# A FULLY WELL-BALANCED SCHEME FOR SHALLOW WATER EQUATIONS WITH CORIOLIS FORCE[*]

VIVIEN DESVEAUX[†] AND ALICE MASSET[‡]

**Abstract.** The present work is devoted to the derivation of a fully well-balanced and positivity-preserving numerical scheme for the shallow water equations with Coriolis force. The first main issue consists in preserving all the steady states. Our strategy relies on a Godunov-type scheme with suitable source term and steady state discretisations. The preservation of moving steady states may lead to ill-defined intermediate states in the Riemann solver. Therefore, a proper correction is introduced in order to obtain a fully well-balanced scheme. The second challenge lies in improving the order of the scheme while preserving the fully well-balanced property. A modification of the classical methods is required since no conservative reconstruction can preserve all the steady states in the case of rotating shallow water equations. A steady state detector is used to overcome this matter. Some numerical experiments are presented to show the relevance and accuracy of both first-order and second-order schemes.

**Keywords.** Shallow water equations; Coriolis force; fully well-balanced schemes; Godunov-type schemes; high-order approximation.

**AMS subject classifications.** 65M08; 65M12.

## 1. Introduction

In the present work we consider the one-dimensional shallow water system with transverse velocity and Coriolis force. This system is also known as 1D rotating shallow-water equations (RSW) and is given by

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \dfrac{gh^2}{2}\right) = fhv - gh\partial_x z, \\ \partial_t(hv) + \partial_x(huv) = -fhu, \end{cases} \tag{1.1}$$

where $h(x,t)$ denotes the fluid height, $u(x,t)$ and $v(x,t)$ are the two components of the horizontal velocity, $z(x)$ designates the topography and is a given function, $g$ is the constant gravitational acceleration and $f$ the Coriolis parameter. This system can be written under the more compact form $\partial_t w + \partial_x F(w) = S(w,z)$ with

$$w = \begin{pmatrix} h \\ hu \\ hv \end{pmatrix}, \quad F(w) = \begin{pmatrix} hu \\ hu^2 + \frac{gh^2}{2} \\ huv \end{pmatrix},$$

and $S(w,z) = S_{cor}(w) + S_{topo}(w)\partial_x z$ where we set

$$S_{cor}(w) = \begin{pmatrix} 0 \\ fhv \\ -fhu \end{pmatrix} \text{ and } S_{topo}(w) = \begin{pmatrix} 0 \\ -gh \\ 0 \end{pmatrix}.$$

The first source term is related to the Coriolis force and the second one to the topography. The vector $w$ must belong to the convex set of admissible states

$$\Omega = \{w = (h, hu, hv)^T \in \mathbb{R}^3; h > 0\}.$$

This 1D system can be obtained from the 2D RSW equations,

$$\begin{cases} \partial_t h + \partial_x(hu) + \partial_y(hv) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \dfrac{gh^2}{2}\right) + \partial_y(huv) = fhv - gh\partial_x z, \\ \partial_t(hv) + \partial_x(huv) + \partial_y\left(hv^2 + \dfrac{gh^2}{2}\right) = -fhu - gh\partial_y z, \end{cases} \tag{1.2}$$

in which the variations in the $y$ direction are neglected.

The RSW system takes into account the force due to the Earth's rotation through the Coriolis term and can therefore model large-scale oceanic or atmospheric fluid flows. One remarkable behaviour of geophysical flows is the geostrophic equilibrium, that received a great attention in the literature these last years, see [3, 4, 9, 11, 18, 22, 26, 31] for instance. Most oceanic and atmospheric circulations are perturbations of the geostrophic equilibrium, which express the balance between the Coriolis force and the horizontal pressure force, as follows in 2D

$$g\nabla(h+z) = f\begin{pmatrix} v \\ -u \end{pmatrix}.$$

In 1D, the geostrophic equilibrium writes

$$\begin{cases} u = 0, \\ g\partial_x(h+z) = fv, \end{cases} \tag{1.3}$$

which is a steady solution of (1.1) with no tangential velocity. Let us notice that in 1D, all the steady solutions of (1.1) with no tangential velocity are described by the geostrophic equilibrium (1.3). With a zero velocity $v$, we recover the lake at rest solution of the classical shallow-water model.

From a numerical point of view, it is well-known since the early works [5, 19, 21], that numerical schemes should capture accurately the steady solutions in order to avoid spurious oscillations, especially on coarse grids. In the few last decades, a large literature was devoted to design such well-balanced schemes that are able to preserve steady solution at rest in different contexts. For the classical shallow-water equations, we can mention the hydrostatic reconstruction method proposed in [1] and numerous other works using various methods, including [13, 14, 24]. Concerning the RSW system, some authors have developed numerical schemes which preserve exactly the geostrophic equilibrium (1.3), for instance in [9, 11, 25, 26].

More recently, some numerical schemes were derived to preserve all the steady states, including the moving ones. Let us emphasize that it is in general a very challenging task to derive such fully well-balanced schemes. The first attempt was in [16] where the author obtains a scheme that preserves all the sonic steady states. In [10], the authors derive a scheme that captures all the steady states of the shallow-water equations with topography. However, this scheme was not able to preserve the positivity of the water height. The first fully well-balanced and positivity-preserving scheme was

derived by Berthon-Chalons [6]. Later, fully well-balanced schemes were also derived for the shallow-water equations with both topography and friction in [28] and for the blood flow equations in [15].

For the 1D RSW equations, the steady solutions are described by

$$
\begin{cases}
\partial_x(hu) = 0, \\
\partial_x\left(hu^2 + \dfrac{gh^2}{2}\right) = fhv - gh\partial_x z, \\
(hu)\partial_x v = -fhu.
\end{cases}
\tag{1.4}
$$

To the best of our knowledge, no fully well-balanced scheme was proposed for the 1D RSW equations. In this system, there is an additional difficulty due to the complex structure of the steady states. Indeed, let us notice that the steady solutions with nonzero tangential velocity satisfy

$$
\begin{cases}
\partial_x(hu) = 0, \\
\partial_x\left(\dfrac{u^2}{2} + g(h+z)\right) = fv, \\
\partial_x v = -f.
\end{cases}
\tag{1.5}
$$

Thus the steady solutions with no tangential velocity described by (1.3) cannot be obtained by setting $u = 0$ in (1.5). It leads to two different families of steady states. This is a discrepancy with the standard shallow-water model, where the lake at rest can be obtained by setting $u = 0$ in the moving steady states equations. The first aim of this paper is therefore to derive a fully well-balanced and positivity-preserving scheme for the one-dimensional RSW equations.

Another issue arises with the 1D RSW equations when we try to increase the order of precision, while preserving the well-balanced property. For other systems with source terms, well-balanced second-order extensions exist. The reader is referred, for instance, to [8,27] for the shallow-water system with topography, [28] for the shallow-water system with both topography and friction and [15] for the blood flow equations. In all these extensions, the main ingredient lies in a reconstruction procedure that preserves the discrete steady states. In order to get the well-balanced and positivity-preserving properties, a standard method is to consider the second-order scheme as a convex combination of first-order schemes on half cells. It requires to use a conservative reconstruction of the unknown variables. Unfortunately, a conservative reconstruction which is fully well-balanced and positivity-preserving is not possible in the case of the 1D RSW, as it will be explained in Section 4.1.

In [28] and [15], a discrete steady state detection procedure is performed. The purpose is to modify the limitation procedure in order to recover the well-balanced first-order scheme near steady states and keep the high-order scheme far from steady state. We propose to adapt this technique for the 1D RSW equations. However, since the numerical flux depends on the space step, we have to complement this technique by adjusting the space step to recover the fully well-balanced first-order scheme at steady states and the second-order scheme far from steady states.

The paper is organized as follows. In Section 2, we start by recalling some general notions about Godunov-type schemes and we choose the discretisation of the continuous steady solutions the scheme will have to preserve. Next, Section 3 is devoted to the derivation of an approximate Riemann solver that leads to a fully well-balanced and

positivity-preserving scheme. In Section 4, we recall the principle of the classical second-order MUSCL extension and we explain why it cannot give a fully well-balanced scheme for the RSW system. Therefore, we present a new strategy based on a discrete steady state detection to recover this property. We also check that this modification does not create non-positive fluid height values. In Section 5, we show some numerical examples that illustrate the fully well-balanced property and the accuracy of both first-order and second-order schemes. Finally, we give some concluding remarks in Section 6.

## 2. Godunov-type scheme with source terms

The numerical scheme we will derive to approximate system (1.1) is a Godunov-type scheme. In this section, we recall the framework of this family of finite volume schemes and we set the notations.

**2.1. Principle.**        We consider a space discretisation made of cells $K_i = (x_{i-1/2}, x_{i+1/2})$, with constant length $\Delta x$. The center of the cell $K_i$ is denoted by $x_i$. The topography is discretized by

$$z_i = \frac{1}{\Delta x} \int_{K_i} z(x) dx.$$

At time $t^n$, we assume that an approximation of the solution of (1.1) is known, which is constant on each cell, and we denote it by

$$w_{\Delta x}(x, t^n) = w_i^n, \text{ if } x \in K_i.$$

In order to simplify the notations, we set the augmented vector $\widetilde{w} = (w, z)$, which belongs to the set

$$\widetilde{\Omega} = \Omega \times \mathbb{R} = \{\widetilde{w} = (h, hu, hv, z)^T \in \mathbb{R}^4; h > 0\}.$$

Since $z$ does not depend on time, we have $\widetilde{w}_i^n = (w_i^n, z_i)$. We aim to update this approximation at time $t^{n+1} = t^n + \Delta t$, with a step $\Delta t$ chosen according to a CFL condition.

Godunov-type schemes are mainly based on Riemann problems, which are Cauchy problems for system (1.1) with an initial data of the form

$$w(x, 0) = \begin{cases} w_L \text{ if } x < 0, \\ w_R \text{ if } x > 0, \end{cases} \tag{2.1}$$

and a topography given by

$$z(x) = \begin{cases} z_L \text{ if } x < 0, \\ z_R \text{ if } x > 0. \end{cases} \tag{2.2}$$

We denote the exact solution of (1.1)–(2.1)–(2.2) by $\mathcal{W}_R(\frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R)$. Let us point out that this solver depends on $(\widetilde{w}_L, \widetilde{w}_R) \in \widetilde{\Omega}^2$ but belongs to the set $\Omega \subset \mathbb{R}^3$ since the variables $h, hu$ and $hv$ evolve through time but the topography $z$ does not.

