

# An explicit residual based approach for shallow water flows

Mario Ricchiuto\*

\* Inria Bordeaux - Sud-Ouest,  
200 Avenue de la Vieille Tour, 33405 Talence cedex - France

## Abstract

We describe fully explicit residual based discretizations of the shallow water equations with friction on unstructured grids. The schemes are obtained by properly adapting the explicit construction proposed in (Ricchiuto and Abgrall, *J.Comput.Phys.* 229, 2010). In particular, previous work on well balanced integration (Ricchiuto, *J.Sci.Comp.* 48, 2011) and preservation of the depth non-negativity (Ricchiuto and Bollermann, *J.Comput.Phys.* 228, 2009) is reformulated in the context of a genuinely explicit time stepping still based on a weighted residual approximation. The paper discusses in depth how to achieve in this context an exact preservation of all the simple known steady equilibria, and how to obtain a super-consistent approximation for smooth non-trivial moving equilibria. The treatment of the wetting/drying interface is also discussed, giving formal conditions for the preservation of the non-negativity of the depth for a particular case, based on a nonlinear variant of a Lax-Friedrichs type scheme. The approach is analyzed and tested thoroughly. The quality of the numerical results shows the interest in the proposed approach over previously proposed schemes, in terms of accuracy and efficiency.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>The shallow water equations</b>	<b>4</b>
2.1	Lake at rest solution . . . . .	6
2.2	Constant energy pseudo-1d flows . . . . .	6
2.3	Steady flows in sloping channels . . . . .	7
<b>3</b>	<b>Explicit residual approach for conservation laws</b>	<b>7</b>
3.1	Generalities . . . . .	7
3.2	The predictor-corrector explicit scheme . . . . .	9
3.3	Basic properties . . . . .	10
<b>4</b>	<b>Application to the shallow water equations</b>	<b>13</b>
4.1	C-property and super-consistency analysis . . . . .	13
4.2	C-property : application to particular steady states . . . . .	16
4.2.1	The lake at rest solution . . . . .	16

4.2.2	Constant energy pseudo-1d flows . . . . .	17
4.2.3	Steady flows in sloping channels . . . . .	17
4.3	Nonlinear Lax-Friedrichs distribution . . . . .	18
4.4	Filtering and streamline dissipation . . . . .	20
4.5	Wet/dry front handling and implementation details . . . . .	21
<b>5</b>	<b>Numerical tests</b>	<b>23</b>
5.1	Flows on flat bathymetry . . . . .	24
5.1.1	Vortex transport, accuracy and efficiency . . . . .	24
5.1.2	Asymmetric break of a dam . . . . .	24
5.2	C-property tests . . . . .	25
5.2.1	Lake at rest solution . . . . .	25
5.2.2	Constant energy flows . . . . .	27
5.2.3	Flows in sloping channels with friction . . . . .	34
5.3	Wetting/drying tests . . . . .	34
5.3.1	Thacker's oscillations in a parabolic bowl . . . . .	34
5.3.2	Runup on a conical island . . . . .	35
5.3.3	Okushiri tsunami experiment . . . . .	37
<b>6</b>	<b>Conclusions</b>	<b>39</b>
<b>A</b>	<b>Proof of proposition 4.2</b>	<b>40</b>
<b>B</b>	<b>Proof of proposition 4.6</b>	<b>51</b>

# 1 Introduction

Free surface flows are relevant in a large number of applications, especially in civil and coastal engineering. The problems concerned are either (relatively) local, such as dam breaks and flooding, overland flows due to rainfall, nearshore wave propagation and interaction with complex bathymetries/structures, and tidal waves in rivers, or global such as in ocean or sea basin models for the study of *e.g.* tsunami generation and propagation.

The simulation of such flows can be carried out by solving directly the three dimensional Navier-Stokes equations. However, for many applications, including *e.g.* nearshore wave propagation and flooding, simplified models obtained by combining vertical averaging and some form of thin layer approximation provide reliable results. The applicability of such models depends on the nature of the flow and on the hypotheses at their basis [15, 44].

The simplest among these models is the so-called shallow water model. The model assumes that the waves developing in the flow are *long* (small ratio amplitude/wavelength), and of a hydrostatic vertical variation of the pressure [35, 48]. More complex nonlinear models can be obtained, by including higher order terms, and depending on the hypotheses on the flow [35, 48, 44, 15]. The first order shallow water approximation constitutes a non-homogeneous hyperbolic system where the effects of the variation of the bathymetry and the viscous friction on the bottom are modeled by the source terms [35, 48].

The amount of literature related to the solution of the shallow water system is vast. This model finds applications in oceanography, hydrology, and meteorology (see *e.g.* [20, 34, 40, 72, 73, 74] and references therein). The main challenges when solving the shallow water system numerically are related to the discretization of the bathymetry and friction terms, and to the numerical treatment of nearly dry regions. For the first issue, one speaks often *asymptotic preserving* character or *well balancedness* of a discretization. The second issue is what is referred to as the wetting/drying strategy.

Well balancing, refers to the ability of the discretization to preserve some steady equilibria involving the existence of a set of invariants exactly, or within some mesh size dependent bounds possibly more favorable than the accuracy of the scheme. The typical example is the so called *lake at rest state* involving a flat still free surface, that should remain flat whatever the shape of the bottom. This property is what one refers to as *Conservation property, or C-property* [14] or well-balancedness [36]. One speaks of approximate C-property when the steady state is kept within an accuracy higher than that of the underlying scheme. This property becomes important when one is interested in flows that, at least locally, are perturbations of one of these steady equilibria, so that numerical perturbations might interfere with the actual flow giving wrong results. There is plenty of literature discussing several different approaches to the preservation of steady equilibria, in particular the so-called lake at rest state. Most of these developments have taken place in the finite volume community, and are thought in terms of one-dimensional flows (see *e.g.* [14, 33, 36, 51] and references therein). The basic approach boils down either to the inclusion of a source term contribution in the FV numerical flux, so that the correct equilibrium is found at the discrete level [14, 36, 41], or to the rewriting of the system in a relaxation form, where an appropriate integral of the source term is added to the physical flux in the Maxwellian on the right hand side [29, 70]. The multidimensional case is often handled by a dimension by dimension extension on structured grids (see [50, 51, 52, 80], for recent examples), or by introducing local pseudo-one dimensional problems along some geometrical directions (*e.g.* normals to grid faces) [12, 28, 41, 49]. These modified FV fluxes are also used in the context of discontinuous Galerkin schemes to retain the C-property (see *e.g.* [32, 78]). A different approach is that of the well balanced wave propagation finite volume schemes of LeVeque and his co-workers [45, 46], the continuous stabilized finite element discretizations proposed by G.Hauke [38], and residual distribution schemes of [17, 58, 60].

On the other hand, the computational treatment of nearly dry areas involves the solution of the following issues : ensuring that in these regions no unphysical negative depths are obtained ; handling some ill-posed problems such as the computation of the local velocity given depth and discharge ; preserving the well balanced character of the method when  $0 < H \ll 1$ .

These three issues are not independent and the large majority of the wetting/drying treatments discussed in literature boil down to : rely on the use of some positivity preserving scheme to be able to keep the depth non-negative ; introducing a cut-off of some sort on the velocity (and mass flux) to avoid zero over zero type divisions ; modify the *numerical* slope of the bathymetry used in the discrete equations ; employ an implicit (split or unsplit) treatment of the friction term to handle the stiffness associated to this term in dry areas. These ideas can be put in practice in various ways, depending on the initial formulation of the method, on the techniques used to reach higher order of accuracy, and on the type of nonlinear mech-

anism used to combine high order and preservation of the positivity. For an overview see [12, 18, 19, 21, 22, 32, 49, 81, 79].

This paper follows the author’s previous work [56, 58, 60, 59] on the construction of residual approximations to the shallow water system. The main objective is to propose a method more efficient than those proposed in the last references, yet retaining all the nice properties of these methods. These properties include the C-property and a generalized C-property for constant energy flows, the preservation of the depth non-negativity, and a robust treatment of moving shorelines. All these properties are achieved on general adaptive unstructured grids. This has a definite advantage as it allows an enhanced resolution of local features, such as e.g. steep variations of the bathymetry leading to complex local flow patterns. More advanced enhancements can be obtained by means of dynamic adaptation for time dependent flows, but this is not considered in this paper. The major limitation of the schemes of [56, 58, 60, 59] is that, while being genuinely implicit and highly nonlinear, they still need to satisfy an explicit type constraint on the allowed time step for the preservation of the depth’s non negativity. Possible routes to overcome this limitation have been suggested in [42, 68], using an unconditionally monotonicity preserving space-time framework, and in [57], where a fully explicit variant of the residual distribution method is proposed. In this paper, we develop the ideas of the last reference, and propose a specific formulation for the shallow water equations. In particular, the main contributions of the present work are : a detailed analysis of the conditions leading to the respect of the C-property for both lake at rest solutions and flows over constant slopes with friction ; a formal characterization of the approximate generalized C-property as referred in [56] (C-property for constant energy flows), here studied in terms of super-consistency with the steady solution both on general grids, and on flow aligned cartesian meshes ; the analysis of the preservation of the non-negativity of the depth for the explicit schemes based on the nonlinear Lax-Friedrichs method of [57] ; an extensive validation and comparison with the implicit scheme of [60]. Few of these results have been presented in [55] (see also the manuscript [54]).

The paper is organized as follows. We recall the form of the shallow system and a number of exact steady equilibria in section §2. The explicit residual discretization approach is then recalled in section §3. Section §4 finally analyzes the properties of the discretization, namely well balancedness (C-properties), accuracy, positivity preservation, and wetting/drying strategy. Lastly, in section §5 we demonstrate the capabilities of the scheme on a large number of numerical tests. Conclusive remarks and future developments end the paper in section §6.

## 2 The shallow water equations

The system of the Nonlinear Shallow Water Equations (NLSW) reads

$$\begin{aligned} \partial_t H + \nabla \cdot (H\vec{v}) + R(x, y, t) &= 0 \\ \partial_t (H\vec{v}) + \nabla \cdot (H\vec{v} \otimes \vec{v} + p(H)\mathbf{I}) + gH(\nabla b + c_f \vec{v}) &= 0' \end{aligned} \tag{1}$$

where (cf. figure 1)  $H$  represents the water depth,  $\vec{v}$  the (vertically averaged) local velocity,  $R$  is a source of mass (*e.g.* associated to rainfall),  $b$  is the bathymetry,  $p(H)$  is given by

$$p(H) = g \frac{H^2}{2},$$

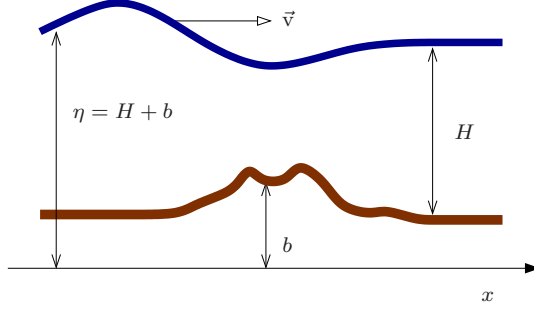


Figure 1: NLSW basic notation

and  $c_f$  is the friction coefficient generally depending on the solution :

$$c_f = c_f(H, \vec{v}). \quad (2)$$

In the following, we will assume that  $R = 0$ , and that the friction coefficient is given by Manning's formula

$$c_f = \frac{n^2 \|\vec{v}\|}{H^{4/3}}, \quad (3)$$

with  $n$  the Manning's coefficient. Introducing the conserved variables  $u$ , conservative fluxes  $\mathcal{F}(u)$  and the source term  $\mathcal{S}$

$$u = \begin{bmatrix} H \\ H\vec{v} \end{bmatrix}, \quad \mathcal{F}(u) = \begin{bmatrix} H\vec{v} \\ H\vec{v} \otimes v + p(H)\mathbf{I} \end{bmatrix}, \quad \mathcal{S}(u, x, y) = gH \begin{bmatrix} 0 \\ \nabla b(x, y) + c_f(u)\vec{v} \end{bmatrix}, \quad (4)$$

system (1) can be recast in the compact form

$$\partial_t u + \nabla \cdot \mathcal{F}(u) + \mathcal{S}(u, x, y) = 0. \quad (5)$$

System (1) is endowed with a mathematical entropy coinciding with the total energy [38, 39, 75, 76], it is hyperbolic, and characterized by the physical constraint of the non-negativity of the depth. Given a direction  $\hat{\xi} \in \mathbb{R}^2$ , with  $\|\hat{\xi}\| = 1$ , and setting for any  $\vec{v} \in \mathbb{R}^2$

$$v_\xi = \vec{v} \cdot \hat{\xi}, \quad (6)$$

the Jacobian matrix

$$K_\xi = \frac{\partial \mathcal{F}_\xi(u)}{\partial u} \quad (7)$$

admits a complete set of real eigenvalues and linearly independent real eigenvectors, with the eigenvalues given by

$$v_\xi - c, \quad v_\xi, \quad v_\xi + c,$$

with  $c$  the celerity

$$c = \sqrt{gH}. \quad (8)$$

It is also useful to introduce the free surface level

$$\eta = H + b, \quad (9)$$

the *specific total energy*

$$\mathcal{E} = g\eta + k, \quad k = \frac{\|\vec{v}\|^2}{2}, \quad (10)$$

with  $k$  the kinetic energy, the discharge

$$\vec{q} = H\vec{v}, \quad (11)$$

and the Froude number

$$\text{Fr} = \frac{\|\vec{v}\|}{c}, \quad (12)$$

playing for (1) the same role as the Mach number in gas dynamics.

System (5) is known to admit a certain number of exact steady solutions whose form depend on the equilibrium between the source term  $\mathcal{S}$  and the remaining terms of the equation. In the following sub-sections we recall some of these solutions which will be used later to test the scheme.

## 2.1 Lake at rest solution

This solution corresponds to the hydrostatic equilibrium

$$\begin{aligned} \vec{q} &= 0 \\ \nabla p(H) + gH\nabla b &= 0 \end{aligned}$$

The last relation is always satisfied by the physical steady state (cf. equation (9))

$$\begin{aligned} \vec{v} &= 0 \\ \eta &= \eta_0 = \text{const} \end{aligned} \quad (13)$$

corresponding to still water on an arbitrary bathymetry.

## 2.2 Constant energy pseudo-1d flows

A pseudo one-dimensional steady equilibrium is readily obtained in the frictionless case by rewriting (1) as (cf. equation (10))

$$\begin{aligned} \partial_t H + \nabla \cdot \vec{q} &= 0 \\ \partial_t \vec{q} + (\vec{v} \cdot \nabla) \vec{q} - (\vec{v}^\perp \cdot \nabla) \vec{q}^\perp + \frac{1}{1 - \text{Fr}^2} \frac{1}{g} (gH\nabla \mathcal{E} - \vec{v}\vec{v} \cdot \nabla \mathcal{E}) &, \\ + \frac{1}{1 - \text{Fr}^2} \left( \frac{\vec{v}}{gH} \vec{v} \cdot (\nabla \vec{q} \cdot \vec{v}) - \text{Fr}^2 (\nabla \vec{q})^t \cdot \vec{v} \right) &= \frac{\vec{v}^\perp \cdot \nabla b}{1 - \text{Fr}^2} \vec{v}^\perp \end{aligned} \quad (14)$$

with  $\vec{v}^\perp = (-v_y, v_x)$ ,  $\vec{q}^\perp = H\vec{v}^\perp$ . The left hand side in the last equations only depends on derivatives of the discharge  $\vec{q}$ , and of the total energy  $\mathcal{E}$ . This shows that, provided we verify the compatibility condition for the bathymetry

$$\vec{v}^\perp \cdot \nabla b = 0, \quad (15)$$

there exist an admissible family of steady solutions characterized by the invariants

$$\begin{aligned}\vec{q} &= \vec{q}_0 = \text{const} \\ \mathcal{E} &= \mathcal{E}_0 = \text{const}\end{aligned}\tag{16}$$

Note that condition (15) allows bathymetry variations only in the direction of the discharge. This makes these solutions basically one-dimensional flows in the  $\vec{v}$  direction.

## 2.3 Steady flows in sloping channels

If the bathymetry has a constant slope, a steady solution is obtained from the equilibrium between friction and the hydrostatic load. If, without loss of generality, we take

$$b = b_0 - \zeta_0 x,$$

with constant slope  $\zeta_0$ , and assume  $v_y = 0$ , the equations reduce to

$$\begin{aligned}q_x &= q_0(y) \\ \partial_x H + c_f(H, \vec{v})v_x &= \zeta_0 \\ q_y &= 0 \\ \partial_y \eta &= 0\end{aligned}\tag{17}$$

A known solution is given by the pseudo one-dimensional state

$$\begin{aligned}H &= H_0 = \text{const} \\ v_x &= v_0 = \text{const} \\ v_y &= 0\end{aligned}\tag{18}$$

When using Manning's formula (3), the values of  $H_0$  and  $v_0$  can be expressed in a general manner as a function of the mass flux  $q_0$  and of the slope as

$$\begin{aligned}H &= H_0 = \left( \frac{n^2 \|\vec{q}_0\|^2}{|\zeta_0|} \right)^{\frac{3}{10}} \\ \vec{v} &= \vec{v}_0 = - \left( \frac{|\zeta_0| \|\vec{q}_0\|^{4/3}}{n^2} \right)^{\frac{3}{10}} \frac{\nabla b}{|\zeta_0|},\end{aligned}\tag{19}$$

representing a pseudo one-dimensional flow in the direction of  $\nabla b$ .

# 3 Explicit residual approach for conservation laws

## 3.1 Generalities

In this section we recall of the explicit residual schemes initially proposed in [57]. Let us first recall that the NLSW can be recast in compact form as the system of nonlinear partial differential equations

$$\partial_t u + \nabla \cdot \mathcal{F}(u) + \mathcal{S}(u, x, y) = 0,\tag{20}$$

defined on a space-time domain  $\Omega \times [0, T]$ . We consider now an unstructured triangulation of  $\Omega$ , which we denote by  $\Omega_h$ ,  $h$  denoting the largest element diameter. We denote by  $K$  the generic element of the mesh, and by  $|K|$  its area. On  $\Omega_h$ , we consider the standard continuous  $P^1$  finite element approximation of  $u$ , which we denote by  $u_h$ . For every node  $i \in \Omega_h$  we define  $K_i$ , the subset of elements containing  $i$  as a node

$$K_i = \bigcup_{K \in \Omega_h | i \in K} K. \quad (21)$$

We denote by  $C_i$  the standard median dual cell obtained by joining the gravity centers of the elements in  $K_i$  with the mid-points of the edges emanating from  $i$  whose area is given by

$$|C_i| = \sum_{K \in K_i} \frac{|K|}{3}. \quad (22)$$

The temporal domain  $[0, T]$  is approximated by a set of time slabs  $[t^n, t^{n+1}]$ . We denote by  $\Delta t^n = t^{n+1} - t^n$ , and we set  $\Delta t = \max_n \Delta t^n$ .

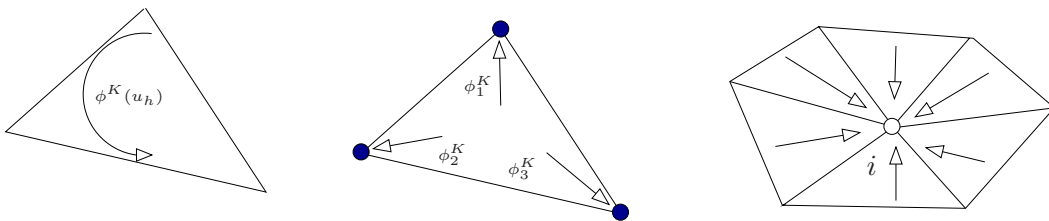


Figure 2: Residual distribution/fluctuation splitting methodology

As summarized on figure 2, the approach proposed follows the Residual Distribution (RD) philosophy inspired by the fluctuation and signals framework initially introduced by P.L. Roe [65]. According to this principle, data are evolved by means of signals proportional to a local error in the approximation of the equation (Roe's fluctuation).

This framework has led to the development of a certain number of numerical schemes known as the fluctuation splitting, or more recently, residual distribution schemes. For a review, the interested reader can consult [4, 5, 27, 66] and references therein. For the time dependent NLSW, the most recent adaptation of this approach is discussed in [59, 60]. The scheme proposed in the last references allows to solve the NLSW with second order of accuracy for smooth solutions on unstructured grids, it satisfies the C-property [14] for the lake at rest state, and an approximate generalized C-property for the pseudo one-dimensional flow of section §2.2 [56], it preserves the positivity of the depth, and, in conjunction with a proper wetting and drying strategy, allows the approximation of solutions involving runup, drying and flooding of complex bathymetries.

The main flaw of the scheme of [59, 60] is that, even though the scheme is highly implicit, positivity preservation is only achieved under an explicit time step restriction dictated by the underlying Crank-Nicholson time integration. In order to overcome this limitations, two



strategies have been proposed. One is the genuinely space-time formulation of [42], extended to the NLSW in [68]. The second is the genuinely explicit Runge-Kutta predictor-corrector variant of [57]. Preliminary results on the application of the latter to the NLSW have been presented in [55], and in the thesis manuscript [54]. Here we will analyze and extend the study of the last reference.

### 3.2 The predictor-corrector explicit scheme

Given the approximation of the initial solution  $u_h^0$ , in its simplest form, the second order scheme of [57] allows to march in time according to the following computational procedure :

1. Predictor step :

- $\forall K \in \Omega_h$  compute the *fluctuation* defined as

$$\phi^K(u_h^n) = \oint_{\partial K} \mathcal{F}_h(u_h^n) \cdot \hat{n} + \int_K \mathcal{S}_h(u_h^n, x, y). \quad (23)$$

- $\forall K \in \Omega_h$  distribute the *fluctuation* to the nodes of  $K$ . If  $\phi_i^K$  denotes the amount of  $\phi^K$  distributed to node  $i \in K$ , then

$$\sum_{j \in K} \phi_j^K(u_h^n) = \phi^K(u_h^n). \quad (24)$$

Equivalently, if there exist *bounded distribution matrix coefficients*  $\beta_i^K$  such that

$$\phi_i^K = \beta_i^K \phi^K, \quad (25)$$

then the *consistency condition* (24) becomes

$$\sum_{j \in K} \beta_j^K = \mathbf{I}. \quad (26)$$

- $\forall i \in \Omega_h$  compute the first order predictor  $u_i^*$  from

$$|C_i| \frac{u_i^* - u_i^n}{\Delta t^n} + \sum_{K \in K_i} \phi_i^K(u_h^n) = 0. \quad (27)$$

2. Corrector step :

- $\forall K \in \Omega_h$  compute the *element residual* defined as (cf. equation (23))

$$\Phi^K(u_h^n, u_h^*) = \int_K \frac{u_h^* - u_h^n}{\Delta t^n} + \frac{1}{2} \phi^K(u_h^n) + \frac{1}{2} \phi^K(u_h^*). \quad (28)$$

Note that, as in [57], we distinguish between the fluctuation (23), containing the integral of the spatial operator, and the element residual (28), defined as the integral of the full semi-discrete equation.

- $\forall K \in \Omega_h$  distribute the *element residual* to the nodes of  $K$ . If  $\Phi_i^K$  denotes the amount of  $\Phi^K$  distributed to node  $i \in K$ , then

$$\sum_{j \in K} \Phi_j^K(u_h^n, u_h^*) = \Phi^K(u_h^n, u_h^*). \quad (29)$$

As before, if there exist *bounded distribution matrix coefficients*  $\beta_i^K$  such that

$$\Phi_i^K = \beta_i^K \Phi^K, \quad (30)$$

then the distribution matrices must verify the consistency condition (26)

- $\forall i \in \Omega_h$  compute the second order correction from

$$|C_i| \frac{u_i^{n+1} - u_i^*}{\Delta t^n} + \sum_{K \in K_i} \Phi_i^K(u_h^n, u_h^*) = 0. \quad (31)$$

In the above expressions it remains to specify not only how the distribution is performed, but also how to define the *discrete approximation of the flux and of the source term*  $\mathcal{F}_h(u_h)$ , and  $\mathcal{S}_h(u_h, x, y)$ , given the nodal values of the solution.

**Remark 3.1** (Fluctuation and signals). *With the notation introduced so far, we can easily recall that the implicit scheme of [59, 60] can be recast as : given the initial solution  $u_h^0$  march in time by solving the nonlinear system*

$$\sum_{K \in K_i} \Phi_i^K(u_h^n, u_h^{n+1}) = 0, \quad \forall i \in \Omega_h$$

where  $\forall K \in \Omega_h$  the quantity  $\Phi_i^K(u_h^n, u_h^{n+1})$  is a *splitting of the element residual*

$$\Phi^K(u_h^n, u_h^{n+1}) = \int_K \frac{u_h^{n+1} - u_h^n}{\Delta t^n} + \frac{1}{2} \phi^K(u_h^n) + \frac{1}{2} \phi^K(u_h^{n+1}) = \sum_{j \in K} \Phi_j^K(u_h^n, u_h^{n+1}).$$

It is immediately clear that the scheme of [59, 60] would be, eventually, obtained using (31) as an iterative scheme to get to the fixed point  $u_i^* = u_i^{n+1}$ .

