



**HAL**  
open science

# WEAK CONSISTENCY OF NON LINEAR CONVECTION OPERATORS ON STAGGERED MESHES. APPLICATION TO A QUASI -SECOND ORDER STAGGERED SCHEME FOR THE TWO-DIMENSIONAL SHALLOW WATER EQUATIONS

T Gallouët, R Herbin, J.-C Latché, Y Nasser, N Therme

► **To cite this version:**

T Gallouët, R Herbin, J.-C Latché, Y Nasser, N Therme. WEAK CONSISTENCY OF NON LINEAR CONVECTION OPERATORS ON STAGGERED MESHES. APPLICATION TO A QUASI -SECOND ORDER STAGGERED SCHEME FOR THE TWO-DIMENSIONAL SHALLOW WATER EQUATIONS. 2020. hal-02940981v1

**HAL Id: hal-02940981**

**<https://hal.science/hal-02940981v1>**

Preprint submitted on 16 Sep 2020 (v1), last revised 17 Nov 2021 (v3)

**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.

# WEAK CONSISTENCY OF NON LINEAR CONVECTION OPERATORS ON STAGGERED MESHES.

## APPLICATION TO A QUASI - SECOND ORDER STAGGERED SCHEME FOR THE TWO-DIMENSIONAL SHALLOW WATER EQUATIONS

T. GALLOUËT<sup>1</sup>, R. HERBIN<sup>2</sup>, J.-C. LATCHÉ<sup>3</sup>, Y. NASSERI<sup>4</sup> AND N. THERME<sup>5</sup>

**Abstract.** In this paper a Lax-Wendroff type result of consistency is given for convection operators on staggered meshes. It is applied to a class of second order finite volume schemes developed to obtain approximate solutions of the shallow water equations with bathymetry. These schemes are based on staggered grids for the space discretization: scalar and vector unknowns are defined on different meshes. MUSCL-like interpolations for the discrete convection operators in the water height and momentum equations are performed in order to improve the precision of the scheme. The time discretization is performed either by a first order segregated forward Euler scheme in time or by the second order Heun scheme. Both schemes are shown to preserve the water height positivity under a CFL condition and an important state equilibrium known as the lake at rest. Using the above mentioned staggered Lax-Wendroff type results, these schemes are shown to be Lax-consistent with the weak formulation of the continuous equations; besides, the forward Euler scheme is shown to be consistent with a weak entropy inequality. Numerical results confirm the efficiency and accuracy of the schemes.

**2010 AMS Subject Classification.** Primary 65M08, 76N15 ; Secondary 65M12, 76N19.

The dates will be set by the publisher.

**Keywords** Finite-volume scheme, MAC grid, shallow water flow.

September 15, 2020

### CONTENTS

1. Introduction	2
2. Space and time discretization	3
2.1. Definitions and notations	3
2.2. The segregated forward Euler scheme	5
2.3. A second order in time Heun scheme	7
3. Stability of the schemes	8
4. Weak consistency of the schemes	10
4.1. Proof of consistency of the forward Euler MAC scheme	11
4.2. Proof of the weak consistency of the Heun scheme	14
4.3. A sufficient condition for the convergence of the intermediate solutions	15
5. Weak entropy consistency of the forward Euler- MAC scheme	18
6. Numerical results	27
6.1. A smooth solution	28

---

*Keywords and phrases:* Finite-volume scheme, MAC grid, shallow water flow.

<sup>1</sup> I2M UMR 7373, Aix-Marseille Université, CNRS, Ecole Centrale de Marseille. 39 rue Joliot Curie. 13453 Marseille, France.

(raphaele.herbin@univ-amu.fr)

<sup>2</sup> I2M UMR 7373, Aix-Marseille Université, CNRS, Ecole Centrale de Marseille. 39 rue Joliot Curie. 13453 Marseille, France.

(raphaele.herbin@univ-amu.fr)

<sup>3</sup> IRSN, BP 13115, St-Paul-lez-Durance Cedex, France (jean-claude.latche@irsn.fr)

<sup>4</sup> I2M UMR 7373, Aix-Marseille Université, CNRS, Ecole Centrale de Marseille. 39 rue Joliot Curie. 13453 Marseille, France.

(yousseouf.nasseri@univ-amu.fr)

<sup>5</sup> CEA/CESTA 33116, Le Barp, France (nicolas.therme@cea.fr)

6.2. A Riemann problem	29
6.3. A circular dam break problem	29
6.4. A so-called partial dam-break problem	31
6.5. Uniform circular motion in a paraboloid	32
Appendix A. Consistency of numerical non linear convection fluxes on staggered meshes	32
Appendix B. Former lemmas	36
B.1. A result on a finite volume convection operator	36
B.2. A result on the space translates	37
References	37

## 1. INTRODUCTION

The shallow water equations form a hyperbolic system of two conservation equations (mass and momentum) which models the flow of an incompressible fluid, assuming that the mean vertical height of the fluid is small compared to the plane scale. It is widely used for the simulation of numerous geophysical phenomena, such as flow in rivers and coastal areas. For a fluid occupying the space-time domain  $\Omega \times (0, T)$ , where  $\Omega$  is an open bounded subset of  $\mathbb{R}^2$  and  $T > 0$ , the shallow water equations with bathymetry solve the water height  $h$  and the (vector) velocity of the fluid  $\mathbf{u} = (u_1, u_2)$  and read:

$$\partial_t h + \operatorname{div}(h\mathbf{u}) = 0 \quad \text{in } \Omega \times (0, T), \quad (1a)$$

$$\partial_t(h\mathbf{u}) + \operatorname{div}(h\mathbf{u} \otimes \mathbf{u}) + \nabla p + gh\nabla z = 0 \quad \text{in } \Omega \times (0, T), \quad (1b)$$

$$p = \frac{1}{2}gh^2 \quad \text{in } \Omega \times (0, T), \quad (1c)$$

$$\mathbf{u} \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \times (0, T), \quad (1d)$$

$$h(\mathbf{x}, 0) = h_0, \quad \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0 \quad \text{in } \Omega. \quad (1e)$$

where  $\partial_t$  is the partial time derivative,  $\operatorname{div}$  denotes the spatial divergence operator,  $g$  is the standard gravity constant and  $z$  the (given) bathymetry, which is supposed to be regular in this paper. The initial conditions are  $h_0 \in L^\infty(\Omega)$  and  $\mathbf{u}_0 = (u_{0,1}, u_{0,2}) \in L^\infty(\Omega, \mathbb{R}^2)$  with  $h_0 \geq 0$ . This system has therefore been intensively studied, both theoretically and numerically, so that it is impossible to give an exhaustive list of references. We refer to the books [6, 29] and to the more recent books or parts of books [3, 8, 30] and the references therein. We recall that it is wellknown that if no dry zone exists, the system is strictly hyperbolic. In all cases, the solution of the system may develop shocks, so that the finite volume method is often preferred for numerical simulations. Two main approaches are found: one is the collocated approach which is usually based on some approximate Riemann solver, see e.g. [6, 8] and references therein; the other one is based on a staggered arrangement of the unknowns on the grid, which is quite classical in the hydraulic and ocean engineering community, see e.g. [2, 5, 28]. These latter staggered schemes have been implemented with an upwind choice for the convection operators and a forward Euler time discretization and analysed in the case of one space dimension [9, 16], following the works on the related barotropic Euler equations, see [22] and references therein. In particular, the weak consistency of the scheme is shown as well as a weak entropy consistency. Let us recall that if  $(h, \mathbf{u})$  is a regular solution of (1), the following elastic potential energy balance and kinetic energy balance are obtained by manipulations on the mass and momentum equations:

$$\partial_t\left(\frac{1}{2}gh^2\right) + \operatorname{div}\left(\frac{1}{2}gh^2\mathbf{u}\right) + \frac{1}{2}gh^2\operatorname{div}\mathbf{u} = 0 \quad (2)$$

$$\partial_t\left(\frac{1}{2}h|\mathbf{u}|^2\right) + \operatorname{div}\left(\frac{1}{2}h|\mathbf{u}|^2\mathbf{u}\right) + \mathbf{u} \cdot \nabla p + gh\mathbf{u} \cdot \nabla z = 0. \quad (3)$$

Summing these equations, we obtain an entropy balance equation:  $\partial_t E + \text{div}\Phi = 0$ , where the entropy-entropy flux pair  $(E, \Phi)$  is given by:

$$E = \frac{1}{2}h|\mathbf{u}|^2 + \frac{1}{2}gh^2 + ghz \text{ and } \Phi = (E + \frac{1}{2}gh^2)\mathbf{u}. \quad (4)$$

For non regular functions the above manipulations are no longer valid, and the entropy inequality  $\partial_t E + \text{div}\Phi \leq 0$  is satisfied in a distributional sense. The weak entropy consistency consists in showing that any possible limit of the scheme satisfies a weak form of the entropy inequality (4) given in (31) below.

In the case of two space dimensions, the consistency of the upwind scheme with respect to the weak formulation and to a weak entropy inequality is stated in [20]; a quasi-second order scheme in time and space using the second order Heun method in time dependent and a MUSCL-like interpolation in space was proposed in [14].

Here, we analyse the former schemes both theoretically and numerically. The framework that is developed here includes three schemes : the first order scheme of [20], the same scheme replacing the upwind choice in the numerical convection operator by a MUSCL-like procedure, and the quasi second order scheme proposed in [14]. Generic properties are shown to be preserved, such as the positivity of the water height and the preservation of the "lake at rest" steady state. The weak consistency of the schemes is proven thanks to a generalisation of Lax-Wendroff type result which is given in an appendix; this consistency result is interesting for its own sake and valid for general convection operators on general collocated or staggered grids in any space dimension. Furthermore, the two first schemes are shown to be entropy-weak consistent in the sense that a weak entropy inequality is satisfied by any possible limit of the scheme as the time and space steps tend to 0, under some CFL condition.

The remainder of the paper is organized as follows: In Section 2 we introduce the space and time discretization. The resulting approximate solutions have some discrete stability and well balance properties which are studied in Section 3. Furthermore, under some convergence and boundedness assumptions, the approximate solutions are shown in Section 4 to converge to a weak solution of (1). This proof of these results heavily relies on the general Lax-Wendroff consistency lemma which is given in the appendix A In Section 5 we consider the first order time discretization and show that any possible limit of the scheme satisfies a weak entropy inequality, again using the consistency result of the appendix. Numerical results comparing the first order scheme of [20], the same scheme replacing the upwind choice in the numerical convection operator by a MUSCL-like procedure, and the quasi second order scheme proposed in [14] are presented in Section 6. Finally, the appendix A contains the general consistency result for a nonlinear convection operator on general meshes with a staggered arrangement of the unknowns, which generalizes the result obtained in [12], while the appendix B contains some technical lemmas which were proved formerly and which are recalled for the sake of completeness.

## 2. SPACE AND TIME DISCRETIZATION

### 2.1. Definitions and notations

We concentrate on the MAC discretization in space, see [17, 18] for some seminal papers and [13] for the convergence analysis of the scheme applied to the incompressible Navier-Stokes equations. This scheme is also widely used by the hydrologist and known as the Arakawa scheme [2].

Let  $\Omega$  be a connected subset of  $\mathbb{R}^2$  consisting in a union of rectangles whose edges are assumed to be orthogonal to the canonical basis vectors, denoted by  $(\mathbf{e}^{(1)}, \mathbf{e}^{(2)})$ .

**Definition 2.1** (MAC discretization). A discretization  $(\mathcal{M}, \mathcal{E})$  of  $\Omega$  with a staggered rectangular grid (or MAC grid), is defined by:

- A primal mesh  $\mathcal{M}$  which consists in a conforming structured, possibly non uniform, rectangular grid of  $\Omega$ . A generic cell of this grid is denoted by  $K$ , and its mass center by  $\mathbf{x}_K$ . The scalar unknowns (water height and pressure) are associated to this mesh.
- A set  $\mathcal{E}$  of all edges of the mesh, with  $\mathcal{E} = \mathcal{E}_{\text{int}} \cup \mathcal{E}_{\text{ext}}$ , where  $\mathcal{E}_{\text{int}}$  (resp.  $\mathcal{E}_{\text{ext}}$ ) are the edges of  $\mathcal{E}$  that lie in the interior (resp. on the boundary) of the domain. The set of edges that are orthogonal to  $\mathbf{e}^{(i)}$

is denoted by  $\mathcal{E}^{(i)}$ , for  $i \in \llbracket 1, 2 \rrbracket$ . We then have  $\mathcal{E}^{(i)} = \mathcal{E}_{\text{int}}^{(i)} \cup \mathcal{E}_{\text{ext}}^{(i)}$ , where  $\mathcal{E}_{\text{int}}^{(i)}$  (resp.  $\mathcal{E}_{\text{ext}}^{(i)}$ ) are the edges of  $\mathcal{E}^{(i)}$  that lie in the interior (resp. on the boundary) of the domain.

For  $\sigma \in \mathcal{E}_{\text{int}}$ , we write  $\sigma = K|L$  if  $\sigma = \partial K \cap \partial L$ . A dual cell  $D_\sigma$  associated to an edge  $\sigma \in \mathcal{E}$  is defined as follows:

- if  $\sigma = K|L \in \mathcal{E}_{\text{int}}$  then  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$ , where  $D_{K,\sigma}$  (resp.  $D_{L,\sigma}$ ) is the half-part of  $K$  (resp.  $L$ ) adjacent to  $\sigma$  (see Fig. 1);
- if  $\sigma \in \mathcal{E}_{\text{ext}}$  is adjacent to the cell  $K$ , then  $D_\sigma = D_{K,\sigma}$ .

For each dimension  $i = 1, 2$ , the domain  $\Omega$  can also be split up in dual cells:  $\Omega = \cup_{\sigma \in \mathcal{E}^{(i)}} \overline{D_\sigma}$ ,  $i \in \llbracket 1, 2 \rrbracket$ ; the  $i^{\text{th}}$  grid is referred to as the  $i^{\text{th}}$  dual mesh; it is associated to the  $i^{\text{th}}$  velocity component, in a sense which is clarified below. The set of the edges of the  $i^{\text{th}}$  dual mesh is denoted by  $\tilde{\mathcal{E}}_i$  (note that these edges may be non-orthogonal to  $\mathbf{e}^{(i)}$ ); the set  $\tilde{\mathcal{E}}_i$  is decomposed into the internal and boundary edges:  $\tilde{\mathcal{E}}_i = \tilde{\mathcal{E}}_{\text{int}}^{(i)} \cup \tilde{\mathcal{E}}_{\text{ext}}^{(i)}$ . The dual edge separating two dual cells  $D_\sigma$  and  $D_{\sigma'}$  is denoted by  $\epsilon = \sigma|\sigma'$ . We denote by  $D_\epsilon$  the cell associated to a dual edge  $\epsilon \in \tilde{\mathcal{E}}$  defined as follows:

- if  $\epsilon = \sigma|\sigma' \in \tilde{\mathcal{E}}_{\text{int}}$  then  $D_\epsilon = D_{\sigma,\epsilon} \cup D_{\sigma',\epsilon}$ , where  $D_{\sigma,\epsilon}$  (resp.  $D_{\sigma',\epsilon}$ ) is the half-part of  $D_\sigma$  (resp.  $D_{\sigma'}$ ) adjacent to  $\epsilon$  (see Fig. 1);
- if  $\epsilon \in \tilde{\mathcal{E}}_{\text{ext}}$  is adjacent to the cell  $D_\sigma$ , then  $D_\epsilon = D_{\sigma,\epsilon}$ .

In order to define the scheme, we need some additional notations. The set of edges of a primal cell  $K$  and of a dual cell  $D_\sigma$  are denoted by  $\mathcal{E}(K) \subset \mathcal{E}$  and  $\tilde{\mathcal{E}}(D_\sigma)$  respectively; note that  $\tilde{\mathcal{E}}(D_\sigma) \subset \tilde{\mathcal{E}}_i$  if  $\sigma \in \mathcal{E}^{(i)}$ . For  $\sigma \in \mathcal{E}$ , we denote by  $\mathbf{x}_\sigma$  the mass center of  $\sigma$ . The vector  $\mathbf{n}_{K,\sigma}$  stands for the unit normal vector to  $\sigma$  outward  $K$ . In some cases, we need to specify the orientation of various geometrical entities with respect to the axis:

- a primal cell  $K$  is denoted  $K = \overrightarrow{[\sigma\sigma']}$  if  $\sigma, \sigma' \in \mathcal{E}^{(i)}(K)$  for some  $i \in \llbracket 1, 2 \rrbracket$  are such that  $(\mathbf{x}_{\sigma'} - \mathbf{x}_\sigma) \cdot \mathbf{e}^{(i)} > 0$ ;
- we write  $\sigma = \overrightarrow{K|L}$  if  $\sigma \in \mathcal{E}^{(i)}$ ,  $\sigma = K|L$  and  $\overrightarrow{\mathbf{x}_K \mathbf{x}_L} \cdot \mathbf{e}^{(i)} > 0$  for some  $i \in \llbracket 1, 2 \rrbracket$ ;
- the dual edge  $\epsilon$  separating  $D_\sigma$  and  $D_{\sigma'}$  is written  $\epsilon = \overrightarrow{\sigma|\sigma'}$  if  $\overrightarrow{\mathbf{x}_\sigma \mathbf{x}_{\sigma'}} \cdot \mathbf{e}^{(i)} > 0$  for some  $i \in \llbracket 1, 2 \rrbracket$ .

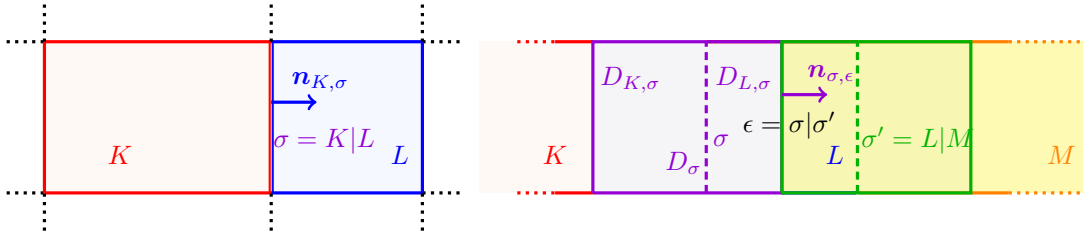


FIGURE 1. Notations for the prima and dual meshes (in two space dimensions, for the first component of the velocity).

The size  $\delta_{\mathcal{M}}$  of the mesh and its regularity  $\theta_{\mathcal{M}}$  are defined by:

$$\delta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \text{diam}(K), \text{ and } \theta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \max_{\sigma \in \mathcal{E}_K} \frac{|D_\sigma|}{|K|}. \quad (5)$$

where  $|\cdot|$  stands for the one (or two) dimensional measure of a subset of  $\mathbb{R}$  (or  $\mathbb{R}^2$ ). Note that in the rectangular case that is considered here, the regularity parameter  $\theta_{\mathcal{M}}$  is also equal to:

$$\theta_{\mathcal{M}} = \frac{1}{2} \left( 1 + \max \left\{ \frac{|\sigma|}{|\sigma'|}, (\sigma, \sigma') \in \mathcal{E}^{(i)2}, i = 1, 2 \right\} \right).$$

The discrete velocity unknowns are associated to the velocity cells and are denoted by  $(u_{i,\sigma})_{\sigma \in \mathcal{E}^{(i)}}$ ,  $i \in \llbracket 1, 2 \rrbracket$ , while the discrete scalar unknowns (water height and pressure) are associated to the primal cells and are denoted respectively by  $(h_K)_{K \in \mathcal{M}}$  and  $(p_K)_{K \in \mathcal{M}}$ .

Let us consider a uniform discretisation  $0 = t_0 < t_1 < \dots < t_N = T$  of the time interval  $(0, T)$ , and let  $\delta t = t_{n+1} - t_n$  for  $n = 0, 1, \dots, N - 1$  be the (constant, for the sake of simplicity) time step.

Here we present two schemes: a first order in time segregated scheme using the forward Euler scheme and the second order in time Heun scheme. Both schemes use a MUSCL-like technique for the computation of the numerical flux, see [27], so that they are quasi second-order in space.

## 2.2. The segregated forward Euler scheme

We propose here a first order in time segregated discretisation and MAC discretization in space of the system (1); the scheme is written in compact form as follows:

$$\mathbf{Initialisation:} \quad u_\sigma^0 = \frac{1}{|D_\sigma|} \int_{D_\sigma} u_{i,0}(\mathbf{x}) \, d\mathbf{x}, \quad h^0 = \frac{1}{|K|} \int_K h_0(\mathbf{x}) \, d\mathbf{x}, \quad p^0 = \frac{1}{2}g(h^0)^2. \quad (6a)$$

**For**  $0 \leq n \leq N - 1$  : solve for  $h^{n+1}$ ,  $p^{n+1}$  and  $\mathbf{u}^{n+1} = (u_i^{n+1})_{i=1,2}$  :

$$\bar{\partial}_t h_K^{n+1} + \operatorname{div}_K(h^n \mathbf{u}^n) = 0, \quad \forall K \in \mathcal{M} \quad (6b)$$

$$p^{n+1} = \frac{1}{2}g(h^{n+1})^2, \quad (6c)$$

$$\bar{\partial}_t(h u_i)_\sigma^{n+1} + \operatorname{div}_{D_\sigma}(h^n \mathbf{u}^n u_i^n) + \bar{\partial}_\sigma p^{n+1} + g h_{\sigma,c}^{n+1} \bar{\partial}_\sigma z = 0, \quad \forall \sigma \in \mathcal{E}_{\text{int}}^{(i)}, i \in \llbracket 1, 2 \rrbracket, \quad (6d)$$

where the different discrete terms and operators introduced here are now defined.

*Discrete time derivative* - In the sequel, we shall denote by  $\bar{\partial}_t v^{n+1}$  the discrete forward time derivative of a given discrete function of time  $v$ , *i.e.*:

$$\bar{\partial}_t v^{n+1} = \frac{v^{n+1} - v^n}{\delta t} \quad (7)$$

*Discrete divergence and gradient operators* - The discrete divergence operator on the primal mesh denoted by  $\operatorname{div}_K$  is defined as follows:

$$\operatorname{div}_K(h\mathbf{u}) = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| \mathbf{F}_\sigma \cdot \mathbf{n}_{K,\sigma}, \quad \text{with } \mathbf{F}_\sigma = h_\sigma \mathbf{u}_\sigma, \quad \text{with } \mathbf{u}_\sigma = u_{i,\sigma} \mathbf{e}^{(i)} \text{ for } \sigma \in \mathcal{E}^{(i)}, i \in \llbracket 1, 2 \rrbracket, \quad (8)$$

and  $h_\sigma$  is approximated by the MUSCL-like interpolation technique with respect to  $\mathbf{u}_\sigma$ ; in the subsequent analysis, we do not need to have an explicit formula for  $h_\sigma$ , but we need the following conditions to be satisfied:

$$\forall K \in \mathcal{M}, \quad \forall \sigma = K|L \in \mathcal{E}_{\text{int}}(K), \\ - \exists \lambda_{K,\sigma} \in [0, 1] : h_\sigma = \lambda_{K,\sigma} h_K + (1 - \lambda_{K,\sigma}) h_L \text{ if } \mathbf{F}_\sigma \cdot \mathbf{n}_{K,\sigma} \geq 0. \quad (9)$$

$$- \exists \alpha_\sigma^K \in [0, 1] \text{ and } M_\sigma^K \in \mathcal{M} : h_\sigma - h_K = \begin{cases} \alpha_\sigma^K (h_K - h_{M_\sigma^K}) & \text{if } \mathbf{u}_\sigma \cdot \mathbf{n}_{K,\sigma} \geq 0, \\ \alpha_\sigma^K (h_{M_\sigma^K} - h_K) & \text{otherwise.} \end{cases} \quad (10)$$

By (9),  $h_\sigma$  is a convex combination of  $h_K$  and  $h_L$ , and if  $\mathbf{u}_\sigma \cdot \mathbf{n}_{K,\sigma} < 0$ , the cell  $M_\sigma^K$  in (10) can be chosen as  $L$  and  $\alpha_{K,\sigma}$  as  $1 - \lambda_{K,\sigma}$ . In the case of a discrete divergence free velocity field  $\mathbf{u}$ , this assumption ensures that  $h_K^{n+1}$  is a convex combination of the values  $h_K^n$  and  $(h_M^n)_{M \in \mathcal{N}_m((K))}$ , where  $\mathcal{N}_m(K)$  denotes the set of cells  $M_\sigma^K$  satisfying (10), see [27, Lemma 3.1], for any structured or unstructured mesh.

