# A DECOUPLED STAGGERED SCHEME FOR THE SHALLOW WATER EQUATIONS

Raphaèle Herbin, Jean-Claude Latché, Youssouf Nasseri and Nicolas Therme

**Abstract.** We present a first order scheme based on a staggered grid for the shallow water equations with topography in two space dimensions, which enjoys several properties: positivity of the water height, preservation of constant states, and weak consistency with the equations of the problem and with the associated entropy inequality.

*Keywords:* Shallow water, finite volumes, staggered grid.

*AMS classification:* 65M08,76B99.

## §1. Introduction

The shallow water equations form a hyperbolic system of two conservation equations (mass and momentum) which are obtained when modelling a flow whose vertical height is considered small with respect to the plane scale. The solution of such a system may develop shocks, so that the finite volume method is usually preferred for numerical simulations. Two main approaches are found: one is the colocated approach which is usually based on some approximate Riemann solver, see e.g. [3] and references therein; the other one is based on a staggered arrangement of the unknowns on the grid. Indeed, staggered schemes have been used for some time in the hydraulic and ocean engineering community, see e.g. [1, 2, 12]. They have been recently analysed in the case of one space dimension [5, 8], following the works on the related barotropic Euler equations, see [11] and references therein. In the present work, we obtain a discrete local entropy inequality; furthermore, we extend the consistency analysis of the scheme to the case of two space dimensions, and we weaken the assumptions on the estimates, namely we no longer require a bound on the $BV$ norm of the approximate solutions, at least for the weak formulation (the passage to the limit in the entropy still necessitates a time $BV$ boundedness).

Let $\Omega$ be an open bounded domain of $\mathbb{R}^2$ and let $T > 0$. We consider the shallow water equations with topography over the space and time domain $\Omega \times (0, T)$:

$$\partial_t h + \mathrm{div}(h\boldsymbol{u}) = 0 \qquad\qquad \text{in } \Omega \times (0, T), \tag{1a}$$

$$\partial_t(h\boldsymbol{u}) + \mathrm{div}(h\boldsymbol{u} \otimes \boldsymbol{u}) + \nabla p + gh\nabla z = 0 \qquad\qquad \text{in } \Omega \times (0, T), \tag{1b}$$

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

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

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

where $t$ stands for the time, $g$ is the standard gravity constant and $z$ the (given) topography, which is supposed to be regular in this paper. These equations solve the water height $h$ and the velocity $\boldsymbol{u}$.

Let us recall that if $(h, \boldsymbol{u})$ is a regular solution of (1), the following elastic potential energy balance and kinetic energy balance is obtained by manipulations on the mass and momentum equations:

$$\partial_t(\frac{1}{2}gh^2) + \mathrm{div}(\frac{1}{2}gh^2\boldsymbol{u}) + \frac{1}{2}gh^2\mathrm{div}\boldsymbol{u} = 0 \tag{2}$$

$$\partial_t(\frac{1}{2}h|\boldsymbol{u}|^2) + \mathrm{div}(\frac{1}{2}h|\boldsymbol{u}|^2\boldsymbol{u}) + \boldsymbol{u} \cdot \nabla p + gh\boldsymbol{u} \cdot \nabla z = 0. \tag{3}$$

Summing these equations, we obtain en entropy equality of the form $\partial_t \eta + \mathrm{div}\Phi = 0$, where the entropy-entropy flux pair $(\eta, \Phi)$ is given by:

$$\eta = \frac{1}{2}h|\boldsymbol{u}|^2 + \frac{1}{2}gh^2 + ghz \text{ and } \Phi = (\eta + \frac{1}{2}gh^2)\boldsymbol{u}. \tag{4}$$

For non regular functions the above manipulations are no longer valid, and the entropy inequality $\partial_t \eta + \mathrm{div}\Phi \leq 0$ is satisfied in a distributional sense.

In this paper, we build a decoupled scheme, involving only explicit steps; the resulting approximate solutions are shown to satisfy some discrete equivalent of (2) and (3); furthermore, under some convergence and boundedness assumptions, the approximate solutions are shown in Section 5 to converge to a weak solution of (1) and to satisfy a weak entropy inequality.

## §2. Mesh and space discretizations

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 $(\boldsymbol{e}^{(1)}, \boldsymbol{e}^{(2)})$.

**Definition 1** (MAC grid). A discretization $(\mathcal{M}, \mathcal{E})$ of $\Omega$ with a staggered rectangular grid (or MAC grid), is defined by:

- A primal grid $\mathcal{M}$ which consists in a conforming structured partition of $\Omega$ in rectangles, possibly non uniform. A generic cell of this grid is denoted by $K$, and its mass center by $\boldsymbol{x}_K$. The scalar unknowns (water height and pressure) are associated to this mesh.
- The set of all edges of the mesh $\mathcal{E}$, with $\mathcal{E} = \mathcal{E}_{\mathrm{int}} \cup \mathcal{E}_{\mathrm{ext}}$, where $\mathcal{E}_{\mathrm{int}}$ (resp. $\mathcal{E}_{\mathrm{ext}}$) 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 $\boldsymbol{e}^{(i)}$ is denoted by $\mathcal{E}^{(i)}$, for $i = 1, 2$. We then have $\mathcal{E}^{(i)} = \mathcal{E}_{\mathrm{int}}^{(i)} \cup \mathcal{E}_{\mathrm{ext}}^{(i)}$, where $\mathcal{E}_{\mathrm{int}}^{(i)}$ (resp. $\mathcal{E}_{\mathrm{ext}}^{(i)}$) are the edges of $\mathcal{E}^{(i)}$ that lie in the interior (resp. on the boundary) of the domain.

  For $\sigma \in \mathcal{E}_{\mathrm{int}}$, 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 $\sigma = K|L \in \mathcal{E}_{\mathrm{int}}$ 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);
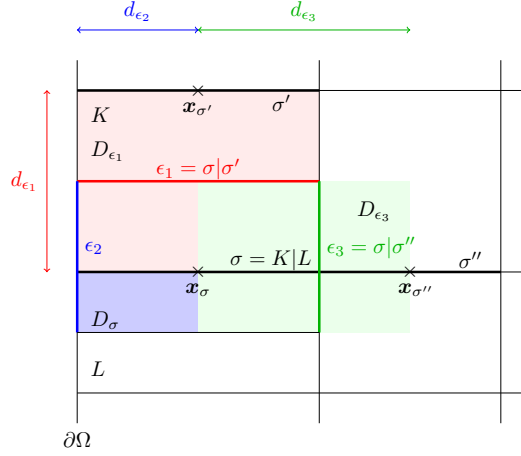
Figure 1: Notations for control volumes and dual cells (in two space dimensions, for the second component of the velocity).