The right hand sides of (27) and (31) are thus easily interpreted as corrections of the nodal values somehow proportional to elemental errors given by the integral of the equations, represented by the residual  $\Phi^K$ . This allows to view the scheme proposed as a truly time dependent generalization of Roe's initial fluctuation splitting idea [65].

### 3.3 Basic properties

The accuracy properties of the explicit scheme described in the previous section are thoroughly discussed in [57]. We recall here the conditions under which scheme (23)-(31) is conservative, second order accurate, and satisfies a discrete maximum principle.

**Conservation.** By conservation we mean the ability to reproduce the correct jump conditions across discontinuous solutions. This property is characterized by a Lax-Wendroff theorem firstly formulated in [7], and then further clarified in [6, 10] and in [26, 62]. Without going

into the details of the theorem, for which we refer the reader to [7, 6, 10], we recall that *provided that the consistency conditions (24) and (29) hold, scheme (23)-(31) is conservative if the discrete approximation of the physical flux  $\mathcal{F}_h(u_h)$  is continuous across edges.*

Note that this continuity condition is satisfied by several definitions of the discrete flux, such as  $\mathcal{F}_h(u_h) = \mathcal{F}_h$  the  $P^1$  finite element interpolation of the nodal values of the flux, or also  $\mathcal{F}_h(u_h) = \mathcal{F}(u_h)$ , and more generally by  $\mathcal{F}_h(u_h) = \mathcal{F}(u(v_h))$  where  $v_h$  is the finite element interpolation of a set of variables different from  $u$ . A similar set of choices is possible for  $\mathcal{S}_h(u_h, x, y)$  as well. This freedom can be exploited to recognize steady equilibria, as we shall see later.

**Second-order.** As discussed in much detail in [57, 54], scheme (23)-(31) can be obtained as a particular case of a mass lumped  $P^1$  Petrov-Galerkin finite element discretization. In particular, in order to formally characterize the accuracy, let us consider a smooth exact pointwise solution  $w(x, y, t)$ , such that

$$\partial_t w + \nabla \cdot \mathcal{F}(w) + \mathcal{S}(w, x, y) = 0 \quad \forall (x, y, t) \in \Omega \times [0, T].$$

Let  $w_i^n$  be its nodal values at the discrete time level  $t^n$ , and  $w_h^n$  the finite element approximation at the same time level. Consider also a  $C^1$  continuous compactly supported function  $\psi \in C_0^1(\Omega \times [0, T])$ , and define the truncation error

$$\epsilon := \sum_{n=0}^N \Delta t^n \sum_{i \in \Omega_h} \psi_i^n \left[ |C_i| \frac{w_i^{n+1} - w_i^*}{\Delta t^n} + \sum_{K \in K_i} \Phi_i^K(w_h^n, w_h^*) \right], \quad (32)$$

with  $\sum_{n=0}^N \Delta t^n = T$ , and  $w_i^*$  obtained from (27) :

$$|C_i| \frac{w_i^* - w_i^n}{\Delta t^n} + \sum_{K \in K_i} \phi_i^K(w_h^n) = 0. \quad (33)$$

In [57, 54] it is proven that *if there exist distribution matrix coefficients  $\{\beta_j^K\}_{j \in K}$  uniformly bounded w.r.t.  $h$ ,  $w_h$ ,  $\phi^K$ ,  $\Phi^K$  and the data of the problem, such that (25) and (30) hold, then*

$$\begin{aligned} \epsilon + \mathcal{O}(\Delta t^2) = & \int_0^T \int_{\Omega_h} \psi_h \partial_t (w_h - w) + \int_0^T \int_{\Omega_h} (\mathcal{F}(w) - \mathcal{F}_h(w_h)) \cdot \nabla \psi_h + \int_0^T \int_{\Omega_h} (\mathcal{S}_h(w_h, x, y) - \mathcal{S}(w, x, y)) \psi_h \\ & + \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i, j \in K} \Delta t^n \frac{\psi_i^n - \psi_j^n}{3} \int_K \gamma_i^K (\partial_t (w_h - w) + \nabla \cdot (\mathcal{F}_h(w_h) - \mathcal{F}(w)) + \mathcal{S}_h(w_h, x, y) - \mathcal{S}(w, x, y)) \end{aligned},$$

with  $\gamma_i^K$  a suitably chosen uniformly bounded Petrov-Galerkin bubble stabilization depending on the distribution coefficients  $\beta_i^K$ . Moreover, under the same hypotheses, and provided that there exist positive bounded constants  $C_1, C_2, C_a, C_b$  such that

$$C_1 \leq \max_{K \in \Omega_h} \frac{h^2}{|K|} \leq C_2, \quad C_a \leq \frac{\Delta t}{h} \leq C_b, \quad (34)$$

then the error can be bounded using classical approximation arguments to obtain the consistency estimate :

$$|\epsilon| \leq C(\Omega_h, T, \psi)h^2. \quad (35)$$

For further details, the interested reader can consult [57] or the manuscript [54]. We limit ourselves to observe that the key to second order of accuracy is the uniform boundedness of the distribution coefficients (25), (30).

**Discrete Maximum Principle.** The non-oscillatory character of the discretization is preserved in the residual distribution framework by making use of the theory of positive coefficient schemes [13, 27, 71]. The theory mainly deals with the case in which (20) reduces to a scalar homogeneous conservation law

$$\partial_t u + \nabla \cdot \mathcal{F}(u) = 0.$$

In this case, scheme (23)-(31) remains formally identical, modulo the fact that  $\mathcal{S} = 0$ . In summary, the idea is to ensure that in (27) and (31) one has  $\forall K$

$$\begin{cases} \phi_i^K = \sum_{j \in K} c_{ij}^n (u_i^n - u_j^n) \\ \Phi_i^K = \frac{1}{2} \sum_{j \in K} c_{ij}^* (u_i^* - u_j^*) + \frac{1}{2} \sum_{j \in K} \bar{c}_{ij}^n (u_i^n - u_j^n) \end{cases} \quad \text{with } c_{ij}^n, c_{ij}^*, \bar{c}_{ij}^n \geq 0. \quad (36)$$

This allows to recast the predicted values as

$$u_i^* = u_i^n - \frac{\Delta t^n}{|C_i|} \sum_{K \in K_i} \sum_{j \in K} c_{ij}^n (u_i^n - u_j^n),$$

which with the hypotheses made, and under the time step restriction  $\Delta t^n \leq |C_i| / \sum_K \sum_{j \in K} c_{ij}^n$ , are easily shown to verify

$$\min_{j \in K_i} u_j^n \leq u_i^* \leq \max_{j \in K_i} u_j^n.$$

For the corrector step we can write

$$\begin{aligned} u_i^{n+1} &= u_i^* - \frac{\Delta t^n}{2|C_i|} \sum_{K \in K_i} \sum_{j \in K} c_{ij}^* (u_i^* - u_j^*) - \frac{\Delta t^n}{2|C_i|} \sum_{K \in K_i} \sum_{j \in K} \bar{c}_{ij}^n (u_i^n - u_j^n) \\ &= \frac{1}{2} \left( u_i^* - \frac{\Delta t^n}{|C_i|} \sum_{K \in K_i} \sum_{j \in K} c_{ij}^* (u_i^* - u_j^*) \right) + \frac{1}{2} \left( u_i^n - \frac{\Delta t^n}{|C_i|} \sum_{K \in K_i} \sum_{j \in K} (\bar{c}_{ij}^n + c_{ij}^n) (u_i^n - u_j^n) \right), \end{aligned}$$

so, provided that

$$\Delta t^n \leq \min \left( \frac{|C_i|}{\sum_{K \in K_i} \sum_{j \in K} c_{ij}^*}, \frac{|C_i|}{\sum_{K \in K_i} \sum_{j \in K} (c_{ij}^n + \bar{c}_{ij}^n)} \right),$$

we have

$$\frac{1}{2} \min_{j \in K_i} u_j^* + \frac{1}{2} \min_{j \in K_i} u_j^n \leq u_i^{n+1} \leq \frac{1}{2} \max_{j \in K_i} u_j^* + \frac{1}{2} \max_{j \in K_i} u_j^n.$$

For scalar homogeneous problems, the theory of positive coefficient schemes allows to give precise conditions leading to local bounds on the numerical solution, related to the local extrema at the old time level. For a hyperbolic system things become much more complicated,

and even in the continuous case the existence of maximum principles is hard to prove in a general way, its definition being unclear even at the continuous level. We mention that a discrete wave decomposition technique has been proposed in [9] as a means to extend positive coefficient schemes theory to systems.

In this paper, we will say that *a scheme is positive if it verifies (36) in the scalar case*. For the NLSW (and for systems in general) practical tests show that indeed these schemes show a non-oscillatory approximation of discontinuities. Using the same theory, however, in section §4 we shall study in more detail the preservation of the physical constraint  $H \geq 0$  for the scheme used in all the numerical applications.

## 4 Application to the shallow water equations

We consider now the properties of scheme (23)-(31) when applied to the NLSW. Two issues are analyzed : the preservation of steady equilibria, and the wetting/drying strategy, including the issue of the preservation of the constraint  $H \geq 0$ . The analysis of the preservation of steady equilibria is general and applies to all schemes that formally verify (25) and (30). On the contrary, so far computations involving dry areas rely on nonlinear Lax-Friedrichs distribution initially introduced in [3], and constituting the basis of the results presented in [60, 59] for the NLSW.

### 4.1 C-property and super-consistency analysis

The C-property or “conservation” property, introduced in of [14], consists of the ability of the discretization in preserving the exact steady state balance of flux divergence and source terms. Originally, a scheme was said to enjoy the C-property if it preserved exactly the steady state (13). However one still speaks of C-property when referring to other steady states [23]. When the conservation of the steady state is no exact but is obtained within error rates below the formal accuracy of the scheme, one often speaks of generalized C-property [56]. Schemes enjoying the C-property, of the generalized C-property, are more often referred to in literature as being “well balanced”. In this paper we will say that a scheme verifies the C-property for a given steady state, if that state is preserved exactly by the scheme. If the preservation is obtained within error bounds decreasing with rates larger than those of the formal accuracy of our scheme, we will speak of *super consistency* with that particular solution.

For the schemes considered here, we start by providing a general result concerning smooth steady equilibria admitting a set of invariants. In the next section, we will analyze the case of the particular steady solutions discussed in section §2.

The underpinning idea is that if for a given state  $u = u^0(x)$  we have  $\phi^K(u^0) = 0$ , then the scheme defined by (23)-(31), (25) and (30) will preserve the initial state indefinitely since

$$|C_i| \frac{u_i^* - u_i^0}{\Delta t^n} = |C_i| \frac{u_i^{n+1} - u_i^*}{\Delta t^n} = 0.$$

The objective of the analysis is to quantify the effect of the approximation choice and of the quadrature errors in the evaluation of  $\phi^K$  on the preservation of smooth equilibria.

To start, we assume here that *both the flux  $\mathcal{F}$  and the source term  $\mathcal{S}$  are at least Lipschitz continuous* :

$$\|\mathcal{F}(u) - \mathcal{F}(w)\| \leq \mathcal{K}_{\mathcal{F}}\|u - w\|, \quad \text{and} \quad \|\mathcal{S}(u, \nabla b) - \mathcal{S}(w, \nabla b)\| \leq \mathcal{K}_{\mathcal{S}}\|u - w\|. \quad (37)$$

Consider now a set of derived variables  $v$  that depend on  $u$ , and that also depend on  $b$ . In particular, we consider now the mappings

$$v : (u, b) \mapsto v(u, b), \quad U = v^{-1} : (v, b) \mapsto u = U(v, b). \quad (38)$$

It is assumed in the following that these mappings are smooth. Examples of such variables are

- Total energy variables (cf. equation (10)) :  $v = [\mathcal{E}, \vec{q}]^t$  ;
- Symmetrizing variables [75, 76, 38, 39] :  $v = [g\eta - k, \vec{v}]^t$  ;

The flux, being independent of  $b$ , can be expressed as the divergence term (cf. equation (38))

$$\begin{aligned} \nabla \cdot \mathcal{F}(u(v, b)) &= \frac{\partial \mathcal{F}}{\partial u}(u(v, b)) \frac{\partial u}{\partial v}(u(v, b)) \cdot \nabla v + \frac{\partial \mathcal{F}}{\partial u}(u(v, b)) \frac{\partial U}{\partial b}(u(v, b)) \cdot \nabla b \\ &= \frac{\partial \mathcal{F}}{\partial v}(u(v, b)) \cdot \nabla v + \mathcal{S}_v(u(v, b), \nabla b), \end{aligned} \quad (39)$$

where  $\mathcal{S}_v(u(v, b), \nabla b)$  is the contribution of all the terms containing derivatives of the bathymetry. For the analysis that follows we make the additional hypothesis of the explicit knowledge of an analytical bathymetry and of the definition of the continuous flux approximation and discrete source term built starting from the nodal values of the steady invariant, and from the local values of  $b(x, y)$ , namely

$$\mathcal{F}_h = \mathcal{F}(u(v_h, b)), \quad \mathcal{S}_h = \mathcal{S}(u(v_h, b), \nabla b). \quad (40)$$

Using this notation, we can prove the following result.

**Lemma 4.1** (Super consistency - local estimate). *Given an analytical bathymetry  $b$ , let  $v(u, b)$  be a set of invariants such that a family of steady equilibria for (1) is completely described by*

$$v = v_0 = \text{const}.$$

*Let  $\mathcal{F}_h = \mathcal{F}(u(v_h, b))$  and  $\mathcal{S}_h = \mathcal{S}(u(v_h, b), \nabla b)$ , with  $v_h$  the piecewise linear continuous  $P^1$  approximation of  $v$ . Assume that  $(v, b) \mapsto u$  is a one to one smooth mapping  $C^l$  with  $l$  sufficiently large, and similarly  $(v, b) \mapsto \mathcal{F}(u(v, b))$  is also  $C^{l'}$  with  $l'$  sufficiently large. Then, for exact integration we have*

$$\phi^K(v_0, b) = \oint_{\partial K} \mathcal{F}_h \cdot \vec{n} + \int_K \mathcal{S}_h = 0.$$

*For approximate integration, let*

$$\phi^K(v_0, b) = \sum_{f \in \partial K} \sum_{q=1}^{f_q} \omega_q \mathcal{F}_h(\vec{x}_q) \cdot \vec{n}_f + \sum_{q=1}^{v_q} \bar{\omega}_q \mathcal{S}_h(\vec{x}_q),$$

and let the line quadrature formula used for the flux be exact for polynomials of degree  $p_f \geq 1$ , and the volume quadrature formula used for the source be exact for polynomials of degree  $p_v \geq 1$ . If  $b \in H^{p+1}(\Omega)$  with  $\nabla b \in H^p(\Omega_h)$  and with  $p \geq \min(p_f, p_v + 1)$ , then

$$\|\phi^K(v_0, b)\| \leq C h^r \quad \text{with } r = \min(p_f + 2, p_v + 3).$$

*Proof.* See appendix A. □

The meaning of the lemma is the following : for a smooth enough bathymetry, provided that the approximation is written in terms of the steady invariants, the fluctuation (23) is small, and its magnitude is dictated by the numerical quadrature used in practice.

Under the same hypotheses, and with the same notation of the lemma, in appendix A we prove the following general result.

**Proposition 4.2** (Super consistency). *Under the regularity assumptions on the mesh and on the time step (34), and provided that (25)-(30) are true for some distribution coefficients  $\beta_i^K$  uniformly bounded w.r.t.  $h$ ,  $u_h$ , element residuals, and w.r.t. to the data of the problem, then scheme (23)-(31) preserves exactly the initial steady equilibrium for exact integration. For approximate integration, under the same hypotheses of lemma 4.1 scheme (23)-(31) verifies a the truncation error estimate*

$$\|\epsilon(v_0, b, \psi)\| \leq C h^l, \quad l = \min(p_f + 1, p_v + 1), \quad (41)$$

with the error  $\epsilon(v_0, b, \psi)$  defined as in (32).

The last proposition shows that for finite time computations, if the bathymetry is regular enough, the discrete solution converges with a rate  $l > 2$ , as soon as the quadrature formulas are at least second order accurate. The proposition explicitly uses the assumption that an analytical bathymetry is used in the discretization, and that the regularity of this expression is such that the full accuracy of the quadrature formulas is recovered.

**Remark 4.3** (Convergence rates observed in practice). *As the results of section §5 will show, the convergence rates obtained in practice for a bathymetry  $b \in H^{p+1}(\Omega)$  are given by*

$$\epsilon(v_0, b, \psi) = \mathcal{O}(h^{\bar{l}}), \quad \bar{l} = \min(p + 1, p_f + 1, p_v + 2). \quad (42)$$

*This means that the proposition fails somehow to capture some error cancelation effects allowing to recover and additional degree of convergence for a given volume quadrature accuracy. Furthermore, the proposition does not take into account the regularity of the bathymetry. In particular, the convergence saturates at a rate given by  $p + 1$ . This result is quite natural, since the estimates, and the form of the scheme, rely on the smoothness of  $b$  in the quadrature set to achieve super convergence. A consequence of this fact is that if the curves across which  $b$  presents jumps in one of its derivatives are aligned with the grid, the local regularity of  $b$  is still high enough to allow a super-consistent behavior, as the numerical results will show.*

*So, while low convergence rates are obtained for low regularity of the bathymetry, if needed this can be cured by using a properly designed adaptive mesh recovering a fully super convergent behavior. For the lake at rest solution, we will also show that a better result is obtained for the lake at rest solution when replacing  $b(x, y)$  by the same finite element interpolation used for the depth.*

## 4.2 C-property : application to particular steady states

In this section we make a case by case analysis of the steady solutions presented in section §2. The principle used will be always to write the continuous approximation of the fluxes and the one of the source term starting from nodal values of the steady state invariants. The discussion made here will be confirmed by the numerical experiments of section §5. As in [14, 23] we will refer to the exact conservation of a steady equilibrium as the C-property, while, following section §4.1, we will talk of super consistency when preservation is only obtained within an error below the theoretical truncation error of the schemes.

Note that, as shown clearly by the results of the previous section, and as it will be clear from the following paragraphs, different approximation choices allow to maintain different type of equilibria. In practice, the interpolation used should take into account both the type of flows one is interested in, and the simplicity and cost of the method, some choices being more expensive than others.

### 4.2.1 The lake at rest solution

The approximation of this solution by means of RD type discretization has been thoroughly studied in [17, 56, 58, 60]. We can distinguish two cases. The first is the one in which the approach of section §4.1 is used, and the discrete flux and source term approximations are written as (40). In this case, the super consistency property of proposition 4.2 is recovered with e.g.  $v = v_0 = [\eta_0, 0]$ . In this case, we note that the components of physical flux  $\mathcal{F}$  are polynomials of degree at most 3 of the components of  $v$ , so that accurate exact numerical quadrature becomes quite easy. However, even for this simple solution, for bathymetries with low regularity this approach will yield errors of  $\mathcal{O}(h)$  which is below the formal accuracy of the scheme.

A more interesting choice is that of [17, 58, 60], and consisting in replacing  $b$  by the same finite element interpolation used for  $H$ . In this case, one simply finds that the discrete approximation of  $H$  obtained as  $H_h = \eta_h - b_h$  is exactly that which would be obtained by approximating directly  $H$ . Moreover, following the above references, we can readily show that along the lake at rest state

$$\begin{aligned} \phi^K(v_0) &= \int_{\partial K} \begin{bmatrix} H_h \vec{v}_h \\ H_h \vec{v}_h \otimes \vec{v}_h + g \frac{H_h^2}{2} \mathbf{I} \end{bmatrix} \cdot \vec{n} + \int_K g H_h \begin{bmatrix} 0 \\ \nabla b_h + c_f \vec{v}_h \end{bmatrix} \\ &= \int_{\partial K} \begin{bmatrix} 0 \\ g \frac{H_h^2}{2} \mathbf{I} \end{bmatrix} \cdot \vec{n} + \int_K g H_h \begin{bmatrix} 0 \\ \nabla b_h \end{bmatrix} = \int_K g H_h \begin{bmatrix} 0 \\ \nabla \eta_h \end{bmatrix} = 0 \end{aligned},$$

where all the terms can be evaluated exactly by means of simple Gauss quadrature formulas. This leads to the following result.

**Proposition 4.4** (Lake at rest - C-property). *For a given bathymetry  $b(x, y)$  let  $b_h$  be its finite element approximation. Let the discrete approximation of the flux be given by  $\mathcal{F}_h = \mathcal{F}([H_h, \vec{q}_h]^t)$ . Provided that the quadrature formulas used in (23) are exact w.r.t  $H_h^2$ , provided that (25)-(30) are true for some distribution coefficients  $\beta_i^K$  uniformly bounded w.r.t.  $h$ ,  $u_h$ , element residuals, and w.r.t. to the data of the problem, then then scheme (23)-(31) preserves exactly the lake at rest solution, independently of the regularity of  $b(x, y)$ .*



### 4.2.2 Constant energy pseudo-1d flows

The constant energy pseudo-1d solution equilibrium fits quite well in the analysis made in section §4.1. Numerical results showing that a super convergent behavior is observed when interpolating directly the total energy have been already reported in [56] for the scheme of [60]. The super consistency property of proposition 4.2 is recovered by interpolating the invariant  $v = v_0 = [\mathcal{E}_0, \vec{q}_0]$ . Note that, similar results have been presented *e.g.* in [33, 51, 80]. However, as the numerical results will show, here we obtain a super convergent behavior on unstructured triangular meshes, while only one dimensional and two dimensional cartesian meshes have been studied previously, with the exception of the author's previous work [56].

Concerning the implementation, as discussed in detail in [51], when writing the approximation in terms of the energy variables  $[\mathcal{E}, \vec{q}]^t$ , the local values of the set  $[H, \vec{v}]^t$  are obtained as the solution of a nonlinear algebraic problem. In particular, given a tuple  $[\mathcal{E}^*, \vec{q}^*]$ , and the local value of the bathymetry  $b^*$ , in order to obtain the corresponding values of  $H$  and  $\vec{v}$  needed to evaluate the flux, one needs to solve the nonlinear system (cf. equation (10))

$$\begin{aligned} H\vec{v} &= \vec{q}^* \\ gH + \frac{\vec{v} \cdot \vec{v}}{2} &= \mathcal{E}^* - gb^* \end{aligned} \quad ,$$

where, using the first equation one obtains a cubic algebraic equation for  $H$  :

$$H^3 + \left(b^* - \frac{\mathcal{E}^*}{g}\right)H^2 + \frac{\vec{q}^* \cdot \vec{q}^*}{2g} = 0.$$

The conditions for the existence of solutions to this equation, and a small summary of the Newton algorithm necessary to obtain them, are given in [51] which we have followed in all the numerical applications discussed in section §5.

Note that last equation needs to be solved every time one needs to go from interpolated to physical values, which means in every mesh node, and in every quadrature point. This renders this approximation quite expensive, and in practice other choices are preferred if constant energy flows are not of interest.

### 4.2.3 Steady flows in sloping channels

These solutions are a particularly simple example of steady equilibrium between friction and a constant slope. In this case,  $H$ ,  $\vec{v}$ , and  $\vec{q}$  are constant. The friction term  $c_f(H_0, \vec{v}_0)$  exactly balances the constant slope  $\nabla b = \zeta_0$  (cf. section §2.3). Simple algebra shows that in this case approximating directly the conserved variables  $[H, \vec{q}]^t$  leads to the identity

$$\phi^K = \oint_{\partial K} \begin{bmatrix} \vec{q}_0 \\ \vec{q}_0 \otimes \vec{v}_0 + g \frac{H_0^2}{2} \mathbf{I} \end{bmatrix} \cdot \vec{n} + \int_K gH_h \begin{bmatrix} 0 \\ \underbrace{\nabla b_0 + c_f \vec{v}_0}_{=0} \end{bmatrix} = \begin{bmatrix} \vec{q}_0 \\ \vec{q}_0 \otimes \vec{v}_0 + g \frac{H_0^2}{2} \mathbf{I} \end{bmatrix} \cdot \overbrace{\oint_{\partial K} \vec{n}}^{=0} = 0.$$

This result is independent of how the integrals are evaluated, thus proving the following property.

**Proposition 4.5** (Steady flows in sloping channels - C-property). *Given the constant slope bathymetry  $b(x) = b_0 - \zeta_0 x$ , let the discrete approximation of the flux be given by  $\mathcal{F}_h =$*

$\mathcal{F}([H_h, \vec{q}_h]^t)$ . Provided that (25)-(30) are true for some distribution coefficients  $\beta_i^K$  uniformly bounded w.r.t.  $h$ ,  $u_h$ , element residuals, and w.r.t. to the data of the problem, then scheme (23)-(31) preserves exactly steady flows on constant sloping channels, independently on the quadrature formulas used in (23).

### 4.3 Nonlinear Lax-Friedrichs distribution

The super-consistency and C- properties discussed in the previous sections are independent of the actual form of the distribution coefficient, and are valid as long as this coefficient is bounded. To complete the presentation of the scheme used in the numerical validation we consider in this section the issue of the preservation of the depth non-negativity, and, in the following, that of the treatment of wet-dry fronts.