This exact solution is usually very difficult to compute. Therefore, we prefer to use an approximate Riemann solver $\widehat{\mathcal{W}}_R\left(\frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R\right)$ instead. According to [20], the approximate Riemann solver has to satisfy the following consistency property:

$$\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \widehat{\mathcal{W}}_R\left(\frac{x}{\Delta t}, \widetilde{w}_L, \widetilde{w}_R\right) dx = \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \mathcal{W}_R\left(\frac{x}{\Delta t}, \widetilde{w}_L, \widetilde{w}_R\right) dx.$$

The average of the exact Riemann solution can be computed and the previous condition is equivalent to

$$
\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \widehat{\mathcal{W}}_R \left( \frac{x}{\Delta t}, \widetilde{w}_L, \widetilde{w}_R \right) dx
$$
$$
= \frac{w_L + w_R}{2} - \frac{\Delta t}{\Delta x} (F(w_R) - F(w_L)) + \frac{1}{\Delta x} \int_0^{\Delta t} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} S \left( \mathcal{W}_R \left( \frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R \right), z(x) \right) dx dt.
$$

(2.3)

In the absence of source terms, we can enforce this equality to ensure the consistency of the approximate Riemann solver. However, it is not always possible to compute exactly the average of the source term. Therefore, it is standard to use a suitable approximation (see for instance [6, 8, 12])

$$
\mathcal{S}(\widetilde{w}_L, \widetilde{w}_R) \approx \frac{1}{\Delta t} \int_0^{\Delta t} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} S \left( \mathcal{W}_R \left( \frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R \right), z(x) \right) dx dt.
$$

This numerical source term should be consistent with the continuous source term $S$ in the following sense.

DEFINITION 2.1.     *The numerical source term $\mathcal{S}$ is consistent with the continuous source term $S(\widetilde{w}) = S_{cor}(w) + S_{topo}(w)\partial_x z$ if it satisfies*

$$
\mathcal{S}((w, z_L), (w, z_R)) = S_{cor}(w)\Delta x + S_{topo}(w)[z].
$$

(2.4)

Provided a consistent numerical source term, the approximate Riemann solver can only satisfy a weaker version of (2.3). It leads to the definition of a weakly consistent approximate Riemann solver.

DEFINITION 2.2.     *The approximate Riemann solver $\widehat{\mathcal{W}}_R$ is weakly consistent if there exists a consistent numerical source term $\mathcal{S}$ such that*

$$
\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \widehat{\mathcal{W}}_R \left( \frac{x}{\Delta t}, \widetilde{w}_L, \widetilde{w}_R \right) dx = \frac{w_L + w_R}{2} - \frac{\Delta t}{\Delta x} (F(w_R) - F(w_L)) + \frac{\Delta t}{\Delta x} \mathcal{S}(\widetilde{w}_L, \widetilde{w}_R).
$$

(2.5)

The following section will be devoted to derive a weakly consistent approximate Riemann solver. For now, we show how we can obtain a numerical scheme from an approximate Riemann solver. A Godunov-type scheme is built in two steps:

- firstly, we consider the juxtaposition of approximate Riemann solvers at each interface $x_{i+1/2}$,

$$
w_{\Delta x}(x, t^n + t) = \widehat{\mathcal{W}}_R \left( \frac{x - x_{i+1/2}}{t}, \widetilde{w}_i^n, \widetilde{w}_{i+1}^n \right), \text{ if } x \in (x_i, x_{i+1});
$$

- secondly, the update at time $t^{n+1}$ is obtained by averaging the previous function on each cell

$$
w_i^{n+1} = \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} w_{\Delta x}(x, t^n + \Delta t) dx,
$$

or equivalently

$$
w_i^{n+1} = \frac{1}{\Delta x} \int_0^{\frac{\Delta x}{2}} \widehat{\mathcal{W}}_R \left( \frac{x}{\Delta t}, \widetilde{w}_{i-1}^n, \widetilde{w}_i^n \right) dx + \frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^0 \widehat{\mathcal{W}}_R \left( \frac{x}{\Delta t}, \widetilde{w}_i^n, \widetilde{w}_{i+1}^n \right) dx.
$$

(2.6)

In order to prevent the approximate Riemann solvers to interact between each other, we must enforce the CFL restriction

$$\frac{\Delta t}{\Delta x}\max_{i\in\mathbb{Z}}|\lambda^{\pm}(\widetilde{w}_i^n,\widetilde{w}_{i+1}^n)|\leq\frac{1}{2}, \tag{2.7}$$

where $\lambda^{\pm}(\widetilde{w}_L,\widetilde{w}_R)$ denotes both the maximum and minimum speed of the waves that appear in $\widehat{\mathcal{W}}_R(\frac{x}{t},\widetilde{w}_L,\widetilde{w}_R)$. Under this condition, we can write the Godunov-type scheme as a finite volume scheme (see for instance [6])

$$w_i^{n+1}=w_i^n-\frac{\Delta t}{\Delta x}(\mathcal{F}_{i+1/2}^n-\mathcal{F}_{i-1/2}^n)+\frac{\Delta t}{2\Delta x}(\mathcal{S}_{i+1/2}^n+\mathcal{S}_{i-1/2}^n), \tag{2.8}$$

with $\mathcal{F}_{i+1/2}^n=\mathcal{F}(\widetilde{w}_i^n,\widetilde{w}_{i+1}^n)$ and $\mathcal{S}_{i+1/2}^n=\mathcal{S}(\widetilde{w}_i^n,\widetilde{w}_{i+1}^n)$, where the numerical flux is given by

$$\mathcal{F}(\widetilde{w}_L,\widetilde{w}_R)=\frac{F(w_L)+F(w_R)}{2}-\frac{\Delta x}{4\Delta t}(w_R-w_L)$$
$$+\frac{1}{2\Delta t}\left(\int_0^{\frac{\Delta x}{2}}\widehat{\mathcal{W}}_R\left(\frac{x}{\Delta t},\widetilde{w}_L,\widetilde{w}_R\right)dx-\int_{-\frac{\Delta x}{2}}^0\widehat{\mathcal{W}}_R\left(\frac{x}{\Delta t},\widetilde{w}_L,\widetilde{w}_R\right)dx\right), \tag{2.9}$$

and the numerical source term $\mathcal{S}(\widetilde{w}_L,\widetilde{w}_R)$ is the same as introduced in Definition 2.2.

At this point, the only property the scheme has to satisfy is the weak consistency of the approximate Riemann solver. We now list some other properties the scheme should satisfy.

**2.2. Numerical scheme properties.**     We present two important features of numerical schemes in this context: robustness and well-balancing. Godunov-type schemes have the advantage of inheriting these properties from the approximate Riemann solver $\widehat{\mathcal{W}}_R$. First, we study the preservation of fluid height positivity.

LEMMA 2.1.  *If the approximate Riemann solver $\widehat{\mathcal{W}}_R$ satisfies the robustness condition*

$$\forall(\widetilde{w}_L,\widetilde{w}_R)\in\widetilde{\Omega}^2,\forall\xi\in\mathbb{R},\widehat{\mathcal{W}}_R(\xi,\widetilde{w}_L,\widetilde{w}_R)\in\Omega, \tag{2.10}$$

*then under the CFL condition* (2.7), *the Godunov-type scheme* (2.8) *preserves the positivity of the fluid height:*

$$\forall i\in\mathbb{Z},h_i^n>0\Rightarrow\forall i\in\mathbb{Z},h_i^{n+1}>0.$$

The reader can refer, for instance, to [8] to find a proof of this result.

Similarly, the Godunov-type scheme is well-balanced as soon as the approximate Riemann solver is. To be more specific, we have to introduce the notion of local steady states. In the following, for any quantity $X$ which has a left value $X_L$ and a right value $X_R$, we will use the following notations

$$[X]=X_R-X_L,\qquad\overline{X}=\frac{X_L+X_R}{2}.$$

DEFINITION 2.3.    *A couple of states $(\widetilde{w}_L,\widetilde{w}_R)$ defines a local steady state for the system* (1.1) *if it satisfies*

$$\begin{cases}h_Ru_R=h_Lu_L=q,\\\left[\dfrac{u^2}{2}+g(h+z)\right]=\Delta xf\overline{v},\\q[v]=-\Delta xfq,\end{cases} \tag{2.11}$$

*or equivalently if the local steady state indicator*

$$\mathcal{E}(\widetilde{w}_L, \widetilde{w}_R, \Delta x) = \sqrt{\left|[hu]\right|^2 + \left|\left[\frac{u^2}{2} + g(h+z)\right] - \Delta x f\overline{v}\right|^2 + \left|\overline{hu}([v] + f\Delta x)\right|^2}, \quad (2.12)$$

*is equal to zero.*

All along this paper, we write $\mathcal{E}_{LR}$ rather than $\mathcal{E}(\widetilde{w}_L, \widetilde{w}_R, \Delta x)$ if no ambiguity is possible. Let us notice that (2.11) is actually a discretization of the Equations (1.4) that define the continuous steady states. Other choices of discretization could be possible, especially in the choice of the mean value $\overline{v}$.

The definition of a well-balanced Riemann solver follows.

DEFINITION 2.4. *An approximate Riemann solver is called well-balanced if*

$$\widehat{\mathcal{W}}_R\left(\frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R\right) = \begin{cases} w_L \ \ if \ x < 0, \\ w_R \ \ if \ x > 0, \end{cases}$$

*as soon as $(\widetilde{w}_L, \widetilde{w}_R)$ is a local steady state.*

Similarly, we define a discrete steady state and a well-balanced scheme.

DEFINITION 2.5.
(1) *A sequence $(\widetilde{w}_i^n)_{i\in\mathbb{Z}}$ defines a discrete steady state if the couples $(\widetilde{w}_i, \widetilde{w}_{i+1})$ are local steady states for all $i \in \mathbb{Z}$.*
(2) *A numerical scheme is called well-balanced if for each discrete steady state $(\widetilde{w}_i^n)_{i\in\mathbb{Z}}$, we have*

$$w_i^{n+1} = w_i^n, \ \forall i \in \mathbb{Z}.$$

Equipped with these definitions, we have the following statement.

LEMMA 2.2. *If the approximate Riemann solver $\widehat{\mathcal{W}}_R$ is well-balanced, then the associated Godunov-type scheme (2.8) is well-balanced.*

*Proof.* Let a sequence $(\widetilde{w}_i^n)_{i\in\mathbb{Z}}$ be a discrete steady state. Then, according to (2.6), the updated approximation satisfies $w_i^{n+1} = w_i^n$ for all $i \in \mathbb{Z}$. □

To summarize, the approximate Riemann solver that we will derive in the next section has to satisfy the weak consistency condition (2.5), the robustness condition (2.10) and the fully well-balanced property given by Definition 2.4.

## 3. A fully well-balanced Godunov-type scheme

Here, we propose an approximate Riemann solver for the system (1.1) that satisfies the three required properties. We adapt to the RSW system the strategy mostly proposed in [15, 27, 28] for different systems.

**3.1. Source term discretisation.** The aim of this section is to propose a numerical source term

$$\mathcal{S}(\widetilde{w}_L, \widetilde{w}_R) = (0, \mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R), \mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R))^T$$

which is consistent with the continuous source term $S$ in the sense of Definition 2.1. Moreover this choice of a numerical source term has to be coherent with the required well-balanced property.

