

Staggered and well-balanced discretization of shallow-water equations

Nguyen Quan Ba Hong

May 9, 2019

Abstract

In this context, we consider a class of finite volume schemes for the shallow water equations with variable bottom topography.

Contents

1	Introduction	2
2	A staggered scheme in 1-D by Doyen and Gunawan	4
3	Three hydrostatic reconstruction schemes based on subcell reconstructions	10
3.1	The original HR method	11
3.2	The HR method of Morales et al.	12
3.3	The third HR method	13
3.4	The numerical flux	13
3.5	Interpretation via subcell reconstructions	14
3.5.1	Splitting the cells into subcells	15
3.5.2	Reconstruction of the bottom $b_\varepsilon(x)$	15
3.5.3	Infinitesimal HR	15
3.5.4	Fluxes and source terms based on subcell reconstructions	19
3.6	Comparison of the HR schemes	20
3.7	Stability analysis	20
4	Offset equilibrium schemes	21

We follow the notations in [4].

1 Introduction

The shallow water equations are a nonlinear hyperbolic system of conservation laws with a source term due to the variable bottom topography.

Balance laws often consist of the conservation laws for the vector $U(x, t)$ of *mass* and *momentum*, accelerated by *conservation advection* and *pressure forces* (denoted by $-\frac{\partial F(U)}{\partial x}$ in (1.1) below), and by additional nonconservative forces $S(U, x)$, also called *source terms*.

The *equations of motion* may be written as

$$\frac{\partial U}{\partial t} + \frac{\partial F(U)}{\partial x} = S(U, x). \quad (1.1)$$

In this context, we consider the case where there is no source term in the equation of mass, so S can be written as $S = (0, s)^T$.

The *residuum* R is defined as

$$R(x, t) := -\frac{\partial F(U)}{\partial x} + S(U, x), \quad (1.2)$$

which indicates *near-equilibrium flows* when it nearly vanishes.

A *semidiscrete, first order accurate finite volume scheme* (abbr., se-dis 1st FVM) may be written as a method of lines

$$\frac{d}{dt}U_i(t) = R_i(t) := -\frac{1}{\Delta x} \left(F_{i+\frac{1}{2}} - F_{i-\frac{1}{2}} \right) + S_i, \quad (1.3)$$

where U_i approximates the *cell average* over the i^{th} cell $C_i := [x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$ at time t , $R_i(t)$ is the cell average of the residuum, $\Delta x := x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$ is the *spatial grid size*, $F_{i\pm\frac{1}{2}}$ is a *conservative numerical flux function*, and S_i approximates the cell average of the source term.

In this context, we focus on one-dimensional (1-D) shallow water equations, given by

$$U = \begin{pmatrix} h \\ hu \end{pmatrix}, \quad F(U) = \begin{pmatrix} hu \\ hu^2 + \frac{1}{2}gh^2 \end{pmatrix}, \quad \text{and } S(U, x) = -\begin{pmatrix} 0 \\ ghb' \end{pmatrix}. \quad (1.4)$$

More explicitly,

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, & g = \text{const}, \\ \partial_t(hu) + \partial_x\left(hu^2 + \frac{1}{2}gh^2\right) = -ghb'(x), \end{cases} \quad (1.5)$$

Here $b(x)$ is the *bottom topography*¹, $h(x, t)$ the *water height*, $u(x, t)$ the *water velocity*, and $g = 9.8\text{m/s}^2$ the *gravitational acceleration*. Thus the source term S models the force of gravity tangential to a sloped bottom². The residuum can be rewritten as

$$R = -\partial_x \begin{pmatrix} hu \\ hu^2 \end{pmatrix} - gh\partial_x \begin{pmatrix} 0 \\ w \end{pmatrix}, \quad (1.6)$$

¹Here the bottom is assumed *rigid* under the effect of gravity, i.e., the bottom is time-independent $\partial_t b = 0$ on Γ_{bot} , see [4, pp. 2–3] for more details.

²The bottom topography b is assumed to be of class C_1 for simplicity, but well-balanced schemes, which preserves the lake at rest discretely, use either continuous or discontinuous topography, see [2, p. 760].

where $w := h + b$ is the *water level*.

For such a problem, where shocks can form in the solution, finite volume methods have proved to be very effective.

We recall two important *equilibria* stated in [2, p. 759]:

i) *still water*, where

$$u = 0 \text{ and } \partial_x w = 0, \text{ and} \quad (1.7)$$

ii) the *lake at rest*, which is still water together with dry boundaries:

$$u = 0 \text{ and } h\partial_x w = 0. \quad (1.8)$$

Thus, the lake at rest residuum combines the dry shore ($h = 0$) with the flat water level ($\partial_x w = 0$) in a single product, which suggests a natural splitting of the nonconservative product at the wet-dry front.

The goal of this context is to modify numerical schemes on staggered grids for a precise consideration of these equilibrium states.

These second equilibria make it possible to take into account the transitions between dry zones and wet areas, whereas the former, defined by relation (??), presuppose water everywhere. We hope, and in fact we see, that preserving exactly these states at the discrete level are also more accurate in the transient case.

We are interested in numerical schemes whose water height and velocity have shifted discretizations. In addition, we restrict ourselves to schemes that allow to capture shock waves. We consider N_x points (or N_x meshes), and the unknown notations from [5].

- water level: h_i^n , $i \in \{1, \dots, N_x\}$,
- velocity: $u_{i+\frac{1}{2}}$, $i \in \{0, \dots, N_x\}$,

and n denotes the number of the iteration in time.

In addition to equation (1.5), considered on $[0, T] \times (0, L)$, T being the final time and L the length of the domain which has dry shore at its both boundaries, we give ourselves an initial condition

$$(h, hu)(0, x) = (h_0, q_0)(x), \quad x \in (0, L),$$

and conditions at the boundaries

$$(hu)(t, 0) = (hu)(t, L) = 0, \quad t \in [0, T].$$

We further define the space step $\Delta x = \frac{L}{N_x}$ and the time step $\Delta t = \frac{T}{N_t}$ which will be subjected to a stability condition.

2 A staggered scheme in 1-D by Doyen and Gunawan

We reuse the notations in Sec. 2, [3, p. 229]: The left end, the center and the right end of the i -th cell are denoted by $x_{i-\frac{1}{2}}$, x_i and $x_{i+\frac{1}{2}}$, respectively. We set $C_i := (x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})$ for all $i \in \mathcal{M} := \{1, \dots, N_x\}$, $\mathcal{E}_{\text{int}} := \{1, \dots, N_x - 1\}$, $\mathcal{E}_b := \{0, N_x\}$, and $\mathcal{E} := \mathcal{E}_{\text{int}} \cup \mathcal{E}_b$. The water height h and the topography b are discretized at the center of the cells. The approximation of h at point x_i and at time t^n is denoted by h_i^n . The approximation of b at point x_i is denoted by b_i . The velocity u is discretized at the interfaces between the cells. The approximation of u at point $x_{i+\frac{1}{2}}$ and at time t^n is denoted by $u_{i+\frac{1}{2}}^n$.

The mass conservation equation is discretized with an explicit upwind scheme:

$$h_i^{n+1} = h_i^n - \frac{\Delta t}{\Delta x} \left(F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n \right), \quad \forall i \in \mathcal{M}, \quad (2.1)$$

where

$$F_{i+\frac{1}{2}}^n := h_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n, \quad \forall i \in \mathcal{M},$$

and where $h_{i+\frac{1}{2}}^n$ is calculated by an upwind shift according to the sign of $u_{i+\frac{1}{2}}^n$:

$$\forall i \in \mathcal{E}, \quad h_{i+\frac{1}{2}}^n = \begin{cases} h_i^n, & \text{if } u_{i+\frac{1}{2}}^n \geq 0, \\ h_{i+1}^n, & \text{otherwise.} \end{cases}$$

Proposition 1 (Conservation of the total water height). *The explicit upwind scheme (2.1) preserves the total water height.*

Proof. Define

$$\mathcal{H}_n := \Delta x \sum_{i=1}^{N_x} h_i^n, \quad \forall n \in \{1, \dots, N_t\},$$

we claim that $\mathcal{H}_{n+1} = \mathcal{H}_n$ for all $n \in \{0, \dots, N_t - 1\}$. Indeed,

$$\begin{aligned} \mathcal{H}_{n+1} &= \Delta x \sum_{i=1}^{N_x} h_i^{n+1} \\ &= \Delta x \sum_{i=1}^{N_x} \left[h_i^n - \frac{\Delta t}{\Delta x} \left(F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n \right) \right] \\ &= \Delta x \sum_{i=1}^{N_x} h_i^n - \Delta t \sum_{i=1}^{N_x} \left(F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n \right) \\ &= \mathcal{H}_n - \Delta t \left(F_{N_x+\frac{1}{2}}^n - F_{\frac{1}{2}}^n \right) \\ &= \mathcal{H}_n, \quad \forall n \in \{1, \dots, N_t - 1\}, \end{aligned}$$

and then $\mathcal{H}_n = \mathcal{H}_1$ for all $n \in \{1, \dots, N_t\}$. This completes our proof. \square