The scheme used in this work is the nonlinear variant of the Lax-Friedrichs scheme in the stabilized form originally proposed in [3], and further adapted for the solution of the time dependent NLSW in [59, 60]. As in the references, the starting point of the construction is the first order Lax-Friedrichs distribution (cf. section §3.2 equations (23), (24), and (27))

$$\phi_i^{\text{LF}}(u_h^n) = \frac{\phi^K(u_h^n)}{3} + \frac{\alpha_{\text{LF}}}{3} \sum_{j \in K} (u_i^n - u_j^n), \quad (43)$$

in the predictor step and (cf. section §3.2 equations (28), (29), and (31))

$$\Phi_i^{\text{LF}} = \frac{|K|}{3} (u_i^* - u_i^n) + \frac{\phi_i^{\text{LF}}(u_h^n) + \phi_i^{\text{LF}}(u_h^*)}{2}, \quad (44)$$

in the corrector step. The usual definition of the LF dissipation coefficient  $\alpha_{\text{LF}}$ , satisfying the positivity requirement of section §3.3 in the scalar case, is some upper bound within the time step to the the largest absolute value of the flux Jacobians evaluated in the nodes of an element. In practice, this upper bound is often approximated by [3, 57]

$$\alpha_{\text{LF}} = \frac{1}{2} \max_{j \in K} (l_j (\|\vec{v}_j^n\| + c_j^n)), \quad (45)$$

with  $c_j$  the NLSW celerity (8), and  $l_j$  the length of the edge opposite node  $j$ . These definitions, combined with equations (27) and (31) give a straightforward two-dimensional generalization of the local Lax-Friedrichs scheme with second order SSP Runge-Kutta time integration, and can be easily shown to preserve the non-negativity of the depth  $H$  (see e.g. [53] and [60]). In the scalar case, the scheme verifies the positivity condition (36), and in general yields non-oscillatory solutions.

Unfortunately, the LF scheme is only first order accurate. Indeed, the LF scheme does not verify the condition for second order of accuracy recalled in section §2.2. In particular, the relations

$$\beta_i^{\text{LF}} \phi^K = \phi_i^{\text{LF}}, \quad \text{and} \quad \beta_i^{\text{LF}} \Phi^K = \Phi_i^{\text{LF}}$$

do not define bounded distribution coefficients, and in the system case do not even give enough conditions to determine them, unless some more assumptions are made. The nonlinear Limited Lax Friedrich's (LLF) scheme is obtained as follows (see e.g. [9] for more details)

- Project the LF elemental contributions (43) (or (44) in the corr. step) and the  $\phi^K$  (or  $\Phi^K$  in the corr. step) onto a basis of the solution space :  $\{\ell_m\}_{m=1}^3$ . Two possibilities are considered in this paper for  $\ell_k$  (see later) : either the left eigenvectors of the linearized flux Jacobian projected onto the local velocity direction (characteristic projection), or the Euclidean  $\mathbb{R}^3$  basis  $\{e_m\}_{m=1}^3$ , with  $e_{mj} = \delta_{mj}$  (limiting equation by equation, no projection). Let

$$\varphi^{K-m} = \ell_m \cdot \phi^K \quad (\text{or } \varphi^{K-m} = \ell_m \cdot \Phi^K \text{ in the corr. step}),$$

and

$$\varphi_i^{\text{LF}-m} = \ell_m \cdot \phi_i^{\text{LF}} \quad (\text{or } \varphi_i^{\text{LF}-m} = \ell_m \cdot \Phi_i^{\text{LF}} \text{ in the corr. step}).$$

- For each  $m \in \{1, 2, 3\}$ , apply a *sign preserving* nonlinear limiter to the otherwise unbounded LF distribution coefficients  $\beta_i^{\text{LF}-m} = \varphi_i^{\text{LF}-m} / \varphi^{K-m}$ . As in [3, 9, 60], the LLF distribution coefficients are computed as

$$\beta_i^{\text{LLF}-m} = \frac{\max(0, \beta_i^{\text{LF}-m})}{\sum_{j \in K} \max(0, \beta_j^{\text{LF}-m})} \quad (46)$$

- Redistribute the projected fluctuations as

$$\varphi_i^{\text{LLF}-m} = \beta_i^{\text{LLF}-m} \varphi^{K-m}.$$

Note that the properties of limiter (46) imply that (see e.g. [9, 60] for details)

$$\varphi_i^{\text{LLF}-m} = \beta_i^{\text{LLF}-m} \varphi^{K-m} = \gamma_i \varphi_i^{\text{LF}-m} \quad \text{with } \gamma_i \in [0, 1]. \quad (47)$$

- Project back to physical space

$$\phi_i^{\text{LLF}} = \sum_{m=1}^3 r_m \beta_i^{\text{LLF}-m} \varphi^{K-m} \quad (\text{or } \Phi_i^{\text{LLF}} = \sum_{m=1}^3 r_m \beta_i^{\text{LLF}-m} \varphi^{K-m} \text{ in the corr. step}),$$

with  $\{r_m\}_{m=1}^3$  a basis orthonormal to  $\{\ell_m\}_{m=1}^3$ . Here, the  $r_m$  will either be the right eigenvectors of the linearized flux Jacobian projected onto the local velocity direction (characteristic projection), or the Euclidean  $\mathbb{R}^3$  basis  $\{e_m\}_{m=1}^3$ , with  $e_{mj} = \delta_{mj}$  (limiting equation by equation, no projection).

The LLF scheme obtained with the above procedure is by construction second order accurate, its distribution coefficients being by construction uniformly bounded. Various theoretical results concerning its stability can be found in [3, 9, 60]. Concerning the explicit predictor corrector variant used in this paper, we prove in appendix B the following result.

**Proposition 4.6.** *Provided that  $\forall K \in \Omega_h$  the Lax-Friedrichs dissipation coefficient verifies*

$$\alpha_{LF} > \max_{x \in K} \|\vec{u}\|_\infty \quad \text{and} \quad \alpha_{LF} \geq C > 0,$$

*that the limiting is applied equation by equation, and that it verifies (47) on the depth equation, for a single  $\gamma_i \in [0, 1]$  in the predictor and corrector steps, then under the time step constraint*

$$\Delta t \leq \min_{i \in \Omega_h} \min \left( \frac{|C_i|}{\sum_{K \in K_i} \alpha_{LF}^n}, \min_{K \in K_i} \frac{|K|}{3\alpha_{LF}^n} \right), \quad (48)$$

*the explicit LLF scheme preserves the positivity of the depth  $H$ .*

The interest in this proposition is that it provides depth non-negativity preservation conditions very similar to those required for the scheme proposed in [60, 59]. In particular, the time step restriction (48) is roughly a half of the one required for the implicit scheme proposed in the references. This restriction has been compared to the one of the node centered first order upwind scheme on structured triangulations in [68]. From the analysis reported in the reference one can see that condition (48) is considerably more constraining than that of the upwind finite volume method. In particular, the results of the reference show that, on a structured triangulation of reference size  $h$ , (48) is equivalent for the constant advection equation to  $a\Delta t/h \leq 1/3$ , while the upwind node centered finite volume scheme requires  $a\Delta t/h \leq 0.53 - 0.62$ , depending on the orientation of the advection speed. It must be remarked that (48) is an extremely conservative condition allowing to prove the preservation of the non-negativity of the depth. In practice, values of CFL up to 1.2-1.3 can be used. This, however, still provides time steps 20%-30% smaller than the finite volume scheme.

On the other hand, compared to the implicit scheme proposed in [8, 60], still bound by a CFL condition, the explicit procedure proposed in this paper represents a definite improvement. In particular, the time step restriction of [60] is only twice as large as (48), leading to a scheme considerably less efficient.

#### 4.4 Filtering and streamline dissipation

While allowing the preservation of depth positivity and a monotone approximation of discontinuous solutions, the LLF scheme suffers from the appearance of weak spurious modes in correspondence of smooth solutions. A thorough analysis of this flaw is made in [3] and is beyond the scope of this paper. Possible solutions to cure this problem are suggested in [3, 60] for the steady case, and in [11, 60, 57] for the time dependent case. In particular, in [3] it is shown that the spurious modes are related to the ill-posedness of the algebraic equations in smooth regions. The reference also shows that the discrete equations are always well posed for schemes with a marked upwind character, which is not the case for the LLF scheme. For this reason, the solution suggested in the reference to filter the spurious modes is to add to the scheme a streamline dissipation term of the form [3, 43]

$$\phi_i^{\text{sd}} = \delta(u_h) \mathbf{K}_{n_i} \tau \phi^K,$$

with  $\vec{n}_i$  the inward pointing normal to the edge facing node  $i$  and scaled by the edge length, with  $\mathbf{K}_{n_i}$  the flux Jacobian projected onto  $\vec{n}_i$  (cf. section §2 equations (6) and (7) for the notation), and where the matrix  $\tau$  is a scaling parameter. Several forms of this parameter are suggested in the literature. The interested reader can refer to [43, 60] and references therein for a discussion. In this paper, we have set [3, 60]

$$\tau = \frac{1}{2} |K| \left( \sum_{j \in K} |\mathbf{K}_{n_j}| \right)^{-1}, \quad (49)$$

where the absolute value of the flux Jacobians  $|\mathbf{K}_{n_j}|$  is computed via standard eigen-decomposition. The parameter  $\delta(u_h)$  is a scalar smoothness sensor ensuring that the correction is only active in smooth parts of the flow. In particular,  $\delta(u_h) < Ch$  across discontinuities, while  $\delta \approx 1$  in smooth regions. So the scheme proposed in [3] reads

$$\phi_i = \phi_i^{\text{LLF}} + \delta(u_h) \mathbf{K}_{n_i} \tau \phi^K.$$

Following [57], in this paper we have used the fact that  $\delta$  is a scalar to define an efficient blending between the linear high order SUPG scheme and the nonlinear high order LLF scheme. The resulting LLFs (Stabilized Limited Lax Friedrichs) distribution reads

$$\phi_i^{\text{LLFs}} = (1 - \delta(u_h))\phi_i^{\text{LLF}} + \delta(u_h) K_{n_i\tau} \phi^K + \frac{\delta(u_h)}{3} \phi^K, \quad (50)$$

in the predictor step, and

$$\Phi_i^{\text{LLFs}} = (1 - \delta(u_h))\Phi_i^{\text{LLF}} + \delta(u_h) K_{n_i\tau} \Phi^K + \frac{\delta(u_h)}{3} \phi^K + \delta(u_h) \sum_{j \in K} m_{ij}^G \frac{u_j^* - u_j^n}{\Delta t^n}, \quad (51)$$

in the corrector step, having denoted by  $m_{ij}^G$  the standard  $P^1$  Galerkin mass matrix. The advantage of (50) and (51) w.r.t. the schemes proposed in [3, 60] is that both  $\delta$  and  $m_{ij}^G$  are scalar quantities, and so at the small additional cost of purely scalar operations a genuinely linear high order scheme is recovered in smooth regions, yielding improved accuracy. Note also that, differently from what is done in [57], for better consistency we have chosen to keep here the full Galerkin mass matrix in the last term instead of using a centered distribution of the integral of the predicted time increment (cf. [57] for more).

Concerning the definition of the smoothness sensor  $\delta(u_h)$  we have used the one proposed in [60] (see also [3]) :

$$\delta(u_h) = \min \left( 1, \frac{h_K^2 \|E\|_{L^\infty(K)} \|\vec{v}\|_{L^\infty(K)}}{|\varphi_E| + 10^{-12}} \right) \quad (52)$$

where  $h_K$  is the local mesh size,  $E$  is the NLSW entropy (see e.g. [38, 39, 75, 76], and cf. equation (10) for the notation)

$$E = H \left( \frac{1}{2} gH + gB + k \right),$$

while  $\varphi_E$  is an approximation of the entropy residual obtained as

$$\varphi_E = v_K \cdot \phi^K \quad (\text{or } \varphi_E = v_K \cdot \Phi^K \text{ in the corr. step}),$$

with  $v_K$  a local averaged value of the symmetrizing variables  $\partial E / \partial u$  (cf. section §4.1 and see [37, 38, 39, 60, 75, 76] for details).

## 4.5 Wet/dry front handling and implementation details

In this last paragraph, we discuss the treatment of the wet/dry fronts, and give some details regarding the implementation of the scheme. As in [60], dry fronts are detected by means of two cut-off constants  $C_H$  and  $C_{\vec{v}}$ . A node is flagged as dry if  $H \leq C_H$ . The value of this threshold is set to  $C_H = 10^{-12}$ . In practice, we have set

$$\phi^K = 0 \quad \text{if } H_j^n \leq C_H \quad \forall j \in K \quad (\text{and } \Phi^K = 0 \quad \text{if } H_j^n, H_j^* \leq C_H \quad \forall j \in K \text{ in the corr. step}).$$

The second cut-off is instead used to avoid division by zero when computing the local speed  $\vec{v}$ . In all computations, we have set  $\forall i \in \Omega_h$

$$\vec{v}_i = \begin{cases} \frac{\vec{q}_i}{H_i} & \text{if } H_i > C_{\vec{v}} \\ 0 & \text{otherwise} \end{cases}.$$

Furthermore, we have enforced  $\vec{q}_i = 0$  if  $\vec{v}_i = 0$ . As in [60], the  $C_{\vec{v}}$  cut-off is also used to benefit from the result of proposition 4.6. In particular, following [60], if  $\phi_H^K$  denotes the first component of  $\phi^K$ , and setting  $H_{\min} = \min_{j \in K} H_j^n - |\phi_H^K|$ , we have set in the computation of the nonlinear LLF distribution coefficients (cf. section §4.3)

$$\ell_m = \begin{cases} \mathbf{l}_m & \text{if } H_{\min} > C_{\vec{v}} \\ e_m & \text{otherwise} \end{cases}, \quad r_m = \begin{cases} \mathbf{r}_m & \text{if } H_{\min} > C_{\vec{v}} \\ e_m & \text{otherwise} \end{cases},$$

with  $\mathbf{l}_m$  and  $\mathbf{r}_m$  the  $m$ -th left and right eigenvector of the locally linearized NLSW flux Jacobian projected on the local velocity direction. In this way, the limiting is performed equation by equation in elements close to the wet/dry fronts, thus falling in the hypotheses of proposition 4.6. Moreover, to avoid  $\alpha_{\text{LF}}$  to be zero or too small in these elements, we have set in practice (cf. equation (45))

$$\alpha_{\text{LF}} = \max_{j \in K} l_j \left( \frac{1}{2} \max_{j \in K} (\|\vec{v}_j\|_{L^\infty} + c_j^n) + \frac{h_K}{L_{\text{ref}}} \right),$$

where the reference length  $L_{\text{ref}}$  is given by

$$L_{\text{ref}} = \max_{i,j \in \Omega_h} \|\vec{x}_i - \vec{x}_j\|.$$

As in [60], to take into account the presence of dry fronts, the smoothness sensor has been modified as

$$\delta^*(u_h) = \delta(u_h) e^{-\frac{a h_K^2}{L_{\text{ref}}^2} \left( \frac{\max_{j \in \Omega_j} H_j^n - C_{\vec{v}}}{\max(C_H, H_{\min} - C_{\vec{v}})} \right)^2},$$

with  $\delta(u_h)$  given by (52). We have experimentally found that the results become insensitive to the choice of  $a$  below the value  $a = 1/10$ , which is the value used in all computations. Concerning the choice of the cut-off  $C_{\vec{v}}$ , the interested reader can refer to the numerical tests reported in [60]. Here, following the reference, we have used the mesh dependent value

$$C_{\vec{v}} = \left( \frac{h}{L_{\text{ref}}} \right)^2.$$

To preserve the C-property in elements containing dry nodes, we have adopted the bathymetry re-definition suggested in [17], namely

$$B_i = \begin{cases} B(x_i, y_i) & \text{if } H_i > C_H \\ \max_{j \in K, H_j > C_H} \eta_j & \text{otherwise} \end{cases}.$$

Lastly, concerning the choice of the approximation and of the quadrature formulas, we have used a simple approximation in conserved variables in all the tests with the exception of those related to the analysis of the constant energy flows. In this case a more expensive interpolation of the energy variables (cf. section §4.2.2) is necessary. Concerning the quadrature formulas, the elemental fluctuations (23) have been computed using standard Gauss line quadrature to evaluate the contour integrals on triangles boundaries. If not stated otherwise, the standard two points formula has been used (see e.g. [31], chapter 8), which allows to exactly evaluate the  $H_h^2$  term, as required by proposition 4.4. The volume integrals of the friction source term

and of the time increment in the residual (28) have been computed by the second order three point formula using the nodes of the triangle, which is exact for the time increment integrals. Concerning the bathymetry, if not stated otherwise, the integral of  $H_h \nabla b_h$  is evaluated exactly for linear  $H_h$  and  $b_h$ , as required by proposition 4.4. In the super-consistency tests, this integral has been computed as (cf. section §4.1, Lemma 4.1)

$$\int_K H_h \nabla b = |K| \sum_{q=1}^{v_q} \bar{\omega}_q H_h(\vec{x}_q) \nabla b(\vec{x}_q),$$

with  $\nabla b$  known analytically and with quadrature formulas of different accuracy (cf. [30]).

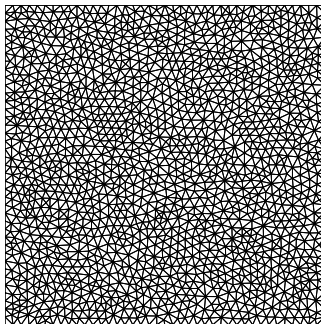


Figure 3: Unstructured grid topology

## 5 Numerical tests

We discuss in this section the numerical results obtained with the LLFs scheme. Three families of tests are considered : flows on flat bathymetries, C-property tests, wetting/drying tests. The first class of problems will be used to assess the accuracy and shock capturing capabilities of the scheme in absence of bathymetric variations, and to compare computational times with the scheme of [60]. The two different schemes are coded in the same platform, allowing fair comparisons. All the computations have been run on a portable 2.66 Ghz Intel Dual Core PC with 4 GB of RAM memory. In all the computations, the time step has been set to (cf. sections §2 and §4.3)

$$\Delta t = \min_{i \in \Omega_h} \min \left( \frac{|C_i|}{\sum_{K \in K_i} (\alpha_{\text{LF}} + \max_{j \in K} g c_f(u_j^n))}, \min_{K \in K_i} \frac{|K|}{3\alpha_{\text{LF}} + \max_{j \in K} g c_f(u_j^n)} \right),$$

for the LLFs proposed here, while for the scheme of [60] we have used the maximum time step allowed for depth non-negativity :

$$\Delta t = 2 \min_{i \in \Omega_h} \min_{K \in K_i} \frac{|K|}{3\alpha_{\text{LF}}}.$$

In both cases, the time step is computed by evaluating  $\alpha_{\text{LF}}$  using the solution at time  $t^n$ .



## 5.1 Flows on flat bathymetry

### 5.1.1 Vortex transport, accuracy and efficiency

We consider the traveling vortex problem proposed in [60] to test the accuracy of the scheme. The exact solution consists of the advection along the  $x$ -direction of a vortex described by an analytical perturbation of  $H$  and  $\vec{v}$  about a constant state  $H = 10 [m]$  and  $\vec{v} = (6, 0)$  (see [60] for details). The spatial domain is the square  $[0, 1]^2$ , with periodic boundary conditions in the  $x$  direction. Numerical solutions are computed at time  $T = 1/6$ , corresponding to the moment in which the vortex has crossed the whole domain and got back to its initial position, due to the periodic boundary conditions. The computations have been performed on 4 unstructured grids with the topology shown on figure 3. The coarsest mesh has size  $h \approx 1/56$ . The other 3 meshes have been generated independently, halving the mesh size at each step.

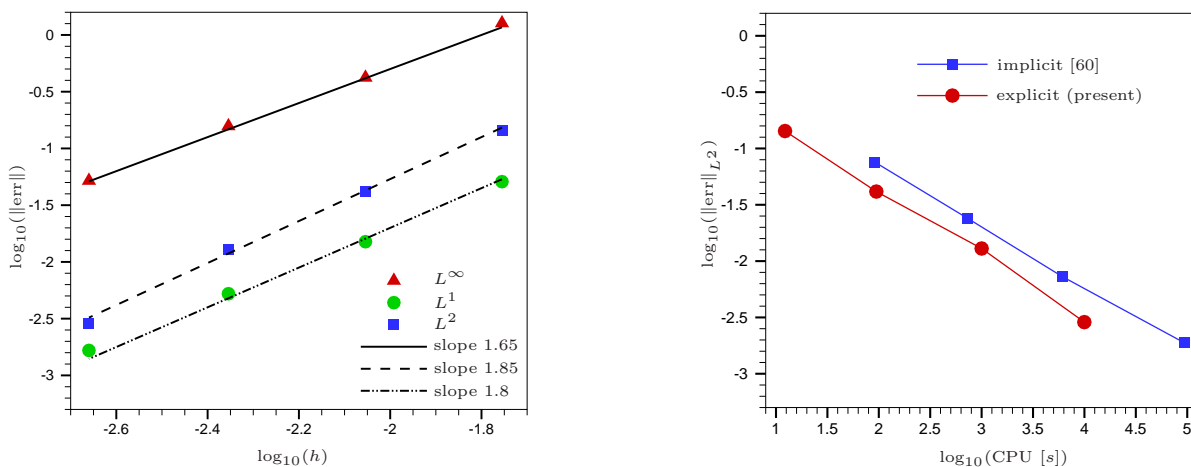


Figure 4: Vortex transport. Grid convergence for the explicit scheme (left), and error-CPU time comparison between explicit and implicit scheme [60]

We present the results on figure 4. The left picture, in particular, shows the grid convergence history of the error, confirming the theoretical second order of accuracy. The right picture, instead, shows the evolution of the  $L^2$  norm of the error w.r.t CPU time, measured in seconds. The plots compares the results of the explicit LLFs scheme proposed here, with those of the implicit scheme of [60]. The plot shows that on a given mesh, the scheme of the reference provides a lower error, however, the computational times required are roughly 6 or 7 times larger than those of the explicit scheme proposed here. In particular, for a given error level, the explicit scheme we propose is roughly 4 times faster than the reference.

### 5.1.2 Asymmetric break of a dam

This test, taken from [58, 70], consists of the asymmetric break of dam separating two basins with water depths of 5 and 10 meters. The dam is contained in the computational domain  $[0, 200]^2$ , and the breaking is initially placed at  $x = 95 [m]$ . Reflective boundary conditions



are used on all boundaries. A sketch of the geometry and a close up view of the mesh are reported on figure 5. The mesh size is  $h \approx 2[m]$ , the mesh contains 19274 triangles and 9899 nodes. We refer to [58, 70] for more details on the test set up. Computations have been run up to time  $T = 7.2 [s]$  with the explicit LLFs scheme proposed here and with the scheme of [60]. Results are summarized on figures 6 and 7.

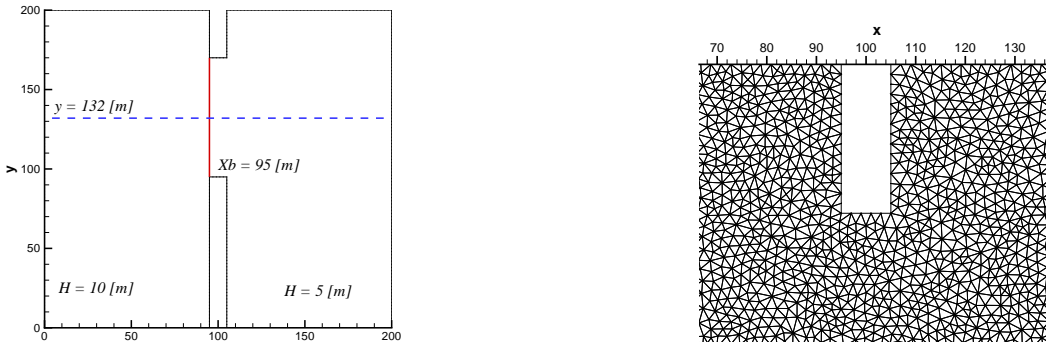


Figure 5: Asymmetric dam break. Left: computational domain. Right: mesh close up ( $h \approx 2$ )

In particular, figure 6 shows 3D visualizations of the free surface level computed with the explicit LLFs scheme proposed here, and the one obtained with the implicit scheme of [60]. The scheme of [60] provides slightly sharper shocks and stronger rarefactions around the corners, as shown by presence of more pronounced kinks in the depth contours in the region of the trough forming in the emptying basin from the interaction of the two corner rarefactions. Similar observation can be made by looking at the depth and Froude number profiles on figure 7. The plots compare the data along the line  $y = 132[m]$ , which is roughly the position of the center of the depth trough due to the interaction of the corner rarefactions (cf. left picture on figure 5). We see again that the implicit scheme provides slightly sharper shocks, while the effect of the stronger corner rarefactions are particularly visible in the deeper trough in depth and higher peak in Froude number at  $x \approx 55[m]$ . The differences between the two solutions are however not striking. On the other hand, the computational times required to obtain the solutions are  $53.3 [s]$  for the explicit scheme proposed here, and  $331 [s]$  for the scheme of [60], which is more than 6 times larger. This shows again the improvement brought by the scheme proposed here.

