



**HAL**  
open science

# Well-Balanced Second-Order Approximation of the Shallow Water Equation with Continuous Finite Elements

Pascal Azerad, Jean-Luc Guermond, Bojan Popov

► **To cite this version:**

Pascal Azerad, Jean-Luc Guermond, Bojan Popov. Well-Balanced Second-Order Approximation of the Shallow Water Equation with Continuous Finite Elements. *SIAM Journal on Numerical Analysis*, 2017, 55 (6), pp.3203 - 3224. 10.1137/17M1122463 . hal-01815500

**HAL Id: hal-01815500**

**<https://hal.science/hal-01815500>**

Submitted on 15 Jun 2018

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

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

## WELL-BALANCED SECOND-ORDER APPROXIMATION OF THE SHALLOW WATER EQUATION WITH CONTINUOUS FINITE ELEMENTS\*

PASCAL AZERAD<sup>†</sup>, JEAN-LUC GUERMOND<sup>‡</sup>, AND BOJAN POPOV<sup>‡</sup>

**Abstract.** This paper investigates a first-order and a second-order approximation technique for the shallow water equation with topography using continuous finite elements. Both methods are explicit in time and are shown to be well-balanced. The first-order method is invariant domain preserving and satisfies local entropy inequalities when the bottom is flat. Both methods are positivity preserving. Both techniques are parameter free, work well in the presence of dry states, and can be made high order in time by using strong stability preserving time stepping algorithms.

**Key words.** shallow water, well-balanced approximation, invariant domain, second-order method, finite element method, positivity preserving

**AMS subject classifications.** 65M08, 65M60, 65M12, 35L50, 35L65, 76M10

**DOI.** 10.1137/17M1122463

**1. Introduction.** The objective of this paper is to develop an invariant domain preserving well-balanced approximation of the shallow water equation with bathymetry using continuous finite elements. There are many finite volume and discontinuous Galerkin (DG) techniques available in the literature that can solve this problem efficiently up to second and higher order in space. Examples of schemes that are well balanced at rest and robust in the presence of dry states can be found, for example, in Audusse and Bristeau [1], Audusse et al. [2], Bollermann, Noelle, and Lukáčová-Medvidová [6], Gallardo, Parés, and Castro [14], Kurganov and Petrova [23], Perthame and Simeoni [27], Ricchiuto and Bollermann [28]. We refer the reader to the book of Bouchut [7] for a review on this topic, to the paper of Xing and Shu [32] for a survey on finite volume and DG methods, and to the paper [23] for a survey of central-upwind schemes. However, to the best of our knowledge, these types of approximations are not developed in the context of continuous finite elements. Or we should say that no robust continuous finite element technique is yet available in the literature that guarantees second-order accuracy, works properly in every regime (subcritical, transcritical, transcritical with hydraulic jumps, wet, and dry regions) and is well-balanced at rest. We propose such a method in the present paper. Two variants of the method are discussed: one variant is first-order accurate in space, positivity preserving, and preserves every convex invariant domain of the system in the absence of bathymetry; the other variant is second-order accurate in space and positivity preserving. Both variants are explicit in time and use continuous finite elements on unstructured meshes.

---

\*Received by the editors March 24, 2017; accepted for publication (in revised form) September 11, 2017; published electronically December 19, 2017.

<http://www.siam.org/journals/sinum/55-6/M112246.html>

**Funding:** The work of the authors was supported in part by the National Science Foundation grants DMS-1619892 and DMS-1620058, by the Air Force Office of Scientific Research, USAF, under grant/contract FA9550-15-1-0257, and by the Army Research Office under grant/contract W911NF-15-1-0517.

<sup>†</sup>Institut Montpellierain Alexander Grothendieck, UMR 5149, Université de Montpellier, 34095 Montpellier, France (azerad@math.univ-montp2.fr).

<sup>‡</sup>Department of Mathematics, Texas A&M University, 3368 TAMU, College Station, TX 77843 (guermond@math.tamu.edu, popov@math.tamu.edu).

The first building block of the method consists of using the methodology introduced in Guermond and Popov [16]. The second building block consists of making the schemes well-balanced with respect to rest states by using the so-called hydrostatic reconstruction from [2, section 2.1] and variations thereof. The technique from [16] is a loose extension of Lax's scheme [24, p. 163] to continuous finite elements; it solves general hyperbolic systems in any space dimension using forward Euler time stepping and continuous finite elements on nonuniform grids. The artificial dissipation is defined so that any convex invariant set containing the initial data is an invariant domain for the method. The solution thus constructed satisfies a discrete entropy inequality for every admissible entropy of the system. The accuracy in space is formally first order and the accuracy in time can be made high order by using strong stability preserving Runge–Kutta time stepping. Some ideas of the method are rooted in the work of Hoff [20, 21], and Frid [13]. The method is made second order and positivity preserving by using techniques introduced in Guermond and Popov [17].

The paper is organized as follows. The model problem and the finite element setting are introduced in section 2. The first-order variant of the method is described in section 3. The main results of this section are Propositions 3.9 and 3.11. The second-order variant of the method is described in section 4. The key results of this section are Propositions 4.2 and 4.4. The performances of the algorithms introduced in the paper are numerically illustrated in section 5 on standard benchmark problems.

**2. Preliminaries.** In this section we introduce the model problem, the finite element setting, and we define (recall) the concept of well-balancing at rest.

**2.1. The model problem.** Let  $D$  be a polygonal domain in  $\mathbb{R}^d$  with  $d \in \{1, 2\}$ , occupied by a body of water evolving in time under the action of gravity. Assuming that the deformations of the free surface are small compared to the water elevation and the bottom topography  $z$  varies slowly, the problem can be well represented by Saint-Venant's shallow water model. This model describes the time and space evolution of the water height  $h$  and flow rate, or discharge,  $\mathbf{q}$  in the direction parallel to the bottom. Using  $\mathbf{u} = (h, \mathbf{q})^\top$  as a dependent variable the model is as follows:

$$(2.1) \quad \partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) + \mathbf{b}(\mathbf{u}, \nabla z) = 0, \quad \mathbf{x} \in D, t \in \mathbb{R}_+,$$

$$(2.2) \quad \mathbf{f}(\mathbf{u}) := \begin{pmatrix} \mathbf{q}^\top \\ \frac{1}{h} \mathbf{q} \otimes \mathbf{q} + \frac{1}{2} g h^2 \mathbb{I}_d \end{pmatrix} \in \mathbb{R}^{(1+d) \times d}, \quad \mathbf{b}(\mathbf{u}, \nabla z) := \begin{pmatrix} 0 \\ g h \nabla z \end{pmatrix}.$$

The quantity  $\mathbf{q}$  is related to the horizontal component of the water velocity  $\mathbf{v}$  by  $\mathbf{q} = \mathbf{v}h$ . The function  $z : D \ni \mathbf{x} \mapsto z(\mathbf{x}) \in \mathbb{R}$  is the given topography.

We assume that either the boundary conditions are periodic or the initial data  $\mathbf{u}_0$  and the bottom topography are constant outside a compact set in  $D$  and the solution to (2.1) is constant outside this compact set over some time interval  $[0, T]$ .

**2.2. The finite element space.** We approximate the solution of (2.2) with continuous finite elements. Let  $(\mathcal{T}_h)_{h>0}$  be a shape-regular family of matching meshes. (Here we slightly abuse notation by denoting the mesh size by  $h$ . For instance we are going to denote by  $h_h$  the finite element approximation of the water height.) The elements in  $\mathcal{T}_h$  are assumed to be generated from a finite number of reference elements denoted  $\{\widehat{K}_r\}_{1 \leq r \leq \varpi}$ . For example, the mesh  $\mathcal{T}_h$  could be composed of a combination of triangles and quadrangles ( $\varpi = 2$  in this case). Given a set of reference finite elements in the sense of Ciarlet  $\{(\widehat{K}_r, \widehat{P}_r, \widehat{\Sigma}_r)\}_{1 \leq r \leq \varpi}$  (the index  $r \in \{1: \varpi\}$  is omitted

in the rest of the paper to simplify the notation) we introduce the finite element space

$$(2.3) \quad P(\mathcal{T}_h) := \left\{ v \in C^0(D; \mathbb{R}) \mid v|_K \circ T_K \in \widehat{P}, \forall K \in \mathcal{T}_h \right\},$$

where for any  $K \in \mathcal{T}_h$ ,  $T_K : \widehat{K} \rightarrow K$  is the geometric bijective transformation that maps the reference element  $\widehat{K}$  to the current element  $K$ . We do not assume that  $T_K$  is affine. The exact nature of the degrees of freedom in  $\widehat{\Sigma}_r$  is not essential, but the reader who is not familiar with finite elements can think of Lagrange elements or Bernstein elements. The reference space  $\widehat{P}$  is assumed to be composed of scalar-valued functions (these are polynomials usually). The reference shape functions are denoted  $\{\widehat{\theta}_i\}_{i \in \{1:n_{\text{sh}}\}}$ ; recall that they form a basis of  $\widehat{P}$ . We assume that the basis  $\{\widehat{\theta}_i\}_{i \in \{1:n_{\text{sh}}\}}$  has the partition of unity property:  $\sum_{i \in \{1:n_{\text{sh}}\}} \widehat{\theta}_i(\widehat{\mathbf{x}}) = 1$  for all  $\widehat{\mathbf{x}} \in \widehat{K}$ . The approximation in space of  $\mathbf{u}$  in (2.2) will be done in  $\mathbf{P}(\mathcal{T}_h) := [P(\mathcal{T}_h)]^{1+d}$ . The approximation of the bathymetry map will be done in  $P(\mathcal{T}_h)$ . The global shape functions in  $P(\mathcal{T}_h)$  are denoted by  $\{\varphi_i\}_{i \in \{1:I\}}$ ; the set  $\{\varphi_i\}_{i \in \{1:I\}}$  is a basis of  $P(\mathcal{T}_h)$ . The partition of unity property on the reference shape functions implies that

$$(2.4) \quad \sum_{i \in \{1:I\}} \varphi_i(\mathbf{x}) = 1 \quad \forall \mathbf{x} \in D.$$

Let  $D_i$  be the support of  $\varphi_i$  and  $|D_i|$  be the measure of  $D_i$ ,  $i \in \{1:I\}$ . For any union of cells  $E \subset \mathcal{T}_h$ , we define  $\mathcal{I}(E) := \{j \in \{1:I\} \mid |D_j \cap E| \neq 0\}$  to be the set that contains the indices of all the shape functions whose support on  $E$  is of nonzero measure. We are going to regularly invoke  $\mathcal{I}(K)$  and  $\mathcal{I}(D_i)$  and the partition of unity property  $\sum_{i \in \mathcal{I}(K)} \varphi_i(\mathbf{x}) = 1$  for all  $\mathbf{x} \in K$ .

Let  $\mathcal{M}$  be the consistent mass matrix with entries  $m_{ij} := \int_D \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) \, dx$ , and let  $\mathcal{M}^L$  be the diagonal lumped mass matrix with entries  $m_i := \int_D \varphi_i(\mathbf{x}) \, dx$ . The partition of unity property implies that  $m_i = \sum_{j \in \mathcal{I}(D_i)} m_{ij}$ . One key assumption that we use in the rest of the chapter is that

$$(2.5) \quad m_i > 0 \quad \forall i \in \{1:I\}.$$

The identities (2.4) are satisfied by all the standard finite elements and (2.5) is satisfied by many Lagrange elements and by the Bernstein–Bezier elements of any degree.

Upon denoting by  $\|\cdot\|_{\ell^2}$  the Euclidean norm in  $\mathbb{R}^d$ , we introduce the following two quantities which will play an important role in the rest of the paper:

$$(2.6) \quad \mathbf{c}_{ij} := \int_D \varphi_i \nabla \varphi_j \, dx, \quad \mathbf{n}_{ij} := \frac{\mathbf{c}_{ij}}{\|\mathbf{c}_{ij}\|_{\ell^2}}, \quad i, j \in \{1:I\}.$$

Note that (2.4) implies  $\sum_{j \in \{1:I\}} \mathbf{c}_{ij} = \mathbf{0}$ . Furthermore, if either  $\varphi_i$  or  $\varphi_j$  is zero on  $\partial D$ , then  $\mathbf{c}_{ij} = -\mathbf{c}_{ji}$ . In particular we have  $\sum_{i \in \{1:I\}} \mathbf{c}_{ij} = \mathbf{0}$  if  $\varphi_j$  is zero on  $\partial D$ . This property will be used to establish conservation.

LEMMA 2.1. *Let  $\mathbf{k} \in C^1(\mathbb{R}^{1+d}; \mathbb{R}^{(1+d) \times d})$ . Let  $\mathbf{u}_h = \sum_{j \in \{1:I\}} \mathbf{U}_j \varphi_j \in \mathbf{P}(\mathcal{T}_h)$ . Then  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{k}(\mathbf{U}_j) \cdot \mathbf{c}_{ij}$  is a second-order approximation of  $\int_D \nabla \cdot (\mathbf{k}(\mathbf{u}_h)) \varphi_i \, dx$ .*

*Proof.* Since we have  $\int_{D_i} \nabla \cdot (\mathbf{k}(\mathbf{u}_h)) \varphi_i \, dx = \sum_{j \in \{1:I\}} \mathbf{k}(\mathbf{U}_j) \int_{D_i} \varphi_i \nabla \varphi_j \, dx$  when  $\mathbf{k}$  is linear, the quantity  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{k}(\mathbf{U}_j) \cdot \mathbf{c}_{ij}$  is a second-order approximation in space of  $\int_D \nabla \cdot (\mathbf{k}(\mathbf{u}_h)) \varphi_i \, dx$ , i.e., the error scales like  $\mathcal{O}(h^2) \|\mathbf{c}_{ij}\|_{\ell^2}$ . □

DEFINITION 2.2 (centrosymmetry). *The mesh  $\mathcal{T}_h$  is said to be centrosymmetric if the following conditions hold true: (i) For all  $i \in \{1:I\}$ , there is a permutation  $\sigma_i : \mathcal{I}(D_i) \rightarrow \mathcal{I}(D_i)$  such that  $\mathbf{c}_{ij} = -\mathbf{c}_{i\sigma_i(j)}$ ; (ii) if the function  $D_i \ni \mathbf{x} \rightarrow \sum_{j \in \mathcal{I}(D_i)} \alpha_j \varphi_j(\mathbf{x}) \in \mathbb{R}$  is linear over  $D_i$  then  $\alpha_i = \frac{1}{2}(\alpha_j + \alpha_{\sigma_i(j)})$  for all  $j \in \mathcal{I}(D_i)$ .*

For instance, in the context of Lagrange elements, the centrosymmetric assumption holds if for any  $i \in \{1:I\}$  the set of the Lagrange nodes with indices in  $\mathcal{I}(D_i)$  can be partitioned into pairs that are symmetric with respect to the Lagrange node of index  $i$ . Although at some point in the paper we will invoke centrosymmetry of the mesh to establish formal consistency of some terms, we do not assume that the mesh is centrosymmetric in the rest of the paper.