- if $\sigma \in \mathcal{E}_{ext}$ is adjacent to the cell $K$, then $D_\sigma = D_{K,\sigma}$.

For each dimension $i = 1, 2$, the domain $\Omega$ is partitioned in dual cells: $\Omega = \cup_{\sigma \in \mathcal{E}^{(i)}} D_\sigma$, $i = 1, 2$; the $i^{th}$ partition is refered to as the $i^{th}$ dual mesh; it is associated to the $i^{th}$ velocity component, in a sense which is clarified below. The set of the edges of the $i^{th}$ dual mesh is denoted by $\widetilde{\mathcal{E}}^{(i)}$ (note that these edges may be orthogonal to any vector of the basis of $\mathbb{R}^2$ and not only $e^{(i)}$) and is decomposed into the internal and boundary edges: $\widetilde{\mathcal{E}}^{(i)} = \widetilde{\mathcal{E}}^{(i)}_{int} \cup \widetilde{\mathcal{E}}^{(i)}_{ext}$. The dual edge separating two duals cells $D_\sigma$ and $D_{\sigma'}$ is denoted by $\epsilon = \sigma|\sigma'$. We denote by $D_\epsilon$ the dual cell associated to a dual edge $\epsilon \in \widetilde{\mathcal{E}}$ defined as follows:

- if $\epsilon = \sigma|\sigma' \in \widetilde{\mathcal{E}}_{int}$ 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 $\epsilon \in \widetilde{\mathcal{E}}_{ext}$ 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 $\mathcal{E}(K)$ and $\widetilde{\mathcal{E}}(D_\sigma)$ respectively. For $\sigma \in \mathcal{E}$, we denote by $x_\sigma$ the mass center of $\sigma$. The vector $n_{K,\sigma}$ 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 $K$ will be denoted $K = [\overrightarrow{\sigma\sigma'}]$ if $\sigma, \sigma' \in \mathcal{E}^{(i)}(K)$ for some $i = 1, 2$ are such that $(x_{\sigma'} - x_\sigma) \cdot e^{(i)} > 0$;
- we write $\sigma = \overrightarrow{K|L}$ if $\sigma \in \mathcal{E}^{(i)}$, $\sigma = K|L$ and $\overrightarrow{x_K x_L} \cdot e^{(i)} > 0$ for some $i = 1, 2$;
- the dual edge $\epsilon$ separating $D_\sigma$ and $D_{\sigma'}$ is written $\epsilon = \overrightarrow{\sigma|\sigma'}$ if $\overrightarrow{x_\sigma x_{\sigma'}} \cdot e^{(i)} > 0$ for some $i = 1, 2$.

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

$$\delta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \text{diam}(K), \text{ and } \eta_{\mathcal{M}} = \max\left\{\frac{|\sigma|}{|\sigma'|}, \ \sigma \in \mathcal{E}^{(i)}, \ \sigma' \in \mathcal{E}^{(j)}, \ i, j = 1, 2, \ i \neq j\right\}, \quad (5)$$

where $|\cdot|$ stands for the one (or two) dimensional measure of a subset of $\mathbb{R}$ (or $\mathbb{R}^2$).

The discrete velocity unknowns are associated to the dual cells and denoted by $(u_{i,\sigma})_{\sigma \in \mathcal{E}^{(i)}}$, $i = 1, 2$, while the scalar unknowns (discrete water height and pressure) are associated to the primal cells and are denoted respectively by $(h_K)_{K \in \mathcal{M}}$ and $(p_K)_{K \in \mathcal{M}}$. The scalar unknown space $L_{\mathcal{M}}$ is defined as the set of piecewise constant functions over each grid cell $K$ of $\mathcal{M}$, and the discrete $i^{th}$ velocity space $H_{\mathcal{E}^{(i)}}$ as the set of piecewise constant functions over each of the grid cells $D_\sigma$, $\sigma \in \mathcal{E}^{(i)}$. As in the continuous case, the Dirichlet boundary conditions are taken into account by defining the subspaces $H_{\mathcal{E}^{(i)},0} \subset H_{\mathcal{E}^{(i)}}$, $i = 1, 2$ as follows

$$H_{\mathcal{E}^{(i)},0} = \left\{u_i \in H_{\mathcal{E}^{(i)}}, \ u_i(\boldsymbol{x}) = 0, \ \forall \boldsymbol{x} \in D_\sigma, \ \sigma \in \mathcal{E}^{(i)}_{\text{ext}}\right\}.$$

We then set $\boldsymbol{H}_{\mathcal{E},0} = H_{\mathcal{E}^{(1)},0} \times H_{\mathcal{E}^{(2)},0}$. Defining the characteristic function $\mathbb{1}_A$ of any subset $A \subset \Omega$ by $\mathbb{1}_A(\boldsymbol{x}) = 1$ if $\boldsymbol{x} \in A$ and $\mathbb{1}_A(\boldsymbol{x}) = 0$ otherwise, the functions $\boldsymbol{u} = (u_1, u_2) \in \boldsymbol{H}_{\mathcal{E},0}$, may then be written:

$$u_i(\boldsymbol{x}) = \sum_{\sigma \in \mathcal{E}^{(i)}} u_{i,\sigma} \mathbb{1}_{D_\sigma}(\boldsymbol{x}), \ i = 1, 2. \quad (6)$$

For $\boldsymbol{u} \in \boldsymbol{H}_{\mathcal{E},0}$, let $[\![u_i]\!]_\epsilon = |u_{i,\sigma} - u_{i,\sigma'}|$, for $\epsilon = \sigma|\sigma' \in \widetilde{\mathcal{E}}^{(i)}_{int}, i = 1, 2$. In the same way the functions $h \in L_{\mathcal{M}}$ are defined by $h(\boldsymbol{x}) = \sum_{K \in \mathcal{M}} h_K \mathbb{1}_K(\boldsymbol{x})$ and the notation $[\![\ ]\!]_\sigma$ refers to $[\![h]\!]_\sigma = |h_K - h_L|$, for $\sigma = K|L \in \mathcal{E}_{\text{int}}(K)$.

## §3. A decoupled explicit scheme

**Description of the scheme** Let us consider a uniform discretisation $0 = t_0 < t_1 < \cdots < t_N = T$ of the time interval $(0, T)$, and let $\delta t = t_{n+1} - t_n$ for $n = 0, 1, \cdots, N - 1$ be the (constant) time step. The discrete velocity $\boldsymbol{u}$ and water height $h$ unknowns are defined by:

$$\boldsymbol{u}(\boldsymbol{x}, t) = \sum_{n=0}^{N-1} \boldsymbol{u}^{n+1}(\boldsymbol{x}) \mathbb{1}_{[t_n, t_{n+1})}(t), \text{ with } \boldsymbol{u}^{n+1} \in \boldsymbol{H}_{\mathcal{E},0},$$

$$h(\boldsymbol{x}, t) = \sum_{n=0}^{N-1} h^{n+1}(\boldsymbol{x}) \mathbb{1}_{[t_n, t_{n+1})}(t), \text{ with } h^{n+1} \in L_{\mathcal{M}},$$

where $\mathbb{1}_{[t_n, t_{n+1})}$ is the characteristic function of the interval $[t_n, t_{n+1})$ and the space functions $\boldsymbol{u}^n$ and $h^n$ take the form defined in the previous section. We propose the following decoupled discretisation of the system (1), written in compact form, with the various discrete operators

defined below.

