impact of these quantities on the energy budget is to integrate them over their respective meshes, setting :

$$E = \sum_{K \in \mathbb{T}} m_K \mathcal{E}_K + \sum_{D \in \mathbb{T}^*} m_D \mathcal{K}_D,$$

and focus on the residual contributions at the level of dual elements. We first remark that, thanks to the duality relation (A.24), the second term of the right hand side of (A.30) furnish a first stabilizing term :

$$\sum_{D \in \mathbb{T}^*} m_D \mathbf{u}_D \cdot \nabla_D \pi = - \sum_{K \in \mathbb{T}} m_K \pi_K \operatorname{div}_K \mathbf{u},$$

according to the definition of  $\pi_K$  (A.28). Then, considering a generic element  $D \in \mathbb{T}^*$ , the non conservative terms in the kinetic energy balance (A.30) result in the following quantities on the left and right hand sides :

$$m_D h_D \mathbf{u}_D \cdot \nabla_D \phi \quad (LHS) \quad , \quad -m_D \omega \mathbf{u}_D^\perp \cdot \mathbf{q}_D. \quad (RHS)$$

On the other hand, the term issuing from (A.31) needs to be counted twice (contributions of the elements  $K$  and  $K_e$  sharing the edge  $e$ ), which yields :

$$-2m_e \delta\phi_e \cdot h_D \mathbf{u}_D \quad (LHS) \quad , \quad -2m_e \delta\phi_e \cdot \mathbf{q}_D. \quad (RHS)$$

Using the definition (A.23), we easily verify that the terms of the left hand side are in exact balance and the remaining term reads :

$$-m_D \mathbf{q}_D \cdot \left( \nabla_D \phi + \omega \mathbf{u}_D^\perp \right).$$

Replacing  $\mathbf{q}_D$  by its expression (A.27), we have thus established the following result :

#### **Theorem A.1.7 (Semi-discrete mechanic balance)**

*The scheme (A.25-A.26-A.27) satisfies the following semi-discrete inequality :*

$$\frac{d}{dt}E \leq -\gamma \Delta t \sum_{D \in \mathbb{T}^*} m_D h_D \left\| \nabla_D \phi + \omega \mathbf{u}_D^\perp \right\|^2 - \nu \Lambda \max\{\Delta x, \Delta y\} \sum_{K \in \mathbb{T}} m_K h_K (\operatorname{div}_K \mathbf{u})^2.$$

As a remark, defining the quantity  $\mathcal{K}_K = \frac{1}{m_K} \sum_{e \subset \partial K} m_D \mathcal{K}_D$ , which can be understood as a kinetic energy associated to the cell  $K$ , the previous result can be reformulated under local energy estimates. Note that the present formalism embeds the implicit and explicit approaches recently proposed in [25] without Coriolis force.

### Steady states

The linearized version of the scheme (A.25) is the following :

$$\begin{cases} \frac{d}{dt} h_K + h_0 \operatorname{div}_K \mathbf{u} - \operatorname{div}_K \mathbf{q} = 0, \\ \frac{d}{dt} \mathbf{u}_D + \nabla_D \phi - \nabla_D \tilde{\pi} = -\omega (\mathbf{u}_D - \mathbf{q}_D)^\perp, \end{cases}$$

where

$$\mathbf{q}_D = \gamma \Delta t \left( \nabla_D \phi + \omega \mathbf{u}_D^\perp \right) \quad \text{and} \quad \tilde{\pi}_K = \nu \Lambda \max\{\Delta x, \Delta y\} \operatorname{div}_K \mathbf{u}.$$

From this, we immediately see that discrete kernel of the system corresponds to the geostrophic balance expressed on the dual grid  $\mathbb{T}^*$  :

$$\nabla_D \phi + \omega \mathbf{u}_D^\perp = 0,$$

and is consequently exactly preserved by the linearized scheme. We may indeed check that  $\operatorname{div}_K(\nabla_D^\perp \phi) = 0$  by construction which implies  $\operatorname{div}_K \mathbf{u} = 0$  when the geostrophic equilibrium holds.