**2.3. Well-balancing properties.** The concept of well-balancing originates in the seminal work of Bermudez and Vazquez [4] and Greenberg and Leroux [15]. The idea is that the scheme should at the very least preserve steady states at rest. Of course, it could be desirable to preserve *general* steady solutions, i.e., not necessarily at rest, but this is beyond the scope of the present paper. We refer the reader to Noelle, Xing, and Shu [26] where this question is addressed. Since at rest  $\mathbf{q} = \mathbf{0}$  the balance of momentum reduces to  $\mathbf{0} = g\nabla(\frac{1}{2}h^2) + gh\nabla z = gh\nabla(h+z)$ , one should have either  $h+z$  is constant (so-called wet state) or  $h$  is zero (so-called dry state). Hence a well-balanced scheme in the context of the shallow water equation is one such that, at rest, dry states remain dry and  $h+z$  remains constant for wet states. This property is not easy to satisfy for approximation techniques that are second order and higher in space. We refer the reader to Bouchut [7] for a concise account and further references on well-balanced schemes. In this paper we are going to adapt to continuous finite elements a methodology proposed in Audusse and Bristeau [1], Audusse et al. [2] known as the “hydrostatic reconstruction” technique.

Let  $z_h = \sum_{i=1}^I Z_i \varphi_i \in P(\mathcal{T}_h)$  be the approximation of the bathymetry map. Let  $h_h = \sum_{i=1}^I H_i \varphi_i \in P(\mathcal{T}_h)$  be the approximation of the water height. Let  $\mathbf{q}_h = \sum_{i=1}^I \mathbf{Q}_i \varphi_i$  be the approximation of the flow rate. Let us now define the rest state. Curiously, defining a rest state is not as trivial as it sounds. We are going to use two definitions. One of them makes use of the following quantity which is known in the literature as the hydrostatic reconstruction of the water height:

$$(2.7) \quad H_i^{*,j} := \max(0, H_i + Z_i - \max(Z_i, Z_j)) \quad \forall i \in \{1:I\}, j \in \mathcal{I}(D_i).$$

To better understand this definition, assume that the water is at rest and consider for instance a dry node  $j$  in the neighborhood of a wet node  $i$ , i.e.,  $j \in \mathcal{I}(D_i)$ , see the left panel of Figure 1. In this case  $H_j = 0$  and  $Z_j \geq H_i + Z_i$ , which then implies  $H_i^{*,j} = H_j^{*,i}$ . Similarly if both  $i$  and  $j$  are dry states we have  $H_i^{*,j} = H_j^{*,i}$ , and if both  $i$  and  $j$  are wet states and are such that  $H_j + Z_j = H_i + Z_i$  we also have  $H_i^{*,j} = H_j^{*,i}$ . These observations motivate the following definition.

DEFINITION 2.3 (rest at large). *A numerical state  $(h_h, \mathbf{q}_h, z_h)$  is said to be at rest at large if the approximate momentum  $\mathbf{q}_h$  is zero, and if the approximate water height  $h_h$  and the approximate bathymetry map  $z_h$  satisfy the following property, for all  $i \in \{1:I\}$ :  $H_i^{*,j} = H_j^{*,i}$  for all  $j \in \mathcal{I}(D_i)$ .*

DEFINITION 2.4 (exact rest). *A numerical state  $(h_h, \mathbf{q}_h, z_h)$  is said to be at exact rest (or exactly at rest) if  $\mathbf{q}_h$  is zero, and if the approximate water height  $h_h$  and the approximate bathymetry map  $z_h$  satisfy the following alternative, for all  $i \in \{1:I\}$ : for all  $j \in \mathcal{I}(D_i)$ , either  $H_j = H_i = 0$  or  $H_j + Z_j = H_i + Z_i$ .*

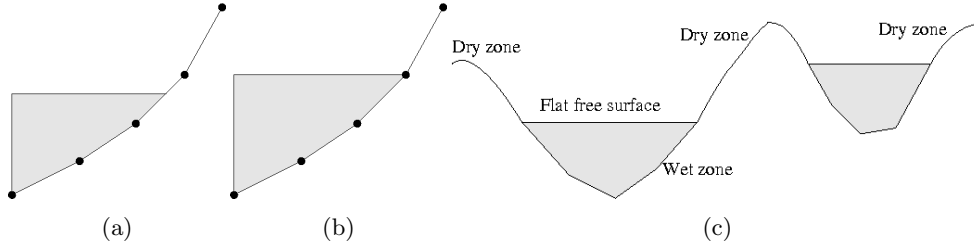


FIG. 1. Configuration (a) is not an exact rest state according to Definition 2.4 whereas configuration (b) is. Both states are at rest at large. Panel (c) shows a typical steady state at rest with wet and dry areas.

The existence of an exact rest state is a compatibility condition between the mesh and the initial data. This compatibility condition is not satisfied by the configuration depicted in the left panel of Figure 1 whereas it is satisfied by the configuration in the center panel. Exact rest implies rest at large. Note in passing that the zone where  $h + z$  is constant may not be connected; that is to say, it is possible to have different free surface heights in disconnected wet zones as shown in the right panel of Figure 1.

DEFINITION 2.5 (well-balancing at large). (i) A function  $\mathbf{K} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbb{R}^I \times (\mathbb{R}^I)^d$  is said to be a well-balanced flux approximation at large if  $\mathbf{K}(\mathbf{u}_h) = 0$  when  $\mathbf{u}_h$  is a rest state at large according to Definition 2.3. (ii) A mapping  $\mathbf{S} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbf{P}(\mathcal{T}_h)$  is a well-balanced scheme at large if  $\mathbf{S}(\mathbf{u}_h) = \mathbf{u}_h$  when  $\mathbf{u}_h$  is a rest state at large.

DEFINITION 2.6 (exact well-balancing). (i) A function  $\mathbf{K} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbb{R}^I \times (\mathbb{R}^I)^d$  is said to be an exactly well-balanced flux approximation if  $\mathbf{K}(\mathbf{u}_h) = 0$  when  $\mathbf{u}_h$  is an exact rest state according to Definition 2.4. (ii) A mapping  $\mathbf{S} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbf{P}(\mathcal{T}_h)$  is an exactly well-balanced scheme if  $\mathbf{S}(\mathbf{u}_h^n) = \mathbf{u}_h^n$  when  $\mathbf{u}_h^n$  is an exact rest state.

DEFINITION 2.7 (conservation). We say that  $\mathbf{u}_h^n \rightarrow \mathbf{u}_h^{n+1}$  is a conservative finite element approximation of (2.1) if  $\sum_{i \in \{1:I\}} m_i \mathbf{H}_i^n = \sum_{i \in \{1:I\}} m_i \mathbf{H}_i^{n+1}$  and if  $\sum_{i \in \{1:I\}} m_i \mathbf{Q}_i^n = \sum_{i \in \{1:I\}} m_i \mathbf{Q}_i^{n+1}$  when the topography map is constant.

**3. First order scheme.** We describe in this section a time and space approximation of (2.2). The scheme is well-balanced at large but approximates the flux to first order in space only. This scheme satisfies local invariant domain properties and local discrete entropy inequalities when the bottom is flat. It is an adaptation of the method presented in Audusse et al. [2] to the continuous finite element setting developed in Guermond and Popov [16]. To the best of our knowledge, this is the first result of this type for continuous finite elements.

**3.1. Flux approximation.** Just like in [2, (2.13)], the key is to consider the hydrostatic reconstruction (2.7) and to observe that  $\sum_{j \in \mathcal{I}(D_i)} \frac{1}{2} ((\mathbf{H}_j^{*,i})^2 - (\mathbf{H}_i^{*,j})^2) \mathbf{c}_{ij}$  is a well-balanced first-order approximation of the flux  $\int_{D_i} (\nabla(\frac{1}{2}h^2) + h\nabla z) \varphi_i \, dx$ .

LEMMA 3.1 (consistency/well-balancing). (i) Assume that  $\{\hat{\theta}_n\}_{n \in \{1:n_{sh}\}}$  consists of Lagrange or Bernstein functions. Then  $\sum_{j \in \mathcal{I}(D_i)} \frac{1}{2} ((\mathbf{H}_j^{*,i})^2 - (\mathbf{H}_i^{*,j})^2) \mathbf{c}_{ij}$  is a first-order approximation of the flux  $\int_{D_i} (\nabla(\frac{1}{2}h^2) + h\nabla z) \varphi_i \, dx$ . (ii) The mapping  $\mathbf{u}_h \rightarrow (0, \sum_{j \in \mathcal{I}(D_i)} \frac{1}{2} ((\mathbf{H}_j^{*,i})^2 - (\mathbf{H}_i^{*,j})^2) \mathbf{c}_{ij})_{i \in \{1:I\}}$  is well-balanced at large.

*Proof.* (i) Let us fix  $i \in \{1:I\}$ . We slightly abuse the notation by using  $h$  to denote the mesh size. For the consistency analysis we assume that the water height and the bathymetry map are smooth and the water height is nonnegative. More

precisely, we assume that there is  $C_z$  such that for all  $i \in \{1:I\}$ ,  $|Z_i - Z_j| \leq C_z h$  for all  $j \in \mathcal{I}(D_i)$ .

Assume first that  $Z_j \geq Z_i$ . We immediately get  $H_j^{*,i} = H_j$ . If, in addition,  $H_i \geq C_z h$ , then  $H_i^{*,j} = \max(0, H_i + (Z_i - Z_j)) = H_i + (Z_i - Z_j)$ , and we have  $\frac{1}{2}((H_j^{*,i})^2 - (H_i^{*,j})^2) = \frac{1}{2}H_j^2 - \frac{1}{2}(H_i + (Z_i - Z_j))^2 = \frac{1}{2}H_j^2 - \frac{1}{2}H_i^2 + H_i(Z_j - Z_i) + \mathcal{O}(h^2)$ . Similarly, if  $H_i \leq C_z h$ , then  $H_i^{*,j} = \mathcal{O}(h)$  and we again have  $\frac{1}{2}((H_j^{*,i})^2 - (H_i^{*,j})^2) = \frac{1}{2}H_j^2 - \frac{1}{2}H_i^2 + H_i(Z_j - Z_i) + \mathcal{O}(h^2)$ . On the other hand, if  $Z_i \leq Z_j$ , we obtain  $\frac{1}{2}((H_j^{*,i})^2 - (H_i^{*,j})^2) = \frac{1}{2}H_j^2 - \frac{1}{2}H_i^2 + H_j(Z_j - Z_i) + \mathcal{O}(h^2)$ . But since  $H_j = H_i + \mathcal{O}(h)$  (we are using continuous finite elements and the water height is assumed to be smooth), we also have  $\frac{1}{2}((H_j^{*,i})^2 - (H_i^{*,j})^2) = \frac{1}{2}H_j^2 - \frac{1}{2}H_i^2 + H_i(Z_j - Z_i) + \mathcal{O}(h^2)$  in this case.

Using Lemma 2.1 we infer that  $\sum_{j \in \mathcal{I}(D_i)} (\frac{1}{2}H_j^2 - \frac{1}{2}H_i^2) \mathbf{c}_{ij}$  is a second-order approximation of  $\int_D (\nabla(\frac{1}{2}h^2)) \varphi_i dx$ . Similarly,  $\sum_{j \in \mathcal{I}(D_i)} (H_i(Z_j - Z_i)) \mathbf{c}_{ij}$  is a second-order approximation of  $H_i \int_D (\nabla z) \varphi_i dx$ . If  $z$  is linear over  $\mathcal{D}_i$  (which is a sufficient assumption for the consistency analysis), then  $H_i \int_D (\nabla z) \varphi_i dx = \nabla z|_{D_i} H_i \int_D \varphi_i dx$ . Since  $H_i \int_D \varphi_i dx$  can be shown to be a second-order approximation of  $\int_{D_i} h \varphi_i dx$  (at least for Lagrange and Bernstein basis functions), we conclude that  $\sum_{j \in \mathcal{I}(D_i)} (H_i(Z_j - Z_i)) \mathbf{c}_{ij}$  is a second-order approximation of  $\int_D (h \nabla z) \varphi_i dx$ . Combining these observations with the above argument and upon observing that  $\|\mathbf{c}_{ij}\|_{\ell^2} \mathcal{O}(h^2) = m_i \mathcal{O}(h)$ , we conclude that  $\sum_{j \in \mathcal{I}(D_i)} \frac{1}{2}((H_j^{*,i})^2 - (H_i^{*,j})^2) \mathbf{c}_{ij}$  is a first-order approximation of  $\int_D (\nabla(\frac{1}{2}h^2) + h \nabla z) \varphi_i dx$ .