Note that if  $K = [\sigma' \sigma]$  with  $\sigma' = J|K$  and  $\sigma = K|L$  and  $\mathbf{u}_\sigma \cdot \mathbf{n}_{K,\sigma} \geq 0$ , the cell  $M_\sigma^K$  in Relation (10) can be chosen as the cell  $J$  and the value  $h_\sigma$  computed using the following limitation procedure:

$$h_\sigma - h_K = \frac{1}{2} \psi(h_L, h_K, h_J), \text{ where}$$

$$\psi(h_L, h_K, h_J) = \begin{cases} \min\text{mod}\left(\frac{h_L - h_J}{2}, \zeta^+(h_L - h_K), \zeta^-(h_K - h_J)\right), & \text{if } (h_L - h_K)(h_K - h_J) > 0, \\ 0, & \text{otherwise,} \end{cases}$$

where the limitation parameters  $\zeta^+, \zeta^-$  are such that  $\zeta^+, \zeta^- \in [0, 2]$ . Observe that if  $\zeta^+ = \zeta^- = 1$ , the classical minmod limiter ( $\min\text{mod}(h_L - h_K, h_K - h_J)$ ) is recovered.

A local discrete derivative applied to a discrete scalar field  $\xi$  (with  $\xi = p, h$  or  $z$ ) is defined by:

$$\bar{\partial}_\sigma \xi = \frac{|\sigma|}{|D_\sigma|} (\xi_L - \xi_K) \text{ for } \sigma = \overrightarrow{K|L} \in \mathcal{E}_{\text{int}}. \quad (11)$$

The above defined discrete divergence and discrete derivatives satisfy the following div-grad duality relationship [13, Lemma 2.4]:

$$\sum_{K \in \mathcal{M}} |K| \xi_K \text{div}_K(h\mathbf{u}) + \sum_{i=1}^2 \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} |D_\sigma| h_\sigma u_{i,\sigma} \bar{\partial}_\sigma \xi = 0. \quad (12)$$

*Discrete water height for the bathymetry term* – In equation (6d) the term  $\bar{\partial}_\sigma z$  denotes the discrete derivative (in the sense of (11)) of the piecewise constant function  $z_\mathcal{M} = \sum_{K \in \mathcal{M}} z(\mathbf{x}_K) \mathbf{1}_K$ , that is:

$$\bar{\partial}_\sigma z = \frac{|\sigma|}{|D_\sigma|} (z(\mathbf{x}_L) - z(\mathbf{x}_K)) \text{ for } \sigma = \overrightarrow{K|L} \in \mathcal{E}_{\text{int}}. \quad (13)$$

The value  $h_{\sigma,c}$  of the water height is defined so as to satisfy:

$$\bar{\partial}_\sigma p + g h_{\sigma,c} \bar{\partial}_\sigma z = 0 \text{ if } \bar{\partial}_\sigma(h + z) = 0, \forall i = 1, 2. \quad (14)$$

This requirement is fulfilled if  $h_{\sigma,c}$  is centered:

$$h_{\sigma,c} = \begin{cases} \frac{1}{2}(h_K + h_L) & \text{for } \sigma = K|L \in \mathcal{E}_{\text{int}}, \\ h_K & \text{for } \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}(K). \end{cases} \quad (15)$$

Indeed, if  $h_{\sigma,c}$  is defined by (15), since  $p = \frac{1}{2}gh^2$ , one has from the definition of the discrete gradient (11), for  $\sigma = K|L$ ,

$$\bar{\partial}_\sigma p + g h_{\sigma,c} \bar{\partial}_\sigma z = \frac{1}{2}g \frac{|\sigma|}{|D_\sigma|} (h_K + h_L) \bar{\partial}_\sigma(h + z)$$

and therefore (14) holds, so that the ‘‘lake at rest’’ steady state is preserved, see Lemma 3.2 below.

*Discrete convection operator* – The term  $(h u_i)_\sigma^{n+1}$  in the discrete time derivative in (6d) is defined by

$$(h u_i)_\sigma^{n+1} = h_{D_\sigma}^{n+1} u_{i,\sigma}^{n+1}, \quad (16a)$$

$$h_{D_\sigma} = \frac{1}{|D_\sigma|} \left( |D_{K,\sigma}| h_K + |D_{L,\sigma}| h_L \right), \text{ with } \sigma = K|L \in \mathcal{E}_{\text{int}}, \quad (16b)$$

where  $D_\sigma$ ,  $D_{K,\sigma}$  and  $D_{L,\sigma}$  are defined in Definition 2.1.

The discrete divergence operator on the dual mesh  $\text{div}_{D_\sigma}$  is given by:

$$\text{div}_{D_\sigma}(h u_i \mathbf{u}) = \frac{1}{|D_\sigma|} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| \mathbf{G}_\epsilon \cdot \mathbf{n}_{\sigma,\epsilon}, \text{ with } \mathbf{G}_\epsilon = \mathbf{F}_\epsilon u_{i,\epsilon}, \quad (17)$$

where

- the flux  $\mathbf{F}_\epsilon$  is computed from the primal numerical mass fluxes; following [19] (see also [21], and [1] for an extension to triangular or quadrangular meshes using low order non-conforming finite element), it is defined as follows:

$$\text{for } \epsilon = \sigma|\sigma', \epsilon \subset K, \quad \mathbf{F}_\epsilon = \frac{1}{2}(\mathbf{F}_\sigma + \mathbf{F}_{\sigma'}), \quad \epsilon \subset K, \quad (\text{left on Figure 2}) \quad (18a)$$

$$\text{for } \epsilon = \sigma|\sigma', \epsilon \not\subset K, \epsilon \subset \tau \cup \tau', \quad \mathbf{F}_\epsilon = \frac{1}{|\epsilon|} \left( \frac{1}{2}|\tau| \mathbf{F}_\tau + \frac{1}{2}|\tau'| \mathbf{F}_{\tau'} \right), \quad (\text{right on Figure 2}), \quad (18b)$$

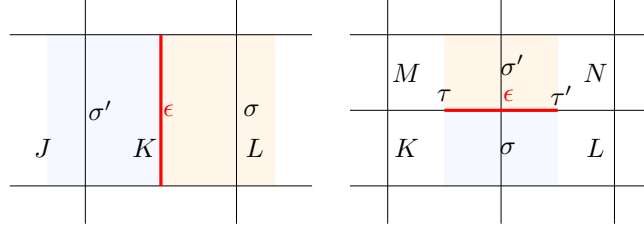


FIGURE 2. Notation for the definition of the momentum flux on the dual mesh for the first component of the velocity- left:  $\epsilon \subset K$  - right:  $\epsilon \subset \tau \cup \tau'$ .

- the value  $u_{i,\epsilon}$  is expressed in terms of the unknowns  $u_{i,\sigma}$ , for  $\sigma \in \mathcal{E}^{(i)}$  by a second order MUSCL-like interpolation scheme with respect to  $\mathbf{F}_\epsilon \cdot \mathbf{n}_{\sigma,\epsilon}$  [27]; the values  $u_{i,\sigma}$  satisfy the following property:

$$\forall \sigma \in \mathcal{E}_{\text{int}}^{(i)}, i = 1, 2, \quad \forall \epsilon = \sigma|\sigma' \in \tilde{\mathcal{E}}(D_\sigma),$$

$$u_{i,\epsilon} \text{ is a convex combination of } u_{i,\sigma} \text{ and } u_{i,\sigma'} : \exists \mu_{\sigma,\epsilon} \in [0, 1] : u_{i,\epsilon} = \mu_{\sigma,\epsilon} u_{i,\sigma} + (1 - \mu_{\sigma,\epsilon}) u_{i,\sigma'} \quad (19)$$

$$\exists \alpha_\epsilon^\sigma \in [0, 1] \text{ and } \tau_\epsilon^\sigma \in \mathcal{E}_{\text{int}}^{(i)} : u_{i,\epsilon} - u_{i,\sigma} = \begin{cases} \alpha_\epsilon^\sigma (u_{i,\sigma} - u_{i,\tau_\epsilon^\sigma}) & \text{if } \mathbf{F}_\epsilon \cdot \mathbf{n}_{\sigma,\epsilon} \geq 0, \\ \alpha_\epsilon^\sigma (u_{i,\tau_\epsilon^\sigma} - u_{i,\sigma}) & \text{otherwise.} \end{cases} \quad (20)$$

Again note that in the case  $\mathbf{F}_\epsilon \cdot \mathbf{n}_{\sigma,\epsilon} < 0$ , the edge  $\tau_\epsilon^\sigma$  may be chosen as  $\sigma'$ .

Let us emphasize that thanks to the definitions (16b) and (18) the following discrete mass balance version on the dual mesh holds:

$$\frac{|D_\sigma|}{\delta t} (h_{D_\sigma}^{n+1} - h_{D_\sigma}^n) + \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} = 0. \quad (21)$$

### 2.3. A second order in time Heun scheme

We retain here the quasi-second order space discretization which we just set up, but consider now a second order time discretization using the Heun (or Runge Kutta 2) scheme.

The initialization of the scheme is the same as that of the forward Euler scheme, see (6a), but the  $n$ -th step now reads:

**Step  $n$**  : For  $h^n$  and  $\mathbf{u}^n = (u_i^n)_{i=1,2}$  known,

$$\hat{h}_K^{n+1} = h_K^n - \delta t \operatorname{div}_K(h^n \mathbf{u}^n), \quad \forall K \in \mathcal{M} \quad (22a)$$

$$\hat{h}_{D_\sigma}^{n+1} \hat{u}_{i,\sigma}^{n+1} = h_{D_\sigma}^n u_{i,\sigma}^n - \delta t \mathcal{F}_{D_\sigma}(h^n, u_i^n), \quad \forall \sigma \in \mathcal{E}_{\text{int}}^{(i)} \quad (22b)$$

$$\tilde{h}_K^{n+1} = \hat{h}_K^{n+1} - \delta t \operatorname{div}_K(\hat{h}^{n+1} \hat{\mathbf{u}}^{n+1}), \quad \forall K \in \mathcal{M} \quad (22c)$$

$$\tilde{h}_{D_\sigma}^{n+1} \tilde{u}_{i,\sigma}^{n+1} = \hat{h}_{D_\sigma}^{n+1} \hat{u}_{i,\sigma}^{n+1} - \delta t \mathcal{F}_{D_\sigma}(\hat{h}^{n+1}, \hat{u}_i^{n+1}), \quad \forall \sigma \in \mathcal{E}_{\text{int}}^{(i)} \quad (22d)$$



$$h_K^{n+1} = \frac{1}{2} (h_K^n + \tilde{h}_K^{n+1}), \quad \forall K \in \mathcal{M} \quad (22e)$$

$$h_{D_\sigma}^{n+1} u_{i,\sigma}^{n+1} = \frac{1}{2} (h_{D_\sigma}^n u_{i,\sigma}^n + \tilde{h}_{D_\sigma}^{n+1} \tilde{u}_{i,\sigma}^{n+1}), \quad \forall \sigma \in \mathcal{E}_{\text{int}}^{(i)} \quad (22f)$$

where

$$\mathcal{F}_{D_\sigma}(h^n, u_i^n) = \text{div}_{D_\sigma}(h^n \mathbf{u}^n u_i^n) + g h_{\sigma,c}^n ((\partial_\sigma h^n) + (\partial_\sigma z)) \quad (23)$$

and the dual cell values  $\hat{h}_{D_\sigma}^{n+1}$ ,  $\tilde{h}_{D_\sigma}^{n+1}$  and  $h_{D_\sigma}^{n+1}$  are computed from the corresponding cell values by the analogue of the formula (16b), so that they satisfy a dual mass balance of the type (21).

The steps (22c)-(22f) of the above scheme (22) may be replaced by the more compact form

$$\partial_t h_K^{n+1} = -\frac{1}{2} \left( \text{div}_K(h^n \mathbf{u}^n) + \text{div}_K(\hat{h}^{n+1} \hat{\mathbf{u}}^{n+1}) \right), \quad \forall K \in \mathcal{M} \quad (24a)$$

$$\partial_t (h_{D_\sigma} u_{i,\sigma})^{n+1} = -\frac{1}{2} \left( \mathcal{F}_{D_\sigma}(h^n, u_i^n) + \mathcal{F}_{D_\sigma}(\hat{h}^{n+1}, \hat{u}_i^{n+1}) \right), \quad \forall \sigma \in \mathcal{E}^{(i)}, \quad (24b)$$

where the dual cell value  $h_{D_\sigma}^{n+1}$  is computed by the formula (16b) and hence satisfies a dual mass balance of the type (21).

### 3. STABILITY OF THE SCHEMES

The positivity of the water height under a CFL like condition is ensured by both the schemes (6) and (22); it is a consequence of the property (10) of the MUSCL choice for the interface values. Indeed, the proof of the positivity in [27, Lemma 3.1] remains valid even if the discrete velocity field is not divergence free, as is the case here.

**Lemma 3.1** (Positivity of the water height). *Let  $n \in \llbracket 0, N-1 \rrbracket$ , let  $(h_K^n)_{K \in \mathcal{M}} \subset \mathbb{R}_+^*$  and  $(\mathbf{u}_\sigma^n)_{\sigma \in \mathcal{E}} \subset \mathbb{R}^d$  be given, and let  $h_K^{n+1}$  be computed by the forward Euler scheme, step (6b). Then  $h_K^{n+1} > 0$ , for all  $K \in \mathcal{M}$  under the following CFL condition,*

$$2 \delta t \leq \frac{|K|}{\sum_{\sigma \in \mathcal{E}(K)} |\sigma| |\mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}|}. \quad (25)$$

If (25) is fulfilled and if furthermore

$$2 \delta t \leq \frac{|K|}{\sum_{\sigma \in \mathcal{E}(K)} |\sigma| |\hat{\mathbf{u}}_\sigma^{n+1} \cdot \mathbf{n}_{K,\sigma}|}, \quad (26)$$

then  $h_K^{n+1}$  computed by the Heun scheme (22) is positive.

Secondly, thanks to the choice (15) for the reconstruction of the water height, the property (14) holds, so that the co-called "lake at rest" steady state is preserved by both schemes.

**Lemma 3.2** (Steady state "lake at rest"). *Let  $n \in \llbracket 0, N-1 \rrbracket$ ,  $C \in \mathbb{R}_+$ ; let  $(h_K^n)_{K \in \mathcal{M}} \subset \mathbb{R}$  such that  $h_K^n + z_K = C$  for all  $K \in \mathcal{M}$  and  $\mathbf{u}_\sigma^n = 0$  for  $\sigma \in \mathcal{E}$ . Then the solution  $(h_K^{n+1})_{K \in \mathcal{M}}$ ,  $(\mathbf{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}}$  of the forward Euler scheme (6) (resp. Heun scheme (22)) satisfies  $h_K^{n+1} + z = C$  for all  $K \in \mathcal{M}$  and  $\mathbf{u}_\sigma^{n+1} = 0$  for  $\sigma \in \mathcal{E}$ .*

As a consequence of the careful discretisation of the convection term, the segregated forward Euler scheme satisfies a discrete kinetic energy balance, as stated in the following lemma. The proof of this result is an easy adaptation of [22, Lemma 3.2].

**Lemma 3.3** (Discrete kinetic energy balance, forward Euler scheme). *A solution to the scheme (6) satisfies the following equality, for  $i = 1, 2$ ,  $\sigma \in \mathcal{E}^{(i)}$  and  $0 \leq n \leq N - 1$ :*

$$\begin{aligned} \frac{|D_\sigma|}{2\delta t} (h_{D_\sigma}^{n+1} (u_{i,\sigma}^{n+1})^2 - h_{D_\sigma}^n (u_{i,\sigma}^n)^2) + \frac{1}{2} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| (u_{i,\epsilon}^n)^2 \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} \\ + |D_\sigma| u_{i,\sigma}^{n+1} (\delta_\sigma p^{n+1}) + |D_\sigma| g h_{\sigma,c}^{n+1} u_{i,\sigma}^{n+1} (\delta_\sigma z) = -R_{i,\sigma}^{n+1}, \end{aligned} \quad (27)$$

with

$$\begin{aligned} R_{i,\sigma}^{n+1} = \frac{1}{2\delta t} |D_\sigma| h_{D_\sigma}^{n+1} (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n)^2 - \frac{1}{2} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \\ + \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n). \end{aligned}$$

The scheme also satisfies the following potential energy balance.

**Lemma 3.4** (Discrete potential balance, forward Euler scheme). *Let, for  $K \in \mathcal{M}$  and  $0 \leq n \leq N$  the potential energy be defined by  $(E_p)_K^n = \frac{1}{2}g(h_K^n)^2 + gh_K^n z_K$ . A solution to the scheme (6) satisfies the following equality, for  $K \in \mathcal{M}$  and  $0 \leq n \leq N - 1$ :*

$$\delta_t (E_p)_K^{n+1} + \operatorname{div}_K \left( \frac{1}{2}g(h^n)^2 \mathbf{u}^n \right) + gz_K \operatorname{div}_K (h^n \mathbf{u}^n) + p_K^n \operatorname{div}_K (\mathbf{u}^n) = -r_K^{n+1}, \quad (28)$$

with  $\delta_t (E_p)_K^{n+1} = \frac{1}{\delta t} ((E_p)_K^{n+1} - (E_p)_K^n)$ ,  $\operatorname{div}_K \left( \frac{1}{2}g(h^n)^2 \mathbf{u}^n \right) = \sum_{\sigma \in \mathcal{E}(K)} |\sigma| \left( \frac{1}{2}g(h_\sigma^n)^2 \right) \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}$  and

$$\begin{aligned} |K| r_K^{n+1} = \frac{1}{2} \frac{|K|}{\delta t} g (h_K^{n+1} - h_K^n)^2 - \frac{1}{2} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| g (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \\ + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| g (h_K^{n+1} - h_K^n) h_\sigma^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}. \end{aligned} \quad (29)$$

*Proof.* Applying [22, Lemma A1], (re-stated in Lemma B.1 below for the sake of completeness), with  $P = K$ ,  $\psi : x \mapsto \frac{1}{2}gx^2$ ,  $\rho_P = h_K^{n+1}$ ,  $\rho_P^* = h_K^n$ ,  $\eta = \sigma$ ,  $\rho_\eta^* = h_\sigma^n$  and  $V_\eta^* = |\sigma| \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}$ , and  $R_K^{n+1} = |K| r_K^{n+1}$ , we get that

$$\begin{aligned} \frac{g}{2} \delta_t (h_K^{n+1})^2 + \operatorname{div}_K \left( \frac{g}{2} (h^n)^2 \mathbf{u}^n \right) + p_K^n \operatorname{div}_K (\mathbf{u}^n) = -\frac{g}{2\delta t} (h_K^{n+1} - h_K^n)^2 + \frac{g}{2} \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \\ - \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| g (h_K^{n+1} - h_K^n) h_\sigma^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}, \end{aligned}$$

Then, multiplying the discrete mass balance equation (6b) by  $gz_K$  yields

$$\delta_t (ghz)_K^{n+1} + \operatorname{div}_K (h^n z \mathbf{u}) + \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)^{(m)}} |\sigma| g (z_K - z_\sigma) h_\sigma^n \mathbf{u}^n \cdot \mathbf{n}_{K,\sigma} = 0$$

Summing the two above equations yields (29).  $\square$

Since the discrete kinetic and potential energies are computed on the dual and primal meshes respectively, the obtention of a discrete entropy inequality is not straightforward. In [20], a kinetic energy inequality on

the primal cell is obtained from the inequality (1d) to get a discrete local entropy inequality. Here, however we proceed otherwise, thanks to a general Lax-Wendroff Lemma for staggered grids (Lemma A.1 in the appendix), which allows to handle each energy inequality on its respective mesh, without any reconstruction, see Section 5 below.

#### 4. WEAK CONSISTENCY OF THE SCHEMES

We now wish to prove the weak consistency of the scheme in the Lax-Wendroff sense, namely to prove that if a sequence of solutions is controlled in suitable norms and converges to a limit, this latter necessarily satisfies a weak formulation of the continuous problem.

The pair of functions  $(\bar{h}, \bar{\mathbf{u}}) \in L^1(\Omega \times [0, T]) \times L^1(\Omega \times [0, T])^2$  is a weak solution to the continuous problem if it satisfies, for any  $\varphi \in C_c^\infty(\Omega \times [0, T])$  ( $\varphi \in C_c^\infty(\Omega \times [0, T])^2$ ):

$$\int_0^T \int_\Omega [\bar{h} \partial_t \varphi + \bar{h} \bar{\mathbf{u}} \cdot \nabla \varphi] \, d\mathbf{x} \, dt + \int_\Omega h_0(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} = 0, \quad (30a)$$

$$\begin{aligned} \int_0^T \int_\Omega [\bar{h} \bar{\mathbf{u}} \cdot \partial_t \boldsymbol{\varphi} + (\bar{h} \bar{\mathbf{u}} \otimes \bar{\mathbf{u}}) : \nabla \boldsymbol{\varphi} + \frac{1}{2} g \bar{h}^2 \operatorname{div} \boldsymbol{\varphi} + g \bar{h} \nabla z \boldsymbol{\varphi}] \, d\mathbf{x} \, dt \\ + \int_\Omega h_0(\mathbf{x}) \mathbf{u}_0(\mathbf{x}) \cdot \boldsymbol{\varphi}(\mathbf{x}, 0) \, d\mathbf{x} = 0. \end{aligned} \quad (30b)$$

A weak solution of (30) is an entropy weak solution if for any nonnegative test function  $\varphi \in C_c^\infty(\Omega \times [0, T], \mathbb{R}_+)$ :

$$\int_0^T \int_\Omega [\bar{E} \partial_t \varphi + \bar{\Phi} \cdot \nabla \varphi] \, d\mathbf{x} \, dt + \int_\Omega E_0(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \geq 0, \quad (31)$$

with

$$\bar{E} = \frac{1}{2} \bar{h} |\bar{\mathbf{u}}|^2 + \frac{1}{2} g \bar{h}^2 + g \bar{h} z \text{ and } \bar{\Phi} = (\bar{E} + \frac{1}{2} g \bar{h}^2) \bar{\mathbf{u}}.$$

Before stating the global weak consistency of the schemes (6) and (22), some definitions and assumptions are needed.

Let  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$  be a sequence of meshes in the sense of Definition 2.1 and let  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  be the associated sequence of solutions of the scheme (6)) defined almost everywhere on  $(\Omega \times [0, T])$  by:

$$\begin{aligned} u_i^{(m)}(\mathbf{x}, t) &= \sum_{n=0}^{N-1} \sum_{\sigma \in (\mathcal{E}^{(i)})^{(m)}} (u_i^{(m)})_\sigma^{n+1} \mathbb{1}_{D_\sigma}(\mathbf{x}) \mathbb{1}_{[t_n, t_{n+1})}(t), \text{ for } i \in [1, 2] \\ h^{(m)}(\mathbf{x}, t) &= \sum_{n=0}^{N-1} \sum_{K \in \mathcal{M}^{(m)}} (h^{(m)})_K^{n+1} \mathbb{1}_K(\mathbf{x}) \mathbb{1}_{[t_n, t_{n+1})}(t), \end{aligned}$$

where  $\mathbb{1}_A$  is the characteristic function of a given set  $A$ , that is  $\mathbb{1}_A(y) = 1$  if  $y \in A$ ,  $\mathbb{1}_A(y) = 0$  otherwise.