The momentum balance equation in (1.5) is discretized with explicit upwind fluxes for the convection term and implicit centered fluxes for the pressure term and topography term:

$$\bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} = \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} \left((h_{i+1}^{n+1})^2 - (h_i^{n+1})^2 \right) + g \bar{h}_{i+\frac{1}{2}}^{n+1} (b_{i+1} - b_i) \right], \quad \forall i \in \mathcal{E}_{\text{int}}, \quad (2.2)$$

where

$$\bar{h}_{i+\frac{1}{2}}^n := \frac{1}{2} (h_i^n + h_{i+1}^n), \quad \forall i \in \mathcal{E}_{\text{int}},$$

and

$$G_i^n = \frac{1}{2} u_i^n \left(F_{i-\frac{1}{2}}^n + F_{i+\frac{1}{2}}^n \right),$$

where u_i^n is calculated by an upwind shift according to the sign of $F_{i-\frac{1}{2}}^n + F_{i+\frac{1}{2}}^n$:

$$\forall i \in \mathcal{M}, \quad u_i^n = \begin{cases} u_{i-\frac{1}{2}}^n, & \text{if } F_{i-\frac{1}{2}}^n + F_{i+\frac{1}{2}}^n \geq 0, \\ u_{i+\frac{1}{2}}^n, & \text{otherwise.} \end{cases}$$

The discrete boundary conditions are

$$u_{i+\frac{1}{2}}^{n+1} = 0, \quad \forall i \in \mathcal{E}_b. \quad (2.3)$$

The computation of the discrete unknowns at each time step is completely explicit. First the discrete water heights $\{h_i^{n+1}\}$ are computed with (2.1), then the discrete velocities $\left\{u_{i+\frac{1}{2}}^{n+1}\right\}$ are computed with (2.2) (if $\bar{h}_{i+\frac{1}{2}}^{n+1} = 0$, by convention, $u_{i+\frac{1}{2}}^{n+1}$ is set to zero).

Proposition 2 (Preserved quantities). *This scheme conserves the mass, and, for a flat topography, the total momentum, provided the space step is small enough for the latter.*

Proof. We will prove the conservation properties for the following quantities:

- *Mass:* Recall in [4] that the mass is defined by

$$\mathcal{Z} := \int_{\mathbb{R}} \zeta(t, x) dx = \int_{\mathbb{R}} (h(t, x) + b(x) - H) dx,$$

this quantity is preserved (see [4, Sec. 3.1, p. 21]):

$$\frac{d}{dt} \mathcal{Z} = 0. \quad (2.4)$$

Now we prove this property in the discrete level. To do this, we define

$$\mathcal{Z}_n := \Delta x \sum_{i=1}^{N_x} (h_i^n + b_i - H), \quad \forall n \in \{0, \dots, N_t\}.$$

We claim that $\mathcal{Z}_{n+1} = \mathcal{Z}_n$ for all $n \in \{0, \dots, N_t - 1\}$. Indeed,

$$\begin{aligned}
\mathcal{Z}_{n+1} &= \Delta x \sum_{i=1}^{N_x} (h_i^{n+1} + b_i - H) \\
&= \Delta x \sum_{i=1}^{N_x} \left(h_i^n - \frac{\Delta t}{\Delta x} (F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n) + b_i - H \right) \\
&= \Delta x \sum_{i=1}^{N_x} (h_i^n + b_i - H) - \Delta t \left(\sum_{i=1}^{N_x} F_{i+\frac{1}{2}}^n - \sum_{i=1}^{N_x} F_{i-\frac{1}{2}}^n \right) \\
&= \mathcal{Z}_n - \Delta t \left(\sum_{i=1}^{N_x} F_{i+\frac{1}{2}}^n - \sum_{i=0}^{N_x-1} F_{i+\frac{1}{2}}^n \right) \\
&= \mathcal{Z}_n - \Delta t (F_{N_x+\frac{1}{2}}^n - F_{\frac{1}{2}}^n) \\
&= \mathcal{Z}_n, \quad \forall n \in \{0, \dots, N_t - 1\},
\end{aligned}$$

where the last equality is deduced from (2.3). We then have $\mathcal{Z}_n = \mathcal{Z}_0$, for all $n \in \{0, \dots, N_t\}$, i.e., \mathcal{Z}_n is independent in time. This is the discrete version of (2.4).

- *Total momentum*: Recall that the total momentum is defined by

$$\mathfrak{M} := \int_{\mathbb{R}} (hu)(t, x) dx,$$

this quantity is preserved (see [4, Sec 3.1, p.21]):

$$\frac{d}{dt} \mathfrak{M} = 0. \quad (2.5)$$

Similarly, we now look at this property in the discrete level. Define

$$\mathfrak{M}_n := \Delta x \sum_{i=1}^{N_x-1} \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n, \quad \forall n \in \{0, \dots, N_t\},$$

we claim that $\mathfrak{M}_{n+1} = \mathfrak{M}_n$ for all $n \in \{0, \dots, N_t - 1\}$ provided that the bottom topography is flat. Indeed,

$$\begin{aligned}
\mathfrak{M}_{n+1} &= \Delta x \sum_{i=1}^{N_x-1} \bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} \\
&= \Delta x \sum_{i=1}^{N_x-1} \left[\bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} \left((h_{i+1}^{n+1})^2 - (h_i^{n+1})^2 \right) \right] \right] \\
&= \Delta x \sum_{i=1}^{N_x-1} \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \Delta t \sum_{i=1}^{N_x-1} \left[G_{i+1}^n - G_i^n + \frac{g}{2} \left((h_{i+1}^{n+1})^2 - (h_i^{n+1})^2 \right) \right] \\
&= \mathfrak{M}_n - \Delta t \left[G_{N_x}^n - G_1^n + \frac{g}{2} \left((h_{N_x}^{n+1})^2 - (h_1^{n+1})^2 \right) \right]
\end{aligned}$$

$$\begin{aligned}
&= \mathfrak{M}_n - \Delta t \left[\frac{1}{2} u_{N_x}^n \left(F_{N_x-\frac{1}{2}}^n + F_{N_x+\frac{1}{2}}^n \right) - \frac{1}{2} u_1^n \left(F_{\frac{1}{2}}^n + F_{\frac{3}{2}}^n \right) \right. \\
&\quad \left. + \frac{g}{2} \left(\left(h_{N_x}^n - \frac{\Delta t}{\Delta x} \left(F_{N_x+\frac{1}{2}}^n - F_{N_x-\frac{1}{2}}^n \right) \right)^2 - \left(h_1^n - \frac{\Delta t}{\Delta x} \left(F_{\frac{3}{2}}^n - F_{\frac{1}{2}}^n \right) \right)^2 \right) \right] \\
&= \mathfrak{M}_n - \Delta t \left[\frac{1}{2} u_{N_x}^n F_{N_x-\frac{1}{2}}^n - \frac{1}{2} u_1^n F_{\frac{3}{2}}^n + \frac{g}{2} \left(\left(h_{N_x}^n + \frac{\Delta t}{\Delta x} F_{N_x-\frac{1}{2}}^n \right)^2 - \left(h_1^n - \frac{\Delta t}{\Delta x} F_{\frac{3}{2}}^n \right)^2 \right) \right] \\
&= \mathfrak{M}_n, \quad \forall n \in \{0, \dots, N_t - 1\},
\end{aligned}$$

where the last inequality is obtain by choosing N_x large enough for which

$$h_1^n = h_{\frac{3}{2}}^n = h_{N_x-\frac{1}{2}}^n = h_{N_x}^n = 0$$

holds. We then have $\mathfrak{M}_n = \mathfrak{M}_0$, for all $n \in \{0, \dots, N_t\}$, i.e., \mathfrak{M}_n is independent in time. This is the discrete version of (2.5).

This completes our proof. \square

We apply the Finite Volume Method for 1-D scalar conservation laws for the mass conservation equation

$$\partial_t h + \partial_x f(h) = 0, \quad \text{in } [0, T] \times (0, L),$$

where $f(h) = (hu)(x, t)$.

We consider the numerical flux $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined piecewisely by

$$g(h, k) := \begin{cases} hu_{i+\frac{1}{2}}^n, & \text{if } u_{i+\frac{1}{2}}^n \geq 0, \\ ku_{i+\frac{1}{2}}^n, & \text{otherwise,} \end{cases} \quad \text{on } C_i \times M_{i+1}.$$

This numerical flux g is not consistent but *monotone* $g(\nearrow, \searrow)$, i.e., g is non-decreasing w.r.t. its first variable and non-increasing w.r.t. its second variable.