## 5.2 C-property tests

### 5.2.1 Lake at rest solution

We start by considering the lake at rest solution. In particular we consider 2 benchmarks allowing to verify numerically proposition 4.4. On the spatial domain  $[0, 2] \times [0, 1]$ , let the bathymetry be defined as :

$$b(x, y) = b_0 e^{\psi(x, y)}.$$

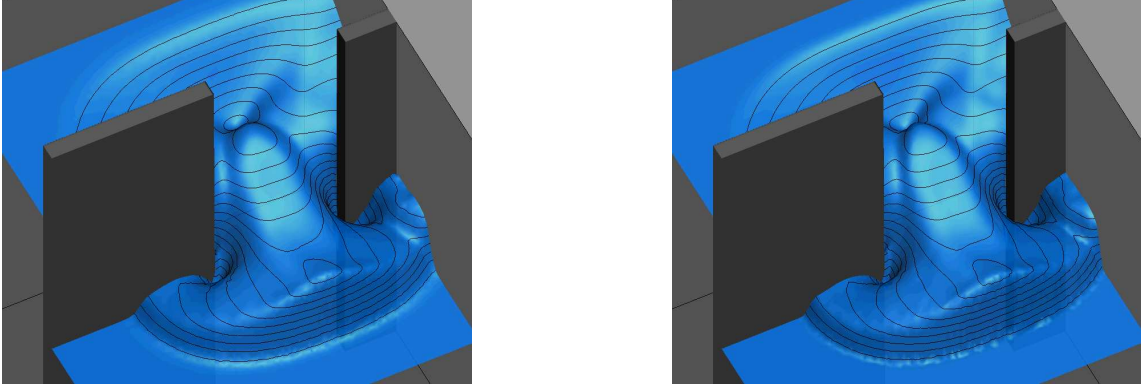


Figure 6: Asymmetric dam break: water height at time  $t = 7.2$ . Left: explicit scheme. Right: implicit scheme of [60]

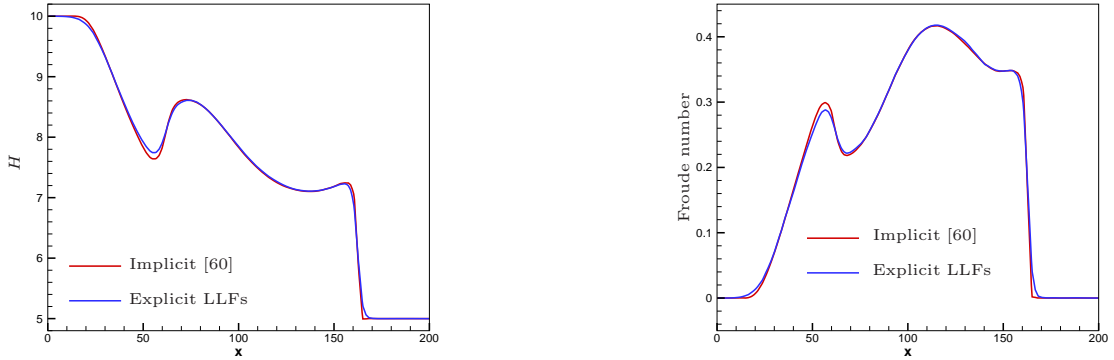


Figure 7: Asymmetric dam break: water height (left) and Froude number (right) at time  $t = 7.2$ . Data extracted along the line  $y = 132[m]$

In particular, we consider two cases. The first definition is obtained with

$$b_0 = 0.8, \quad \text{and} \quad \psi = -5(x - 0.9)^2 - 50(y - 0.5)^2.$$

This definition is used in a large number of references to verify the C-property (see e.g. [78, 70, 45] and references therein). We also consider the following non-smooth case, proposed in [58] :

$$b_0 = 0.6, \quad \text{and} \quad \psi = \begin{cases} \sqrt{(x - 0.9)^2 + (y - 0.5)^2} & \text{if } \vec{x} \in [0.9, 1.1] \times [0.3, 0.7] \\ -5(x - 0.9)^2 - 50(y - 0.5)^2 & \text{otherwise} \end{cases}.$$

For both definitions, we have run the explicit LLFs scheme with the initial lake at rest solution  $(\eta, q) = (1, 0) \forall \vec{x} \in \Omega$ . As foreseen by proposition 4.4 the results are constant up to machine accuracy. To visualize this fact, we consider a perturbation of the lake at rest state given by  $\vec{q} = 0$ , and

$$\eta_0(x, y) = \begin{cases} 1.01 & \text{if } x \in [0.05, 0.15] \\ 1 & \text{otherwise} \end{cases}.$$

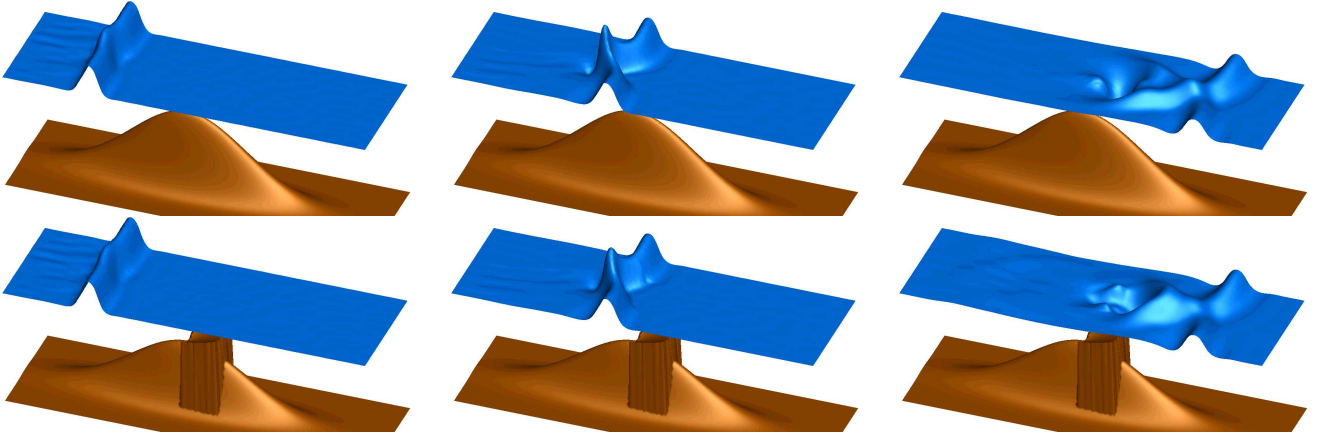


Figure 8: Small perturbation of lake at rest: 3d view of total water height at times  $t = 0.12$  (left),  $t = 0.24$  (middle), and  $t = 0.48$  (right). Explicit LLFs scheme with smooth (top row) and non-smooth (bottom row) bathymetry (rescaled for better visualization)

With this initial condition we have run the LLFs scheme on a mesh with typical size  $h \approx 1/100$ , containing 20037 nodes and 39472 triangles. The results obtained on the smooth and non-smooth bathymetry are visualized in terms of 3d view of the free surface  $\eta$  (and bathymetry  $B$ ) on figure 8, of 1d line plots of the data extracted along the line  $y = 0.5$  (both  $\eta$  and bathymetry) on figure 9, and of contours of the free surface  $\eta$  on figure 10. In all pictures, the bathymetry is rescaled to allow properly visualizing the free surface.

The results show the perfect preservation of the initial lake at rest state away from the perturbation. The features captured by the scheme proposed here match quite well those of other results presented in published literature (see e.g. [45, 58, 60, 70, 78] and references therein), confirming both the well balanced character and the accuracy of the scheme proposed.

### 5.2.2 Constant energy flows

In this section we check numerically the super-consistency property of proposition 4.2. Following [56], we consider, on the square  $[0, 25]^2$ , the pseudo one-dimensional bathymetry

$$b(x, y) = \begin{cases} f(x) & \text{if } x \in [8, 12] \\ 0 & \text{otherwise} \end{cases} .$$

The function  $f(x)$  is chosen with increasing regularity. We start with the  $C^0$  definition

$$f(x) = 0.2 - (x - 10)^2/20, \quad (53)$$

giving a  $H^1$  bathymetry, and then consider

$$f(x) = 0.2 \sin^{2l}(\pi(x - 8)/4), \quad l \in \{1, 2, 4\}, \quad (54)$$

yielding  $H^{2l}$  bathymetries with  $C^{2l-1}$  continuity. We compute initial nodal values from (16), with  $\mathcal{E}_0 = 22.06605$  and  $\vec{q}_0 = (4.42, 0)$ , and run unsteady computations until time  $T = 0.1$  [s] on four nested unstructured grids with the topology shown on figure 3. To fit into the hypotheses of proposition 4.2, in all computations we use the spatial approximation  $\mathcal{F}_h = \mathcal{F}(u(v_h, b))$  and

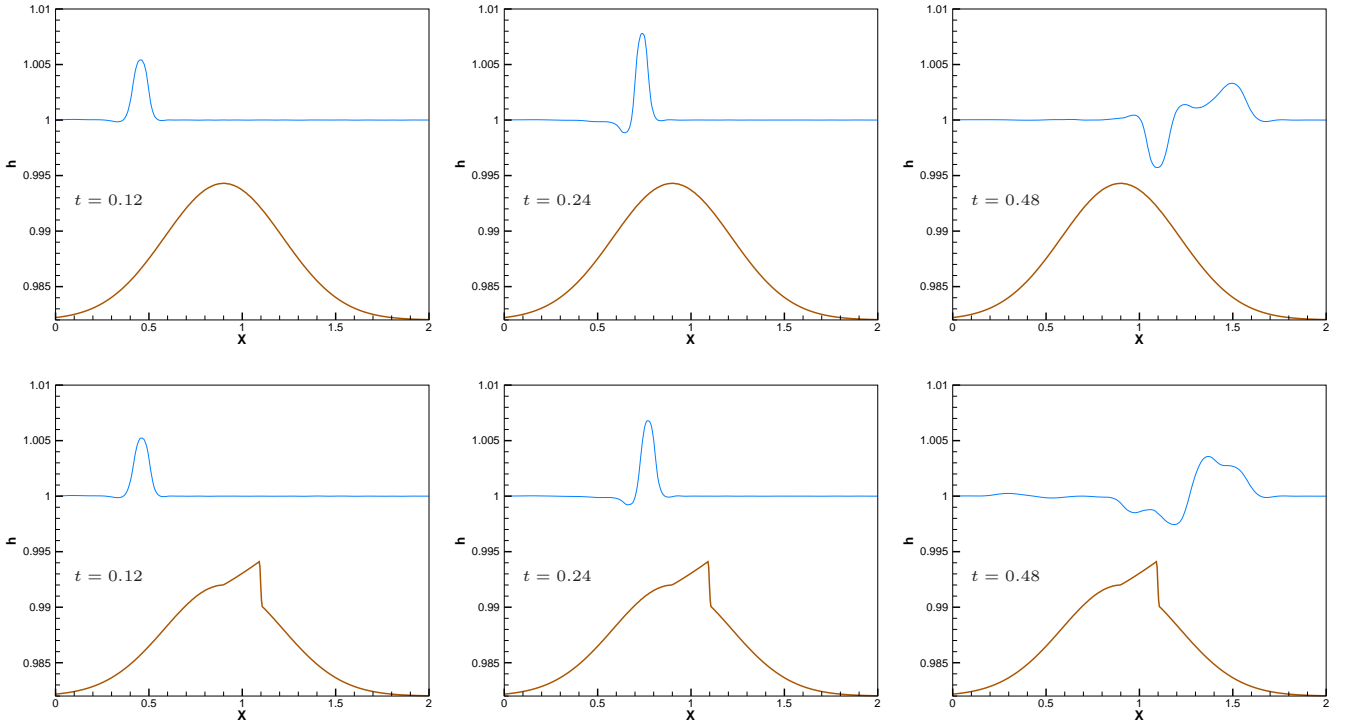


Figure 9: Small perturbation of lake at rest: line plot of total water height at times  $t = 0.12$  (left),  $t = 0.24$  (middle), and  $t = 0.48$  (right). Explicit LLFs scheme with smooth (top row) and non-smooth (bottom row) bathymetry. Data extracted along the  $y = 0.5$  line

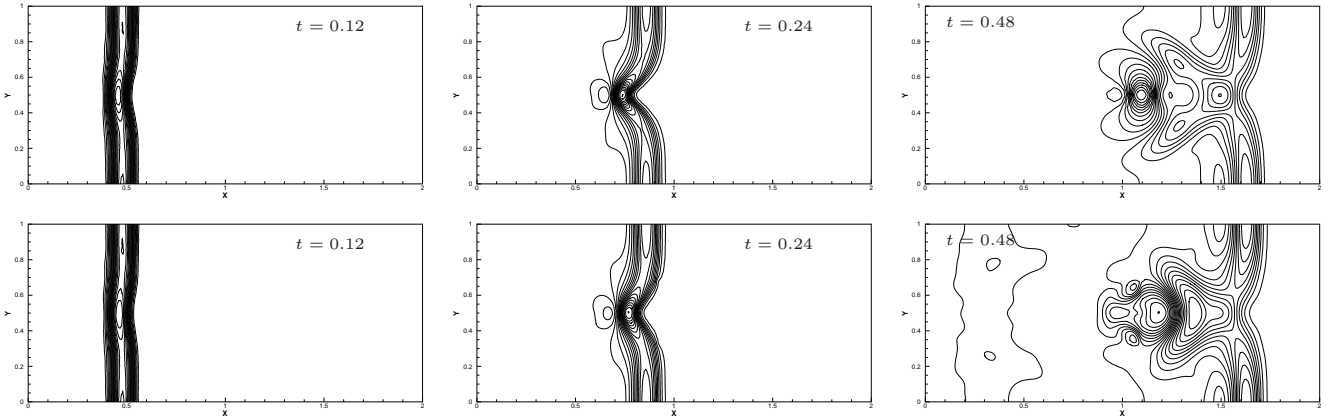


Figure 10: Small perturbation of lake at rest: contour plot of total water height at times  $t = 0.12$  (left),  $t = 0.24$  (middle), and  $t = 0.48$  (right). Explicit LLFs scheme with smooth (top row) and non-smooth (bottom row) bathymetry.

$\mathcal{S}_h = \mathcal{S}(u(v_h, b), \nabla b)$ , with  $b$  the analytical bathymetry, and  $v_h = (\mathcal{E}_h, \vec{q}_h)$  linear. The runs are repeated with different quadrature strategies. Edge integrals are computed with formulas exact for polynomial degrees  $p_f = 1, 3,$  and  $11$ . Two-dimensional integration formulas on triangles exact for polynomial degrees  $p_v = 1, 4,$  and  $6$  are taken from [30]. We discuss the

results obtained with the four strategies :

$$Q_1 : p_f = 1, p_v = 1; \quad Q_2 : p_f = 3, p_v = 4; \quad Q_3 : p_f = 11, p_v = 4; \quad Q_4 : p_f = 11, p_v = 6.$$

Recall that for a  $H^{p+1}$  bathymetry, with sufficiently large  $p$ , the accuracy measured for a finite time computation should be, according to proposition 4.2

$$\epsilon = \mathcal{O}(h^l) \quad \text{with } l = \min(p_f + 1, p_v + 1), \quad (55)$$

which gives for the four quadrature strategies considered the theoretical slopes

$$l_1 = 2, \quad l_2 = 4, \quad l_3 = 5, \quad l_4 = 7.$$

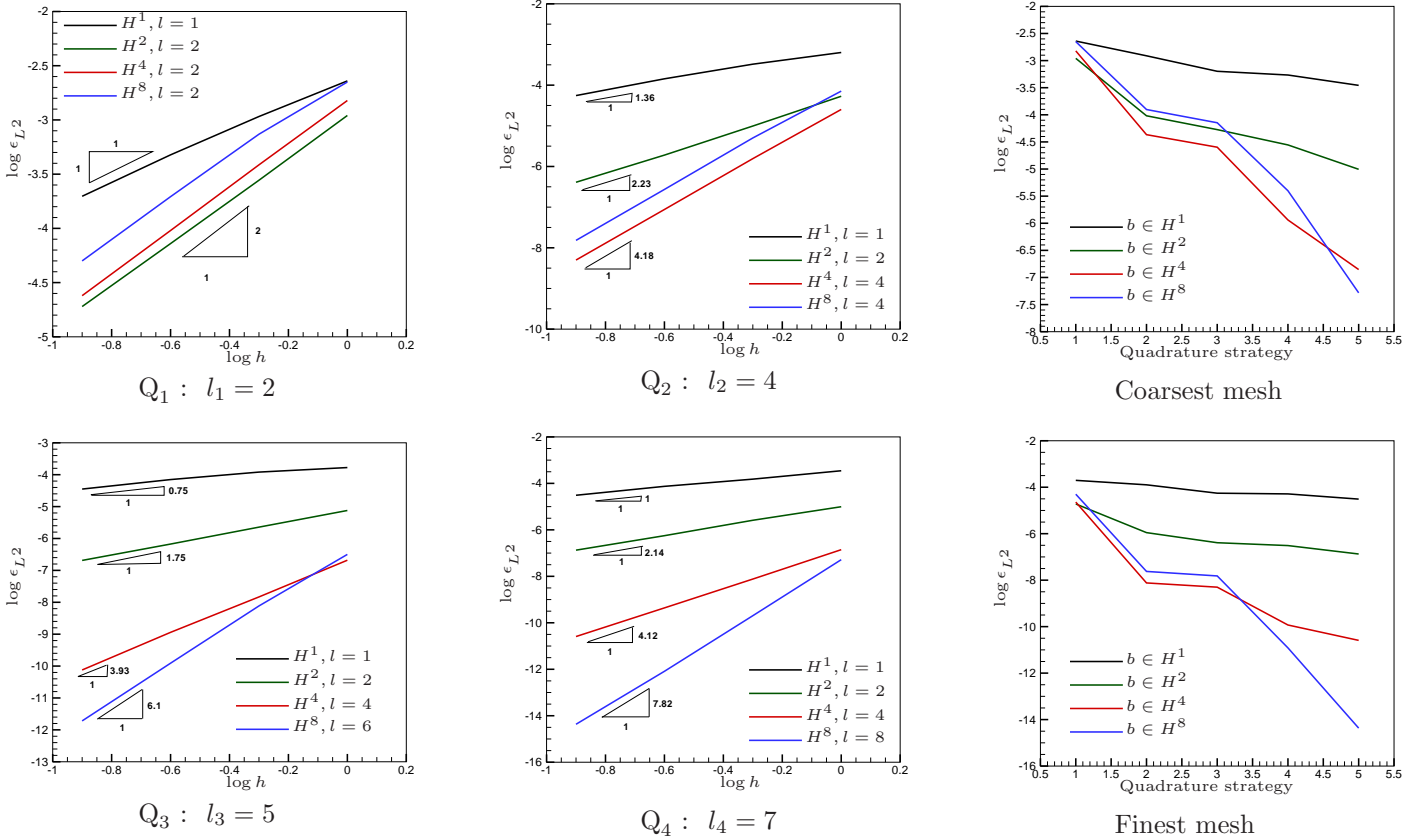


Figure 11: Super consistency for constant energy flows : grid convergence (left and middle columns) and quadrature convergence (rightmost column) of the depth error.

The numerical results are shown on figure 11. In the figure, the first four pictures on the first and second column represent the grid convergence of the depth at  $t = 0.1$  [s]. The last two pictures show the error convergence on a fixed grid (the coarsest on top, the finest on the bottom) when increasing the accuracy of the quadrature. Below each picture, we have reported the quadrature strategy used, and the expected theoretical slope. We recall that the explicit LLFs scheme used is formally second order accurate.

We can see that the discrete solution at time  $T = 0.1$  [s] converges with rates higher than two if the bathymetry is regular enough. In particular, in this case the convergence rates observed are even better than those foreseen by proposition 4.2 and given by  $\mathcal{O}(h^{\min(p_f+1, p_v+2)})$  instead of just  $p_v + 1$ . Moreover, as anticipated in remark 4.3, when the regularity of  $b(x, y)$  gets lower the accuracy saturates at an order  $l$  given by (42). In particular, only first order of accuracy is observed if  $b \in H^1$ . The degree of super-consistency depends also on the quadrature formula. In particular the rightmost column on figure 11 shows the error reduction for the different quadrature formulas on the coarsest and finest meshes used in the grid convergence study. Note that the  $x$ -axis in these pictures has no quantitative meaning, except that the accuracy of the edge, and/or volume quadrature increases from left to right. These plots confirm that the error indeed converges to zero (towards exact preservation) if the quadrature accuracy is increased. The smoother the bathymetry, the faster the convergence.

As already observed in remark 4.3, the disappointing fact is that with this approach only first order of accuracy is obtained when the bathymetry is only continuous, despite the fact that the scheme is formally second order accurate. This is a direct consequence of the fact that, in order to provide the super-consistent results, this approach relies directly on the availability of an analytical bathymetry which is sufficiently regular within the element. As remarked earlier, given the one-dimensional nature of these flows, a simple way around the problem is to use structured grids, or local structured layers, or any other (possibly anisotropic) adaptation technique placing element edges right across the discontinuities in the bathymetry.

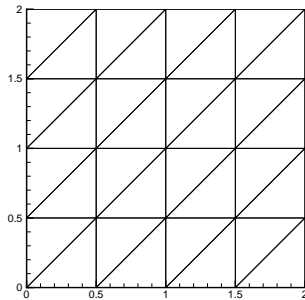


Figure 12: Structured grid for the super-consistency test

To show the potential for such type of technique, we repeat the previous test for the case in which the bathymetry is only  $C^0$  on the structured grids shown on figure 12. We repeat the grid convergence tests taking care of aligning the mesh edges with the lines of discontinuity of the first derivative of  $b(x)$ . We test the quadrature strategies  $Q_2$  and  $Q_3$  which should give fourth and sixth order super-consistent results when interpolating the steady invariants.

The results of the test are summarized on table 1 where we have reported the errors and convergence rates for the interpolation in steady invariants, and in conservative variables. Note that in the second case, the strategy  $Q_2$  provides already exact integration of the polynomials involved, so only one column is reported, the results in the second case being identical. The results in the table show that in this case the expected super-consistency is still observed and



	$Q_2 - u(v_h, b)$	$Q_3 - u(v_h, b)$	$Q_2 - u_h$
25/50	1.452714e-07	3.698282e-10	3.35738e-04
25/100	9.508237e-09	4.450410e-12	8.85116e-05
rate	3.947	6.399	1.930
25/200	6.584230e-10	4.688134e-14	2.36592e-05
rate	3.865	6.591	1.913

Table 1: Super consistency for constant energy flows : structured grid results for a  $C^0$  bathymetry.  $L^2$  error on the total energy at time  $t = 0.1$  for interpolation in steady invariants (first two columns), and in conserved variables (last column).

with the rates foreseen by the theoretical analysis presented.

Finally, to show *visually* the benefits of using our residual approach in conjunction with the approximation in total energy variables, we consider the perturbation of a pseudo one dimensional flow on a bathymetry representative of a ribbed channel. In particular, we consider on the square  $[0, 25]^2$  a series of 5 ribs starting at  $x = 1.5, 6, 10.5, 15, \text{ and } 19.5 [m]$ . Two rib shapes are considered : a  $C^0$  piecewise parabolic definition obtained by shifting (53), and a  $C^1$  definition obtained shifting (54) with  $l = 1$  ( $\sin^2$  ribs). The bathymetries obtained are visualized on figure 13. In the pictures we also show the mesh used in the computations, containing 6553 nodes and 12784 elements ( $h \approx 25/80$  ).

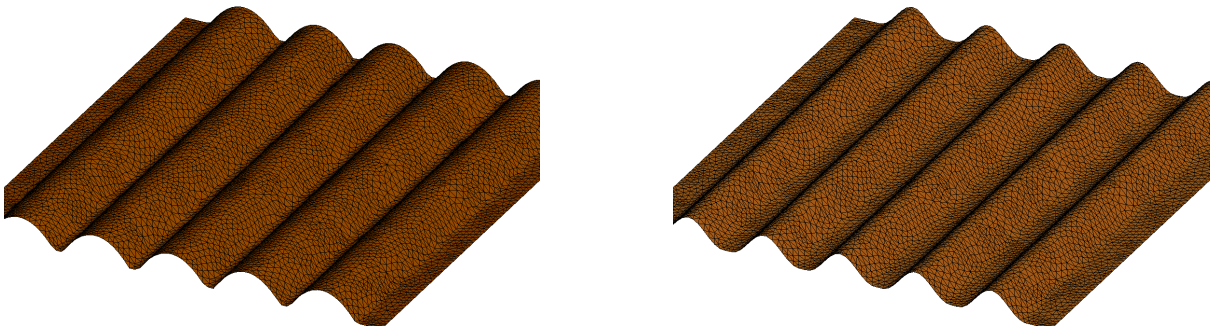


Figure 13: Ribbed channel bathymetries. Left :  $C^0$  piecewise parabolic. Right :  $C^1$  piecewise  $\sin^2$

We then consider a supercritical solution obtained from (16) setting the value of the gravity acceleration to  $g = 1 [m/s^2]$ , and with  $\mathcal{E}_0 = 6$ , and  $\vec{q}_0 = (5.65685, 0)$ . We denote by  $\eta_{\text{steady}}$  the free surface level associated to this initial solution. Following [51], we perturb the total energy  $\mathcal{E}$  in the box  $[6.5, 7.5] \times [12, 13]$ , adding the perturbation  $\delta\mathcal{E} = 0.1$ . We then compute the evolution of the perturbation until time  $t = 2 [s]$  (recall that we have set  $g = 1 [m/s^2]$ ) with the explicit LLFs scheme proposed here, using either the standard approximation in  $[H_h, \vec{q}_h]$  variables with a piecewise linear approximation of the bathymetry  $b_h$ , or using the approximation of proposition 4.2 :  $\mathcal{F}_h = \mathcal{F}(u(v_h, b))$  and  $\mathcal{S}_h = \mathcal{S}(u(v_h, b), \nabla b)$ , with  $b$  the analytical bathymetry, and  $v_h = (\mathcal{E}_h, \vec{q}_h)$  linear.