**Assumed estimates** - Some boundedness and compactness assumptions on the sequence of discrete solutions  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  are needed in order to prove the Lax-Wendroff type consistency result. First of all we assume that  $h^{(m)} > 0, \forall m \in \mathbb{N}$  which can be obtained under uniform versions of the CFL conditions (25) and (26), thanks to Lemma 3.1. Furthermore, we assume that:

- the water height  $h^{(m)}$  and its inverse are uniformly bounded in  $L^\infty(\Omega \times (0, T))$ , *i.e.* there exists  $C_{\mathcal{M}}^h \in \mathbb{R}_+^*$  such that for  $m \in \mathbb{N}$  and  $0 \leq n < N^{(m)}$ :

$$\frac{1}{C^h} < (h^{(m)})_K^n \leq C^h, \quad \forall K \in \mathcal{M}^{(m)}, \quad (32)$$

– the velocity  $\mathbf{u}^{(m)}$  is also uniformly bounded in  $L^\infty(\Omega \times (0, T))^2$ , i.e. there exists  $C^u \in \mathbb{R}_+^*$  such that

$$|(\mathbf{u}^{(m)})_\sigma^n| \leq C^u, \quad \forall \sigma \in \mathcal{E}^{(m)}. \quad (33)$$

**Theorem 4.1** (Weak consistency of the schemes). *Let  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$  be a sequence of meshes such that  $\delta t^{(m)}$  and  $\delta_{\mathcal{M}^{(m)}} \rightarrow 0$  as  $m \rightarrow +\infty$ ; assume that there exists  $\theta > 0$  such that  $\theta_{\mathcal{M}^{(m)}} \leq \theta$  for any  $m \in \mathbb{N}$  (with  $\theta_{\mathcal{M}^{(m)}}$  defined by (5)).*

*Let  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  be a sequence of solutions to the scheme (6) satisfying (32) and (33) converging to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$ . Then  $(\bar{h}, \bar{\mathbf{u}})$  satisfies the weak formulation (30) of the shallow water equations.*

*Similarly, if  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}, (\hat{h}^{(m)}, \hat{\mathbf{u}}^{(m)})_{m \in \mathbb{N}}$  are sequences of solutions to the scheme (22) both uniformly bounded in the sense of (32) and (33) and converging to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$ , then the limit  $(\bar{h}, \bar{\mathbf{u}})$  satisfies (30).*

The proof of this theorem is the object of the following paragraphs; it relies on some general consistency lemmas which generalize the results of [12] to staggered meshes; these results are independent of the problem at hand and are given in the Appendix A. The proof of the consistency of the schemes is given in Section 4.1 for the forward Euler time discretization and in Section 4.2 for the Heun time discretization.

Note that because the convergence and boundedness of the approximate solutions are assumed, no CFL condition is required in Theorem 4.1. However, recall that a CFL condition is for instance already needed to show the positivity of the water height, see Lemma 3.1.

Finally, in Section 4.3, we give some conditions that imply the boundedness and convergence of the sequence  $(\hat{h}^{(m)}, \hat{\mathbf{u}}^{(m)})_{m \in \mathbb{N}}$  if the boundedness and convergence of the sequence  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  is assumed. One of this condition is a rather strong CFL-like condition.

## 4.1. Proof of consistency of the forward Euler MAC scheme

### 4.1.1. Consistency, mass equation

Under the assumptions of Theorem 4.1, the aim here is to prove that the limit  $(\bar{h}, \bar{\mathbf{u}})$  of the scheme (6) satisfies the weak form of the mass equation (30a). In order to do so, we apply the consistency result of Lemma A.1 in the appendix A, with  $U = (h, \mathbf{u})$ ,  $\beta(U) = h$ ,  $\mathbf{f}(U) = h\mathbf{u}$ ,  $\mathcal{P}^{(m)} = \mathcal{M}^{(m)}$ ,  $\mathfrak{F}^{(m)} = \mathcal{E}^{(m)}$ , and

$$\begin{aligned} \mathcal{C}_{\text{MASS}}^{(m)}(U^{(m)}) : \quad \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \delta_t (h^{(m)})_K^{n+1} + \operatorname{div}_K ((h^{(m)})^n \mathbf{u}^n) \text{ for } \mathbf{x} \in K \text{ and } t \in (t_n, t_{n+1}) \end{aligned} \quad (34)$$

We first note that the assumptions (32) and (33) imply that (77) holds. Furthermore, the assumption of Theorem 4.1 that  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  is a sequence of solutions to the scheme (6) converging to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$  implies that (78) holds.

By the initialisation (6a) of the scheme, it is clear that

$$\sum_{K \in \mathcal{M}^{(m)}} \int_K |(h^{(m)})_K^0 - h_0(\mathbf{x})| d\mathbf{x} = 0,$$

so that the assumption (80) is satisfied.

Since for any  $n \in \llbracket 0, N_m - 1 \rrbracket$  and  $K \in \mathcal{M}$ , one has  $\beta(U^{(m)}(\mathbf{x}, t)) = h_K^n$  for any  $(\mathbf{x}, t) \in K \times [t_n, t_{n+1})$  and  $(\beta^{(m)})_K^n = h_K^n$ ,

$$\sum_{n=0}^{N_m-1} \sum_{K \in \mathcal{M}^{(m)}} \int_{t_n}^{t_{n+1}} \int_K |(h^{(m)})_K^n - h^{(m)}(\mathbf{x}, t)| d\mathbf{x} dt = 0,$$

and therefore the assumption (81) is also clearly satisfied. Now  $(\mathbf{F}^{(m)})_\sigma^n = h_\sigma^n \mathbf{u}_\sigma^n$  and, because the velocity components are piecewise constant on different grids,

$$\begin{aligned} \mathbf{f}(U^m(\mathbf{x}, t)) &= (f_1(U^m(\mathbf{x}, t)), f_2(U^m(\mathbf{x}, t))), \text{ with} \\ f_i(U^m(\mathbf{x}, t)) &= \begin{cases} h_K^n u_{i,\sigma}^n & \text{if } \mathbf{x} \in D_{K,\sigma} \\ h_K^n u_{i,\sigma'}^n & \text{if } \mathbf{x} \in D_{K,\sigma'}, \end{cases} \quad \text{with } K = [\sigma\sigma'] \text{ and where } \sigma \text{ and } \sigma' \perp \mathbf{e}^{(i)}. \end{aligned}$$

For  $\mathbf{x} \in K$ ,  $\sigma = K|L$  and  $t \in [t_n, t_{n+1})$ ,

$$\begin{aligned} \left| \left( (\mathbf{F}^{(m)})_\sigma^n - \mathbf{f}(U^m(\mathbf{x}, t)) \right) \cdot \mathbf{n}_{K,\sigma} \right| &= \left| \left( h_\sigma^n \mathbf{u}_\sigma^n - h_K^n \mathbf{u}_\sigma^n + h_K^n \mathbf{u}_\sigma^n - h_K^n \mathbf{u}(\mathbf{x}, t) \right) \cdot \mathbf{n}_{K,\sigma} \right| \\ &\leq C^u |h_K - h_L| + C^h |\mathbf{u}_\sigma - \mathbf{u}_{\sigma'}|. \end{aligned}$$

Thanks to Lemma B.2 ( [12, Lemma 4.2], recalled in Lemma B.2 in the appendix below) we have

$$\sum_{n=0}^{N_m-1} \sum_{K \in \mathcal{M}^{(m)}} \int_{t_n}^{t_{n+1}} \frac{\text{diam}(K)}{|K|} \int_K |\sigma| \left| \left( h_\sigma^n \mathbf{u}_\sigma^n - h(\mathbf{x}, t) \mathbf{u}(\mathbf{x}, t) \right) \cdot \mathbf{n}_{K,\sigma} \right| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty.$$

so that the assumption (82) is also satisfied.

Hence, by Lemma A.1,

$$\begin{aligned} \forall \varphi \in C_c^\infty(\Omega \times [0, T]), \int_0^T \int_\Omega \mathfrak{e}_{\text{MASS}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) d\mathbf{x} dt \rightarrow \\ - \int_\Omega h_0(\mathbf{x}) \varphi(\mathbf{x}, 0) d\mathbf{x} - \int_0^T \int_\Omega \left[ \bar{h}(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + \bar{h}(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right] d\mathbf{x} dt \text{ as } m \rightarrow +\infty. \end{aligned} \quad (35)$$

From (6b) and (35), we conclude that the limit  $(\bar{h}, \bar{\mathbf{u}})$  of the approximate solutions defined by the forward Euler scheme (6) satisfies (30a).

#### 4.1.2. Consistency, momentum equation

Let  $\varphi = (\varphi_1, \dots, \varphi_d) \in (C_c^\infty(\Omega \times [0, T]))^d$  be a test function and let  $\varphi_{i,\sigma}^{n+1}$  denote the mean value of  $\varphi_i$  over  $\sigma \times (t_n, t_{n+1})$ . Multiplying the equation (6d) by  $|D_\sigma| \varphi_{i,\sigma}^{n+1}$ , summing the result over  $\sigma \in \mathcal{E}^{(i)}$  and then summing over  $n \in [0, N-1]$  and  $i = 1, 2$  yields:

$$\sum_{i=1}^2 Q_{1,i}^{(m)} + Q_{2,i}^{(m)} + Q_{3,i}^{(m)} + Q_{4,i}^{(m)} = 0, \quad (36)$$

with (dropping the exponents  $(m)$  in the summations for the sake of simplicity)

$$Q_{1,i}^{(m)} = \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in (\mathcal{E}^{(m)})^{(i)}} |D_\sigma| \bar{\partial}_t (h u_i)_{\sigma}^{n+1}, \quad (37)$$

$$Q_{2,i}^{(m)} = \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in (\mathcal{E}^{(m)})^{(i)}} |D_\sigma| \text{div}_{D_\sigma} (h^n \mathbf{u}^n u_i^n) \varphi_{i,\sigma}^{n+1}, \quad (38)$$

$$Q_{3,i}^{(m)} = \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in (\mathcal{E}^{(m)})^{(i)}} |D_\sigma| \bar{\partial}_\sigma p^{n+1} \varphi_{i,\sigma}^{n+1}, \quad (39)$$

$$Q_{4,i}^{(m)} = \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in (\mathcal{E}^{(m)})^{(i)}} |D_\sigma| g h_{\sigma,c}^{n+1} \bar{\partial}_\sigma z \varphi_{i,\sigma}^{n+1}. \quad (40)$$

**The nonlinear convection operator.** In order to study the limit of the discrete non linear convection operator defined by  $Q_i^{(m)} = Q_{1,i}^{(m)} + Q_{2,i}^{(m)}$ , we apply Lemma A.1 with  $U = (h, \mathbf{u})$ ,  $\beta(U) = hu_i$ ,  $\mathbf{f}(U) = h\mathbf{u}u_i$ , with  $\mathcal{P}^{(m)}$  the set of dual cells associated with  $u_i$  (that is with the cells corresponding to the vertical edges for  $i = 1$  and the horizontal edges for  $i = 2$ ), with  $\mathfrak{F} = \tilde{\mathcal{E}}_i^{(m)}$  and with the dual fluxes  $(\mathbf{G})_\epsilon^n$  defined by (18). The discrete non linear convection operator thus reads

$$\begin{aligned} [\mathcal{C}_{\text{MOM}}^{(m)}(U^{(m)})]_i : \quad & \Omega \times (0, T) \rightarrow \mathbb{R}, \\ (\mathbf{x}, t) \mapsto & \partial_t (hu_i)_\sigma^{n+1} - \operatorname{div}_{D_\sigma} (h^n u_i \mathbf{u}^n) \text{ for } \mathbf{x} \in D_\sigma \text{ and } t \in (t_n, t_{n+1}) \end{aligned} \quad (41)$$

(again dropping the exponents  $^{(m)}$  for the sake of simplicity).

Again, by the initialisation of the scheme (6a) and by the definition of  $(h\mathbf{u})_{i,\sigma}^0$  (see (16)), it is clear that

$$\sum_{\sigma \in \tilde{\mathcal{E}}_i^{(m)}} \int_{D_\sigma} |(h\mathbf{u})_{i,\sigma}^0 - h_0(\mathbf{x})u_{i,0}(\mathbf{x})| d\mathbf{x} = 0 \text{ and } \sum_{n=0}^{N_m-1} \sum_{\sigma \in \mathcal{E}^{(m)}} \int_{t_n}^{t_{n+1}} \int_{D_\sigma} |(h\mathbf{u})_{i,\sigma}^n - h(\mathbf{x}, t)u_i(\mathbf{x}, t)| d\mathbf{x} dt = 0, \quad i = 1, 2.$$

so that the assumptions (80) and (81) are satisfied.

In order to show that the assumption (82) is satisfied, we need to show that

$$\sum_{n=0}^{N_m-1} \sum_{\sigma \in \mathcal{E}^{(m)}} \int_{t_n}^{t_{n+1}} \frac{\operatorname{diam}(D_\sigma)}{|D_\sigma|} \int_{D_\sigma} |\epsilon| \left| \sum_{\epsilon \in \tilde{\mathcal{E}}_i^{(m)}} \left( (\mathbf{G}^{(m)})_\epsilon^n - h(\mathbf{x}, t)u_i(\mathbf{x}, t)\mathbf{u}(\mathbf{x}, t) \right) \cdot \mathbf{n}_{\sigma,\epsilon} \right| d\mathbf{x} dt \rightarrow 0$$

as  $m \rightarrow +\infty$ . (42)

Let us then estimate, for any  $\epsilon \in \tilde{\mathcal{E}}_i^{(m)}$ ,  $n \in \llbracket 0, N_m - 1 \rrbracket$  and  $\mathbf{x} \in D_\sigma$  the quantity  $Y_\epsilon^n$  defined by:

$$Y_\epsilon^n(\mathbf{x}) = \left| \left( (\mathbf{G}^{(m)})_\epsilon^n - h(\mathbf{x}, t)u_i(\mathbf{x}, t)\mathbf{u}(\mathbf{x}, t) \right) \cdot \mathbf{n}_{\sigma,\epsilon} \right|.$$

Let  $L$  be the (primal) cell such that  $\sigma = K|L$ .

- (1) If  $\epsilon = \sigma'| \sigma \subset K$ , then  $(\mathbf{G}^{(m)})_\epsilon^n$  is defined by (18a). By the triangular inequality and thanks to the assumptions (9), (19),(32), and (33), we get that

$$Y_\epsilon^n(\mathbf{x}) \leq \frac{1}{2}(C^u)^2 |h_K - h_L| + \frac{1}{2}(C^u)^2 |h_K - h_J| + C^h C^u |u_{\sigma,i} - u_{\sigma',i}|, \quad \forall \mathbf{x} \in D_\sigma,$$

where  $J$  is the (primal) cell such that  $\sigma' = J|K$ , see Figure 2, left.

- (2) If  $\epsilon \subset K$ , then  $(\mathbf{G}^{(m)})_\epsilon^n$  is defined by (18b). Again by the triangular inequality and thanks to the assumptions (9), (19),(32), and (33), we get that

$$Y_\epsilon^n(\mathbf{x}) \leq \frac{1}{2}(C^u)^2 |h_K - h_M| + \frac{1}{2}(C^u)^2 |h_K - h_N| + C^h C^u |u_{\sigma,i} - u_{\sigma',i}|, \quad \forall \mathbf{x} \in D_\sigma,$$

where  $M$  and  $N$  are the two (primal) cells such that  $\tau = K|M$  and  $\tau' = L|N$ , as depicted on Figure 2, right.

Now recall that the sequence of meshes is assumed to be regular in the sense that  $\theta^{(m)} \leq \theta$  with  $\theta^{(m)}$  defined by (5); therefore, since the sequences  $h^{(m)}$  and  $\mathbf{u}^{(m)}$  converge in  $L^1$  as  $m$  tends to  $+\infty$ , we may again apply

Lemma B.2 below, to get that (42) holds. Hence, owing to Lemma A.1, we get that

$$\begin{aligned}
Q_i^{(m)} &= Q_{1,i}^{(m)} + Q_{2,i}^{(m)} = \int_0^T \int_{\Omega} \mathfrak{C}_{\text{MOM}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \rightarrow \\
&\int_0^T \int_{\Omega} \left[ \bar{h}(\mathbf{x}, t) \bar{u}_i(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + \bar{h}(\mathbf{x}, t) \bar{u}_i(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right] d\mathbf{x} \, dt \\
&\quad + \int_{\Omega} h_0(\mathbf{x}) u_{i,0}(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \text{ as } m \rightarrow +\infty. \quad (43)
\end{aligned}$$

**Pressure gradient and bathymetry.** Let us now study the terms  $Q_{3,i}^{(m)}$  and  $Q_{4,i}^{(m)}$  defined by (39) and (40). By the definition (11) of  $\bar{\partial}_{\sigma} p$  and by conservativity, we have (again dropping the exponents  $(m)$ )

$$\begin{aligned}
\sum_{i=1}^2 Q_{3,i}^{(m)} &= - \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{K \in \mathcal{M}^{(m)}} \sum_{\sigma \in \mathcal{E}(K)} p_K^{n+1} \int_{\sigma} \varphi_{\sigma}^{n+1} \cdot \mathbf{n}_{K,\sigma} \\
&= - \int_0^T \int_{\Omega} p^{(m)}(\mathbf{x}, t) \operatorname{div} \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt
\end{aligned}$$

Since the sequence  $(h^{(m)})_{m \in \mathbb{N}}$  is bounded in  $L^{\infty}(\Omega \times (0, T))$  and converges to  $\bar{h}$  in  $L^1(\Omega \times (0, T))$ , the sequence  $(p^{(m)})_{m \in \mathbb{N}}$  converges to  $\bar{p} = \frac{1}{2} g \bar{h}^2$  in  $L^1(\Omega \times (0, T))$  as  $m \rightarrow +\infty$ . Hence we get

$$\sum_{i=1}^2 Q_{3,i}^{(m)} \rightarrow \int_0^T \int_{\Omega} \bar{p}(\mathbf{x}, t) \operatorname{div} \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty. \quad (44)$$

Let us now turn to the bathymetry term  $Q_{4,i}^{(m)}$ , which may be written

$$Q_{4,i}^{(m)} = \int_0^T \int_{\Omega} \tilde{h}^{(m)}(\mathbf{x}, t) \bar{\partial}_i^{(m)} z(\mathbf{x}) \tilde{\varphi}_i^{(m)}(\mathbf{x}, t) \, d\mathbf{x} \, dt,$$

where

- the function  $\tilde{h}^{(m)} : \Omega \times (0, T) \rightarrow \mathbb{R}$  is defined by  $\tilde{h}(\mathbf{x}, t) = h_{\sigma,c}^{n+1} = \frac{1}{2}(h_K^{n+1} + h_L^{n+1})$  for  $\mathbf{x} \in D_{\sigma}$  and  $t \in (t_n, t_{n+1})$ ; the sequence  $(\tilde{h}^{(m)})_{m \in \mathbb{N}}$  is therefore bounded in  $L^{\infty}(\Omega \times (0, T))$  and converges to  $\bar{h}$  in  $L^1(\Omega \times (0, T))$ ;
- the function  $\tilde{\varphi}_i^{(m)} : \Omega \times (0, T) \rightarrow \mathbb{R}$  is defined by  $\tilde{\varphi}_i^{(m)}(\mathbf{x}, t) = \varphi_{\sigma}^{n+1}$  for  $\mathbf{x} \in D_{\sigma}$  and  $t \in (t_n, t_{n+1})$ ; by the regularity of  $\varphi$ , the sequence  $(\tilde{\varphi}_i^{(m)})_{m \in \mathbb{N}}$  converges to  $\varphi_i$  uniformly.
- by (13), the function  $\bar{\partial}_i^{(m)} z : \Omega \rightarrow \mathbb{R}$  is defined by  $\bar{\partial}_i^{(m)} z = \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \frac{|\sigma|}{|D_{\sigma}|} (z(\mathbf{x}_L) - z(\mathbf{x}_K)) \mathbb{1}_{D_{\sigma}}$ . Since  $z$  is a regular function, the sequence of functions  $(\bar{\partial}_i^{(m)} z)_{m \in \mathbb{N}}$  converges uniformly to the derivative  $\partial_i z$  of  $z$  with respect to the  $i$ -th variable as  $m \rightarrow +\infty$ .

Hence

$$Q_{4,i}^{(m)} \rightarrow \int_0^T \int_{\Omega} \bar{h}(\mathbf{x}, t) \partial_i z(\mathbf{x}) \varphi_i(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow \infty. \quad (45)$$

**Limit of the momentum equation.** Passing to the limit in (36) as  $m \rightarrow +\infty$ , using (43), (44) and (45), we get that the limit  $(\bar{h}, \bar{\mathbf{u}})$  of the approximate solutions defined by the forward Euler scheme (6) satisfies (30b), which concludes first part of the proof of Theorem 4.1.

## 4.2. Proof of the weak consistency of the Heun scheme

### 4.2.1. Mass balance

Under the assumptions of Theorem 4.1, the aim here is to prove that the limit  $(\bar{h}, \bar{\mathbf{u}})$  of the scheme (22a)-(22f) satisfies the weak form of the mass equation (30a). In order to do so, we consider the equivalent mass

equation (24a). Because of the structure of the scheme, we cannot use here Lemma A.1 straightforwardly as in the case of the forward Euler scheme. We apply Lemma A.2 with  $U = (h, \mathbf{u})$ ,  $\beta(U) = h$ ,  $\mathbf{f}(U) = h\mathbf{u}$ ,  $\mathcal{P}^{(m)} = \mathcal{M}^{(m)}$ ,  $\mathfrak{F}^{(m)} = \mathcal{E}^{(m)}$  and then Lemma A.3 twice: once with  $U = (h, \mathbf{u})$ ,  $\mathbf{f}(U) = h\mathbf{u}$ ,  $\mathcal{P}^{(m)} = \mathcal{M}^{(m)}$ ,  $\mathfrak{F}^{(m)} = \mathcal{E}^{(m)}$ , and then with  $U = (\hat{h}, \hat{\mathbf{u}})$ ,  $\mathbf{f}(U) = \hat{h}\hat{\mathbf{u}}$ . Thanks to the arguments developed in Section 4.1.1, it is easy to check that in each case, the assumptions of the lemmas are satisfied, so that we can conclude that  $(\bar{h}, \bar{\mathbf{u}})$  satisfies (30a).

#### 4.2.2. Momentum balance

Still under the assumptions of Theorem 4.1, we now prove that the limit  $(\bar{h}, \bar{\mathbf{u}})$  of the scheme (22a)-(22f) satisfies the weak form of the mass equation (30b). Again we consider the equivalent momentum equation (24b). Multiplying the equation (24b) by  $|D_\sigma|\varphi_{i,\sigma}^{n+1}$ , summing the result over  $\sigma \in \mathcal{E}^{(i)}$  and then summing over  $n \in \llbracket 0, N-1 \rrbracket$  and  $i = 1, 2$  yields:

$$\sum_{i=1}^2 \left[ Q_{1,i}^{(m)} + \frac{1}{2}(Q_{2,i}^{(m)} + \widehat{Q}_{2,i}^{(m)} + Q_{3,i}^{(m)} + \widehat{Q}_{3,i}^{(m)}) + Q_{4,i}^{(m)} + \widehat{Q}_{4,i}^{(m)} \right] = 0, \quad (46)$$

where  $Q_{1,i}^{(m)}, \dots, Q_{4,i}^{(m)}$  are defined by (37)-(40), and  $\widehat{Q}_{2,i}^{(m)}, \widehat{Q}_{3,i}^{(m)}, \widehat{Q}_{4,i}^{(m)}$  are defined by (38)-(40), replacing the unknowns  $h, p, u$  by  $\hat{h}, \hat{p}, \hat{u}$ .

Again, because of the structure of the scheme, we cannot use Lemma A.1 directly: we use Lemma A.2 for the time derivative term  $Q_{1,i}^{(m)}$  and Lemma A.3 for the terms  $Q_{2,i}^{(m)}$  and  $\widehat{Q}_{2,i}^{(m)}$ , with  $\mathcal{P}^{(m)}$  the set of dual cells associated with  $u_i$  (that is with the vertical edges for  $i = 1$  and the horizontal edges for  $i = 2$ ), with  $\mathfrak{F} = \widehat{\mathcal{E}}_i^{(m)}$  and with the dual fluxes  $(\mathbf{G})_\varepsilon^n$  defined by (18). We first apply Lemma A.2 with  $U = (h, \mathbf{u})$ ,  $\beta(U) = hu_i$ ,  $\mathbf{f}(U) = h\mathbf{u}u_i$ , and then Lemma A.3, once with  $U = (h, \mathbf{u})$ ,  $\beta(U) = hu_i$ ,  $\mathbf{f}(U) = h\mathbf{u}u_i$  and then with  $U = (\hat{h}, \hat{\mathbf{u}})$ ,  $\beta(U) = \hat{h}\hat{u}_i$ ,  $\mathbf{f}(U) = \hat{h}\hat{\mathbf{u}}\hat{u}_i$ . Thanks to the arguments developed in Section 4.1.2, it is easy to check that in each case, the assumptions of the lemmas are satisfied, so that

$$Q_{1,i}^{(m)} + \frac{1}{2}(Q_{2,i}^{(m)} + \widehat{Q}_{2,i}^{(m)}) \rightarrow \int_0^T \int_\Omega \left[ \bar{h}(\mathbf{x}, t) \bar{u}_i(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + \bar{h}(\mathbf{x}, t) \bar{u}_i(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right] d\mathbf{x} dt + \int_\Omega h_0(\mathbf{x}) u_{i,0}(\mathbf{x}) \varphi(\mathbf{x}, 0) d\mathbf{x} \text{ as } m \rightarrow +\infty. \quad (47)$$

The proof of convergence of the pressure gradient and bathymetry terms  $(Q_{3,i}^{(m)}, Q_{4,i}^{(m)}, \widehat{Q}_{3,i}^{(m)})$  and  $\widehat{Q}_{4,i}^{(m)}$  follow the exact same lines as that of the terms  $(Q_{3,i}^{(m)})$  and  $Q_{4,i}$  in Section 4.1.2. Hence

$$\sum_{i=1}^2 \frac{1}{2} (Q_{3,i}^{(m)} + \widehat{Q}_{3,i}^{(m)} + Q_{4,i}^{(m)} + \widehat{Q}_{4,i}^{(m)}) \rightarrow \int_0^T \int_\Omega \left( \bar{p}(\mathbf{x}, t) \operatorname{div} \varphi(\mathbf{x}, t) + \bar{h}(\mathbf{x}, t) \nabla z(\mathbf{x}) \cdot \varphi(\mathbf{x}, t) \right) d\mathbf{x} dt \text{ as } m \rightarrow +\infty. \quad (48)$$

Therefore, owing to (47) and (48), we may pass to the limit in (46) and conclude that  $(\bar{h}, \bar{\mathbf{u}})$  satisfies (30b). This concludes the proof of Theorem 4.1.