Proposition 3. *If the Courant-Friedrichs-Levy-like (CFL-like) condition:*

$$\Delta t \leq \frac{\Delta x}{\left(-u_{i+\frac{1}{2}}^n \right)_- + \left(u_{i+\frac{1}{2}}^n \right)_+}, \quad \forall i \in \mathcal{M} \quad (2.6)$$

holds (where, for any $a \in \mathbb{R}$, $(a)_+ := \max(a, 0)$ and $(a)_- := \min(a, 0)$) then the numerical scheme (2.1) is monotone: Let $H : \mathbb{R}^3 \rightarrow \mathbb{R}$ defined by $h_i^{n+1} = H(h_{i-1}^n, h_i^n, h_{i+1}^n)$ then H is non-decreasing w.r.t. its three variables $H(\nearrow, \nearrow, \nearrow)$. Therefore

$$A \leq h_0 \leq B \Rightarrow A \leq h_\Delta \leq B \quad \text{a.e.,}$$

where

$$h_\Delta(t, x) := \sum_{i=1}^{N_x} \sum_{n=1}^{N_t} h_i^n \mathbf{1}_{[t^n, t^{n+1}] \times C_i}(t, x), \quad \text{in } [0, T] \times (0, L).$$

Proof. Define

$$H(h, k, l) := k - \frac{\Delta t}{\Delta x} (g(k, l) - g(h, k)), \quad \text{in } \mathbb{R}^3,$$

its first-order partial derivatives are given by

$$\begin{aligned} \partial_1 H &= \frac{\Delta t}{\Delta x} \partial_1 g(h, k) \geq 0, \\ \partial_2 H &= 1 - \frac{\Delta t}{\Delta x} (\partial_1 g(k, l) - \partial_2 g(h, k)), \\ \partial_3 H &= -\frac{\Delta t}{\Delta x} \partial_2 g(k, l) \geq 0, \end{aligned}$$

Notice that the weak first-order partial derivatives of g are given by

$$\partial_1 g(h, k) = \begin{cases} u_{i+\frac{1}{2}}^n, & \text{if } u_{i+\frac{1}{2}}^n \geq 0, \\ 0, & \text{otherwise,} \end{cases} \quad \text{on } C_i \times M_{i+1},$$

and

$$\partial_2 g(h, k) = \begin{cases} 0, & \text{if } u_{i+\frac{1}{2}}^n \geq 0, \\ u_{i+\frac{1}{2}}^n, & \text{otherwise,} \end{cases} \quad \text{on } C_i \times M_{i+1},$$

i.e.,

$$\begin{aligned} \partial_1 g(h, k) &= \sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)_+ \mathbf{1}_{M_i \times M_{i+1}}(h, k), \\ \partial_2 g(h, k) &= \sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)^- \mathbf{1}_{M_i \times M_{i+1}}(h, k), \end{aligned}$$

then we can estimate $\partial_2 H$ as follows,

$$\begin{aligned} \partial_2 H &= 1 - \frac{\Delta t}{\Delta x} (\partial_1 g(k, l) - \partial_2 g(h, k)) \\ &= 1 - \frac{\Delta t}{\Delta x} \left(\sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)_+ \mathbf{1}_{M_i \times M_{i+1}}(k, l) - \sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)^- \mathbf{1}_{M_i \times M_{i+1}}(h, k) \right) \\ &= 1 - \frac{\Delta t}{\Delta x} \left(\sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)_+ \mathbf{1}_{C_i}(k) \times \mathbf{1}_{M_{i+1}}(l) - \sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)^- \mathbf{1}_{C_i}(h) \times \mathbf{1}_{M_{i+1}}(k) \right) \\ &= 1 - \frac{\Delta t}{\Delta x} \left(\sum_{i=1}^{N_x-1} \left(u_{i+\frac{1}{2}}^n \right)_+ \mathbf{1}_{C_i}(k) \times \mathbf{1}_{M_{i+1}}(l) - \sum_{i=2}^{N_x} \left(u_{i-\frac{1}{2}}^n \right)^- \mathbf{1}_{M_{i-1}}(h) \times \mathbf{1}_{C_i}(k) \right) \\ &= 1 - \frac{\Delta t}{\Delta x} \left[\sum_{i=2}^{N_x-1} \left[\mathbf{1}_{C_i}(k) \times \left(\left(u_{i+\frac{1}{2}}^n \right)_+ \mathbf{1}_{M_{i+1}}(l) - \left(u_{i-\frac{1}{2}}^n \right)^- \mathbf{1}_{M_{i-1}}(h) \right) \right] \right. \\ &\quad \left. + \left(u_{\frac{3}{2}}^n \right)_+ \mathbf{1}_{M_1}(k) \times \mathbf{1}_{M_2}(l) - \left(u_{N_x-\frac{1}{2}}^n \right)^- \mathbf{1}_{M_{N_x-1}}(h) \times \mathbf{1}_{M_{N_x}}(k) \right]. \end{aligned}$$

Hence, if the CFL-like condition (2.6) holds, then $\partial_2 H \geq 0$. Thus, $H(\nearrow, \nearrow, \nearrow)$.

Now, we assume that $A \leq h_i^0 \leq B$, for all $i \in \{1, \dots, N_x\}$. Let us remark that

$$\forall \kappa \in \mathbb{R}, \quad H(\kappa, \kappa, \kappa) = \kappa - \frac{\Delta t}{\Delta x} (g(\kappa, \kappa) - g(\kappa, \kappa)) = \kappa.$$

As a consequence: if $h_0 \equiv \kappa$ then $h_\Delta \equiv \kappa$.

Now, if we take $\kappa = A$ and $\kappa = B$, we can obtain the L^∞ -stability by monotonicity: Assume that $\forall i \in \mathbb{Z}$, $A \leq h_i^n \leq B$, then

$$\begin{aligned} h_i^{n+1} &= H(h_{i-1}^n, h_i^n, h_{i+1}^n) \leq H(B, h_i^n, h_{i+1}^n) \leq H(B, B, h_{i+1}^n) \leq H(B, B, B) = B, \\ h_i^{n+1} &= H(h_{i-1}^n, h_i^n, h_{i+1}^n) \geq H(A, h_i^n, h_{i+1}^n) \geq H(A, A, h_{i+1}^n) \geq H(A, A, A) = A, \end{aligned}$$

This ends our proof. \square

Applying this proposition with $A = 0$, we obtain immediately the following direct consequence.

Corollary 1 (Nonnegativity of water height). *The water height remains nonnegative at time t^{n+1} under the CFL-like condition (2.6).*

The next proposition indicates the well-balanced property of the proposed numerical scheme.

Proposition 4 (Well-balanced property). *The steady states at rest are preserved under the numerical scheme (2.1)-(2.2), i.e., if $u_{i+\frac{1}{2}}^n = 0$ and $h_i^n + b_i = \text{const}$ for all $i \in \mathcal{M}$, then $u_{i+\frac{1}{2}}^{n+1} = 0$ and $h_i^{n+1} + b_i = \text{const}$ for all $i \in \mathcal{M}$.*

Proof. We recall (2.1) and rewrite (2.2) as

$$\begin{aligned} h_i^{n+1} &= h_i^n - \frac{\Delta t}{\Delta x} \left(h_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - h_{i-\frac{1}{2}}^n u_{i-\frac{1}{2}}^n \right), \quad \forall i \in \mathcal{M}, \\ \bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} &= \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} (h_i^{n+1} + h_{i+1}^{n+1}) ((h_{i+1}^{n+1} + b_{i+1}) - (h_i^{n+1} + b_i)) \right], \quad \forall i \in \mathcal{E}_{\text{int}}. \end{aligned}$$

Suppose that $u_{i+\frac{1}{2}}^n = 0$ and $h_i^n + b_i = \text{const}$ for all $i \in \mathcal{E}$, then these equations become

$$\begin{aligned} h_i^{n+1} &= h_i^n, \quad \forall i \in \mathcal{M}, \\ \frac{h_i^{n+1} + h_{i+1}^{n+1}}{2} u_{i+\frac{1}{2}}^{n+1} &= -\frac{\Delta t}{\Delta x} \left[\frac{g}{2} (h_i^{n+1} + h_{i+1}^{n+1}) ((h_{i+1}^{n+1} + b_{i+1}) - (h_i^{n+1} + b_i)) \right], \quad \forall i \in \mathcal{E}_{\text{int}}. \end{aligned} \quad (2.7)$$

As a direct consequence, $h_i^{n+1} + b_i = h_i^n + b_i = \text{const}$.

Plugging the former into the latter yields

$$\frac{h_i^n + h_{i+1}^n}{2} u_{i+\frac{1}{2}}^{n+1} = -\frac{\Delta t}{\Delta x} \left[\frac{g}{2} (h_i^n + h_{i+1}^n) ((h_{i+1}^n + b_{i+1}) - (h_i^n + b_i)) \right], \quad \forall i \in \mathcal{E}_{\text{int}}.$$

Since $h_i^n + b_i = \text{const}$ for all $i \in \mathcal{E}$, the last equation implies $u_{i+\frac{1}{2}}^{n+1} = 0$ ($h_i^n \geq 0$ and $h_{i+1}^n \geq 0$ by Corollary 1). This finishes our proof. \square

Consider the case which dry shore are permitted in our domain Ω , we decompose \mathcal{M} as $\mathcal{M} = \mathcal{M}_{\text{dry}}^n \cup \mathcal{M}_{\text{wet}}^n$. The *dry-component* \mathcal{M}_{dry} consists of all indices i such that $h_i^n = 0$ and the *wet-component* \mathcal{M}_{wet} consists of all indices i such that $h_i^n > 0$. This is reasonable because the water height at a particular cell obviously indicates whether that cell is “dry” or “wet”.