#### A.1.3 Time discretisation

For the discretisation in time, fluxes are taken explicit. We only remind the reader of the standard discretisation of the right hand side (Coriolis force). The ODE

$$\frac{du}{dt} = \omega v, \quad \frac{dv}{dt} = -\omega u,$$

is solved by means of the  $\theta$  scheme

$$\frac{u^{n+1} - u^n}{\Delta t} = \omega(\theta_u v^n + (1 - \theta_u)v^{n+1}), \quad \frac{v^{n+1} - v^n}{\Delta t} = -\omega(\theta_v u^n + (1 - \theta_v)u^{n+1})$$

with  $\theta_u + \theta_v \leq 1$ . Here we choose  $\theta_u = 1$  and  $\theta_v = 0$  so that the system is solved explicitly.

The time step is chosen following [5] such that

$$\Delta t^n \leq \min \left\{ \frac{2}{\omega}, \frac{\min(\Delta x, \Delta y)}{\max_D \|\mathbf{u}_D^n\| + \sqrt{gh_D^n}} \right\}.$$

## A.2 Numerical assessments of the schemes

### A.2.1 Presentation of the test cases

To assess the numerical schemes designed in this work, we consider the stationary vortex described in [7]. Set in the domain  $\Omega = [-0.5, 0.5] \times [-0.5, 0.5]$ , the equations are supplemented with the initial velocity field given in polar coordinates

$$\mathbf{u}^0(r, \theta) = \nu_\theta(r) \mathbf{e}_\theta \quad \text{where} \quad \nu_\theta(r) = \varepsilon \left[ 5r \mathbf{1}_{\{r < 1/5\}} + (2 - 5r) \mathbf{1}_{\{1/5 \leq r < 2/5\}} \right].$$

The initial water height is then computed so that the steady version of the nonlinear equations is satisfied, *i.e.* such that

$$\partial_r h^0(r, \theta) = \frac{\nu_\theta(r)}{g} \left( \omega + \frac{\nu_\theta(r)}{r} \right).$$

As mentioned in [5], this case induces a Froude number of order  $\mathcal{O}(\varepsilon)$ .