### 4.3. A sufficient condition for the convergence of the intermediate solutions

In Theorem 4.1, we assumed the boundedness and convergence of both sequences  $(h^{(m)}, u^{(m)})$  and  $(\hat{h}^{(m)}, \hat{u}^{(m)})$ . In fact, under a restricted CFL condition, we may prove that the convergence and boundedness of the sequence  $(h^{(m)}, u^{(m)})$  implies the convergence and boundedness of the sequence  $(\hat{h}^{(m)}, \hat{u}^{(m)})$ .

**Lemma 4.2** (Bound on the intermediate step, Heun scheme). *Let  $n \in \llbracket 0, N-1 \rrbracket$ , let  $(h_K^n)_{K \in \mathcal{M}} \subset \mathbb{R}_+^*$  and  $(u_\sigma^n)_{\sigma \in \mathcal{E}} \subset \mathbb{R}^d$  be given. Assume that there exists  $\zeta \in (0, 1)$  such that the following restricted CFL-like*



condition holds (note that it slightly more restrictive than (25)):

$$2 \delta t \leq \zeta \frac{|K|}{\sum_{\sigma \in \mathcal{E}(K)} |\sigma| |\mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}|} \text{ for all } K \in \mathcal{M}. \quad (49)$$

Let  $C_{\mathcal{M}}^{\delta t}$ ,  $C_{\mathcal{M}}^h$  and  $C_{\mathcal{M}}^u \in \mathbb{R}_+^*$  be such that

$$\delta t \leq C_{\mathcal{M}}^{\delta t} \min_{\sigma \in \mathcal{E}} |\sigma|, \quad (50a)$$

$$\frac{1}{C_{\mathcal{M}}^h} \leq h_K^n \leq C_{\mathcal{M}}^h, \forall n \in \llbracket 0, N-1 \rrbracket, \forall K \in \mathcal{M}, \quad (50b)$$

$$\max_{\sigma \in \mathcal{E}} |\mathbf{u}_\sigma^n| \leq C_{\mathcal{M}}^u, \forall n \in \llbracket 0, N-1 \rrbracket. \quad (50c)$$

Then the solutions  $(\widehat{h}_K^{n+1})_{K \in \mathcal{M}}$   $(\widehat{\mathbf{u}}_\sigma^n)_{\sigma \in \mathcal{E}}$  of the Heun steps (22a)-(22b) satisfy:

$$\frac{1-\zeta}{C_{\mathcal{M}}^h} \leq \widehat{h}_K^{n+1} \leq 2C_{\mathcal{M}}^h \quad \forall K \in \mathcal{M}, \quad (51a)$$

$$|\widehat{\mathbf{u}}_\sigma^{n+1}| \leq C_{\mathcal{M}}^u + C_{\mathcal{M}}^{\delta t} \frac{(C_{\mathcal{M}}^h)^2}{1-\zeta} \left( 4(C_{\mathcal{M}}^u)^2 + g(C_{\mathcal{M}}^h + \|z\|_\infty) \right). \quad \forall \sigma \in \mathcal{E}, \quad (51b)$$

*Proof.* From (22a) and by the definition (8) of the discrete divergence, we have

$$\widehat{h}_K^{n+1} = h_K^n - \sum_{\sigma \in \mathcal{E}(K)} (\omega_{K,\sigma}^n)^+ h_\sigma^n + \sum_{\sigma \in \mathcal{E}(K)} (\omega_{K,\sigma}^n)^- h_\sigma^n \text{ with } \omega_{K,\sigma}^n = \delta t \frac{|\sigma|}{|K|} \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}.$$

Thanks to the condition (10),  $\exists \alpha_\sigma^K \in [0, 1]$  and  $M_\sigma^K \in \mathcal{M}$  such that  $h_\sigma^n - h_K^n = \alpha_\sigma^K (h_K^n - h_{M_\sigma^K}^n)$  if  $\omega_{K,\sigma}^n \geq 0$ , and therefore

$$h_\sigma^n = h_K^n (1 + \alpha_\sigma^K) - \alpha_\sigma^K h_{M_\sigma^K}^n.$$

Hence

$$\widehat{h}_K^{n+1} = \left( 1 - \sum_{\sigma \in \mathcal{E}(K)} (\omega_{K,\sigma}^n)^+ (1 + \alpha_\sigma^K) \right) h_K^n + \sum_{\sigma \in \mathcal{E}(K)} (\omega_{K,\sigma}^n)^- h_\sigma^n + \sum_{\sigma \in \mathcal{E}(K)} \alpha_\sigma^K (\omega_{K,\sigma}^n)^+ h_{M_\sigma^K}^n.$$

Therefore, thanks to the condition (49), we get (51a).

Let us now prove (51b); from (22b) we have

$$\begin{aligned} \widehat{u}_{i,\sigma}^{n+1} &= \frac{1}{\widehat{h}_{D_\sigma}^{n+1}} \left( h_{D_\sigma}^n u_{i,\sigma}^n - \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} u_{i,\epsilon}^n \right) - \frac{\delta t |\sigma|}{|D_\sigma|} \frac{gh_{\sigma,c}^n}{\widehat{h}_{D_\sigma}^{n+1}} (h_L^n - h_K^n + z_L - z_K) \\ &= \frac{1}{\widehat{h}_{D_\sigma}^{n+1}} \left[ \left( h_{D_\sigma}^n - \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} \right) u_{i,\sigma}^n - \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) \right] \\ &\quad - \frac{\delta t |\sigma|}{|D_\sigma|} \frac{gh_{\sigma,c}^n}{\widehat{h}_{D_\sigma}^{n+1}} (h_L^n - h_K^n + z_L - z_K). \end{aligned}$$

Since the values  $\widehat{h}_{D_\sigma}^{n+1}$  and  $\widehat{h}_{D_\sigma}^n$  are computed by an equivalent formula to (17), they satisfy a discrete dual mass balance of the type (21), and therefore:

$$\widehat{u}_{i,\sigma}^{n+1} = u_{i,\sigma}^n - \frac{1}{\widehat{h}_{D_\sigma}^{n+1}} \left[ \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) \right] - \frac{\delta t |\sigma|}{|D_\sigma|} \frac{gh_{\sigma,c}^n}{\widehat{h}_{D_\sigma}^{n+1}} (h_L^n - h_K^n + z_L - z_K).$$

Thanks to the CFL condition (50a) and to the bounds on  $\mathbf{u}_\sigma^n$  and  $\widehat{h}_{D_\sigma}^{n+1}$  for all  $\sigma$  (recall that for  $\sigma = K|L$ ,  $\widehat{h}_{D_\sigma}^{n+1}$  is a convex combination of  $\widehat{h}_K^{n+1}$  and  $\widehat{h}_L^{n+1}$ ),

$$\frac{1}{\widehat{h}_{D_\sigma}^{n+1}} \left| \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) \right| \leq \frac{4C_{\mathcal{M}}^{\delta t} t (C_{\mathcal{M}}^u)^2 (C_{\mathcal{M}}^h)^2}{1 - \zeta}.$$

Furthermore, since  $2h_{\sigma,c}^n = h_K^n + h_L^n$  and again owing to (50a),

$$\frac{\delta t |\sigma|}{|D_\sigma|} \frac{gh_{\sigma,c}^n}{\widehat{h}_{D_\sigma}^{n+1}} (h_L^n - h_K^n + z_L - z_K) \leq C_{\mathcal{M}}^{\delta t} g \frac{(C_{\mathcal{M}}^h)^2}{1 - \zeta} \left( \max_{K \in \mathcal{M}} (h_K^n) + \max_{K \in \mathcal{M}} (z_K) \right) \leq \frac{C_{\mathcal{M}}^{\delta t} (C_{\mathcal{M}}^h)^2 g}{1 - \zeta} (C_{\mathcal{M}}^h + \|z\|_\infty).$$

Therefore,

$$|\widehat{u}_{i,\sigma}^{n+1}| \leq C_{\mathcal{M}}^u + \frac{4C_{\mathcal{M}}^{\delta t} t (C_{\mathcal{M}}^u)^2 (C_{\mathcal{M}}^h)^2}{1 - \zeta} + \frac{C_{\mathcal{M}}^{\delta t} (C_{\mathcal{M}}^h)^2 g}{1 - \zeta} (C_{\mathcal{M}}^h + \|z\|_\infty),$$

which concludes the proof that (51b) holds.  $\square$

**Lemma 4.3** ( $L^1$  convergence of the intermediate step, Heun scheme). *Let  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$  be a sequence of meshes such that  $\delta t^{(m)}$  and  $\delta_{\mathcal{M}^{(m)}} \rightarrow 0$  as  $m \rightarrow +\infty$ ; assume that  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$  is uniformly regular, in the sense that there exists  $\theta > 0$  such that  $\theta_{\mathcal{M}^{(m)}} \leq \theta$  for any  $m \in \mathbb{N}$  (with  $\theta_{\mathcal{M}^{(m)}}$  defined by (5)).*

*Let  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  be a sequence of solutions to the scheme (6) satisfying (32) and (33) converging to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$ .*

*Assume that there exists  $\zeta \in (0, 1)$  such that the following restricted CFL-like condition holds:*

$$2 \delta t^{(m)} \leq \zeta \frac{|K|}{\sum_{\sigma \in \mathcal{E}(K)} |\sigma| |(\mathbf{u}^{(m)})_\sigma^n \cdot \mathbf{n}_{K,\sigma}|}, \quad \forall K \in \mathcal{M}^{(m)}, \forall m \in \mathbb{N}, \quad (52)$$

and assume that there exists  $C^{\delta t} \in \mathbb{R}_+^*$  not depending on  $m$  such that

$$\delta t^{(m)} \leq C^{\delta t} \min_{\sigma \in \mathcal{E}^{(m)}} |\sigma|, \quad \forall m \in \mathbb{N}. \quad (53)$$

Then there exists  $\widehat{C}^u, \widehat{C}^h \in \mathbb{R}_+^*$  such that

$$\frac{1}{\widehat{C}^h} < (\widetilde{h}^{(m)})_K^n \leq \widehat{C}^h, \quad \forall K \in \mathcal{M}^{(m)}, \quad (54a)$$

$$|(\widetilde{\mathbf{u}}^{(m)})_\sigma^n| \leq \widehat{C}^u, \quad \forall \sigma \in \mathcal{E}^{(m)}. \quad (54b)$$

Furthermore, the sequence  $(\widehat{h}^{(m)}, \widehat{\mathbf{u}}^{(m)})_{m \in \mathbb{N}}$  converges to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$ .

*Proof.* Under the above assumptions, the hypotheses (49) and (50) hold uniformly with respect to  $m$ , so that the bounds (54) are a direct consequence of Lemma 4.2.

Now, from equation (22a), we get that

$$\widehat{h}_K^{n+1} - h_K^n = -\frac{\delta t}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} (h_\sigma^n - h_K^n) - \frac{\delta t}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| h_K^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}, \quad \forall K \in \mathcal{M}^{(m)}.$$

For  $K \in \mathcal{M}^{(m)}$ , let us denote by  $\sigma_{K,i}$  and  $\sigma'_{K,i}$  the edges of  $K$  in the direction  $i \in \llbracket 1, 2 \rrbracket$ , so that  $K = \overrightarrow{[\sigma_{K,i}, \sigma'_{K,i}]}$  for  $i \in \llbracket 1, 2 \rrbracket$ ; noting that  $\mathbf{n}_{K,\sigma_{K,i}} = -\mathbf{n}_{K,\sigma'_{K,i}}$  and that  $|\sigma_{K,i}| = |\sigma'_{K,i}|$ , and owing to (32), we get that

$$\left| \sum_{\sigma \in \mathcal{E}(K)} |\sigma| h_K^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \right| \leq C^h \sum_{i=1}^d |\sigma_{K,i}| |u_{i,\sigma_{K,i}}^n - u_{i,\sigma'_{K,i}}^n|, \quad \forall K \in \mathcal{M}^{(m)}.$$

Since  $h_\sigma$  is a convex combination of  $h_K$  and  $h_L$ , with  $K$  and  $L$  such that  $\sigma = K|L$ , we get:

$$|\widehat{h}_K^{n+1} - h_K^n| \leq \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} \frac{\delta t^{(m)}}{|K|} |\sigma| |\mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}| |h_L^n - h_K^n| + C^h \sum_{i=1}^2 \delta t \frac{|\sigma|}{|K|} |u_{i,\sigma_{K,i}}^n - u_{i,\sigma'_{K,i}}^n|, \quad \forall K \in \mathcal{M}^{(m)}.$$

Noting that (53) implies that  $\frac{\delta t^{(m)}}{|K|} |\sigma| \leq 1$  and thanks to the condition (33), we thus get that there exists  $C \in \mathbb{R}_+$  depending on  $C^h$ ,  $C^u$ ,  $C^{\delta t}$  such that

$$|\widehat{h}_K^{n+1} - h_K^n| \leq C \left[ \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |h_L^n - h_K^n| + \sum_{i=1}^d |u_{i,\sigma_{K,i}}^n - u_{i,\sigma'_{K,i}}^n| \right], \quad \forall K \in \mathcal{M}^{(m)}.$$

Multiplying this latter inequality by  $|K| \delta t^{(m)}$  and summing over  $K \in \mathcal{M}^{(m)}$  and  $n \in \llbracket 0, N \rrbracket$ , using the uniform regularity of the mesh and owing again to the convergence result on the space translates given in Lemma B.2, we conclude that

$$\int_0^T \int_\Omega |\widehat{h}^{(m)} - h^{(m)}| \, d\mathbf{x} \, dt \rightarrow 0 \text{ as } m \rightarrow +\infty.$$

Let us now turn to the intermediate velocities. Owing to (22b), (23) and since  $\widehat{u}$  satisfies a dual mass balance of the form (21), we have for  $\sigma = K|L \in (\mathcal{E}_{\text{int}}^{(i)})^{(m)}$ ,  $i \in \llbracket 1, 2 \rrbracket$ :

$$\begin{aligned} \widehat{h}_{D_\sigma}^{n+1} (\widehat{u}_{i,\sigma}^{n+1} - u_{i,\sigma}^n) &= -(\widehat{h}_{D_\sigma}^{n+1} - h_{D_\sigma}^n) u_{i,\sigma}^n - \frac{\delta t}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} u_{i,\epsilon}^n - \delta t g h_{\sigma,c}^n ((\partial_\sigma h^n) + (\partial_\sigma z)) \\ &= - \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} \left[ \frac{\delta t |\epsilon|}{|D_\sigma|} \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) + \frac{\delta t |\sigma|}{|D_\sigma|} g h_{\sigma,c}^n (h_L^n - h_K^n + z_L - z_K) \right]. \end{aligned}$$

Hence, owing to (32), (33), (53) and to the fact that for  $\epsilon = \sigma|\sigma'$ ,  $u_{i,\epsilon}^n$  is a convex combination of  $u_{i,\sigma}^n$  and  $u_{i,\sigma'}^n$ , there exists  $C \in \mathbb{R}_+$  depending only on  $C^h$ ,  $C^u$ ,  $C^{\delta t}$  and  $g$  such that

$$|\widehat{u}_{i,\sigma}^{n+1} - u_{i,\sigma}^n| \leq C \left[ \sum_{\substack{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma) \\ \epsilon = \sigma|\sigma'}} |u_{i,\sigma'}^n - u_{i,\sigma}^n| + |h_L^n - h_K^n| + |z_L - z_K| \right], \quad \text{for } i = 1, 2.$$

Multiplying this latter inequality by  $|D_\sigma| \delta t^{(m)}$  and summing over  $\sigma \in \mathcal{M}^{(m)}$  and  $n \in \llbracket 0, N \rrbracket$ , using the uniform regularity of the mesh and again thanks to Lemma B.2 we conclude that

$$\int_0^T \int_\Omega |\widehat{u}_i^{(m)} - u_i^{(m)}| \, d\mathbf{x} \, dt \rightarrow 0 \text{ as } m \rightarrow +\infty, \text{ for } i = 1, 2.$$

□

## 5. WEAK ENTROPY CONSISTENCY OF THE FORWARD EULER- MAC SCHEME

**Theorem 5.1** (Weak entropy consistency of the forward Euler MAC scheme). *Let  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$  be a sequence of meshes such that  $\delta t^{(m)}$  and  $\delta_{\mathcal{M}^{(m)}} \rightarrow 0$  as  $m \rightarrow +\infty$ ; assume that there exists  $\theta > 0$  such that  $\theta_{\mathcal{M}^{(m)}} \leq \theta$  for any  $m \in \mathbb{N}$  (with  $\theta_{\mathcal{M}^{(m)}}$  defined by (5)). Let  $(h^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$  be a sequence of solutions to the*

scheme (6) converging to  $(\bar{h}, \bar{\mathbf{u}})$  in  $L^1(\Omega \times (0, T)) \times L^1(\Omega \times (0, T))^2$ , such that (32), (33) hold. Assume the following CFL-like condition:

$$\delta t^{(m)} \leq \frac{|D_\sigma| h_{D_\sigma}^{n+1}}{\sum_{\substack{\epsilon \in \mathcal{E} D_\sigma \\ \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma, \epsilon} > 0}} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma, \epsilon}}. \quad (55)$$

Assume furthermore that

$$\exists C_{BVt} \in \mathbb{R}_+ : \sum_{K \in \mathcal{M}^{(m)}} |K| |(h^{(m)})_K^{n+1} - (h^{(m)})_K^n| \sum_{K \in \mathcal{M}^{(m)}} |K| |(h^{(m)})_K^{n+1} - (h^{(m)})_K^n| \leq C_{BVt} \quad \forall m \in \mathbb{N}, \quad (56a)$$

$$\frac{\delta t^{(m)}}{\inf_{K \in \mathcal{M}^{(m)}} \text{diam}(K)} \sum_{n=0}^{N_m-1} \rightarrow 0 \text{ as } m \rightarrow +\infty, \quad (56b)$$

and that the coefficients  $\lambda_{K, \sigma}$  and  $\mu_{\sigma, \epsilon}$  in (9) and (19) satisfy:

$$\lambda_{K, \sigma} \in \left[\frac{1}{2}, 1\right] : \text{ if } \mathbf{F}_\sigma \cdot \mathbf{n}_{K, \sigma} \geq 0, \quad (57)$$

$$\mu_{\sigma, \epsilon} \in \left[\frac{1}{2}, 1\right] : \text{ if } \mathbf{F}_\sigma \cdot \mathbf{n}_{K, \sigma} \geq 0. \quad (58)$$

Then  $(\bar{h}, \bar{\mathbf{u}})$  satisfies the entropy inequality (31).

Note also that the condition (9) implies that

$$\forall K \in \mathcal{M}, \forall \sigma = K|L \in \mathcal{E}_{\text{int}}(K) \text{ with } K \text{ such that } \mathbf{F}_\sigma \cdot \mathbf{n}_{K, \sigma} \geq 0,$$

$$h_\sigma = \frac{1}{2}(h_K + h_L) + (\lambda_{K, \sigma} - \frac{1}{2})(h_K - h_L), \text{ with } \lambda_{K, \sigma} - \frac{1}{2} \geq 0. \quad (59)$$

Also Note that the condition (57) is rather restrictive. Indeed, it is satisfied by the usual two slopes minmod limiter [15] only in the case of a uniform Cartesian mesh [27], and it is not satisfied by the three slopes minmod limiter.

*Proof.* Let  $\varphi \in C_c^\infty(\Omega \times [0, T], \mathbb{R}_+)$ , and for a given discretization  $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})$  let  $\varphi_K^n$  (resp.  $\varphi_\sigma^n$ ) denote the mean value of  $\varphi$  on  $K \times (t_n, t_{n+1})$  (resp.  $D_\sigma \times (t_n, t_{n+1})$ ), for any  $K \in \mathcal{M}^{(m)}$  (resp.  $\sigma \in \mathcal{E}^{(m)}$ ) and  $n \in \llbracket 0, N_m - 1 \rrbracket$ . Let us multiply the discrete kinetic energy balance (27) by  $\delta t \varphi_\sigma^{n+1}$  and sum over  $\sigma \in \mathcal{E}^{(m)}$  and  $i \in \{1, 2\}$ ; let us then multiply the discrete potential energy balance (28) by  $\delta t |K| \varphi_K^n$  and sum over  $K \in \mathcal{M}^{(m)}$ . Summing the two resulting equations and summing over  $n \in \llbracket 0, N_m - 1 \rrbracket$ , we get, owing to lemmas 3.3 and 3.4,

$$\int_0^T \int_\Omega \mathfrak{C}_{\text{KIN}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt + \int_0^T \int_\Omega \mathfrak{C}_{\text{POT}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt + \mathcal{P}^{(m)} + \mathcal{Z}^{(m)} = -\mathcal{R}_k^{(m)} - \mathcal{R}_p^{(m)}, \quad (60)$$

with

$$\begin{aligned}
\mathfrak{C}_{\text{KIN}}^{(m)}(U^{(m)})|_{D_\sigma} &= \sum_{i=1}^2 \mathfrak{C}_{\text{KIN}}^{(m,i)}(U^{(m)})|_{D_\sigma} \text{ with } \mathfrak{C}_{\text{KIN}}^{(m,i)}(U^{(m)})|_{D_\sigma} = (\partial_t E_{k,i})_\sigma^{n+1} + \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} |\epsilon| \frac{(u_{i,\epsilon}^n)^2}{2} \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon}, \\
&\text{and } (E_{k,i})_\sigma^n = \frac{1}{2} h_{D_\sigma}^n (u_{i,\sigma}^n)^2, \\
\mathfrak{C}_{\text{POT}}^{(m)}(U^{(m)})|_K &= \frac{1}{2} (\partial_t h_K^2)^n + \text{div}_K \left( \frac{1}{2} g (h^n)^2 \mathbf{u}^n \right), \\
\mathfrak{P}^{(m)} &= \sum_{n=0}^{N_m-1} \delta t^{(m)} \left[ \sum_{\sigma \in \mathcal{E}^{(m)}} |D_\sigma| u_\sigma^{n+1} \partial_\sigma p^{n+1} \varphi_\sigma^{n+1} + \sum_{K \in \mathcal{M}^{(m)}} |K| p_K^n \text{div}_K \mathbf{u}^n \varphi_K^n \right], \\
\mathfrak{Z}^{(m)} &= + \sum_{n=0}^{N_m-1} \delta t^{(m)} \left[ \sum_{\sigma \in \mathcal{E}^{(m)}} |D_\sigma| h_{\sigma,c}^{n+1} u_\sigma^{n+1} \partial_\sigma z \varphi_\sigma^{n+1} + \sum_{K \in \mathcal{M}^{(m)}} g \left( z_K (\partial_t h_K)^n + g z_K \text{div}_K (h^n \mathbf{u}^n) \right) \varphi_K^n \right], \\
\mathfrak{R}_k^{(m)} &= \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i,m)}} \left[ \frac{1}{2} \frac{|D_\sigma|}{\delta t^{(m)}} h_{D_\sigma}^{n+1} (u_\sigma^{n+1} - u_\sigma^n)^2 \right. \\
&\quad \left. + \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} \left( -\frac{1}{2} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 + (u_{i,\epsilon}^n - u_{i,\sigma}^n) (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n) \right) \right] \varphi_\sigma^{n+1} \\
\mathfrak{R}_p^{(m)} &\geq \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{K \in \mathcal{M}^{(m)}} \left[ -\frac{1}{2} g \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| g (h_K^{n+1} - h_K^n) h_\sigma^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \right] \varphi_K.
\end{aligned}$$

*Kinetic energy convection term*

Let us check that the above defined convection operator  $\mathfrak{C}_{\text{KIN}}^{(m,i)}$  satisfies the hypotheses (80)–(82) of Lax-Wendroff type consistency Lemma A.1 given in the appendix which we apply here with  $d = 2$ ,  $\mathcal{P}^{(m)}$  and  $\mathfrak{F}^{(m)}$  the  $i$ -th dual mesh and its set of edges,  $U = (h, \mathbf{u})$ ,  $\beta(U) = E_{k,i}(U) = \frac{1}{2} h u_i^2$ , for  $i = 1, 2$ .