With allowing dry shore in our domain Ω , a well-balanced property still holds in this scenario, with a little modification.

Proposition 5 (Ill-balanced property with dry shore). *When the domain Ω contains dry shore, the above well-balanced property does not hold in general.*

Proof. Given n , since the domain Ω contains dry shore, there exists an index $i_0 \in \mathcal{M}$ such that $i_0 \in \mathcal{M}_{\text{dry}}$ and $i_0 + 1 \in \mathcal{M}_{\text{wet}}$, or $i_0 \in \mathcal{M}_{\text{wet}}$ and $i_0 + 1 \in \mathcal{M}_{\text{dry}}$. Without loss of generality, suppose that the former holds, i.e., $h_{i_0}^n = 0$ and $h_{i_0+1}^n > 0$, we follow the proof of Proposition 4 to arrive at $h_{i_0}^{n+1} = h_{i_0}^n = 0$, $h_{i_0+1}^{n+1} = h_{i_0+1}^n > 0$, and (2.7) gives us

$$\frac{h_{i_0+1}^n}{2} u_{i_0+\frac{1}{2}}^{n+1} = -\frac{g\Delta t}{2\Delta x} h_{i_0+1}^n (h_{i_0+1}^n + b_{i_0+1} - b_{i_0}).$$

So if $h_{i_0+1}^n + b_{i_0+1} < b_{i_0}$ (a figure of this case can be easily drawn) then this equality implies $u_{i_0+\frac{1}{2}}^{n+1} > 0$. Therefore, the previous well-balanced property falls in this case. \square

At each time step, the Courant number is defined by

$$\nu := \frac{\Delta t}{\Delta x} \max_{i \in \mathcal{M}} \left(\frac{|q_{i+\frac{1}{2}}^n + q_{i-\frac{1}{2}}^n|}{2h_i^n} + \sqrt{gh_i^n} \right).$$

The numerical simulations in [3, pp. 232–234] show that the staggered scheme is stable under the CFL condition $\nu < 1$.

3 Three hydrostatic reconstruction schemes based on subcell reconstructions

We briefly rewrite the three hydrostatic reconstruction schemes in [2] in our notations. All of them introduce reconstructed values $U_{i+\frac{1}{2}\pm}$ of the unknowns to the left and right of the interface $x_{i+\frac{1}{2}}$, and define the numerical flux via a Riemann solver \mathcal{F} :

$$F_{i+\frac{1}{2}} = \mathcal{F}\left(U_{i+\frac{1}{2}-}, U_{i+\frac{1}{2}+}\right). \quad (3.1)$$

Those HR schemes split the singular source term at the interface into a left and a right part, $S_{i+\frac{1}{2}-}$ and $S_{i+\frac{1}{2}+}$, and compute the source term in (1.3) as

$$S_i = S_{i-\frac{1}{2}+} + S_{i+\frac{1}{2}-} = \left(0, s_{i-\frac{1}{2}+}\right)^T + \left(0, s_{i+\frac{1}{2}-}\right)^T. \quad (3.2)$$

There are two main types of interface, depending on how the water covers the bottom to the left and right sides of the interfaces (see Fig. 1³):

- *Fully wet interface*: where the water level on each side is higher than the higher side of the bottom topography:

$$\min(w_i, w_{i+1}) > \max(b_i, b_{i+1}), \quad (3.3)$$

where b_i and w_i are the bottom topography and the water level, respectively, in cell C_i .

- *Partially wet interface*: where the water level on one side is equal to or below the topography on the other side:

$$\min(w_i, w_{i+1}) \leq \max(b_i, b_{i+1}). \quad (3.4)$$

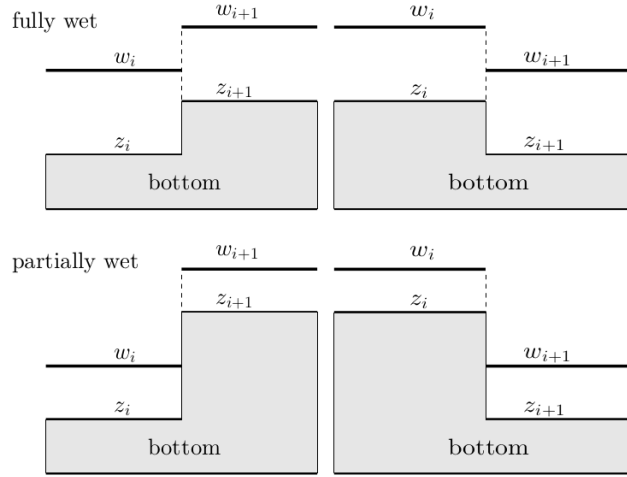


Figure 1: Examples of interfaces: fully wet (top), and partially wet (bottom).

After the water heights $h_{i-\frac{1}{2}+}$ and $h_{i+\frac{1}{2}-}$ are reconstructed in the following three HR schemes, the conservative variables are defined as

$$U_{i+\frac{1}{2}-} := \begin{pmatrix} h_{i+\frac{1}{2}-} \\ h_{i+\frac{1}{2}-} u_{i+\frac{1}{2}-} \end{pmatrix}, \quad U_{i-\frac{1}{2}+} := \begin{pmatrix} h_{i-\frac{1}{2}+} \\ h_{i-\frac{1}{2}+} u_{i-\frac{1}{2}+} \end{pmatrix} \quad \text{with } u_{i+\frac{1}{2}-} = u_{i-\frac{1}{2}+} = u_i. \quad (3.5)$$

3.1 The original HR method

Audusse et al. [1] introduce their first order HR scheme by choosing the *intermediate bottom* as

$$b_{i+\frac{1}{2}}^{\text{AUD}} := \max(b_i, b_{i+1}), \quad (3.6)$$

³This figure is extracted in [2, p. 761].

and the *interface water height* as

$$h_{i+\frac{1}{2}-}^{\text{AUD}} := \left(w_i - b_{i+\frac{1}{2}}^{\text{AUD}}\right)^+, \quad h_{i+\frac{1}{2}+}^{\text{AUD}} := \left(w_{i+1} - b_{i+\frac{1}{2}}^{\text{AUD}}\right)^+. \quad (3.7)$$

The source term is discretized as

$$s_{i+\frac{1}{2}-}^{\text{AUD}} := \frac{g}{2\Delta x} \left(\left(h_{i+\frac{1}{2}-}^{\text{AUD}}\right)^2 - h_i^2 \right), \quad (3.8)$$

$$s_{i+\frac{1}{2}+}^{\text{AUD}} := \frac{g}{2\Delta x} \left(h_{i+1}^2 - \left(h_{i+\frac{1}{2}+}^{\text{AUD}}\right)^2 \right). \quad (3.9)$$

3.2 The HR method of Morales et al.

The HR scheme of Morales et al. [6] is identical to that of Audusse's scheme,

$$b_{i+\frac{1}{2}}^{\text{MOR}} := b_{i+\frac{1}{2}}^{\text{AUD}}, \quad h_{i+\frac{1}{2}-}^{\text{MOR}} := h_{i+\frac{1}{2}-}^{\text{AUD}}, \quad h_{i+\frac{1}{2}+}^{\text{MOR}} := h_{i+\frac{1}{2}+}^{\text{AUD}}, \quad (3.10)$$

and the source term is defined as

$$s_{i+\frac{1}{2}-}^{\text{MOR}} := s_{i+\frac{1}{2}-}^{\text{AUD}}, \quad s_{i+\frac{1}{2}+}^{\text{MOR}} := s_{i+\frac{1}{2}+}^{\text{AUD}}, \quad (3.11)$$

except for the partially wet interfaces (3.4), where water either flows downhill or flows uphill with enough kinetic energy to climb the jump of the bottom at the interface. This results in the following two cases.