**Initialisation:** $\quad u^0 = \mathcal{P}_\mathcal{E} u_0,\ h^0 = \mathcal{P}_\mathcal{M} h_0,\ p^0 = \frac{1}{2} g(h^0)^2.$ $\qquad$ (7a)

**Iteration** $n,\ 0 \le n \le N-1 :$ solve for $u^{n+1} \in \boldsymbol{H}_{\mathcal{E},0}, h^{n+1} \in L_\mathcal{M}$ and $p^{n+1} \in L_\mathcal{M} :$

$$\eth_t h^{n+1} + \mathrm{div}_\mathcal{M}(h^n u^n) = 0, \qquad (7b)$$

$$p^{n+1} = \frac{1}{2} g(h^{n+1})^2, \qquad (7c)$$

$$\eth_t(h u)^{n+1} + \boldsymbol{C}_\mathcal{E}(h^n u^n) u^n + \nabla_\mathcal{E} p^{n+1} + g\, \boldsymbol{I}_\mathcal{E} h^{n+1}\, \nabla_\mathcal{E} z = 0, \qquad (7d)$$

*Projection operators* - The operators $\mathcal{P}_\mathcal{E}$ and $\mathcal{P}_\mathcal{M}$ used in the initialisation step are defined by $\mathcal{P}_\mathcal{E} = (\mathcal{P}_{\mathcal{E}^{(i)}})_{i=1,\cdots,d}$ with

$$\mathcal{P}_{\mathcal{E}^{(i)}} : \quad \left|\ \begin{aligned} & L^1(\Omega) \longrightarrow H_{\mathcal{E}^{(i)},0} \\ & v \longmapsto \mathcal{P}_{\mathcal{E}^{(i)}} v = \sum_{\sigma \in \mathcal{E}_{\mathrm{int}}^{(i)}} v_\sigma\, \mathbb{1}_{D_\sigma}\ \text{with}\ v_\sigma = \frac{1}{|D_\sigma|} \int_{D_\sigma} v(\boldsymbol{x})\, \mathrm{d}\boldsymbol{x},\ \text{for}\ \sigma \in \mathcal{E}_{\mathrm{int}}^{(i)}. \end{aligned} \right. \quad (8)$$

For $q \in L^2(\Omega), \mathcal{P}_\mathcal{M} q \in L_\mathcal{M}$ is defined by:

$$\mathcal{P}_\mathcal{M} q = \sum_{K \in \mathcal{M}} q_K \mathbb{1}_K\ \text{with}\ q_K = \frac{1}{|K|} \int_K q(\boldsymbol{x})\, \mathrm{d}\boldsymbol{x}\ \text{for}\ K \in \mathcal{M}. \qquad (9)$$

*Discrete time derivative* - The symbol $\eth_t$ denotes the discrete time derivative for both water height and momentum:

$$\eth_t h^{n+1} = \sum_{K \in \mathcal{M}} \eth_t h_K^{n+1} \mathbb{1}_K,\ \eth_t h_K^{n+1} = \frac{1}{\delta t}(h_K^{n+1} - h_K^n),\ \text{and}\ \eth_t(h u)^{n+1} = (\eth_t(h u_1)^{n+1}, \eth_t(h u_2)^{n+1})$$

$$\text{with}\ \eth_t(h u_i)^{n+1} = \sum_{\sigma \in \mathcal{E}^{(i)}} \eth_t(h u_i)_\sigma^{n+1} \mathbb{1}_{D_\sigma},\ \text{and}\ \eth_t(h u_i)_\sigma^{n+1} = \frac{1}{\delta t}(h_{D_\sigma}^{n+1} u_{i,\sigma}^{n+1} - h_{D_\sigma}^n u_{i,\sigma}^n),\ i = 1, 2,$$

where $h_{D_\sigma}$ is the discrete water height in the dual cell, which is computed from the primal unknowns $(h_K^n)_{n \in \mathbb{N}, K \in \mathcal{M}}$ and defined so as to satisfy a discrete mass balance, see below.

*Discrete divergence and gradient operators* - The discrete divergence operator $\mathrm{div}_\mathcal{M}$ is defined by:

$$\mathrm{div}_\mathcal{M} : \quad \left|\ \begin{aligned} & \boldsymbol{H}_{\mathcal{E},0} \longrightarrow L_{\mathcal{M},0} \\ & u \longmapsto \mathrm{div}_\mathcal{M}(h u) = \sum_{K \in \mathcal{M}} \mathrm{div}_K(h u) \mathbb{1}_K,\ \text{with}\ \mathrm{div}_K(h u) = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}, \end{aligned} \right. \quad (10)$$

where $F_{K,\sigma}$ is the (conservative) numerical mass flux, defined by $F_{K,\sigma} = |\sigma|\, h_\sigma u_{K,\sigma}$ with $u_{K,\sigma} = u_{i,\sigma} \boldsymbol{n}_{K,\sigma} \cdot \boldsymbol{e}^{(i)}$ for $\sigma \in \mathcal{E}_{\mathrm{int}}^{(i)}, i = 1, 2$, while $h_\sigma$ is approximated by the first order upwind scheme, namely, for $\sigma = K|L \in \mathcal{E}_{int}, h_\sigma = h_K$ if $u_{K,\sigma} \ge 0$ and $h_\sigma = h_L$ otherwise.