Let us start with the assumption (80). For a given function  $\psi \in L^1(\Omega)$ , and any subset  $A$  of  $\Omega$  we denote by  $\langle \psi \rangle_A$  the mean value of  $\psi$  on  $A$ . By definition of the kinetic energy, we have  $(E_{k,i})_\sigma^0 = \frac{1}{2} h_{D_\sigma}^0 |u_{i,\sigma}^0|^2 = \frac{1}{2} \langle h_0 \rangle_{D_\sigma} (\langle u_{i,0} \rangle_{D_\sigma})^2$  and  $E_{k,i}(U_0) = E_{k,i}(h_0, \mathbf{u}_0) = \frac{1}{2} h_0 u_{i,0}^2$ . Therefore, owing to the assumptions (32)–(33) on the functions  $h^{(m)}$  and  $\mathbf{u}^{(m)}$  and to the fact that these sequences converge in  $L^1$

$$\begin{aligned}
\sum_{P \in \mathcal{P}^{(m)}} \int_P |(\beta^{(m)})_P^0 - \beta(U_0(\mathbf{x}))| d\mathbf{x} &= \sum_{\sigma \in \mathcal{E}^{(m)}} \int_{D_\sigma} |(E_{k,i})_\sigma^0 - E_{k,i}(h_0, \mathbf{u}_0)| d\mathbf{x} \\
&= \frac{1}{2} \sum_{\sigma \in \mathcal{E}^{(m)}} |D_\sigma| \left| \langle h_0 \rangle_{D_\sigma} \langle u_{i,0} \rangle_{D_\sigma}^2 - \langle h_0 u_{i,0}^2 \rangle_{D_\sigma} \right| \\
&\rightarrow 0 \text{ as } m \rightarrow +\infty.
\end{aligned}$$

The assumption (80) is thus satisfied.

Let us then note that the assumption (81), which reads

$$\sum_{n=0}^{N_m-1} \sum_{\sigma \in \mathcal{E}^{(m)}} \int_{t_n}^{t_{n+1}} \int_{D_\sigma} |E_{k,i}^{(m)}|_\sigma^n - E_{k,i}(U^{(m)}(\mathbf{x}, t))| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

is satisfied, again thanks to the assumptions (32)–(33) on the functions  $h^{(m)}$  and  $\mathbf{u}^{(m)}$  and to the fact that these sequences converge in  $L^1$ .

Let us now turn to the assumption (82), which reads

$$\sum_{n=0}^{N_m-1} \sum_{\sigma \in \mathcal{E}(m)} \int_{t_n}^{t_{n+1}} \frac{\text{diam}(D_\sigma)}{|D_\sigma|} \int_{D_\sigma} \left| \sum_{\epsilon \in \tilde{\mathcal{E}}(m)} |\epsilon| \left( (\mathbf{F}^{(m)})_\epsilon^n - \mathbf{f}(U^m(\mathbf{x}, t)) \right) \cdot \mathbf{n}_{\sigma, \epsilon} \right| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

with  $(\mathbf{F}^{(m)})_\epsilon^n = \frac{1}{2}(u_{i, \epsilon^n})^2 \mathbf{F}_\epsilon^n$  and  $\mathbf{f}(U) = \frac{1}{2}h|u|^2 u_i$ . This assumption is indeed satisfied since  $\mathbf{F}_\epsilon^n$  is a convex combination of  $h_\sigma \mathbf{u}_\sigma$  and  $h_{\sigma'} \mathbf{u}_{\sigma'}$  for  $\epsilon = \sigma | \sigma'$ , and thanks to the boundedness and convergence assumptions on the sequences  $(h^{(m)})_{m \in \mathbb{N}}$  and  $(\mathbf{u}^{(m)})_{m \in \mathbb{N}}$ .

By Lemma A.1, we thus get that

$$\int_0^T \int_\Omega [C_{\text{KIN}}^{(m)}(U^{(m)})]_i \varphi(\mathbf{x}, t) d\mathbf{x} dt \rightarrow - \int_\Omega E_{k,i}(U^{(0)}) \varphi(x, 0) d\mathbf{x} - \int_0^T \int_\Omega E_{k,i}(\bar{U}) \partial_t \varphi + \frac{1}{2} E_{k,i}(\bar{U}) \bar{u}_i \partial_i \varphi d\mathbf{x} dt$$

with  $E_{k,i}(\bar{U}) = \frac{1}{2} g \bar{h} \bar{u}_i^2$ ; summing over  $i = 1, 2$ , we get that

$$\begin{aligned} \int_0^T \int_\Omega \mathcal{C}_{\text{KIN}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) d\mathbf{x} dt \rightarrow \\ - \int_\Omega E_k(U^{(0)}) \varphi(x, 0) d\mathbf{x} - \int_0^T \int_\Omega \left[ E_k(\bar{U}) \partial_t \varphi + \frac{1}{2} E_k(\bar{U}) \mathbf{u} \cdot \nabla \varphi \right] d\mathbf{x} dt \text{ as } m \rightarrow +\infty. \end{aligned} \quad (61)$$

with  $E_k(\bar{U}) = \frac{1}{2} g \bar{\mathbf{u}}^2$ .

*Potential energy convection terms*

Let us now check that the above defined convection operator  $\mathcal{C}_{\text{POT}}^{(m)}$  satisfies the hypotheses (80)–(82) of Lemma A.1 which we now apply with  $d = 2$ ,  $\mathcal{P}^{(m)}$  and  $\mathfrak{F}^{(m)}$  the primal mesh and its set of edges,  $U = (h, \mathbf{u})$ ,  $\beta(U) = \frac{1}{2} g h^2$  and  $\mathbf{f}(U) = \frac{1}{2} g h^2 \mathbf{u}$ .

Indeed,

$$\sum_{K \in \mathcal{M}} \int_K \left| \langle h(\cdot, 0) \rangle_K - h(x, 0) \right|^2 d\mathbf{x} \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

so that the hypothesis (80) is satisfied. Next,

$$\sum_{n=0}^{N_m-1} \int_{t_n}^{t_{n+1}} \sum_{K \in \mathcal{M}} \int_K \left| (h_K^n)^2 - h^2(\mathbf{x}, t) \right| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

thanks to the boundedness and convergence assumptions on the sequence  $(h^{(m)})_{m \in \mathbb{N}}$ . so that the hypothesis (81) is satisfied. Finally, the left hand side of (82) reads

$$\begin{aligned} X_F &= \sum_{n=0}^{N_m-1} \int_{t_n}^{t_{n+1}} \sum_{K \in \mathcal{M}} \frac{\text{diam}(K)}{|K|} \int_K \left| \sum_{\sigma \in \mathcal{E}(K)} |\sigma| \left( \frac{1}{2} g (h_\sigma^n)^2 \mathbf{u}_\sigma^n - \frac{1}{2} g h^2(\mathbf{x}, t) \mathbf{u}(\mathbf{x}, t) \right) \cdot \mathbf{n}_{K, \sigma} \right| d\mathbf{x} dt \\ &= \sum_{n=0}^{N_m-1} \int_{t_n}^{t_{n+1}} \sum_{K \in \mathcal{M}} \frac{\text{diam}(K)}{|K|} \left| \sum_{\sigma \in \mathcal{E}(K)} |\sigma| \int_{D_{K, \sigma}} \left( \frac{1}{2} g (h_\sigma^n)^2 - \frac{1}{2} g (h_K^n)^2 \right) \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K, \sigma} \right| d\mathbf{x} dt \\ &\rightarrow 0 \text{ as } m \rightarrow +\infty \end{aligned}$$

thanks to the fact that  $h_\sigma^n$  is a convex combination of  $h_K^n$  and  $h_L^n$  for  $\sigma = K | L$ , and thanks to the boundedness and convergence assumptions on the sequences  $(h^{(m)})_{m \in \mathbb{N}}$  and  $(\mathbf{u}^{(m)})_{m \in \mathbb{N}}$ . Therefore, the assumption (82) is also satisfied.

Hence by Lemma A.1,

$$\int_0^T \int_{\Omega} \mathcal{C}_{\text{POT}}^{(m)}(U^{(m)}) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \rightarrow -\frac{1}{2} \int_{\Omega} g h^2(\mathbf{x}, 0) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_{\Omega} \left[ E_p(\bar{U}) \partial_t \varphi + E_p(\bar{U}) \bar{\mathbf{u}} \cdot \nabla \varphi \right] \, d\mathbf{x} \, dt, \quad (62)$$

with  $E_p(\bar{U}) = \frac{1}{2} g \bar{h}^2$ .

*Pressure terms*

Let us rewrite  $\mathcal{P}^{(m)}$  as

$$\begin{aligned} \mathcal{P}^{(m)} &= \sum_{n=0}^{N_m-1} \delta t^{(m)} (A^{n+1} + B^{n+1}) - \delta t^{(m)} B^0, \\ \text{with } A^n &= \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}}^{(m)} \\ \sigma = K|L}} |D_{\sigma}| u_{\sigma}^n \partial_{\sigma} p^n \varphi_{\sigma}^n, \text{ and } B^n = \sum_{K \in \mathcal{M}^{(m)}} |K| p_K^n \operatorname{div}_K(\mathbf{u}^n) \varphi_K^n. \end{aligned}$$

By Lemma 5.2 below,

$$A^{n+1} + B^{n+1} = \sum_{K \in \mathcal{M}^{(m)}} \sum_{\sigma \in \mathcal{E}(K)} |D_{K,\sigma}| p_K^{n+1} \mathbf{u}_{\sigma}^{n+1} \cdot \frac{|\sigma|(\varphi_K^{n+1} - \varphi_{\sigma}^{n+1})}{|D_{K,\sigma}|} \mathbf{n}_{K,\sigma}$$

On each subcell  $D_{K,\sigma}$  the quantity  $\frac{|\sigma|(\varphi_K^{n+1} - \varphi_{\sigma}^{n+1})}{|D_{K,\sigma}|} \mathbf{n}_{K,\sigma}$  is, up to higher order terms, a discrete differential quotient of  $\varphi$  between  $\mathbf{x}_K$  and  $\mathbf{x}_{\sigma}$ , in the direction  $i$  if  $\sigma \in \mathcal{E}^{(i)}$ , which uniformly converges to  $\partial_i \varphi \mathbf{e}_i$  in the case of a rectangular grid, and therefore,

$$\sum_{n=0}^{N_m-1} \delta t^{(m)} (A^{n+1} + B^{n+1}) \rightarrow - \int_0^T \int_{\Omega} \bar{p}(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty.$$

Now, since we assume  $u_0 \in W^{1,1}(\Omega)$ ,

$$\begin{aligned} \delta t^{(m)} |B^0| &= \delta t^{(m)} \left| \sum_{K \in \mathcal{M}^{(m)}} |K| p_K^0 \operatorname{div}_K(\mathbf{u}^0) \varphi_K^0 \right| \\ &\leq g \delta t^{(m)} \|h_0\|_{\infty}^2 \|\varphi\|_{\infty} \sum_{K \in \mathcal{M}^{(m)}} |K| |\operatorname{div}_K(\mathbf{u}^0)| \\ &\leq 2g \frac{\delta t^{(m)}}{\inf_{K \in \mathcal{M}^{(m)}} \operatorname{diam}(K)} \|h_0\|_{\infty}^2 \|\varphi\|_{\infty} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| d_{\sigma} \|\mathbf{u}_0\|_{\infty}, \end{aligned}$$

so that, by the assumption (56b),  $\sum_{n=0}^{N_m-1} \delta t^{(m)} B^0 \rightarrow 0$  as  $m \rightarrow +\infty$ ; and therefore,

$$\mathcal{P}^{(m)} \rightarrow - \int_0^T \int_{\Omega} \bar{p}(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty. \quad (63)$$

In the above bound, we used the assumption (56b); this could be avoided if we assume  $u_0 \in W^{1,1}(\Omega)$  ) or  $u_0 \in L^1(0, T; BV(\Omega))$ ; indeed, in this case we have

$$|B^0| \leq g \|h_0\|_{\infty}^2 \|\varphi\|_{\infty} \|\mathbf{u}_0\|_{W^{1,1}(\Omega)}.$$

However, the assumption (56) seems unavoidable to deal with the remainder term appearing in the discrete potential energy, see below.

*Bathymetry terms*

Let us introduce the following piecewise constant functions:

- $\tilde{h}^{(m)}$  is the piecewise constant function equal to  $h_{\sigma,c}^{n+1} = \frac{1}{2}(h_K^{n+1} + h_L^{n+1})$  on each set  $D_\sigma \times (t_n, t_{n+1})$ , for  $\sigma = K|L \in \mathcal{E}_{\text{int}}^{(m)}$  and  $n \in \llbracket 0, N_m - 1 \rrbracket$ ;
- $\nabla^{(m)} z^{(m)}$  is the piecewise constant function equal to  $\frac{|\sigma|}{|D_\sigma|}(z_L - z_K)$  on each set  $D_\sigma$ , for  $\sigma = K|L \in \mathcal{E}_{\text{int}}^{(m)}$ ;
- $\tilde{\varphi}^{(m)}$  is the piecewise constant function equal to  $\varphi_\sigma$  on each set on each set  $D_\sigma \times (t_n, t_{n+1})$ , for  $\sigma = K|L \in \mathcal{E}_{\text{int}}^{(m)}$  and  $n \in \llbracket 0, N_m - 1 \rrbracket$ ;

With these notations, we get that

$$\begin{aligned} \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}^{(m)}} |D_\sigma| h_{\sigma,c}^{n+1} u_\sigma^{n+1} \partial_\sigma z \varphi_\sigma^{n+1} &= \int_\Omega \tilde{h}^{(m)}(\mathbf{x}, t) \mathbf{u}^{(m)}(\mathbf{x}, t) \cdot \nabla z^{(m)}(\mathbf{x}) \tilde{\varphi}^{(m)}(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } n \rightarrow +\infty \\ &\rightarrow \int_\Omega h(\mathbf{x}, t) \mathbf{u}(\mathbf{x}, t) \cdot \nabla z(\mathbf{x}) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty, \end{aligned} \quad (64)$$

thanks to the convergence assumptions on  $h^{(m)}$  and  $\mathbf{u}^{(m)}$  and owing to the strong convergence of the discrete gradient  $\nabla^{(m)}$  (which would be only a weak convergence in the case of a non rectangular mesh, see [12, Lemma 3.1]).

Now let  $T_K^n = g \partial_t h_K^{n+1} z_K$  and  $Z_K^n = \frac{1}{|K|} g \sum_{\sigma \in \mathcal{E}(K)} |\sigma| h_\sigma^n u_{K,\sigma}^n z_K$ . Using a discrete summation by parts in

time and thanks to the convergence assumption on  $h^{(m)}$ , we get that

$$\sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{K \in \mathcal{M}^{(m)}} |K| T_K^n \rightarrow - \int_\Omega g z(\mathbf{x}) h(\mathbf{x}, 0) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_\Omega g z(\mathbf{x}) h(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt. \quad (65)$$

Using next a discrete summation by parts in space, we get

$$\begin{aligned} \sum_{K \in \mathcal{M}^{(m)}} |K| Z_K^n &= \sum_{K \in \mathcal{M}^{(m)}} g z_K \varphi_K^{n+1} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| h_\sigma^n u_{K,\sigma}^n z_K \\ &= \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}}^{(m)} \\ \sigma = K|L}} |\sigma| h_\sigma^n u_{K,\sigma}^n (z_K \varphi_K^{n+1} - z_L \varphi_L^{n+1}) \\ &= - \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}}^{(m)} \\ \sigma = K|L}} |D_\sigma| h_\sigma^n \mathbf{u}_\sigma^n \cdot (\nabla^{(m)}(z\varphi))_\sigma^{n+1}, \end{aligned}$$

where  $\nabla^{(m)}(z\varphi)$  is the piecewise constant discrete gradient defined by:

$$\begin{aligned} \forall \sigma = K|L \in \mathcal{E}_{\text{int}}^{(m)}, \forall n \in \llbracket 0, N_m - 1 \rrbracket, \forall (\mathbf{x}, t) \in D_\sigma \times [t_n, t_{n+1}), \\ \nabla^{(m)}(z\varphi)^{n+1}(\mathbf{x}, t) = (\nabla^{(m)}(z\varphi))_\sigma =^{n+1} \frac{|\sigma|}{|D_\sigma|} (z_K \varphi_K^{n+1} - z_L \varphi_L^{n+1}) \mathbf{n}_{K,\sigma}, \end{aligned}$$

which converges to  $\nabla(z\varphi)$  uniformly in the case of a rectangular mesh, and weakly in the case of a general mesh, see [12, Lemma 3.1].

Therefore, thanks to the convergence assumptions on  $h$  and  $\mathbf{u}$ ,

$$\sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{K \in \mathcal{M}^{(m)}} |K| Z_K^n \rightarrow - \int_0^T \int_\Omega g \bar{h}(\mathbf{x}, t) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla(z\varphi)(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty. \quad (66)$$



Owing to (64), (65) and (66), we thus get that

$$\begin{aligned} \mathcal{Z}^{(m)} \rightarrow & - \int_{\Omega} g z(\mathbf{x}) h(\mathbf{x}, 0) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_{\Omega} g z(\mathbf{x}) \bar{h}(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ & - \int_0^T \int_{\Omega} g \bar{h}(\mathbf{x}, t) z(\mathbf{x}) \bar{\mathbf{u}}(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty. \end{aligned} \quad (67)$$

*Remainder terms*

The remainder term  $\mathcal{R}_k^{(m)}$  in (60) satisfies

$$\mathcal{R}_k^{(m)} = \mathcal{R}_{k,1}^{(m)} + \mathcal{R}_{k,2}^{(m)} + \mathcal{R}_{k,3}^{(m)} \quad (68)$$

with

$$\begin{aligned} \mathcal{R}_{k,1}^{(m)} &= \frac{1}{2} \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \frac{1}{\delta t} |D_{\sigma}| h_{D_{\sigma}}^{n+1} (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n)^2 \varphi_{\sigma}^{n+1}, \\ \mathcal{R}_{k,2}^{(m)} &= -\frac{1}{2} \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_{\sigma})} |\epsilon| \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \varphi_{\sigma}^{n+1} \\ \mathcal{R}_{k,3}^{(m)} &= \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_{\sigma})} |\epsilon| \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n) (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n) \varphi_{\sigma}^{n+1}. \end{aligned}$$

The term  $\mathcal{R}_{k,3}^{(m)}$  satisfies

$$\mathcal{R}_{k,3}^{(m)} \geq \mathcal{R}_{k,3,1}^{(m)} + \mathcal{R}_{k,3,2}^{(m)} \quad (69)$$

with

$$\begin{aligned} \mathcal{R}_{k,3,1}^{(m)} &= -\frac{1}{2} \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_{\sigma}) \\ \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} > 0}} |\epsilon| \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\sigma}^{n+1} - u_{i,\sigma}^n)^2 \varphi_{\sigma}^{n+1}, \\ \mathcal{R}_{k,3,2}^{(m)} &= -\frac{1}{2} \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_{\sigma})} |\epsilon| \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \varphi_{\sigma}^{n+1}. \end{aligned}$$

Thanks to the CFL condition (55), we get that

$$\mathcal{R}_{k,1}^{(m)} + \mathcal{R}_{k,3,1}^{(m)} \geq 0. \quad (70)$$

Let us now study the term

$$\tilde{\mathcal{R}}_{k,2}^{(m)} = \mathcal{R}_{k,3,2}^{(m)} + \mathcal{R}_{k,2}^{(m)} = - \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_{\sigma})} |\epsilon| \mathbf{F}_{\epsilon}^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \varphi_{\sigma}^{n+1},$$

which we decompose as:  $\tilde{\mathcal{R}}_{k,2}^{(m)} \geq \tilde{\mathcal{R}}_{k,2,1}^{(m)} + \tilde{\mathcal{R}}_{k,2,2}^{(m)}$ , with

$$\begin{aligned}\tilde{\mathcal{R}}_{k,2,1}^{(m)} &= - \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma) \\ \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} > 0}} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \varphi_\epsilon^{n+1}, \\ \tilde{\mathcal{R}}_{k,2,2}^{(m)} &= - \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 (\varphi_\epsilon^{n+1} - \varphi_\sigma^{n+1}),\end{aligned}$$

and, by conservativity,

$$\begin{aligned}\tilde{\mathcal{R}}_{k,2,1}^{(m)} &\geq \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\substack{\epsilon=\sigma|\sigma' \in \tilde{\mathcal{E}}_{\text{int}}^{(i)} \\ \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} > 0}} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} \left[ (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 - (u_{i,\epsilon}^n - u_{i,\sigma'}^n)^2 \right] \varphi_\epsilon^{n+1} \\ &\geq \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \|u_i\|_\infty \sum_{\substack{\epsilon=\sigma|\sigma' \in \tilde{\mathcal{E}}_{\text{int}}^{(i)} \\ \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} > 0}} |D_\epsilon| (2u_{i,\epsilon} - u_{i,\sigma} - u_{i,\sigma'}) (u_{i,\sigma'} - u_{i,\sigma}).\end{aligned}$$

Therefore, thanks to (19) and (58),

$$\tilde{\mathcal{R}}_{k,2,1}^{(m)} \geq \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\substack{\epsilon=\sigma|\sigma' \in \tilde{\mathcal{E}}_{\text{int}}^{(i)} \\ \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} > 0}} |\epsilon| \mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon} (2\mu_{\sigma,\epsilon} - 1) (u_{i,\sigma}^n - u_{i,\sigma'}^n)^2 \varphi_\epsilon^{n+1} \geq 0. \quad (71)$$

Let us then write that, thanks to the regularity of  $\varphi$ ,

$$\begin{aligned}|\tilde{\mathcal{R}}_{k,2,2}^{(m)}| &\leq C_\varphi \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |D_\epsilon| |\mathbf{F}_\epsilon^n \cdot \mathbf{n}_{\sigma,\epsilon}| (u_{i,\epsilon}^n - u_{i,\sigma}^n)^2 \\ &\leq C_\varphi \|h\|_\infty \|\mathbf{u}\|_\infty \sum_{i=1}^2 \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\sigma \in \mathcal{E}_{\text{int}}^{(i)}} \sum_{\epsilon \in \tilde{\mathcal{E}}^{(i)}(D_\sigma)} |D_\epsilon| |u_{i,\epsilon}^n - u_{i,\sigma}^n|\end{aligned}$$

so that, thanks to the  $L^1$  convergence of  $\mathbf{u}^{(m)}$  and to the regularity of the mesh, we may again apply Lemma B.2 to obtain

$$|\tilde{\mathcal{R}}_{k,2,2}^{(m)}| \rightarrow 0 \text{ as } m \rightarrow +\infty \quad (72)$$

Therefore, owing to (68)-(72)

$$\lim_{m \rightarrow +\infty} \mathcal{R}_k^{(m)} \geq 0. \quad (73)$$

Let us now turn to the remainder  $\mathcal{R}_p^{(m)}$ . We have  $\mathcal{R}_p^{(m)} \geq \mathcal{R}_{p,1}^{(m)} + \mathcal{R}_{p,2}^{(m)}$ , with

$$\begin{aligned}\mathcal{R}_{p,1}^{(m)} &= - \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{K \in \mathcal{M}^{(m)}} g \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_K^n, \\ \mathcal{R}_{p,2}^{(m)} &= \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{K \in \mathcal{M}^{(m)}} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| g (h_K^{n+1} - h_K^n) h_\sigma^n \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_K^n.\end{aligned}$$