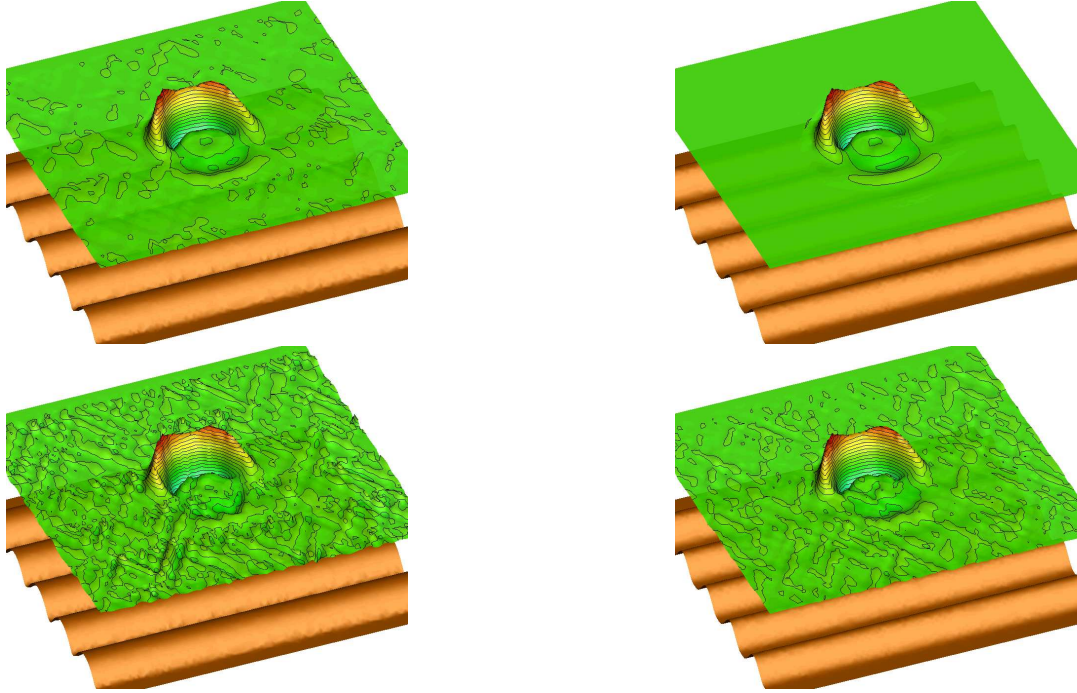


Figure 14: Perturbation of constant energy state states. Free-surface perturbation  $\eta - \eta_{\text{steady}}$  at time  $t = 0.4s$ . Left pictures :  $C^0$  bathymetry (discontinuous derivative). Right pictures :  $C^1$  bathymetry. Top : approximation in  $[\mathcal{E}, \vec{q}]$  with analytical  $b$ . Bottom : approximation in  $[H, \vec{q}]$  with linear  $b_h$

The results obtained are reported on figures 14 and 15. In particular, on figure 14 we report an exaggerated 3d view of the free surface level perturbation  $\eta - \eta_{\text{steady}}$  for the  $C^0$  bathymetry (left column) and for the  $C^1$  case (right column). The top results are obtained with the approximation in total energy variables, verifying the hypotheses of proposition 4.2, while the results on the bottom row are obtained with a standard approximation in  $[H_h, \vec{q}_h]$  variables. The figures clearly show that even if formally only first order accurate, even for a  $C^0$  bathymetry, the approximation in total energy variables provide a much better preservation of the undisturbed steady state. The results obtained for the  $C^1$  bathymetry show that the preservation of the steady solution is practically perfect when using total energy variables, the deviation being of the order of  $10^{-8}$ . To confirm these observations, on figure 15 we report in the top row the one-dimensional distribution of the free surface perturbation  $\eta - \eta_{\text{steady}}$  along the line  $y = 12.5$ , and in the bottom row the 1d plot of  $\eta - \eta_{\text{steady}}$  in *all the mesh points in the box*  $[0, 6] \times [0, 25]$ . The left column shows the results obtained on the  $C^0$  bathymetry, while the results of the  $C^1$  case are reported on the right column. The top pictures show that the result obtained using total energy variables (red curves) is smoother in both cases, with much smaller deviation from zero in the undisturbed region (left most and rightmost ends of the domain) in the  $C^0$  case, and clearly no visible deviation from the steady state in the  $C^1$  case. The data reported in the pictures on the bottom confirm this last observation : in the  $C^0$  case (left picture) the results in total energy variables (red curve) are much more clustered around zero, the largest peaks being of the order of one third of those of the standard approximation. On the other hand, the right picture confirms that in the  $C^1$  case the preservation of the



steady state is practically perfect.

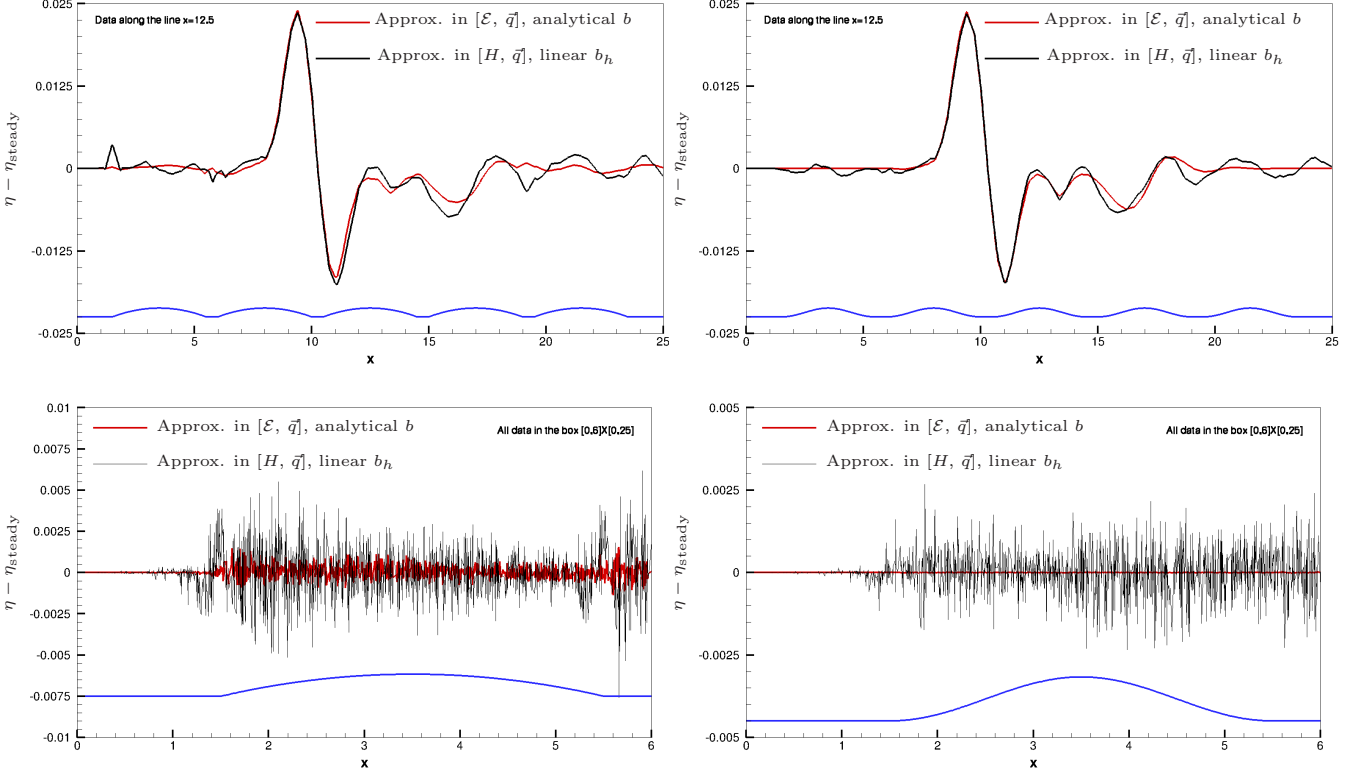


Figure 15: Perturbation of constant energy state states. Free-surface perturbation  $\eta - \eta_{\text{steady}}$  at time  $t = 0.4s$ . Left :  $C^0$  bathymetry (discontinuous derivative). Right :  $C^1$  bathymetry. Top : data extracted along the line  $y = 12.5$ . Bottom : all data in the “unperturbed” box  $[0, 6] \times [0, 25]$ .

**Remark 5.1** (Choice of the approximation). *The results shown, as well as the theoretical developments, are based on the assumption that an analytical bathymetry is available, and that its exact form is used in the discretization in conjunction with a direct interpolation of flux and total energy.*

*The use of this interpolation leads to a scheme which is substantially more expensive, due to the need of recovering the depth and velocity from the total energy in every mesh node and quadrature point (cf. section §4.2.2). So in practice, simpler choices, such as interpolating the conserved variables, should be preferred if these equilibria are not of interest.*

**Remark 5.2** (Bathymetry representation). *For situations in which the preservation of constant energy flows is important, there is the question of the availability of a smooth analytical bathymetry. This is of course questionable. Nevertheless, given the uncertainties in bathymetric data, it would not be unthinkable, when this type of flow is relevant, to replace irregular experimental bathymetries with a regularized  $C^1$  approximation (obtained under some physical constraints of e.g. equal total water volume at rest).*

### 5.2.3 Flows in sloping channels with friction

We present a verification of proposition 4.5. On the square  $[0, 1]^2$ , we consider a rotated solution of the type (19), with a bathymetry obtained as  $b = b_0 - \xi x^*$ , with  $x^* = (2x + y)/\sqrt{5}$ . We then compute two solutions from (19) corresponding to a sub-critical ( $\text{Fr} \approx 0.638$ ) and super-critical ( $\text{Fr} \approx 2.536$ ) flow. We then perturb the free surface level by  $\delta\eta = 0.01$  within the circle centered in  $(0.5, 0.5)$  and of radius  $r = 0.1$ . We compute the evolution of the perturbation with the explicit LLFs proposed here, using a standard approximation in  $[H_h, \vec{q}_h]$  variables, on an unstructured triangulation similar to that shown on figure 3, and containing 6553 nodes and 12784 elements ( $h \approx 1/80$ ).

A three-dimensional visualization of the evolution of the perturbation  $\eta - \eta_{\text{steady}}$  is reported on figure 16. The results clearly show the perfect preservation of the underlying steady state.

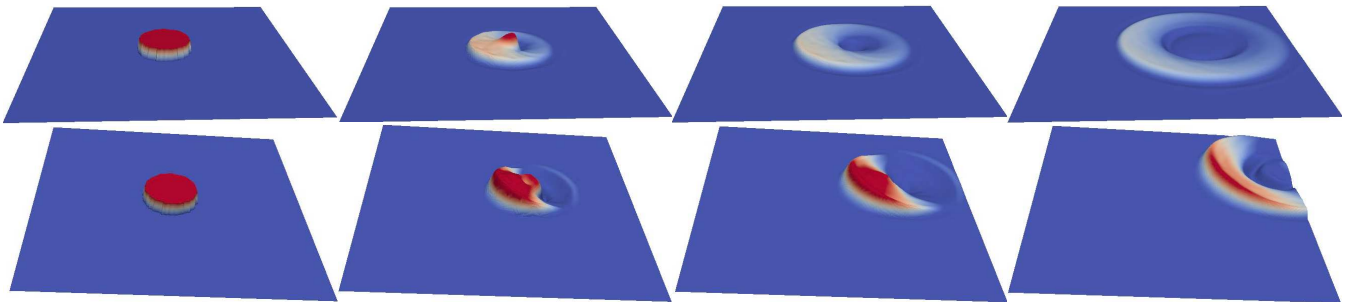


Figure 16: Flows in sloping channels with friction. Evolution of a free surface perturbation in the sub-critical (top row) and super-critical (bottom row) cases.

## 5.3 Wetting/drying tests

### 5.3.1 Thacker's oscillations in a parabolic bowl

To verify the capability of the scheme to provide an accurate and stable approximation of moving shorelines we consider the periodic oscillations of a curved free surface in a paraboloid [77]. The spatial domain considered is  $[-1.2, 1.2]^2$ , which we discretize with an unstructured triangulation with the topology shown on figure 3 containing 10113 points and 19824 elements, and with size  $h \approx 1/40$ , giving roughly 50 cells the oscillating region. Details concerning the setup and exact solutions can be found in [77]. As a first test, here we set the free surface to the analytical solution at time  $t = 0$  and let it oscillate for three full periods. We then look at the solution at times  $3T + \delta t$  for  $\delta t \in \{T/6, T/3, T/2, 2T/3, 5T/6, T\}$ .

The data along the line  $y = 0$  computed by the LLFs scheme proposed here is compared to the exact solution on figure 17. The computed solutions are nicely close to the analytical ones, even on this relatively coarse mesh. The close ups of the wetting/drying region reported on the left column also show a clean and oscillation free capturing of the moving shoreline.

Next we have compared the results obtained with the explicit LLFs scheme proposed here with those obtained with the scheme of [60]. The accuracy of the two schemes is compared in terms of  $L^1$  norm of the error on the free surface after one period in figure 18. We can see that, while the slope obtained with the scheme of [60] is closer to 2, the absolute value of the error of the scheme proposed here is significantly lower. Moreover, the CPU time required for one

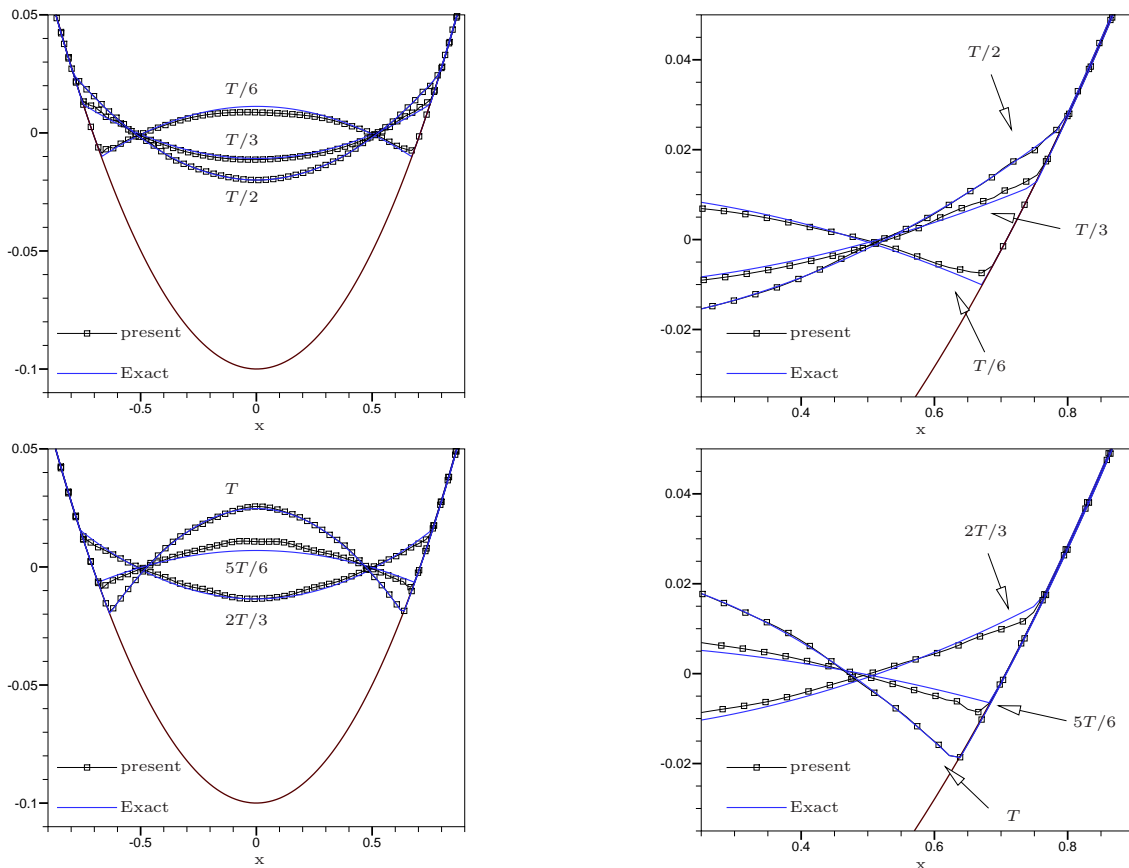


Figure 17: Periodic oscillations in a parabolic bowl [77]. Free surface level at  $3T + \delta t$ ,  $\delta t \in \{T/6, T/3, T/2, 2T/3, 5T/6, T\}$ . Data along the line  $y = 0$  compared with exact solution [77]

period on the  $h = 1/40$  mesh, containing 10113 points and 19824 elements, is of 103.470 [s] for the explicit LLFs proposed here, and of 13 minutes and 33.110 [s] for the implicit scheme of [60]. This further confirms the significant improvement brought by the present work.

### 5.3.2 Runup on a conical island

This is a standard test to validate the ability of a scheme to correctly predict long wave run up. The test aims at reproducing the experiments performed in [16]. A sketch of the test is depicted on the left on figure 19 : a solitary wave travels over an island of conical shape. The experiments of [16] have provided both point wise time series of the water level in the gauge points indicated in figure 19, and the maximum run up heights over the island. For more details the interested reader is referred to [16].

The computational domain used to reproduce the test is the rectangle of  $[-12.96, 12.4] \times [-13.8, 16.2]$  with the center of the island placed in the origin of the axes. The island's lower radius is 3.6 [m], the upper one is 1.1 [m] and the slope 1/4 with a peak height of 0.625 [m]. We have considered the case in which the water depth far from the island is  $h_0 = 0.32$  [m].

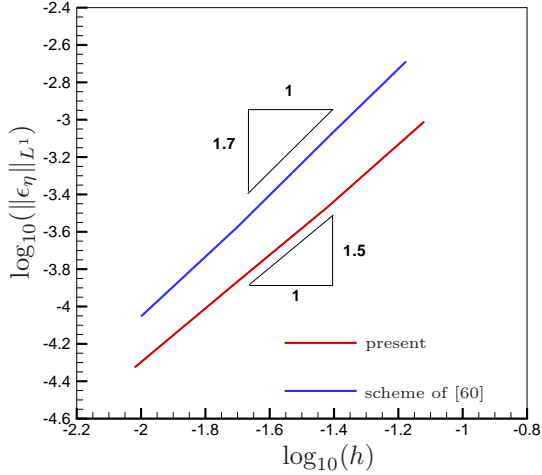


Figure 18: Periodic oscillations in a parabolic bowl [77]. Grid convergence and comparison with the scheme of [60]

The solitary wave shape imposed at the left hand of the domain is defined by the free surface perturbation

$$d\eta = A \operatorname{sech}^2(\sqrt{3A/(4h_0^3)}x),$$

with a corresponding velocity perturbation obtained from the linearized shallow water equations :  $\vec{v} = (\sqrt{g/h_0}, 0)d\eta$ . We consider here the case of a soliton of amplitude  $A = 0.2$ . The spatial domain is discretized with an unstructured triangulation. On the right on figure 19 we report a view of the mesh which is refined around the island. The largest mesh size is of 50 [cm], while the finest is of 7.5 [cm]. For better understanding, the picture also shows the lower and upper circles of the island, and the positions of the four gauges which will be used for validation :  $g6 = (-3.6, 0)$  ,  $g9 = (-2.6, 0)$  ,  $g16 = (0, -2.58)$  ,  $g22 = (2.6, 0)$ . The mesh dependent cut-off coefficients needed for the wet-dry treatment (cf. section §4.5) are computed using a local mesh size.

An exaggerated three-dimensional visualization of the run up process is presented in figures 20 and 21. The pictures show the soliton run up first on the front side of the island, then the secondary waves running around the island and meeting behind it giving the rear side run up visible in the leftmost picture on the bottom row. The rear wave then splits again into two smaller waves running back around the island.

On figure 22 we report the comparison with the experimental data [16] of the computed time history of the water height deviation  $\eta - \eta_0$  ,  $\eta_0$  being the free surface level at still water. The results of the explicit LLFs scheme proposed match quite well the experimental data, within at least the limits of the capabilities of the NLSW model. Non-hydrostatic terms are needed to better match the oscillations seen e.g. in gauge g9 after the backwash phase (around time  $t = 10$  [s]). To further confirm the soundness of the wetting/drying procedure, on figure 23 we compare the maximum run up with the data of [16]. To obtain the figure we have superposed the solutions at all times and then blanked the cells in which the minimum depth

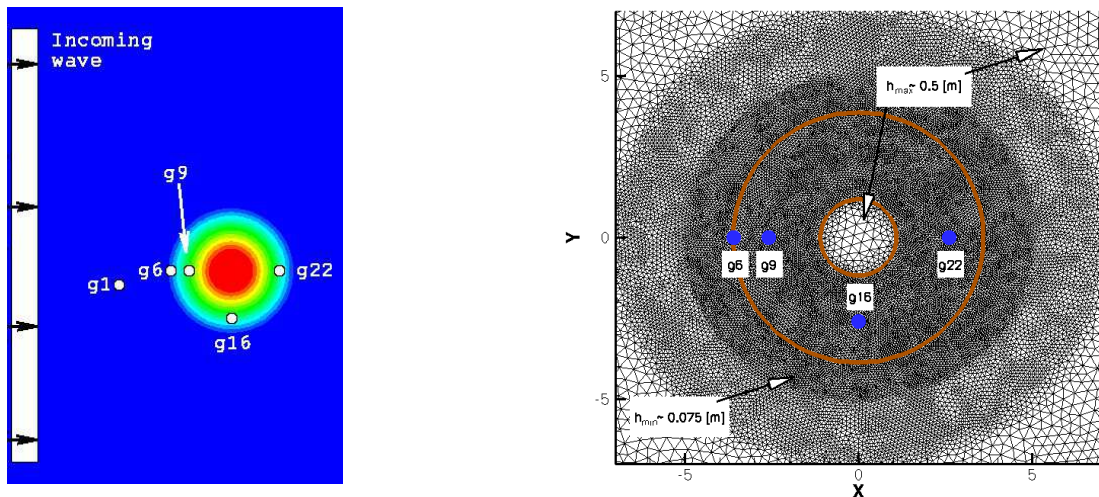


Figure 19: Run up on a conical island [16] : problem sketch (left) and computational grid (right)

is below  $10^{-5}$  [m]. This value has been set by trial and error. Using smaller values, results in the random appearance of “wet” cells not connected to the rest of the domain, as already visible on the top part of the figure. The boundary of the wet region computed by our scheme is in excellent agreement with the experiments.

### 5.3.3 Okushiri tsunami experiment

As a final application we consider the second benchmark of the third international workshop on long wave runup models : Tsunami runup onto a complex three-dimensional beach. The test is thoroughly described on the web pages [1, 2] and in [47] to which we refer for details. The test is a scaled down laboratory reproduction of the tsunami wave that hit the Okushiri island in Japan in 1993. The web page provides data files for the bathymetry of the coast of the island in the region of the Monai village, which is the one where the most damage has been observed. A three-dimensional view of the bathymetry is reported on the left picture on figure 24. On the figure, the highest point (about 50 [m] in real life) is the Monai village region while the small island in front of the coast is the Muen island reaching 10 [m] height in real life. The web site gives a 400 times scaled down geometry together with the shape of the wave used in the experiment, which is reported on the right on figure 24. In the observations [1, 2, 47] the highest runup is of 32 [m], and it occurs in the region of the Monay valley where the bathymetry is steepest. For clarity, this region is encircled in most of the two- and three-dimensional results presented in the following.

The spatial domain  $[0, 5.448] \times [0, 3.402]$  has been discretized with two meshes, both adapted to the bathymetric variations. A close up view of the meshes is reported on figure 25. The coarse mesh (left picture) contains 7000 nodes and 13720 triangles, with maximum and minimum mesh sizes given roughly by 0.1 [m] and 0.025 [m]. The fine mesh (right picture) contains 18711 nodes and 36911 triangles, with maximum and minimum mesh sizes given roughly by 0.05 [m] and 0.01 [m]. Note that the mesh size recommended in [1, 2, 47] for this test is of 0.014 [m] which would give, in absence of mesh adaptation, a number of triangles of 189000, roughly five times more than the finest mesh used here.