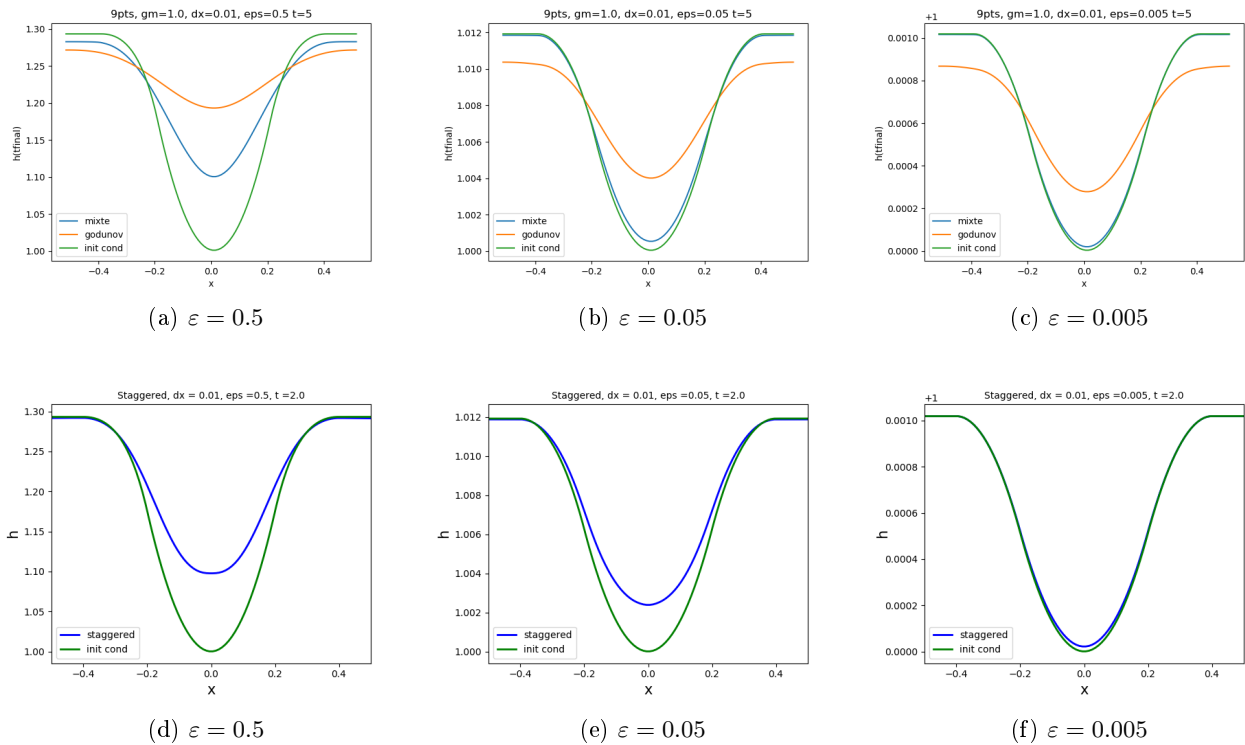


FIGURE A.3 – Stationary vortex : cross-section of the water height at final time for different values of  $\epsilon$ ; comparisons between the exact solution (*green curves*), the numerical solution obtained with the classic Godunov scheme (*orange curve*), the collocated scheme (A.12) (A-B-C) and the staggered scheme (A.25) (D-E-F) (*blue curves*)

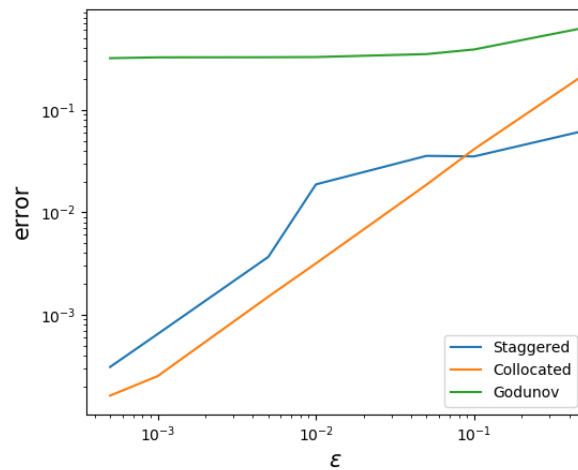


FIGURE A.4 – Error in  $L^2$  norm between the numerical solution and the exact solution at final time

### A.2.2 Numerical results

Given the initial condition prescribed in Section A.2.1, we compare the robustness of several numerical schemes. The data are a steady solution of the nonlinear equations but their approximation on the grid is not necessarily preserved by any numerical scheme. We remind that the present paper was aimed at designing a scheme that is robust in the low Froude number regime by means of diffusion terms based on linear equilibria.

We observe on Figure A.3 that the smaller  $\varepsilon$ , the more accurate the solution obtained with the collocated scheme (A.12) while the accuracy of the standard Godunov scheme is barely improved.

This is confirmed by Fig. A.4 where a kind of relative error is plotted :

$$E(\varepsilon) = \frac{\sum_i (h_K^N - h_K^0)^2}{\sum_i (h_K^0 - \bar{h})^2}$$

where  $N$  is the index of the final time and  $\bar{h} = \max_K h_K^0$ . The Godunov scheme broadly provides results of which the accuracy is independent from  $\varepsilon$  while the modifications proposed in this paper leads to a second order scheme.

These results tend to show that the target is reached insofar as the resulting scheme is robust in the low Froude number regime where the standard Godunov scheme was inaccurate. Considering the presented simulations, it does not seem that the two constants  $\gamma$  and  $\nu$  must necessarily be both strictly positive to get stability. A linear stability analysis and other numerical investigations are currently in progress to clarify this point.