(ii) Let us prove the well-balancing at large. Assuming that  $\mathbf{u}_h$  is a rest state at large, according to Definition 2.3 we have  $H_j^{*,i} = H_i^{*,j}$ , hence  $(H_j^{*,i})^2 - (H_i^{*,j})^2 = 0$ . The conclusion follows immediately.  $\square$

Let us introduce the gas dynamics flux  $\mathbf{g}(\mathbf{u}) := (\mathbf{q}, \frac{1}{h} \mathbf{q} \otimes \mathbf{q})^\top$ . We now need to approximate  $\int_{D_i} \mathbf{g}(\mathbf{u}) \varphi_i dx$ . Since we have seen above that using  $H^*$  is a good idea to guarantee well-balancing at large, one could imagine working with the pair  $(H_i^{*,j}, \mathbf{Q}_i)^\top$ . The problem with this choice is that if it happens that  $H_i^{*,j}$  is zero (because  $H_i + Z_i \leq \max(Z_i, Z_j)$ ), there is no reason for the approximate flow rate  $\mathbf{Q}_i$  to be zero; hence the quantity  $\mathbf{Q}_i/H_i^{*,j}$  which approximates the velocity could be unbounded. To avoid this problem, we proceed as in [2] by working with the quantities

$$(3.1) \quad \mathbf{Q}_i^{*,j} := \mathbf{Q}_i \frac{H_i^{*,j}}{H_i}, \quad \mathbf{U}_i^{*,j} := \left( H_i^{*,j}, \mathbf{Q}_i^{*,j} \right)^\top$$

with the convention that  $\mathbf{Q}_i^{*,j} := 0$  if  $H_i = 0$ . Note that we have  $\|\mathbf{Q}_i^{*,j}\|_{\ell^2} \leq \|\mathbf{Q}_i\|_{\ell^2}$  since  $0 \leq H_i^{*,j} \leq H_i$  by definition. We now face the question of constructing a consistent approximation of  $\int_{D_i} \mathbf{g}(\mathbf{u}) \varphi_i dx$  using the state variable  $\mathbf{U}_i^{*,j}$ . To simplify the notation let us introduce the approximate velocity  $\mathbf{v}_h = \sum_{i \in \{1:I\}} \mathbf{V}_i \varphi_i$  with

$$(3.2) \quad \mathbf{V}_i := \frac{\mathbf{Q}_i}{H_i}, \quad i \in \{1:I\}.$$

DEFINITION 3.2 (shoreline). *We say that a degree of freedom  $i$  is away from the shoreline if either  $H_j = 0$  for all  $j \in \mathcal{I}(D_i)$  or  $\min(H_j, H_i) > |Z_i - Z_j|$  for all  $j \in \mathcal{I}(D_i)$ .*

Note that if the bottom topography is smooth, i.e., there is  $C_z$  such that for all  $i \in \{1:I\}$ ,  $|Z_i - Z_j| \leq C_z h$ , then any degree of freedom  $i$  such that  $H_j \geq C_z h$  for all  $j \in \mathcal{I}(D_i)$ , is away from the shoreline according to the above definition. Roughly speaking, a degree of freedom  $i$  is said to be away from the shoreline if either all the

degrees of freedom around  $i$  are dry or the water depth around  $i$  is at least  $C_z h$  if the bottom topography is smooth ( $h$  being the mesh size).

LEMMA 3.3. *The quantity  $\sum_{j \in \mathcal{I}(D_i)} (\mathbf{g}(\mathbf{U}_j^{*,i}) + \mathbf{g}(\mathbf{U}_i^{*,j})) \cdot \mathbf{c}_{ij}$  is a first-order approximation of  $\int_{D_i} \nabla \cdot \mathbf{g}(\mathbf{u}) \varphi_i \, dx$  away from the shoreline if the mesh is centrosymmetric.*

*Proof.* Let  $i \in \{1: I\}$  be a degree of freedom away from the shoreline. The approximation of the flux is  $\sum_{j \in \mathcal{I}(D_i)} (\mathbf{V}_j \mathbf{H}_j^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij}$  for the mass conservation equation and  $\sum_{j \in \mathcal{I}(D_i)} ((\mathbf{V}_j \otimes \mathbf{V}_j) \mathbf{H}_j^{*,i} + (\mathbf{V}_i \otimes \mathbf{V}_i) \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij}$  for the flow rate conservation. Let us start with the mass conservation equation. We proceed as in the proof of Lemma 3.1 and again assume that the water height and the bathymetry map are smooth and the water height is nonnegative. Since the mesh is centrosymmetric by hypothesis, we can assume without loss of generality that  $Z_j \geq Z_i \geq Z_{\sigma_i(j)}$ . Then  $\mathbf{H}_j^{*,i} = \mathbf{H}_j$  and since  $i$  is away from the shoreline we have either  $\mathbf{H}_i^{*,j} = \mathbf{H}_i + Z_i - Z_j$  if  $\mathbf{H}_i \neq 0$ , or  $\mathbf{H}_i^{*,j} = 0$  if  $\mathbf{H}_i = 0$ . Similarly,  $\mathbf{H}_i^{*,\sigma_i(j)} = \mathbf{H}_i$  and since  $i$  is away from the shoreline we have either  $\mathbf{H}_{\sigma_i(j)}^{*,i} = \mathbf{H}_{\sigma_i(j)} + Z_{\sigma_i(j)} - Z_i$  if  $\mathbf{H}_{\sigma_i(j)} \neq 0$ , or  $\mathbf{H}_{\sigma_i(j)}^{*,i} = 0$  if  $\mathbf{H}_{\sigma_i(j)} = 0$ . Hence, if  $i$  is a wet state (and all the states in  $\mathcal{I}(D_i)$  are wet since  $i$  is away from the shoreline), we have

$$\begin{aligned} & (\mathbf{V}_j \mathbf{H}_j^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij} + (\mathbf{V}_{\sigma_i(j)} \mathbf{H}_{\sigma_i(j)}^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,\sigma_i(j)}) \cdot \mathbf{c}_{i\sigma_i(j)} \\ &= (\mathbf{V}_j \mathbf{H}_j + \mathbf{V}_i (\mathbf{H}_i + Z_i - Z_j) - (\mathbf{V}_{\sigma_i(j)} (\mathbf{H}_{\sigma_i(j)} + Z_{\sigma_i(j)} - Z_i) + \mathbf{V}_i \mathbf{H}_i)) \cdot \mathbf{c}_{ij} \\ &= (\mathbf{V}_j \mathbf{H}_j - \mathbf{V}_i \mathbf{H}_i) \cdot \mathbf{c}_{ij} + (\mathbf{V}_{\sigma_i(j)} \mathbf{H}_{\sigma_i(j)} - \mathbf{V}_i \mathbf{H}_i) \cdot \mathbf{c}_{i\sigma_i(j)} \\ & \quad + \mathbf{V}_i (Z_i - Z_j) \cdot \mathbf{c}_{ij} + \mathbf{V}_{\sigma_i(j)} (Z_{\sigma_i(j)} - Z_i) \cdot \mathbf{c}_{i\sigma_i(j)}, \end{aligned}$$

where we have used the centrosymmetry property  $\mathbf{c}_{ij} = -\mathbf{c}_{i\sigma_i(j)}$ . If  $i$  is a dry state (recall that  $j$  and  $\sigma_i(j)$  are also dry states since  $i$  is away from the shoreline) then

$$\begin{aligned} & (\mathbf{V}_j \mathbf{H}_j^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij} + (\mathbf{V}_{\sigma_i(j)} \mathbf{H}_{\sigma_i(j)}^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,\sigma_i(j)}) \cdot \mathbf{c}_{i\sigma_i(j)} \\ &= (\mathbf{V}_j \mathbf{H}_j - \mathbf{V}_i \mathbf{H}_i) \cdot \mathbf{c}_{ij} + (\mathbf{V}_{\sigma_i(j)} \mathbf{H}_{\sigma_i(j)} - \mathbf{V}_i \mathbf{H}_i) \cdot \mathbf{c}_{i\sigma_i(j)}. \end{aligned}$$

Since according to Lemma 2.1,  $\sum_{j \in \mathcal{I}(D_i)} (\mathbf{V}_j \mathbf{H}_j - \mathbf{V}_i \mathbf{H}_i) \cdot \mathbf{c}_{ij} = \sum_{j \in \mathcal{I}(D_i)} \mathbf{V}_j \mathbf{H}_j \cdot \mathbf{c}_{ij}$  is a second-order approximation of  $\int_D \nabla \cdot (\mathbf{v}_h h_h) \varphi_i \, dx$ , we have to show that the contribution of the extra term  $\mathbf{V}_i (Z_i - Z_j) \cdot \mathbf{c}_{ij} - \mathbf{V}_{\sigma_i(j)} (Z_{\sigma_i(j)} - Z_i) \cdot \mathbf{c}_{ij}$  that arises when  $i$  is a wet state is small. Assuming that the velocity is smooth, we have  $\mathbf{V}_{\sigma_i(j)} = \mathbf{V}_i + \mathcal{O}(h)$ , which shows that  $\mathbf{V}_i (Z_i - Z_j) \cdot \mathbf{c}_{ij} - \mathbf{V}_{\sigma_i(j)} (Z_{\sigma_i(j)} - Z_i) \cdot \mathbf{c}_{ij} = \mathbf{V}_i (2Z_i - Z_j - Z_{\sigma_i(j)}) \cdot \mathbf{c}_{ij} + \|\mathbf{c}_{ij}\|_{\ell^2} \mathcal{O}(h^2)$ . The centrosymmetry assumption implies that  $2Z_i - Z_j - Z_{\sigma_i(j)} = \mathcal{O}(h^2)$  if the bathymetry map is smooth. In conclusion  $\sum_{j \in \mathcal{I}(D_i)} (\mathbf{V}_j \mathbf{H}_j^{*,i} + \mathbf{V}_i \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij} = \sum_{j \in \mathcal{I}(D_i)} \mathbf{V}_j \mathbf{H}_j \cdot \mathbf{c}_{ij} + m_i \mathcal{O}(h)$  away from the shoreline. Using the same argument one proves that

$$\sum_{j \in \mathcal{I}(D_i)} ((\mathbf{V}_j \otimes \mathbf{V}_j) \mathbf{H}_j^{*,i} + (\mathbf{V}_i \otimes \mathbf{V}_i) \mathbf{H}_i^{*,j}) \cdot \mathbf{c}_{ij} = \sum_{j \in \mathcal{I}(D_i)} (\mathbf{V}_j \otimes \mathbf{V}_j) \mathbf{H}_j + m_i \mathcal{O}(h).$$

This concludes the proof. □

Remark 3.4 (hydrostatic reconstruction). The lack of consistency of the hydrostatic reconstruction at the shoreline or in the presence of large gradients in the topography map has been identified in Delestre et al. [10, Prop. 2.1]. Various alternatives to the hydrostatic reconstruction have since been proposed like in Berthon and Foucher [5], Bryson et al. [9], Duran, Liang, and Marche [12], where the authors propose to work with the free surface elevation instead of the water height.



**3.2. Full time and space approximation.** Let  $\mathbf{u}_h^0 = \sum_{i=1}^I \mathbf{U}_i^0 \varphi_i \in \mathcal{P}(\mathcal{T}_h)$  be a reasonable approximation of  $\mathbf{u}_0$ . Let  $n \in \mathbb{N}$ ,  $\tau$  be the time step,  $t_n$  be the current time, and let us set  $t_{n+1} = t_n + \tau$ . Let  $\mathbf{u}_h^n = \sum_{i=1}^I \mathbf{U}_i^n \varphi_i \in \mathcal{P}(\mathcal{T}_h)$  be the space approximation of  $\mathbf{u}$  at time  $t_n$ . Upon denoting  $H_i^{*,j,n} := \max(0, H_i^n + Z_i - \max(Z_i, Z_j))$ , we propose to estimate  $\mathbf{U}_i^{n+1}$  as follows:

$$(3.3) \quad m_i \frac{\mathbf{U}_i^{n+1} - \mathbf{U}_i^n}{\tau} + \sum_{j \in \mathcal{I}(D_i)} (\mathbf{g}(\mathbf{U}_j^{*,i,n}) + \mathbf{g}(\mathbf{U}_i^{*,j,n})) \cdot \mathbf{c}_{ij} + \left( \frac{1}{2} g((H_j^{*,i,n})^2 - (H_i^{*,j,n})^2) \mathbf{c}_{ij} \right) - \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) = 0,$$

where the artificial viscosity coefficient  $d_{ij}^n$  is defined by

$$(3.4) \quad d_{ij}^n := \max(d_{ij}^{\mathbf{f},n}, d_{ji}^{\mathbf{f},n}),$$

$$(3.5) \quad d_{ij}^{\mathbf{f},n} := \max\left(\lambda_{\max}^{\mathbf{f}}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^{*,i,n}), \lambda_{\max}^{\mathbf{f}}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_i^{*,j,n})\right) \|\mathbf{c}_{ij}\|_{\ell^2},$$

and  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)$  is the maximum wave speed in the Riemann problem:

$$(3.6) \quad \partial_t \mathbf{u} + \partial_x(\mathbf{f}(\mathbf{u}) \cdot \mathbf{n}) = 0, \quad \mathbf{u}(x, 0) = (1 - H(x))\mathbf{U}_L + H(x)\mathbf{U}_R,$$

where  $H(x)$  is the Heaviside function. Note that  $d_{ij}^n \geq 0$  and  $d_{ij}^n = d_{ji}^n$  for all  $j \neq i$  in  $\mathcal{I}(D_i)$ . For convenience we denote  $d_{ii}^n := -\sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n$ . Therefore we have  $\sum_{j \in \mathcal{I}(D_i)} d_{ij}^n = \sum_{j \in \mathcal{I}(D_i)} d_{ji}^n = 0$ ; this property will be used in the rest of the paper.

**3.3. Reduction to the one-dimensional Riemann problem.** For completeness, we show how the estimation of  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)$  can be reduced to estimating the maximum wave speed in a one-dimensional Riemann problem independent of  $\mathbf{n}$ . Similarly to [16], we make a change of basis and introduce  $\mathbf{t}_1, \dots, \mathbf{t}_{d-1} \in \mathbb{R}^d$  so that  $\{\mathbf{n}, \mathbf{t}_1, \dots, \mathbf{t}_{d-1}\}$  is an orthonormal basis of  $\mathbb{R}^d$ . With respect to this basis we have that  $\mathbf{q} = (q, \mathbf{q}^\perp)$ , where  $q := \mathbf{q} \cdot \mathbf{n}$  and  $\mathbf{q}^\perp := (\mathbf{q} \cdot \mathbf{t}_1, \dots, \mathbf{q} \cdot \mathbf{t}_{d-1})^\top$ . Then, with the notation  $v = q/h$ , the Riemann problem (3.6) can be rewritten in the new orthonormal basis as follows:

$$(3.7) \quad \partial_t \mathbf{u} + \partial_x(\mathbf{n} \cdot \mathbf{f}(\mathbf{u})) = \mathbf{0}, \quad \mathbf{u} = \begin{pmatrix} h \\ q \\ \mathbf{q}^\perp \end{pmatrix}, \quad \mathbf{f}(\mathbf{u}) \cdot \mathbf{n} = \begin{pmatrix} q \\ vq + \frac{q}{2}h^2 \\ v\mathbf{q}^\perp \end{pmatrix}$$

with data  $\mathbf{U}_L = (h_L, q_L, \mathbf{q}_L^\perp)^\top$ ,  $\mathbf{U}_R = (h_R, q_R, \mathbf{q}_R^\perp)^\top$ . The solution to (3.7) is henceforth denoted  $\mathbf{u}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)(x, t)$ . Following [16], we introduce the following definition.