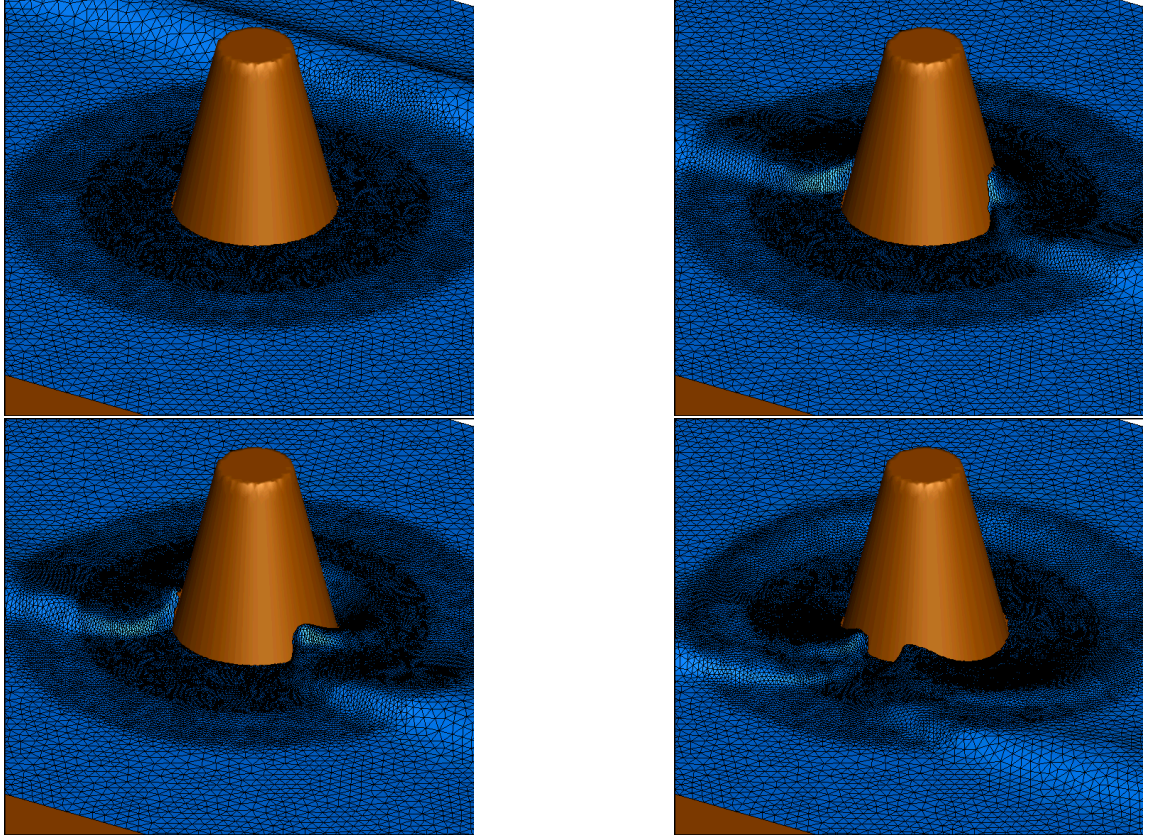


Figure 20: Run up on a conical island [16] : three-dimensional visualizations of free surface evolution (from left to right and from top to bottom).

Simulations have been run with the LLFs scheme for 25 [s], using local mesh sizes to compute the cut-off thresholds for the wetting/drying treatment (cf. section §4.5). Three-dimensional visualizations of the computed flow are reported on figures 26 and 27.

The pictures (from top to bottom and from left to right) show the initial withdrawing of the water followed by the arrival of the main wave (top row). After hitting the beach, the wave reflects, and a large wave travels toward the right to hit the steepest slopes in the region of the Monai village (bottom row, third picture from the left). The reflected wave eventually reaches and inundates the Muen island (bottom row last picture from the left). As already said, the highest runup observed is about 32 [m] (corresponding to 32/400 [m] in the scaled model), and it has been observed in the region of the Monai valley, highlighted by a yellow circle in figures 26, 27, and 29.

During the experiment, probes have been set to measure the water height history in three locations shown in the top left picture on figure 28. On the same figure, we report the comparison of the computed water height deviation from its initial value, with the experimental data provided on the web page of the workshop. Two remarks can be made. Firstly, the agreement between measured and computed heights is quite satisfactory. Second, there is no remarkable difference between the results obtained on the coarse and fine mesh in the probes.

As a last verification, we present on figure 29 the maximum runup plot obtained on the coarse and fine meshes. The plot has been obtained as described in section §5.3.2 for the conical

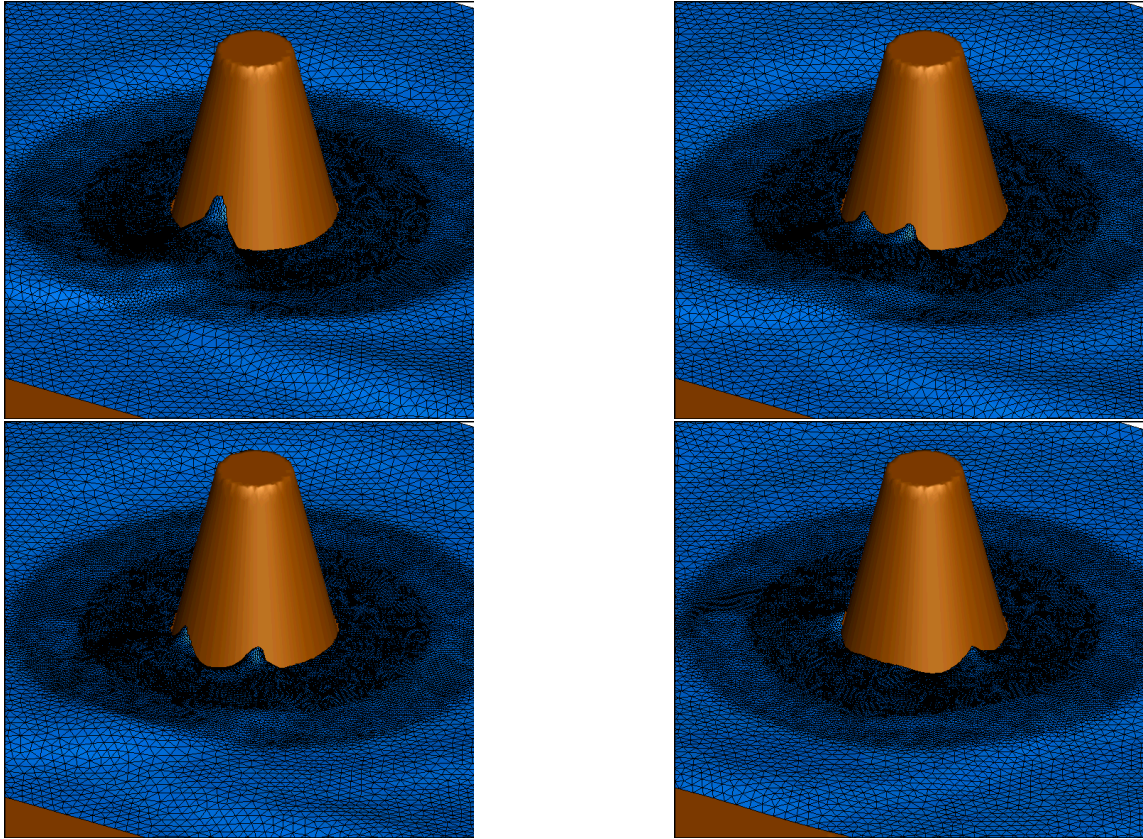


Figure 21: Run up on a conical island [16] (continued) : three-dimensional visualizations of free surface evolution (from left to right and from top to bottom).

island test. In the runup plots we have reported as a reference, the curve corresponding to the maximum experimental runup of 32 [m] (real life scale), which we recall is observed in the Monay valley (encircled region). The pictures clearly show that the higher resolution is necessary to obtain an accurate prediction of the maximum runup region. In particular, the coarse mesh results underestimate the maximum runup by roughly 10 [m], while the fine mesh results overshoot the line of the 32 [m] by one row of giving a more conservative prediction of the maximum runup of about 36 meters. As remarked already, these results are obtained with a grid containing five times less elements than what a uniform mesh following the prescriptions of [1, 2, 47] would contain, thus showing the interest of the use of unstructured adaptive meshes.

## 6 Conclusions

In this paper we have discussed a genuinely explicit residual discretization of the shallow water equations based on an improved and adapted formulation of the nonlinear stabilized explicit limited Lax-Friedrichs scheme of [57]. The scheme has been shown to enjoy all the most interesting properties relevant for shallow water applications, namely the C-property and its generalization to moving equilibria, positivity preservation, and a robust handling of

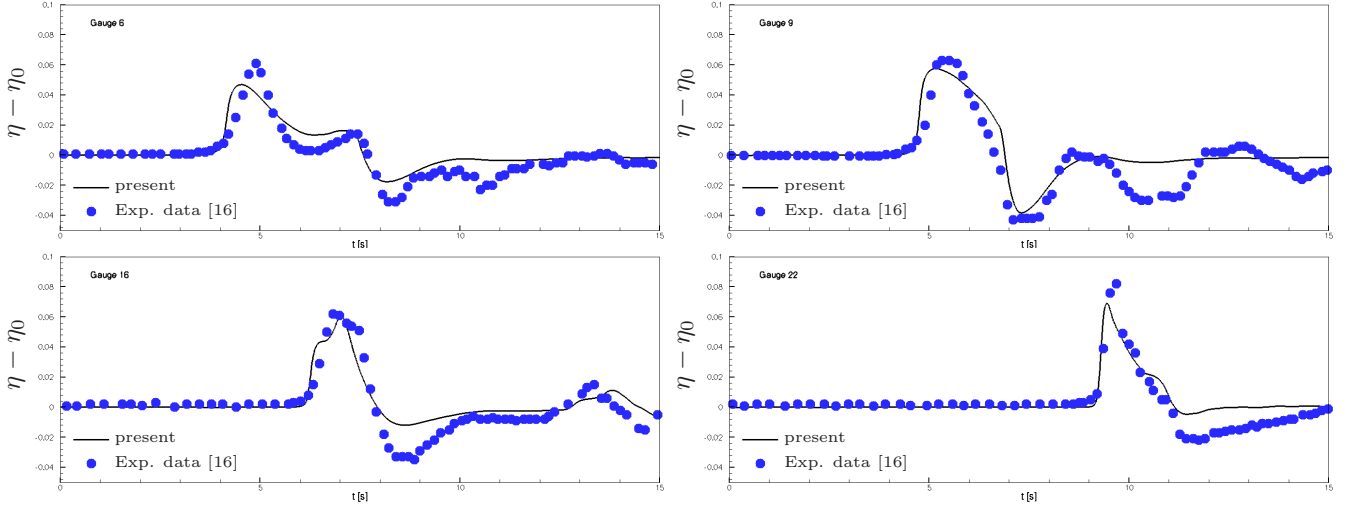


Figure 22: Run up on a conical island [16] : gauge data (left)

the wetting/drying front. As shown in the numerical results section, this work represents a considerable improvement over previous work by the author and his collaborators [56, 58, 60, 59], providing the same properties and similar accuracy with a much reduced computational cost. Theoretical results and thorough benchmarking have confirmed this fact. Following the initial work of [57], this paper finally brings residual distribution schemes to a cost similar to that of explicit second order Godunov type schemes (including DG), retaining all the advantages of the continuous residual based formulation, including the potential to capture both steady and moving equilibria on unstructured grids.

The scheme proposed has been already combined with uncertainty quantification techniques to study the sensitivity of long wave runup simulations to variations in physical parameters [61], and for robust code-to-code validation [25]. In this last reference, in particular, it has been shown that, in long wave runup simulations, the scheme proposed here provides accuracy levels very close to state of the art high order finite volume schemes.

Concerning the scheme, foreseen improvements are the design of higher (at least third) order formulations, and the use of ALE based moving mesh techniques for mesh adaptation, in particular to follow moving shorelines. Concerning the extension to other models, future work will involve the inclusion of Coriolis terms as already shown in [69], and eventually a formulation of the schemes on manifolds as in [67]. Current work also includes the investigation of residual based discretizations of non-hydrostatic models [63, 64].

## A Proof of proposition 4.2

In all following proofs, we consider the error in approximating an initial solution which is also a steady exact solution  $u^0 = u^0(x, y)$ . This state is assumed to be characterized by the existence of an invariant  $v(u, b)$  verifying  $v(u^0, b) = v_0 = \text{const}$ . In particular, for the finite element approximation of this exact solution we trivially have  $v_h = v_0$ . When interpolating



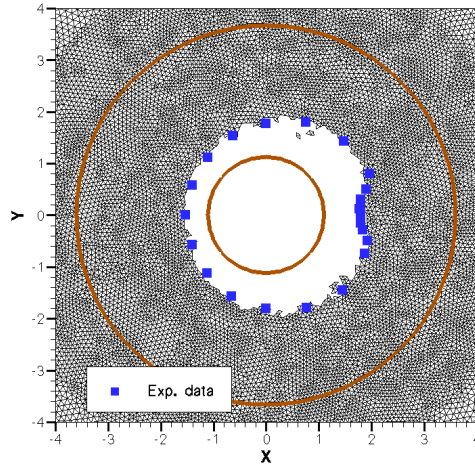


Figure 23: Run up on a conical island [16] : maximum run up plot

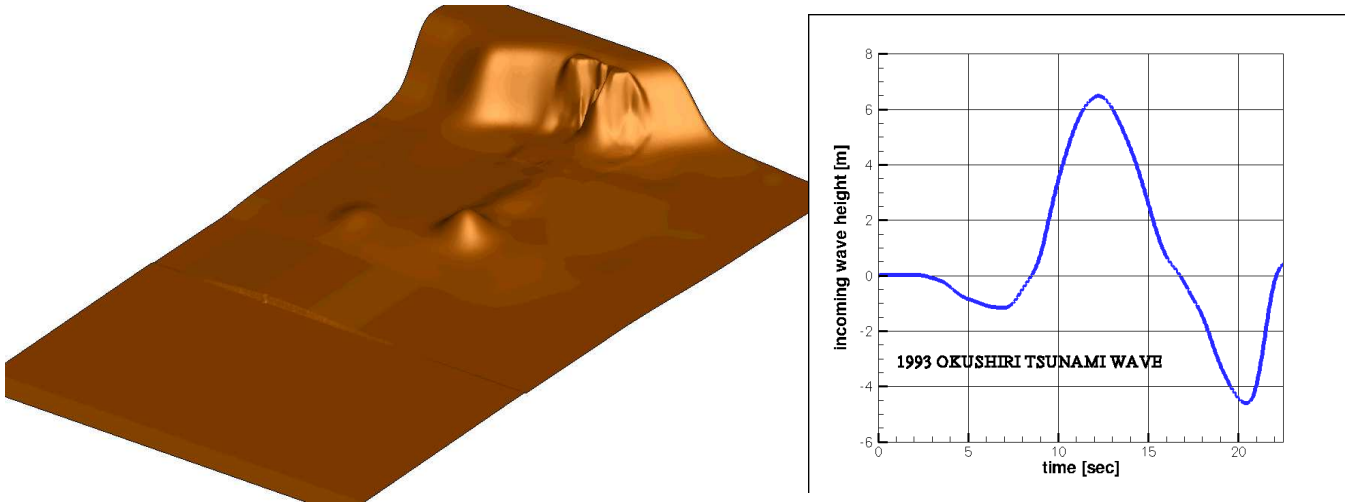


Figure 24: The Okushiri Tsunami experiment. Left : bathymetry. Right : inlet wave

the steady invariant with the analytical bathymetry  $b = b(x, y)$  we have moreover

$$u_h^0 = u(v_h, b) = u(v_0, b) = u^0 = u^0(x, y).$$

So, under the hypotheses made  $u_h^0 = u^0$  is the exact steady solution in conservative variables, with  $v_h = v_0$  being the exact steady state invariant. Also note that, as done in standard truncation error analysis, in the following we will formally replace  $u_h = u(v_h, b)$  by the approximation of the exact solution, and estimate the rest. In particular we will use everywhere the hypothesis

$$w_h^n = w_h^{n+1} = u_h^0 = u^0.$$

The analysis that follows only considers the two-dimensional case.

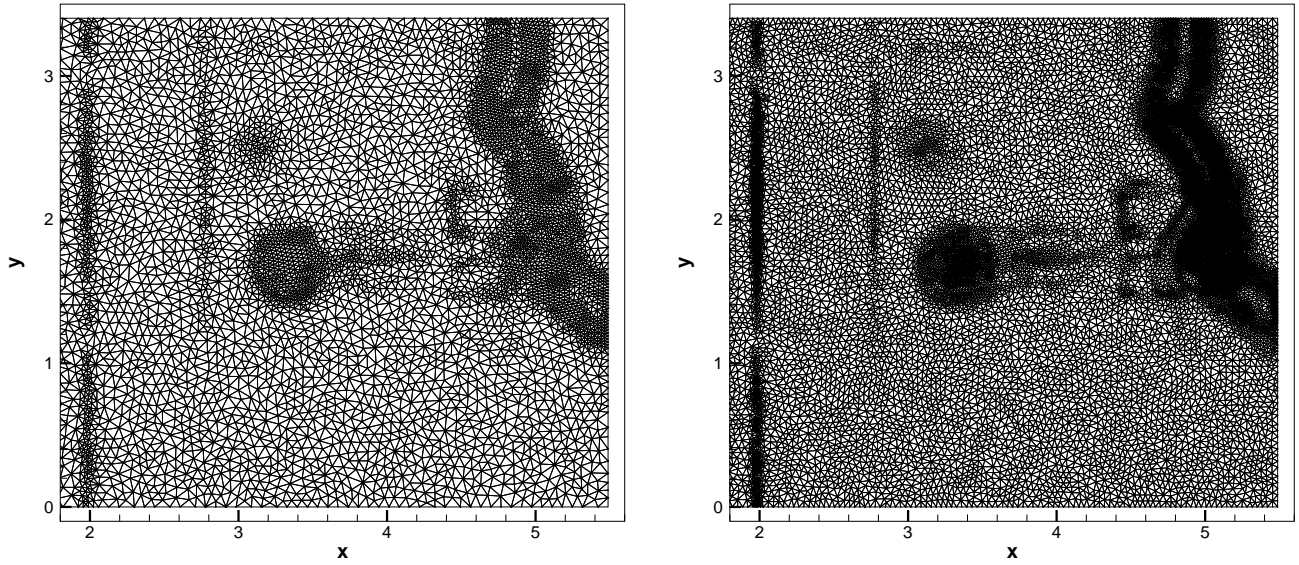


Figure 25: The Okushiri Tsunami experiment. Left : coarse adaptive mesh ( $h_{\min} \approx 0.025$  and  $h_{\max} \approx 0.1$ ) bathymetry. Right : fine adaptive mesh ( $h_{\min} \approx 0.01$  and  $h_{\max} \approx 0.05$ )

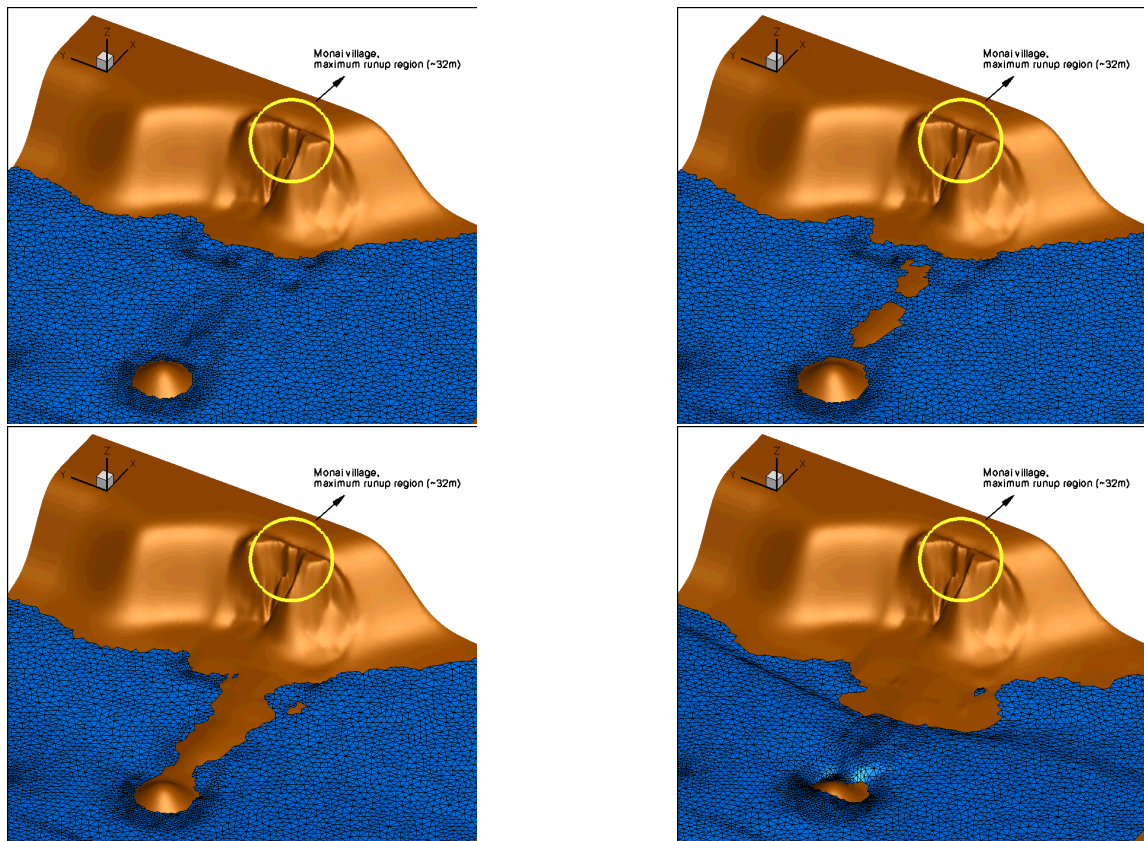


Figure 26: The Okushiri Tsunami experiment. From from left to right and from top to bottom : 3D visualization of the inundation and reflection process

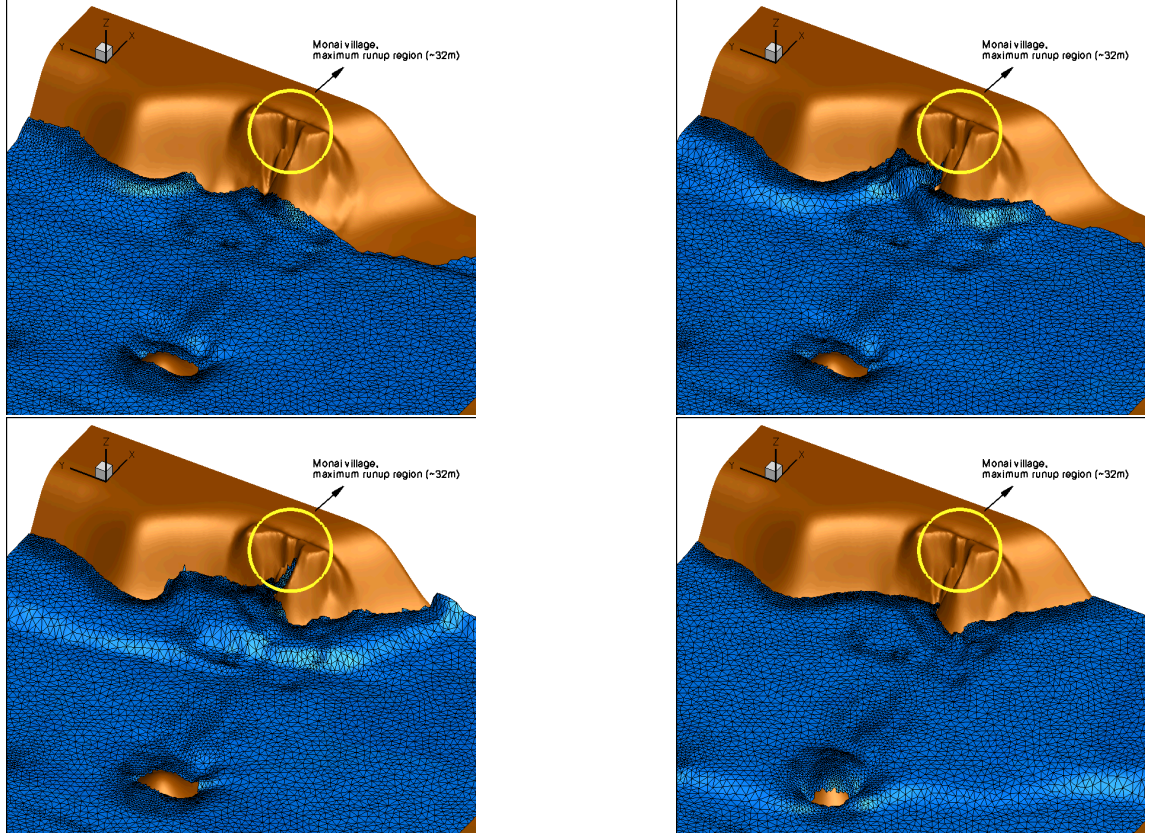


Figure 27: The Okushiri Tsunami experiment (continued). From from left to right and from top to bottom : 3D visualization of the inundation and reflection process

### Proof of Lemma 4.1

To prove the lemma we start by noting that for  $b \in H^{p+1}$  with  $p \geq \min(p_f, p_v + 1) \geq 1$  we can write for exact integration (cf. equation (39))

$$\oint_{\partial K} \mathcal{F}(u(v_h, b)) \cdot \vec{n} = \int_K \nabla \cdot \mathcal{F}(u(v_h, b)) = \int_K \left( \frac{\partial \mathcal{F}}{\partial v}(u(v_h, b)) \cdot \nabla v_h + \mathcal{S}_v(u(v_h, b), \nabla b) \right).$$

and so

$$\phi^K = \int_K \left( \frac{\partial \mathcal{F}}{\partial v}(u(v_h, b)) \cdot \nabla v_h + \mathcal{S}_v(u(v_h, b), \nabla b) + \mathcal{S}(u(v_h, b), \nabla b) \right).$$

Since  $v_0$  is an invariant, and it describes a steady equilibrium, then we deduce immediately that (cf. (39) and (20))

$$\nabla v_0 = 0, \quad \mathcal{S}_v(u(v_0, b), \nabla b) + \mathcal{S}(u(v_0, b), \nabla b) = 0,$$

and, as a consequence, we deduce that  $\phi^K(v_0, b) = 0$ .