## A.3 Conclusion

We managed to construct some corrections to the standard Godunov scheme so that the resulting collocated and staggered schemes are accurate and energy-decreasing around the geostrophic equilibrium in the low Froude number regime. The stationary vortex test case go to show the improvement due to these new terms. The proofs are given in the semi-discrete case. In particular, the semi-discrete energy inequality requires  $\gamma \geq 0$  and  $\nu \geq 0$  so that it decreases. Future works will thus deal with the stability analysis of the fully discrete version. Moreover, further numerical investigations are needed to ensure the schemes behave well in a large range of applications.

*Acknowledgements.* The authors express their deep gratitude to the French National Research Agency project NABUCO, grant ANR-17-CE40-0025 and the french INSU-CNRS (Institut National des Sciences de l'Univers – Centre National de la Recherche Scientifique) program LEFE-MANU (Méthodes Mathématiques et Numériques), project DWAVE for the funding of the project, as well as to the organizing committee of the CEMRACS'19 event. The authors also thank Martin Parisot for fruitful discussions.

# Bibliographie

- [1] N. Aïssiouene. *Numerical analysis and discrete approximation of a dispersive shallow water model*. Theses, Université Pierre et Marie Curie - Paris VI, Dec. 2016.
- [2] N. Aïssiouene, M.-O. Bristeau, E. Godlewski, A. Mangeney, C. Parés, and J. Sainte-Marie. A two-dimensional method for a family of dispersive shallow water models. *SMAI Journal of Computational Mathematics*, 6 :39, Sept. 2020.
- [3] G. Ansanay-Alex, F. Babik, J.-C. Latché, and D. Vola. An  $L^2$ -stable approximation of the Navier–Stokes convection operator for low-order non-conforming finite elements. *International Journal for Numerical Methods in Fluids*, 66(5) :555–580, 2011.
- [4] E. Audusse, F. Bouchut, M. Bristeau, R. Klein, and B. Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM Journal on Scientific Computing*, 25(6) :2050–2065, 2004.
- [5] E. Audusse, M. Do, P. Omnes, and Y. Penel. Analysis of modified godunov type schemes for the two-dimensional linear wave equation with coriolis source term on cartesian meshes. *Journal of Computational Physics*, 373 :91–129, 2018.
- [6] E. Audusse, V. Dubos, A. Duran, N. Gaveau, Y. Nasser, and Y. Penel. Numerical approximation of the shallow water equations with Coriolis source term. *ESAIM : Proceedings*, 70 :31–44, June 2021.
- [7] E. Audusse, R. Klein, and A. Owinoh. Conservative discretization of coriolis force in a finite volume framework. *Journal of Computational Physics*, 228(8) :2934–2950, 2009.
- [8] J. Behrens and F. Dias. New computational methods in tsunami science. *Philosophical Transactions of the Royal Society A : Mathematical, Physical and Engineering Sciences*, 373(2053) :20140382, 2015.
- [9] A. Bermudez and M. Vazquez-Cendon. Upwind Methods for Hyperbolic Conservation Laws with Source Terms. *Computers and Fluids*, 23(8) :1049–1071, 1994.
- [10] C. Berthon, A. Duran, A., F. Foucher, K. Saleh, and J. De Dieu Zabsonré. Improvement of the hydrostatic reconstruction scheme to get fully discrete entropy inequalities. *Journal of Scientific Computing*, May 2019.
- [11] C. Berthon, M. M’Baye, M. H. Le, and D. Seck. A well-defined moving steady states capturing Godunov-type scheme for Shallow-water model. *International Journal on Finite Volumes*, 15, 2020.