To this end, we start by considering a Riemann data $(\widetilde{w}_L, \widetilde{w}_R)$ which is a local steady state according to Definition 2.3. Since we want the approximate Riemann solver to be both weakly consistent and well-balanced, the condition (2.5) enforces

$$\mathcal{S}(\widetilde{w}_L, \widetilde{w}_R) = F(w_R) - F(w_L), \tag{3.1}$$

or equivalently

$$\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R) = h_R u_R^2 + \frac{g h_R^2}{2} - h_L u_L^2 - \frac{g h_L^2}{2},$$
$$\mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R) = h_R u_R v_R - h_L u_L v_L.$$

Based on the chosen Definition 2.3 of the local steady states, these relations can be written by

$$\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R) = \left( g\overline{h} - \frac{q^2}{h_L h_R} \right) [h], \tag{3.2}$$

$$\mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R) = -\Delta x f q. \tag{3.3}$$

The expression (3.2) cannot be used to define the numerical source term in the general case since it would not be consistent in the sense of Definition 2.1. Hence, we continue to develop this expression for a local steady state. First, from the second equality of (2.11), we get

$$\frac{q^2}{2h_R^2} + g(h_R + z_R) - \frac{q^2}{2h_L^2} - g(h_L + z_L) = \Delta x f \overline{v}, \tag{3.4}$$

which leads to

$$[h] \left( 1 - \frac{q^2 \overline{h}}{g h_L^2 h_R^2} \right) = \Delta x f \overline{v}/g - (z_R - z_L). \tag{3.5}$$

It follows

$$[h] = \frac{\Delta x f \overline{v}/g - [z]}{1 - \mathrm{Fr}}, \tag{3.6}$$

where $\mathrm{Fr} = \frac{\overline{h}|u_L u_R|}{g h_L h_R}$ is a discrete Froude number. Injecting this relation into (3.2), we get

$$\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R) = \Delta x f \overline{h}\,\overline{v} - g\overline{h}[z] + \frac{g \mathrm{Fr}[h]^2}{4\overline{h}} \frac{(\Delta x f \overline{v}/g - [z])}{(1 - \mathrm{Fr})}. \tag{3.7}$$

We inject one more time (3.6) in the above equality to obtain a more convenient expression

$$\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R) = \Delta x f \overline{h}\,\overline{v} - g\overline{h}[z] + \frac{g \mathrm{Fr}[h]}{4\overline{h}} \frac{(\Delta x f \overline{v}/g - [z])^2}{(1 - \mathrm{Fr})^2}. \tag{3.8}$$

This expression is *a priori* not well-defined if $\mathrm{Fr} = 1$. However, the combination of (3.6) and (3.7) leads to a paraphrase of $\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R)$ in the form

$$\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R) = g\overline{h}[h](1 - \mathrm{Fr}) + \frac{g}{4\overline{h}} \mathrm{Fr}[h]^3.$$

Therefore $\mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R)$ admits the following limit when Fr goes to 1

$$\lim_{\text{Fr}\to 1} \mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R) = \frac{g}{4\overline{h}}[h]^3. \tag{3.9}$$

At this point, (3.8) and (3.3) are suitable definitions for the numerical source terms $\mathcal{S}^{hu}$ and $\mathcal{S}^{hv}$ when a local steady state is considered. However, let us point out that the limit (3.9) is only valid for a local steady state. Therefore, the right-hand side of (3.8) is not well-defined when $\varepsilon_{LR}\neq 0$ and $\text{Fr}=1$. To deal with this issue, we add a nonnegative term $\mathcal{E}_{LR}$ to the denominator as follows

$$\mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R) = \Delta x f\overline{h}\,\overline{v} - g\overline{h}[z] + \frac{g\text{Fr}[h]}{4\overline{h}}\frac{(\Delta x f\overline{v}/g - [z])^2}{(1-\text{Fr})^2 + \mathcal{E}_{LR}}. \tag{3.10}$$

Indeed, the denominator in (3.10) can only vanish when $\varepsilon_{LR}=0$, which means the Riemann data $(\widetilde{w}_L,\widetilde{w}_R)$ is a local steady state. But then the source term can be defined by the limit (3.9) as mentioned before.

To generalise (3.3) away from local steady states, we need to define a general discharge $\widetilde{q}$ which coincides with $q$ as soon as a local steady state is considered or as soon as $w_L = w_R = w$. There are several possible definitions, $\widetilde{q} = \overline{hu}$ for instance.

We finally obtain the following definitions for the numerical source terms

$$\mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R) = \begin{cases} \Delta x f\overline{h}\,\overline{v} - g\overline{h}[z] + \dfrac{g\text{Fr}[h]}{4\overline{h}}\dfrac{(\Delta x f\overline{v}/g - [z])^2}{(1-\text{Fr})^2 + \mathcal{E}_{LR}} & \text{if } \text{Fr}\neq 1 \text{ or } \mathcal{E}_{LR}\neq 0, \\[4mm] \dfrac{g}{4\overline{h}}[h]^3 & \text{if } \text{Fr}=1 \text{ and } \mathcal{E}_{LR}=0. \end{cases} \tag{3.11}$$

$$\mathcal{S}^{hv}(\widetilde{w}_L,\widetilde{w}_R) = -\Delta x f\widetilde{q}. \tag{3.12}$$

To conclude, we prove that these numerical source terms are consistent.

LEMMA 3.1. *The numerical source term*

$$\mathcal{S}(\widetilde{w}_L,\widetilde{w}_R) = (0,\mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R),\mathcal{S}^{hv}(\widetilde{w}_L,\widetilde{w}_R))^T$$

*defined by* (3.11) *and* (3.12) *is consistent in the sense of Definition* 2.1.

*Proof.* The consistency is immediate for $\mathcal{S}^{hv}$, as for $\mathcal{S}^{hu}$ in the case $\text{Fr}\neq 1$ or $\mathcal{E}_{LR}\neq 0$. In the case $\text{Fr}=1$ and $\mathcal{E}_{LR}=0$, let us notice that according to (3.6), we have $\Delta x f\overline{v}/g = [z]$. Therefore the source term $\mathcal{S}^{hu}$ can be written under the form

$$\mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R) = \Delta x f\overline{h}\,\overline{v} - g\overline{h}[z] + \frac{g[h]^3}{4\overline{h}},$$

and the consistency follows. □

**3.2. Approximate Riemann solver.** The numerical source term being well-defined, we now turn to build a weakly consistent approximate Riemann solver which is fully well-balanced and preserves the positivity of the fluid height.

Let us notice that the well-balanced property of the approximate Riemann solver strongly depends on the choice that was made in Definition 2.3 to discretise the steady states. However, the following procedure stands for any discretisation of the steady

states. This is not the case for the source term discretisation which was done in the previous section and should be adapted to the steady state discretisation.

We consider a Riemann data $(\widetilde{w}_L, \widetilde{w}_R) \in \widetilde{\Omega}^2$. We choose to build an approximate Riemann solver $\widehat{\mathcal{W}}_R$ with four constant states separated by three discontinuities with respective speeds $\lambda_L < 0$, $\lambda_0 = 0$ and $\lambda_R > 0$, as described in Figure 3.1. This approximate Riemann solver writes

$$
\widehat{\mathcal{W}}_R \left( \frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R \right) = \begin{cases} w_L & \text{if } \dfrac{x}{t} < \lambda_L, \\ w_L^{\star} & \text{if } \lambda_L < \dfrac{x}{t} < 0, \\ w_R^{\star} & \text{if } 0 < \dfrac{x}{t} < \lambda_R, \\ w_R & \text{if } \dfrac{x}{t} > \lambda_R. \end{cases} \tag{3.13}
$$

This leads to two intermediate states $w_L^{\star}$ and $w_R^{\star}$ and thus six unknowns. In order
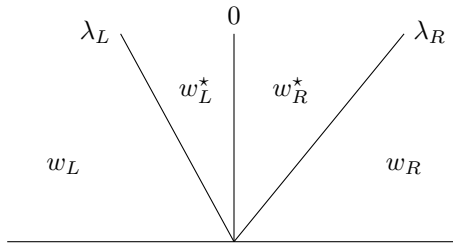


FIG. 3.1. *Approximate Riemann solver* $\widehat{\mathcal{W}}_R$

to simplify the subsequent notations, we introduce the intermediate state of the HLL approximate Riemann solver (see [20])

$$
w^{HLL} = \frac{\lambda_R w_R - \lambda_L w_L}{\lambda_R - \lambda_L} - \frac{F(w_R) - F(w_L)}{\lambda_R - \lambda_L}. \tag{3.14}
$$

Let us notice that its first component can be written as

$$
h^{HLL} = \frac{u_L - \lambda_L}{\lambda_R - \lambda_L} h_L + \frac{\lambda_R - u_R}{\lambda_R - \lambda_L} h_R.
$$

As a consequence, as soon as the speeds $\lambda_L$ and $\lambda_R$ satisfy

$$
\lambda_L < u_L \qquad \text{and} \qquad \lambda_R > u_R, \tag{3.15}
$$

we have $h^{HLL} > 0$.

We are searching for six relations to define the unknowns $w_L^{\star}$ and $w_R^{\star}$. The first three relations come from the weak consistency condition. Noticing that the average of the approximate Riemann solver $\widehat{\mathcal{W}}_R$ is given by

$$
\frac{1}{\Delta x} \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} \widehat{\mathcal{W}}_R \left( \frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R \right) = \frac{w_L + w_R}{2} - \frac{\Delta t}{\Delta x} (\lambda_R w_R - \lambda_L w_L) + \frac{\Delta t}{\Delta x} (\lambda_R w_R^{\star} - \lambda_L w_L^{\star}),
$$

the weak consistency condition (2.5) writes

$$
\lambda_R h_R^{\star} - \lambda_L h_L^{\star} = (\lambda_R - \lambda_L) h^{HLL}, \tag{3.16}
$$

$$\lambda_R h_R^\star u_R^\star - \lambda_L h_L^\star u_L^\star = (\lambda_R - \lambda_L)(hu)^{HLL} + \mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R), \tag{3.17}$$

$$\lambda_R h_R^\star v_R^\star - \lambda_L h_L^\star v_L^\star = (\lambda_R - \lambda_L)(hv)^{HLL} + \mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R). \tag{3.18}$$

The three missing relations will come from the fully well-balanced constraint given by Definition 2.4. In other words, we have to choose three additional relations such that the solution of the system formed by these relations and Equations (3.16), (3.17) and (3.18) satisfies

$$w_R^\star = w_R \quad \text{and} \quad w_L^\star = w_L,$$

as soon as $(\widetilde{w}_L, \widetilde{w}_R)$ is a local steady state. We will deal with each variable separately.