Note that if  $h_\sigma^n$  is the upwind choice for any  $\sigma \in \mathcal{E}^{(m)}$ , then  $\mathcal{R}_{p,1}^{(m)} \geq 0$ . In the general case, we may write that

$$\mathcal{R}_{p,1}^{(m)} = \mathcal{R}_{p,1,1}^{(m)} + \mathcal{R}_{p,1,2}^{(m)}$$

with

$$\begin{aligned} \mathcal{R}_{p,1,1}^{(m)} &= - \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{K \in \mathcal{M}^{(m)}} g \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_\sigma^n \\ \mathcal{R}_{p,1,2}^{(m)} &= - \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{K \in \mathcal{M}^{(m)}} g \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (h_\sigma^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} (\varphi_K^n - \varphi_\sigma^n) \end{aligned}$$

By conservativity,

$$\begin{aligned} \mathcal{R}_{p,1,1}^{(m)} &\geq -g \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{\substack{\sigma=K|L \in \mathcal{E} \\ \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} > 0}} |\sigma| \left[ (h_\sigma^n - h_K^n)^2 - (h_\sigma^n - h_L^n)^2 \right] \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_\sigma^n \\ &= -g \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{\substack{\sigma=K|L \in \mathcal{E} \\ \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} > 0}} |\sigma| (h_L^n - h_K^n)(2h_\sigma^n - h_K^n - h_L^n) \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_\sigma^n \end{aligned}$$

Owing to the assumption (9), one has

$$(h_L^n - h_K^n)(2h_\sigma^n - h_K^n - h_L^n) = -2\lambda_{K,\sigma}(h_K^n - h_L^n)$$

and since by (57),  $\lambda_{K,\sigma} \geq \frac{1}{2}$  if  $\mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} > 0$ ,

$$\mathcal{R}_{p,1,1}^{(m)} = 2g \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{\substack{\sigma=K|L \in \mathcal{E} \\ \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} > 0}} |\sigma| \left( \lambda_{K,\sigma} - \frac{1}{2} \right) (h_L^n - h_K^n)^2 \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma} \varphi_\sigma^n \geq 0.$$

Now

$$|\mathcal{R}_{p,1,2}^{(m)}| \leq \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{2} \sum_{K \in \mathcal{M}^{(m)}} g \|h\|_\infty \|\mathbf{u}\|_\infty C_\varphi \sum_{\sigma \in \mathcal{E}(K)} |D_{K,\sigma}| |h_L^n - h_K^n| \rightarrow 0 \text{ as } m \rightarrow +\infty$$

so that

$$\lim_{m \rightarrow +\infty} \mathcal{R}_{p,1}^{(m)} \geq 0.$$

Let us now turn to  $\mathcal{R}_{p,2}^{(m)}$ . Since for all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}(K)$  we have  $|\sigma| \leq \frac{|K|}{\inf_{K \in \mathcal{M}} \text{diam}(K)}$ , we have

$$\begin{aligned} |\mathcal{R}_{p,2}^{(m)}| &\leq g \|h\|_\infty \|\mathbf{u}\|_\infty \|\varphi\|_\infty \sum_{n=0}^{N_m-1} \frac{\delta t^{(m)}}{\inf_{K \in \mathcal{M}^{(m)}} \text{diam}(K)} \sum_{n=0}^{N_m-1} \sum_{K \in \mathcal{M}^{(m)}} |K| |(h^{(m)})_K^{n+1} - (h^{(m)})_K^n| \\ &\rightarrow 0 \text{ as } m \rightarrow +\infty, \end{aligned}$$

thanks to the assumption (56). Hence

$$\lim_{m \rightarrow +\infty} \mathcal{R}_p^{(m)} \geq 0. \quad (74)$$

*Conclusion of the proof*

Owing to (73) and (74), passing to the limit in (60) as  $m \rightarrow +\infty$  yields, together with (61), (62), (63) and (67), that the limit  $(\bar{h}, \bar{u})$  satisfies the weak entropy inequality (31).  $\square$

The next lemma, used to pass to the limit in the pressure terms of the entropy is the discrete equivalent, on a staggered grid, of the formal equality  $\int_{\Omega} (\mathbf{u} \cdot \nabla p \varphi + p \operatorname{div} \mathbf{u} \varphi) \, d\mathbf{x} = - \int_{\Omega} p \mathbf{u} \cdot \nabla \varphi \, d\mathbf{x}$ .

**Lemma 5.2** (Pressure terms). *Let  $(\mathcal{M}, \mathcal{E})$  be a MAC discretization of  $\Omega$  in the sense of Definition 2.1 ; Let  $(p_K)_{K \in \mathcal{M}} \subset \mathbb{R}$  and  $(\mathbf{u}_\sigma)_{\sigma \in \mathcal{E}} \subset \mathbb{R}^d$  be some discrete unknowns associated to  $\mathcal{M}$  and  $\mathcal{E}$  respectively. Let  $\varphi \in C_c^\infty(\Omega)$ , and let  $\varphi_K$  (resp.  $\varphi_\sigma$ ) denote the mean value of  $\varphi$  on  $K$  (resp.  $D_\sigma$ ), for any  $K \in \mathcal{M}$  (resp.  $\sigma \in \mathcal{E}^{(m)}$ ). Then*

$$\sum_{\substack{\sigma \in \mathcal{E}^{\text{int}} \\ \sigma = K|L}} |D_\sigma| u_\sigma \partial_\sigma p \varphi_\sigma + \sum_{K \in \mathcal{M}} |K| p_K \operatorname{div}_K \mathbf{u} \varphi_K = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} |D_{K,\sigma}| p_K \mathbf{u}_\sigma \cdot \frac{|\sigma|(\varphi_K - \varphi_\sigma)}{|D_{K,\sigma}|} \mathbf{n}_{K,\sigma}.$$

*Proof.* Let us denote by  $A$  and  $B$  the first and second terms of the right hand side. Then, with the notations of Definition 2.1,

$$\begin{aligned} A &= \sum_{K \in \mathcal{M}^{(m)}} \sum_{\sigma \in \mathcal{E}(K)} |D_{K,\sigma}| u_\sigma \partial_\sigma p \varphi_\sigma = \sum_{K \in \mathcal{M}^{(m)}} \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |D_{K,\sigma}| u_{K,\sigma} \frac{p_L - p_K}{|D_\sigma|} |\sigma| \varphi_\sigma \\ &= \sum_{K \in \mathcal{M}^{(m)}} \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |D_{K,\sigma}| u_{K,\sigma} \frac{p_\sigma - p_K}{|D_{K,\sigma}|} |\sigma| \varphi_\sigma, \end{aligned}$$

where  $p_\sigma$  is defined by  $\frac{p_\sigma - p_K}{|D_{K,\sigma}|} = \frac{p_L - p_K}{|D_\sigma|}$ . By conservativity,  $\sum_{K \in \mathcal{M}^{(m)}} \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} u_{K,\sigma} p_\sigma |\sigma| \varphi_\sigma = 0$ , so that

$$A = - \sum_{K \in \mathcal{M}^{(m)}} \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |D_{K,\sigma}| u_{K,\sigma} \frac{p_K}{|D_{K,\sigma}|} |\sigma| \varphi_\sigma$$

Now

$$B = \sum_{K \in \mathcal{M}^{(m)}} |K| p_K \operatorname{div}_K \mathbf{u} \varphi_K = \sum_{K \in \mathcal{M}^{(m)}} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| p_K |D_{K,\sigma}| u_\sigma \varphi_K.$$

Adding the results for  $A$  and  $B$  concludes the proof.  $\square$

## 6. NUMERICAL RESULTS

This section is devoted to numerical tests: we first check the order of convergence of the proposed scheme on a two-dimensional regular solution (Section 6.1); then we turn to one-dimensional and two-dimensional shock solutions on a plane topography (Sections 6.2 and 6.3); in Section 6.4, we address a two-dimensional dam-break problem in a closed computational domain with a variable topography, which, in particular, shows the ability of staggered scheme to "natively" cope with reflection boundary conditions; finally, we compute the motion of a liquid slug over a partly dry support (6.5).

In this section, we compare three schemes: the second-order scheme developed here, the scheme referred to in Section 2.2 as the segregated forward Euler scheme (combining a segregated forward Euler scheme in time and the proposed MUSCL-like discretization of the convection fluxes) and a first order scheme which still features the segregated forward Euler scheme in time but with first-order upwind convection fluxes. These schemes are referred to in the following as the *second-order*, *segregated* and *first-order* scheme respectively.

The schemes have been implemented within the CALIF<sup>3S</sup> open-source software [7] of the French Institut de Sûreté et de Radioprotection Nucléaire (IRSN); this software is used for the following tests.

### 6.1. A smooth solution

We begin here by checking the accuracy of the scheme on a known regular solution consisting in a travelling vortex. This solution is obtained through the following steps: we first derive a compact-support  $H^2$  solution consisting in a standing vortex which becomes time-dependent by adding a constant velocity motion. The velocity field of the standing vortex and the pressure are sought under the form:

$$\hat{\mathbf{u}} = f(\xi) \begin{bmatrix} -x_2 \\ x_1 \end{bmatrix}, \quad \hat{p} = \wp(\xi),$$

with  $\xi = x_1^2 + x_2^2$ . A simple derivation of these expressions yields:

$$\hat{\mathbf{u}} \cdot \nabla \hat{\mathbf{u}} = -f(\xi)^2 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

and

$$\nabla \hat{p} = 2 \wp'(\xi) \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

Using the relation  $p = \frac{1}{2}gh^2$ , we thus obtain a stationary solution of the shallow water equations (1) with a topography  $z = 0$  if  $\wp$  satisfies  $8g\wp = (F + c)^2$ , where  $F$  is such that  $F' = f^2$ ,  $F(0) = 0$  and  $c$  is a positive real number. For the present numerical study, we choose  $f(\xi) = 10\xi^2(1 - \xi)^2$  if  $\xi \in (0, 1)$ ,  $f = 0$  otherwise, which indeed yields an  $H^2(\mathbb{R}^2)$  velocity field (note that as a consequence, the pressure and the water height are also regular), and  $c = 1$ . The problem is made unsteady by a time translation: given a constant vector field  $\mathbf{a}$ , the pressure  $p$  and the velocity  $\mathbf{u}$  are deduced from the steady state solution  $\hat{p}$  and  $\hat{\mathbf{u}}$ :

$$h(\mathbf{x}, t) = \hat{h}(\mathbf{x} - \mathbf{a}t), \quad \mathbf{u}(\mathbf{x}, t) = \hat{\mathbf{u}}(\mathbf{x} - \mathbf{a}t) + \mathbf{a}.$$

The center of the vortex is initially located at  $\mathbf{x}_0 = (0, 0)^t$ , the translation velocity  $\mathbf{a}$  is set to  $\mathbf{a} = (1, 1)^t$ , the computational domain is  $\Omega = (-1.2, 2.)^2$  and the computation is run on the time interval  $(0, 0.8)$ .

Computations are performed with successively refined meshes with square cells, and the time step is  $\delta t = \delta_{\mathcal{M}}/8$ , and corresponds to a Courant (or CFL) number with respect to the celerity of the fastest waves close to  $1/3$ . The discrete  $L^1$ -norm of the difference between the exact solution and the solution obtained by the second-order scheme is given in Table 1. The observed order of convergence over the whole sequence is 2 for the water height and 1.5 for the velocity. Results with the first-order scheme are given in Table 2; one observes that the second-order scheme is much more accurate. Finally, the segregated scheme yields good results on coarse meshes (it is the most accurate scheme on the  $32 \times 32$  mesh); unfortunately, when refining the mesh, oscillations appear, and the convergence is lost. This results confirms a behaviour already observed for the transport operator in [27]: for multi-dimensional problems, the smoothing produced by the Heun time-stepping seems to be necessary to compensate the oscillatory character of the MUSCL scheme (which, for the transport operator, does not lead, of course, to violate the local maximum principle warranted by construction of the limitation process).

mesh	error( $h$ )	ord( $h$ )	error( $u$ )	ord( $u$ )
$32 \times 32$	$3.61 \cdot 10^{-3}$	/	$2.93 \cdot 10^{-1}$	/
$64 \times 64$	$1.15 \cdot 10^{-3}$	1.65	$1.14 \cdot 10^{-1}$	1.36
$128 \times 128$	$2.58 \cdot 10^{-4}$	2.16	$4.06 \cdot 10^{-2}$	1.49
$256 \times 256$	$5.85 \cdot 10^{-5}$	2.14	$1.49 \cdot 10^{-2}$	1.45
$512 \times 512$	$1.53 \cdot 10^{-5}$	1.93	$4.67 \cdot 10^{-3}$	1.68

TABLE 1. Measured numerical errors for the travelling vortex – Discrete  $L^1$ -norm of the difference between the numerical and exact solution at  $t = 0.8$ , for the height and the velocity, and corresponding order of convergence.

mesh	error( $h$ )	ord( $h$ )	error( $u$ )	ord( $u$ )
$32 \times 32$	$8.04 \cdot 10^{-3}$	/	$6.55 \cdot 10^{-1}$	/
$64 \times 64$	$5.56 \cdot 10^{-3}$	0.53	$4.84 \cdot 10^{-1}$	0.44
$128 \times 128$	$3.53 \cdot 10^{-3}$	0.66	$3.22 \cdot 10^{-1}$	0.59
$256 \times 256$	$2.08 \cdot 10^{-3}$	0.76	$1.96 \cdot 10^{-1}$	0.72
$512 \times 512$	$1.15 \cdot 10^{-3}$	0.85	$1.16 \cdot 10^{-1}$	0.76

TABLE 2. Measured numerical errors for the travelling vortex with the first order scheme - Discrete  $L^1$ -norm of the difference between the numerical and exact solution at  $t = 0.8$ , for the height and the velocity, and corresponding order of convergence.

mesh	error( $h$ )	error( $u$ )
$32 \times 32$	$2.06 \cdot 10^{-3}$	$2.33 \cdot 10^{-1}$
$64 \times 64$	$1.37 \cdot 10^{-3}$	$1.18 \cdot 10^{-1}$
$128 \times 128$	$1.24 \cdot 10^{-3}$	$8.50 \cdot 10^{-2}$
$256 \times 256$	$1.26 \cdot 10^{-3}$	$6.16 \cdot 10^{-2}$
$512 \times 512$	$1.56 \cdot 10^{-3}$	$4.85 \cdot 10^{-2}$

TABLE 3. Measured numerical errors for the travelling vortex with the segregated scheme - Discrete  $L^1$ -norm of the difference between the numerical and exact solution at  $t = 0.8$ , for the height and the velocity.

## 6.2. A Riemann problem

We now turn to a one-dimensional shock solution, corresponding to a Riemann problem posed over  $\Omega = (0, 1)$ . The initial height is  $h = 1$  if  $x < 0.5$  and  $h = 0.2$  otherwise, and the topography  $z$  is set to zero over the computational domain; the fluid is initially at rest. The solution consists in a 1-rarefaction wave and a 2-shock.

We plot on Figure 3 and Figure 4 the results obtained at  $t = 0.1$  with the second-order scheme, the segregated scheme and the first-order scheme. The space step is  $\delta x = 1/200$  and the time step is chosen as  $\delta t = \delta x/10$ , which corresponds to a CFL number lower than 0.5 with respect to the waves celerity (the maximal speed of sound is close to 3 and the maximal velocity is close to 2). As expected, the first order scheme is more diffusive than the other ones. As in the previous test, the segregated forward Euler scheme (with MUSCL fluxes) exhibits some oscillations, which are damped by the Heun time discretization (see the Figure 4). In this test case, for both the second-order and the segregated scheme, the shock is captured with only one intermediate cell between the left and the right state.

## 6.3. A circular dam break problem

The objective of this test-case is to check the capability of the scheme to capture a multi-dimensional shock solution. The fluid is initially at rest and the height is given by:

$$h = 2.5 \text{ if } r < 2.5, \quad h = 0.5 \text{ otherwise, with } r^2 = x_1^2 + x_2^2.$$

The computational domain is  $\Omega = (-20, 20) \times (-20, 20)$  and the final time is  $T = 4.7$ .

We plot on Figure 5 the results obtained with a  $800 \times 800$  uniform mesh, with the second-order scheme. The time-step is  $\delta t = h_M/10$  (with a maximal velocity in the range of 3.5 and a maximal speed of sound in the range of 5). In addition, to cure some oscillations (see Figure 7), we add a slight stabilization in the momentum balance equation which consists in adding to the discrete momentum equation associated to an

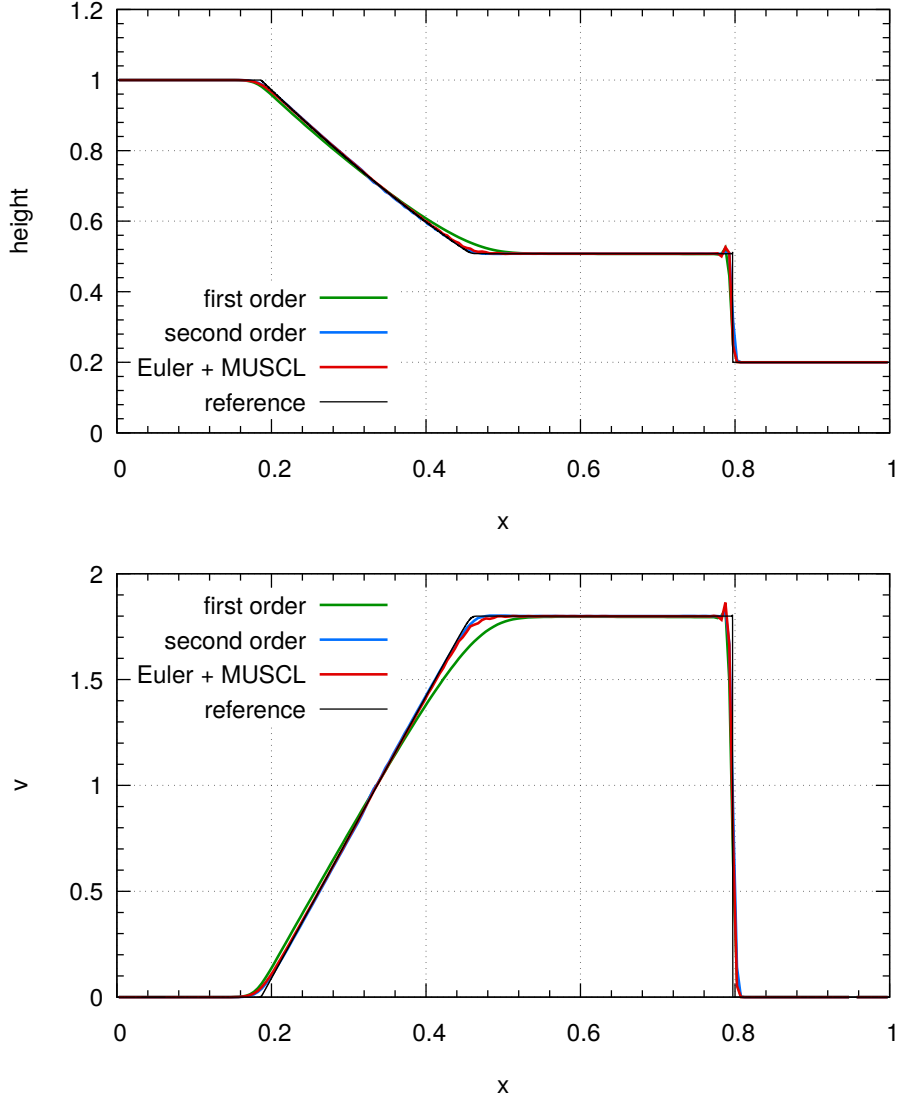


FIGURE 3. Riemann problem. Top: flow height – Bottom: velocity.

edge  $\sigma$  included in a cell  $K$  the following flux through a dual edge  $\epsilon = D_\sigma |D_{\sigma'}$ :

$$F_{\text{stab},\sigma,\epsilon} = \zeta h_K \text{diam}(K)^{d-1} (u_\sigma - u'_\sigma),$$

where  $\zeta$  is a user-defined parameter. Here,  $\zeta = 0.1$ , which is significantly lower than the diffusion generated by the use of an upwind scheme in the momentum balance equation; indeed, the upwind scheme may be seen as the centered one complemented by a diffusion taking the same expression as  $F_{\text{stab},\sigma,\epsilon}$  with  $\zeta h_K$  replaced by  $|F_{\sigma,\epsilon}|/2$ . The interest of this stabilization stems from the fact that the numerical diffusion introduced in the present family of schemes depends on the material velocity (and not on the waves celerity as, for instance, in collocated schemes based on Riemann solvers), and is sometimes too low in the zones where the fluid is almost at rest [22]. Note that, as a counterpart, the scheme does not become overdiffusive for low-Mach number flows. For the same computation, we give on Figure 6 the height and the radial velocity along the axis  $x_2 = 0$  (*i.e.* the first component of the velocity) at different times.

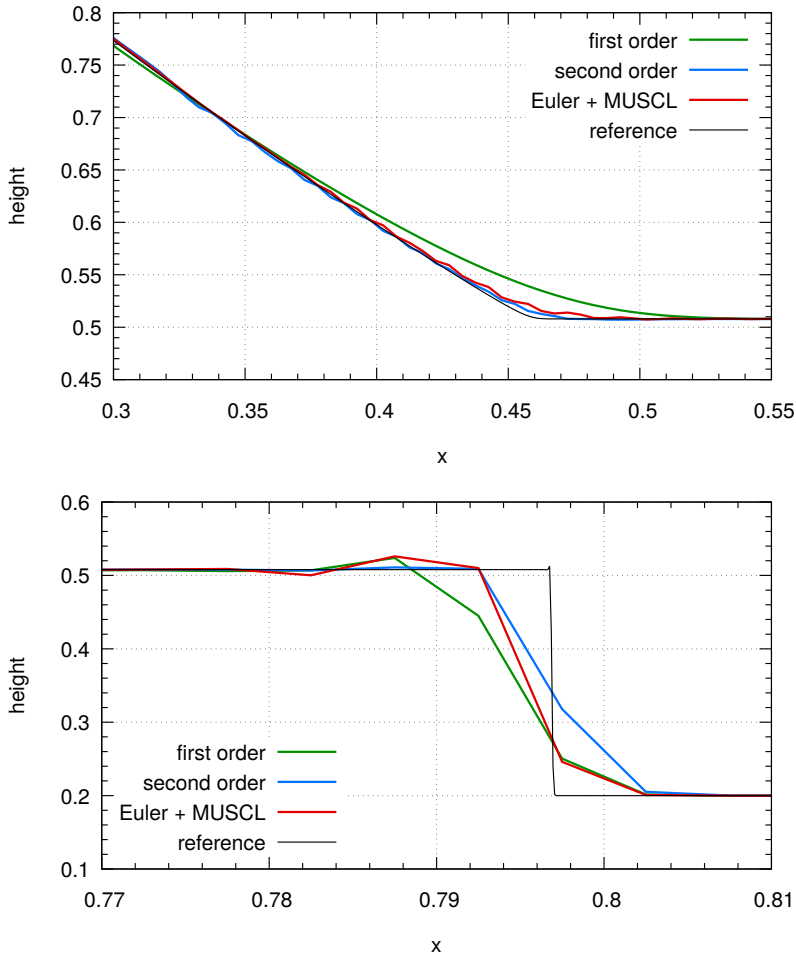


FIGURE 4. Riemann problem. Details of the flow height.

This computation is also used as "reference computation" on Figure 7, where we compare the results obtained at  $t = 3T/5$  with a  $200 \times 200$  mesh with the second-order scheme, the second-order scheme with stabilization and the first-order scheme. This latter is significantly more diffusive, and we observe how the stabilization (even if added to the momentum balance only and not on the mass balance) damps the oscillations obtained with the second-order scheme for both the flow height and the velocity.

#### 6.4. A so-called partial dam-break problem

We now turn to a test consisting in a partial dam-break problem with reflection phenomena, and with a non-flat bathymetry. In this test, the computational domain is  $\Omega = (0, 200) \times (0, 200) \setminus \Omega_w$  with  $\Omega_w = (95, 105) \times (0, 95) \cup (95, 105) \times (170, 200)$ . The fluid is supposed to be initially at rest, the initial water height is  $h = 10$  for  $x_1 \leq 100$  and  $h = 5 - 0.04(x_1 - 100)$  otherwise, and the bathymetry is  $z = 0$  if  $x_1 \leq 100$  and  $z = 0.04(x_1 - 100)$  otherwise. A zero normal velocity is prescribed at all the boundaries of the computational domain. The computation is performed with a mesh obtained from a  $1000 \times 1000$  regular grid by removing the cells included in  $\Omega_w$ . The time step is  $\delta t = \delta_{\mathcal{M}}/40$  (the maximal speed of sound and the maximal velocity are both close to 10). A stabilization with  $\zeta = 0.25$  (so two orders of magnitude lower than the artificial viscosity generated by the upwind scheme in high momentum zones) is added to damp oscillations appearing in the zones at rest, where no numerical diffusion is generated by our schemes. Results obtained at  $t = 20$  with the first order in time and space and the present scheme are compared on Figure 8. One can observe that the second-order scheme is clearly less diffusive. In addition, these results illustrate the capacity of the



staggered scheme to deal with reflection conditions by simply imposing the normal velocity to the boundary at zero.

