

A SUBSONIC-WELL-BALANCED RECONSTRUCTION SCHEME FOR SHALLOW WATER FLOWS

FRANÇOIS BOUCHUT AND TOMÁS MORALES

ABSTRACT. We consider the Saint-Venant system for shallow water flows with non-flat bottom. In the past years, efficient well-balanced methods have been proposed in order to well resolve solutions close to steady states at rest. Here we describe a strategy based on a local subsonic steady-state reconstruction that allows to derive a subsonic-well-balanced scheme, preserving exactly all the subsonic steady states. It generalizes the now well-known hydrostatic solver, and as the latter it preserves nonnegativity of water height and satisfies a semi-discrete entropy inequality. An application to the Euler-Poisson system is proposed.

1. INTRODUCTION

We consider the classical Saint-Venant system for shallow water flows with topography. It is a hyperbolic system of conservation laws that approximately describes various geophysical flows, such as rivers, coastal areas, oceans when completed with a Coriolis term, and granular flows when completed with friction. Numerical approximate solutions to this system can be generated using conservative finite volume methods, which are known to properly handle shocks and contact discontinuities. As is now well-known, in the occurrence of source terms such as topography, a classical centered discretization does not allow precise computations for near steady states. One has then to use the so called well-balanced schemes, that properly balance the fluxes and the source at the level of each interface. Such schemes have been proposed in [13], [14], [4], [18], [12], [17], [15], [8], [16], [3], [5], [11], [10], [2], [6], [1], [9].

Additionally to the well-balanced property, the difficulty is to also have schemes that satisfy very natural properties such as conservativity of the water height ρ , nonnegativity of ρ , the ability to compute dry states $\rho = 0$ and transcritical flows when the Jacobian matrix F' of the flux function becomes singular, and eventually to satisfy a discrete entropy inequality. The solvers satisfying all these requirements are very few, they are those obtained by exact resolution in [10], by the kinetic method of [17], by the hydrostatic reconstruction method of [1], and by the Suliciu relaxation method of [7].

Nevertheless, apart from [10], these solvers preserve only the steady states at rest, for which $u \equiv 0$. The object of this paper is to go further in well-balanced schemes by building a solver with all the above requirements, and overall the property to be preserve exactly

Date: May 7, 2009.

Key words and phrases. shallow water, subsonic reconstruction, subsonic steady states, well-balanced scheme, semi-discrete entropy inequality.

all subsonic steady states. Note that a solver proposed in [9] is able to maintain all the steady states, but however it is not satisfying a discrete entropy inequality.

2. SAINT VENANT SYSTEM AND WELL-BALANCED SCHEMES

The Saint Venant system describes the evolution of the water height $\rho(t, x)$ and the velocity $u(t, x)$ in the horizontal direction, of a thin layer of water flowing over a slowly varying topography. In one space dimension, the systems writes as

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) + \rho g z_x = 0, \end{cases} \quad (2.1)$$

where $g > 0$ is the gravitational constant and $z(x)$ is the topography. We shall denote

$$Z = gz. \quad (2.2)$$

The physically relevant case is $p(\rho) = g\rho^2/2$, but we shall deal with the general case $p(\rho)$. We shall assume as usual that $p' > 0$, and we suppose that

$$\rho^2 p'(\rho) \text{ is strictly increasing,} \quad \rho^2 p'(\rho) \rightarrow \infty \text{ as } \rho \rightarrow \infty, \quad (2.3)$$

$$\int_0^1 \frac{p'(\rho)}{\rho} d\rho < \infty, \quad p'(\rho) \rightarrow 0 \text{ as } \rho \rightarrow 0, \quad \int_1^\infty \frac{p'(\rho)}{\rho} d\rho = \infty. \quad (2.4)$$

These assumptions are satisfied in particular for the pressure law of isentropic gas dynamics $p(\rho) = \kappa\rho^\gamma$, with $\gamma > 1$ and $\kappa > 0$. We define as usual the internal energy $e(\rho)$ by

$$e'(\rho) = \frac{p(\rho)}{\rho^2}. \quad (2.5)$$

Note that the integrability conditions in (2.4) imply that $e(\rho) + p(\rho)/\rho$ has a finite limit as $\rho \rightarrow 0$, and tends to ∞ as $\rho \rightarrow \infty$. For future reference we denote the flux by

$$F(U) = (\rho u, \rho u^2 + p(\rho)), \quad U = (\rho, \rho u). \quad (2.6)$$

The Saint Venant model is very robust, being hyperbolic and admitting an entropy inequality related to the physical energy,

$$\partial_t \tilde{\eta}(U, Z) + \partial_x \tilde{G}(U, Z) \leq 0, \quad (2.7)$$

where

$$\begin{aligned} \eta(U) &= \rho u^2/2 + \rho e(\rho), & G(U) &= (\rho u^2/2 + \rho e(\rho) + p(\rho)) u, \\ \tilde{\eta}(U, Z) &= \eta(U) + \rho Z, & \tilde{G}(U, Z) &= G(U) + \rho u Z. \end{aligned} \quad (2.8)$$

The steady states of (2.1) can be described as follows. We subtract u times the first equation in (2.1) to the second, and divide the result by ρ . We get

$$\partial_t u + \partial_x \left(u^2/2 + e(\rho) + \frac{p(\rho)}{\rho} + Z \right) = 0. \quad (2.9)$$

Therefore, the steady states are exactly the functions $\rho(x)$, $u(x)$ satisfying

$$\begin{cases} \rho u = Cst, \\ \frac{u^2}{2} + \left(e + \frac{p}{\rho}\right) (\rho) + Z = Cst. \end{cases} \quad (2.10)$$

In particular, we have the so-called steady state at rest

$$u = 0, \quad e + \frac{p}{\rho} + Z = Cst. \quad (2.11)$$

As exposed in [7], a first-order finite volume method for solving (2.1) writes generically with $U = (\rho, \rho u)$

$$U_i^{n+1} - U_i + \frac{\Delta t}{\Delta x_i} (F_{i+1/2-} - F_{i-1/2+}) = 0, \quad (2.12)$$

$$F_{i+1/2-} = \mathcal{F}_l(U_i, U_{i+1}, \Delta Z_{i+1/2}), \quad F_{i+1/2+} = \mathcal{F}_r(U_i, U_{i+1}, \Delta Z_{i+1/2}), \quad (2.13)$$

for some left/right numerical fluxes $\mathcal{F}_l(U_l, U_r, \Delta Z)$, $\mathcal{F}_r(U_l, U_r, \Delta Z)$, with

$$\Delta Z_{i+1/2} = Z_{i+1} - Z_i, \quad (2.14)$$

and where Δx_i denotes a possibly variable mesh size, $\Delta x_i = x_{i+1/2} - x_{i-1/2}$. We have then the following characterizations (see [7]).