- $b_i < b_{i+1}$ (*ascending bottom*). If $u_i < 0$, or $u_i > 0$ and

$$\frac{|u_i|^2}{2} + g(w_i - b_{i+1}) \geq \frac{3}{2} \sqrt{g(h_i |u_i|)^3}, \quad (3.12)$$

then the left interface source term is redefined as

$$s_{i+\frac{1}{2}-}^{\text{MOR}} := -\frac{g}{\Delta x} \frac{h_i}{2} \left(b_{i+\frac{1}{2}}^{\text{MOR}} - b_i \right). \quad (3.13)$$

- $b_i > b_{i+1}$ (*descending bottom*). If $u_{i+1} > 0$, or $u_{i+1} < 0$ and

$$\frac{|u_{i+1}|^2}{2} + g(w_{i+1} - b_i) \geq \frac{3}{2} \sqrt{g(h_{i+1} |u_{i+1}|)^3}, \quad (3.14)$$

then the right interface source term is redefined as

$$s_{i+\frac{1}{2}+}^{\text{MOR}} := -\frac{g}{\Delta x} \frac{h_{i+1}}{2} \left(b_{i+1} - b_{i+\frac{1}{2}}^{\text{MOR}} \right). \quad (3.15)$$

3.3 The third HR method

For the third HR method, the intermediate bottom is defined as

$$b_{i+\frac{1}{2}}^{\text{CN}} := \min(\max(b_i, b_{i+1}), \min(w_i, w_{i+1})). \quad (3.16)$$

The interface water heights are given by

$$h_{i+\frac{1}{2}-}^{\text{CN}} := \min(w_i - b_{i+\frac{1}{2}}^{\text{CN}}, h_i), \quad h_{i+\frac{1}{2}+}^{\text{CN}} := \min(w_{i+1} - b_{i+\frac{1}{2}}^{\text{CN}}, h_{i+1}), \quad (3.17)$$

and the interface source terms are defined as

$$s_{i+\frac{1}{2}-}^{\text{CN}} := -\frac{g}{\Delta x} \frac{h_i + h_{i+\frac{1}{2}-}^{\text{CN}}}{2} (b_{i+\frac{1}{2}}^{\text{CN}} - b_i), \quad (3.18)$$

$$s_{i+\frac{1}{2}+}^{\text{CN}} := -\frac{g}{\Delta x} \frac{h_{i+\frac{1}{2}+}^{\text{CN}} + h_{i+1}}{2} (b_{i+1} - b_{i+\frac{1}{2}}^{\text{CN}}). \quad (3.19)$$

Remark 1. *For fully wet interfaces, all schemes are identical. The differences for the partially wet case are rather subtle.*

3.4 The numerical flux

In [2, p. 763], the authors used Harten-Lax-van Leer (HLL)-type Riemann solvers,

$$\mathcal{F}_{\text{HLL}}(U_-, U_+) = \frac{s^+ F(U_-) - s^- F(U_+) + s^+ s^- (U_+ - U_-)}{s^+ - s^-}, \quad (3.20)$$

where the smallest and largest wave speed s^- and s^+ are chosen as

$$s^- = \min(u_- - a_-, u_+ - a_+, 0), \quad s^+ = \max(u_- + a_-, u_+ + a_+, 0), \quad (3.21)$$

with gravitational wave speed $a = \sqrt{gh}$.

The following inequalities satisfied by the HLL flux are crucial for the stability analysis. The first one states that there is no numerical mass flux out of an empty cell, and is used to prove positivity of the water height later.

Lemma 1. *The first component of the HLL flux satisfies*

$$\mathcal{F}^h((0, 0), (h, hu)) \leq 0, \quad \mathcal{F}^h((h, hu), (0, 0)) \geq 0. \quad (3.22)$$

Proof. From the definition of the wave speeds (3.21), the value of HLL flux at $U_- = (0, 0)$ and $U_+ = (h, hu)$ is

$$\mathcal{F}_{\text{HLL}}((0, 0), (h, hu)) = \frac{-s^- (hu, hu^2 + \frac{1}{2}gh^2) + s^+ s^- (h, hu)}{s^+ - s^-}. \quad (3.23)$$

Hence, the first component of the HLL flux is given by

$$\mathcal{F}_{HLL}^h((0,0), (h, hu)) = \frac{-s^- hu + s^+ s^- h}{s^+ - s^-} = h s^- \frac{s^+ - u}{s^+ - s^-}. \quad (3.24)$$

Since $h \geq 0$, $s^- \leq 0$, $s^+ \geq u$, $s^+ > s^-$, the last equality gives us $\mathcal{F}_{HLL}^h \leq 0$.

Similarly,

$$\mathcal{F}_{HLL}((h, hu), (0,0)) = \frac{s^+ (hu, hu^2 + \frac{1}{2}gh^2) - s^+ s^- (h, hu)}{s^+ - s^-}, \quad (3.25)$$

and thus

$$\mathcal{F}_{HLL}^h((h, hu), (0,0)) = s^+ h \frac{u - s^-}{s^+ - s^-}. \quad (3.26)$$

Since $h \geq 0$, $s^+ \geq 0$, $u \geq s^-$, $s^+ > s^-$, the last equality implies that $\mathcal{F}_{HLL}^h \geq 0$. \square

The second inequality is used to prove the semidiscrete entropy inequality. It states that the numerical mass flux into a vacuum cell is at least as large as the physical mass flux.

Lemma 2. *The first component of the HLL flux satisfies*

$$\mathcal{F}_{HLL}^h((h, hu), (0,0)) - hu \geq 0, \quad hu - \mathcal{F}_{HLL}^h((0,0), (h, hu)) \geq 0. \quad (3.27)$$

Proof. The proof follows from (3.20) and (3.21) via

$$\begin{aligned} \mathcal{F}_{HLL}^h((h, hu), (0,0)) - hu &= \frac{s^+ hu - s^+ s^- h}{s^+ - s^-} - hu = -\frac{s^- h (s^+ - u)}{s^+ - s^-} \geq 0, \\ hu - \mathcal{F}_{HLL}^h((0,0), (h, hu)) &= hu - \frac{-s^- hu + s^+ s^- h}{s^+ - s^-} = \frac{s^+ h (u - s^-)}{s^+ - s^-} \geq 0. \end{aligned}$$

This completes our proof. \square

3.5 Interpretation via subcell reconstructions

Recall that in [2] the *singular layers* (or *internal boundary layers*) are defined by

$$\widehat{C}_{i+\frac{1}{2}}^\varepsilon := [x_{i+\frac{1}{2}} - \varepsilon, x_{i+\frac{1}{2}} + \varepsilon], \forall i \in \mathcal{E}. \quad (3.28)$$

Over each of these infinitesimal layers the bottom is reconstructed continuously by a function $b_\varepsilon(x)$. The flow variables are reconstructed by piecewise continuous functions $h_\varepsilon(x)$, $w_\varepsilon(x)$, and $u_\varepsilon(x)$ over the singular subcells

$$\widehat{C}_{i+\frac{1}{2}-}^\varepsilon := [x_{i+\frac{1}{2}} - \varepsilon, x_{i+\frac{1}{2}}] \text{ and } \widehat{C}_{i+\frac{1}{2}+}^\varepsilon := [x_{i+\frac{1}{2}}, x_{i+\frac{1}{2}} + \varepsilon], \quad \forall i \in \mathcal{E}. \quad (3.29)$$

These reconstructions provide the data of the Riemann problem at the interface, with an approximate Riemann solver $F_\varepsilon(x_{i+\frac{1}{2}})$. The source term is computed over the singular subcells. Together, this gives the *residuum*

$$R_i^\varepsilon := -\frac{1}{\Delta x} \left(F_\varepsilon(x_{i+\frac{1}{2}}) - F_\varepsilon(x_{i-\frac{1}{2}}) \right) + \frac{1}{\Delta x} \int_{C_i} S(U_\varepsilon(x), b_\varepsilon(x)) dx. \quad (3.30)$$

3.5.1 Splitting the cells into subcells

Let us denote the *interior subcell* by $C_i^\varepsilon := [x_{i-\frac{1}{2}} + \varepsilon, x_{i+\frac{1}{2}} - \varepsilon]$. Then

$$C_i = \widehat{C}_{i-\frac{1}{2}+}^\varepsilon \cup C_i^\varepsilon \cup \widehat{C}_{i+\frac{1}{2}-}^\varepsilon. \quad (3.31)$$

The *piecewise continuous reconstruction* is defined as follows.

Definition 1 (subcell reconstruction). *Given values φ_i and $\varphi_{i+\frac{1}{2}\pm}$ for $i \in \mathbb{Z}$, let $\widehat{\varphi}_{i+\frac{1}{2}\pm}^\varepsilon : \widehat{C}_{i+\frac{1}{2}\pm}^\varepsilon \rightarrow \mathbb{R}$ be Lipschitz continuous functions with boundary values*

$$\varphi_{i+\frac{1}{2}\pm}(x_{i+\frac{1}{2}}) = \varphi_{i+\frac{1}{2}\pm}, \quad \varphi_{i+\frac{1}{2}\pm}(x_{i+\frac{1}{2}} \pm \varepsilon) = \varphi_{i+\frac{1}{2}\pm 1}. \quad (3.32)$$

Then $\varphi_\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ is the piecewise continuous function given by

$$\varphi_\varepsilon(x) := \begin{cases} \varphi_i, & \text{if } x \in C_i^\varepsilon, \\ \widehat{\varphi}_{i+\frac{1}{2}\pm}^\varepsilon, & \text{if } x \in \widehat{C}_{i+\frac{1}{2}\pm}^\varepsilon. \end{cases}$$

Remark 2. If $\widehat{\varphi}_{i+\frac{1}{2}-}^\varepsilon$ and $\widehat{\varphi}_{i+\frac{1}{2}+}^\varepsilon$ are both linear, we call them the standard subcell reconstruction at interface $x_{i+\frac{1}{2}}$. The only exception from the standard subcell reconstruction will occur in the definition of the water level for Audusse's and Morales' schemes for partially wet interfaces.