First, we are going to determine the intermediate discharges $h_L^\star u_L^\star$ and $h_R^\star u_R^\star$. Since $hu$ is the Riemann invariant associated with the characteristic field of speed $\lambda_0$, it is natural to enforce the relation

$$h_L^\star u_L^\star = h_R^\star u_R^\star = q^\star, \tag{3.19}$$

The system (3.17)–(3.19) can be solved immediately to obtain the intermediate discharge

$$q^\star = (hu)^{HLL} + \frac{\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R)}{\lambda_R - \lambda_L}. \tag{3.20}$$

Now we are going to determine the intermediate fluid heights $h_L^\star$ and $h_R^\star$. Let us notice that when $(\widetilde{w}_L, \widetilde{w}_R)$ is a local steady state, we have according to (3.2)

$$\alpha_{LR}(h_R - h_L) = \mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R), \tag{3.21}$$

where $\alpha_{LR} = g\overline{h} - |u_L u_R|$. Therefore, a simple choice for the additional equation would be a linearisation of (3.21) as

$$\alpha_{LR}(h_R^\star - h_L^\star) = \mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R).$$

Together with Equation (3.16), this leads to a simple linear system. However this system does not admit a unique solution when $\alpha_{LR}$ vanishes. We suggest the following modification

$$(\alpha_{LR}^2 + \mathcal{E}_{LR})(h_R^\star - h_L^\star) = \alpha_{LR}\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R). \tag{3.22}$$

The coefficient $\alpha_{LR}^2 + \mathcal{E}_{LR}$ can still vanish if both $\alpha_{LR}$ and $\mathcal{E}_{LR}$ vanish. However, it is quite natural in that case to enforce $h_R^\star - h_L^\star = h_R - h_L$ because of the fully well-balanced constraint. Moreover, we can notice that for all $\alpha_{LR} \neq 0$ and according to relations (3.22) and (3.21) we have

$$\lim_{\varepsilon_{LR} \to 0} h_R^\star - h_L^\star = h_R - h_L.$$

Thus we introduce the following quantity

$$\Delta_{LR}^h = \begin{cases} \dfrac{\alpha_{LR}\mathcal{S}^{hu}(\widetilde{w}_L, \widetilde{w}_R)}{\alpha_{LR}^2 + \mathcal{E}_{LR}} & \text{if } \mathcal{E}_{LR} \neq 0, \\ h_R - h_L & \text{if } \mathcal{E}_{LR} = 0, \end{cases}$$

and we choose the following additional equation

$$h_R^\star - h_L^\star = \Delta_{LR}^h. \tag{3.23}$$

Solving the system (3.16)–(3.23), we obtain

$$h_L^\star = h^{HLL} - \frac{\lambda_R}{\lambda_R - \lambda_L} \Delta_{LR}^h,$$

$$h_R^\star = h^{HLL} - \frac{\lambda_L}{\lambda_R - \lambda_L} \Delta_{LR}^h.$$

Nothing ensures these intermediate fluid heights to be positive. To address this issue, we adapt the cut-off procedure suggested in [2, 15, 27]. Let us introduce the threshold

$$\delta = \min(\varepsilon, h_L, h_R, h^{HLL}), \tag{3.24}$$

where $\varepsilon > 0$ is a small parameter. We recall that $h^{HLL}$ is positive as soon as condition (3.15) is satisfied, and therefore we have $\delta > 0$. If $h_L^\star < \delta$, we set $h_L^\star = \delta$, and $h_R^\star$ is modified according to (3.16). In this case, we have

$$h_R^\star = \left(1 - \frac{\lambda_L}{\lambda_R}\right) h^{HLL} + \frac{\lambda_L}{\lambda_R} h_L^\star \geq \delta,$$

so both intermediate water heights $h_L^\star$ and $h_R^\star$ are positive. We proceed similarly if $h_R^\star < \delta$. Taking this procedure into account, the intermediate fluid heights write

$$h_L^\star = \min\left(\max\left(h^{HLL} - \frac{\lambda_R}{\lambda_R - \lambda_L} \Delta_{LR}^h, \delta\right), \left(1 - \frac{\lambda_R}{\lambda_L}\right) h^{HLL} + \frac{\lambda_R}{\lambda_L} \delta\right), \tag{3.25}$$

$$h_R^\star = \min\left(\max\left(h^{HLL} - \frac{\lambda_L}{\lambda_R - \lambda_L} \Delta_{LR}^h, \delta\right), \left(1 - \frac{\lambda_L}{\lambda_R}\right) h^{HLL} + \frac{\lambda_L}{\lambda_R} \delta\right). \tag{3.26}$$

Finally, we proceed similarly in order to determine the intermediate transverse speeds $v_L^\star$ and $v_R^\star$. We introduce the quantity

$$\Delta_{LR}^v = \begin{cases} \dfrac{\widetilde{q}\, \mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R)}{\widetilde{q}^2 + \mathcal{E}_{LR}} & \text{if } \mathcal{E}_{LR} \neq 0, \\ v_R - v_L & \text{if } \mathcal{E}_{LR} = 0, \end{cases}$$

and we enforce the following additional equation

$$v_R^\star - v_L^\star = \Delta_{LR}^v. \tag{3.27}$$

The system (3.18)–(3.27) then leads to

$$v_L^\star = \frac{(hv)^{HLL}}{h^{HLL}} + \frac{1}{(\lambda_R - \lambda_L) h^{HLL}} \left(\mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R) - \lambda_R h_R^\star \Delta_{LR}^v\right), \tag{3.28}$$

$$v_R^\star = \frac{(hv)^{HLL}}{h^{HLL}} + \frac{1}{(\lambda_R - \lambda_L) h^{HLL}} \left(\mathcal{S}^{hv}(\widetilde{w}_L, \widetilde{w}_R) - \lambda_L h_L^\star \Delta_{LR}^v\right). \tag{3.29}$$

The approximate Riemann solver is then completely defined by relations (3.19), (3.20), (3.25), (3.26), (3.28) and (3.29). Let us notice that it is automatically weakly consistent by the choice of the three first Equations (3.16), (3.17) and (3.18). The cut-off procedure does not alter the weak consistency, since Equation (3.16) is still enforced when it is applied.

Moreover, thanks to the cut-off procedure, we can prove that the approximate Riemann solver is robust.

LEMMA 3.2.    *If the initial fluid heights $h_L$ and $h_R$ are positive and if the speeds $\lambda_L$ and $\lambda_R$ satisfy (3.15), then both intermediate fluid heights $h_L^\star$ and $h_R^\star$ defined by (3.25) and (3.26) are positive.*

*Proof.*    Under these asumptions, the intermediate fluid heights $h_R^\star$ and $h_L^\star$ are greater or equal to $\delta$, which is positive.                                              □

We now prove that the approximate Riemann solver is also well-balanced.

LEMMA 3.3.    *The approximate Riemann solver $\widehat{\mathcal{W}}_R$ is well-balanced.*

*Proof.*    Let us consider a local steady state $(\widetilde{w}_L, \widetilde{w}_R)$. First, we state that the cut-off procedure cannot apply in this case. Indeed, thanks to the Definition (3.24), the intermediate fluid heights computed before the cut-off procedure satisfy $h_L^\star = h_L \geq \delta$ and $h_R^\star = h_R \geq \delta$.

Therefore, we only need to show that $w_L^\star = w_L$ and $w_R^\star = w_R$ to prove the result. Since $(w_L^\star, w_R^\star)$ is defined as the unique solution of the system of Equations (3.16), (3.17), (3.18), (3.19), (3.23), (3.27), it is sufficient to prove that $(w_L, w_R)$ is a solution of this system.

Since we have $\mathcal{E}_{LR} = 0$, it is immediate that $(w_L, w_R)$ is a solution of (3.19), (3.23), (3.27). The approximate solver being weakly consistent and $(\widetilde{w}_L, \widetilde{w}_R)$ being a local steady state, Equation (2.5) enforces

$$\mathcal{S}(\widetilde{w}_L, \widetilde{w}_R) = f(w_R) - f(w_L),$$

so the intermediate state of the HLL solver defined by (3.14) satisfies

$$(\lambda_R - \lambda_L) w^{HLL} = \lambda_R w_R - \lambda_L w_L - \mathcal{S}(\widetilde{w}_L, \widetilde{w}_R).$$

As a consequence, Equations (3.16), (3.17) and (3.18) rewrite

$$\lambda_R h_R^\star - \lambda_L h_L^\star = \lambda_R h_R - \lambda_L h_L,$$
$$\lambda_R h_R^\star u_R^\star - \lambda_L h_L^\star u_L^\star = \lambda_R h_R u_R - \lambda_L h_L u_L,$$
$$\lambda_R h_R^\star v_R^\star - \lambda_L h_L^\star v_L^\star = \lambda_R h_R v_R - \lambda_L h_L v_L.$$

We deduce $(w_L, w_R)$ is a solution of these equations and thus $w_L^\star = w_L$ and $w_R^\star = w_R$.    □

The approximate Riemann solver $\widehat{\mathcal{W}}_R$ thus satisfies all the required properties.

**3.3. The final scheme.**    We summarize in this section the full scheme and its properties.

THEOREM 3.1.    *The approximate Riemann solver (3.13) where the intermediate states are given by (3.19), (3.20), (3.25), (3.26), (3.28) and (3.29) leads to a Godunov-type scheme that can be written under the form (2.8). The numerical flux*

$$\mathcal{F}(\widetilde{w}_L, \widetilde{w}_R) = \left( \mathcal{F}^h(\widetilde{w}_L, \widetilde{w}_R), \mathcal{F}^{hu}(\widetilde{w}_L, \widetilde{w}_R), \mathcal{F}^{hv}(\widetilde{w}_L, \widetilde{w}_R) \right)^T,$$

*is given by*

$$\mathcal{F}^h(\widetilde{w}_L,\widetilde{w}_R) = \overline{hu} + \frac{\lambda_R}{2}(h_R^\star - h_R) + \frac{\lambda_L}{2}(h_L^\star - h_L),$$

$$\mathcal{F}^{hu}(\widetilde{w}_L,\widetilde{w}_R) = \overline{hu^2 + \frac{gh^2}{2}} + \frac{\lambda_R}{2}(h_R^\star u_R^\star - h_R u_R) + \frac{\lambda_L}{2}(h_L^\star u_L^\star - h_L u_L),$$

$$\mathcal{F}^{hv}(\widetilde{w}_L,\widetilde{w}_R) = \overline{huv} + \frac{\lambda_R}{2}(h_R^\star v_R^\star - h_R v_R) + \frac{\lambda_L}{2}(h_L^\star v_L^\star - h_L v_L),$$

*and the numerical source term*

$$\mathcal{S}(\widetilde{w}_L,\widetilde{w}_R) = \left(0, \mathcal{S}^{hu}(\widetilde{w}_L,\widetilde{w}_R), \mathcal{S}^{hv}(\widetilde{w}_L,\widetilde{w}_R)\right)^T$$

*is defined by* (3.11) *and* (3.12).

*Under the CFL restriction* (2.7) *and if the speeds $\lambda_L$ and $\lambda_R$ are chosen according to* (3.15), *this scheme is fully well-balanced and preserves the positivity of $h$.*

*Proof.*  The expression of the numerical flux is obtained from a straightforward computation in (2.9).

Assume $(w_L, w_R)$ are in $\Omega$. Since $h^{HLL} > 0$, the cut-off procedure ensures $h_L^\star > 0$ and $h_R^\star > 0$. Thus the variable $h$ remains positive in the approximate Riemann solver. According to Lemma 2.1, the scheme preserves the positivity of $h$.

The well-balanced property of the scheme is a direct consequence of Lemmas 2.2 and 3.3.                                                                    □

### 4. Second-order scheme

In this section, we propose to improve the scheme precision using the MUSCL method. Our goal is to build a second-order scheme in space that preserves the good properties of the first-order one, namely the positivity of $h$ and the well-balanced property. The second-order in time is obtained with the usual Runge-Kutta method. We do not describe it here, but the reader can refer to [7, 17, 29].

We start by a description of the standard MUSCL method, and we explain why it is not adapted to get the fully well-balanced property for the RSW system. Indeed, no conservative reconstruction can preserve the structure of all the steady states defined by (2.11) as it will be explained in Subsection 4.1. We explain in Subsection 4.2 how to recover the fully well-balanced property by adapting the ideas proposed in [28] and [15] to our generalised MUSCL scheme.