▷ The conservativity or density writes, with $\mathcal{F}_{l/r} = (\mathcal{F}_{l/r}^\rho, \mathcal{F}_{l/r}^{\rho u})$,

$$\mathcal{F}_l^\rho = \mathcal{F}_r^\rho \equiv \mathcal{F}^\rho. \quad (2.15)$$

▷ The consistency-conservativity can be written

$$\begin{cases} \mathcal{F}^\rho(U, U, 0) = \rho u, \\ \mathcal{F}_l^{\rho u}(U, U, 0) = \mathcal{F}_r^{\rho u}(U, U, 0) = \rho u^2 + p(\rho), \end{cases} \quad (2.16)$$

$$\mathcal{F}_r^{\rho u}(U_l, U_r, \Delta Z) - \mathcal{F}_l^{\rho u}(U_l, U_r, \Delta Z) = -\rho \Delta Z + o(\Delta Z), \quad \text{as } U_l, U_r \rightarrow U, \Delta Z \rightarrow 0. \quad (2.17)$$

▷ The well-balancing property can be stated as the property to have, for the considered steady-states,

$$\mathcal{F}_l(U_l, U_r, \Delta Z) = F(U_l), \quad \mathcal{F}_r(U_l, U_r, \Delta Z) = F(U_r). \quad (2.18)$$

▷ The property to satisfy a semi-discrete entropy inequality (i.e. related to the limit $\Delta t \rightarrow 0$) is characterized by the existence of a numerical entropy flux $\tilde{\mathcal{G}}(U_l, U_r, Z_l, Z_r)$ consistent with the exact flux $\tilde{G}(U, Z)$, such that

$$\tilde{G}(U_r, Z_r) + \tilde{\eta}'(U_r, Z_r)(\mathcal{F}_r(U_l, U_r, \Delta Z) - F(U_r)) \leq \tilde{\mathcal{G}}(U_l, U_r, Z_l, Z_r), \quad (2.19)$$

$$\tilde{\mathcal{G}}(U_l, U_r, Z_l, Z_r) \leq \tilde{G}(U_l, Z_l) + \tilde{\eta}'(U_l, Z_l)(\mathcal{F}_l(U_l, U_r, \Delta Z) - F(U_l)), \quad (2.20)$$

where $\tilde{\eta}'(U, Z)$ is the derivative of $\tilde{\eta}(U, Z)$ with respect to U .

The hydrostatic reconstruction scheme satisfies all the above, and is defined as

$$\begin{aligned} \mathcal{F}_l(U_l, U_r, \Delta Z) &= \mathcal{F}(U_l^*, U_r^*) + \begin{pmatrix} 0 \\ p(\rho_l) - p(\rho_l^*) \end{pmatrix}, \\ \mathcal{F}_r(U_l, U_r, \Delta Z) &= \mathcal{F}(U_l^*, U_r^*) + \begin{pmatrix} 0 \\ p(\rho_r) - p(\rho_r^*) \end{pmatrix}, \end{aligned} \quad (2.21)$$

where $\mathcal{F}(U_l, U_r)$ is a numerical flux for the shallow water problem without source ($Z = cst$), and the reconstructed states U_l^* , U_r^* are defined by

$$U_l^* = (\rho_l^*, \rho_l^* u_l), \quad U_r^* = (\rho_r^*, \rho_r^* u_r), \quad (2.22)$$

$$\begin{aligned} (e + p/\rho)(\rho_l^*) &= \left((e + p/\rho)(\rho_l) - (\Delta Z)_+ \right)_+, \\ (e + p/\rho)(\rho_r^*) &= \left((e + p/\rho)(\rho_r) - (-\Delta Z)_+ \right)_+, \end{aligned} \quad (2.23)$$

where we use the notation $X_+ = \max(0, X)$, and we assumed that $e(\rho) + p(\rho)/\rho \rightarrow 0$ as $\rho \rightarrow 0$.

3. WELL-BALANCED SCHEME WITH SUBSONIC RECONSTRUCTION

We would like now to explain how it is possible to extend the hydrostatic reconstruction scheme (2.21)-(2.23) in order to obtain a scheme that satisfies the above requirements and preserves some more general steady-states than the rest steady states. We shall obtain in particular a scheme that preserves all subsonic steady states, that is the steady states that verify $u^2 < p'(\rho)$. This property will be called subsonic-well-balanced. Note in particular that the steady-states at rest (with $u = 0$) are subsonic.

3.1. Parametrization of numerical fluxes. Following [1], [7], we propose and analyze finite volume schemes defined by (2.12), (2.13) with numerical fluxes

$$\begin{aligned} \mathcal{F}_l(U_l, U_r, \Delta Z) &= \mathcal{F}(U_l^*, U_r^*) + \begin{pmatrix} 0 \\ p(\rho_l) - p(\rho_l^*) + T_l(U_l, U_r, \Delta Z) \end{pmatrix}, \\ \mathcal{F}_r(U_l, U_r, \Delta Z) &= \mathcal{F}(U_l^*, U_r^*) + \begin{pmatrix} 0 \\ p(\rho_r) - p(\rho_r^*) + T_r(U_l, U_r, \Delta Z) \end{pmatrix}, \end{aligned} \quad (3.1)$$

where \mathcal{F} stands for a numerical flux for the homogeneous problem ($Z = cst$), and the interface values U_l^*, U_r^* are derived from a local reconstruction procedure. They should satisfy at least that $U_l^* = U_l$, $U_r^* = U_r$ when $\Delta Z = 0$.

The extra terms T_l , T_r appear here in order to balance the advection term $\partial_x(\rho u^2)$ in (2.1), that was not considered in the hydrostatic scheme. Taking into account the constant discharge condition in (2.10), this balancing requirement suggests the relations

$$T_l = \rho_l u_l (u_l - u_l^*), \quad T_r = \rho_r u_r (u_r - u_r^*). \quad (3.2)$$

It is obvious from the characterization (2.18) that a steady state is maintained exactly by the scheme (2.12), (2.13), (3.1), if for such a state, the reconstructed states satisfy $U_l^* = U_r^*$, $\rho_l^* u_l^* = \rho_l u_l$, $\rho_r^* u_r^* = \rho_r u_r$, and (3.2) is satisfied.

Taking into account the steady states equation (2.10), one could guess a reconstruction of the states U_l^* , U_r^* as

$$\begin{cases} \frac{(u_l^*)^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_l^*) + Z^* = \frac{u_l^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_l) + Z_l, \\ \rho_l^* u_l^* = \rho_l u_l, \end{cases} \quad (3.3)$$

$$\begin{cases} \frac{(u_r^*)^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_r^*) + Z^* = \frac{u_r^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_r) + Z_r, \\ \rho_r^* u_r^* = \rho_r u_r, \end{cases} \quad (3.4)$$

with

$$Z^* = \max(Z_l, Z_r). \quad (3.5)$$

There exist solutions to the previous system if ΔZ is small enough (recall that $\Delta Z = Z_r - Z_l$), but this is not true for arbitrary ΔZ , as we shall see later on. The idea is thus to consider generalized T_l, T_r , to be defined later on. Their definition is motivated by the entropy inequality.