For the second part of the proof, due to the assumed regularity of  $(v, b) \mapsto \mathcal{F}(u(v, b))$ , and since  $v_h = v_0$  which is constant, then we deduce that  $\mathcal{F}_h = \mathcal{F}(u(v_0, b))$  has the same regularity



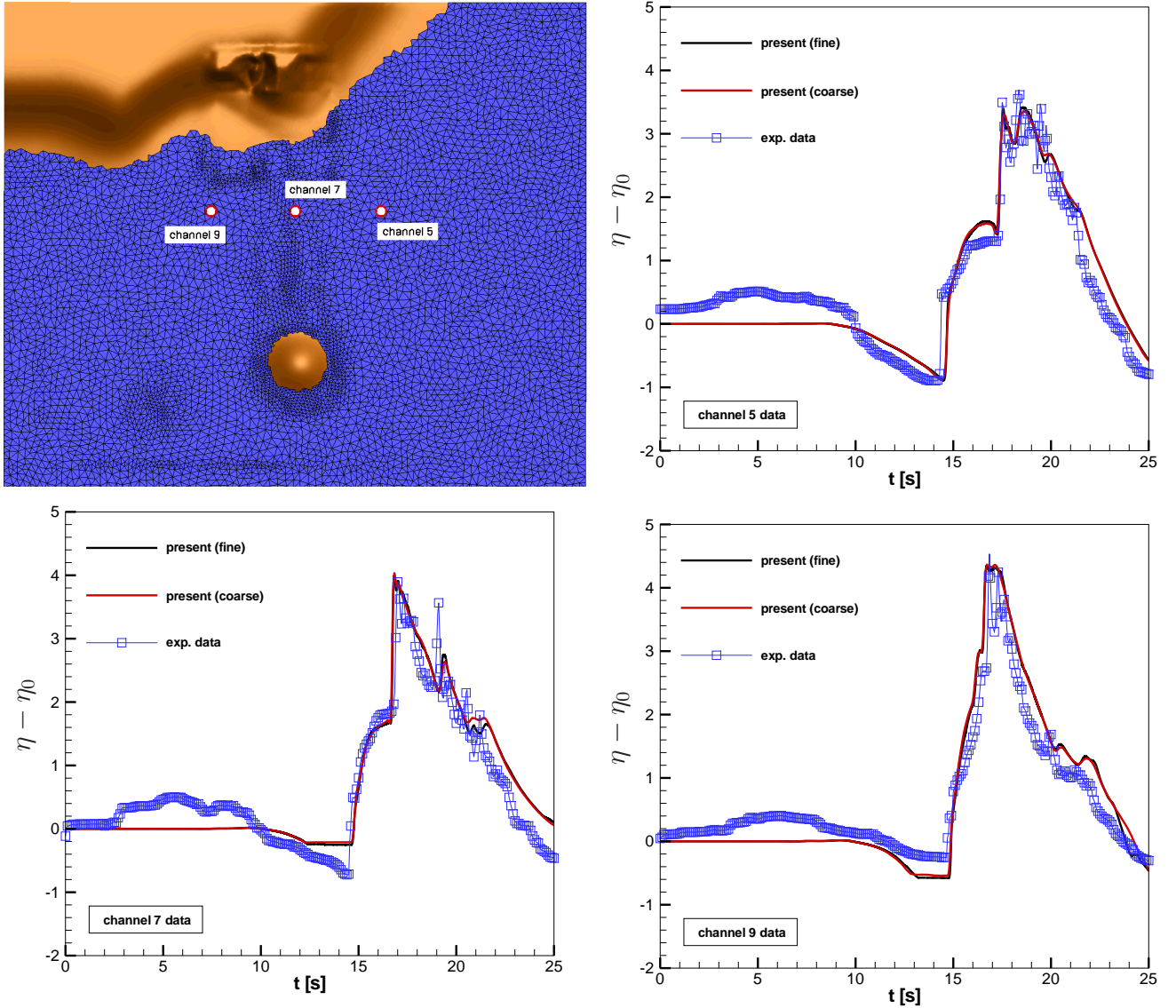


Figure 28: The Okushiri Tsunami experiment. Experimental gauge positions (top left) and comparisons with experiments in channels 5 (top right), 7 (bottom left) and 9 (bottom right)

as  $b$ , which means that  $\mathcal{F}(u(v_0, b)) \in H^{p+1}(\Omega_h)$ . Similarly, we argue that  $\mathcal{S}_h = \mathcal{S}(u(v_0, b), \nabla b)$  is in  $H^p$ . We now consider on each  $K$ , the polynomials  $\widehat{\mathcal{F}}_h$  of degree  $p_f$ , and the polynomials  $\widetilde{\mathcal{S}}_h$  of degree  $p_v$  such that, denoting by  $f$  the generic face of  $\partial K$  (we omit the additional superscript  $^K$ )

$$\sum_{q=1}^{f_q} \omega_q \mathcal{F}_h(\vec{x}_q) \cdot \vec{n}_f = \int_f \widehat{\mathcal{F}}_h \cdot \vec{n}_f, \quad \text{and} \quad \sum_{q=1}^{v_q} \bar{\omega}_q \mathcal{S}_h(\vec{x}_q) = \int_K \widetilde{\mathcal{S}}_h.$$

For conservation reasons, and without loss of generality, the polynomial approximation  $\widehat{\mathcal{F}}_h$  is assumed to be continuous across element edges. With this notation, we can write, subtracting

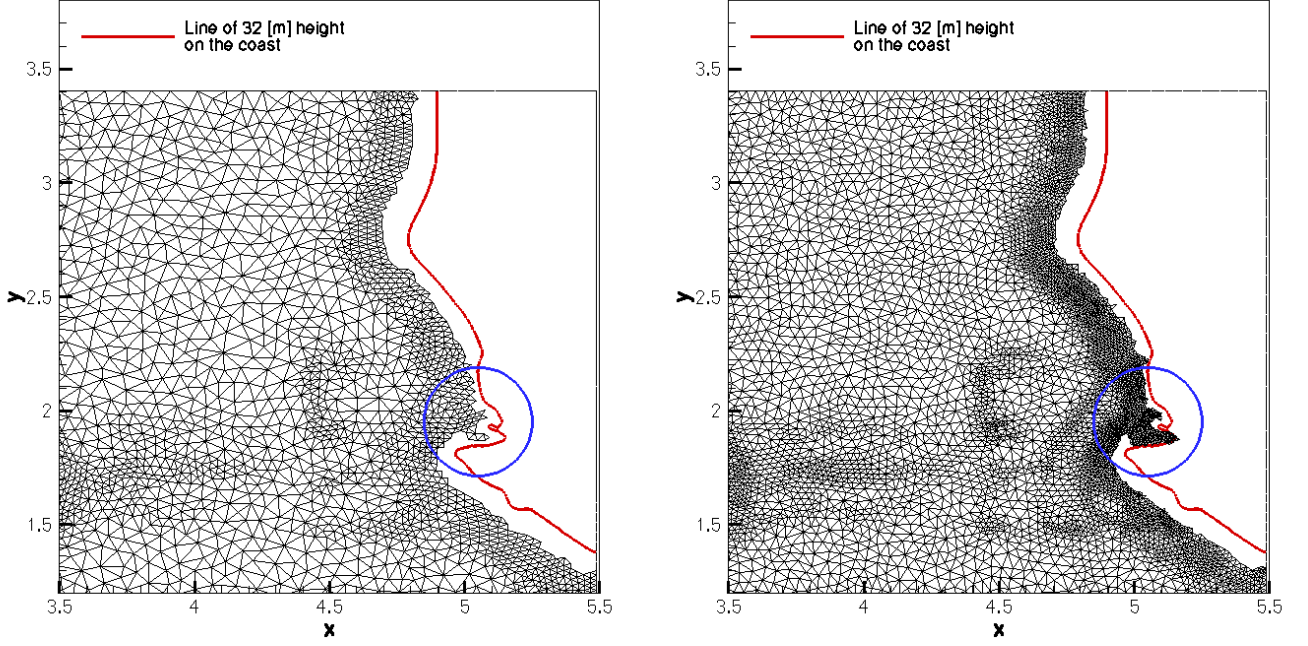


Figure 29: The Okushiri Tsunami experiment. Run up plots on the coarse (left) and fine (right) mesh. Circle : region of max runup (32 [m]) in observations.

the exact integral which is zero :

$$\begin{aligned}
\|\phi^K(v_0, b)\| &= \left\| \sum_{f \in \partial K} \int_f \widehat{\mathcal{F}}_h \cdot \vec{n}_f + \int_K \widetilde{\mathcal{S}}_h \right\| \\
&= \left\| \sum_{f \in \partial K} \int_f (\widehat{\mathcal{F}}_h - \mathcal{F}(u(v_0, b))) \cdot \vec{n}_f + \int_K (\widetilde{\mathcal{S}}_h - \mathcal{S}(u(v_0, b), \nabla b)) \right\| \\
&\leq \sum_{f \in \partial K} \int_f \|(\widehat{\mathcal{F}}_h - \mathcal{F}(u(v_0, b))) \cdot \vec{n}_f\| + \int_K \|\widetilde{\mathcal{S}}_h - \mathcal{S}(u(v_0, b), \nabla b)\|.
\end{aligned}$$

For the given regularity of  $b$ , we can write using standard approximation arguments [24, 31]

$$\begin{aligned}
\|\widehat{\mathcal{F}}_h - \mathcal{F}(u(v_0, b))\| &\leq C(v_0, b) h^{p_f+1} \Rightarrow \int_f \|(\widehat{\mathcal{F}}_h - \mathcal{F}(u(v_0, b))) \cdot \vec{n}_f\| = \mathcal{O}(h^{p_f+2}) \\
\|\widetilde{\mathcal{S}}_h - \mathcal{S}(u(v_0, b), \nabla b)\| &\leq C'(v_0, b) h^{p_v+1} \Rightarrow \int_K \|\widetilde{\mathcal{S}}_h - \mathcal{S}(u(v_0, b), \nabla b)\| = \mathcal{O}(h^{p_v+3})
\end{aligned}$$

This leads to the final estimate  $\|\phi^K(v_0, b)\| \leq C'' \max(h^{p_f+2}, h^{p_v+3})$ . Note that, even if  $b$  has lower regularity, this proof only uses the assumed regularity within the element  $K$ . As a consequence, if the edges of  $K$  are aligned with lines across which some derivatives of  $b$  are discontinuous, this will not affect the final estimate which will be the same.

## Proof of Proposition 4.2

In order to prove the proposition, we rewrite the truncation error (32) for a steady smooth equilibrium  $v = v_0$ . First of all, the predictor step (33) provides an estimate on the error  $\|w^* - u^0\|$ . Indeed, setting  $u_i^0 = u^0(x_i, y_i)$ , from

$$|C_i| \frac{w_i^* - u_i^0}{\Delta t^n} + \sum_{K \in K_i} \beta_i^K \phi^K(v_0, b) = 0,$$

and using Lemma 4.1 and (34) we immediately deduce that for exact integration  $w^* = u^0$ , and it is trivial to see that (32) is identically zero due to Lemma 4.1.

For approximate integration, a crude estimate based on Lemma 4.1 gives

$$\|w_h^* - u^0\| = \mathcal{O}(h^l), \quad l = \min(p_f + 1, p_v + 2), \quad (56)$$

since  $|C_i|/\Delta t^n = \mathcal{O}(h)$ . A sharper estimate can be obtained as follows. First, with the notation introduced in the previous paragraph, we consider the following local Galerkin projection :

$$\phi_i^G(v_0, b) = \int_K \varphi_i (\nabla \cdot \widehat{\mathcal{F}} + \widetilde{\mathcal{S}}_h), \quad (57)$$

with  $\varphi_i$  the  $P^1$  basis functions, and where using the properties of the shape functions we have that by construction

$$\sum_{j \in K} \phi_j^G(v_0, b) = \phi^K(v_0, b). \quad (58)$$

Proceeding as in the last paragraph, we can easily show that

$$\phi_i^G(v_0, b) = \mathcal{O}(h^{l+1}). \quad (59)$$

Consider now the quantity

$$e_0 = \sum_{K \in \Omega_h} \sum_{j \in K} \psi_j \beta_j^K \phi^K(v_0, b), \quad (60)$$

with  $\psi$  a smooth compactly supported function as in (32). Clearly, for exact integration we have  $e_0=0$ . For approximate integration, proceeding as in e.g. [58] we recast it as<sup>1</sup>

$$e_0 = - \int_{\partial\Omega_h} (\widehat{\mathcal{F}} - \mathcal{F}) \cdot \nabla \psi_h + \int_{\Omega_h} \psi_h (\widetilde{\mathcal{S}}_h - \mathcal{S}) + \sum_{K \in \Omega_h} \sum_{i, j \in K} \frac{\psi_i - \psi_j}{2} (\beta_i^K \phi^K - \phi_i^G),$$

having also used the fact that for the exact solution we have  $\nabla \cdot \mathcal{F} + \mathcal{S} = 0$  (cf. [58] for more). We can now easily estimate  $e_0$  as

$$\begin{aligned} \|e_0\| &\leq C_{\Omega_h} \left\{ \|\widehat{\mathcal{F}}_h - \mathcal{F}\| \|\nabla \psi\| + \|\psi\| \|\widetilde{\mathcal{S}}_h - \mathcal{S}\| \right. \\ &\quad \left. + \frac{|\Omega_h|}{h^2} \|\nabla \psi\| h \sup_{K \in \Omega_h} \sup_{j \in K} (\|\beta_j^K\| \|\phi^K\| + \|\phi_i^G\|) \right\} \leq C_0 h^{\min(p_f+1, p_v+1)}, \end{aligned} \quad (61)$$

<sup>1</sup>dependence on  $(v_0, b)$ ,  $u(v_0, b)$ , and  $\nabla b$  omitted for simplicity

having used (59), the estimates on  $\phi^K$ , the hypotheses on the quadrature and on the polynomials  $\widehat{\mathcal{F}}_h$  and  $\widehat{\mathcal{S}}_h$ , and the fact that the number of elements in the mesh is of order  $\mathcal{O}(|\Omega_h|/h^2)$ .

The error of the predictor step is now immediately estimated using the identity [10, 58]

$$\sum_{i \in \Omega_h} |C_i| \psi_i(w_i^* - u_i^0) = -\Delta t^n e_0, \quad (62)$$

which, by virtue of the estimate on  $e_0$ , and of hypotheses (34), gives

$$\left\| \sum_{i \in \Omega_h} |C_i| \psi_i(w_i^* - u_i^0) \right\| \leq C_0 h^{\min(p_f+2, p_v+2)}. \quad (63)$$

To use this result we now consider the integral

$$\int_{\Omega_h} \psi_h(w_h^* - u_h^0) = \sum_{K \in \Omega_h} \int_K \psi_h(w_h^* - u_h^0),$$

and apply the standard second order formula using edge midpoints to evaluate it exactly. This gives on each  $K$  :

$$\int_K \psi_h(w_h^* - u_h^0) = \frac{|K|}{3} \left( \sum_{j \in K} \left( \frac{\psi_j(w_j^* - u_j^0)}{2} + \sum_{l \neq j} \frac{\psi_l(w_j^* - u_j^0)}{4} \right) \right).$$

Considering that for each  $j$  there are two indices  $l \neq j$ , and that with the hypotheses made on  $\psi$  we have  $\psi_l = \psi_j + \mathcal{O}(h)$  we readily obtain that

$$\int_K \psi_h(w_h^* - u_h^0) = \frac{|K|}{3} \sum_{j \in K} (\psi_j + C_j(\psi) h) (w_j^* - u_j^0),$$

for some bounded constants  $C_j$  depending of the derivatives of  $\psi$ . This leads to the estimate

$$\begin{aligned} \left\| \int_{\Omega_h} \psi_h(w_h^* - u_h^0) \right\| &\leq \left\| \sum_{i \in \Omega_h} |C_i| \psi_i(w_i^* - u_i^0) \right\| + \sum_{K \in \Omega_h} C_K |K| h \sup_{j \in K} |w_j^* - u_j^0| \\ &\leq \left\| \sum_{i \in \Omega_h} |C_i| \psi_i(w_i^* - u_i^0) \right\| + C_{\Omega_h} h \sup_{i \in \Omega_h} |w_i^* - u_i^0|, \end{aligned}$$

having used the fact that the number of elements is of  $\mathcal{O}(h^{-2})$  and that  $|K| = \mathcal{O}(h^2)$ . Lastly, using (56) and (63), we can write

$$\left\| \int_{\Omega_h} \psi_h(w_h^* - u_h^0) \right\| \leq C_1 h^{\min(p_f+2, p_v+2)}. \quad (64)$$

The objective is now to bound the truncation error, using the results obtained for the predictor step. Recalling that for the truncation error analysis  $w_h^{n+1} = w_h^n = u^0$ , and explicitly

using the predictor step to replace the values  $w_i^*$ , we start by writing

$$\begin{aligned}
& |C_i| \frac{w_i^{n+1} - w_i^*}{\Delta t^n} + \sum_{K \in K_i} \Phi_i^K(w_h^n, w_h^*) = \\
& |C_i| \frac{u_i^0 - w_i^*}{\Delta t^n} + \sum_{K \in K_i} \beta_i^K \left( \int_K \frac{w_h^* - u_h^0}{\Delta t^n} + \frac{1}{2} \phi^K(w_h^*) + \frac{1}{2} \phi^K(u_h^0) \right) = \quad (65) \\
& \sum_{K \in K_i} \beta_i^K \int_K \frac{w_h^* - u_h^0}{\Delta t^n} + \sum_{K \in K_i} \beta_i^K \phi^K(u_h^0) + \frac{1}{2} \sum_{K \in K_i} \beta_i^K (\phi^K(w_h^*) - \phi^K(u_h^0))
\end{aligned}$$

Note that here, consistently with the assumptions of the proposition,  $\phi^K(w_h^*) = \phi^K(u(v_h^*, b))$  is obtained by evaluating exactly

$$\phi^K(u(v_h^*, b)) = \int_K \left( \nabla \cdot \widehat{\mathcal{F}}_h^* + \widetilde{\mathcal{S}}_h^* \right),$$

with  $\widehat{\mathcal{F}}_h^*$  and  $\widetilde{\mathcal{S}}_h^*$  the polynomials of degree respectively  $p_f$ , and  $p_v$ , obtained by interpolating the invariant  $v(w^*, b)$ , exactly as  $\widehat{\mathcal{F}}_h$  and  $\widetilde{\mathcal{S}}_h$  are obtained by using  $v_0$ .

Injecting (65) in the definition of the error (32), and following [10, 58, 57], we can write the consistency error as

$$\epsilon = \text{I} + \text{II} + \text{III} + \text{IV},$$

with

$$\begin{aligned}
\text{I} &= \sum_{n=0}^N \Delta t^n \int_{\Omega_h} \psi_h \frac{w_h^* - u_h^0}{\Delta t^n} - \sum_{n=0}^N \Delta t^n \int_{\partial\Omega_h} \widehat{\mathcal{F}}_h \cdot \nabla \psi_h + \sum_{n=0}^N \Delta t^n \int_{\Omega_h} \psi_h \widetilde{\mathcal{S}}_h, \\
\text{II} &= - \sum_{n=0}^N \frac{\Delta t^n}{2} \int_{\partial\Omega_h} (\widehat{\mathcal{F}}_h^* - \widehat{\mathcal{F}}_h) \cdot \nabla \psi_h + \sum_{n=0}^N \frac{\Delta t^n}{2} \int_{\Omega_h} (\widetilde{\mathcal{S}}_h^* - \widetilde{\mathcal{S}}_h) \psi_h, \\
\text{III} &= \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) (w_h^* - u_h^0), \\
&+ \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \Delta t^n \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) \nabla \cdot \widehat{\mathcal{F}}_h \\
&+ \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \Delta t^n \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) \widetilde{\mathcal{S}}_h \\
\text{IV} &= \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \frac{\Delta t^n}{2} \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) \nabla \cdot (\widehat{\mathcal{F}}_h^* - \widehat{\mathcal{F}}_h) \\
&+ \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \frac{\Delta t^n}{2} \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) (\widetilde{\mathcal{S}}_h^* - \widetilde{\mathcal{S}}_h).
\end{aligned}$$

We can now estimate each term using the results and hypotheses available.



The estimate of term I is almost identical to that of  $e_0$ . In particular, using of the fact that the exact solution verifies  $\nabla \cdot \mathcal{F} + \mathcal{S} = 0$ , and that the number of time steps is of order  $\mathcal{O}(\Delta t^{-1}) = \mathcal{O}(h^{-1})$ , we can write

$$\|\text{I}\| \leq C_0 \left( h^{-1} \left\| \int_{\Omega_h} \psi_h (w_h^* - u_h^0) \right\| + \left\| \int_{\partial\Omega_h} (\widehat{\mathcal{F}}_h - \mathcal{F}) \cdot \nabla \psi_h \right\| + \left\| \int_{\Omega_h} \psi_h (\widetilde{\mathcal{S}}_h - \mathcal{S}) \right\| \right).$$

Making now use of hypotheses (34), of the assumptions made on the polynomials  $\widehat{\mathcal{F}}_h$  and  $\widetilde{\mathcal{S}}_h$ , of estimate (64), and proceeding as for  $e_0$ , we end with

$$\|\text{I}\| \leq C_1 h^{\min(p_f+1, p_v+1)}. \quad (66)$$

Similar arguments can be used to estimate the term III. In particular, we can use the fact that for the analytical solution  $\nabla \cdot \mathcal{F} + \mathcal{S} = 0$  to recast this term as

$$\begin{aligned} \text{III} &= \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) (w_h^* - u_h^0) \\ &+ \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \Delta t^n \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) \nabla \cdot (\widehat{\mathcal{F}}_h - \mathcal{F}) \\ &+ \sum_{n=0}^N \sum_{K \in \Omega_h} \sum_{i,j \in K} \Delta t^n \frac{\psi_i - \psi_j}{2} \int_K (\beta_i^K - \varphi_i) (\widetilde{\mathcal{S}}_h - \mathcal{S}). \end{aligned}$$

For a smooth solution, provided that (34) holds we can initially write

$$\begin{aligned} \|\text{III}\| &\leq C \left\{ \frac{T}{\Delta t} \frac{|\Omega_h| h}{h^2} \frac{1}{2} \|\nabla \psi\| h^2 \left( 1 + \sup_{\substack{K \in \Omega_h \\ i \in K}} \|\beta_i^K\|_K \right) \|w_h^* - u_h^0\|_K \right. \\ &+ \frac{|\Omega_h| h}{h^2} \frac{1}{2} \|\nabla \psi\| h^2 \left( 1 + \sup_{\substack{K \in \Omega_h \\ i \in K}} \|\beta_i^K\|_K \right) \|\nabla \cdot (\widehat{\mathcal{F}}_h - \mathcal{F})\|_K \\ &\left. + \frac{|\Omega_h| h}{h^2} \frac{1}{2} \|\nabla \psi\| h^2 \left( 1 + \sup_{\substack{K \in \Omega_h \\ i \in K}} \|\beta_i^K\|_K \right) \|\widetilde{\mathcal{S}}_h - \mathcal{S}\|_K \right\}. \end{aligned}$$

Last expression can be recast as

$$\|\text{III}\| \leq \overline{C}_{\text{III}} \|\nabla \psi\| \left( 1 + \sup_{\substack{K \in \Omega_h \\ i \in K}} \|\beta_i^K\|_K \right) \left( \|w_h^* - u_h^0\|_K + h \|\nabla \cdot (\widehat{\mathcal{F}}_h - \mathcal{F})\|_K + h \|\widetilde{\mathcal{S}}_h - \mathcal{S}\|_K \right)$$

Using the regularity of  $\psi$ , and the hypothesis on the boundedness of the distribution coefficients, we can now bound this term using the estimate (56), and standard approximation arguments [24, 31] to estimate the remaining terms. In particular, recalling that  $\widehat{\mathcal{F}}_h$  is a polynomial of degree  $p_f$ , and that  $\widetilde{\mathcal{S}}_h$  is a polynomial of degree  $p_v$ , we have [24, 31]

$$\|\nabla \cdot (\widehat{\mathcal{F}}_h - \mathcal{F})\|_K = \mathcal{O}(h^{p_f}) \quad \text{and} \quad \|\widetilde{\mathcal{S}}_h - \mathcal{S}\|_K = \mathcal{O}(h^{p_v+1}),$$

leading to

$$\|\text{III}\| \leq C_{\text{III}} h^{\min(p_f+1, p_v+2)}. \quad (67)$$

To end the proof we need to estimate II and IV. To do this, we use the hypothesis on the Lipschitz continuity of the flux and source (37). In particular, using the edge continuity of the flux approximation, we start recasting II as

$$\text{II} = - \sum_{n=0}^N \frac{\Delta t^n}{2} \int_{\partial\Omega_h} (\widehat{\mathcal{F}}_h^* - \widehat{\mathcal{F}}_h) \cdot \nabla \psi_h + \sum_{n=0}^N \frac{\Delta t^n}{2} \int_{\Omega_h} (\widetilde{\mathcal{S}}_h^* - \widetilde{\mathcal{S}}_h) \psi_h.$$

We next use (37) to obtain

$$\|\widehat{\mathcal{F}}_h^* - \widehat{\mathcal{F}}_h\| \leq \mathcal{K}_{\mathcal{F}} \|w_h^* - u^0\| \leq C h^{\min(p_f+1, p_v+2)}$$

$$\|\widetilde{\mathcal{S}}_h^* - \widetilde{\mathcal{S}}_h\| \leq \mathcal{K}_{\mathcal{S}} \|w_h^* - u^0\| \leq C' h^{\min(p_f+1, p_v+2)}$$

and, thus we can readily bound this term as

$$\|\text{II}\| \leq C_{\text{II}} h^{\min(p_f+1, p_v+2)}. \quad (68)$$

The last term is estimated in a similar way. In particular, we first express the polynomial  $\widehat{\mathcal{F}}_h$  using a local high order finite element basis

$$\widehat{\mathcal{F}}_h = \sum_{\sigma} \mathcal{F}_{\sigma} \widehat{\varphi}_{\sigma},$$

with the  $\widehat{\varphi}_{\sigma}$  the kernel of a higher degree (at least  $p_f$ ) Lagrange approximation. Next we observe that using (56) we can write

$$\|\nabla \cdot (\widehat{\mathcal{F}}_h^* - \widehat{\mathcal{F}}_h)\|_K = \left\| \sum_{\sigma} (\mathcal{F}(w_{\sigma}^*) - \mathcal{F}(u_{\sigma}^0)) \cdot \nabla \widehat{\varphi}_{\sigma} \right\|_K \leq \frac{\widehat{C}_K \mathcal{K}_{\mathcal{F}}}{h} \sum_{\sigma} \|w_{\sigma}^* - u_{\sigma}^0\| \leq C h^{\min(p_f, p_v+1)}.$$

Similarly one obtains the estimate

$$\|\widetilde{\mathcal{S}}_h^* - \widetilde{\mathcal{S}}_h\| \leq C h^{\min(p_f+1, p_v+2)}.$$

Finally, term IV is estimated as

$$\|\text{IV}\| \leq C \frac{|\Omega_h|}{h^2} \|\nabla \psi\| h h^2 \left(1 + \sup_{\substack{K \in \Omega_h \\ i \in K}} \|\beta_i^K\|_K\right) h^{\min(p_f, p_v+1)} \leq C_{\text{IV}} h^{\min(p_f+1, p_v+2)}, \quad (69)$$

which together with (66), (68), and (67) achieves the proof.