The discrete gradient operator applies to the pressure and the topography and is defined by:

$$\nabla_{\mathcal{E}} : \quad \begin{vmatrix} L_{\mathcal{M}} \longrightarrow \boldsymbol{H}_{\mathcal{E},0} \\[4pt] p \longmapsto \nabla_{\mathcal{E}} p, \end{vmatrix}$$

with for $i = 1, 2$:

$$(\nabla_{\mathcal{E}} p)_i = \sum_{\sigma \in \mathcal{E}_{int}^{(i)}} \eth_\sigma p \, \mathbb{1}_{D_\sigma} \text{ with for } \sigma = \overrightarrow{K|L}, \ \eth_\sigma p = \frac{|\sigma|}{|D_\sigma|} (p_L - p_K). \tag{11}$$

The above defined discrete divergence and gradient operators satisfy the following div-grad duality relationship [7, Lemma 2.5]:

$$\text{for } p \in L_{\mathcal{M}}, \ \boldsymbol{u} \in \boldsymbol{H}_{\mathcal{E},0}, \quad \int_\Omega p \operatorname{div}_{\mathcal{M}}(\boldsymbol{u}) \, \mathrm{d}\boldsymbol{x} + \int_\Omega \nabla_{\mathcal{E}} p \cdot \boldsymbol{u} \, \mathrm{d}\boldsymbol{x} = 0.$$

*Discrete convection operator* – The discrete nonlinear convection operator $\boldsymbol{C}_{\mathcal{E}}(h\boldsymbol{u})$ is linked to the discrete divergence operator on the dual mesh by the relation $\boldsymbol{C}_{\mathcal{E}}(h\boldsymbol{u})\boldsymbol{u} = \operatorname{div}_{\mathcal{E}}(h\boldsymbol{u} \otimes \boldsymbol{u})$, where the full discrete convection operator $\boldsymbol{C}_{\mathcal{E}}(h\boldsymbol{u})$ is defined by:

$$\boldsymbol{C}_{\mathcal{E}}(h\boldsymbol{u})\,\boldsymbol{u} = (C_{\mathcal{E}^{(1)}}(h\boldsymbol{u})\,u_1, C_{\mathcal{E}^{(2)}}(h\boldsymbol{u})\,u_2),$$

and the $i$-th component $C_{\mathcal{E}^{(i)}}(h\boldsymbol{u})$ of the convection operator is defined by:

$$C_{\mathcal{E}^{(i)}}(h\boldsymbol{u}) : \quad \begin{vmatrix} H_{\mathcal{E}^{(i)},0} \longrightarrow H_{\mathcal{E}^{(i)},0} \\[4pt] u_i \longmapsto C_{\mathcal{E}^{(i)}}(h\boldsymbol{u})\,u_i = \sum_{\sigma \in \mathcal{E}_{int}^{(i)}} \operatorname{div}_{\mathcal{E}^{(i)}}(hu_i\boldsymbol{u}) \ \mathbb{1}_{D_\sigma}, \\[10pt] \text{with } \operatorname{div}_{\mathcal{E}^{(i)}}(hu_i\boldsymbol{u}) = \frac{1}{|D_\sigma|} \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon} u_{i,\epsilon}, \end{vmatrix} \tag{12}$$

where for $\epsilon = \sigma|\sigma'$, $u_{i,\epsilon}$ is the upwind choice between $u_\sigma$ and $u_{\sigma'}$ with respect to the sign of $F_{\sigma,\epsilon}$. The quantity $F_{\sigma,\epsilon}$ is the numerical mass flux through $\epsilon$ outward $D_\sigma$; it must be chosen carefully to ensure some stability properties of the scheme as in [7, 11]. Indeed we recall that in order to derive a discrete kinetic energy balance (Lemma 3 below), it is necessary that a discrete equation of the mass balance holds in the dual mesh, namely:

$$\frac{|D_\sigma|}{\delta t}(h_{D_\sigma}^{n+1} - h_{D_\sigma}^n) + \operatorname{div}_{\mathcal{E}}(h^n \boldsymbol{u}^n) = 0, \quad \text{with} \quad |D_\sigma| \operatorname{div}_{\mathcal{E}}(h^n \boldsymbol{u}^n) = \sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon}^n. \tag{13}$$

The water height $h_{D_\sigma}$ and the flux $F_{\sigma,\epsilon}$ are computed from the primal unknowns and fluxes so as to satisfy this latter relation thanks to the discrete mass balance on the primal mesh (7b). For $\sigma = K|L \in \mathcal{E}_{int}$, the water height $h_{D_\sigma}$ is defined as a weighted average between $h_K$ and $h_L$:

$$|D_\sigma|\, h_{D_\sigma} = |D_{K,\sigma}|\, h_K + |D_{L,\sigma}|\, h_L, \tag{14}$$

where $D_\sigma$, $D_{K,\sigma}$ and $D_{L,\sigma}$ are defined in Definition 1. The numerical flux $F_{\sigma,\epsilon}$ on the internal dual edges, is defined according to the location of the edges as follows:
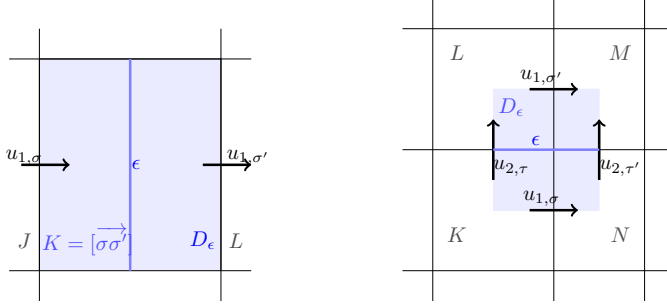
Figure 2: Notations for the definition of the momentum flux on the dual mesh for the first component of the velocity- left: first case - right: second case.

- First case – The vector $\boldsymbol{e}^{(i)}$ is normal to $\epsilon$, and $\epsilon$ is included in a primal cell $K$, with $K = [\overrightarrow{\sigma\sigma'}]$ (see Definition 1 and Figure 2 on the left for $i = 1$). Then for a dual edge $\epsilon \in \widetilde{\mathcal{E}}^{(i)}$ such that $\epsilon = \overrightarrow{\sigma|\sigma'}$, the flux $F_{\sigma,\epsilon}$ through the edge $\epsilon$ is given by:

$$F_{\sigma,\epsilon} = \frac{1}{2}(F_{K,\sigma'} - F_{K,\sigma}) = \frac{1}{2}|\epsilon|\,(h_\sigma u_{i,\sigma} + h_{\sigma'} u_{i,\sigma'}), \tag{15}$$

  since $|\sigma| = |\sigma'| = |\epsilon|$.

- Second case – The vector $\boldsymbol{e}^{(i)}$ is tangent to $\epsilon$, and $\epsilon$ is the union of the halves of two primal edges $\tau$ and $\tau'$ such that $\tau = \overrightarrow{K|L}$, $\tau \in \mathcal{E}(K)$ and $\tau' = \overrightarrow{N|M} \in \mathcal{E}(N)$ (see Definition 1 and Figure 2 on the right for $i = 1$). Let $j \in \{1, 2\}$, $j \neq i$: the numerical flux through $\epsilon$ is then given by:

$$F_{\sigma,\epsilon} = \frac{1}{2}\,(F_{K\tau} + F_{L\tau'}) = \frac{1}{2}\,(|\tau|\,h_\tau u_{j,\tau} + |\tau'|\,h_{\tau'} u_{j,\tau'}). \tag{16}$$

Note that the numerical momentum flux on a dual edge is conservative. It is easy to check that the unknowns $h_{D_\sigma}^n$ and $F_{\sigma,\epsilon}^n$ thus defined satisfy the discrete dual mass balance (13).