### 6.5. Uniform circular motion in a paraboloid

We address in this section a classical test which admits a closed-form solution and corresponds to the uniform rotation of a drop of liquid on a paraboloid-shaped support. The solution is very regular (at a given time, the velocity field is constant and  $h + z$  is affine outside the dry zones), and the essential interest of this test is to check whether the scheme is able to cope with dry zones, *i.e.* zones where the height is zero (in the continuous setting) or very close to zero, as we shall use numerically. The computational domain is  $\Omega = (0, L) \times (0, L)$  and the topography is given by

$$z = -\frac{h_0}{a^2} \left( a^2 - \left( x - \frac{L}{2} \right)^2 - \left( y - \frac{L}{2} \right)^2 \right),$$

with  $h_0$  and  $a$  parameters which are given below. The height is:

$$h = \max(0, \bar{h}) \text{ with } \bar{h} = \eta \frac{h_0}{a^2} \left( 2 \left( x - \frac{L}{2} \right) \cos(\omega t) + 2 \left( y - \frac{L}{2} \right) \sin(\omega t) - \eta \right) - z,$$

with  $\eta$  a parameter and  $\omega$  (the angular rotation velocity of the drop) given by

$$\omega = \frac{(2gh_0)^{1/2}}{a}.$$

Finally, the velocity is

$$\mathbf{u} = \eta \omega \begin{bmatrix} -\sin(\omega t) \\ \cos(\omega t) \end{bmatrix}.$$

The computation is run up to  $T = 6\pi/\omega$ , so the drop is supposed to perform 3 turns and to lie at the final time at its initial position. The parameters are fixed here to  $L = 4$ ,  $h_0 = 0.1$ ,  $a = 1$  and  $\eta = 0.5$ .

For numerical tests, we bound  $h$  from below by  $10^{-8}$ , *i.e.* we set  $h = \max(10^{-8}, \bar{h})$ , in particular to avoid divisions by zero in the averaging steps of the Heun scheme (Equations (22e) and (22f)). The computation are performed with a uniform  $100 \times 100$  mesh, with  $\delta t = \delta_{\mathcal{M}}/16$ , without changing anything to the numerical fluxes to cope with dry zones. This is clearly dangerous, since a non-upwind approximation of the water height at a face separating two cells with a large ratio of water height may lead to a huge outflow mass flux in view of the cell mass inventory (or, in other words, a very large CFL number). This probably explains the rather small time step used here (the CFL number with respect to the celerity of the fastest waves is in the range of 1/8); the first-order scheme, which uses upwind fluxes, works with time steps four times larger. This problem would be probably cured by a more careful limitation of the mass fluxes outward an almost dry cell.

Results obtained with the first order, the segregated and the second order scheme at  $t = 6\pi/\omega$  are plotted on Figure 9. All schemes give good results, which, for the first-order scheme, is probably due to the regularity of the solution. For the momentum, one observes that the second-order scheme is less accurate than the other ones; this seems to be due to the time-stepping procedure, which perhaps generates some diffusion at the interface between dry and wet zones, especially in the last averaging step, since the segregated scheme is the most accurate one (and superimposed to the exact solution on Figure 9).

## APPENDIX A. CONSISTENCY OF NUMERICAL NON LINEAR CONVECTION FLUXES ON STAGGERED MESHES

We give here some general lemmas which generalise the Lax-Wendroff theorem to multidimensional staggered meshes, and which we state for any space dimension  $d = 1, 2$  or  $3$ . The well-known Lax-Wendroff theorem [25] states that, on uniform 1D grids, a flux-consistent and conservative cell-centered finite-volume scheme for a system of conservation laws is weakly consistent, in the sense that the limit of any a.e. convergent sequence of  $L^\infty$ -bounded numerical solutions, obtained with a sequence of grids with mesh and time

steps tending to zero, is a weak solution of the conservation law; it is also stated in a different form [26, Section 12.10], with a BV bound assumption on the scheme. It is generalised to non uniform 1D or Cartesian meshes in [11, Theorem 21.2]. In a recent work [4], the Lax-Wendroff theorem is extended to obtain some error estimates for higher order schemes on uniform 1D meshes. The case of general (and, in particular, unstructured) discretizations has been also been tackled over the past decades: [24], [15, Section 4.2.2] [10], [12]. In this latter work, the quasi-uniformity assumption that is required in [10] is relaxed, but while in [10] the flux is only required to be continuous, it is supposed to be Lipschitz continuous or at least “lip-diag”. In all these works, the scheme is supposed to be colocated, in the sense that the discrete unknowns are associated to the cells of the mesh; these results may not be used directly on staggered meshes, and for instance, in [23], the consistency of an explicit staggered scheme for the full compressible Euler equations is proven recovering the kinetic energy inequality on the primal mesh.

The consistency result that we give here is valid for general polygonal or polyhedral grids with a colocated or staggered arrangement of the unknowns. The main new idea is that in the proof of consistency, rather than using a convergence result for the discrete gradient, which is only weak and demands some regularity on the mesh, we use the actual mean value of the gradient of the test function on each cell, which converges strongly to the gradient, and does not require any regularity of the mesh. As in [12], the proof also relies on the control of some residual terms, involving the difference between the numerical solution and a space or time translate of this latter, and we use the estimate on the translates given [12, Lemma 4.2] to this purpose, which we recall in the appendix B for the sake of completeness.

Let us suppose that:

$$\Omega \subset \mathbb{R}^d, \quad d = 1, 2, 3, \quad T \in (0, +\infty), \quad (75a)$$

$$p \in \mathbb{N}^*, \quad \beta \in C^1(\mathbb{R}^p, \mathbb{R}), \quad \mathbf{f} \in C^1(\mathbb{R}^p, \mathbb{R}^d), \quad U \in L^\infty(\Omega \times (0, T), \mathbb{R}^p), \quad (75b)$$

and consider the conservative convection operator defined (in the distributional sense) by:

$$\begin{aligned} \mathcal{C}(U) : \quad \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \partial_t(\beta(U))(\mathbf{x}, t) + \operatorname{div}(\mathbf{f}(U))(\mathbf{x}, t). \end{aligned} \quad (76)$$

**Lemma A.1** (Weak consistency for a multi-dimensional conservative convection operator). *Under the assumptions (75), let  $(U^{(m)})_{m \in \mathbb{N}} \subset L^\infty(\Omega \times (0, T), \mathbb{R}^p)$  be a sequence of functions such that:*

$$\exists C^u \in \mathbb{R}_+^* : \|U^{(m)}\|_\infty \leq C^u \quad \forall m \in \mathbb{N}, \quad (77)$$

$$\exists \bar{U} \in L^\infty(\Omega \times (0, T), \mathbb{R}^p) : \|U^{(m)} - \bar{U}\|_{L^1(\Omega \times (0, T), \mathbb{R}^p)} \rightarrow 0 \text{ as } m \rightarrow +\infty. \quad (78)$$

Let  $(\mathcal{P}_m)_{m \in \mathbb{N}}$  be a sequence of polygonal or polyhedral conforming mesh of  $\Omega$  such that

$$\delta(\mathcal{P}_m) = \max_{P \in \mathcal{P}_m} \operatorname{diam}(P) \rightarrow 0 \text{ as } m \rightarrow +\infty.$$

Let  $\mathfrak{F}^{(m)}$  denote the set of faces (or edges) of the mesh, and for a given polyhedron (or polygon)  $P \in \mathcal{P}^{(m)}$ , let  $\mathfrak{F}^{(m)}(P)$  be the set of faces (or edges) of  $P$ . For  $m \in \mathbb{N}$ , let  $t_0^{(m)} = 0 < t_1^{(m)} < \dots < t_{N_m}^{(m)} = T$  be a discretization of  $(0, T)$  with  $\delta t^{(m)} = t_{k+1}^{(m)} - t_k^{(m)} \rightarrow 0$  as  $m \rightarrow +\infty$ , and consider the discrete convection operator

$$\begin{aligned} \mathcal{C}^{(m)}(U^{(m)}) : \quad \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \bar{\partial}_t(\beta^{(m)})_P^n + \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| |(\mathbf{F}^{(m)})_\zeta^n \cdot \mathbf{n}_{P,\zeta}| \text{ for } \mathbf{x} \in P \text{ and } t \in (t_n, t_{n+1}) \end{aligned} \quad (79)$$

with  $\bar{\partial}_t(\beta^{(m)})_P^n = \frac{1}{\delta t}((\beta^{(m)})_P^{n+1} - (\beta^{(m)})_P^n)$  and where the families  $\{(\beta^{(m)})_P^n, P \in \mathcal{P}^{(m)}, n \in \llbracket 0, N_m - 1 \rrbracket\}$  of real numbers and  $\{(\mathbf{F}^{(m)})_\zeta^n, \zeta \in \mathfrak{F}^{(m)}, n \in \llbracket 0, N_m - 1 \rrbracket\}$  of real vectors are such that

$$\sum_{P \in \mathcal{P}^{(m)}} \int_P |(\beta^{(m)})_P^0 - \beta(U_0(\mathbf{x}))| d\mathbf{x} \rightarrow 0 \text{ as } m \rightarrow +\infty, \text{ with } U_0 \in L^\infty(\Omega, \mathbb{R}^p), \quad (80)$$

$$\sum_{n=0}^{N_m-1} \sum_{P \in \mathcal{P}^{(m)}} \int_{t_n}^{t_{n+1}} \int_P |(\beta^{(m)})_P^n - \beta(U^{(m)}(\mathbf{x}, t))| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty, \quad (81)$$

$$\sum_{n=0}^{N_m-1} \sum_{P \in \mathcal{P}^{(m)}} \int_{t_n}^{t_{n+1}} \frac{\text{diam}(P)}{|P|} \int_P \left| \sum_{\zeta \in \mathfrak{F}^{(m)}} |\zeta| \left( (\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \right) \cdot \mathbf{n}_{P,\zeta} \right| d\mathbf{x} dt \rightarrow 0 \text{ as } m \rightarrow +\infty. \quad (82)$$

Let  $\varphi \in C_c^\infty(\Omega \times [0, t])$ , then

$$\begin{aligned} \int_0^T \int_\Omega \mathfrak{C}^{(m)}(U^{(m)})(\mathbf{x}, t) \varphi(\mathbf{x}, t) d\mathbf{x} dt &\rightarrow - \int_\Omega \beta(U_0(\mathbf{x})) \varphi(\mathbf{x}, 0) d\mathbf{x} \\ &\quad - \int_0^T \int_\Omega \left( \beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right) d\mathbf{x} dt \text{ as } m \rightarrow +\infty. \end{aligned} \quad (83)$$

*Proof.* The result of this lemma is the consequence of the two following lemmas, which prove respectively the convergence of the time derivative part and the space derivative part. Indeed, let us decompose

$$\int_0^T \int_\Omega \mathfrak{C}^{(m)}(U^{(m)})(\mathbf{x}, t) \varphi(\mathbf{x}, t) d\mathbf{x} dt = X_1^{(m)} + X_2^{(m)}, \text{ with} \quad (84)$$

$$X_1^{(m)} = \sum_{n=0}^{N_m-1} \delta t \sum_{P \in \mathcal{P}^{(m)}} |P| \bar{\partial}_t^n \beta_P^{(m)} \varphi_P^n \quad (85)$$

$$X_2^{(m)} = \sum_{n=0}^{N_m-1} \delta t \sum_{P \in \mathcal{P}^{(m)}} \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| (\mathbf{F}^{(m)})_\zeta^n \cdot \mathbf{n}_{P,\zeta} \varphi_P^n \quad (86)$$

where  $\varphi_P^n$  denotes the mean value of  $\varphi$  on  $P \times (t_n, t_{n+1})$ . Then, by Lemma A.2 below,

$$X_1^{(m)} \rightarrow - \int_\Omega \beta(U_0(\mathbf{x})) d\mathbf{x} - \int_0^T \int_\Omega \beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) d\mathbf{x} dt \text{ as } m \rightarrow +\infty,$$

and by Lemma A.3 below,

$$X_2^{(m)} \rightarrow - \int_0^T \int_\Omega \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) d\mathbf{x} dt \text{ as } m \rightarrow +\infty,$$

which concludes the proof □

**Lemma A.2** (Weak consistency, time derivative). *Under the assumptions and notations of Lemma A.1,*

$$\int_0^T \int_\Omega \bar{\partial}_t(\beta^{(m)})_P^n \varphi(\mathbf{x}, t) d\mathbf{x} dt \rightarrow - \int_\Omega \beta(U_0(\mathbf{x})) \varphi(\mathbf{x}, 0) d\mathbf{x} - \int_0^T \int_\Omega \beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) d\mathbf{x} dt$$

*Proof.* By definition of  $\bar{\partial}_t^n \beta_P^{(m)}(\mathbf{x}, t)$  and thanks to a discrete integration by parts,

$$\begin{aligned} X_1^{(m)} &= \int_0^T \int_{\Omega} \bar{\partial}_t(\beta^{(m)})_P^n \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ &= - \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^0 \varphi_P^0 - \sum_{n=1}^{N_m} \delta t \sum_{P \in \mathcal{P}^{(m)}} |P| \beta_P^{(m)}(\mathbf{x}, t) \frac{1}{\delta t} (\varphi_P^n - \varphi_P^{n-1}). \end{aligned}$$

Thanks to the assumptions (77), (78) (80) and (81), we get that

$$\lim_{m \rightarrow +\infty} X_1^{(m)} = - \int_{\Omega} \beta(U_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_{\Omega} (\beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) \, dt \, d\mathbf{x}. \quad (87)$$

□

**Lemma A.3** (Weak consistency, space derivative). *Under the assumptions and notations of Lemma A.1,*

$$\int_0^T \int_{\Omega} \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| (\mathbf{F}^{(m)})_{\zeta}^n \cdot \mathbf{n}_{P,\zeta} \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \rightarrow - \int_0^T \int_{\Omega} \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty.$$

*Proof.* Let  $X_2^{(m)}$  denote the left-hand-side of the above assertion. Since for a face  $\zeta$  separating  $P$  and  $P'$ , one has  $\mathbf{n}_{P,\zeta} = -\mathbf{n}_{P',\zeta}$ , we may rewrite  $X_2^{(m)}$  as

$$\begin{aligned} X_2^{(m)} &= \int_0^T \int_{\Omega} \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| (\mathbf{F}^{(m)})_{\zeta}^n \cdot \mathbf{n}_{P,\zeta} \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ &= \sum_{n=0}^{N_m-1} \delta t^{(m)} \sum_{P \in \mathcal{P}^{(m)}} A_P^n \text{ with } A_P^n = \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| (\mathbf{F}^{(m)})_{\zeta}^n \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_{\zeta}^n), \end{aligned}$$

where  $\varphi_P^n$  (resp.  $\varphi_{\zeta}^n$ ) denotes the mean value of  $\varphi$  over  $P \times (t_n, t_{n+1})$  (resp.  $\zeta \times (t_n, t_{n+1})$ ). Now for any  $\mathbf{x} \in P$ ,  $t \in [t_n, t_{n+1})$ , we can decompose  $A_P^n$  as

$$\begin{aligned} A_P^n &= B_P^n(\mathbf{x}, t) + R_P^n(\mathbf{x}, t), \text{ with } B_P^n(\mathbf{x}) = \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_{\zeta}^n), \text{ and} \\ R_P^n(\mathbf{x}, t) &= \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| (\mathbf{F}_{\zeta}^n - \mathbf{f}(U^{(m)}(\mathbf{x}, t))) \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_{\zeta}^n) \, d\mathbf{x}. \end{aligned}$$

Since  $\sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| \mathbf{n}_{P,\zeta} = 0$ , we have

$$\begin{aligned} B_P^n(\mathbf{x}, t) &= - \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \cdot \mathbf{n}_{P,\zeta} \varphi_{\zeta}^n = -|P| \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \cdot (\nabla \varphi)_P^n, \\ &\text{with } (\nabla \varphi)_P^n = \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}^{(m)}(P)} |\zeta| \varphi_{\zeta}^n \mathbf{n}_{P,\zeta} = \frac{1}{|P|} \nabla \varphi(\mathbf{x}) \, d\mathbf{x}. \end{aligned} \quad (88)$$

Note that the piecewise function  $\nabla^{(m)} \varphi : \Omega \times (0, T) \rightarrow \mathbb{R}^d$  defined by  $\nabla^{(m)} \varphi(\mathbf{x}, t) = (\nabla \varphi)_P^n$  for  $(\mathbf{x}, t) \in P \times (t_n, t_{n+1})$  converges uniformly to  $\nabla \varphi$  in  $L^\infty(\Omega \times (0, T))^d$ . Integrating (88) over  $\mathbf{x} \in P$ ,

$$\int_P B_P^n(\mathbf{x}, t) \, d\mathbf{x} = |P| \int_P \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \cdot \nabla^{(m)} \varphi(\mathbf{x}, t) \, d\mathbf{x},$$

Since  $A_P^n = \frac{1}{\delta t^{(m)}|P|} \left( \int_{t_n}^{t_{n+1}} \int_P B_P^n(\mathbf{x}, t) d\mathbf{x} dt + \int_{t_n}^{t_{n+1}} \int_P R_P^n(\mathbf{x}, t) d\mathbf{x} dt \right)$ , we get

$$\begin{aligned} X_2^{(m)} &= \sum_{n=0}^{N_m-1} \sum_{P \in \mathcal{P}^{(m)}} \left( \int_{t_n}^{t_{n+1}} \int_P B_P^n(\mathbf{x}, t) d\mathbf{x} dt + \int_{t_n}^{t_{n+1}} \frac{1}{|P|} \int_P R_P^n(\mathbf{x}, t) d\mathbf{x} dt \right) \\ &= - \int_0^T \int_{\Omega} \mathbf{f}(U^{(m)}(\mathbf{x}, t)) \nabla^{(m)} \varphi(\mathbf{x}, t) d\mathbf{x} dt + \sum_{n=0}^{N_m-1} \sum_{P \in \mathcal{P}^{(m)}} \int_{t_n}^{t_{n+1}} \frac{1}{|P|} \int_P R_P^n(\mathbf{x}, t) d\mathbf{x} dt. \end{aligned}$$

Owing to the boundedness and convergence assumptions on  $U^{(m)}$  and to the uniform convergence of  $\nabla^{(m)} \varphi$  to  $\nabla \varphi$ , the first term tends to  $-\int_0^T \int_{\Omega} \mathbf{f}(U(\mathbf{x}, t)) \nabla \varphi(\mathbf{x}, t) d\mathbf{x} dt$  as  $m \rightarrow +\infty$ . Since  $|\varphi_{\zeta}^n - \varphi_P^n| \leq C_{\varphi} \text{diam}(P)$ , with  $C_{\varphi}$  depending only on  $\varphi$ , the second term tends to 0 thanks to the assumption (82). Therefore

$$\lim_{m \rightarrow +\infty} X_2^{(m)} \rightarrow - \int_0^T \int_{\Omega} \mathbf{f}(U(\mathbf{x}, t)) \cdot \nabla \varphi(\mathbf{x}, t) d\mathbf{x} dt. \quad (89)$$

□

## APPENDIX B. FORMER LEMMAS

### B.1. A result on a finite volume convection operator

We begin with a property of the convection operator  $\mathcal{C} : \rho \mapsto \partial_t(\rho) + \text{div}(\rho \mathbf{u})$ ; at the continuous level, this property may be formally obtained as follows (see [22] for the detailed derivation). Let  $\psi$  be a regular function from  $(0, +\infty)$  to  $\mathbb{R}$ ; then:

$$\psi'(\rho) \mathcal{C}(\rho) = \partial_t(\psi(\rho)) + \text{div}(\psi(\rho) \mathbf{u}) + (\rho \psi'(\rho) - \psi(\rho)) \text{div} \mathbf{u}. \quad (90)$$

This computation is of course completely formal and only valid for regular functions  $\rho$  and  $\mathbf{u}$ . The following lemma states a discrete analogue to (90), and its proof follows the formal computation which we just described.

**Lemma B.1.** *[On the discrete convection operator, [21, Lemma A1]] Let  $P$  be a polygonal (resp. polyhedral) bounded set of  $\mathbb{R}^2$  (resp.  $\mathbb{R}^3$ ), and let  $\mathcal{E}(P)$  be the set of its edges (resp. faces). Let  $\psi$  be a twice continuously differentiable function defined over  $(0, +\infty)$ . Let  $\rho_P^* > 0$ ,  $\rho_P > 0$ ,  $\delta t > 0$ ; consider three families  $(\rho_{\eta}^*)_{\eta \in \mathcal{E}(P)} \subset \mathbb{R}_+ \setminus \{0\}$ ,  $(V_{\eta}^*)_{\eta \in \mathcal{E}(P)} \subset \mathbb{R}$  and  $(F_{\eta}^*)_{\eta \in \mathcal{E}(P)} \subset \mathbb{R}$  such that*

$$\forall \eta \in \mathcal{E}(P), \quad F_{\eta}^* = \rho_{\eta}^* V_{\eta}^*.$$

Let  $R_{P, \delta t}$  be defined by:

$$\begin{aligned} R_{P, \delta t} &= \left[ \frac{|P|}{\delta t} (\rho_P - \rho_P^*) + \sum_{\eta \in \mathcal{E}(P)} F_{\eta}^* \right] \psi'(\rho_P) \\ &\quad - \left[ \frac{|P|}{\delta t} [\psi(\rho_P) - \psi(\rho_P^*)] + \sum_{\eta \in \mathcal{E}(P)} \psi(\rho_{\eta}^*) V_{\eta}^* + [\rho_P^* \psi'(\rho_P^*) - \psi(\rho_P^*)] \sum_{\eta \in \mathcal{E}(P)} V_{\eta}^* \right]. \end{aligned}$$

Then this quantity may be expressed as follows:

$$R_{P, \delta t} = \frac{1}{2} \frac{|P|}{\delta t} (\rho_P - \rho_P^*)^2 \psi''(\bar{\rho}_P^{(1)}) - \frac{1}{2} \sum_{\eta \in \mathcal{E}(P)} V_{\eta}^* (\rho_P^* - \rho_{\eta}^*)^2 \psi''(\bar{\rho}_{\eta}^*) + \sum_{\eta \in \mathcal{E}(P)} V_{\eta}^* \rho_{\eta}^* (\rho_P - \rho_P^*) \psi''(\bar{\rho}_P^{(2)}),$$

where  $\bar{\rho}_P^{(1)}, \bar{\rho}_P^{(2)} \in [\rho_P, \rho_P^*]$  and  $\forall \eta \in \mathcal{E}(P)$ ,  $\bar{\rho}_{\eta}^* \in [\rho_P^*, \rho_{\eta}^*]$ . We recall that, for  $a, b \in \mathbb{R}$ , we denote by  $\llbracket a, b \rrbracket$  the interval  $\llbracket a, b \rrbracket = \{\theta a + (1 - \theta)b, \theta \in [0, 1]\}$ .