**DEFINITION 3.5** (invariant set). *A convex set  $A \subset \mathcal{A}$  is said to be invariant for the flat bottom system, i.e., (2.1) with  $\mathbf{b} = 0$ , if for any admissible pair  $(\mathbf{U}_L, \mathbf{U}_R) \in A \times A$  and any unit vector  $\mathbf{n} \in \mathbb{R}^d$ , we have  $\mathbf{u}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)(x, t) \in A$  for a.e.  $x \in \mathbb{R}$ ,  $t > 0$ .*

Let  $\bar{\mathbf{u}}(t, \mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) := \int_{-\frac{1}{2}}^{\frac{1}{2}} \mathbf{u}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)(x, t) dx$ . Then, the following result is a consequence of  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)$  being finite; see [16, Lem. 2.1].

**LEMMA 3.6** (invariant set and average). (i) *Let  $A \subset \mathcal{A}$  be an invariant set for the flat bottom system. If  $(\mathbf{U}_L, \mathbf{U}_R) \in A$ , then  $\bar{\mathbf{u}}(t, \mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) \in A$ .* (ii) *Assume that  $2t \lambda_{\max}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) \leq 1$ , then  $\bar{\mathbf{u}}(t, \mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) = \frac{1}{2}(\mathbf{U}_L + \mathbf{U}_R) - t(\mathbf{f}(\mathbf{U}_R) - \mathbf{f}(\mathbf{U}_L)) \cdot \mathbf{n}$ .*

This lemma is the key motivation for the definition of the viscosity coefficients  $d_{ij}^{\mathbf{f},n}$  in (3.5) (see [16, section 3.3] for more details).

The maximum wave speed in the Riemann problem (3.7) is determined by the one-dimensional shallow water system for the component  $(h, q)^T$  because the last component is just passively transported and does not influence the first two equations of the system. That is to say (3.7) reduces to solving the Riemann problem

$$(3.8) \quad \partial_t(h, q)^T + \partial_x(\mathbf{f}_{1D}(h, q)) = 0$$

with data  $\mathbf{u}_L := (h_L, q_L)$ ,  $\mathbf{u}_R := (h_R, q_R)$  and flux  $\mathbf{f}_{1D}(h, q) := (q, vq + \frac{g}{2}h^2)^T$ . This establishes the following result which will be useful to estimate  $d_{ij}^{\mathbf{f};n}$  in (3.5).

**PROPOSITION 3.7** (maximum wave speed). *Let  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R)$ ,  $\lambda_{\max}^{\mathbf{f}_{1D}}(\mathbf{u}_L, \mathbf{u}_R)$  be the maximum wave speed in the Riemann problems (3.7) and (3.8), respectively. Then  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) = \lambda_{\max}^{\mathbf{f}_{1D}}(\mathbf{u}_L, \mathbf{u}_R)$ .*

In order to estimate  $\lambda_{\max}^{\mathbf{f}_{1D}}(\mathbf{u}_L, \mathbf{u}_R)$  from above, we introduce

$$(3.9) \quad \lambda_1^-(h_*) := v_L - \sqrt{gh_L} \left( 1 + \left( \frac{h_* - h_L}{2h_L} \right)_+ \right)^{\frac{1}{2}} \left( 1 + \left( \frac{h_* - h_L}{h_L} \right)_+ \right)^{\frac{1}{2}},$$

$$(3.10) \quad \lambda_2^+(h_*) := v_R + \sqrt{gh_R} \left( 1 + \left( \frac{h_* - h_R}{2h_R} \right)_+ \right)^{\frac{1}{2}} \left( 1 + \left( \frac{h_* - h_R}{h_R} \right)_+ \right)^{\frac{1}{2}}.$$

The following result is proved in Guermond and Popov [18]:

**LEMMA 3.8.** *Let  $h_{\min} = \min(h_L, h_R)$ ,  $h_{\max} = \max(h_L, h_R)$ ,  $x_0 = (2\sqrt{2} - 1)^2$ , and*

$$\bar{h}_* := \begin{cases} \frac{(v_L - v_R + 2\sqrt{gh_L} + 2\sqrt{gh_R})_+^2}{16g} & \text{if case 1,} \\ \left( -\sqrt{2h_{\min}} + \sqrt{3h_{\min} + 2\sqrt{2}h_{\min}h_{\max}} + \sqrt{\frac{2}{g}(v_L - v_R)\sqrt{h_{\min}}} \right)^2 & \text{if case 2,} \\ \sqrt{h_{\min}h_{\max}} \left( 1 + \frac{\sqrt{2}(v_L - v_R)}{\sqrt{gh_{\min} + \sqrt{gh_{\max}}}} \right) & \text{if case 3,} \end{cases}$$

where case 1 is  $0 \leq f(x_0h_{\min})$ , case 2 is  $f(x_0h_{\min}) < 0 \leq f(x_0h_{\max})$ , and case 3 is  $f(x_0h_{\max}) < 0$ . Then  $\lambda_{\max}^{\mathbf{f}}(\mathbf{n}, \mathbf{U}_L, \mathbf{U}_R) = \lambda_{\max}^{\mathbf{f}_{1D}}(\mathbf{u}_L, \mathbf{u}_R) \leq \max(|\lambda_1^-(\bar{h}_*)|, |\lambda_2^+(\bar{h}_*)|)$ .

**3.4. Stability properties.** We collect in this section some remarkable stability properties of the scheme defined by (3.3)–(3.5).

**PROPOSITION 3.9** (well-balancing/conservation). *The scheme defined in (3.3) is well-balanced at large, and it is conservative in the sense of Definition 2.7.*

*Proof.* Let  $\mathbf{u}_h^n$  be a rest state at large, then  $\mathbf{H}_j^{*,i,n} = \mathbf{H}_i^{*,j,n}$  for all  $i \in \{1:I\}$  and all  $j \in \mathcal{I}(D_i)$ ; this identity implies well-balancing at large. Let us now establish conservation. Since  $\mathbf{c}_{ij} = -\mathbf{c}_{ji}$  and  $d_{ij}^n = d_{ji}^n$  we have

$$\sum_{i \in \{1:I\}} \sum_{j \in \mathcal{I}(D_i)} \mathbf{c}_{ji} \alpha_{ij} = 0, \quad \sum_{i \in \{1:I\}} \sum_{j \in \mathcal{I}(D_i)} d_{ji}^n \beta_{ij} = 0$$

for any symmetric field  $\alpha_{ij} = \alpha_{ji}$  and any skew-symmetric field  $\beta_{ij} = -\beta_{ji}$ . Hence, we only have to deal with the nonconservative flux in (3.3),  $\frac{1}{2}g((\mathbf{H}_j^{*,i,n})^2 - (\mathbf{H}_i^{*,j,n})^2)\mathbf{c}_{ij}$ . This quantity is zero for a constant topography map. This concludes the proof.  $\square$

Since the shallow water system makes sense only for nonnegative water heights, and the water discharge should be zero in dry states, we are led to consider the following definition for the admissibility of shallow water states.

DEFINITION 3.10 (admissible water states). *A shallow water state  $\mathbf{U} = (\mathbf{H}, \mathbf{Q})^\top$  is admissible if  $\mathbf{H} \geq 0$  and  $\mathbf{Q} = \mathbf{0}$  if  $\mathbf{H} = 0$ . The set of admissible states is denoted  $\mathcal{A}$ .*

Note that a convex combination of admissible states is always an admissible state.

PROPOSITION 3.11 (invariant domain). *Let  $\mathbf{u}_h^{n+1}$  be given by (3.3)–(3.5),  $n \geq 0$ . Let  $\ell \in \{1: I\}$ . Assume that  $1 + 4\frac{\tau}{m_i}d_{ii}^n \geq 0$ . Let  $A_i^n$  be an invariant set of the shallow water equation that contains  $\{\mathbf{U}_j^n\}_{j \in \mathcal{I}(D_i)}$ . Then the following properties hold true:*

- (i) *If the bathymetry map is constant then  $\mathbf{U}_i^{n+1} \in A_i^n$ .*
- (ii) *If the bathymetry is not constant, let*

$$\Delta \mathbf{Z}_i^n := \frac{\tau}{m_i} \sum_{i \neq j \in \mathcal{I}(D_i)} g((\mathbf{H}_i^n)^2 - (\mathbf{H}_i^{*,j,n})^2) \mathbf{c}_{ij}$$

and  $\Delta \mathbf{U}_i^{*,n} := \frac{2\tau}{m_i} \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n (1 - \frac{\mathbf{H}_i^{*,j,n}}{\mathbf{H}_i^n}) \mathbf{U}_i^n$  then  $\mathbf{U}_i^{n+1} \in \text{conv}(A_i^n, \mathbf{0}) + (0, \Delta \mathbf{Z}_i^n)^\top + \Delta \mathbf{U}_i^{*,n}$ ; in particular the scheme preserves the nonnegativity of the water height.

- (iii) *If the states  $\{\mathbf{U}_i^n\}$  are admissible then the states  $\{\mathbf{U}_i^{n+1}\}$  are also admissible.*

*Proof.* Recalling that  $\mathbf{f}(\mathbf{u}) = \mathbf{g}(\mathbf{u}) + (0, \frac{1}{2}gh^2\mathbb{1}_d)^\top$ , then (3.3) can also be rewritten

$$\begin{aligned} \frac{m_i}{\tau} (\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) + \sum_{j \in \mathcal{I}(D_i)} \mathbf{f}(\mathbf{U}_j^{*,i,n}) \cdot \mathbf{c}_{ij} - d_{ij}^m \mathbf{U}_j^{*,i,n} + \mathbf{f}(\mathbf{U}_i^{*,j,n}) \cdot \mathbf{c}_{ij} - d_{ij}^m \mathbf{U}_i^{*,j,n} \\ + \sum_{j \in \mathcal{I}(D_i)} \left( 0, -g(\mathbf{H}_i^{*,j,n})^2 \mathbf{c}_{ij} \right)^\top + (d_{ij}^m + d_{ij}^m) \mathbf{U}_i^{*,j,n} = \mathbf{0}. \end{aligned}$$

Using conservation, i.e.,  $\mathbf{c}_{ii} = -\sum_{i \neq j \in \mathcal{I}(D_i)} \mathbf{c}_{ij}$ , this equation can be recast into

$$\begin{aligned} \frac{m_i}{\tau} (\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) \\ = \sum_{i \neq j \in \mathcal{I}(D_i)} - \left( \mathbf{f}(\mathbf{U}_j^{*,i,n}) - \mathbf{f}(\mathbf{U}_i^n) \right) \cdot \mathbf{c}_{ij} + d_{ij}^m (\mathbf{U}_j^{*,i,n} + \mathbf{U}_i^n) \\ + \sum_{i \neq j \in \mathcal{I}(D_i)} - \left( \mathbf{f}(\mathbf{U}_i^{*,j,n}) - \mathbf{f}(\mathbf{U}_i^n) \right) \cdot \mathbf{c}_{ij} + d_{ij}^m (\mathbf{U}_i^{*,j,n} + \mathbf{U}_i^n) \\ + \sum_{i \neq j \in \mathcal{I}(D_i)} \left( 0, g \left( (\mathbf{H}_i^n)^2 - (\mathbf{H}_i^{*,j,n})^2 \right) \mathbf{c}_{ij} \right)^\top - (d_{ij}^m + d_{ij}^m) (\mathbf{U}_i^{*,j,n} + \mathbf{U}_i^n). \end{aligned}$$

Upon introducing the vectors  $\overline{\mathbf{U}}_{ij}^n \in \mathbb{R}^{1+d}$ ,  $\overline{\mathbf{W}}_{ij}^n \in \mathbb{R}^{1+d}$ , and  $\Delta \mathbf{Z}_i^n \in \mathbb{R}^d$  defined by

$$\begin{aligned} \overline{\mathbf{U}}_{ij}^n &:= -\frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2d_{ij}^n} \left( \mathbf{f}(\mathbf{U}_j^{*,i,n}) - \mathbf{f}(\mathbf{U}_i^n) \right) \cdot \mathbf{n}_{ij} + \frac{1}{2} (\mathbf{U}_j^{*,i,n} + \mathbf{U}_i^n), \\ \overline{\mathbf{W}}_{ij}^n &:= -\frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2d_{ij}^n} \left( \mathbf{f}(\mathbf{U}_i^{*,j,n}) - \mathbf{f}(\mathbf{U}_i^n) \right) \cdot \mathbf{n}_{ij} + \frac{1}{2} (\mathbf{U}_i^{*,j,n} + \mathbf{U}_i^n), \\ \Delta \mathbf{Z}_i^n &:= \sum_{i \neq j \in \mathcal{I}(D_i)} g \left( (\mathbf{H}_i^n)^2 - (\mathbf{H}_i^{*,j,n})^2 \right) \mathbf{c}_{ij}, \end{aligned}$$

we finally obtain

$$\begin{aligned} \mathbf{U}_i^{n+1} &= \left(1 - \sum_{i \neq j \in \mathcal{I}(D_i)} \frac{4\tau}{m_i} d_{ij}^n\right) \mathbf{U}_i^n + \sum_{i \neq j \in \mathcal{I}(D_i)} \frac{2\tau}{m_i} d_{ij}^n (\overline{\mathbf{U}}_{ij}^n + \overline{\mathbf{W}}_{ij}^n) \\ &\quad + \frac{\tau}{m_i} (0, \Delta \mathbf{Z}_i^n)^\top + \frac{2\tau}{m_i} \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n \left(1 - \frac{H_i^{*,j,n}}{H_i^n}\right) \mathbf{U}_i^n. \end{aligned}$$