Lemma 3.1. *Let $\mathcal{F}(U_l, U_r)$ be a given consistent numerical flux for the Saint Venant problem without source that verifies a semi-discrete entropy inequality for the entropy pair (η, G) given by (2.8), and denote $\mathcal{F} = (\mathcal{F}^\rho, \mathcal{F}^{\rho u})$.*

*A sufficient condition for the scheme (2.12), (2.13), (3.1), to be semi-discrete entropy satisfying for the entropy pair $(\tilde{\eta}, \tilde{G})$ in (2.8) is that for some Z^**

$$\begin{aligned} & G(U_l^*) + \eta'(U_l^*)(\mathcal{F}(U_l^*, U_r^*) - F(U_l^*)) + \mathcal{F}^\rho(U_l^*, U_r^*)Z^* \\ & \leq G(U_l) + \eta'(U_l)(\mathcal{F}_l(U_l, U_r, \Delta Z) - F(U_l)) + \mathcal{F}^\rho(U_l^*, U_r^*)Z_l, \end{aligned} \quad (3.6)$$

and

$$\begin{aligned} & G(U_r) + \eta'(U_r)(\mathcal{F}_r(U_l, U_r, \Delta Z) - F(U_r)) + \mathcal{F}^\rho(U_l^*, U_r^*)Z_r \\ & \leq G(U_r^*) + \eta'(U_r^*)(\mathcal{F}(U_l^*, U_r^*) - F(U_r^*)) + \mathcal{F}^\rho(U_l^*, U_r^*)Z^*. \end{aligned} \quad (3.7)$$

Proof. The numerical flux \mathcal{F} satisfies a semidiscrete entropy inequality associated to the entropy pair (η, G) , thus

$$\begin{aligned} & G(U_r) + \eta'(U_r)(\mathcal{F}(U_l, U_r) - F(U_r)) \\ & \leq \mathcal{G}(U_l, U_r) \leq G(U_l) + \eta'(U_l)(\mathcal{F}(U_l, U_r) - F(U_l)), \end{aligned} \quad (3.8)$$

for a numerical flux \mathcal{G} consistent with G . For the scheme (2.12), (2.13), (3.1) to be semi-discrete entropy satisfying for the entropy pair $(\tilde{\eta}, \tilde{G})$, (2.19), (2.20) should hold. Let

$$\tilde{\mathcal{G}}(U_l, U_r, Z_l, Z_r) = \mathcal{G}(U_l^*, U_r^*) + \mathcal{F}^\rho(U_l^*, U_r^*)Z^*. \quad (3.9)$$

As \mathcal{G} is consistent with G , $\tilde{\mathcal{G}}$ is consistent with \tilde{G} . The comparison between (3.8) evaluated at (U_l^*, U_r^*) and (2.19), (2.20) gives that (3.6), (3.7) are sufficient conditions. \square

Lemma 3.2. *Denote $(\mathcal{F}^\rho, \mathcal{F}^{\rho u}) \equiv \mathcal{F}(U_l^*, U_r^*)$, $T_l \equiv T_l(U_l, U_r, \Delta Z)$, $T_r \equiv T_r(U_l, U_r, \Delta Z)$, and define the quantities*

$$\begin{aligned} W_l & \equiv \mathcal{F}^\rho \cdot \left(\left(e + \frac{p}{\rho}\right)(\rho_l) - \left(e + \frac{p}{\rho}\right)(\rho_l^*) + Z_l - Z^* + \frac{(u_l^*)^2}{2} - \frac{u_l^2}{2} \right) \\ & \quad + (u_l - u_l^*)(\mathcal{F}^{\rho u} - p(\rho_l^*)) + u_l T_l, \end{aligned} \quad (3.10)$$

$$\begin{aligned} W_r & \equiv \mathcal{F}^\rho \cdot \left(\left(e + \frac{p}{\rho}\right)(\rho_r) - \left(e + \frac{p}{\rho}\right)(\rho_r^*) + Z_r - Z^* + \frac{(u_r^*)^2}{2} - \frac{u_r^2}{2} \right) \\ & \quad + (u_r - u_r^*)(\mathcal{F}^{\rho u} - p(\rho_r^*)) + u_r T_r. \end{aligned} \quad (3.11)$$

A necessary and sufficient condition for (3.6), (3.7) to hold is that

$$W_l \geq 0, \quad W_r \leq 0. \quad (3.12)$$

Proof. From the explicit value of F , G , and computing $\eta'(U) = (e(\rho) + p(\rho)/\rho - u^2/2, u)$, one gets the identity $G(U) - \eta'(U)F(U) = -up(\rho)$. Plugging this into (3.6), (3.7) yields the result. \square

We remark that in the particular case when U_l^* , U_r^* verify (3.3)-(3.4), we have

$$\begin{aligned} W_l &= (u_l^* - u_l) \left((u_l + u_l^*) \mathcal{F}^\rho - \mathcal{F}^{\rho u} + p(\rho_l^*) \right) + u_l T_l, \\ W_r &= (u_r^* - u_r) \left((u_r + u_r^*) \mathcal{F}^\rho - \mathcal{F}^{\rho u} + p(\rho_r^*) \right) + u_r T_r. \end{aligned} \quad (3.13)$$

Thus, the choice of T_l , T_r given by (3.2) makes $W_l = W_r = 0$ whenever $U_l^* = U_r^*$.

3.2. Subsonic reconstruction. We intend here to define reconstructed states that verify (3.3), (3.4) as far as possible. The possibility of achieving these relations is related to the following definitions.

Definition 3.3. Let $\rho \geq 0$ and $u \in \mathbb{R}$. We say that (ρ, u) is a sonic, subsonic or supersonic point for the Saint-Venant system (2.1) if we have respectively $u^2 = p'(\rho)$, $u^2 < p'(\rho)$ or $u^2 > p'(\rho)$.

Definition 3.4. Let $q \in \mathbb{R}$. We define $\rho_s(q)$ as the solution to

$$\rho_s^2 p'(\rho_s) = q^2, \quad \rho_s \geq 0, \quad (3.14)$$

and

$$m_s(q) = \left(e + \frac{p}{\rho} + \frac{p'}{2} \right) (\rho_s(q)). \quad (3.15)$$

Because of the assumptions (2.3), (2.4) on p , there exists a unique ρ_s solution to (3.14). Moreover, $m_s(0)$ is well defined,

$$m_s(0) = \left(e + \frac{p}{\rho} \right) (0). \quad (3.16)$$

In the particular case when $p(\rho) = \kappa \rho^\gamma$, we have

$$\rho_s(q) = \left(\frac{q^2}{\kappa \gamma} \right)^{\frac{1}{\gamma+1}}, \quad m_s(q) = \left(\frac{1}{2} + \frac{1}{\gamma-1} \right) (\kappa \gamma)^{\frac{2}{\gamma+1}} |q|^{2\frac{\gamma-1}{\gamma+1}}. \quad (3.17)$$

The following proposition gives the necessary and sufficient conditions for the existence of a solution to (3.3) and (3.4).