## B.2. A result on the space translates

**Lemma B.2** (Convergence of the space translates [12, Lemma 4.2]). *For a given mesh  $\mathcal{M}$ , let*

$$\theta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \max_{\sigma \in \mathcal{E}_K} \frac{|D_{\sigma}|}{|K|}.$$

*Let  $\theta > 0$  and  $(\mathcal{M}^{(m)})_{m \in \mathbb{N}}$  be a sequence of meshes such that  $\theta_{\mathcal{M}^{(m)}} \leq \theta$  for all  $m \in \mathbb{N}$  and  $\lim_{m \rightarrow +\infty} h_{\mathcal{M}^{(m)}} = 0$ . We suppose that the number of faces of a cell  $K \in \mathcal{M}^{(m)}$  is bounded by  $\mathcal{N}_{\mathcal{E}}$ , for any  $m \in \mathbb{N}$ . Let  $\psi \in L^1(\Omega)$ , let  $\langle \psi \rangle_K$  denote the mean value of  $\psi$  on a cell  $K$ . Then,*

$$\lim_{m \rightarrow +\infty} \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} |D_{\sigma}| |\langle \psi \rangle_K - \langle \psi \rangle_L| = 0. \quad (91)$$

## REFERENCES

- [1] 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*, 66:555–580, 2011.
- [2] A. Arakawa and V. Lamb. A potential enstrophy and energy conserving scheme for the shallow water equations. *Monthly Weather Review*, 109:18–36, 1981.
- [3] E. Audusse. *Autour du système de Saint-Venant : Méthodes numériques pour le transport sédimentaire, les fluides en rotation et les équations primitives*. Habilitation à diriger des recherches, Université Paris 13 Villetaneuse, Nov. 2018.
- [4] M. Ben-Artzi and J. Li. Consistency and convergence of finite volume approximations to nonlinear hyperbolic balance laws, 2019.
- [5] L. Bonaventura and T. Ringler. Analysis of discrete shallow-water models on geodesic delaunay grids with c-type staggering. *Monthly Weather Review*, 133(8):2351–2373, 2005.
- [6] F. Bouchut. *Nonlinear Stability of finite volume methods for hyperbolic conservation laws*. Birkhauser, 2004.
- [7] CALIF<sup>3</sup>S. A software components library for the computation of fluid flows. <https://gforge.irsnn.fr/gf/project/califs>.
- [8] M. J. Castro, T. Morales de Luna, and C. Parés. Well-balanced schemes and path-conservative numerical methods. In *Handbook of numerical methods for hyperbolic problems*, volume 18 of *Handb. Numer. Anal.*, pages 131–175. Elsevier/North-Holland, Amsterdam, 2017.
- [9] D. Doyen and H. Gunawan. An explicit staggered finite volume scheme for the shallow water equations. In *Finite volumes for complex applications. VII. Methods and theoretical aspects*, volume 77 of *Springer Proc. Math. Stat.*, pages 227–235. Springer, Cham, 2014.
- [10] V. Elling. A Lax-Wendroff type theorem for unstructured quasi-uniform grids. *Mathematics of Computation*, 76:251–272, 2007.
- [11] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. In P. Ciarlet and J. Lions, editors, *Handbook of Numerical Analysis, Volume VII*, pages 713–1020. North Holland, 2000, <https://hal.archives-ouvertes.fr/>.
- [12] T. Gallouët, R. Herbin, and J.-C. Latché. On the weak consistency of finite volumes schemes for conservation laws on general meshes. *SeMA J.*, 76(4):581–594, 2019.
- [13] T. Gallouët, R. Herbin, J.-C. Latché, and K. Mallem. Convergence of the marker-and-cell scheme for the incompressible Navier-Stokes equations on non-uniform grids. *Foundations of Computational Mathematics*, 18:249–289, 2018.
- [14] T. Gallouët, R. Herbin, J.-C. Latché, and Y. Nasser. A second order consistent MAC scheme for the shallow water equations on non uniform grids. In *Finite Volumes for Complex Applications IX - Methods, Theoretical Aspects, Examples*, pages 123–131. Springer, 6 2020.
- [15] E. Godlewski and P.-A. Raviart. Numerical approximation of hyperbolic systems of conservation laws. In *Springer*, page 118. Applied Mathematical Sciences, New York, 1996.
- [16] H. Gunawan. *Numerical simulation of shallow water equations and related models*. PhD thesis, Université Paris-Est and Institut Teknologi Bandung, 2015.
- [17] F. Harlow and A. Amsden. A numerical fluid dynamics calculation method for all flow speeds. *Journal of Computational Physics*, 8:197–213, 1971.
- [18] F. Harlow and J. Welsh. Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface. *Physics of Fluids*, 8:2182–2189, 1965.
- [19] R. Herbin and J.-C. Latché. Kinetic energy control in the MAC discretization of the compressible Navier-Stokes equations. *International Journal on Finite Volumes*, 7, 2010.
- [20] R. Herbin, J.-C. Latché, Y. Nasser, and N. Therme. A decoupled staggered scheme for the shallow water equations. *Monografías Matemáticas García de Galdeano*, 52:1–16, 2019.
- [21] R. Herbin, J.-C. Latché, and T. Nguyen. Explicit staggered schemes for the compressible Euler equations. *ESAIM: Proceedings*, 40:83–102, 2013.

- [22] R. Herbin, J.-C. Latché, and T. Nguyen. Consistent segregated staggered schemes with explicit steps for the isentropic and full Euler equations. *ESAIM: Mathematical Modelling and Numerical Analysis*, 52:893–944, 2018.
- [23] R. Herbin, J.-C. Latch, S. Minjeaud, and N. Therme. Conservativity and weak consistency of a class of staggered finite volume methods for the euler equations. *Mathematics of Computation*, on line 2020.
- [24] D. Kroner, M. Rokyta, and M. Wierse. A Lax-Wendroff type theorem for upwind finite volume schemes in 2-D. *East-West Journal of Numerical Mathematics*, 4:279–292, 1996.
- [25] P. Lax and B. Wendroff. Systems of conservation laws. *Communications in Pure and Applied Mathematics*, 13:217–237, 1960.
- [26] R. Leveque. *Finite Volume Methods for Hyperbolic Problems*. Cambridge texts in applied mathematics. Cambridge University Press, 2002.
- [27] L. Piar, F. Babik, R. Herbin, and J.-C. Latché. A formally second-order cell centred scheme for convection-diffusion equations on general grids. *Internat. J. Numer. Methods Fluids*, 71(7):873–890, 2013.
- [28] G. Stelling and S. Duinmeijer. A staggered conservative scheme for every Froude number in rapidly varied shallow water flows. *International Journal for Numerical Methods in Fluids*, 43:1329–1354, 2003.
- [29] W.-Y. Tan. *Shallow water hydrodynamics: Mathematical theory and numerical solution for a two-dimensional system of shallow-water equations*. Elsevier, 1992.
- [30] Y. Xing. Numerical methods for the nonlinear shallow water equations. In *Handbook of numerical methods for hyperbolic problems*, volume 18 of *Handb. Numer. Anal.*, pages 361–384. Elsevier/North-Holland, Amsterdam, 2017.

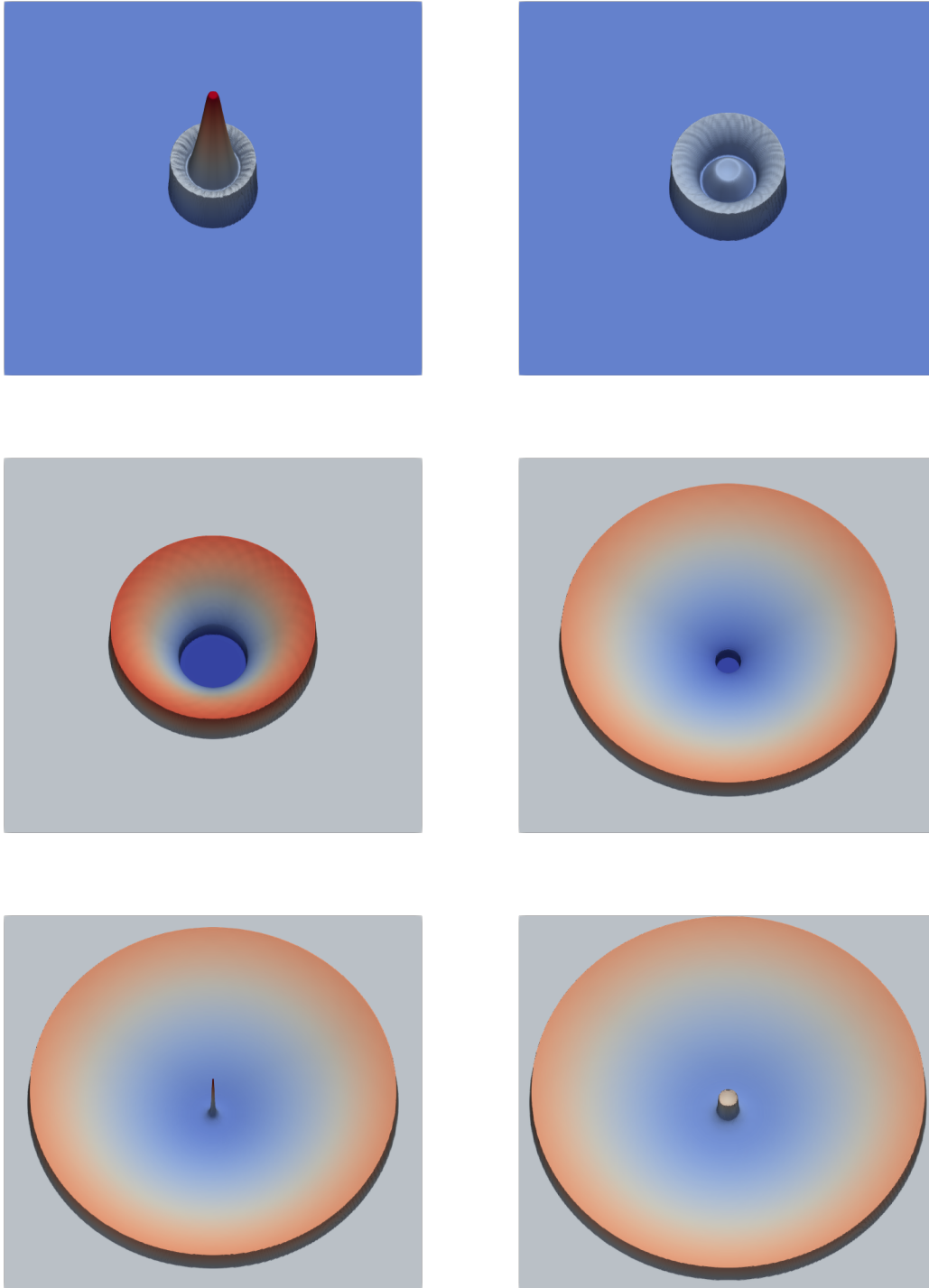


FIGURE 5. Circular dam-break problem. Height obtained at  $t = 0.38$ ,  $t = 0.705$ ,  $t = 1.88$ ,  $t = 3.76$ ,  $t = 4.28$  and  $t = T = 4.7$  with the stabilized second-order scheme and a  $800 \times 800$  mesh. The color range corresponds to the  $(0.1, 2.5)$  interval for the first two plots, and to the  $(0.1, 1)$  interval for the last four ones.



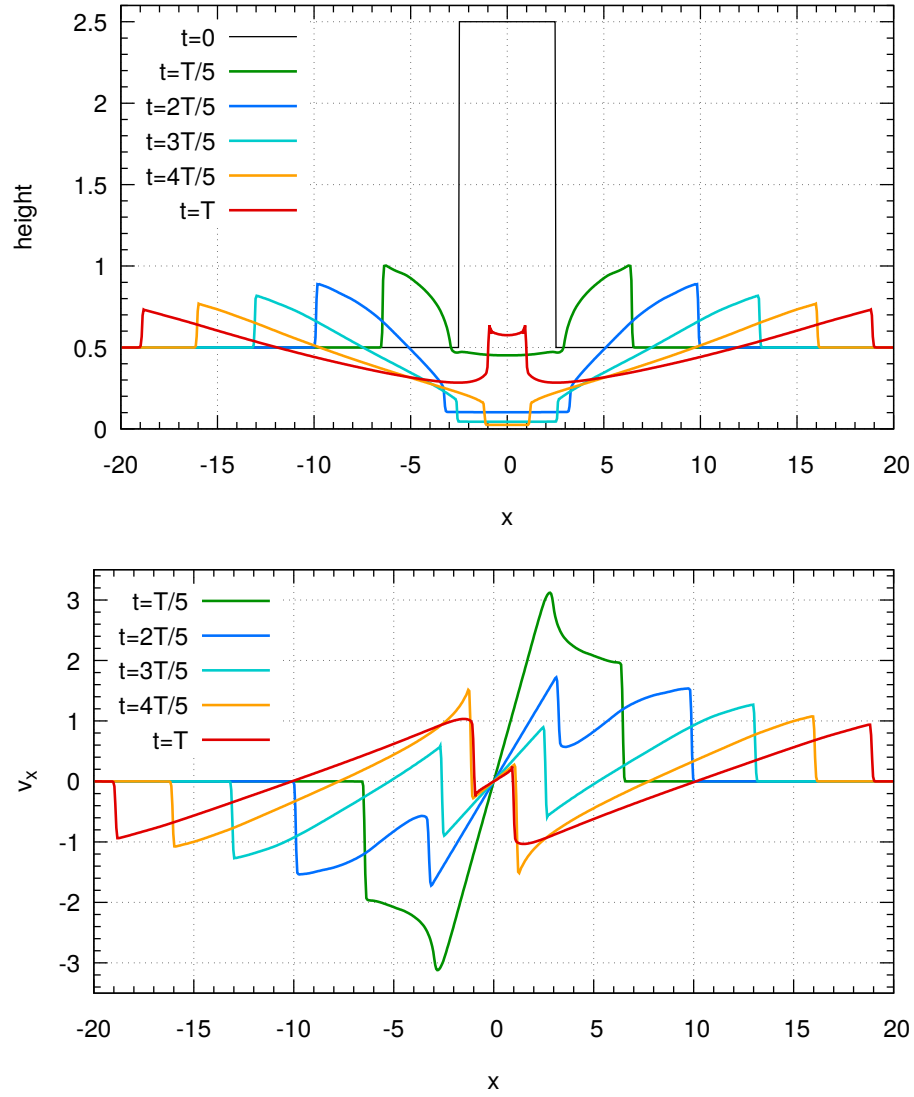


FIGURE 6. Circular dam-break problem. Height and radial velocity obtained at different times along the line  $x_2 = 0$  with the stabilized second-order scheme and a  $800 \times 800$  mesh.

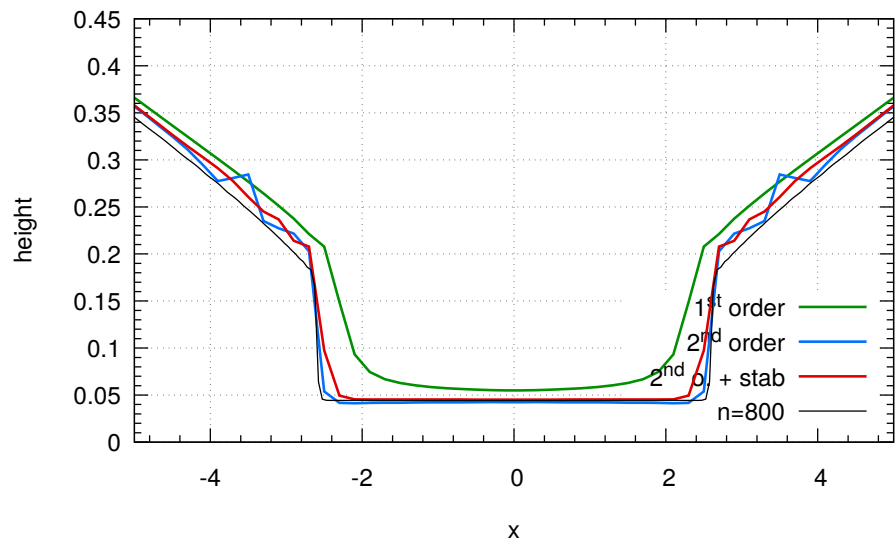
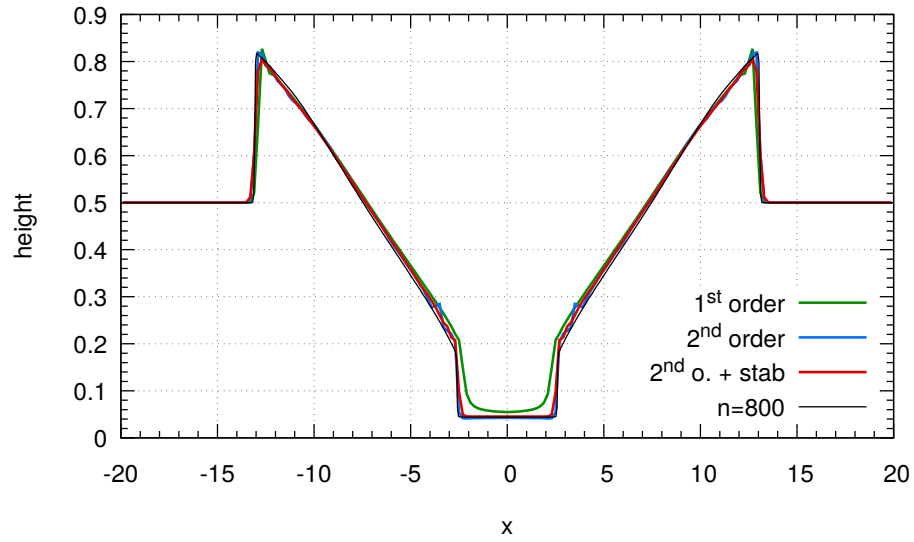


FIGURE 7. Circular dam-break problem. Height obtained at  $t = 3T/5$  with the first-order scheme and the second-order scheme with and without stabilization, with a  $200 \times 200$  mesh.

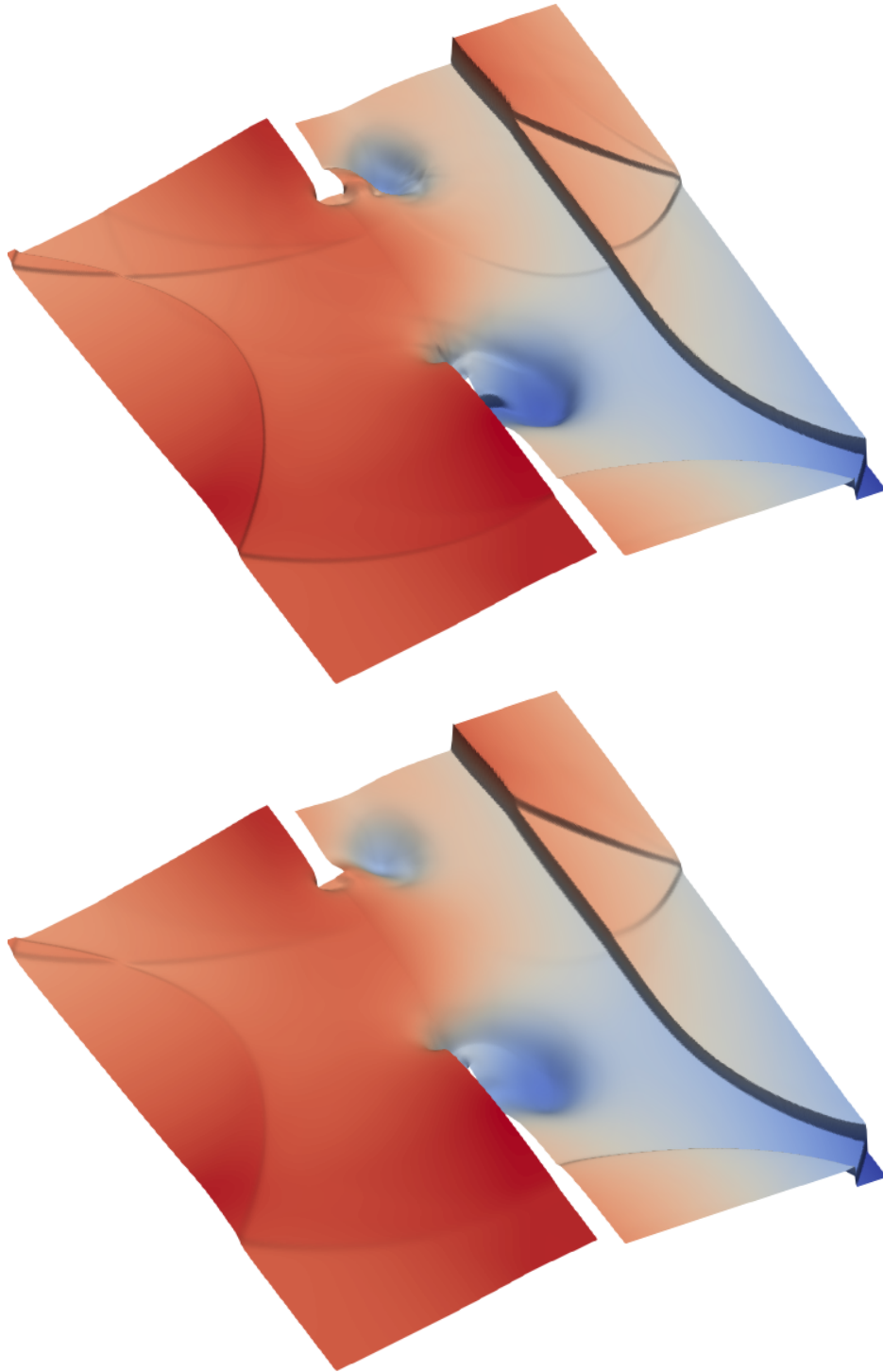


FIGURE 8. Partial dam-break flow. Top: MUSCL scheme – Bottom: upwind scheme.

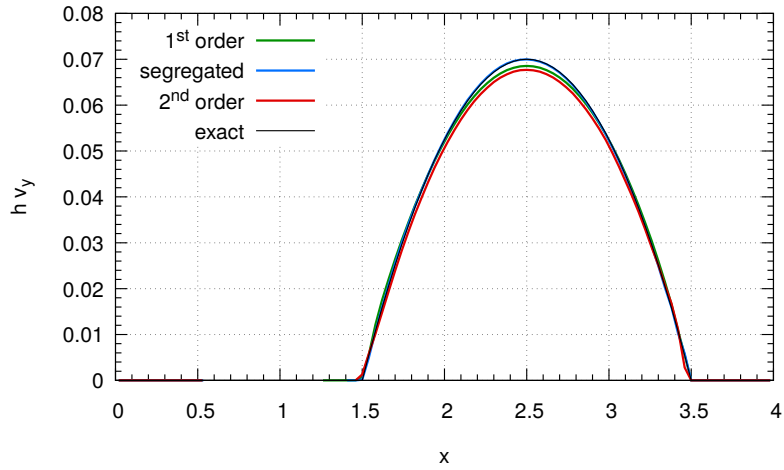
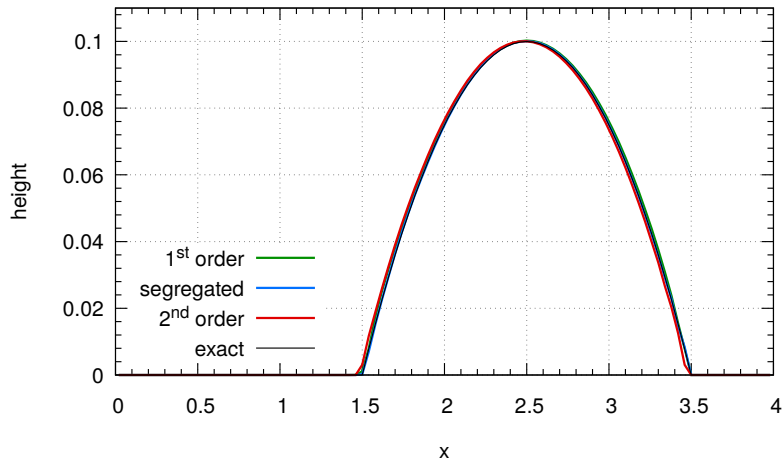
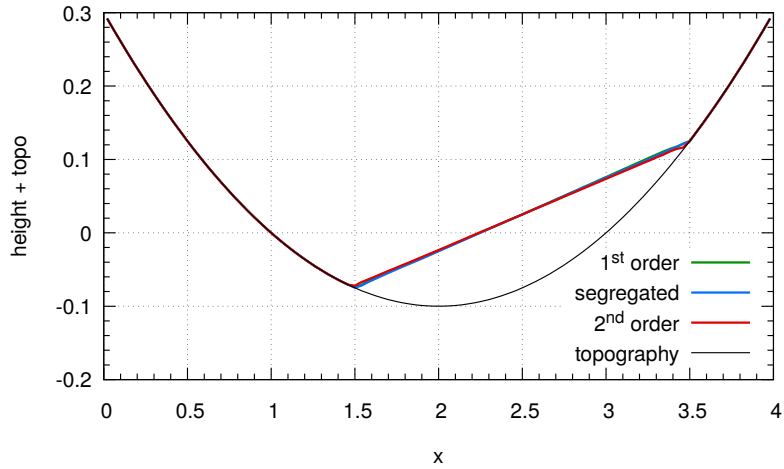


FIGURE 9. Circular motion of a drop over a paraboloid-shaped topography. Sum of the height and the topography, height alone and second component of the momentum along the  $y = L/2$  line at  $t = 6\pi$ .