Upon introducing the fake time  $t = \frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2d_{ij}^n}$  and observing that the definition of  $d_{ij}^n$  implies that  $2t\lambda_{\max}^f(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^{*,i,n}) \leq 1$  and  $2t\lambda_{\max}^f(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_i^{*,j,n}) \leq 1$ , we infer from Lemma 3.6 that  $\overline{\mathbf{U}}_{ij}^n \in \text{conv}_{j \in \mathcal{I}(D_i)}(\mathbf{U}_j^{*,i,n})$  and  $\overline{\mathbf{W}}_{ij}^n \in \text{conv}_{j \in \mathcal{I}(D_i)}(\mathbf{U}_i^{*,j,n})$ ; hence,  $\frac{\overline{\mathbf{U}}_{ij}^n + \overline{\mathbf{W}}_{ij}^n}{2} \in \text{conv}_{j \in \mathcal{I}(D_i)}(\mathbf{U}_j^{*,i,n}, \mathbf{U}_i^{*,j,n})$ . In conclusion, under the CFL condition  $1 + 4\frac{\tau}{m_i}d_{ii}^n \geq 0$ , the state  $\tilde{\mathbf{U}}_i^{n+1} := (1 + \frac{4\tau}{m_i}d_{ii}^n)\mathbf{U}_i^n + \sum_{i \neq j \in \mathcal{I}(D_i)} \frac{2\tau}{m_i}d_{ij}^n(\overline{\mathbf{U}}_{ij}^n + \overline{\mathbf{W}}_{ij}^n)$  belongs to  $\text{conv}_{j \in \mathcal{I}(D_i)}(\mathbf{U}_j^{*,i,n}, \mathbf{U}_i^{*,j,n})$ . If the bathymetry map is flat then  $H_i^n = H_i^{*,j,n}$  and we obtain  $\mathbf{U}_i^{n+1} = \tilde{\mathbf{U}}_i^{n+1} \in \text{conv}_{j \in \mathcal{I}(D_i)}(\mathbf{U}_j^n) \subset A_i^n$  and this proves (i). If the bathymetry is not flat, then  $\mathbf{U}_j^{*,i,n}$  is in the convex hull of  $\mathbf{U}_j^n$  and  $\mathbf{0}$  for all  $j \in \mathcal{I}(D_i)$  and  $\mathbf{U}_i^{*,j,n}$  is in the convex hull of  $\mathbf{U}_i^n$  and  $\mathbf{0}$  for all  $j \in \mathcal{I}(D_i)$ ; this proves that  $\tilde{\mathbf{U}}_i^{n+1} \in \text{conv}(A_i^n, \mathbf{0})$ . Hence, if the bathymetry is not flat we get  $\mathbf{U}_i^{n+1} \in \text{conv}(A_i^n, \mathbf{0}) + (0, \Delta \mathbf{Z}_i^n)^\top + \Delta \mathbf{U}_i^{*,n}$  as announced. The water height in  $\Delta \mathbf{U}_i^{*,n}$  is  $\frac{2\tau}{m_i} \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n (H_i^n - H_i^{*,j,n}) \geq 0$ . Since all the states in  $A_i^n$  have nonnegative water height, we conclude that  $H_i^{n+1} \geq 0$  and this proves (ii). Finally, fix  $n \geq 0$  and assume that all states  $\{\mathbf{U}_j^n\}$  are admissible in the sense of Definition 3.10. If  $H_i^n > 0$  then we have that

$$H_i^{n+1} \geq \left(1 - \sum_{i \neq j \in \mathcal{I}(D_i)} \frac{4\tau}{m_i} d_{ij}^n\right) H_i^n > 0,$$

and this proves that  $\mathbf{U}_i^{n+1}$  is admissible. In the remaining case  $H_i^n = 0$ , we have that  $H_i^{*,j,n} = 0$  for all  $j \in \mathcal{I}(D_i)$  and  $\Delta \mathbf{Z}_i^n = 0$ . Hence  $\mathbf{U}_j^{n+1} = \tilde{\mathbf{U}}_i^{n+1}$  and using that  $\tilde{\mathbf{U}}_i^{n+1}$  is a convex combination of admissible states we conclude that the state  $\mathbf{U}_i^{n+1}$  is admissible and this proves (iii).  $\square$

We finish with a discrete inequality which reduces to a standard discrete entropy inequality when the bottom topography is flat. The proof is omitted for brevity.

**PROPOSITION 3.12.** *Let  $\mathbf{u}_h^{n+1}$  be given by (3.3)–(3.5). Assume the CFL condition  $1 + 4\frac{\tau}{m_i}d_{ii}^n \geq 0$ . Then for any flat bed shallow water entropy pair  $(\eta, \mathbf{G})$ , we have the following discrete entropy inequality:*

$$\begin{aligned} (3.11) \quad &\frac{m_i}{\tau} (\eta(\mathbf{U}_i^{n+1}) - \eta(\mathbf{U}_i^n)) + \sum_{i \neq j \in \mathcal{I}(D_i)} \left(\mathbf{G}(\mathbf{U}_j^{*,i,n}) + \mathbf{G}(\mathbf{U}_i^{*,j,n})\right) \cdot \mathbf{c}_{ij} \\ &\leq \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n \left(\eta(\mathbf{U}_j^{*,i,n}) + \eta(\mathbf{U}_i^{*,j,n}) - 2\eta(\mathbf{U}_i^n)\right) \\ &\quad + \left((0, \Delta \mathbf{Z}_i^n)^\top + \sum_{i \neq j \in \mathcal{I}(D_i)} 2d_{ij}^n \left(1 - \frac{H_i^{*,j,n}}{H_i^n}\right) \mathbf{U}_i^n\right) \cdot \nabla \eta(\mathbf{U}_i^{n+1}). \end{aligned}$$

*Remark 3.13* (literature). We refer the reader to Bouchut and Frid [8, section 2] for an alternative point of view to derive the invariant domain property and entropy inequality obtained above.

**4. Second-order extension.** In this section we propose a scheme that is second-order accurate in space, is exactly well-balanced, and is positivity preserving.

**4.1. Flux approximation.** We start by constructing a well-balanced second-order approximation of the quantity  $\int_{D_i} (\nabla(\frac{1}{2}h^2) + h\nabla z)\varphi_i \, dx$ .

LEMMA 4.1 (consistency/well-balancing). (i) Assume that  $\{\widehat{\theta}_n\}_{n \in \{1:n_{\text{sh}}\}}$  consists of Lagrange or Bernstein basis functions. The expression  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_i(\mathbf{H}_j + \mathbf{Z}_j) \mathbf{c}_{ij}$  is a second-order approximation of  $\int_{D_i} (\nabla(\frac{1}{2}h^2) + h\nabla z)\varphi_i \, dx$ . (ii) The mapping  $\mathbf{u}_h \rightarrow (0, \sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_i(\mathbf{H}_j + \mathbf{Z}_j) \mathbf{c}_{ij})_{i \in \{1:I\}}$  is an exactly well-balanced flux.

*Proof.* (i) If  $h + z$  is linear over  $K \in \mathcal{T}_h$  then  $\int_K h\nabla(h + z)\varphi_i \, dx = \nabla(h + z)|_K \int_K h\varphi_i \, dx$  and the approximation  $\int_K h\varphi_i \, dx \approx \mathbf{H}_i \frac{1}{d} |K|$  is second-order accurate, at least for Lagrange and Bernstein basis functions. Hence, upon noticing that  $\sum_{K \subset D_i} \nabla(h + z)|_K \frac{1}{d} |K| = \int_{D_i} \nabla(h + z)\varphi_i \, dx = \sum_{j \in \mathcal{I}(D_i)} (\mathbf{H}_j + \mathbf{Z}_j) \mathbf{c}_{ij}$ , the expression  $\int_{D_i} h\nabla(h + z)\varphi_i \, dx \approx \sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_i(\mathbf{H}_j + \mathbf{Z}_j) \mathbf{c}_{ij}$  is formally second-order accurate.

(ii) Let us now prove well-balancing. Let us assume exact rest. Let us fix  $i \in \{1:I\}$ . Notice that owing to the partition of unity property we have  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{c}_{ij} = 0$ ; hence  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_i(\mathbf{H}_j + \mathbf{Z}_j) \mathbf{c}_{ij} = \sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_i(\mathbf{H}_j + \mathbf{Z}_j - \mathbf{H}_i - \mathbf{Z}_i) \mathbf{c}_{ij}$ . Consider  $j \in \mathcal{I}(D_i)$ . According to our definition of the exact rest state (see Definition 2.4), either  $\mathbf{H}_i = 0$  and  $\mathbf{H}_j = 0$ , or  $\mathbf{H}_j + \mathbf{Z}_j - \mathbf{H}_i - \mathbf{Z}_i = 0$ ; whence the conclusion.  $\square$

Let us introduce the gas dynamics flux  $\mathbf{g}(\mathbf{u}) := (\mathbf{q}, \frac{1}{h} \mathbf{q} \otimes \mathbf{q})^\top$ ; then upon invoking Lemma 2.1,  $\sum_{j \in \mathcal{I}(D_i)} \mathbf{g}(\mathbf{U}_j) \cdot \mathbf{c}_{ij}$  is a second-order approximation of  $\int_{D_i} \nabla \cdot (\mathbf{g}(\mathbf{u}))\varphi_i \, dx$ .

**4.2. Full time and space approximation.** Let  $\mathbf{u}_h^0 = \sum_{i=1}^I \mathbf{U}_i^0 \varphi_i \in \mathcal{P}(\mathcal{T}_h)$  be a reasonable approximation of  $\mathbf{u}_0$ . Let  $n \in \mathbb{N}$ ,  $\tau$  be the time step,  $t_n$  be the current time, and  $t_{n+1} := t_n + \tau$ . Let  $\mathbf{u}_h^n = \sum_{i=1}^I \mathbf{U}_i^n \varphi_i \in \mathcal{P}(\mathcal{T}_h)$  be the space approximation of  $\mathbf{u}$  at time  $t_n$  and let  $\mathbf{u}_h^{n+1} := \sum_{i=1}^I \mathbf{U}_i^{n+1} \varphi_i$ . We estimate  $\mathbf{U}_i^{n+1}$  as follows:

$$(4.1) \quad \frac{m_i}{\tau} (\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) = \sum_{j \in \mathcal{I}(D_i)} -\mathbf{g}(\mathbf{U}_j^n) \cdot \mathbf{c}_{ij} - (0, g\mathbf{H}_i^n (\mathbf{H}_j^n + \mathbf{Z}_j) \mathbf{c}_{ij})^\top + \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) + \mu_{ij}^n (\mathbf{U}_j^n - \mathbf{U}_j^{*,i,n} - (\mathbf{U}_i^n - \mathbf{U}_i^{*,j,n}))$$

$$(4.2) \quad \mu_{ij}^n := \max((\mathbf{V}_i \cdot \mathbf{n}_{ij})_-, (\mathbf{V}_j \cdot \mathbf{n}_{ij})_+) \|\mathbf{c}_{ij}\|_{\ell^2}, \quad d_{ij}^n \geq \mu_{ij}^n, \quad i \neq j.$$

Here we use the notation  $a_+ := \max(a, 0)$  and  $a_- = -\min(a, 0)$ . In the above scheme  $d_{ij}^n = d_{ji}^n$  can be any nonnegative number larger than  $\mu_{ij}^n$  when  $i \neq j$ . One could just take  $d_{ij}^n = \mu_{ij}^n$ , but a more robust choice consists of using  $d_{ij}^n = \max(d_{ij}^{f,n}, d_{ji}^{f,n})$ ; note that in this case the local maximum wave speed formulas (3.9) and (3.10) used with  $\mathbf{u}_L := (\mathbf{H}_i^n, \mathbf{Q}_i^n \cdot \mathbf{n}_{ij})$  and  $\mathbf{u}_R := (\mathbf{H}_j^n, \mathbf{Q}_j^n \cdot \mathbf{n}_{ij})$  imply that  $d_{ij}^n \geq \mu_{ij}^n$ . Notice that  $\mu_{ij}^n = \mu_{ji}^n$  because  $\mathbf{n}_{ij} = -\mathbf{n}_{ji}$  owing to the assumed boundary condition. We adopt again the convention  $d_{ii}^n := -\sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n$ .

PROPOSITION 4.2. The scheme (4.1)–(4.2) is exactly well-balanced and conservative. It is positivity preserving provided  $1 + 2d_{ii}^n \frac{\tau}{m_i} \geq 0$  for all  $i \in \{1:I\}$ .

*Proof.* The artificial viscosity term on the right-hand side of (4.1) at exact rest is  $\sum_{i \neq j \in \mathcal{I}(D_i)} -\mu_{ij}^n (-\mathbf{H}_j^n + \mathbf{H}_i^n, \mathbf{0})^\top = 0$ , since  $\mu_{ij}^n = 0$  at rest state (at large). The remainder of the proof is a consequence of Lemma 4.1, which establishes exact well-balancing. Since  $\sum_{j \in \mathcal{I}(D_i)} -\mathbf{g}(\mathbf{U}_j^n) \cdot \mathbf{c}_{ij} = \sum_{j \in \mathcal{I}(D_i)} (\mathbf{g}(\mathbf{U}_i^n) - \mathbf{g}(\mathbf{U}_j^n)) \cdot \mathbf{c}_{ij}$ , the conservation can be shown like in the proof of Proposition 3.9. Finally, to prove positivity, let us fix  $i$  and assume that  $\mathbf{H}_j^n \geq 0$  for all  $j \in \mathcal{I}(D_i)$ . The water height update is

$$\begin{aligned} \mathbf{H}_i^{n+1} &= \mathbf{H}_i^n - \frac{\tau}{m_i} \sum_{i \neq j} \left( \mu_{ij}^n \mathbf{H}_i^n + (d_{ij}^n - \mu_{ij}^n) \mathbf{H}_i^{*,j,n} \right) \\ &\quad + \frac{\tau}{m_i} \sum_{i \neq j} \left( (\mu_{ij}^n - \mathbf{c}_{ij} \cdot \mathbf{V}_j^n) \mathbf{H}_j^n + (d_{ij}^n - \mu_{ij}^n) \mathbf{H}_j^{*,i,n} \right). \end{aligned}$$

Using that  $d_{ij}^n - \mu_{ij}^n \geq 0$ ,  $\mu_{ij}^n \geq 0$ ,  $\mathbf{H}_i^n \geq \mathbf{H}_i^{*,j,n} \geq 0$ , and  $\mathbf{H}_j^{*,i,n} \geq 0$  we obtain

$$\mathbf{H}_i^{n+1} \geq \mathbf{H}_i^n \left( 1 - \frac{\tau}{m_i} \sum_{i \neq j} d_{ij}^n \right) + \frac{\tau}{m_i} \sum_{i \neq j} (\mu_{ij}^n - \mathbf{c}_{ij} \cdot \mathbf{V}_j^n) \mathbf{H}_j^n.$$

The conclusion follows from the assumption on the CFL number and the definition of  $\mu_{ij}^n$  which implies that  $\mu_{ij}^n - \mathbf{c}_{ij} \cdot \mathbf{V}_j^n \geq ((\mathbf{V}_j^n \cdot \mathbf{n}_{ij})_+ - \mathbf{V}_j^n \cdot \mathbf{n}_{ij}) \|\mathbf{c}_{ij}\|_{\ell^2} \geq 0$ .  $\square$

*Remark 4.3.* Note that the approximation of the flux in the scheme (4.1) is formally second-order accurate in space and contrary to (3.3) does not suffer from the small inconsistency of the hydrostatic reconstruction, since the hydrostatic reconstruction is used only in the artificial viscosity. In particular (4.1) is formally second-order accurate in space when the artificial viscosity is set to zero.