Up to this point, for the sake of conciseness, we neither mentioned explicitly the dependence on $\Delta x$ in the numerical fluxes and the source terms, nor in the definition of local steady states. However in the following, we will consider half-cells, which will impose to make these dependencies appear, in particular to determine if the scheme is fully well-balanced. It will also be useful to consider the $\Delta x$ that appears in the approximate Riemann solver $\widehat{\mathcal{W}}_R$ and the $\Delta x$ that appears in the numerical scheme Definition (2.8) as two separated parameters. For the sake of clarity, the first one will be denoted by $d > 0$, and therefore the approximate Riemann solver writes

$$\widehat{\mathcal{W}}_R\left(\frac{x}{t}, \widetilde{w}_L, \widetilde{w}_R, d\right) = \begin{cases} w_L & \text{if } \frac{x}{t} < \lambda_L, \\ w_L^\star(d) & \text{if } \lambda_L < \frac{x}{t} < 0, \\ w_R^\star(d) & \text{if } 0 < \frac{x}{t} < \lambda_R, \\ w_R & \text{if } \frac{x}{t} > \lambda_R. \end{cases}$$

According to Subsection [2.1], the resulting Godunov-type scheme writes

$$w_i^{n+1} = w_i^n - \frac{\Delta t}{\Delta x} \left( \mathcal{F}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, d\right) - \mathcal{F}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, d\right) \right)$$
$$+ \frac{\Delta t}{2\Delta x} \left( \mathcal{S}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, d\right) + \mathcal{S}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, d\right) \right), \quad (4.1)$$

provided the CFL condition [2.7] is satisfied. Notice that for $d = \Delta x$, we recover the fully well-balanced scheme derived in Section [3]. Moreover, we establish the following lemma, that will be useful for the forthcoming proof.

LEMMA 4.1.    *Under the CFL condition* [2.7] *and if the approximate Riemann solver wavespeeds $\lambda_L$ and $\lambda_R$ satisfy* [3.15], *then the Godunov-type scheme* [4.1] *preserves the positivity of $h$, for all $d > 0$.*

*Proof.*    Independently of the parameter $d$, the cut-off procedure leads to positive intermediate states $h_L^\star$ and $h_R^\star$ according to Definition [3.25]-[3.26], since the wavespeeds $\lambda_L$ and $\lambda_R$ satisfy the condition [3.15]. Then, we apply the Lemma [2.1] to conclude the scheme [4.1] preserves the positivity of $h$ for all $d > 0$.                    □

**4.1. Standard MUSCL method.**    The main idea of the MUSCL method is to reach second-order by considering a linear reconstruction of the solution on each cell, instead of a constant one. We recall here the standard reconstruction procedure.

Starting from a piecewise constant approximation at time $t^n$,

$$\widetilde{w}_{\Delta x}(x, t^n) = \widetilde{w}_i^n \text{ if } x \in K_i,$$

we reconstruct states at the interfaces of each cells as

$$\widetilde{w}_i^{n, \pm} = \widetilde{w}_i^n \pm \frac{\Delta x}{2} \sigma_i^n(\widetilde{w}), \quad (4.2)$$

where $\sigma_i^n(\widetilde{w})$ is a slope vector to determine. Let us emphasize that this procedure includes the topography.

To avoid spurious oscillations, it is well-known that a limitation procedure must be applied to the slopes. In this paper, we consider the minmod limiter function defined by

$$\mathrm{minmod}(\sigma_L, \sigma_R) = \begin{cases} \min(\sigma_L, \sigma_R) & \text{if } \sigma_L > 0 \text{ and } \sigma_R > 0, \\ \max(\sigma_L, \sigma_R) & \text{if } \sigma_L < 0 \text{ and } \sigma_R < 0, \\ 0 & \text{otherwise.} \end{cases}$$

Then the slope vector is defined by $\sigma_i^n(\widetilde{w}) = \mathrm{minmod}\left( \frac{\widetilde{w}_i^n - \widetilde{w}_{i-1}^n}{\Delta x}, \frac{\widetilde{w}_{i+1}^n - \widetilde{w}_i^n}{\Delta x} \right)$. Other limiters can be considered, see [23, 30] for instance. We enforce an additional limitation procedure on the first component $\sigma_i^{n,h}(\widetilde{w})$ in order to ensure that the fluid heights $h_i^{n, \pm}$ remain positive. An immediate computation shows that the condition

$$|\sigma_i^{n,h}(\widetilde{w})| \leq \frac{2h_i^n}{\Delta x},$$

is sufficient.

The standard MUSCL extension is obtained as follows. For a first-order scheme under the form [2.8], the second-order scheme is defined by

$$w_i^{n+1} = w_i^n - \frac{\Delta t}{\Delta x}(\mathcal{F}(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, d) - \mathcal{F}(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, d))$$
$$+ \frac{\Delta t}{2\Delta x}(\mathcal{S}(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, d) + \mathcal{S}_c(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, d) + \mathcal{S}(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, d)), \quad (4.3)$$

where $\mathcal{S}_c(\widetilde{w}_L, \widetilde{w}_R, d)$ is a centered source term, added to take into account the source term when its jumps at interfaces are small (see [8]). Several possibilities exist to define it.

In the present work, it will be determined naturally by considering the MUSCL scheme (4.3) as a convex combination between two first-order schemes applied on the reconstructed states and on half-cells, as described in Figure 4.1. More precisely, we define

$$w_i^{n+1,+} = w_i^{n,+} - \frac{2\Delta t}{\Delta x}\left(\mathcal{F}\left(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \frac{\Delta x}{2}\right) - \mathcal{F}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right)\right)$$
$$+ \frac{\Delta t}{\Delta x}\left(\mathcal{S}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right) + \mathcal{S}\left(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \frac{\Delta x}{2}\right)\right),$$

and

$$w_i^{n+1,-} = w_i^{n,-} - \frac{2\Delta t}{\Delta x}\left(\mathcal{F}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right) - \mathcal{F}\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \frac{\Delta x}{2}\right)\right)$$
$$+ \frac{\Delta t}{\Delta x}\left(\mathcal{S}\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \frac{\Delta x}{2}\right) + \mathcal{S}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right)\right).$$

Taking the arithmetic mean of these two states, we obtain

$$w_i^{n+1} = \frac{w_i^{n,-} + w_i^{n,+}}{2} - \frac{\Delta t}{\Delta x}\left(\mathcal{F}\left(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \frac{\Delta x}{2}\right) - \mathcal{F}\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \frac{\Delta x}{2}\right)\right)$$
$$+ \frac{\Delta t}{2\Delta x}\left(\mathcal{S}\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \frac{\Delta x}{2}\right) + 2\mathcal{S}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right) + \mathcal{S}\left(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \frac{\Delta x}{2}\right)\right). \quad (4.4)$$

Assuming the reconstruction is conservative, namely $w_i^n = \frac{w_i^{n,-} + w_i^{n,+}}{2}$, we notice that the scheme (4.4) can be written under the form (4.3) with $d = \frac{\Delta x}{2}$ and by defining the centered source term as

$$\mathcal{S}_c(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, d) = 2\mathcal{S}\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \frac{\Delta x}{2}\right).$$

An advantage of this procedure is that the MUSCL scheme (4.4) automatically preserves the positivity of $h$ as soon as the associated first-order scheme does, up to a half CFL restriction.

However, the well-balanced property is not reached as easily. Indeed, in order for the MUSCL scheme (4.4) to be well-balanced, the reconstruction would have to satisfy for any discrete steady state $(\widetilde{w}_i^n)_{i\in\mathbb{Z}}$,

$$\mathcal{E}(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \Delta x/2) = \mathcal{E}(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \Delta x/2) = 0, \text{ for all } i\in\mathbb{Z}.$$

Unfortunately, we cannot provide such a reconstruction in our case. Indeed, we have to reconstruct four variables, including the conservative ones $h, hu, hv$ which leaves only
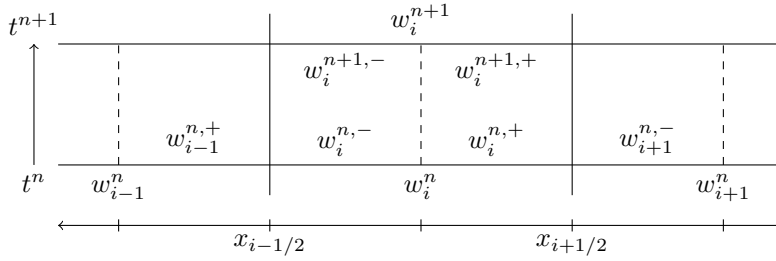
FIG. 4.1. *MUSCL second-order scheme*

one free variable to reconstruct. Moreover, according to Definition (2.11), the moving steady states involve three expressions among which two are not conservative quantities. Therefore, we would need to reconstruct two free variables in order to preserve steady states.

We propose in the next section, to modify the MUSCL method and the reconstruction to get around that problem and recover a fully well-balanced second-order scheme.

**4.2. Fully well-balanced recovering.** We suggest a modification based on an idea introduced in [28] and [15]. The main principle is to consider the second-order scheme (4.3) far from steady states and recover the first-order scheme (2.8) near a steady state, which guarantees the scheme to be well-balanced.

The difficulty lies in the definition of being far from/close to a steady state.

For this purpose, we consider a smooth increasing function $\theta$, valued in $[0,1]$ and such that $\theta(0)=0$ and $\theta(x)\approx 1$ far from 0. We choose the following function

$$\theta(x)=\frac{x^2}{x^2+\Delta x^2}.$$

We set $\theta_i^n=\theta(\mathcal{E}_i^n)$, where $\mathcal{E}_i^n=\mathcal{E}(\widetilde{w}_{i-1}^n,\widetilde{w}_i^n,\Delta x)+\mathcal{E}(\widetilde{w}_i^n,\widetilde{w}_{i+1}^n,\Delta x)$ detects if both couples $(\widetilde{w}_{i-1}^n,\widetilde{w}_i^n)$ and $(\widetilde{w}_i^n,\widetilde{w}_{i+1}^n)$ are local steady states simultaneously.

The reconstructed states are now defined as a convex combination between the linear reconstructed states and the first-order states

$$\widetilde{w}_i^{n,\pm}=(1-\theta_i^n)\widetilde{w}_i^n+\theta_i^n\left(\widetilde{w}_i^n\pm\frac{\Delta x}{2}\sigma_i^n(\widetilde{w})\right)=\widetilde{w}_i^n\pm\theta_i^n\frac{\Delta x}{2}\sigma_i^n(\widetilde{w}). \tag{4.5}$$

This reconstruction amounts to considering an additional limitation that involves the steady state detector $\theta_i^n$. For a discrete steady state, we have $\theta_i^n=0$ and we recover the first-order states $\widetilde{w}_i^{n,\pm}=\widetilde{w}_i^n$. Far from steady states and for smooth solutions, a mere computation shows that $\widetilde{w}_i^{n,\pm}=\widetilde{w}_i^n\pm\frac{\Delta x}{2}\sigma_i^n(\widetilde{w})+O(\Delta x^3)$ when $\Delta x$ tends to 0, which means the perturbation added to the usual second-order reconstruction is small enough to recover the seeking order.