- [12] P. Bonneton, E. Barthelemy, F. Chazel, R. Cienfuegos, D. Lannes, F. Marche, and M. Tissier. Recent advances in serre–green naghdi modelling for wave transformation, breaking and runup processes. *European Journal of Mechanics-B/Fluids*, 30(6) :589–597, 2011.
- [13] F. Bouchut. *Nonlinear stability of finite volume methods for hyperbolic conservation laws, and well-balanced schemes for sources*, volume 2/2004. Birkhäuser Basel, 2004.
- [14] F. Bouchut, J. Le Sommer, and V. Zeitlin. Frontal geostrophic adjustment and nonlinear wave phenomena in one-dimensional rotating shallow water. Part 2. High-resolution numerical simulations. *Journal of Fluid Mechanics*, 514 :35–63, 2004.
- [15] M.-O. Bristeau, A. Mangeney, J. Sainte-Marie, and N. Seguin. An energy-consistent depth-averaged euler system : derivation and properties. *Discrete and Continuous Dynamical Systems - Series B*, 20(4) :28, 2015.
- [16] M. Castro, J. Antonio López, and C. Pares. Finite volume simulation of the geostrophic adjustment in a rotating shallow-water system. *SIAM Journal on Scientific Computing*, 31 :444–477, 01 2008.
- [17] F. Chazel, D. Lannes, and F. Marche. Numerical simulation of strongly nonlinear and dispersive waves using a green–naghdi model. *Journal of Scientific Computing*, 48(1) :105–116, 2011.
- [18] A. Chertock, M. Dudzinski, A. Kurganov, and M. Lukáčová-Medvid’ová. Well-balanced schemes for the shallow water equations with coriolis forces. *Numerische Mathematik*, 138 :939–973, 2018.
- [19] A. Chorin. Numerical solution of the navier–stokes equations. *Mathematics of Computation*, 22, 10 1968.
- [20] F. Couderc, A. Duran, and J.-P. Vila. An explicit asymptotic preserving low froude scheme for the multilayer shallow water model with density stratification. *Journal of Computational Physics*, 343 :235–270, 2017.
- [21] M. W. Dingemans. *Water Wave Propagation Over Uneven Bottoms*. World Scientific Publishing Company, 1997.
- [22] J. Droniou, R. Eymard, T. Gallouët, C. Guichard, and R. Herbin. *The gradient discretisation method*, volume 82 of *Mathématiques et Applications*. Springer International Publishing AG, Aug. 2018.
- [23] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Gradient schemes : a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. *Mathematical Models and Methods in Applied Sciences*, pages 23(13), 2395–2432, 2013.
- [24] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. A unified analysis of elliptic problems with various boundary conditions and their approximation. *Czechoslovak Mathematical Journal*, 70 :339–368, 2020.
- [25] A. Duran, J.-P. Vila, and R. Baraille. Energy-stable staggered schemes for the shallow water equations. <https://hal.archives-ouvertes.fr/hal-01988382/>, Jan. 2019.
- [26] J. R. Edwards. Reflections on the early development of the “ausm family” of riemann solvers. *Shock Waves*, 29(5) :601–609, 2019.
- [27] A. Ern and J.-L. Guermond. *Theory and Practice of Finite Elements*. Springer-Verlag New York, 2004.