**4.3. Second-order positivity preserving viscosity.** In order to make the proposed method fully second-order accurate in space, we now propose a new definition of the viscosity along the line of Guermond and Popov [17]. Namely, we choose the viscous terms  $d_{ij}^n$  and  $\mu_{ij}^n$  in the scheme (4.1) to be  $d_{ij}^n := \alpha_{ij}^n d_{ij}^{v,n}$  and  $\mu_{ij}^n := \alpha_{ij}^n \mu_{ij}^{v,n}$ , where  $d_{ij}^{v,n} := \max(d_{ij}^{f,n}, d_{ij}^{f_i,n})$  is the first-order viscosity based on the maximum wave speed,  $\mu_{ij}^{v,n} := \max((\mathbf{V}_i \cdot \mathbf{n}_{ij})_-, (\mathbf{V}_j \cdot \mathbf{n}_{ij})_+) \|\mathbf{c}_{ij}\|_{\ell^2}$  and  $\alpha_{ij}^n \in [0, 1]$  is appropriately chosen. More precisely, the proposed second-order scheme is

$$\begin{aligned} \frac{m_i}{\tau} (\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) &= \sum_{j \in \mathcal{I}(D_i)} -\mathbf{g}(\mathbf{U}_j^n) \cdot \mathbf{c}_{ij} - (0, g\mathbf{H}_i^n (\mathbf{H}_j^n + Z_j) \mathbf{c}_{ij})^\top \\ (4.3) \quad &+ \sum_{i \neq j \in \mathcal{I}(D_i)} d_{ij}^n (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) + \mu_{ij}^n (\mathbf{U}_j^n - \mathbf{U}_j^{*,i,n} - (\mathbf{U}_i^n - \mathbf{U}_i^{*,j,n})), \end{aligned}$$

$$(4.4) \quad \mu_{ij}^n := \max(\psi_i^n, \psi_j^n) \mu_{ij}^{v,n}, \quad i \neq j,$$

$$(4.5) \quad d_{ij}^n := \max(\psi_i^n, \psi_j^n) d_{ij}^{v,n}, \quad i \neq j,$$

with  $\psi_i^n \in [0, 1]$  yet to be determined. One possible choice for the second-order coefficient  $\psi_i^n$  consists of setting  $\psi_i^n = \psi(\alpha_i^n)$ , where we define

$$(4.6) \quad \alpha_i^n := \frac{|\sum_{j \in \mathcal{I}(D_i)} \mathbf{H}_j^n - \mathbf{H}_i^n|}{\sum_{j \in \mathcal{I}(D_i)} |\mathbf{H}_j^n - \mathbf{H}_i^n|}.$$

It is shown in Guermond and Popov [19] that any function  $\psi$  in  $C^{0,1}([0, 1]; [0, 1])$  with  $\psi(1) = 1$  gives an algorithm that is positivity preserving up to a CFL condition, (see also [17] for the scalar version of the method and other possible choices for  $\psi_i^n$ ). We take  $\psi(\alpha) = \alpha^2$  in all the numerical simulations reported at the end of the paper.

PROPOSITION 4.4. *Let  $k_\psi$  be the Lipschitz constant of  $\psi$ . The scheme (4.3)–(4.4)–(4.5) is positivity preserving provided that  $\frac{\tau}{m_i}(-d_{ii}^n + \sum_{j \in \mathcal{I}(D_i)} (\mathbf{c}_{ij} \cdot \mathbf{V}_j^n)_-) \leq \frac{1}{2}$  and  $\frac{\tau}{m_i} \max_{i \neq j \in \mathcal{I}(D_i)} (\mathbf{c}_{ij} \cdot \mathbf{V}_j)_- \leq \frac{1}{4k_\psi c_\#^2}$ , where  $c_\# = \max_{i \in \{1: I\}} \text{card}(\mathcal{I}(D_i))$ .*

*Proof.* By proceeding as in the proof of Proposition 4.2, we obtain

$$\begin{aligned} H_i^{n+1} &= H_i^n - \frac{\tau}{m_i} \sum_{i \neq j} \left( \mu_{ij}^n H_i^n + (d_{ij}^n - \mu_{ij}^n) H_i^{*,j,n} \right) \\ &\quad + \frac{\tau}{m_i} \sum_{i \neq j} \left( (\mu_{ij}^n - \mathbf{c}_{ij} \cdot \mathbf{V}_j^n) H_j^n + (d_{ij}^n - \mu_{ij}^n) H_j^{*,i,n} \right). \end{aligned}$$

Using that  $d_{ij}^n \geq \mu_{ij}^n$  and  $H_j^{*,i,n} \geq 0$ ,  $H_i^n \geq H_i^{*,j,n}$  for all  $j$ , we obtain

$$H_i^{n+1} \geq H_i^n \left( 1 - \frac{\tau}{m_i} \sum_{i \neq j} d_{ij}^n \right) + \frac{\tau}{m_i} \sum_{i \neq j} (\mu_{ij}^n - \mathbf{c}_{ij} \cdot \mathbf{V}_j^n) H_j^n.$$

To finish the proof, it remains to show that the right-hand side is nonnegative under the appropriate CFL condition. The reader is referred to [19] for the proof of this result and other choices for  $\alpha_{ij}^n$  that also make the scheme (4.3) positivity preserving.  $\square$

*Remark 4.5* (linearity preserving). It is possible to modify the definition of  $\alpha_i^n$  in (4.6) to make the method linearity preserving (the reader is referred to Berger, Aftosis, and Murman [3] for a review on linearity-preserving limiters in the finite volume literature). More precisely, when the shape functions are Lagrange based, one can set  $\alpha_i^n := |\sum_{j \in \mathcal{I}(D_i)} \beta_{ij} (H_j^n - H_i^n)| / \sum_{j \in \mathcal{I}(D_i)} \beta_{ij} |H_j^n - H_i^n|$ , where the coefficients  $\beta_{ij}$  are generalized barycentric coordinates; see Guermond and Popov [17] for details. We take  $\beta_{ij} = 1$  in all the numerical simulations reported at the end of the paper.

**5. Numerical illustrations.** In this section we illustrate the performance of the various algorithms introduced in the paper. Most of the test cases are taken from the so-called SWASHES suite from Delestre et al. [11].

**5.1. Technical details.** All the numerical simulations are done in two space dimensions even when the problem under consideration has a one-dimensional solution. In order to avoid extraneous superconvergence effects we use unstructured, nonnested, Delaunay meshes composed of triangles. The computations are done with continuous Lagrange  $\mathbb{P}_1$  finite elements. The time stepping is done with the SSP RK(3,3) method (three stages, third-order), see Shu and Osher [30, (2.18)] and Kraaijevanger [22, Thm. 9.4]. All the computations reported in this section have been done with the upper bound on  $\lambda_{\max}^{\mathbf{f}_1^D}(\mathbf{v}_L, \mathbf{v}_R)$  given by Lemma 3.8.

To avoid division by zero in the presence of dry states we introduce  $h_\epsilon := \epsilon \max_{\mathbf{x} \in D} h_0(\mathbf{x})$  with  $\epsilon = 10^{-16}$ , where  $h_0$  is the initial water height. That is to say, we approximate the 0 water height by  $10^{-16}$  times the maximum water height at the initial time. Then we regularize the gas dynamics flux  $\mathbf{g}$  as follows:  $\mathbf{g}_\epsilon(\mathbf{u}) := (\mathbf{q}, \frac{2h}{h^2 + \max(h, h_\epsilon)^2} \mathbf{q} \otimes \mathbf{q})^\top$ . That is to say the speed  $\mathbf{v} := \mathbf{g}/h$  is regularized by setting  $\mathbf{v}_\epsilon := \frac{2h}{h^2 + \max(h, h_\epsilon)^2} \mathbf{q}$ . Note that we obtain  $\mathbf{g}(\mathbf{u}) = \mathbf{g}_\epsilon(\mathbf{u})$  and  $\mathbf{v}_\epsilon = \mathbf{v}$  when  $h \geq h_\epsilon$ ; that is, the regularization is active only when  $h \leq h_\epsilon$ .

All the schemes proposed in this paper are positivity preserving on the water height provided they are programmed correctly. Hence, provided the initial water is nonnegative, the water height should never become negative up to roundoff errors.

We have observed that it is possible to avoid the effects of roundoff errors in the presence of dry regions by programming the update of the water height as follows:

$$(5.1) \quad \mathbf{H}_i^{n+1} = \mathbf{H}_i^n \left( 1 - \frac{\tau}{m_i} \left( \mathbf{c}_{ii} \cdot \mathbf{V}_i^n + \sum_{i \neq j} \mu_{ij}^n + (d_{ij}^n - \mu_{ij}^n) \frac{\mathbf{H}_i^{*,j,n}}{\mathbf{H}_i^n} \right) \right) + \frac{\tau}{m_i} \sum_{i \neq j} \left( -\mathbf{c}_{ij} \cdot \mathbf{Q}_j^n + \mu_{ij}^n \mathbf{H}_j^n + (d_{ij}^n - \mu_{ij}^n) \mathbf{H}_j^{*,i,n} \right)$$

instead of setting  $H_i^{n+1} = H_i^n + \frac{\tau}{m_i} \Delta R_i^n$  with

$$\Delta R_i^n := \sum_{j \in \mathcal{I}(D_i)} -\mathbf{c}_{ij} \cdot \mathbf{Q}_j^n + \mu_{ij}^n (\mathbf{H}_j^n - \mathbf{H}_i^n) + (d_{ij}^n - \mu_{ij}^n) (\mathbf{H}_j^{*,i,n} - \mathbf{H}_i^{*,j,n}).$$

When doing convergence tests over meshes of different mesh size, the convergence rates are estimated as follows: given two errors  $e_1, e_2$  obtained on two meshes  $\mathcal{T}_{h1}, \mathcal{T}_{h2}$ , and denoting  $I_1 := \dim P(\mathcal{T}_{h1})$   $I_2 := \dim P(\mathcal{T}_{h2})$ , the convergence rate is defined to be the ratio  $d \log(e_1/e_2) / \log(I_2/I_1)$  since the quantity  $I^{-\frac{1}{d}}$  scales like the mesh size. In all the test cases we take  $g = 9.81 \text{ m s}^{-1}$  and  $d = 2$ .

**5.2. Well-balancing.** We have verified on various tests, not reported here for brevity, that the proposed methods are well-balanced. More precisely, the first-order algorithm (3.3)–(3.5) is well-balanced irrespective of the structure of the mesh, i.e., the discharge stays close to the roundoff error indefinitely. The well-balancing of the second-order algorithm depends on whether or not exact rest is possible as defined in Definition 2.4. If the mesh is such that exact rest is possible, then the algorithm is well-balanced up to machine accuracy indefinitely. If exact rest is not supported by the mesh, approximate well-balancing is achieved up to truncation error indefinitely.

**5.3. Flows over a bump.** We consider in this section several classical test cases detailed in [11, section 3.1]. The domain is a one-dimensional channel  $[0, L]$  with length  $L = 25 \text{ m}$ . The bathymetry profile proposed in [11, section 3.1] is flat with a parabolic bump, but to increase the smoothness of the solution in order to estimate the convergence rate properly, we modify a little bit the profile as follows:

$$(5.2) \quad z(x) = \begin{cases} \frac{0.2}{64}(x - 8)^3(12 - x)^3 & \text{if } 8 \leq x \leq 12, \\ 0 & \text{otherwise.} \end{cases}$$

Steady solutions satisfy mass conservation  $q(x) = q(0)$  and the Bernoulli relation

$$(5.3) \quad \frac{q^2}{2gh^2} + h(x) + z(x) = C_{\text{Ber}},$$

where the Bernoulli constant  $C_{\text{Ber}}$  depends on the data. All the computations in section 5.3 are done in two dimensions in the channel  $D = [0, L] \times [0, 1]$ .

**5.3.1. Subcritical flow.** We now consider a steady state solution with the inflow discharge  $-\mathbf{q} \cdot \mathbf{n} = q_{\text{in}} = 4.42 \text{ m}^2 \text{ s}^{-1}$  imposed at  $\{x = 0\}$  and  $\mathbf{q} \cdot \mathbf{n} = 0$  on the sides of the channel  $\{y = 0\} \cup \{y = 1\}$ . The water height is enforced to be equal to  $h_L = 2 \text{ m}$  at  $\{x = L\}$ ; hence  $C_{\text{Ber}} := \frac{q_{\text{in}}^2}{2gh_L^2} + h_L$ . The initial condition is  $\mathbf{q}_0(\mathbf{x}) = 0$  and  $h_0(\mathbf{x}) = h_L - z(x)$ . We look for the solution at  $t = 80 \text{ s}$  which should be close to steady state. From Bernoulli’s relation (5.3),  $z(x) + h(x) + \frac{q_{\text{in}}^2}{2gh^2(x)} = C_{\text{Ber}}$  one gets



TABLE 1

Subcritical flow over a bump with  $h$  given by (5.6). Computation done at  $t = 80$ s with initial data at rest;  $CFL = 1.25$ .  $L^1$ -norm is given in rows 2–6,  $L^2$ -norm is given in rows 7–11. Viscosities are  $\psi(\alpha) = \alpha^2$  (columns 3–4) and first-order viscosity (columns 5–6).