3.5.2 Reconstruction of the bottom $b_\varepsilon(x)$

For all three HR schemes, a continuous bottom is defined by the standard subcell reconstruction (Def. 1) with

$$b_{i-\frac{1}{2}+} := b_{i-\frac{1}{2}}, \quad b_{i+\frac{1}{2}-} := b_{i+\frac{1}{2}}. \quad (3.33)$$

The reconstructed bottom is globally continuous for fixed $\varepsilon > 0$, but will have steep layers in $\widehat{C}_{i+\frac{1}{2}-}^\varepsilon$ and $\widehat{C}_{i+\frac{1}{2}+}^\varepsilon$.

3.5.3 Infinitesimal HR

The water level and height are now reconstructed. Several modern well-balanced schemes, as well as the third HR schemes, use the fact that the water level $w(x)$ is constant for still water. Thus, the piecewise constant reconstruction becomes exact for this equilibrium state.

Given the bottom topography $b_\varepsilon(x)$, these schemes reconstruct the water level $w_\varepsilon(x)$. The reconstructed water height is then reconstructed simply as

$$h_\varepsilon(x) := w_\varepsilon(x) - b_\varepsilon(x). \quad (3.34)$$

The *conservative variables* are given by

$$U_\varepsilon(x) := \sum_{i=1}^{N_x} \begin{pmatrix} h_\varepsilon(x) \\ h_\varepsilon(x) u_i \end{pmatrix} \mathbf{1}_{C_i}(x), \quad (3.35)$$

The following integral averages of $h_\varepsilon(x)$ in subcells $\widehat{C}_{i+\frac{1}{2}-}^\varepsilon$ and $\widehat{C}_{i+\frac{1}{2}+}^\varepsilon$ are used to calculate the flux and the source term later.

$$\bar{h}_{i+\frac{1}{2}-} := \frac{1}{\varepsilon} \int_{\widehat{C}_{i+\frac{1}{2}-}^\varepsilon} h_\varepsilon(x) dx, \quad \bar{h}_{i+\frac{1}{2}+} := \frac{1}{\varepsilon} \int_{\widehat{C}_{i+\frac{1}{2}+}^\varepsilon} h_\varepsilon(x) dx. \quad (3.36)$$

3.5.3.1. The original HR method. Define $w_\varepsilon(x)$ as in Def. 1, with

$$w_i^{\text{AUD}} := h_i + b_i, \quad w_{i+1}^{\text{AUD}} := h_{i+1} + b_{i+1}, \quad (3.37)$$

$$w_{i+\frac{1}{2}-}^{\text{AUD}} := \max(b_{i+\frac{1}{2}}^{\text{AUD}}, w_i), \quad w_{i+\frac{1}{2}+}^{\text{AUD}} := \max(b_{i+\frac{1}{2}}^{\text{AUD}}, w_{i+1}). \quad (3.38)$$

In the fully wet case (see Fig. 2),

$$\widehat{w}_{i+\frac{1}{2}-}^{\text{AUD}} \equiv w_i, \quad \widehat{w}_{i+\frac{1}{2}+}^{\text{AUD}} \equiv w_{i+1}, \quad (3.39)$$

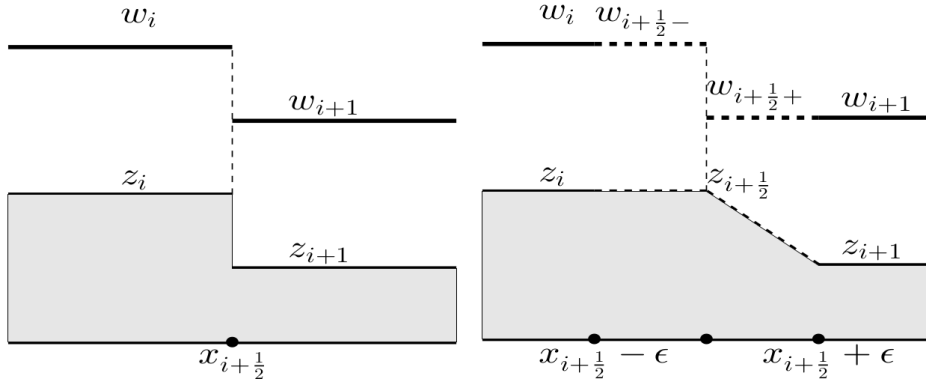


Figure 2: Subcell reconstruction of topography and water level for the fully wet case. Left: Riemann data. Right: reconstructed $b_\varepsilon(x)$ and $w_\varepsilon(x)$ for all three HR schemes.

while in the partially wet case (see Fig. 3),

$$\widehat{w}_{i+\frac{1}{2}-}^{\text{AUD}}(x) := \max(b_\varepsilon^{\text{AUD}}(x), w_i), \quad \widehat{w}_{i+\frac{1}{2}+}^{\text{AUD}}(x) := \max(b_\varepsilon^{\text{AUD}}(x), w_{i+1}). \quad (3.40)$$

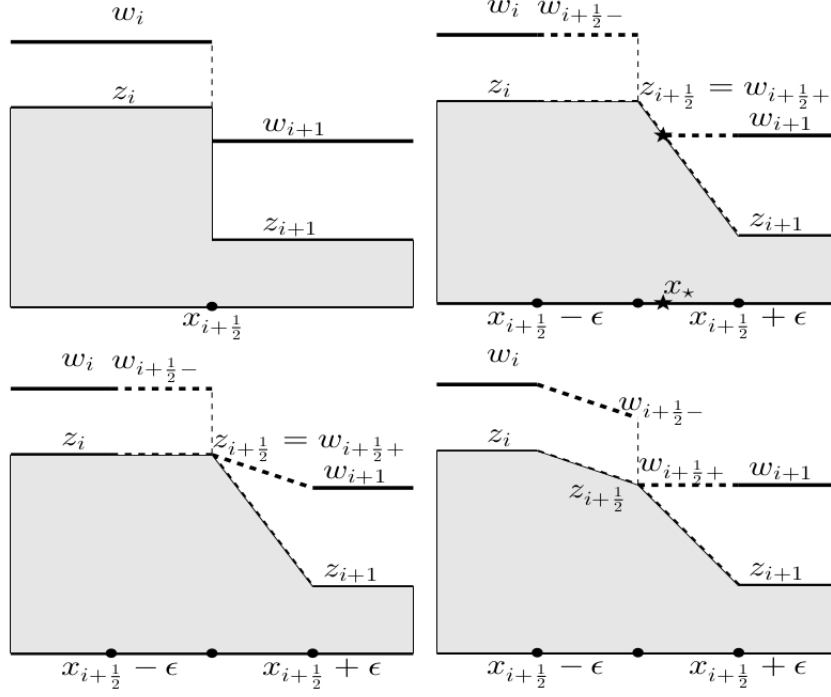


Figure 3: Subcell reconstruction of topography and water level for the partially wet case. Top left: Riemann data. Top right: Audusse or Morales scheme for slow uphill flow (the vacuum region $[x_{i+\frac{1}{2}}, x_\star]$). Bottom left: Morales scheme for fast uphill flow. Bottom right: CN scheme.

This is the only instance where the subcell reconstruction may differ from the standard definition. This will happen if and only if the wet-dry front is contained in one of the cells $\hat{C}_{i+\frac{1}{2}-}^\varepsilon$ or $\hat{C}_{i+\frac{1}{2}+}^\varepsilon$.

The average of $h_\varepsilon^{\text{AUD}}(x)$ over subcell $\hat{C}_{i+\frac{1}{2}-}^\varepsilon$ is defined as follows (that over subcell $\hat{C}_{i+\frac{1}{2}+}^\varepsilon$ is defined similarly). There are two cases to be discussed.

- i) In the fully wet case (in which all three schemes coincide) the water height $\hat{h}_{i+\frac{1}{2}+}^{\text{AUD}}$ is linear, so

$$\bar{h}_{i+\frac{1}{2}+}^{\text{AUD}} = \frac{h_{i+\frac{1}{2}+}^{\text{AUD}} + h_{i+1}}{2}. \quad (3.41)$$

- ii) In the partially wet case, $h_{i+\frac{1}{2}+}^{\text{AUD}} = 0$ (see top right of Fig. 3). Assume that the wet front is located at $x_\star \in \hat{C}_{i+\frac{1}{2}+}^\varepsilon$. Then the average water height is

$$\bar{h}_{i+\frac{1}{2}+}^{\text{AUD}} = \frac{x_{i+\frac{1}{2}} + \varepsilon - x_\star}{\varepsilon} \frac{h(x_\star) + h_{i+1}}{2} = \frac{h_{i+1}}{b_{i+\frac{1}{2}} - b_{i+1}} \frac{h_{i+1}}{2}, \quad (3.42)$$

where we have used the intercept theorem in the last inequality.