## B Proof of proposition 4.6

To prove proposition 4.7 we use the properties of the limiter and the definition of the LF distribution recalled in section §4.3, in particular (47) (see [57, 60] for more details), The explicit limited Lax-Friedrichs (LLF) scheme obtained by applying the limiter equation by equation leads to the following updates for the water height

1. Predictor step

$$|C_i|(H_i^* - H_i^n) = -\Delta t \sum_{K \in K_i} \gamma_i \sum_{j \in K, j \neq i} \frac{1}{3} (\alpha_{\text{LF}} - k_j^n) (H_i^n - H_j^n).$$

2. Corrector step

$$\begin{aligned} |C_i|(H_i^{n+1} - H_i^*) = & -\Delta t \sum_{K \in K_i} \gamma_i \left( \frac{|K|}{3} \frac{H_i^* - H_i^n}{\Delta t} + \sum_{j \in K, j \neq i} \frac{1}{6} (\alpha_{\text{LF}} - k_j^n) (H_i^n - H_j^n) \right. \\ & \left. + \sum_{j \in K, j \neq i} \frac{1}{6} (\alpha_{\text{LF}} - k_j^*) (H_i^* - H_j^*) \right), \end{aligned}$$

where

$$k_j = \frac{\vec{u}_j \cdot \vec{n}_j}{2},$$

with  $\vec{n}_j$  the inward normal to the edge facing node  $j$ , scaled by its length  $l_j$ . Note that, by definition of  $\alpha_{\text{LF}}$ , we have

$$\alpha_{\text{LF}} - k_j^n \geq 0 \quad \text{and} \quad \alpha_{\text{LF}} - k_j^* \geq 0.$$

If  $H_i^n \geq 0$ ,  $\forall i$ , the positivity of  $H_i^*$  is easily shown to lead to the condition

$$\Delta t \sum_{K \in K_i} \frac{\gamma_i}{3} (2\alpha_{\text{LF}} + k_i^n) \leq |C_i|, \quad (70)$$

which, using  $\gamma_i^n \in [0, 1]$  and  $\alpha_{\text{LF}} \geq k_j$  (cf. (47) in section §4.5, and see [60]), can be replaced by the constraint in the first slot of the  $\min(\cdot, \cdot)$  operator in (48).

We now set on every element

$$\omega_i^n = \frac{2\alpha_{\text{LF}} + k_i^n}{3}, \quad \omega_i^* = \frac{2\alpha_{\text{LF}} + k_i^*}{3},$$

and also

$$\Omega_i^n = \sum_{K \in K_i} \omega_i^n, \quad \Omega_i^* = \sum_{K \in K_i} \omega_i^*.$$

The hypotheses made allow to show easily that (see e.g. [60] for details)

$$0 \leq \Delta t \Omega_i^n \leq |C_i|, \quad 0 \leq \Delta t \omega_i^n \leq \frac{|K|}{3}. \quad (71)$$

and

$$0 \leq \Delta t \Omega_i^* \leq |C_i|, \quad 0 \leq \Delta t \omega_i^* \leq \frac{|K|}{3}. \quad (72)$$

We now analyze the second iteration which can be recast as

$$\begin{aligned}
|C_i|H_i^{n+1} &= \sum_{K \in K_i} (1 - \gamma_i) \frac{|K|}{3} H_i^* - \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \frac{2\alpha_{\text{LF}} + k_i^*}{3} H_i^* + \sum_{K \in K_i} \gamma_i \left( \frac{|K|}{3} - \frac{\Delta t}{2} \frac{2\alpha_{\text{LF}} + k_i^n}{3} \right) H_i^n \\
&\quad + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^*}{3} H_j^* + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^n}{3} H_j^n \\
&= \sum_{K \in K_i} (1 - \gamma_i) \frac{|K|}{3} H_i^* - \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \omega_i^* H_i^* + \sum_{K \in K_i} \gamma_i \left( \frac{|K|}{3} - \frac{\Delta t}{2} \omega_i^n \right) H_i^n \\
&\quad + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^*}{3} H_j^* + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^n}{3} H_j^n.
\end{aligned}$$

We now add and remove  $H_i^n$  in the second term on the right hand side, and use the first iteration to replace the value of  $H_i^* - H_i^n$ . This leads to

$$\begin{aligned}
|C_i|H_i^{n+1} &= \sum_{K \in K_i} (1 - \gamma_i) \frac{|K|}{3} H_i^* + \sum_{K \in K_i} \gamma_i \left( \frac{|K|}{3} - \Delta t \frac{\omega_i^n + \omega_i^*}{2} \right) H_i^n \\
&\quad + \frac{\Delta t^2}{2} \frac{\Omega_i^*}{|C_i|} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^n}{3} (H_i^n - H_j^n) \\
&\quad + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^*}{3} H_j^* + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^n}{3} H_j^n \\
&= \sum_{K \in K_i} (1 - \gamma_i) \frac{|K|}{3} H_i^* + \sum_{K \in K_i} \gamma_i \left( \frac{|K|}{3} - \Delta t \frac{\omega_i^n + \omega_i^*}{2} + \frac{\Delta t^2}{2} \frac{\Omega_i^*}{|C_i|} \omega_i^n \right) H_i^n \\
&\quad + \frac{\Delta t}{2} \sum_{K \in K_i} \gamma_i \sum_{j \neq i} \frac{\alpha_{\text{LF}} - k_j^*}{3} H_j^* + \frac{\Delta t}{2} \sum_{K \in K_i} \sum_{j \neq i} \gamma_i \frac{\alpha_{\text{LF}} - k_j^n}{3} \left( 1 - \frac{\Delta t}{|C_i|} \frac{\Omega_i^*}{3} \right) H_j^n.
\end{aligned}$$

Using (71), (72), and the hypotheses on  $\alpha_{\text{LF}}$ , the last two lines of the expression obtained are a positive coefficient combination of the nodal values of  $H^n$  and  $H^*$ . Being  $H^n$  positive by hypothesis, and being  $H^*$  positive under the hypotheses made, this achieves the proof.

## References

- [1] Benchmark problem #2, Tsunami runup onto a complex three-dimensional beach. The third international workshop on long-wave runup models, <http://isec.nacse.org/workshop/2004.cornell/bmark2.html>.
- [2] Tsunami runup onto a complex three-dimensional beach; Monai valley. Benchmarks of the NOAA Center for tsunami research, <http://nctr.pmel.noaa.gov/benchmark/Laboratory/LaboratoryMonaiValley/>.
- [3] R. Abgrall. Essentially non oscillatory residual distribution schemes for hyperbolic problems. *J. Comput. Phys*, 214(2):773–808, 2006.

- [4] R. Abgrall. Residual distribution schemes: Current status and future trends. *Computers and Fluids*, 35(7):641 – 669, 2006.
- [5] R. Abgrall. A review of residual distribution schemes for hyperbolic and parabolic problems: the July 2010 state of the srt. *Comm.Comput.Phys.*, 11(4):1043–1080, 2012.
- [6] R. Abgrall and T.J. Barth. Residual distribution schemes for conservation laws via adaptive quadrature. *SIAM J. Sci. Comput.*, 24(3):732–769, 2002.
- [7] R. Abgrall, K. Mer, and B. Nkonga. A Lax–Wendroff type theorem for residual schemes. In M. Hafez and J.J. Chattot, editors, *Innovative methods for numerical solutions of partial differential equations*, pages 243–266. World Scientific, 2002.
- [8] R. Abgrall and M. Mezine. Construction of second-order accurate monotone and stable residual distribution schemes for unsteady flow problems. *J. Comput. Phys.*, 188:16–55, 2003.
- [9] R. Abgrall and M. Mezine. Construction of second-order accurate monotone and stable residual distribution schemes for steady flow problems. *J. Comput. Phys.*, 195:474–507, 2004.
- [10] R. Abgrall and P.L. Roe. High-order fluctuation schemes on triangular meshes. *J. Sci. Comput.*, 19(3):3–36, 2003.
- [11] R. Abgrall and J. Treflik. An example of high order residual distribution scheme using non-lagrange elements. *Journal of Scientific Computing*, 45:3–25, 2010.
- [12] E. Audusse and M.-O. Bristeau. A well-balanced positivity preserving second-order scheme for shallow water flows on unstructured meshes. *Journal of Computational Physics*, 206(1):311 – 333, 2005.
- [13] T.J. Barth. Numerical methods for conservation laws on structured and unstructured meshes. *VKI LS 2003-05, 33<sup>rd</sup> Computational Fluid dynamics Course, von Karman Institute for Fluid Dynamics*, 2003.
- [14] A. Bermudez and M.E. Vazquez. Upwind methods for hyperbolic conservation laws with source terms. *Computers & Fluids*, 23(8):1049 – 1071, 1994.
- [15] P. Bonneton, E. Barthelemy, F. Chazel, R. Cienfuegos, D. Lannes, F. Marche, and M. Tissier. Recent advances in Serre-Green Naghdi modelling for wave transformation, breaking and runup processes. *European Journal of Mechanics - B/Fluids*, 30(6):589 – 597, 2011.
- [16] M.J. Briggs, C.E. Synolakis, G.S. Harkins, and D.R. Green. Laboratory experiments of tsunami runup on a circular island. *Pure and Applied Geophysics*, 144:569–593, 1995.
- [17] P. Brufau and P. Garcia-Navarro. Unsteady free surface flow simulation over complex topography with a multidimensional upwind technique. *Journal of Computational Physics*, 186(2):503 – 526, 2003.
- [18] P. Brufau, P. Garcia-Navarro, and M.E. Vazquez-Cendon. Zero mass error using unsteady wetting-drying conditions in shallow flows over dry irregular topography. *Int. J. Numer. Meth. Fluids*, 45:1047–1082, 2004.
- [19] P. Brufau, M.E. Vazquez-Cendon, and P. Garcia-Navarro. A numerical model for the flooding and drying of irregular domains. *Int. J. Numer. Meth. Fluids*, 39:247–275, 2002.

- [20] V. Caleffi, A. Valliani, and A. Zanni. Finite volume method for simulating extreme flood events in natural channels. *J.of Hydraulic Research*, 41(2):167–177, 2003.
- [21] M.J. Castro, A.M. Ferreiro Ferreiro, J.A. Garcia-Rodriguez, J.M. Gonzalez-Vida, J. Macias, C. Pares, and M.E. Vazquez-Cendon. The numerical treatment of wet/dry fronts in shallow flows: application to one-layer and two-layer systems. *Mathematical and Computer Modelling*, 42(3-4):419 – 439, 2005.
- [22] M.J. Castro, J. Gonzalez-Vida, and C. Pares. Numerical treatment of wet/dry fronts in shallow flows with a Roe scheme. *Mathematical Models and Methods in Applied Sciences*, 16(6):897–931, 2006.
- [23] L. Cea and M.E. Vazquez-Cendon. Unstructured finite volume discretization of bed friction and convective flux in solute transport models linked to the shallow water equations. *J.Comput.Phys.*, 231(8):3317–3339, 2011.
- [24] P.G. Ciarlet and P.A. Raviart. General Lagrange and Hermite interpolation in  $\mathbb{R}^n$  with applications to finite element methods. *Arch.Ration.Mech.Anal.*, 46:177–199, 1972.
- [25] P.M. Congedo, A.I. Delis, and M. Ricchiuto. Robust code-to-code comparison for long wave run up. SIAM Conf. on Mathematical and Computational Issues in the Geosciences, Padova (Italy), June 2013.
- [26] Á. Csík, M. Ricchiuto, and H. Deconinck. A conservative formulation of the multidimensional upwind residual distribution schemes for general nonlinear conservation laws. *J. Comput. Phys*, 179(2):286–312, 2002.
- [27] H. Deconinck and M. Ricchiuto. Residual distribution schemes: foundation and analysis. In E. Stein, R. de Borst, and T.J.R. Hughes, editors, *Encyclopedia of Computational Mechanics*. John Wiley & Sons, Ltd., 2007. DOI: 10.1002/0470091355.ecm054.
- [28] A.I. Delis, M.Kazolea, and N.A.Kampanis. A robust high-resolution finite volume scheme for the simulation of long waves over complex domains. *Int. J. for Numerical Methods in Fluids*, 56:419–452, 2008.
- [29] A.I. Delis and N.Katsaounis. Relaxation schemes for the shallow water equations. *Int. J. for Numerical Methods in Fluids*, 41:695–719, 2003.
- [30] D.A. Dunavant. High degree efficient symmetrical gaussian quadrature rules for the triangle. *Int. J. Numer. Methods in Engrg.*, 21:1129–1148, 1985.
- [31] A. Ern and J.-C. Guermond. *Theory and practice of finite elements*, volume 159 of *Applied Mathematical Sciences*. Springer, 2004.
- [32] A. Ern, S. Piperno, and K. Djadel. A well-balanced Runge–Kutta discontinuous Galerkin method for the shallow-water equations with flooding and drying. *Int. J. for Numerical Methods in Fluids*, 58(1):1–25, 2008.
- [33] T. Gallouët, J.-M. Hérard, and N. Seguin. Some approximate Godunov schemes to compute shallow-water equations with topography. *Computers & Fluids*, 32(4):479–513, 2003.
- [34] P. Garcia-Navarro, J. Burguete, and R. Aliod. Numerical simulation of runoff over dry beds. *Monografias del Semin. Matem. Garcia de Galdeano*, 27:307–314, 2003.

- [35] J.-F. Gerbeau and B. Perthame. Derivation of viscous Saint-Venant system for laminar shallow water ; numerical validation. *Discrete and Continuous Dynamical Systems, Ser. B*, 1(1):89–102, 2001.
- [36] J.M. Greenberg and A.Y. Leroux. A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM J. Numer. Anal.*, 33:1–16, 1996.
- [37] A. Harten. On the symmetric form of systems of conservation laws with entropy. *J. Comput. Phys.*, 49:151–164, 1983.
- [38] G. Hauke. A symmetric formulation for computing transient shallow water flows. *Computer Methods in Applied Mechanics and Engineering*, 163(1-4):111–122, 1998.
- [39] G. Hauke, A. Landaberea, I. Garmendia, and J. Canales. On the thermodynamics, stability and hierarchy of entropy functions in fluid flow. *Computer Methods in Applied Mechanics and Engineering*, 195(33-36):4473–4489, 2006.
- [40] M. Hubbard and N. Dodd. A 2d numerical model of wave run-up and overtopping. *Coastal Engineering*, 47(1):1–26, 2002.
- [41] M. Hubbard and P. Garcia-Navarro. Flux difference splitting and the balancing of source terms and flux gradients. *J. Comp. Phys.*, 165(1):89–125, 2000.
- [42] M. Hubbard and M. Ricchiuto. Discontinuous upwind residual distribution: A route to unconditional positivity and high order accuracy. *Computers and Fluids*, 46(1):263 – 269, 2011.
- [43] T.J.R. Hughes, G. Scovazzi, and T. Tezduyar. Stabilized methods for compressible flows. *J. Sci. Comp.*, 43:343–368, 2010.
- [44] D. Lannes and P. Bonneton. Derivation of asymptotic two-dimensional time-dependent equations for surface water wave propagation. *Physics of Fluids*, 21, 2009. 016601 doi:10.1063/1.3053183.
- [45] R.J. LeVeque. Balancing source terms and flux gradients in high-resolution Godunov methods: the quasi-steady wave-propagation algorithm. *J. Comput. Phys.*, 146(1):346–365, 1998.
- [46] R.J. LeVeque. A well-balanced path-integral f-wave method for hyperbolic problems with source terms. *J. Sci. Comp.*, 48(1-3):209–226, 2011.
- [47] P. L.-F. Liu, H. Yeh, and C. Synolakis, editors. *Advanced Numerical Models for Simulating Tsunami Waves and Runup*, volume 10 of *Advances in Coastal and Ocean Engineering*. World Scientific, 2008.
- [48] F. Marche. Derivation of a new two-dimensional viscous shallow water model with varying topography, bottom friction and capillary effects. *European Journal of Mechanics - B/Fluids*, 26(1):49 – 63, 2007.
- [49] I.K. Nikolos and A.I. Delis. An unstructured node-centered finite volume scheme for shallow water flows with wet/dry fronts over complex topography. *Computer Methods in Applied Mechanics and Engineering*, 198(47-48):3723 – 3750, 2009.
- [50] S. Noelle, N. Pankratz, G. Puppo, and J.R. Natvig. Well-balanced finite volume schemes of arbitrary order of accuracy for shallow water flows. *J. Comput. Phys.*, 213(2):474–499, 2006.



- [51] S. Noelle, Y. Xing, and C.-W. Shu. High order well-balanced finite volume weno schemes for shallow water equation with moving water. *J.Comput.Phys*, 226:29–58, 2007.
- [52] S. Noelle, Y. Xing, and C.-W. Shu. High order well-balanced schemes. In: Puppo, G., Russo, G. (eds.) *Numerical Methods for Relaxation Systems and Balance Equations*, Quaderni di Matematica, volume 24, Seconda Università di Napoli, pp. 1–66, 2009.
- [53] B. Perthame and C.-W. Shu. On positivity-preserving finite-volume schemes for euler equations. *Numerische Mathematik*, 73(1):119–130, 1996.
- [54] M. Ricchiuto. *Contributions to the development of residual discretizations for hyperbolic conservation laws with application to shallow water flows*. HDR, Université Sciences et Technologies - Bordeaux I, December 2011.
- [55] M. Ricchiuto. Explicit residual discretizations for shallow water flows. *Aip Conference Proceedings*, 1389(1):919–922, 2011.
- [56] M. Ricchiuto. On the C-property and generalized C-property of residual distribution for the shallow water equations. *Journal of Scientific Computing*, 48:304–318, 2011.
- [57] M. Ricchiuto and R. Abgrall. Explicit Runge-Kutta residual distribution schemes for time dependent problems: Second order case. *Journal of Computational Physics*, 229(16):5653 – 5691, 2010.
- [58] M. Ricchiuto, R. Abgrall, and H. Deconinck. Application of conservative residual distribution schemes to the solution of the shallow water equations on unstructured meshes,. *J. Comput. Phys.*, 222:287–331, 2007.
- [59] M. Ricchiuto and A. Bollermann. Accuracy of stabilized residual distribution for shallow water flows including dry beds. In E.Tadmor, J.G.Liu, and A.Tzavaras, editors, *HYP08: 12th international conference on hyperbolic problems : theory, numerics, applications*, volume 67(2). AMS, American Mathematical Society, 2009.
- [60] M. Ricchiuto and A. Bollermann. Stabilized residual distribution for shallow water simulations. *J. Comput. Phys*, 228(4):1071–1115, 2009.
- [61] M. Ricchiuto, P.M. Congedo, G. Geraci, and R. Abgrall. Uncertainty propagation in shallow water long wave run up simulations. First International Conference on Frontiers of Comput. Physics: Modelling the Earth System, Boulder (CO), USA, December 2012.
- [62] M. Ricchiuto, Á. Csík, and H. Deconinck. Residual distribution for general time dependent conservation laws. *J. Comput. Phys*, 209(1):249–289, 2005.
- [63] M. Ricchiuto and A.G. Filippini. Upwind residual discretizations of the enhanced Boussinesq equations for wave propagation over complex bathymetries. SIAM Conf. on Mathematical and Computational Issues in the Geosciences, Padova (Italy), June 2013.
- [64] M. Ricchiuto and A.G. Filippini. Upwind residual discretization of enhanced boussinesq equations for wave propagation over complex bathymetries. *Journal of Computational Physics*, 271(0):306 – 341, 2014.
- [65] P.L. Roe. Fluctuations and signals - a framework for numerical evolution problems. In K.W. Morton and M.J. Baines, editors, *Numerical Methods for Fluids Dynamics*, pages 219–257. Academic Press, 1982.
- [66] P.L. Roe. Computational fluid dynamics, retrospective and prospective. *International Journal of Computational Fluid Dynamics*, 19(8):581–594, 2005.



- [67] J.A. Rossmannith. Residual distribution schemes for hyperbolic balance laws in generalized coordinates. *Numerical Modeling of Space Plasma Flows, ASP Conference Series*, 359:213–219, 2006.
- [68] D. Sarmany, M. Hubbard, and M. Ricchiuto. Unconditionally stable space-time discontinuous residual distribution for shallow-water flows. *Journal of Computational Physics*, 253:86 – 113, 2013.
- [69] D. Sarmany and M.E. Hubbard. Upwind residual distribution for shallow-water ocean modelling. *Ocean Modelling*, 64:1–11, 2013.
- [70] M. Seaïd. Non-oscillatory relaxation methods for the shallow-water equations in one and two space dimensions. *International Journal for Numerical Methods in Fluids*, 46(5):457–484, 2004.
- [71] S.P. Spekreijse. Multigrid solution of monotone second-order discretizations of hyperbolic conservation laws. *Math. Comp.*, 49:135–155, 1987.
- [72] C.E. Synolakis, E. Bernard, V. Titov, U. Kanoglu, and F. Gonzalez. Validation and verification of tsunami numerical models. *Pure and Applied Geophysics*, 165:2197–2228, 2008.
- [73] J. Szmelter and P.K. Smolarkiewicz. An edge-based unstructured mesh discretisation in geospherical framework. *Journal of Computational Physics*, 229(13):4980 – 4995, 2010.
- [74] J. Szmelter and P.K. Smolarkiewicz. An edge-based unstructured mesh framework for atmospheric flows. *Computers & Fluids*, 46(1):455 – 460, 2011.
- [75] E. Tadmor. Skew-selfadjoint form for systems of conservation laws. *J. Math. Anal. Appl.*, 103:428–442, 1984.
- [76] E. Tadmor. Entropy functions for symmetric systems of conservation laws. *J. Math. Anal. Appl.*, 122:355–359, 1987.
- [77] W.C. Thacker. Some exact solutions to the nonlinear shallow-water wave equations. *J. Fluid Mechanics*, 107:499–508, 1981.
- [78] Y. Xing and C.-W. Shu. High-order well-balanced finite volume WENO schemes and discontinuous Galerkin methods for a class of hyperbolic systems with source terms. *J. Comput. Phys.*, 214(2):567–598, 2006.
- [79] Y. Xing and C.-W. Shu. High-order finite volume weno schemes for the shallow water equations with dry states. *Advances in Water Resources*, 34(8):1026 – 1038, 2011.
- [80] Y. Xing, C.-W. Shu, and S. Noelle. On the advantage of well-balanced schemes for moving-water equilibria of the shallow water equations. *Journal of Scientific Computing*, 48:339–349, 2011.
- [81] Y. Xing, X. Zhang, and C.-W. Shu. Positivity-preserving high order well-balanced discontinuous Galerkin methods for the shallow water equations. *Advances in Water Resources*, 33(12):1476 – 1493, 2010.