Norm	$I$	$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$	
$L^1$	248	1.46E-03	Rate	4.99E-03	Rate
	885	2.57E-04	2.73	3.39E-03	0.61
	3069	3.44E-05	3.08	1.95E-03	0.84
	12189	1.21E-06	3.09	1.03E-03	0.98
	48053	7.47E-07	2.66	5.19E-04	1.00
$L^2$	248	2.91E-3	Rate	9.57E-03	Rate
	885	6.48E-04	2.35	6.36E-03	0.64
	3069	1.25E-04	2.52	3.62E-03	0.86
	12189	2.31E-05	2.59	1.90E-03	0.99
	48053	4.04E-06	2.55	9.57E-04	1.00

that the exact steady state solution  $h(x)$  solves the algebraic equation

$$(5.4) \quad h^3(x) + (z(x) - C_{\text{Ber}}) h^2(x) + \frac{q_{\text{in}}^2}{2g} = 0 \quad \forall x \in [0, L].$$

Let  $b(x) := z(x) - C_{\text{Ber}}$  and  $d := \frac{q_{\text{in}}^2}{2g}$ . With the considered data, the cubic equation  $h^3 + bh^2 + d = 0$  has three real zeros. The one that corresponds to the steady state solution is the largest root. Upon defining

$$(5.5) \quad Q(x) := -\frac{b^2(x)}{9}, \quad R(x) := -\frac{27d + 2b^3(x)}{54}, \quad \cos(\theta(x)) = (-Q(x))^{-\frac{3}{2}} R(x),$$

the water height is given by the trigonometric form of Cardano’s formula:

$$(5.6) \quad h(x) = 2\sqrt{-Q(x)} \cos\left(\frac{\theta(x)}{3}\right) - \frac{b(x)}{3}.$$

Two types of computations are performed with the scheme (4.3)–(4.5) using either the second-order viscosity  $\psi(\alpha) = \alpha^2$  or the first-order viscosity  $\psi(\alpha) = 1$ . We use  $CFL = 1.25$ . In order to speed up the convergence to steady state we additionally impose the exact water height at  $x = 0$ . This artifact is used only to observe the theoretical convergence rate in space at  $t = 80$ . We show in Table 1 the error on the water height measured in the  $L^1$ -norm and in the  $L^2$ -norm. All the errors are relative to the corresponding norm of the exact solution. We observe that the convergence rate exceeds 2 both in the  $L^1$ -norm and in the  $L^2$ -norm for the viscosity  $\psi(\alpha) = \alpha^2$ . This is a superconvergence effect that we do not really understand at the moment. Let us recall that the meshes that are used here are nonnested, unstructured, and the initial condition is rest. As expected the asymptotic convergence rate of the solution obtained with the first-order viscosity  $\psi(\alpha) = 1$  is 1 irrespective of the norm.

**5.3.2. Transcritical flow.** We run again the above test in the transcritical regime. Given  $q_{\text{in}}$ , we set the Bernoulli constant  $C_{\text{Ber}}$  so that the Bernoulli relation (5.4) has two identical positive roots at the top of the bump, meaning that the discriminant of (5.4),  $Q^3 + R^2$ , is zero, where  $Q$  and  $R$  are defined in (5.5). This fixes the Bernoulli constant  $C_{\text{Ber}}$  to be equal to  $z_M + \frac{3}{2}\left(\frac{q_{\text{in}}^2}{g}\right)^{\frac{1}{3}}$ , where  $z_M$  is the height of the bump. The flow is fluvial (subsonic) upstream and becomes torrential (supersonic) at the top of the bump. The exact water height is the largest root of (5.4) when  $x \leq x_M$

TABLE 2

Transcritical flow over a bump with  $h$  given by (5.7). Computation done at  $t = 80$  s with initial data at rest;  $CFL = 0.95$ .  $L^1$ -norm is given in rows 2–6,  $L^2$ -norm is given in rows 7–11. Viscosities are  $\psi(\alpha) = \alpha^2$  (columns 3–4) and first-order viscosity (columns 5–6).

Norm	$I$	$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$	
$L^1$	248	2.03E-02	Rate	1.63E-01	Rate
	885	3.49E-03	2.77	9.09E-02	0.92
	3069	4.71E-04	3.08	4.67E-02	1.02
	12189	9.86E-05	2.40	2.35E-02	1.05
	48053	1.95E-05	2.38	1.17E-02	1.02
$L^2$	248	2.28E-02	Rate	1.57E-01	Rate
	885	4.41E-03	2.58	8.73E-02	0.93
	3069	6.40E-04	2.96	4.49E-02	1.02
	12189	1.30E-04	2.44	2.27E-02	1.05
	48053	2.49E-05	2.42	1.13E-02	1.02

and is the other positive root of (5.4) in the other case:

$$(5.7) \quad h(x) = \begin{cases} 2\sqrt{-Q(x)} \cos\left(\frac{\theta(x)}{3}\right) - \frac{b(x)}{3} & \text{if } x \leq x_M, \\ 2\sqrt{-Q(x)} \cos\left(\frac{4\theta(x)}{3}\right) - \frac{b(x)}{3} & \text{otherwise,} \end{cases}$$

where  $\theta(x)$  is defined in (5.5) and  $x_M$  is such that  $z(x_M)$  is the maximum of  $z(x)$ .

We take  $q_{in} = 1.53 \text{ m}^2 \text{ s}^{-1}$ . With the bottom topography defined in (5.2), we have  $x_M = 10$  m and  $z_M = 0.2$  m. The flow rate is enforced at  $\{x = 0\}$  and the exact water height (given by (5.7)) is enforced at the outflow  $\{x = L\}$ . We start with the initial condition  $q(x) = 0 \text{ m}^2 \text{ s}^{-1}$  and  $h(x) + z(x) = 0.66$  m. The errors are measured at  $t = 80$  s. All the errors are relative to the corresponding norm of the exact solution. The computational domain is again  $D = [0, 25] \times [0, 1]$ . Two types of computations are done with the scheme (4.3)–(4.5) using either the second-order viscosity  $\psi(\alpha) = \alpha^2$  or the first-order viscosity  $\psi(\alpha) = 1$ . We use  $CFL = 0.95$ . We show in Table 2 the error on the water height measured in the  $L^1$ -norm and in the  $L^2$ -norm.

**5.3.3. Transcritical flow over a bump with shock.** We run again the above test in the transcritical regime with a hydraulic jump (i.e., a shock). To get a shock the flow must at some point become sonic and the water height at the outflow boundary must be larger than the water height at the sonic point. At the sonic point the discriminant of the Bernoulli relation (5.4) is zero. Just like in the test in section 5.3.2 we position the sonic point at the top of the bump, i.e., the Bernoulli constant  $C_{Ber}$  is equal to  $z_M + \frac{3}{2} \left(\frac{q_{in}^2}{g}\right)^{\frac{1}{3}}$ , where  $z_M$  is the height of the bump. The flow is fluvial (subsonic) upstream and becomes torrential (supersonic) at the top of the bump and stays supersonic up to the hydraulic jump. Now we fix the location of the shock  $x_S \in (x_M, 12)$ . The water height before the hydraulic jump is the second largest root of (5.4):  $h(x_S^-) = 2\sqrt{-Q(x_S^-)} \cos\left(\frac{4\theta(x_S^-)}{3}\right) - \frac{b(x_S^-)}{3}$ . The water height after the jump is determined by the Rankine–Hugoniot relation:  $h(x_S^+) = 0.5(-h(x_S^-) + \sqrt{\Delta})$ , where  $\Delta = (h(x_S^-))^2 + \frac{8q_{in}^2}{gh(x_S^-)}$ . In conclusion the exact solution for the water height is

$$(5.8) \quad h(x) = \begin{cases} 2\sqrt{-Q(x)} \cos\left(\frac{\theta(x)}{3}\right) - \frac{b(x)}{3} & \text{if } x \leq x_M, \\ 2\sqrt{-Q(x)} \cos\left(\frac{4\theta(x)}{3}\right) - \frac{b(x)}{3} & \text{if } x_M \leq x < x_S, \\ h(x_S^+) + z(x_S) - z(x) & \text{if } x_S < x. \end{cases}$$

TABLE 3

*Transcritical flow with a shock, (5.8). Computation done at  $t = 80$  s with initial data at rest; CFL = 0.95.  $L^1$ -norm is given in rows 2–6,  $L^2$ -norm is given in rows 7–11. Viscosities are  $\psi(\alpha) = \alpha^2$  (columns 3–4) and first-order viscosity  $\psi(\alpha) = 1$  (columns 5–6).*

Norm	$I$	$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$	
$L^1$	248	2.79E-02	Rate	7.40E-02	Rate
	885	7.97E-03	1.97	4.43E-02	0.81
	3069	4.03E-03	1.05	2.71E-02	0.75
	12189	2.69E-03	0.62	1.74E-02	0.68
	48053	1.54E-03	0.82	1.15E-02	0.61
$L^2$	248	6.70E-02	Rate	1.12E-01	Rate
	885	4.81E-02	0.52	8.60E-02	0.42
	3069	3.75E-02	0.38	7.71E-02	0.17
	12189	3.37E-02	0.17	7.19E-02	0.11
	48053	2.55E-02	0.41	6.54E-02	0.14

The bottom topography defined in (5.2) gives  $x_M = 10$  m,  $z_M = 0.2$  m. In our computations we take  $q_{\text{in}} = 0.18 \text{ m}^2 \text{ s}^{-1}$  to be consistent with the literature, Delestre et al. [11], Noelle, Xing, and Shu [26], but we could take any value for  $q_{\text{in}}$ . We use  $x_S = 11.7$  m and compute the water height at the outflow boundary  $h_L := h(x_S^+) + z(x_S) - z(L)$  (using  $g = 9.81 \text{ m s}^{-2}$ , this gives  $h_L = 0.282\,052\,798\,138\,021\,81$  m). Note that in [11, 26] the topography is different ( $z(x) = \max(0, 0.2 - 0.05(x - 10)^2)$ ), the gravity constant is also different ( $g = 9.812 \text{ m s}^{-2}$ ), and the shock location is also different ( $x_S = 11.665\,504\,281\,554\,291$  m). We insist on using our smooth bottom topography (5.2) instead of the parabolic profile, since it allows us to properly estimate the convergence rate of the method. With the nonsmooth topography used in the literature ( $z(x) = \max(0, 0.2 - 0.05(x - 10)^2)$ ), the distance between the shock and the kink in the bottom topography is 0.3 m, which represents 1.2% of the length of the domain. To start observing a meaningful convergence rate with this topography, using a quasi-uniform mesh would require one to have at least 10 grid points between the two singularities, which would require one to have at least 833 grid point in the  $x$ -direction and 33 points in the  $y$ -direction (since  $D = [0, 25] \times [0, 1]$ ). The asymptotic convergence range is reached with far fewer grid points with our smooth topography.

The flow rate is enforced at  $\{x = 0\}$  and the exact water height  $h_L$  is enforced at the outflow  $\{x = L\}$ . The initial condition is  $q(x) = q_{\text{in}}$  and  $h(x) + z(x) = h_L$ . The errors are measured at  $t = 80$  s. Two types of computations are done with the scheme (4.3)–(4.5) using either the second-order viscosity  $\psi(\alpha) = \alpha^2$  or the first-order viscosity  $\psi(\alpha) = 1$ . We use CFL = 0.95. We show in Table 3 the relative error on the water height measured in the  $L^1$ -norm and in the  $L^2$ -norm. Once again the superiority of the second-order viscosity  $\psi(\alpha) = \alpha^2$  is evident.

**5.4. Unsteady flows.** In the preceding sections, we went through steady state solutions of increasing difficulties. These solutions are useful to check well-balancing and accuracy in space, but they do not give information about the transient behavior. Thus, in this section, we test transient solutions with wet/dry transitions.

**5.4.1. Dam break on a dry bottom.** We start with an ideal dam break called Ritter's solution; see [29]. This is a Riemann problem with the initial condition

$$(5.9) \quad h(x) = \begin{cases} h_l & \text{if } 0 \leq x < x_0, \\ 0 & \text{if } x_0 \leq x < L, \end{cases}$$

TABLE 4

Problem (5.9) at  $t = 6$  with data (5.10)–(5.11) at  $t = 1$  (columns 3–6) and  $t = 0$  (columns 7–10); CFL = 0.5.  $L^1$ -norm is given in rows 2–6,  $L^2$ -norm is given in rows 7–11. Viscosities are  $\psi(\alpha) = \alpha^2$  (columns 3–4; 7–8) and first-order viscosity (columns 5–6; 9–10).

Norm	$I$	Initialization time $t = 1$				Initialization time $t = 0$			
		$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$		$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$	
$L^1$	248	1.52E-02	Rate	3.64E-02	Rate	3.33E-02	Rate	4.82E-02	Rate
	816	7.41E-03	1.20	2.17E-02	0.81	1.82E-02	1.01	3.38E-02	0.56
	3069	3.03E-03	1.35	1.22E-02	0.88	1.08E-02	0.79	2.39E-02	0.53
	12189	1.21E-03	1.34	6.70E-03	0.92	4.81E-03	1.16	1.52E-02	0.69
	48053	4.73E-04	1.37	3.54E-03	0.93	2.65E-03	0.87	9.61E-03	0.67
$L^2$	248	2.00E-01	–	4.65E-02	–	4.31E-01	–	6.14E-02	–
	816	1.10E-02	1.01	2.97E-02	0.70	2.45E-02	0.95	4.36E-02	0.54
	3069	5.42E-03	1.06	1.82E-02	0.76	1.40E-02	0.84	3.11E-02	0.52
	12189	2.65E-03	1.04	1.11E-02	0.76	7.13E-03	0.98	2.06E-02	0.63
	48053	1.28E-03	1.06	6.64E-03	0.75	3.83E-03	0.91	1.34E-02	0.63

where  $h_l > 0$  and  $v(x) = 0$  m/s. The analytical solution is

$$(5.10) \quad h(x, t) = \begin{cases} h_l & \text{if } 0 \leq x \leq x_A(t), \\ \frac{4}{9g} (\sqrt{gh_l} - \frac{x-x_0}{2t})^2 & \text{if } x_A(t) \leq x \leq x_B(t), \\ 0 & \text{if } x_B(t) \leq x \leq L, \end{cases}$$

$$(5.11) \quad v(x, t) = \begin{cases} 0 & \text{if } 0 \leq x \leq x_A(t), \\ \frac{2}{3} (\frac{x-x_0}{t} + \sqrt{gh_l}) & \text{if } x_A(t) \leq x \leq x_B(t), \\ 0 & \text{if } x_B(t) \leq x \leq L, \end{cases}$$

where  $x_A(t) = x_0 - t\sqrt{gh_l}$ ,  $x_B(t) = x_0 + 2t\sqrt{gh_l}$ . This test is used to check if the scheme preserves positivity of the water height and is able to locate and treat correctly the wet/dry transition. As in SWASHES [11], we consider  $h_l = 0.005$  m,  $x_0 = 5$  m,  $L = 10$  m, and  $t = 6$  s. The computational domain in  $D = [0, L] \times [0, 1]$ .