Summarizing (3.41) and (3.42), we obtain

$$\bar{h}_{i+\frac{1}{2}+}^{\text{AUD}} = \begin{cases} \frac{h_{i+1} + h_{i+\frac{1}{2}+}^{\text{AUD}}}{2}, & \text{if } h_{i+\frac{1}{2}+}^{\text{AUD}} > 0, \\ \frac{h_{i+1}}{2} \frac{h_{i+1}}{b_{i+\frac{1}{2}}^{\text{AUD}} - b_{i+1}}, & \text{if } h_{i+\frac{1}{2}+}^{\text{AUD}} = 0. \end{cases} \quad (3.43)$$

Similarly,

$$\bar{h}_{i+\frac{1}{2}-}^{\text{AUD}} = \begin{cases} \frac{h_i + h_{i+\frac{1}{2}-}^{\text{AUD}}}{2}, & \text{if } h_{i+\frac{1}{2}-}^{\text{AUD}} > 0, \\ \frac{h_i}{2} \frac{h_i}{b_{i+\frac{1}{2}}^{\text{AUD}} - b_i}, & \text{if } h_{i+\frac{1}{2}-}^{\text{AUD}} = 0. \end{cases} \quad (3.44)$$

3.5.3.2. The HR scheme of Morales et al. The continuous bottom of Morales' scheme coincides with that of the original HR scheme:

$$b_\varepsilon^{\text{MOR}}(x) \equiv b_\varepsilon^{\text{AUD}}(x). \quad (3.45)$$

The water level coincides with that of the original scheme, except for the partially wet interface (3.4). If the water flows downhill, or uphill with enough kinetic energy to climb the discrete jump of the bottom, i.e., (3.12) (resp., (3.14)) holds, then over $\hat{C}_{i+\frac{1}{2}-}$ (resp. $\hat{C}_{i+\frac{1}{2}+}$) the reconstructed water level $w_\varepsilon^{\text{MOR}}(x)$ is given by the standard subcell reconstruction (see Def. 1) instead of Audusse's piecewise linear reconstruction (3.38) (see Fig. 3). Then the local averages of $h_\varepsilon^{\text{MOR}}(x)$ over subcells $\hat{C}_{i+\frac{1}{2}-}^\varepsilon$ (resp., $\hat{C}_{i+\frac{1}{2}+}^\varepsilon$) are simply

$$\bar{h}_{i+\frac{1}{2}-}^{\text{MOR}} = \frac{h_i + h_{i+\frac{1}{2}-}^{\text{MOR}}}{2} = \frac{h_i}{2}, \quad \bar{h}_{i+\frac{1}{2}+}^{\text{MOR}} = \frac{h_{i+1} + h_{i+\frac{1}{2}+}^{\text{MOR}}}{2} = \frac{h_{i+1}}{2}, \quad (3.46)$$

since $h_{i+\frac{1}{2}-}^{\text{MOR}} = 0$ (resp., $h_{i+\frac{1}{2}+}^{\text{MOR}} = 0$).

3.5.3.3. The third HR scheme. The continuous bottom of the third HR scheme in [2], denoted by $b_\varepsilon^{\text{CN}}$ is defined by the standard subcell reconstruction with

$$b_{i+\frac{1}{2}-}^{\text{CN}} = b_{i+\frac{1}{2}+}^{\text{CN}} = b_{i+\frac{1}{2}}^{\text{CN}}. \quad (3.47)$$

The reference values for the water surface are given by

$$w_{i+\frac{1}{2}-}^{\text{CN}} := \min(w_i, b_{i+\frac{1}{2}}^{\text{CN}} + h_i), \quad w_{i+\frac{1}{2}+}^{\text{CN}} := \min(w_{i+1}, b_{i+\frac{1}{2}}^{\text{CN}} + h_{i+1}), \quad (3.48)$$

and the reference values for the water depth are given by (3.17). Due to the linearity of $\hat{h}_{i+\frac{1}{2}-}$ (resp., $\hat{h}_{i+\frac{1}{2}+}$), the average values of h_ε over the singular subcells $\hat{C}_{i+\frac{1}{2}-}^\varepsilon$ (resp., $\hat{C}_{i+\frac{1}{2}+}^\varepsilon$) are

$$\bar{h}_{i+\frac{1}{2}-}^{\text{CN}} = \frac{h_i + h_{i+\frac{1}{2}-}^{\text{CN}}}{2}, \quad \bar{h}_{i+\frac{1}{2}+}^{\text{CN}} = \frac{h_i + h_{i+\frac{1}{2}+}^{\text{CN}}}{2}. \quad (3.49)$$

Remark 3. From (3.43), (3.44), (3.46) and (3.49), the subcell averages $\bar{h}_{i+\frac{1}{2}-}$ and $\bar{h}_{i+\frac{1}{2}+}$ of $h_\varepsilon(x)$ over subcells $\widehat{C}_{i+\frac{1}{2}-}^\varepsilon$ and $\widehat{C}_{i+\frac{1}{2}+}^\varepsilon$ obtained by the three schemes are in fact independent of ε .

Proof. □

3.5.4 Fluxes and source terms based on subcell reconstructions

For all three hydrostatic schemes the flux vector $F_\varepsilon(x)$ is reconstructed by the standard subcell reconstruction (see Def. 1) with reference values

$$F_i := F(U_i), \quad F_{i+\frac{1}{2}-} := F_{i+\frac{1}{2}+} := F_{i+\frac{1}{2}}, \quad (3.50)$$

where

$$F_{i+\frac{1}{2}} := \mathcal{F}_{\text{HLL}} \left(U_{i+\frac{1}{2}-}, U_{i+\frac{1}{2}+} \right) \quad (3.51)$$

is an approximate Riemann solver. Note that $F_\varepsilon(x)$ is globally continuous. The definition of the reconstructed source term $S_\varepsilon(x) := (0, s_\varepsilon(x))^T$ takes the natural form

$$s_\varepsilon(x) := -gh_\varepsilon(x) b_\varepsilon'(x), \quad (3.52)$$

and hence corresponds directly to (1.4).

Given (3.50)-(3.52), we now introduce the reconstructed, cell-averaged residuum by

$$R_i^\varepsilon := -\frac{1}{\Delta x} \left(F_\varepsilon \left(x_{i+\frac{1}{2}} \right) - F_\varepsilon \left(x_{i-\frac{1}{2}} \right) \right) + \frac{1}{\Delta x} \int_{C_i} S(U_\varepsilon(x), b_\varepsilon(x)) dx. \quad (3.53)$$

Depending on the choice of hydrostatic scheme, we denote the residuum by $R_i^{\varepsilon, \text{AUD}}$, $R_i^{\varepsilon, \text{MOR}}$, and $R_i^{\varepsilon, \text{CN}}$.

Theorem 1. For each of the three hydrostatic schemes, and for each cell C_i , the reconstructed residuums are independent of ε , $R_i^\varepsilon = \bar{R}_i$ for all $\varepsilon > 0$, and coincide with the original definitions given above:

$$\bar{R}_i^{\text{AUD}} = R_i^{\text{AUD}}, \bar{R}_i^{\text{MOR}} = R_i^{\text{MOR}}, \bar{R}_i^{\text{CN}} = R_i^{\text{CN}}. \quad (3.54)$$

Proof. See [2, pp. 768–769] □

Remark 4. The main advantage of the third HR scheme is its accuracy for shallow downhill flows.

3.6 Comparison of the HR schemes

The third HR scheme only differs from the previous methods in the partially wet case (3.4). Additionally, it differs only in $\bar{h}_{i+\frac{1}{2}\pm}$ and $b_{i+\frac{1}{2}}$.

Proposition 6. *i) For all interfaces $x_{i+\frac{1}{2}}$,*

$$h_{i+\frac{1}{2}\pm}^{\text{AUD}} = h_{i+\frac{1}{2}\pm}^{\text{MOR}} = h_{i+\frac{1}{2}\pm}^{\text{CN}} \quad (3.55)$$

and

$$0 \leq h_{i+\frac{1}{2}-} \leq h_i, \quad 0 \leq h_{i+\frac{1}{2}+} \leq h_{i+1}. \quad (3.56)$$

ii) For fully wet interfaces $x_{i+\frac{1}{2}}$,

$$b_{i+\frac{1}{2}}^{\text{AUD}} = b_{i+\frac{1}{2}}^{\text{MOR}} = b_{i+\frac{1}{2}}^{\text{CN}} \text{ and } \bar{h}_{i+\frac{1}{2}\pm}^{\text{AUD}} = \bar{h}_{i+\frac{1}{2}\pm}^{\text{MOR}} = \bar{h}_{i+\frac{1}{2}\pm}^{\text{CN}}. \quad (3.57)$$

iii) For partially wet interfaces $x_{i+\frac{1}{2}}$ (see (3.4)), and if water flows downhill, or if it flows uphill with too little kinetic energy to climb the discrete jump of the bottom (i.e., neither (3.12) nor (3.14) holds), then

$$b_{i+\frac{1}{2}}^{\text{AUD}} = b_{i+\frac{1}{2}}^{\text{MOR}} \text{ and } \bar{h}_{i+\frac{1}{2}\pm}^{\text{AUD}} = \bar{h}_{i+\frac{1}{2}\pm}^{\text{MOR}}. \quad (3.58)$$

Proof. The proof is a direct computation based on the definition of the interface values (3.6), (3.7), (3.16), and (3.17), and the integral averages of the subcell water heights defined in (3.36).