Proposition 3.5. Let $u_0 \in \mathbb{R}$, $\rho_0 \geq 0$, $\delta \in \mathbb{R}$. Consider the system

$$\begin{cases} \frac{(u^*)^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho^*) = \frac{u_0^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_0) + \delta, \\ \rho^* u^* = \rho_0 u_0, \\ \rho^* \geq 0, \quad u^* \in \mathbb{R}, \end{cases} \quad (3.18)$$

and denote $\rho_s \equiv \rho_s(\rho_0 u_0)$. There exists a solution (ρ^*, u^*) to (3.18) if and only if

$$\frac{u_0^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_0) + \delta \geq m_s(\rho_0 u_0). \quad (3.19)$$

Moreover,

(i) If we have equality in (3.19), there is only one solution (ρ^*, u^*) to (3.18), it is given by

$$\rho^* = \rho_s, \quad u^* = \begin{cases} \frac{\rho_0 u_0}{\rho_s} & \text{if } \rho_0 u_0 \neq 0, \\ 0 & \text{if } \rho_0 u_0 = 0. \end{cases} \quad (3.20)$$

(ii) If we have a strict inequality in (3.19), then there are exactly two different solutions $(\rho_{sup}^*, u_{sup}^*)$ and $(\rho_{sub}^*, u_{sub}^*)$ to (3.18), with $\rho_{sup}^* \leq \rho_s < \rho_{sub}^*$, and $\rho_{sup}^* < \rho_s$ for $\rho_0 u_0 \neq 0$.

(iii) A solution (ρ^*, u^*) to (3.18), with $\rho^* u^* \neq 0$, is a sonic (resp. subsonic or supersonic) point if and only if $\rho^* = \rho_s$ (resp. $\rho^* > \rho_s$ or $\rho^* < \rho_s$).

Proof. Let us suppose first that $\rho_0 u_0 \neq 0$, and consider the function

$$f : \mathbb{R} \times (0, \infty) \longrightarrow \mathbb{R}, \quad (q, \rho) \longmapsto f(q, \rho) = \frac{q^2}{2\rho^2} + \left(e + \frac{p}{\rho} \right) (\rho). \quad (3.21)$$

Then (ρ^*, u^*) is a solution to (3.18) if and only if $\rho^* u^* = \rho_0 u_0$ and $f(\rho_0 u_0, \rho^*) = f(\rho_0 u_0, \rho_0) + \delta$. We have $\frac{\partial f}{\partial \rho}(q, \rho) = (\rho^2 p'(\rho) - q^2)/\rho^3$, thus according to (2.3), (2.4), f is strictly increasing in $(\rho_s(q), \infty)$ and strictly decreasing in $(0, \rho_s(q))$. Therefore, $\rho_s(q)$ is a minimum point of f with minimum value $f(q, \rho_s(q)) = m_s(q)$. Figure 1 shows a sketch of the function $f(\rho_0 u_0, \cdot)$ and the solutions $\rho_{sub}^*, \rho_{sup}^*$ in the case $\delta < 0$. Thus, condition (3.19) follows, as well as (i) and (ii).

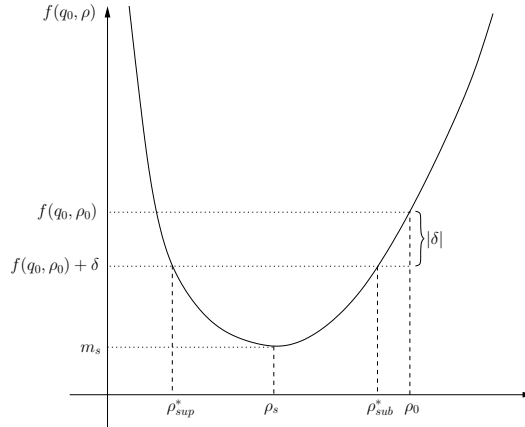


FIGURE 1. Function $f(q_0, \cdot)$

Now, if we consider (ρ^*, u^*) a solution of (3.18), we have

$$(u^*)^2 > p'(\rho^*) \Leftrightarrow (\rho_0 u_0)^2 > (\rho^*)^2 p'(\rho^*). \quad (3.22)$$

Since $\rho^2 p'(\rho)$ is a strictly increasing function with $\rho_s^2 p'(\rho_s) = (\rho_0 u_0)^2$,

$$(u^*)^2 > p'(\rho^*) \Leftrightarrow \rho^* < \rho_s, \quad (3.23)$$

which proves (iii).

In the case $\rho_0 u_0 = 0$, the second condition in (3.18) simplifies to $\rho^* u^* = 0$. Thus the system can have solutions with either $\rho^* = 0$, giving $(u^*)^2/2 = r.h.s - m_s(0)$, or $u^* = 0$, giving $(e + p/\rho)(\rho^*) = r.h.s$. One sees easily that a solution exists if and only if $r.h.s. \geq m_s(0)$, proving condition (3.19). In case of equality, the only solution is $\rho^* = 0$, $u^* = 0$, which is sonic. In case of inequality, the solutions are given by

$$\rho_{sup}^* = 0, \quad \frac{(u_{sup}^*)^2}{2} = \frac{u_0^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_0) + \delta - m_s(0), \quad (3.24)$$

$$u_{sub}^* = 0, \quad \left(e + \frac{p}{\rho}\right)(\rho_{sub}^*) = \frac{u_0^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_0) + \delta. \quad (3.25)$$

This gives the result, with the convention that for (3.24) we identify the two solutions having density $\rho_{sup}^* = 0$. Note that these solutions are supersonic, while the one corresponding to (3.25) are subsonic. \square

Corollary 3.6. *Let $\rho_0 > 0$ and $u_0 \in \mathbb{R} \setminus \{0\}$. Then (ρ_0, u_0) is a sonic (respectively subsonic or supersonic) point if and only if $\rho_0 = \rho_s(\rho_0 u_0)$ (respectively $\rho_0 > \rho_s(\rho_0 u_0)$ or $\rho_0 < \rho_s(\rho_0 u_0)$).*

Proof. The point (ρ_0, u_0) is a trivial solution to (3.18) with $\delta = 0$. Thus the result follows from (iii) in the previous proposition. \square

Lemma 3.7. *Suppose that in Proposition 3.5 we are in the case of two solutions $\rho_{sub}^*, \rho_{sup}^*$. Then one has the following ordering.*

- (A) Case (ρ_0, u_0) subsonic,
 - (A.1) if $\delta > 0$, then $\rho_{sup}^* < \rho_0 < \rho_{sub}^*$,
 - (A.2) if $\delta \leq 0$, then $\rho_{sup}^* < \rho_{sub}^* \leq \rho_0$,
- (B) Case (ρ_0, u_0) supersonic,
 - (B.1) if $\delta \geq 0$, then $\rho_{sup}^* \leq \rho_0 < \rho_{sub}^*$,
 - (B.2) if $\delta < 0$, then $\rho_0 \leq \rho_{sup}^* < \rho_{sub}^*$,
- (C) Case (ρ_0, u_0) sonic,
 - $\rho_{sup}^* \leq \rho_0 < \rho_{sub}^*$.