Next, we define the scheme as

$$w_i^{n+1}=w_i^n-\frac{\Delta t}{\Delta x}\left(\mathcal{F}\left(\widetilde{w}_i^{n,+},\widetilde{w}_{i+1}^{n,-},\Delta x_1\right)-\mathcal{F}\left(\widetilde{w}_{i-1}^{n,+},\widetilde{w}_i^{n,-},\Delta x_1\right)\right)$$

$$+\frac{\Delta t}{2\Delta x}\left(\mathcal{S}\left(\widetilde{w}_{i-1}^{n,+},\widetilde{w}_i^{n,-},\Delta x_1\right)+2\mathcal{S}\left(\widetilde{w}_i^{n,-},\widetilde{w}_i^{n,+},\Delta x_2\right)+\mathcal{S}\left(\widetilde{w}_i^{n,+},\widetilde{w}_{i+1}^{n,-},\Delta x_1\right)\right), \tag{4.6}$$

where the coefficients $\Delta x_1$ and $\Delta x_2$ have to be adapted, depending if we apply the first or second-order scheme. Far from steady states we need

$$\Delta x_1 = \Delta x_2 = \frac{\Delta x}{2}, \tag{4.7}$$

in (4.6) to recover the second-order scheme (4.3). For a discrete steady state, the scheme reads as

$$w_i^{n+1} = w_i^n - \frac{\Delta t}{\Delta x} \left( \mathcal{F}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, \Delta x_1\right) - \mathcal{F}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, \Delta x_1\right)\right)$$
$$+ \frac{\Delta t}{2\Delta x}\left( \mathcal{S}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, \Delta x_1\right) + 2\mathcal{S}\left(\widetilde{w}_i^n, \widetilde{w}_i^n, \Delta x_2\right) + \mathcal{S}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, \Delta x_1\right)\right).$$

We notice that $\mathcal{S}(\widetilde{w}, \widetilde{w}, 0) = 0$ according to the source term consistency (2.4). Therefore, we have to set

$$\Delta x_1 = \Delta x \text{ and } \Delta x_2 = 0, \tag{4.8}$$

to recover the first-order scheme (2.8) at steady states.

In order to satisfy (4.7) far form steady states and (4.8) at steady states, coefficients $\Delta x_1$ and $\Delta x_2$ are set as convex combinations as follows

$$\Delta x_1 = \Delta x \left( 1 - \frac{\theta_i^n}{2}\right) \quad \text{and} \quad \Delta x_2 = \theta_i^n \frac{\Delta x}{2}. \tag{4.9}$$

We prove in the following theorem that the resulting second-order scheme is fully well-balanced, and that it preserves the positivity of $h$ under the classical second-order CFL restriction.

THEOREM 4.1.    *Under the CFL condition*

$$\frac{\Delta t}{\Delta x} \max_{i \in \mathbb{Z}} \left( |\lambda^\pm(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+})|, |\lambda^\pm(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-})| \right) \leq \frac{1}{4},$$

*and if the speeds $\lambda_L$ and $\lambda_R$ of the approximate Riemann solver satisfy the condition* (3.15), *then the second-order scheme* (4.5)-(4.6)-(4.9) *is fully well-balanced and preserves the positivity of $h$.*

*Proof.*    First, we consider a discrete steady state $(\widetilde{w}_i^n)_{i\in\mathbb{Z}}$. By definition, we have $\theta_i^n = 0$ for all $i \in \mathbb{Z}$. Hence, the scheme (4.5)-(4.6)-(4.9) gives

$$w_i^{n+1} = w_i^n - \frac{\Delta t}{\Delta x}\left( \mathcal{F}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, \Delta x\right) - \mathcal{F}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, \Delta x\right)\right)$$
$$+ \frac{\Delta t}{2\Delta x}\left( \mathcal{S}\left(\widetilde{w}_{i-1}^n, \widetilde{w}_i^n, \Delta x\right) + \mathcal{S}\left(\widetilde{w}_i^n, \widetilde{w}_{i+1}^n, \Delta x\right)\right), \quad (4.10)$$

which is nothing but the fully well-balanced first-order scheme (2.8).

Now we prove the positivity-preserving property. We assume that $h_i^n$ is positive for all $i \in \mathbb{Z}$. The update of variable $h$ with the scheme (4.5)-(4.6)-(4.9) writes

$$h_i^{n+1} = h_i^n - \frac{\Delta t}{\Delta x}\left( \mathcal{F}^h\left(\widetilde{w}_i^{n,+}, \widetilde{w}_{i+1}^{n,-}, \Delta x_1\right) - \mathcal{F}^h\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \Delta x_1\right)\right)$$
$$= \frac{1}{2}\left( h_i^{n,-} - \frac{\Delta t}{\Delta x/2}\left( \mathcal{F}^h\left(\widetilde{w}_i^{n,-}, \widetilde{w}_i^{n,+}, \Delta x_1\right) - \mathcal{F}^h\left(\widetilde{w}_{i-1}^{n,+}, \widetilde{w}_i^{n,-}, \Delta x_1\right)\right)\right)$$

$$+ \frac{1}{2}\left( h_i^{n,+} - \frac{\Delta t}{\Delta x/2}\left(\mathcal{F}^h\left(\widetilde{w}_i^{n,+},\widetilde{w}_{i+1}^{n,-},\Delta x_1\right) - \mathcal{F}^h\left(\widetilde{w}_i^{n,-},\widetilde{w}_i^{n,+},\Delta x_1\right)\right)\right).$$

Then $h_i^{n+1}$ is a convex combination between first-order schemes applied on half-cells with parameter $d = \Delta x_1$. As proved in Lemma 4.1, the first-order scheme preserves the positivity of $h$ independently of the value of the parameter $d$. Therefore, we conclude $h_i^{n+1} > 0$ for all $i \in \mathbb{Z}$. □

## 5. Numerical results

This section is devoted to numerical experiments. For the sake of simplicity, the initial discretisation will be defined as

$$w_i^0 = w_0(x_i).$$

Considering a continuous steady solution, the initial discretisation can satisfy exactly the Definition 2.5 of the discrete steady states. In this case, both our first-order and second-order schemes were proved to preserve the initial condition. This will be illustrated in Subsection 5.1.

However, it is also possible that the initial discretisation of a continuous steady state does not lead to a discrete steady state according to Definition 2.5. The behaviour of our numerical schemes in such a case will be investigated in Subsection 5.2. In order to measure how close a given discretisation $(w_i^n)_{i\in\mathbb{Z}}$ at time $t^n$ is to a discrete steady state, we will use the steady state distance

$$\mathcal{E}_{\infty,j}^n = \max_{1\leq i\leq N}\mathcal{E}(w_i^n,w_{i+1}^n),$$

where $j = 1$ for the first-order scheme and $j = 2$ for the second-order scheme.

In Subsection 5.3, we test the long-time convergence towards a steady state on a topography with a bump, using the same distance $\mathcal{E}_{\infty,j}^n$.

Finally, in Subsection 5.4, we consider a particular solution constant in space, but not in time, and for which we compute the errors in space and time.

**5.1. Moving steady state.** We consider here a simple moving steady state. As initial data, we take (see Figure 5.1)

$$h_0(x) = \exp^{2x}, \ u_0(x) = \exp^{-2x} \text{ and } v_0(x) = -fx,$$

and the topography is given by

$$z(x) = -\frac{1}{2}f^2x^2 - \exp^{2x} - \frac{1}{2}\exp^{-4x}.$$

We compute this test on the domain $[0,1]$ with $N = 200$ cells and the parameters $f = g = 1$.

The initial discretisation is a discrete steady state in the sense of Definition 2.5. Indeed the steady state distance at time $t_0 = 0$ is

$$\mathcal{E}_{\infty,1}^0 = \mathcal{E}_{\infty,2}^0 = 8.87\times 10^{-16}.$$

At final time $T_{\max} = 0.5$, the steady state is still preserved by both first-order and second-order schemes, even if small computationnal errors have spread. Indeed, the computation of the steady state distance at the end of the simulations gives

$$\mathcal{E}_{\infty,1}^{T_{\max}} = 5.19\times 10^{-14} \quad \text{and} \quad \mathcal{E}_{\infty,2}^{T_{\max}} = 8.86\times 10^{-15}.$$
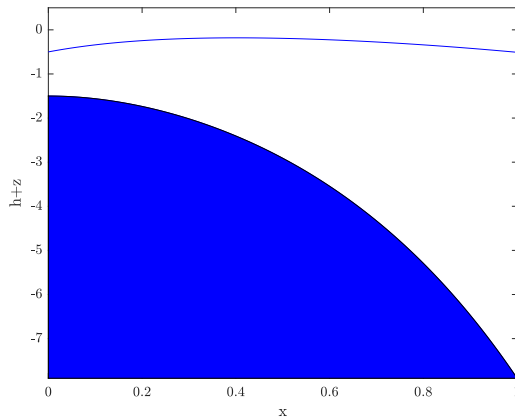
FIG. 5.1. *Initial variable h+z for the moving steady state.*

**5.2. Geostrophic steady state.** Next, we test the numerical schemes on another geostrophic steady state introduced in [11]. The computational domain is $[-5,5]$ with a flat topography $(z \equiv 0)$. We set $f = 10$, $g = 1$, and we consider

$$h_0(x) = \frac{2}{g} - e^{-x^2}, \quad u_0(x) = 0, \quad v_0(x) = \frac{2g}{f} x e^{-x^2},$$

as initial condition, which is a continuous steady state according to (1.4), see Figure 5.2. However, the initial data discretisation is not exactly a discrete steady state because the second relation in Definition 2.3 is not satisfied at each interface. This is confirmed numerically since we have for $N = 200$ discretisation points

$$\mathcal{E}_{\infty,1}^0 = \mathcal{E}_{\infty,2}^0 = 4.06 \times 10^{-5}.$$

Therefore, Theorems 3.1 and 4.1 do not guarantee the behaviour of the numerical schemes on this test case. We present in Table 5.1 the steady state distance at different times. We can observe it slowly decreases through time for both schemes, which means that the schemes slowly converge to a steady solution.
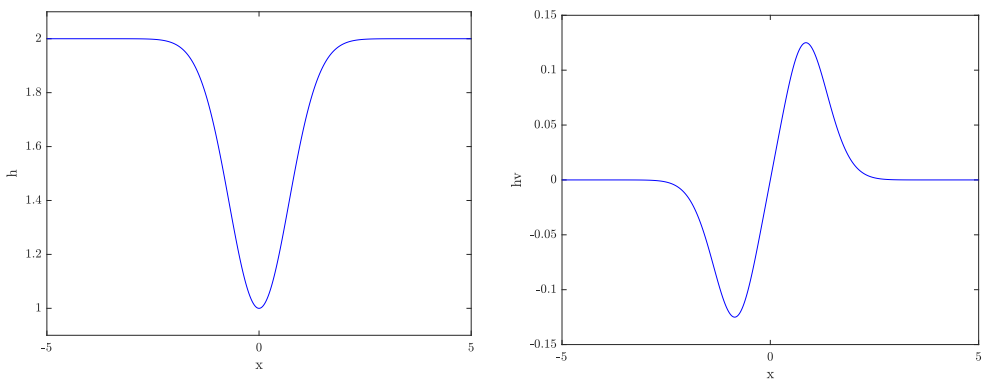