*Discrete water height on the dual mesh, for the topography term* – In equation (7d) the interpolation operator $\boldsymbol{I}_{\mathcal{E}}$ is defined as the mean value of the water height:

$$\boldsymbol{I}_{\mathcal{E}}h = \sum_{\sigma \in \mathcal{E}_{int}} h_{\sigma,c}\,\mathbb{1}_{D_\sigma} \text{ with } h_{\sigma,c} = \begin{cases} \frac{1}{2}(h_K + h_L) \text{ for } \sigma = K|L \in \mathcal{E}_{int}, \\ h_K \text{ for } \sigma \in \mathcal{E}_{ext} \cap \mathcal{E}(K). \end{cases} \tag{17}$$

This choice is important to preserve steady states, see Lemma 2.

## §4. Properties of the scheme

The scheme (7) enjoys some interesting properties, which we now state. First of all, thanks to the upwind choice for $h^n$ in (1a), the positivity of the water height is preserved under a CFL like condition.

**Lemma 1** (Positivity of the water height). *Let $n \in [\![0, N-1]\!]$, let $(h_K^n, u_{i,\sigma}^n)_{K \in \mathcal{M}, \sigma \in \mathcal{E}^{(i)}}$ be given and such that $h_K^n \geq 0$, for all $K \in \mathcal{M}$, and let $h_K^{n+1}$ be computed by (7b). Then $h_K^{n+1} \geq 0$, for all $K \in \mathcal{M}$ under the following CFL condition,*

$$\forall K \in \mathcal{M}, \delta t \leq \frac{|K|}{\displaystyle\sum_{\sigma \in \mathcal{E}(K)} |\sigma| \, |\boldsymbol{u}_{K,\sigma}^n|}. \tag{18}$$

Second, thanks to the choice (17) for the reconstruction of the water height, the "lake at rest" steady state is preserved by the scheme.

**Lemma 2** (Steady state "lake at rest"). *Let $n \in [\![0, N-1]\!]$, $C \in \mathbb{R}_+$; let $\boldsymbol{u}^{n+1} \in \boldsymbol{H}_{\mathcal{E},0}$ and $h^{n+1} \in L_{\mathcal{M}}$ be a solution to (7b)-(7d) with $\boldsymbol{u}^n = 0$ and $h^n + z = C$, where $C$ is a given real number. Then $\boldsymbol{u}^{n+1} = 0$ and $h^{n+1} + z = C$.*

As a consequence of the careful discretisation of the convection term, the scheme satisfies a discrete kinetic energy balance, as stated in the following lemma. The proof of this result is an easy adaptation of [10, Lemma 3.2].

**Lemma 3** (Discrete kinetic energy balance). *A solution to the scheme (7) satisfies the following equality, for $i = 1, 2$, $\sigma \in \mathcal{E}^{(i)}$ and $0 \leq n \leq N - 1$:*

$$\frac{1}{2\,\delta t}(h_{D_\sigma}^{n+1}(u_{i,\sigma}^{n+1})^2 - h_{D_\sigma}^n(u_{i,\sigma}^n)^2) + \frac{1}{2\,|D_\sigma|}\sum_{\epsilon \in \mathcal{E}^{(i)}(D_\sigma)} F_{\sigma,\epsilon}^n(u_{i,\epsilon}^n)^2$$
$$+ u_{i,\sigma}^{n+1}\eth_\sigma p^{n+1} + g\,h_{\sigma,c}^{n+1}\,u_{i,\sigma}^{n+1}\eth_\sigma z = -R_{i,\sigma}^{n+1}, \quad (19)$$

*with $R_{i,\sigma}^{n+1} \geq 0$ under the CFL like restriction:*

$$\forall \sigma \in \mathcal{E}^{(i)}, \qquad \delta t \leq \frac{|D_\sigma| \, h_{D_\sigma}^{n+1}}{\displaystyle\sum_{\epsilon \in \widetilde{\mathcal{E}}(D_\sigma)} (F_{\sigma,\epsilon}^n)^-}. \tag{20}$$

The scheme also satisfies the following potential energy balance [10, Lemma 3.3].

**Lemma 4** (Discrete elastic potential balance). *Let, for $K \in \mathcal{M}$ and $0 \leq n \leq N$ the potential energy be defined by $(E_p)_K^n = \frac{1}{2}g\,(h_K^n)^2$. A solution to the scheme (7) satisfies the following equality, for $K \in \mathcal{M}$ and $0 \leq n \leq N - 1$:*

$$(\eth_t E_p)_K^{n+1} + \text{div}_K(E_p^n \boldsymbol{u}^n) + p_K^n \text{div}_K \boldsymbol{u}^n = -R_K^{n+1}, \tag{21}$$

*with*

$$R_K^{n+1} \geq \frac{1}{|K|}g\sum_{\sigma \in \mathcal{E}(K)} |\sigma| \, u_{K,\sigma}^n h_\sigma^n (h_K^{n+1} - h_K^n). \tag{22}$$

Note that the right-hand side of Equation (22) may be negative, and thus also the quantities $R_K^{n+1}$. This is specific to explicit schemes (for implicit or pressure-correction schemes, this residual is non-negative [9]) and prevents getting a stability estimate for the scheme.

However, combining the two previous lemmas allows to prove that convergent sequences of solutions to the scheme satisfy an entropy inequality, as depicted in the next section. To this purpose, we will pass to the limit in a discrete entropy balance which is built as follows. Let $K \in \mathcal{M}$ and let us denote by $(E_k)_K^n$ the following quantity, which may be seen as a kinetic energy associated to $K$:

$$(E_k)_K^n = \frac{1}{4\,|K|} \sum_{i=1}^{2} \sum_{\sigma \in \mathcal{E}(K) \cap \mathcal{E}^{(i)}} |D_\sigma|\, h_{D_\sigma}^n (u_{i,\sigma}^n)^2.$$