The proof is left to the reader. Now, in order to define a scheme that preserves the non-negativity of the density, the reconstructed states need to verify $\rho_l^* \leq \rho_l$ and $\rho_r^* \leq \rho_r$ (see Theorem 3.12 (i)). We notice that the equations (3.3), (3.4), correspond to the problem (3.18) with successively $\delta = Z_l - Z^*$, $\delta = Z_r - Z^*$. Since it is natural to try to choose a solution (ρ^*, u^*) of the same sonicity as (ρ_0, u_0) , one sees with the previous Lemma that in order to have $\rho^* \leq \rho_0$ one needs $\delta \leq 0$ in the subsonic case, and $\delta \geq 0$ in the supersonic

case. Recalling the values $\delta = Z_l - Z^*$ and $\delta = Z_r - Z^*$, this gives that one should have $Z^* \geq \max(Z_l, Z_r)$ in subsonic regions, while $Z^* \leq \min(Z_l, Z_r)$ in supersonic regions. We observe then that it is not possible to satisfy both conditions with Z^* depending continuously on the data. Thus we make the choice of solving exactly the subsonic steady states and disregard supersonic steady states, which justifies $Z^* = \max(Z_l, Z_r)$, i.e. (3.5).

Consider now the function f given by (3.21). For $q \in \mathbb{R}$ fixed, $f(q, \cdot)$ is strictly increasing in $[\rho_s(q), \infty)$ (recall that $\rho_s(q)$ and $m_s(q)$ are defined as (3.14), (3.15)),

$$\begin{aligned} f(q, \cdot)|_{[\rho_s(q), \infty)} : [\rho_s(q), \infty) &\rightarrow [m_s(q), \infty) \\ \rho &\mapsto \frac{q^2}{2\rho^2} + \left(e + \frac{p}{\rho}\right)(\rho). \end{aligned} \quad (3.26)$$

We consider its inverse and denote it by $f_r^{-1}(q, \cdot)$,

$$f_r^{-1}(q, \cdot) : [m_s(q), \infty) \rightarrow [\rho_s(q), \infty). \quad (3.27)$$

This inverse function corresponds to choosing the subsonic solution to (3.18). With the choice (3.5), the equations (3.3), (3.4) correspond to the problem (3.18) with successively $\delta = Z_l - Z^* = -(\Delta Z)_+$, $\delta = Z_r - Z^* = -(-\Delta Z)_+$. Therefore, we define the reconstructed states U_l^* , U_r^* by

$$\begin{aligned} \rho_l^* &= \min \left\{ \rho_l, f_r^{-1} \left(\rho_l u_l, \max \left\{ f(\rho_l u_l, \rho_l) - (\Delta Z)_+, m_s(\rho_l u_l) \right\} \right) \right\}, \\ u_l^* &= \rho_l u_l / \rho_l^* \quad (u_l^* = u_l \text{ if } \rho_l^* = 0), \\ \rho_r^* &= \min \left\{ \rho_r, f_r^{-1} \left(\rho_r u_r, \max \left\{ f(\rho_r u_r, \rho_r) - (-\Delta Z)_+, m_s(\rho_r u_r) \right\} \right) \right\}, \\ u_r^* &= \rho_r u_r / \rho_r^* \quad (u_r^* = u_r \text{ if } \rho_r^* = 0). \end{aligned} \quad (3.28)$$

Note that these definitions imply

$$\rho_l^* u_l^* = \rho_l u_l, \quad \rho_r^* u_r^* = \rho_r u_r. \quad (3.29)$$

Lemma 3.8. *The definitions (3.28) can be interpreted as follows.*

(A) *Case $\Delta Z \leq 0$:*

we have the trivial solution to (3.3) $U_l^ = U_l$.*

(B) *In the general case:*

by Proposition 3.5, we know that the system (3.3) has a solution if and only if

$$\frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) - (\Delta Z)_+ \geq m_s(\rho_l u_l). \quad (3.30)$$

(B.1) *If (ρ_l, u_l) is a supersonic point or a sonic point, then $U_l^* = U_l$.*

(B.2) *If (ρ_l, u_l) is a subsonic point and we have strict inequality in (3.30), then (ρ_l^*, u_l^*) is the subsonic solution to (3.3).*

(B.3) If (ρ_l, u_l) is a subsonic point and we have equality in (3.30) or the inequality is not satisfied, then $\rho_l^* = \rho_s(\rho_l u_l)$.

Similar statements hold for (ρ_r^*, u_r^*) .

The case $\rho_l^* = 0$ could pose some problems in the previous definition of u_l^* , but as we consider conservative variables, the product $\rho_l^* u_l^*$ is well defined. The following result shows that there is indeed no problem of continuity.

Lemma 3.9. *The reconstructed states (3.28) verify:*

$$(i) \quad \min\left\{\rho_l, \rho_s(\rho_l u_l)\right\} \leq \rho_l^* \leq \rho_l, \quad \min\left\{\rho_r, \rho_s(\rho_r u_r)\right\} \leq \rho_r^* \leq \rho_r, \quad (3.31)$$

(ii) *independently of the other arguments, one has*

$$\lim_{\rho_l \rightarrow 0} \rho_l^* = 0, \quad \lim_{\rho_r \rightarrow 0} \rho_r^* = 0, \quad (3.32)$$

(iii) *for $\rho_l \geq 0$, $\rho_r \geq 0$, ΔZ fixed, we have*

$$\lim_{u_l \rightarrow 0} \rho_l^* = \left(e + \frac{p}{\rho}\right)^{-1} \left(\max\left\{\left(e + \frac{p}{\rho}\right)(\rho_l) - (\Delta Z)_+, m_s(0)\right\}\right), \quad (3.33)$$

$$\lim_{u_r \rightarrow 0} \rho_r^* = \left(e + \frac{p}{\rho}\right)^{-1} \left(\max\left\{\left(e + \frac{p}{\rho}\right)(\rho_r) - (-\Delta Z)_+, m_s(0)\right\}\right), \quad (3.34)$$

(iv) *for ρ_l, ρ_r bounded, one has*

$$\lim_{u_l \rightarrow 0} u_l^* = 0, \quad \lim_{u_r \rightarrow 0} u_r^* = 0. \quad (3.35)$$

Proof. Note that (iii) means the continuity of ρ_l^* and ρ_r^* in this asymptotics. We shall only give the proof for ρ_l^* , the proof for ρ_r^* being similar.