FIG. 5.2. *Initial data for the geostrophic steady state.*

| $t^n$ | $\mathcal{E}_{\infty,1}^n$ | $\mathcal{E}_{\infty,2}^n$ |
|---|---|---|
| 200 | $1.67 \times 10^{-4}$ | $1.67 \times 10^{-4}$ |
| 1000 | $1.15 \times 10^{-4}$ | $1.14 \times 10^{-4}$ |
| 2000 | $2.70 \times 10^{-5}$ | $2.68 \times 10^{-5}$ |
| 5000 | $3.14 \times 10^{-5}$ | $3.31 \times 10^{-5}$ |
| 10000 | $3.17 \times 10^{-5}$ | $3.26 \times 10^{-5}$ |
| 50000 | $9.62 \times 10^{-6}$ | $9.38 \times 10^{-6}$ |

TABLE 5.1. *Steady state distance computed for the gesotrophic steady state at different times.*

(a) first-order scheme

| N | h | | hv | |
|---|---|---|---|---|
| 200 | $6.15 \times 10^{-1}$ | | $1.89 \times 10^{-1}$ | |
| 400 | $4.01 \times 10^{-1}$ | 0.618 | $1.35 \times 10^{-1}$ | 0.485 |
| 800 | $2.43 \times 10^{-1}$ | 0.720 | $8.81 \times 10^{-2}$ | 0.615 |
| 1600 | $1.38 \times 10^{-1}$ | 0.820 | $5.31 \times 10^{-2}$ | 0.730 |
| 3200 | $7.50 \times 10^{-2}$ | 0.877 | $3.09 \times 10^{-2}$ | 0.780 |
| 6400 | $3.99 \times 10^{-2}$ | 0.910 | $1.79 \times 10^{-2}$ | 0.786 |
| 12800 | $2.09 \times 10^{-2}$ | 0.937 | $1.00 \times 10^{-2}$ | 0.838 |

(b) second-order scheme

| N | h | | hv | |
|---|---|---|---|---|
| 200 | $6.16 \times 10^{-1}$ | | $1.89 \times 10^{-1}$ | |
| 400 | $4.02 \times 10^{-1}$ | 0.614 | $1.35 \times 10^{-1}$ | 0.483 |
| 800 | $2.45 \times 10^{-1}$ | 0.713 | $8.88 \times 10^{-2}$ | 0.608 |
| 1600 | $1.40 \times 10^{-1}$ | 0.813 | $5.39 \times 10^{-2}$ | 0.722 |
| 3200 | $7.63 \times 10^{-2}$ | 0.873 | $3.16 \times 10^{-2}$ | 0.768 |
| 6400 | $4.06 \times 10^{-2}$ | 0.908 | $1.85 \times 10^{-2}$ | 0.775 |
| 12800 | $2.11 \times 10^{-2}$ | 0.944 | $1.03 \times 10^{-2}$ | 0.843 |

TABLE 5.2. $L^1$ *error in space for the geostrophic steady state at time* $T_{\max} = 200$ *for first and second-order schemes.*

Let us now check the convergence of both schemes when $\Delta x$ tends to 0. We define the $L_1$ discrete error in space at time $t^n$ between the exact solution $w_0$ and the numerical approximation by

$$E^n = \Delta x \sum_{i=1}^{N} |w_0(x_i) - w_i^n|.$$

We present in Table 5.2 the discrete errors for variables $h$ and $hv$ at final time $T_{\max} = 200$ for the first-order and second-order schemes. We observe that order two is not reached by the second-order scheme. This can be explained by the fact that the initial discretisation is close to a discrete steady state. Therefore, the parameter $\theta_i^n$ in (4.5) is close to zero and the reconstructed states are close to the first-order state. Let us emphasize that far away from steady states, the second-order scheme actually reaches second-order precision, as will be illustrated in Section 5.4.
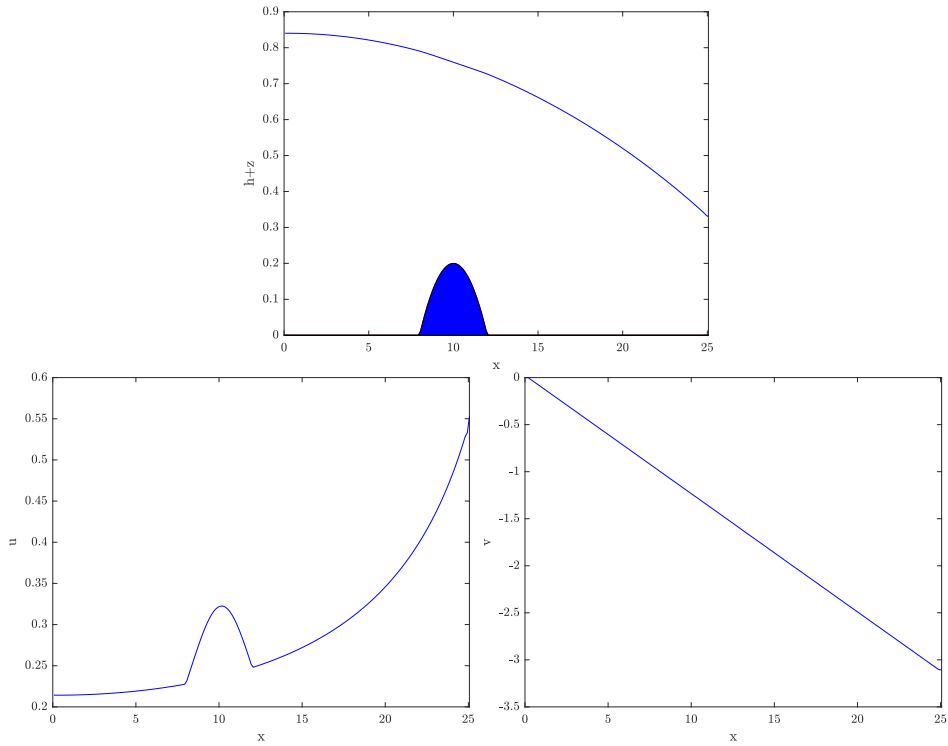
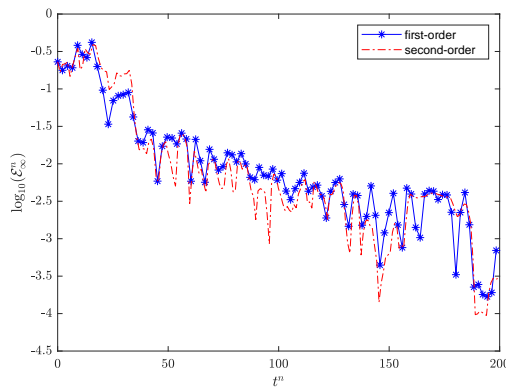FIG. 5.3. *Approximate solution of the steady flow over a bump test case at time* $T_{\max} = 200$.



FIG. 5.4. *Steady flow over a bump, steady state distance* $\mathcal{E}_\infty^n$ *in logarithmic scale.*

**5.3. Convergence towards a steady flow over a bump.** This test case aims to study the convergence towards a steady flow over a bump. It is a classical test for the shallow water equations adapted with the Coriolis source term in [31]. The topography is given by

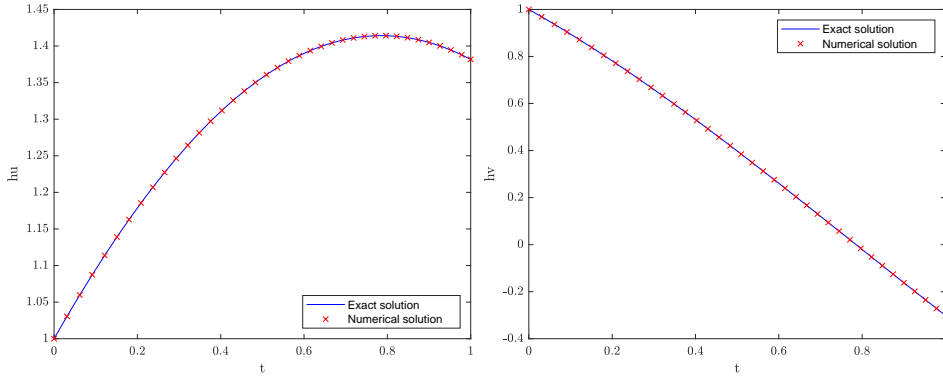$$z(x) = \begin{cases} 0.2 - 0.05(x-10)^2 & \text{if } 8 < x < 12, \\ 0 & \text{otherwise.} \end{cases}$$

FIG. 5.5. *Stationary state in space at time $T_{\max} = 1$*

(a) first-order scheme

| $N$ | hu | | hv | |
|---|---|---|---|---|
| 200 | $7.57 \times 10^{-4}$ | | $1.64 \times 10^{-4}$ | |
| 400 | $3.77 \times 10^{-4}$ | 1.00 | $8.21 \times 10^{-5}$ | 1.00 |
| 800 | $1.88 \times 10^{-4}$ | 1.00 | $4.10 \times 10^{-5}$ | 1.00 |
| 1600 | $9.42 \times 10^{-5}$ | 1.00 | $2.05 \times 10^{-5}$ | 1.00 |
| 3200 | $4.71 \times 10^{-5}$ | 1.00 | $1.03 \times 10^{-5}$ | 1.00 |
| 6400 | $2.35 \times 10^{-5}$ | 1.00 | $5.13 \times 10^{-6}$ | 1.00 |

(b) second-order scheme

| $N$ | hu | | hv | |
|---|---|---|---|---|
| 200 | $1.50 \times 10^{-8}$ | | $6.89 \times 10^{-8}$ | |
| 400 | $3.74 \times 10^{-9}$ | 2.00 | $1.72 \times 10^{-8}$ | 2.00 |
| 800 | $9.35 \times 10^{-10}$ | 2.00 | $4.29 \times 10^{-9}$ | 2.00 |
| 1600 | $2.34 \times 10^{-10}$ | 2.00 | $1.07 \times 10^{-9}$ | 2.00 |
| 3200 | $5.84 \times 10^{-11}$ | 1.99 | $2.68 \times 10^{-10}$ | 2.00 |
| 6400 | $1.46 \times 10^{-11}$ | 1.99 | $6.70 \times 10^{-11}$ | 2.00 |

TABLE 5.3. *$L_1$ error in space for the stationary in space test case computed at time $T_{\max} = 1$.*

We consider the following initial data

$$h_0(x) = 0.33, \quad u_0(x) = 0.18/0.33, \quad v_0(x) = 0.$$

We compute the solution on the domain $[0,25]$ with $N = 200$ cells and we set $f = \frac{2\pi}{50}$ and $g = 9.81$. The boundary conditions are set as

$$(hu)(x=0) = 0.18, \quad h(x=25) = 0.33, \quad v(x=0) = 0.$$

The numerical solution at time $T_{\max} = 200$ is represented in Figure 5.3. The time evolution of the steady state distance $\mathcal{E}_{\infty,j}^n$ is shown for both schemes in Figure 5.4. We can see these distances diminishing through time, which means both schemes actually converge towards a steady state.