- [28] L. Euler. Principes généraux du mouvement des fluides. *Mémoires de l'académie des sciences de Berlin, Volume 11, pp. 274-315.*, 1757.
- [29] R. Eymard, T. Gallouët, and R. Herbin. *Finite Volume Methods*, volume 7 of *Handbook of Numerical Analysis*. Elsevier, 2000.
- [30] R. Eymard, T. Gallouët, and R. Herbin. RTk mixed finite elements for some nonlinear problems. *Mathematics and Computers in Simulation*, pages 1–20, Dec. 2014.
- [31] R. Eymard, C. Guichard, and R. Herbin. Small-stencil 3D schemes for diffusive flows in porous media. *ESAIM : Mathematical Modelling and Numerical Analysis*, 46 :265–290, 2012.
- [32] V. Girault and P.-A. Raviart. *Finite Element Methods for Navier-Stokes Equations*. 1986.
- [33] S. Glimsdal, G. K. Pedersen, C. B. Harbitz, and F. Lovholt. Dispersion of tsunamis : does it really matter ? *Natural Hazards and Earth System Sciences*, 13(6) :1507–1526, 2013.
- [34] L. Gosse. *Computing qualitatively correct approximations of balance laws*, volume 2. Springer, 2013.
- [35] A. E. Green and P. M. Naghdi. A derivation of equations for wave propagation in water of variable depth. *Journal of Fluid Mechanics*, 78(2) :237–246, 1976.
- [36] J. Greenberg and A. LeRoux. A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM Journal on Numerical Analysis*, 33(1) :1–16, 1996.
- [37] R. Herbin, W. Kheriji, and J.-C. Latché. On some implicit and semi-implicit staggered schemes for the shallow water equations. *ESAIM :M2AN*, 48 :1807–1857, 2014.
- [38] S. Jin. Efficient asymptotic-preserving (AP) schemes for some multiscale kinetic equations. *SIAM Journal on Scientific Computing*, 21(2) :441–454, 1999.
- [39] R. J. LeVeque. *Finite volume methods for hyperbolic problems*. Cambridge Texts in Applied Mathematics. Cambridge University Press, Cambridge, 2002.
- [40] M.-S. Liou. A sequel to aum : Ausm+. *Journal of Computational Physics*, 129(2) :364–382, 1996.
- [41] M.-S. Liou and C. J. Steffen. A new flux splitting scheme. *Journal of Computational Physics*, 107(1) :23–39, 1993.
- [42] X. Liu, A. Chertock, and A. Kurganov. An asymptotic preserving scheme for the two-dimensional shallow water equations with Coriolis forces. *Journal of Computational Physics*, 391 :259–279, 2019.
- [43] M. Parisot and J.-P. Vila. Centered-potential regularization for the advection upstream splitting method. *SIAM Journal on Numerical Analysis*, 54(5) :3083–3104, 2015.
- [44] R. Rannacher. On chorin’s projection method for the incompressible navier-stokes equations. In J. G. Heywood, K. Masuda, R. Rautmann, and V. A. Solonnikov, editors, *The Navier-Stokes Equations II — Theory and Numerical Methods*, pages 167–183, Berlin, Heidelberg, 1992. Springer Berlin Heidelberg.
- [45] T. Ringler, J. Thuburn, J. Klemp, and W. Skamarock. A unified approach to energy conservation and potential vorticity dynamics for arbitrarily-structured c-grids. *Journal of Computational Physics*, 229(9) :3065–3090, 2010.



- 
- [46] A. Saint-Venant. Théorie du mouvement non permanent des eaux, avec application aux crues des rivières et a l'introduction de marées dans leurs lits. *Comptes Rendus des Séances de Académie des Sciences*, 73, 147, 237., 1871.
- [47] R. Temam. Une méthode d'approximation de la solution des équations de navier-stokes. *Bulletin de la Société Mathématique de France*, 96 :115–152, 1968.
- [48] J. Thuburn, T. Ringler, W. Skamarock, and J. Klemp. Numerical representation of geostrophic modes on arbitrarily structured c-grids. *Journal of Computational Physics*, 228(22) :8321–8335, 2009.
- [49] E. F. Toro. *Riemann Solvers and Numerical Methods for Fluid Dynamics*. Springer, 2009.
- [50] R. A. Walters, E. Hanert, J. Pietrzak, and D. Le Roux. Comparison of unstructured, staggered grid methods for the shallow water equations. *Ocean Modelling*, 28(1) :106–117, 2009. The Sixth International Workshop on Unstructured Mesh Numerical Modelling of Coastal, Shelf and Ocean Flows.
- [51] R. A. Walters, E. M. Lane, and E. Hanert. Useful time-stepping methods for the coriolis term in a shallow water model. *Ocean Modelling*, 28(1) :66–74, 2009. The Sixth International Workshop on Unstructured Mesh Numerical Modelling of Coastal, Shelf and Ocean Flows.
- [52] H. Zakerzadeh. The rs-imex scheme for the rotating shallow water equations with the coriolis force. In C. Cancas and P. Omnes, editors, *Finite Volumes for Complex Applications VIII - Methods and Theoretical Aspects*, pages 199–207. Springer, 2017.
- [53] V. Zeitlin. *Geophysical fluid dynamics : understanding (almost) everything with rotating shallow water models*. Oxford University Press, 2018.



**Résumé :** Nous étudions dans cette thèse des méthodes numériques pour les écoulements en eaux peu profondes à surface libre. D’une part, nous nous intéressons aux modèles non-hydrostatique de Saint-Venant prenant en compte des effets dispersifs. D’autre part, nous étudions le modèle de Saint-Venant avec terme source de Coriolis et ses équilibres géostrophiques.