According to Lemma 3.8, the only case when (3.31) is nontrivial is (B.2) with $\Delta Z \geq 0$. Then Corollary 3.6 allows to conclude the proof of (i). Then, (ii) is a consequence of (i).

In order to prove (iii), for $\rho_l \geq 0$ and ΔZ fixed, we shall denote

$$\begin{aligned} \beta_l(u_l) &= \max\left\{f(\rho_l u_l, \rho_l) - (\Delta Z)_+, m_s(\rho_l u_l)\right\}, \\ \alpha_l(u_l) &= f_r^{-1}(\rho_l u_l, \beta_l(u_l)). \end{aligned} \quad (3.36)$$

We have

$$\lim_{u_l \rightarrow 0} \beta_l(u_l) = \max\left\{\left(e + \frac{p}{\rho}\right)(\rho_l) - (\Delta Z)_+, m_s(0)\right\}. \quad (3.37)$$

Since $\beta_l(u_l) \geq m_s(\rho_l u_l)$, we have $\alpha_l(u_l) \geq \rho_s(\rho_l u_l)$, which yields

$$0 \leq \frac{(\rho_l u_l)^2}{\alpha_l(u_l)^2} \leq \frac{(\rho_l u_l)^2}{\rho_s(\rho_l u_l)^2} = p'(\rho_s(\rho_l u_l)). \quad (3.38)$$

Therefore, $(\rho_l u_l)^2 / \alpha_l(u_l)^2$ tends to 0 as $u_l \rightarrow 0$. Then, using the identity

$$\frac{(\rho_l u_l)^2}{2(\alpha_l(u_l))^2} + \left(e + \frac{p}{\rho}\right)(\alpha_l(u_l)) = \beta_l(u_l), \quad (3.39)$$

we get

$$\lim_{u_l \rightarrow 0} \alpha_l(u_l) = \left(e + \frac{p}{\rho}\right)^{-1} \left(\max \left\{ \left(e + \frac{p}{\rho}\right)(\rho_l) - (\Delta Z)_+, m_s(0) \right\} \right). \quad (3.40)$$

Since the right-hand side is at most ρ_l , we deduce that $\rho_l^* = \min\{\rho_l, \alpha_l(u_l)\}$ has the same limit, which concludes (iii).

The last statement (iv) is also consequence of (i) since either $\rho_l^* \geq \rho_l$ giving $|u_l^*| \leq |u_l|$, or $\rho_l^* \geq \rho_s(\rho_l u_l)$, giving

$$(u_l^*)^2 = \frac{(\rho_l u_l)^2}{(\rho_l^*)^2} \leq \frac{(\rho_l u_l)^2}{\rho_s(\rho_l u_l)^2} = p'(\rho_s(\rho_l u_l)). \quad (3.41)$$

Since the right-hand side tends to 0, this concludes the proof. \square

We would like to end this subsection by giving an iterative procedure in order to solve the system (3.18). According to Lemma 3.8, we need to solve this system only in the case (B.2) with $\Delta Z > 0$. This is done as follows.

Proposition 3.10. *Let $u_0 \in \mathbb{R} \setminus \{0\}$, $\rho_0 > 0$, and $\delta < 0$. We suppose that $u_0^2 < p'(\rho_0)$ and that (3.19) is strictly satisfied. Let $V_0 = \frac{u_0^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_0) + \delta$ and $\psi(\rho) = \rho^\alpha(f(q_0, \rho) - V_0)$, where f is the function given by (3.21), and $q_0 = \rho_0 u_0$. Then for $\alpha \geq 3/2$, the relation*

$$\rho_{n+1} = \rho_n - \frac{\psi(\rho_n)}{\psi'(\rho_n)} \quad (3.42)$$

(starting from ρ_0) defines a decreasing sequence that converges to ρ_{sub}^* , the subsonic solution to (3.18).

Proof. According to Proposition 3.5, there is a unique solution $\rho_{sub}^* \in (\rho_s, \rho_0)$ to the equation $\psi(\rho) = 0$. We have

$$\begin{aligned} \psi'(\rho) &= \alpha \rho^{\alpha-1} \left(f(q_0, \rho) - V_0 \right) + \rho^\alpha \frac{\partial f}{\partial \rho}(q_0, \rho), \\ \psi''(\rho) &= \alpha(\alpha-1) \rho^{\alpha-2} \left(f(q_0, \rho) - V_0 \right) + 2\alpha \rho^{\alpha-1} \frac{\partial f}{\partial \rho}(q_0, \rho) + \rho^\alpha \frac{\partial^2 f}{\partial \rho^2}(q_0, \rho), \end{aligned} \quad (3.43)$$

and for $\rho \geq \rho_{sub}^*$,

$$\begin{aligned} f(q_0, \rho) - V_0 &\geq 0, \\ \frac{\partial f}{\partial \rho}(q_0, \rho) &= \frac{\rho^2 p'(\rho) - q_0^2}{\rho^3} > 0, \\ \frac{\partial^2 f}{\partial \rho^2}(q_0, \rho) &= -\frac{3}{\rho} \frac{\partial f}{\partial \rho}(q_0, \rho) + \frac{1}{\rho^3} \frac{\partial}{\partial \rho} \left(\rho^2 p'(\rho) \right) \geq -\frac{3}{\rho} \frac{\partial f}{\partial \rho}(q_0, \rho). \end{aligned} \quad (3.44)$$

Thus for $\alpha \geq 3/2$, ψ is strictly increasing and convex. Therefore, the Newton method converges to the zero of the function and this proves the result. \square

3.3. Definition of left and right fluxes. In order to define the left and right numerical fluxes as (3.1), taking the definitions (3.28) for U_l^* , U_r^* (recall that Z^* is given by (3.5)), we need still to define T_l and T_r . We denote $(\mathcal{F}^\rho, \mathcal{F}^{\rho u}) \equiv \mathcal{F}(U_l^*, U_r^*)$, and define

$$\begin{aligned} T_l(U_l, U_r, \Delta Z) &= \frac{\rho_l - \rho_l^*}{\rho_l^*} \left(\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^\rho \right) - (u_l^* - u_l) \mathcal{F}^\rho \\ &\quad + \left(\left(e + \frac{p}{\rho} \right) (\rho_l^*) - \left(e + \frac{p}{\rho} \right) (\rho_l) + (\Delta Z)_+ + \frac{(u_l^*)^2}{2} - \frac{u_l^2}{2} \right) \frac{\mathcal{F}^\rho}{u_l}, \end{aligned} \quad (3.45)$$

$$\begin{aligned} T_r(U_l, U_r, \Delta Z) &= \frac{\rho_r - \rho_r^*}{\rho_r^*} \left(\mathcal{F}^{\rho u} - p(\rho_r^*) - u_r^* \mathcal{F}^\rho \right) - (u_r^* - u_r) \mathcal{F}^\rho \\ &\quad + \left(\left(e + \frac{p}{\rho} \right) (\rho_r^*) - \left(e + \frac{p}{\rho} \right) (\rho_r) + (-\Delta Z)_+ + \frac{(u_r^*)^2}{2} - \frac{u_r^2}{2} \right) \frac{\mathcal{F}^\rho}{u_r}. \end{aligned} \quad (3.46)$$

The definition (3.45) (respectively (3.46)) is ambiguous in the case when $u_l = 0$ or $\rho_l^* = 0$ (respectively when $u_r = 0$ or $\rho_r^* = 0$). In order to overcome this difficulty, we make here the convention that $0/0 = 0$. Then, one has to take into account the following remarks. They are stated for T_l , but of course similar statements hold for T_r .