Then, for $\sigma_0 \in \mathcal{E}(K)$, we define a kinetic energy flux, which we denote by $G_{K,\sigma_0}^n$, as follows. Let us suppose, for instance, that $\sigma_0 \in \mathcal{E}^{(1)}$. We denote by $\epsilon$ the face of $D_{\sigma_0}$ parallel to $\sigma_0$ and included in $K$ and by $\epsilon'$ the opposite face of $D_{\sigma_0}$. In addition, $\sigma_0$ is the union of two half-faces of the dual mesh associated to the second component of the velocity, which we denote by $\tau$ and $\tau'$, and we denote by $\sigma$ and $\sigma'$ the two faces of $K$ belonging to $\mathcal{E}^{(2)}$ such that $\tau \in \widetilde{\mathcal{E}}(D_\sigma)$ and $\tau' \in \widetilde{\mathcal{E}}(D_{\sigma'})$. We then set:

$$G_{K,\sigma_0}^n = \frac{1}{4}\Big[-F_{\sigma_0,\epsilon}^n (u_{1,\epsilon}^n)^2 + F_{\sigma_0,\epsilon'}^n (u_{1,\epsilon'}^n)^2 + F_{\sigma,\tau}^n (u_{2,\tau}^n)^2 + F_{\sigma,\tau'}^n (u_{2,\tau'}^n)^2\Big]$$

Multiplying the kinetic energy balance equation (19) associated to each face $\sigma$ of $K$ by $\frac{1}{2}\,|D_\sigma|$ and summing the four obtained relations with Equation (21) multiplied by $|K|$, we get

$$(\delta_t E_k)_K^{n+1} + (\delta_t E_p)_K^{n+1} + \sum_{\sigma \in \mathcal{E}(K)} [G_{K,\sigma}^n + \frac{1}{2} g h_\sigma^n F_{K,\sigma}^n] + \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |\sigma| p_K^n u_{K,\sigma}^n$$

$$+ \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |\sigma| \frac{1}{2}\, (p_L^{n+1} - p_K^{n+1})\, u_{K,\sigma}^{n+1} + \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |\sigma| \frac{1}{4} g\, (h_K^{n+1} + h_L^{n+1})\, (z_L - z_K)\, u_{K,\sigma}^{n+1} = -T_K^{n+1},$$

where $T_K^{n+1}$ collects the residual terms in (19) and (21), and thus $T_K^{n+1} \geq R_K^{n+1}$. We now remark that, thanks to the discrete mass balance equation and the fact that the topography does not depend on time,

$$\frac{1}{2} \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}^n (z_L - z_K) = \frac{|K|}{\delta t}\, (h_K^{n+1} z_K - h_K^n z_K) + \frac{1}{2} \sum_{\sigma \in \mathcal{E}(K),\, \sigma=K|L} F_{K,\sigma}^n\, (z_K + z_L),$$

and we finally obtain the following discrete entropy balance:

$$(\delta_t E_k)_K^{n+1} + (\delta_t E_p)_K^{n+1} + g z_K (\delta_t h)_K^{n+1} + \sum_{\sigma \in \mathcal{E}(K)} [G_{K,\sigma}^n + \frac{1}{2} g h_\sigma^n F_{K,\sigma}^n + \frac{1}{2} F_{K,\sigma}^n (z_K + z_L)]$$

$$+ \sum_{\substack{\sigma \in \mathcal{E}(K) \\ \sigma = K|L}} |\sigma| \frac{1}{2}\, (p_K^n + p_L^n)\, u_{K,\sigma}^{n+1} = -(R_e)_K^{n+1}, \quad (23)$$

with

$$
(R_e)_K^{n+1} \geq T_K^{n+1} + g \sum_{\sigma \in \mathcal{E}(K)} \left[ \frac{1}{2} F_{K,\sigma}^n - \frac{1}{4}|\sigma| \, (h_K^{n+1} + h_L^{n+1}) \, u_{K,\sigma}^{n+1} \right] (z_L - z_K)
$$

$$
+ \sum_{\sigma \in \mathcal{E}(K), \, \sigma = K|L} |\sigma| \, \frac{1}{2} \, (p_K^{n+1} u_{K,\sigma}^{n+1} - p_K^n u_{K,\sigma}^n + p_L^{n+1} u_{L,\sigma}^{n+1} - p_L^n u_{L,\sigma}^n). \quad (24)
$$

Note that the above remainder term $R_e$ cannot be proven to be non negative; however, the relation (24) allows to obtain a weak entropy consistency property when passing to the limit, as stated in the next section.

## §5. Consistency analysis

The objective of this section is to show that the schemes are consistent in the Lax-Wendroff sense, namely 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.

A weak solution to the continuous problem satisfies, for any $\varphi \in C_c^\infty(\Omega \times [0, T))$ ($\boldsymbol{\varphi} \in C_c^\infty(\Omega \times [0, T))^2$):

$$
\int_0^T \int_\Omega \left[ h \, \partial_t \varphi + h \, \boldsymbol{u} \cdot \nabla \varphi \right] \mathrm{d}\boldsymbol{x} \, \mathrm{d}t + \int_\Omega h_0(\boldsymbol{x}) \, \varphi(\boldsymbol{x}, 0) \, \mathrm{d}\boldsymbol{x} = 0, \quad (25\mathrm{a})
$$

$$
- \int_0^T \int_\Omega \left[ h \, \boldsymbol{u} \cdot \partial_t \boldsymbol{\varphi} + (h\boldsymbol{u} \otimes \boldsymbol{u}) : \boldsymbol{\varphi} + \frac{1}{2} \, g \, h^2 \mathrm{div}(\boldsymbol{\varphi}) + g \, h \, \nabla(z) \boldsymbol{\varphi} \right] \mathrm{d}\boldsymbol{x} \, \mathrm{d}t \quad (25\mathrm{b})
$$

$$
- \int_\Omega h_0(\boldsymbol{x}) \, \boldsymbol{u}_0(\boldsymbol{x}) \cdot \varphi(\boldsymbol{x}, 0) \, \mathrm{d}\boldsymbol{x} = 0.
$$

This system is supplemented with a weak entropy inequality, for any nonnegative test functions $\varphi \in C_c^\infty(\Omega \times [0, T), \mathbb{R}_+)$ :

$$
- \int_0^T \int_\Omega \left[ \eta \, \partial_t \varphi + \boldsymbol{\Phi} \cdot \nabla \varphi \right] \mathrm{d}\boldsymbol{x} \, \mathrm{d}\mathbf{t} - \int_\Omega \eta_0(\boldsymbol{x}) \, \varphi(\boldsymbol{x}, \mathbf{0}) \, \mathrm{d}\boldsymbol{x} \leq \mathbf{0}, \quad (26)
$$

with $\eta$ and $\boldsymbol{\Phi}$ defined by (4).

Before stating the global weak consistency of the scheme (7), some definitions and estimate assumptions are needed.

Let $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$ be a sequence of meshes in the sense of Definition 1 and let $(h^{(m)}, \boldsymbol{u}^{(m)})_{m \in \mathbb{N}}$ be the associated sequence of solutions of the scheme (7)).