.

□

3.7 Stability analysis

Recall that the well-known convex decomposition of the semidiscrete finite volume scheme (1.3) is defined by

$$\frac{d}{dt}U_i(t) = R_{i-\frac{1}{2}+} + R_{i-\frac{1}{2}+} \quad (3.59)$$

$$:= -\frac{1}{\Delta x} \left(F(U_i) - F_{i-\frac{1}{2}} \right) + S_{i-\frac{1}{2}+} - \frac{1}{\Delta x} \left(F_{i+\frac{1}{2}} - F(U_i) \right) + S_{i+\frac{1}{2}-}. \quad (3.60)$$

The following theorem states that the third HR scheme preserves the positivity of the water height under the same condition as Audusse's scheme.

Theorem 2 (Positivity of water height). *Under condition (3.22), the new semidiscrete HR scheme guarantees nonnegative water height for the homogeneous shallow water equations.*

Proof.

□

Before proving that the third HR scheme is well-balanced, we would like to distinguish the follow two classes of equilibria.

Definition 2. *i) Given a constant water level w_{eq} , the still water equilibrium is given by $u(x) \equiv 0$ and*

$$h(x) + b(x) \equiv w_{eq}. \quad (3.61)$$

The cell averages are consistent with the still water equilibrium, if for all i , $u_i = 0$ and

$$h_i + b_i = w_{eq}. \quad (3.62)$$

ii) The lake at rest equilibrium is given by $u(x) \equiv 0$ and

$$h(x) \partial_x (h(x) + b(x)) \equiv 0, \quad (3.63)$$

for some constant $w_{eq} \geq \max_{x \in \mathbb{R}} b(x)$. Moreover, near a wet-dry interface, the dry part of b should not be lower than the adjacent water level.

The cell averages are locally (at interface $x_{i+\frac{1}{2}}$) consistent with the lake at rest, if $u_i = u_{i+1} = 0$ and either $x_{i+\frac{1}{2}}$ is an interior interface (the still water case)

$$h_i > 0, h_{i+1} > 0, \text{ and } h_i + b_i = h_{i+1} + b_{i+1}, \quad (3.64)$$

or a dry-wet front

$$h_i = 0, h_{i+1} > 0, \text{ and } b_i \geq h_{i+1} + b_{i+1}, \quad (3.65)$$

or a wet-dry front

$$h_i > 0, h_{i+1} = 0, \text{ and } b_{i+1} \geq h_i + b_i, \quad (3.66)$$

or dry

$$h_i = h_{i+1} = 0. \quad (3.67)$$

The cell averages are globally consistent with the lake at rest, if they are locally consistent with the lake at rest for all interfaces $x_{i+\frac{1}{2}}$.

iii) Suppose that the cell averages of the semidiscrete finite volume scheme (1.3) are consistent with a given equilibrium state. Then we call the scheme well balanced for this equilibrium state if $R_i = 0$ for all i .

4 Offset equilibrium schemes

It suffices to modify in the upwind decenter scheme (2.1) only the discretization of velocity (unlike collocated schemes).

In [3], the author proposed the following modification:

$$\bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} = \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} \left((h_{i+1}^n)^2 - (h_i^n)^2 \right) + \frac{g}{2} (h_i^n + h_{i+1}^n) (b_{i+1} - b_i) \right],$$

which can be rewritten as

$$\bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} = \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} (h_i^n + h_{i+1}^n) ((h_{i+1}^n + b_{i+1}) - (h_i^n + b_i)) \right]. \quad (4.1)$$

We can easily verify that this numerical scheme preserves the balance-type states *still water*.

Proposition 7 (Well-balanced property). *The numerical scheme (2.1)-(4.1) preserved the steady states at rest stated in Proposition 4.*

Proof. Suppose that $u_{i+\frac{1}{2}}^n = 0$ and $h_i^n + b_i = \text{const}$ for all $i \in \mathcal{E}$, as the proof of Proposition 4, we arrive at

$$\begin{aligned} h_i^{n+1} &= h_i^n, \quad \forall i \in \mathcal{M}, \\ \frac{h_i^n + h_{i+1}^n}{2} u_{i+\frac{1}{2}}^{n+1} &= -\frac{\Delta t}{\Delta x} \frac{g}{2} (h_i^n + h_{i+1}^n) ((h_{i+1}^n + b_{i+1}) - (h_i^n + b_i)), \quad \forall i \in \mathcal{E}_{\text{int}}. \end{aligned}$$

It is straightforward to obtain $h_i^{n+1} + b_i = \text{const}$ and $u_{i+\frac{1}{2}}^{n+1} = 0$ for all $i \in \mathcal{M}$. □

In the case where the computational domain includes transitions between zones dry and wet areas, we can draw inspiration from the scheme developed in [2].

We define

$$\begin{cases} w_i^n = h_i^n + b_i, \\ b_{i+\frac{1}{2}}^n = \min(\max(b_i, b_{i+1}), \min(w_i^n, w_{i+1}^n)), \\ h_{i+\frac{1}{2}-}^n = \min(w_i^n - b_{i+\frac{1}{2}}^n, h_i^n), \\ h_{i+\frac{1}{2}+}^n = \min(w_{i+1}^n - b_{i+\frac{1}{2}}^n, h_{i+1}^n), \end{cases} \quad (4.2)$$

from which we deduce the discretization of the equation in velocity

$$\bar{h}_{i+\frac{1}{2}}^{n+1} u_{i+\frac{1}{2}}^{n+1} = \bar{h}_{i+\frac{1}{2}}^n u_{i+\frac{1}{2}}^n - \frac{\Delta t}{\Delta x} \left[G_{i+1}^n - G_i^n + \frac{g}{2} (h_i^n + h_{i+1}^n) (h_{i+\frac{1}{2}+}^n - h_{i+\frac{1}{2}-}^n) \right]. \quad (4.3)$$

We can check that this new scheme preserves the steady states of *the lake at rest*. We can assimilate the calculations of water depths $h_{i+\frac{1}{2}\pm}^n$ to a local *runup* algorithm. It remains to be seen whether this scheme practice on unsteady cases.

References

- [1] Emmanuel Audusse et al. “A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows”. In: *SIAM J. Sci. Comput.* 25.6 (2004), pp. 2050–2065. ISSN: 1064-8275. DOI: 10.1137/S1064827503431090. URL: <https://doi.org/10.1137/S1064827503431090>.
- [2] Guoxian Chen and Sebastian Noelle. “A new hydrostatic reconstruction scheme based on subcell reconstructions”. In: *SIAM J. Numer. Anal.* 55.2 (2017), pp. 758–784. ISSN: 0036-1429. DOI: 10.1137/15M1053074. URL: <https://doi.org/10.1137/15M1053074>.
- [3] D. Doyen and P. H. Gunawan. “An explicit staggered finite volume scheme for the shallow water equations”. In: *Finite volumes for complex applications VII. Methods and theoretical aspects*. Vol. 77. Springer Proc. Math. Stat. Springer, Cham, 2014, pp. 227–235. DOI: 10.1007/978-3-319-05684-5_21. URL: https://doi.org/10.1007/978-3-319-05684-5_21.
- [4] Vincent Duchêne. *Shallow-water models for water waves (Lectures in Master 2 Fundamental Mathematics and Applications, Université de Rennes 1)*. 2019.
- [5] R. Herbin, J.-C. Latché, and T. T. Nguyen. “Explicit staggered schemes for the compressible Euler equations”. In: *Applied mathematics in Savoie—AMIS 2012: Multiphase flow in industrial and environmental engineering*. Vol. 40. ESAIM Proc. EDP Sci., Les Ulis, 2013, pp. 83–102. DOI: 10.1051/proc/201340006. URL: <https://doi.org/10.1051/proc/201340006>.
- [6] T. Morales de Luna, M. J. Castro Díaz, and C. Parés. “Reliability of first order numerical schemes for solving shallow water system over abrupt topography”. In: *Appl. Math. Comput.* 219.17 (2013), pp. 9012–9032. ISSN: 0096-3003. DOI: 10.1016/j.amc.2013.03.033. URL: <https://doi.org/10.1016/j.amc.2013.03.033>.