1. If (ρ_l^*, u_l^*) solves the system (3.3) (with (3.5)), the factor of $\frac{\mathcal{F}^\rho}{u_l}$ in T_l vanishes.
2. If (ρ_l, u_l) is a supersonic point, we have $U_l^* = U_l$ and $T_l \equiv \frac{\mathcal{F}^\rho}{u_l} (\Delta Z)_+$, which is well defined since $u_l^2 > p'(\rho_l) \geq 0$.
3. If (ρ_l, u_l) is a subsonic point and $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_l) - (\Delta Z)_+ > m_s(\rho_l u_l)$, according to (B.2) in Lemma 3.8 we have that $\rho_l^* > 0$, and the factor of $\frac{\mathcal{F}^\rho}{u_l}$ in T_l vanishes, thus T_l is well defined. This is true in particular when (ρ_l, u_l) is a subsonic point if we consider a continuous bottom $z(x)$, which implies that ΔZ is small for a sufficiently fine grid.
4. If (ρ_l, u_l) is a subsonic point and $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_l) - (\Delta Z)_+ \leq m_s(\rho_l u_l)$, but $u_l \neq 0$, according to (B.3) in Lemma 3.8 we have $\rho_l^* = \rho_s(\rho_l u_l) > 0$, thus T_l is well defined.
5. Some difficulties may arise for (ρ_l, u_l) sonic close to $(0, 0)$, or for (ρ_l, u_l) subsonic with $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho} \right) (\rho_l) - (\Delta Z)_+ \leq m_s(\rho_l u_l)$ and u_l close to 0 (which implies also that (ρ_l, u_l) is close to $(0, 0)$). In these cases, the flux \mathcal{F} has to verify some conditions in order to define $\frac{\mathcal{F}^\rho}{u_l}$ and $(\rho_l^*)^{-1} (\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^\rho)$. As \mathcal{F} is consistent with F , we expect these quantities to be unambiguously defined.

Lemma 3.11. *The definitions (3.45), (3.46) of T_l and T_r imply that the conditions (3.12) are satisfied (with Z^* given by (3.5)).*

Proof. Consider the case of W_l . We have

$$u_l \frac{\rho_l - \rho_l^*}{\rho_l^*} = u_l^* - u_l, \quad (3.47)$$

thus

$$\begin{aligned}
u_l T_l &= \left(u_l^* - u_l \right) \left(\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^\rho - u_l \mathcal{F}^\rho \right) \\
&\quad + \left(\left(e + \frac{p}{\rho} \right) (\rho_l^*) - \left(e + \frac{p}{\rho} \right) (\rho_l) + (\Delta Z)_+ + \frac{(u_l^*)^2}{2} - \frac{u_l^2}{2} \right) \mathcal{F}^\rho \\
&= \left(u_l^* - u_l \right) \left(\mathcal{F}^{\rho u} - p(\rho_l^*) \right) \\
&\quad + \left(\left(e + \frac{p}{\rho} \right) (\rho_l^*) - \left(e + \frac{p}{\rho} \right) (\rho_l) + (\Delta Z)_+ - \frac{(u_l^*)^2}{2} + \frac{u_l^2}{2} \right) \mathcal{F}^\rho.
\end{aligned} \tag{3.48}$$

Putting this value in (3.10) gives $W_l = 0$. \square

3.4. Properties of the subsonic reconstruction scheme.

Theorem 3.12. *Let $\mathcal{F}(U_l, U_r)$ be a given consistent numerical flux for the Saint Venant problem without source that preserves nonnegativity of ρ by interface and satisfies a semi-discrete entropy inequality for the entropy pair (η, G) given by (2.8). Then the scheme (2.12), (2.13), with numerical fluxes (3.1), (3.28), (3.45), (3.46)*

(0) *is conservative in density,*

(i) *preserves the nonnegativity of ρ by interface,*

(ii) *preserves the discrete subsonic steady-states,*

(iii) *is consistent with the Saint Venant system away from sonic points,*

(iv) *satisfies a semi-discrete entropy inequality associated to the entropy pair $(\tilde{\eta}, \tilde{G})$ in (2.8).*

Proof. Notice first that when $\Delta Z = 0$ we have $U_l^* = U_l$, $U_r^* = U_r$, $T_l = 0$, $T_r = 0$, so that the scheme (3.1) reduces to the conservative scheme with numerical flux \mathcal{F} .

Property (0) is obvious from the definition (3.1) and the characterization (2.15). For (i), the assumption that $\mathcal{F}(U_l, U_r)$ preserves nonnegativity of ρ by interface means that (see [7]) there exists some $\sigma_l(U_l, U_r) < 0 < \sigma_r(U_l, U_r)$ such that

$$\rho_l + \frac{\mathcal{F}^\rho(U_l, U_r) - \rho_l u_l}{\sigma_l(U_l, U_r)} \geq 0, \quad \rho_r + \frac{\mathcal{F}^\rho(U_l, U_r) - \rho_r u_r}{\sigma_r(U_l, U_r)} \geq 0, \tag{3.49}$$

for any U_l and U_r (with nonnegative densities ρ_l, ρ_r). This implies in particular that

$$\rho_l^* + \frac{\mathcal{F}^\rho(U_l^*, U_r^*) - \rho_l^* u_l^*}{\sigma_l(U_l^*, U_r^*)} \geq 0, \quad \rho_r^* + \frac{\mathcal{F}^\rho(U_l^*, U_r^*) - \rho_r^* u_r^*}{\sigma_r(U_l^*, U_r^*)} \geq 0. \tag{3.50}$$

According to (3.28), one has $\rho_l^* u_l^* = \rho_l u_l$, $\rho_r^* u_r^* = \rho_r u_r$, and $\rho_l^* \leq \rho_l$, $\rho_r^* \leq \rho_r$, thus

$$\rho_l + \frac{\mathcal{F}^\rho(U_l^*, U_r^*) - \rho_l u_l}{\sigma_l(U_l^*, U_r^*)} \geq 0, \quad \rho_r + \frac{\mathcal{F}^\rho(U_l^*, U_r^*) - \rho_r u_r}{\sigma_r(U_l^*, U_r^*)} \geq 0, \tag{3.51}$$

proving that the scheme preserves the nonnegativity of ρ by interface. The associated speeds involved in the CFL condition are $\sigma_l(U_l^*, U_r^*)$, $\sigma_r(U_l^*, U_r^*)$.