**Assumed estimates** - We need also some a priori estimates on the sequence of discrete solutions $(h^{(m)}, \boldsymbol{u}^{(m)})_{m \in \mathbb{N}}$ in order to prove the consistency result we are seeking. First of all we assume that $h^{(m)} > 0$, $\forall m \in \mathbb{N}$ which can be obtained under the CFL condition (18). Furthermore:

– The water height $h^{(m)}$ and its inverse are uniformly bounded in $L^\infty(\Omega \times (0,T))$, *i.e.* there exists some constants $C, C' \in \mathbb{R}_+^*$ such that for $m \in \mathbb{N}$ and $0 \le n < N^{(m)}$:

$$1/C < (h^{(m)})_K^n \le C, \quad 1/C' < 1/(h^{(m)})_K^n \le C' \quad \forall K \in \mathcal{M}^{(m)} \tag{27}$$

– The velocity $\boldsymbol{u}^{(m)}$ is also uniformly bounded in $L^\infty(\Omega \times (0,T))^2$:

$$|(\boldsymbol{u}^{(m)})_\sigma^n| \le C, \quad \forall \sigma \in \mathcal{E}^{(m)}. \tag{28}$$

Finally, the weak consistency to the entropy inequality is only proved under additional assumptions. First we need the following condition on the space and time steps, which is stronger than a CFL condition:

$$\frac{\delta t^{(m)}}{\delta_{\mathcal{M}^{(m)}}} \to 0 \text{ as } m \to +\infty \tag{29}$$

Second, the $L^1(\Omega, BV)$ norm of the height is required to be bounded, *i.e.* there exists one constant $C$ such that, for $m \in \mathbb{N}$,

$$\sum_{n=0}^{N-1} \sum_{K \in \mathcal{M}} |K| |(h^{(m)})_K^{n+1} - (h^{(m)})_K^n| \le C. \tag{30}$$

We are now in position to state the following consistency result.

**Theorem 5** (Weak consistency of the scheme). *Let $(\mathcal{M}^{(m)}, \mathcal{E}^{(m)})_{m \in \mathbb{N}}$ be a sequence of meshes such that $\delta t^{(m)}$ and $\delta_{\mathcal{M}^{(m)}} \to 0$ as $m \to +\infty$; assume that there exists $\eta > 0$ such that $\eta_{\mathcal{M}^{(m)}} \le \eta$ for any $m \in \mathbb{N}$ (with $\eta_{\mathcal{M}^{(m)}}$ defined by (5)); assume moreover that (27) and (28) hold. Let $(h^{(m)}, \boldsymbol{u}^{(m)})_{m \in \mathbb{N}}$ be a sequence of solutions to the scheme (7) converging to $(\bar{h}, \bar{\boldsymbol{u}})$ in $L^1(\Omega \times (0,T)) \times L^1(\Omega \times (0,T))^2$. Then $(\bar{h}, \bar{\boldsymbol{u}})$ satisfies the weak formulation (25) of the shallow water equations.*

*If we furthermore assume the space and time steps satisfy (29) and that the sequence of heights is uniformly bounded in $L^1(\Omega, BV)$, i.e. satisfy (30), then $(\bar{h}, \bar{\boldsymbol{u}})$ satisfies the entropy inequality (26).*

*Proof.* The proof is obtained by passing to the limit in the scheme and in the discrete entropy balance (23), using the tool of [6] (or, more precisely speaking, simplified versions of these tools adapted to Cartesian grids). The additional assumptions required for the entropy condition are used to prove that the residual term appearing in the discrete potential energy balance, given by (22), tends to zero. □

## §6. Numerical results

We now assess the behaviour of the scheme on some numerical experiments. The computations presented here are performed with the CALIF³S free software developed at IRSN [4].
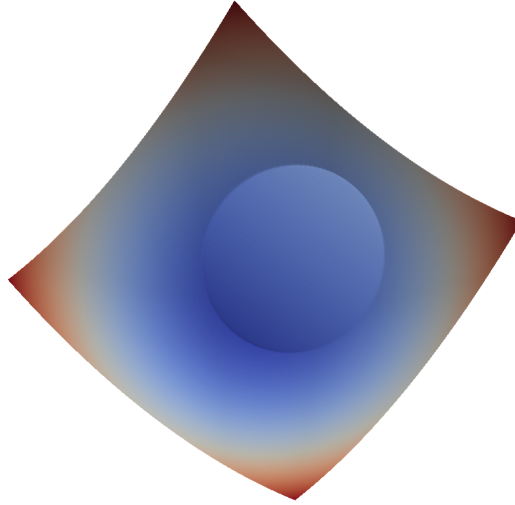
Figure 3: Sloshing of a drop on a parabolic support – State obtained after one revolution (very close to the initial state).

## 6.1. Rotation in a paraboloid

This first test case consists in calculating the uniform rotation of a circular drop on a support of parabolic shape (see Figure 3). The computational domain is $(0, L) \times (0, L)$ and the elevation of the support is:

$$z = -h_0 \left(1 - (x - \frac{L}{2})^2 - (y - \frac{L}{2})^2\right),$$

with $L = 4$ and $h_0 = 0.1$. The fluid height is given by

$$h = h_0 \; \max(0, \; (x - \frac{L}{2}) \cos(\omega t) + (y - \frac{L}{2}) \sin(\omega t) - z - 0.5),$$

and the velocity is

$$\boldsymbol{u} = \frac{1}{2} \omega \begin{bmatrix} -\sin(\omega t) \\ \cos(\omega t) \end{bmatrix}.$$

It is then easy to check that the mass and momentum balance equations are verified provided that $\omega^2 = 2 g h_0$. The solution is thus regular, and this test features a regular topography and dry zones (*i.e.* zones where $h = 0$). We compare the numerical and theoretical height obtained after one rotation (*i.e.* $\omega t = 2\pi$), for different uniform grids and with a time step $\delta t = \delta x / 8$. (the maximal speed of sound and the maximal velocity are both close to 1); results are gathered in the following table:

| grid | error (discrete $L^1$ norm) |
|---|---|
| $100 \times 100$ | $3.02 \ 10^{-3}$ |
| $200 \times 200$ | $1.54 \ 10^{-3}$ |
| $400 \times 400$ | $0.896 \ 10^{-3}$ |
| $800 \times 800$ | $0.511 \ 10^{-3}$ |

We observe an order of convergence between 0.8 and 1, which is consistent with a first-order approximation of the fluxes and the time derivative.

## 6.2. A dam-break problem

In this test, the computational domain is:

$$\Omega = (0, 200) \times (0, 200) \setminus \Omega_w \text{ with } \Omega_w = (95, 105) \times (0, 95) \cup (95, 170) \times (0, 200).$$

The fluid is supposed to be initially at rest, and the initial height is $h = 10$ for $x_1 \leq 100$ and $h = 5$ for $x_1 > 100$. 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 $\delta t = \delta x / 25$ (the maximal speed of sound and the maximal velocity are both close to 10). The obtained fluid height is shown at different times on Figure 4; they confirm the efficiency of the scheme, and its capability to deal with reflexion phenomena very simply (*i.e.* just by setting the normal velocity at the boundary to zero, by contrast with schemes based on Riemann solvers which need to implement fictitious cells techniques).