We show in Table 4 convergence results on the water height for the solution to the above problem at  $t = 6$  s with two different initializations. The results in columns 3–6 have been obtained with the initial data given by (5.10)–(5.11) with the initial time  $t = 1$  s. This test is meant to estimate the accuracy of the method with a solution whose partial derivatives are in  $BV(D)$ . We observe the rates  $\frac{4}{3}$  in the  $L^1$ -norm and 1 in the  $L^2$ -norm with the viscosity  $\psi(\alpha) = \alpha^2$ . The rates are 1 and  $\frac{3}{4}$  for the first-order viscosity,  $\psi(\alpha) = 1$ . The results on the discharge (not shown) give exactly the same convergence rates. The results in columns 7–10 have been obtained by using the Riemann data (5.9) at  $t = 0$  s. There is a loss of accuracy since the initial data are now only in  $BV(D)$ . We observe the convergence rate 1 in the  $L^1$ -norm and the  $L^2$ -norm for the viscosity  $\psi(\alpha) = \alpha^2$  and  $\frac{2}{3}$  in the  $L^1$ -norm and the  $L^2$ -norm with the first-order viscosity  $\psi(\alpha) = 1$ . The results on the discharge (not shown) give exactly the same convergence rates. Note that with both initializations the  $\psi(\alpha) = \alpha^2$  viscosity performs better than the first-order viscosity  $\psi(\alpha) = 1$ . We have also performed the above tests with the first-order scheme (3.3)–(3.5) and the results (not shown) are almost indistinguishable from those given by the scheme (4.3)–(4.5) with the first-order viscosity  $\psi(\alpha) = 1$ .

**5.5. Planar surface in a paraboloid.** We now consider a two-dimensional solution with moving shoreline developed by Thacker; see [31]. It is periodic in time with moving wet/dry transitions. It provides a perfect test for shallow water codes as it deals with bed slope and wetting/drying with two-dimensional effects. Moreover, as

TABLE 5

Planar free surface in a paraboloid vessel with exact solution (5.12). Computations done at  $t = 3 \times 2\pi/\omega$  with initial data (5.12) at  $t = 0$ ;  $CFL = 0.3$ .  $L^1$ -norm is given in rows 2–6; second-order method with  $\psi(\alpha) = \alpha^2$  (columns 3–4); second-order method with  $\psi(\alpha) = 1$  (columns 5–6); first-order method from section 3 (columns 7–8).

Norm	$I$	$\psi(\alpha) = \alpha^2$		$\psi(\alpha) = 1$		Mthd. from section 3	
$L^1$	508	2.71E-01	Rate	6.25E-01	Rate	7.85E-01	Rate
	1926	6.51E-02	2.13	4.27E-01	0.57	7.44E-01	0.08
	7553	1.58E-02	2.08	2.54E-01	0.76	5.46E-01	0.45
	29870	4.46E-03	1.83	1.49E-01	0.88	3.33E-01	0.72
	118851	1.50E-03	1.58	7.26E-02	0.94	1.82E-01	0.87

the gradient of the solution has BV regularity, it is appropriate to verify the accuracy of a numerical method up to second order in  $L^1(D)$ . The topography is a paraboloid of revolution defined by

$$z(\mathbf{x}) = -h_0 \left( 1 - \left( \frac{r(\mathbf{x})}{a} \right)^2 \right)$$

with  $r(\mathbf{x}) = \sqrt{(x - L/2)^2 + (y - L/2)^2}$  for each  $\mathbf{x} := (x, y) \in [0, L] \times [0, L]$ . When the water is at rest,  $h_0$  is the water height at the central point of the domain and  $a$  is the radius of the circular free surface. An analytical solution with a moving shoreline and a free surface that remains planar in time is given by

$$(5.12) \quad \begin{cases} h(\mathbf{x}, t) = \max\left(\frac{\eta h_0}{a^2} \left(2(x - \frac{L}{2}) \cos(\omega t) + 2(y - \frac{L}{2}) \sin(\omega t)\right) - z(x, y), 0\right), \\ v_x(\mathbf{x}, t) = -\eta \omega \sin(\omega t), \\ v_y(\mathbf{x}, t) = \eta \omega \cos(\omega t), \end{cases}$$

where the frequency is defined by  $\omega = \sqrt{2gh_0}/a$  and  $\eta$  is a free parameter. To visualize this case, one can think of a glass with some liquid in rotation inside.

The initial condition is the analytic solution at  $t = 0$ . Boundary conditions are natural, i.e., nothing is enforced. Typical values of parameters are the same as in SWASH [11],  $a = 1$  m,  $h_0 = 0.1$  m,  $L = 4$  m,  $\eta = 0.5$ . The solution is computed up to time  $t = 3 \times 2\pi/\omega$ . The computational domain is  $D = [0, L] \times [0, L]$ . The results are reported in Table 5.

**5.6. Tidal wave over an island.** We finish with a simulation of an experiment reported in Liu et al. [25], which consists of a water tank  $D = [0, 30] \times [0, 25]$  with a conical island. The topography is

$$(5.13) \quad z(\mathbf{x}) := \min(h_{\text{top}}, (h_{\text{cone}} - r(\mathbf{x})/s_{\text{cone}})_+), \quad r(\mathbf{x}) := \sqrt{(x - 15)^2 + (y - 13)^2},$$

where  $h_{\text{top}} = 0.625$  m,  $h_{\text{cone}} = 0.9$  m,  $s_{\text{cone}} = 4$  m. All the dimensions are in meters. We do not use the experimental setup for the initial conditions since there is no real consensus in the literature on the setup of the initial data. Instead, we set the initial condition to be a (solitary) wave big enough to overtop the island to demonstrate that the method is robust with respect to the presence of dry states. Moreover, we impose transparent boundary conditions to show that they are easy to enforce in the finite element setting. Essentially, imposing transparent boundary conditions consist

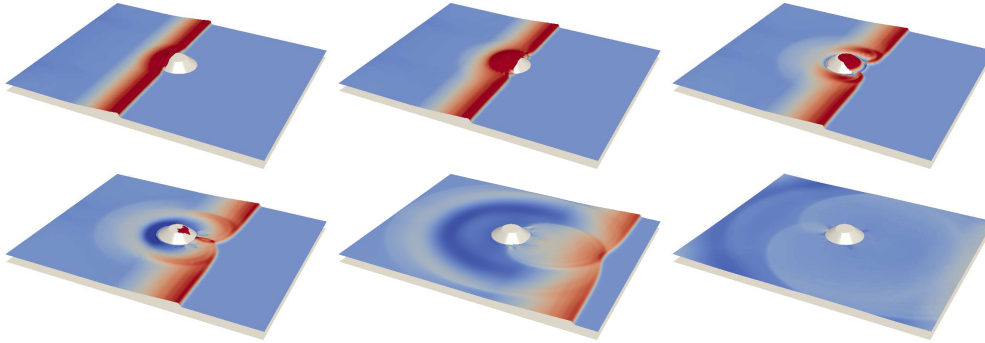


FIG. 2. Tidal wave overtopping a conical island.

of not doing anything (these are the so-called natural boundary conditions). The initial condition is given by  $h(\mathbf{x}, 0) = h_{\text{init}}(\mathbf{x})$ ,  $\mathbf{q}(\mathbf{x}, 0) = (u_{\text{init}}(\mathbf{x})h_{\text{init}}(\mathbf{x}), 0)$ , where

$$(5.14) \quad h_{\text{init}}(\mathbf{x}) := \left( h_0 + \frac{A}{\cosh^2 \left( \sqrt{\frac{3A}{4h_0^3}} (x - x_s) \right)} - z(\mathbf{x}) \right)_+,$$

$$(5.15) \quad u_{\text{init}}(\mathbf{x}) := \frac{A}{\cosh^2 \left( \sqrt{\frac{3A}{4h_0^3}} (x - x_s) \right)} \sqrt{\frac{g}{h_0}}$$

with  $h_0 = 0.32$  m,  $A = h_0$ , and  $x_s = 2.04$  m. The computations are done on an unstructured Delaunay mesh composed of 174432 triangles and 87767 grid points. The average mesh size is 0.1 m. We report in Figure 2 the water elevation at 6 different times, 4.08 s, 4.92 s, 5.88 s, 6.96 s, 9.72 s, 14.52 s showing the various stages of the overtopping of the island. To visualize properly the dry areas, the water height is set to zero in the images (not in the computations) when  $h \leq 10^{-3}h_0$ . For rendering purposes, the elevation map and the water height in the images are scaled by 3.

## REFERENCES

- [1] E. AUDUSSE AND M.-O. BRISTEAU, *A well-balanced positivity preserving “second-order” scheme for shallow water flows on unstructured meshes*, J. Comput. Phys., 206 (2005), pp. 311–333.
- [2] E. AUDUSSE, F. BOUCHUT, M.-O. BRISTEAU, R. KLEIN, AND B. PERTHAME, *A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows*, SIAM J. Sci. Comput., 25 (2004), pp. 2050–2065.
- [3] M. BERGER, M. J. AFTOSMIS, AND S. M. MURMAN, *Analysis of Slope Limiters on Irregular Grids*, AIAA Paper 2005-0490, American Institute for Aeronautics and Astronautics, Reno, NV, 2005.
- [4] A. BERMUDEZ AND M. E. VAZQUEZ, *Upwind methods for hyperbolic conservation laws with source terms*, Comput. Fluids, 23 (1994), pp. 1049–1071.
- [5] C. BERTHON AND F. FOUCHER, *Efficient well-balanced hydrostatic upwind schemes for shallow-water equations*, J. Comput. Phys., 231 (2012), pp. 4993–5015.
- [6] A. BOLLERMANN, S. NOELLE, AND M. LUKÁČOVÁ-MEDVIDOVÁ, *Finite volume evolution Galerkin methods for the shallow water equations with dry beds*, Commun. Comput. Phys., 10 (2011), pp. 371–404.
- [7] F. BOUCHUT, *Nonlinear Stability of Finite Volume Methods for Hyperbolic Conservation Laws and Well-Balanced Schemes for Sources*, Front. Math., Birkhäuser, Basel, 2004.
- [8] F. BOUCHUT AND H. FRID, *Finite difference schemes with cross derivatives correctors for multidimensional parabolic systems*, J. Hyperbolic Differ. Equ., 3 (2006), pp. 27–52.

- [9] S. BRYSON, Y. EPSHTEYN, A. KURGANOV, AND G. PETROVA, *Well-balanced positivity preserving central-upwind scheme on triangular grids for the Saint-Venant system*, ESAIM Math. Model. Numer. Anal., 45 (2011), pp. 423–446.
- [10] O. DELESTRE, S. CORDIER, F. DARBOUX, AND F. JAMES, *A limitation of the hydrostatic reconstruction technique for shallow water equations*, C. R. Math., 350 (2012), pp. 677–681.
- [11] O. DELESTRE, C. LUCAS, P.-A. K SINANT, S. CORDIER, F. DARBOUX, C. LAGUERRE, T.-N.-T. VO, AND F. JAMES, *SWASHES: A compilation of shallow water analytic solutions for hydraulic and environmental studies*, Internat. J. Numer. Methods Fluids, 72 (2013), pp. 269–300.
- [12] A. DURAN, Q. LIANG, AND F. MARCHE, *On the well-balanced numerical discretization of shallow water equations on unstructured meshes*, J. Comput. Phys., 235 (2013), pp. 565–586.
- [13] H. FRID, *Maps of convex sets and invariant regions for finite-difference systems of conservation laws*, Arch. Ration. Mech. Anal., 160 (2001), pp. 245–269.
- [14] J. M. GALLARDO, C. PARÉS, AND M. CASTRO, *On a well-balanced high-order finite volume scheme for shallow water equations with topography and dry areas*, J. Comput. Phys., 227 (2007), pp. 574–601.
- [15] 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 (1996), pp. 1–16.
- [16] J.-L. GUERMOND AND B. POPOV, *Invariant domains and first-order continuous finite element approximation for hyperbolic systems*, SIAM J. Numer. Anal., 54 (2016), pp. 2466–2489.
- [17] J.-L. GUERMOND AND B. POPOV, *Invariant domains and second-order continuous finite element approximation for scalar conservation equations*, SIAM J. Numer. Anal., to appear.
- [18] J.-L. GUERMOND AND B. POPOV, *Estimation From Above of the Maximum Wave Speed in the Riemann Problem for the Euler Equations and Related Problems*, manuscript.
- [19] J.-L. GUERMOND AND B. POPOV, *Second-Order Positivity Preserving Continuous Finite Element Approximation of Hyperbolic Systems of Conservation Laws*, manuscript.
- [20] D. HOFF, *A finite difference scheme for a system of two conservation laws with artificial viscosity*, Math. Comp., 33 (1979), pp. 1171–1193.
- [21] D. HOFF, *Invariant regions for systems of conservation laws*, Trans. Amer. Math. Soc., 289 (1985), pp. 591–610.
- [22] J. F. B. M. KRAAIJEVANGER, *Contractivity of Runge-Kutta methods*, BIT, 31 (1991), pp. 482–528.
- [23] A. KURGANOV AND G. PETROVA, *A second-order well-balanced positivity preserving central-upwind scheme for the Saint-Venant system*, Commun. Math. Sci., 5 (2007), pp. 133–160.
- [24] P. D. LAX, *Weak solutions of nonlinear hyperbolic equations and their numerical computation*, Comm. Pure Appl. Math., 7 (1954), pp. 159–193.
- [25] P. L. F. LIU, Y.-S. CHO, M. J. BRIGGS, U. KANOGLU, AND C. E. SYNOLAKIS, *Runup of solitary waves on a circular island*, J. Fluid Mech., 302 (1995), pp. 259–285.
- [26] 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 (2007), pp. 29–58.
- [27] B. PERTHAME AND C. SIMEONI, *A kinetic scheme for the Saint-Venant system with a source term*, Calcolo, 38 (2001), pp. 201–231.
- [28] M. RICCHIUTO AND A. BOLLERMANN, *Stabilized residual distribution for shallow water simulations*, J. Comput. Phys., 228 (2009), pp. 1071–1115.
- [29] A. RITTER, *Die fortpflanzung der wasserwellen*, Z. Vereines Deutsch. Ingen., 36 (1892), pp. 947–954.
- [30] C.-W. SHU AND S. OSHER, *Efficient implementation of essentially non-oscillatory shock-capturing schemes*, J. Comput. Phys., 77 (1988), pp. 439–471.
- [31] W. C. THACKER, *Some exact solutions to the nonlinear shallow-water wave equations*, J. Fluid Mech., 107 (1981), pp. 499–508.
- [32] Y. XING AND C.-W. SHU, *A survey of high order schemes for the shallow water equations*, J. Math. Study, 47 (2014), pp. 221–249.