Dans un premier temps nous considérons des méthodes numériques pour une famille de modèles d’Euler moyennés sur la hauteur d’eau provenant de la littérature. Il s’agit de modèles de Saint-Venant avec une pression non-hydrostatique. Les méthodes numériques considérées sont basées sur une méthode de prédiction-correction consistant à décomposer le problème en deux étapes à chaque itération en temps. L’étape de prédiction mène à résoudre un système de Saint-Venant qui est habituellement résolu par une méthode de volumes finis, alors que l’étape de correction mène à résoudre un problème elliptique régissant la pression non-hydrostatique. Notre but premier est d’analyser la convergence d’un schéma mixte avec condensation de masse introduit dans la littérature en utilisant la Méthode de Discrétisation du Gradient (GDM) qui forme un cadre englobant des schémas de discrétisation usuels et récents pour des problèmes de diffusion. Cette méthode nous permet d’obtenir la convergence du schéma proposé pour le problème elliptique. Par la suite, nous proposons une nouvelle formulation conforme du problème en pression et donnons les estimateurs d’erreur correspondant via la GDM. Enfin, nous donnons un exemple d’application avec la méthode des éléments finis  $\mathbb{P}_1$ .

Dans un second temps nous souhaitons avoir des schémas explicites colocalisés volumes finis pour les équations de Saint-Venant non-linéaires avec un terme source de Coriolis qui seraient précis aux alentours de l’équilibre géostrophique et stables dans le cadre non-linéaire. Nous construisons plusieurs schémas volumes finis et étudions les deux propriétés suivantes : la décroissance de l’énergie semi-discrète et la préservation de l’équilibre géostrophique par la version linéarisée. Nous proposons également une version conservative de l’un de nos schémas. Enfin, nous observons le comportement de ces schémas à travers plusieurs cas tests et obtenons de meilleurs résultats en comparaison d’un schéma volumes finis classique.

**Mots-clés :** écoulements à surface libre, équations de Saint-Venant, pression non-hydrostatique, méthode de discrétisation du gradient, éléments finis, force de Coriolis, équilibre géostrophique, schéma équilibre, diffusion numérique, volumes finis.

**Abstract :** In this work we study some numerical methods for free surface shallow water flows. On one hand, we are interested in non-hydrostatic shallow water models which take into account dispersive effects. On the other hand, we study the shallow water model with Coriolis source term and the associated geostrophic equilibrium.

First, we consider numerical methods for a family of depth-averaged Euler models from the literature. These models are shallow water models with non-hydrostatic pressure. The considered numerical methods are based on a prediction-correction method which leads to decompose the problem into two parts at each time step. The prediction part consists in solving a shallow water system which is usually solved by finite volume methods, while the correction part consists in solving an elliptic problem governing the non-hydrostatic pressure. Our first goal is to analyse the convergence of a mixed mass-lumped scheme introduced in the literature through the Gradient Discretisation Method (GDM) which is a framework comprising classic and recent discretisation schemes for diffusion problems. This method allows us to obtain a convergence result on the proposed scheme for the elliptic problem. Then, we propose a new conforming formulation of the problem on the pressure and give the corresponding error estimate provided by the GDM. Finally, an application to the  $\mathbb{P}_1$  conforming finite element method is given.

Then, we aim at designing explicit collocated finite volume schemes for the nonlinear Shallow Water equations with Coriolis source term that are proved to be accurate around the geostrophic equilibrium and stable in the nonlinear framework. We define some finite volume schemes and study the two properties we are interested in : the decrease of the semi-discrete energy and the preservation of the geostrophic equilibrium in the linearized version. We also propose a conservative version of one of our scheme. Finally, we illustrate the behaviour of the schemes for some standard test cases and we exhibit a great improvement when compared to a classic finite volume scheme.

**Keywords :** free surface flows, shallow water equations, non-hydrostatic pressure, Gradient Discretisation Method, finite element, Coriolis force, geostrophic equilibrium, well-balanced scheme, numerical diffusion, finite volume.