## Acknowledgements

## References

[1] ARAKAWA, A., AND LAMB, V. A potential enstrophy and energy conserving scheme for the shallow water equations. *Monthly Weather Review 109* (1981), 18–36.

[2] BONAVENTURA, L., AND RINGLER, T. Analysis of discrete shallow-water models on geodesic delaunay grids with c-type staggering. *Monthly Weather Review 133*, 8 (2005), 2351–2373.

[3] BOUCHUT, F. *Nonlinear Stability of finite volume methods for hyperbolic conservation laws*. Birkhauser, 2004.

[4] CALIF$^3$S. A software components library for the computation of fluid flows. `https://gforge.irsn.fr/gf/project/califs`.

[5] DOYEN, D., AND GUNAWAN, H. An explicit staggered finite volume scheme for the shallow water equations. In *Finite volumes for complex applications. VII. Methods and theoretical aspects*, vol. 77 of *Springer Proc. Math. Stat.* Springer, Cham, 2014, pp. 227–235.
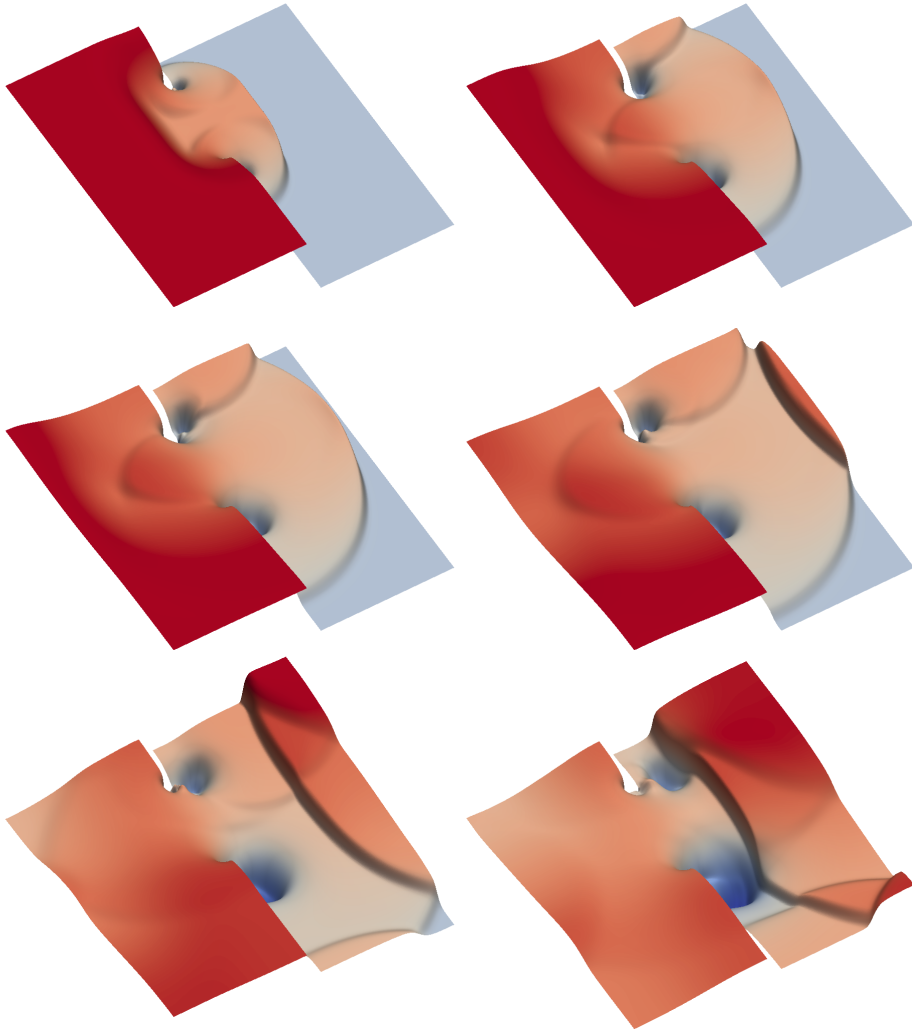
Figure 4: Partial dam break – Height obtained at $t = 4$, $t = 8$, $t = 10$, $t = 12$, $t = 16$ and $t = 20$ with a mesh obtained (by supression of the zones associated to the obstacles) from a $1000 \times 1000$ regular grid. In the last Figure ($t = 20$), the obtained minimal and maximal heights are $h = 2.149$ and $h = 9.306$ respectively.

[6] Gallouët, T., Herbin, R., and Latché. On the weak consistency of finite volumes schemes for conservation laws on general meshes. *under revision* (2019). Available from: `https://hal.archives-ouvertes.fr/hal-02055794`.

[7] Gallouët, T., Herbin, R., Latché, J.-C., and Mallem, K. Convergence of the marker-and-cell scheme for the incompressible Navier-Stokes equations on non-uniform grids. *Foundations of Computational Mathematics 18* (2018), 249–289.

[8] Gunawan, H. *Numerical simulation of shallow water equations and related models*. PhD thesis, Université Paris-Est and Institut Teknologi Bandung, 2015.

[9] Herbin, R., Kheriji, W., and Latché, J.-C. On some implicit and semi-implicit staggered schemes for the shallow water and Euler equations. *ESAIM: Mathematical Modelling and Numerical Analysis 48* (2014), 1807–1857.

[10] Herbin, R., Latché, J.-C., and Nguyen, T. Explicit staggered schemes for the compressible Euler equations. *ESAIM: Proceedings 40* (2013), 83–102.

[11] Herbin, R., Latché, J.-C., and Nguyen, T. Consistent segregated staggered schemes with explicit steps for the isentropic and full Euler equations. *ESAIM: Mathematical Modelling and Numerical Analysis 52* (2018), 893–944.

[12] Stelling, G., and Duinmeijer, S. A staggered conservative scheme for every Froude number in rapidly varied shallow water flows. *International Journal for Numerical Methods in Fluids 43* (2003), 1329–1354.

R. Herbin and Y. Nasseri
Aix-Marseille Université, Institut de Mathématiques de Marseille,
39 rue Joliot Curie
13453 Marseille
`raphaele.herbin@univ-amu.fr` and `youssouf.nasseri@univ-amu.fr`

J.-C. Latché                                          N. Therme
Institut de Radioprotection et Sûreté Nucléaire,     CEA/CESTA
13115, Saint-Paul-lez-Durance                        33116, Le Barp, France
`jean-claude.latche@irsn.fr`                          `nicolas.therme@cea.fr`