(a) first-order scheme

| $N$ | hu | | hv | |
|---|---|---|---|---|
| 200 | $3.82 \times 10^{-4}$ | | $8.06 \times 10^{-5}$ | |
| 400 | $1.91 \times 10^{-4}$ | 0.99 | $4.03 \times 10^{-5}$ | 0.99 |
| 800 | $9.56 \times 10^{-5}$ | 0.99 | $2.01 \times 10^{-5}$ | 0.99 |
| 1600 | $4.78 \times 10^{-5}$ | 0.99 | $1.01 \times 10^{-5}$ | 0.99 |
| 3200 | $2.39 \times 10^{-5}$ | 0.99 | $5.04 \times 10^{-6}$ | 0.99 |
| 6400 | $1.20 \times 10^{-5}$ | 0.99 | $2.52 \times 10^{-6}$ | 0.99 |

(b) second-order scheme

| $N$ | hu | | hv | |
|---|---|---|---|---|
| 200 | $7.71 \times 10^{-9}$ | | $3.58 \times 10^{-8}$ | |
| 400 | $1.92 \times 10^{-9}$ | 1.99 | $8.95 \times 10^{-9}$ | 1.99 |
| 800 | $4.82 \times 10^{-10}$ | 1.99 | $2.24 \times 10^{-9}$ | 1.99 |
| 1600 | $1.20 \times 10^{-10}$ | 2.00 | $5.60 \times 10^{-10}$ | 1.99 |
| 3200 | $3.01 \times 10^{-11}$ | 2.00 | $1.40 \times 10^{-10}$ | 1.99 |
| 6400 | $7.52 \times 10^{-12}$ | 2.00 | $3.50 \times 10^{-11}$ | 1.99 |

TABLE 5.4. $L_1$ error in time for the stationary in space test case computed from time $t_0 = 0$ until time $T_{\max} = 1$.

**5.4. Stationary state in space.**     This test case is based on a particular exact solution of the RSW equations without topography. For a constant initial condition $(h_0, u_0, v_0)$ fixed, the exact solution of RSW equations writes

$$h(x,t) = h_0,$$
$$u(t) = u_0 \cos(ft) + v_0 \sin(ft),$$
$$v(t) = v_0 \cos(ft) - u_0 \sin(ft).$$

For any fixed time $t \geq 0$, the solution remains constant in space. We compute the scheme on domain $[0,1]$ until time $T_{\max} = 1$. We choose

$$h_0 = 1, \quad u_0 = 1, \quad v_0 = 1$$

as initial data, with the parameters $f = g = 1$ and we use periodic boundary conditions.

The solution is well-captured by the scheme as one can see in Figure 5.5, where we represent $hu$ and $hv$ with respect to time. According to Table 5.3 that shows the discrete $L_1$ error in space $E^{T_{\max}}$, both schemes reach the expected accuracy in space. Since the exact solution is known and constant in space, we can also check the scheme's accuracy in time. We introduce the discrete $L^1$ error in time between the exact solution $w_{ex}$ and the numerical approximation at point $x_i$

$$E_i = \sum_n (t^{n+1} - t^n) |w_{ex}(x_i, t^n) - w_i^n|.$$

Let us notice that the choice of the point $x_i$ is irrelevant since the solution is constant in space. We recover the expected order of accuracy in time as one can see in error Table 5.4.

## 6. Conclusions

In this work, we have built a second-order fully well-balanced scheme for the RSW system. In the first part, we have developed a fully well-balanced approximate Riemann solver by selecting carefully the numerical source term definitions and the relations used to define the intermediate states $w_L^\star$ and $w_R^\star$. The positivity of the variable $h$ has been recovered thanks to a cut-off procedure. We have proved in Theorem 3.1 that the resulting Godunov-type scheme satisfies all of the required features: consistency, positivity-preserving and fully well-balanced property.

In the second part, we have proposed a way to extend the Godunov-type scheme to second-order. We have explained the limitations of the classical MUSCL method in view of the fully well-balanced property in the case of the RSW equations. Then we have adapted an idea proposed by [15, 27], which consists in getting the standard MUSCL second-order scheme far from steady states and recovering the first-order fully well-balanced scheme near steady states. That procedure preserves the positivity of $h$ as proved in Theorem 4.1.

Finally, we have presented some numerical experiments that illustrate the robustness and the efficiency of both first-order and second-order schemes.

This work can be easily extended to the two-dimensional RSW equations by involving a standard convex combination of 1D schemes by interface. Additionally, the Coriolis parameter has been assumed constant all along this paper. It would be an interesting development of this work to consider a space-dependent Coriolis force, since it would be more realistic for large-scale simulations.

REFERENCES

[1] 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(6):2050–2065, 2004. 1

[2] E. Audusse, C. Chalons, and P. Ung, *A simple well-balanced and positive numerical scheme for the shallow-water system*, Commun. Math. Sci., 13(5):1317–1332, 2015. 3.2

[3] E. Audusse, S. Dellacherie, M.H. Do, P. Omnes, and Y. Penel, *Godunov type scheme for the linear wave equation with Coriolis source term*, ESAIM Proc. Surveys, 58:1–26, 2017. 1

[4] E. Audusse, R. Klein, D.D. Nguyen, and S. Vater, *Preservation of the discrete geostrophic equilibrium in shallow water flows*, in J. Fořt, J. Fürst, J. Halama, R. Herbin, and F. Hubert (eds.), Finite Volumes for Complex Applications VI Problems & Perspectives, Springer, 59–67, 2011. 1

[5] A. Bermudez and M.E. Vazquez, *Upwind methods for hyperbolic conservation laws with source terms*, Comput. Fluids, 23(8):1049–1071, 1994. 1

[6] C. Berthon and C. Chalons, *A fully well-balanced, positive and entropy-satisfying G-type method for the shallow-water equations*, Math. Comput., 85(299):1281–1307, 2016. 1, 2.1, 2.1

[7] C. Berthon, *Stability of the MUSCL schemes for the Euler equations*, Commun. Math. Sci., 3(2):133–157, 2005. 4

[8] F. Bouchut, *Nonlinear Stability of Finite Volume Methods for Hyperbolic Conservation Laws: and Well-Balanced Schemes for Sources*, Springer Science & Business Media, 2004. 1, 2.1, 2.2, 4.1

[9] F. Bouchut, J. Le Sommer, and V. Zeitlin, *Frontal geostrophic adjustment and nonlinear wave phenomena in one-dimensional rotating shallow water. Part 2. High-resolution numerical simulations*, J. Fluid Mech., 514:35–63, 2004. 1, 1

[10] M.J. Castro, A. Pardo Milanés, and C. Parés, *Well-balanced numerical schemes based on a generalized hydrostatic reconstruction technique*, Math. Model. Meth. Appl. Sci., 17(12):2055–2113, 2007. 1

[11] A. Chertock, M. Dudzinski, A. Kurganov, and M. Lukáčová-Medvid'ová, *Well-balanced schemes*

*for the shallow water equations with Coriolis forces*, Numer. Math., 138(4):939–973, 2018. 1, 1, 5.2

[12] V. Desveaux, M.Zenk, C. Berthon, and C. Klingenberg, *Well-balanced schemes to capture non-explicit steady states: Ripa model*, Math. Comput., 85(300):1571–1602, 2016. 2.1

[13] E.D. Fernandez-Nieto, D. Bresch, and J. Monnier, *A consistent intermediate wave speed for a well-balanced HLLC solver*, C.R. Math., 346(13-14):795–800, 2008. 1

[14] U.S. Fjordholm, S. Mishra, and E. Tadmor, *Well-balanced and energy stable schemes for the shallow water equations with discontinuous topography*, J. Comput. Phys., 230(14):5587–5609, 2011. 1

[15] B. Ghitti, C. Berthon, M.H. Le, and E.F. Toro, *A fully well-balanced scheme for the 1D blood flow equations with friction source term*, J. Comput. Phys., 421:109750, 2020. 1, 1, 3, 3.2, 4, 4.2, 6

[16] L. Gosse, *A well-balanced flux-vector splitting scheme designed for hyperbolic systems of conservation laws with source terms*, Comput. Math. Appl., 39(9-10):135–159, 2000. 1

[17] S. Gottlieb and C.-W. Shu, *Total variation diminishing Runge-Kutta schemes*, Math. Comput., 67(221):73–85, 1998. 4

[18] E. Gouzien, N. Lahaye, V. Zeitlin, and T. Dubos, *Thermal instability in rotating shallow water with horizontal temperature/density gradients*, Phys. Fluids, 29(10):101702, 2017. 1

[19] 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):1–16, 1996. 1

[20] A. Harten, P.D. Lax, and B. van Leer, *On upstream differencing and Godunov-type schemes for hyperbolic conservation laws*, SIAM Rev., 25(1):35–61, 1983. 2.1, 3.2

[21] S. Jin, *A steady-state capturing method for hyperbolic systems with geometrical source terms*, ESAIM: Math. Model. Numer. Anal., 35(04):631–645, 2001. 1

[22] N. Lahaye, *Dynamique, interactions et instabilités de structures cohérentes agéostrophiques dans les modèles en eau peu profonde*, PhD thesis, Université Pierre et Marie Curie, Sorbonne Université, 2014. 1

[23] R.J. LeVeque, *Finite Volume Methods for Hyperbolic Problems*, Cambridge Texts in Applied Mathematics, Cambridge University Press, Cambridge, 2004. 4.1

[24] Q. Liang and F. Marche, *Numerical resolution of well-balanced shallow water equations with complex source terms*, Adv. Water Resour., 32(6):873–884, 2009. 1

[25] X. Liu, A. Chertock, and A. Kurganov, *An asymptotic preserving scheme for the two-dimensional shallow water equations with Coriolis forces*, J. Comput. Phys., 391:259–279, 2019. 1

[26] M. Lukáčová-Medvid'ová, S. Noelle, and M. Kraft, *Well-balanced finite volume evolution Galerkin methods for the shallow water equations*, J. Comput. Phys., 221(1):122–147, 2007. 1, 1

[27] V. Michel-Dansac, C. Berthon, S. Clain, and F. Foucher, *A well-balanced scheme for the shallow-water equations with topography*, Comput. Math. Appl., 72(3):568–593, 2016. 1, 3, 3.2, 6

[28] V. Michel-Dansac, C. Berthon, S. Clain, and F. Foucher, *A well-balanced scheme for the shallow-water equations with topography or Manning friction*, J. Comput. Phys., 335:115–154, 2017. 1, 1, 3, 4, 4.2

[29] C.-W. Shu and S. Osher, *Efficient implementation of essentially non-oscillatory shock-capturing schemes*, J. Comput. Phys., 77(2):439–471, 1988. 4

[30] E.F. Toro, *Riemann Solvers and Numerical Methods for Fluid Dynamics: A Practical Introduction*, Springer Science & Business Media, 2013. 4.1

[31] V. Zetilin, G.M. Reznik, S.B. Medvedev, F. Bouchut, and A. Stegner, *Nonlinear Dynamics of Rotating Shallow Water Methods and Advances*, Edited Series on Advances in Nonlinear Science and Complexity, Elsevier, 2, 2007. 1, 5.3