In order to prove (ii), consider left and right states U_l, U_r , such that the steady state equations (2.10) are satisfied,

$$\begin{cases} \frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) + Z_l = \frac{u_r^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_r) + Z_r, \\ \rho_l u_l = \rho_r u_r, \end{cases} \quad (3.52)$$

and such that both are subsonic,

$$u_l^2 < p'(\rho_l), \quad u_r^2 < p'(\rho_r). \quad (3.53)$$

Note in particular that $\rho_l > 0, \rho_r > 0$. Recall that $\Delta Z \equiv Z_r - Z_l$. If $U_l = U_r$, then $\Delta Z = 0$ and according to the remark above one has $\mathcal{F}_l = \mathcal{F}(U_l, U_r) = F(U_l)$, $\mathcal{F}_r = \mathcal{F}(U_l, U_r) = F(U_r)$, proving (2.18). Assume now that $U_l \neq U_r$. Then $\Delta Z \neq 0$ otherwise (3.52) would give two subsonic solutions to a system (3.18), which is not possible by Proposition 3.5. Consider first the case when $\Delta Z > 0$. Then according to Lemma 3.8 we have

$$U_l^* = U_r^* = U_r, \quad (3.54)$$

and

$$\mathcal{F}(U_l^*, U_r^*) = F(U_r), \quad \mathcal{F}^{\rho u}(U_l^*, U_r^*) - p(\rho_l^*) - u_l^* \mathcal{F}^\rho(U_l^*, U_r^*) = 0, \quad (3.55)$$

thus

$$T_l = -(u_l^* - u_l) \rho_l^* u_l^*, \quad T_r = 0. \quad (3.56)$$

According to the remark after (3.2), the relations (2.18) are satisfied. The case $\Delta Z < 0$ is similar, with $U_l^* = U_r^* = U_l$. This proves (ii).

Property (iv) follows from Lemmas 3.1, 3.2, and 3.11.

It remains to prove the consistency (iii). The first property (2.16) is obvious according to the remark above on the case $\Delta Z = 0$. About (2.17), taking into account (3.1), we have to prove that

$$p(\rho_r) - p(\rho_r^*) + T_r - p(\rho_l) + p(\rho_l^*) - T_l = -\rho \Delta Z + o(\Delta Z), \quad (3.57)$$

as $U_l, U_r \rightarrow U$ and $\Delta Z \rightarrow 0$. As stated, we consider only the case when U is not sonic. Let us assume that $\Delta Z \geq 0$, the complementary case being similar. Then $U_r^* = U_r$ and $T_r = 0$.

(a) Case (ρ, u) supersonic. Then (ρ_l, u_l) is also supersonic if close enough to (ρ, u) (and in particular $u_l \neq 0$), and we have

$$\rho_l^* = \rho_l, \quad u_l^* = u_l, \quad \rho_r^* = \rho_r, \quad u_r^* = u_r, \quad (3.58)$$

$$\mathcal{F}_r^{\rho u} - \mathcal{F}_l^{\rho u} = -T_l = -\mathcal{F}^\rho(U_l, U_r) \frac{\Delta Z}{u_l} = -\rho \Delta Z + o(\Delta Z). \quad (3.59)$$

(b) Case (ρ, u) subsonic. Then $\rho > 0$, and $\frac{u^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho) > m_s(\rho u)$. Therefore, for U_l close enough to U and ΔZ small enough, we have (ρ_l, u_l) subsonic, $\rho_l > 0$, $\frac{u_l^2}{2} +$

$\left(e + \frac{p}{\rho}\right)(\rho_l) - \Delta Z > m_s(\rho_l u_l)$. From Lemma 3.8 (B.2) we have that $\rho_l^* > 0$, and we compute

$$\begin{aligned}\rho_l^* &= f_r^{-1}\left(\rho_l u_l, \frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) - \Delta Z\right) \\ &= f_r^{-1}\left(\rho_l u_l, f(\rho_l u_l, \rho_l) - \Delta Z\right) \\ &= \rho_l + O(\Delta Z),\end{aligned}\tag{3.60}$$

$$u_l^* = \frac{\rho_l u_l}{\rho_l^*} = u_l + O(\Delta Z).\tag{3.61}$$

Now, let us denote $\phi(\rho) = \left(e + \frac{p}{\rho}\right)(\rho)$. Since ϕ is a strictly increasing function, we can consider its inverse ϕ^{-1} . According to (3.60),

$$\left(e + \frac{p}{\rho}\right)(\rho_l^*) = \left(e + \frac{p}{\rho}\right)(\rho_l) + \frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} - \Delta Z,\tag{3.62}$$

thus using (3.61) we get

$$\begin{aligned}\rho_l^* &= \phi^{-1}\left(\phi(\rho_l) + \frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} - \Delta Z\right) \\ &= \rho_l + (\phi^{-1})'(\phi(\rho_l))\left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} - \Delta Z\right) + O\left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} - \Delta Z\right)^2 \\ &= \rho_l + \frac{\rho_l}{p'(\rho_l)}\left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} - \Delta Z\right) + O(\Delta Z)^2.\end{aligned}\tag{3.63}$$

Then,

$$\begin{aligned}p(\rho_l^*) &= p(\rho_l) + p'(\rho_l)(\rho_l^* - \rho_l) + O(\Delta Z)^2 \\ &= p(\rho_l) - \rho_l \Delta Z + \rho_l \left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2}\right) + O(\Delta Z)^2.\end{aligned}\tag{3.64}$$

We also have

$$\frac{\rho_l - \rho_l^*}{\rho_l^*} \left(\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^\rho\right) = O(\Delta Z) \cdot o(1) = o(\Delta Z).\tag{3.65}$$

Therefore,

$$\begin{aligned}\mathcal{F}_r^{\rho u} - \mathcal{F}_l^{\rho u} &= p(\rho_l^*) - p(\rho_l) - T_l \\ &= p(\rho_l^*) - p(\rho_l) - \frac{\rho_l - \rho_l^*}{\rho_l^*} \left(\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^\rho\right) + (u_l^* - u_l) \mathcal{F}^\rho \\ &= -\rho \Delta Z + \rho_l \left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2}\right) + (u_l^* - u_l) \mathcal{F}^\rho + o(\Delta Z) \\ &= -\rho \Delta Z + (u_l - u_l^*) \left(\rho_l \frac{u_l + u_l^*}{2} - \mathcal{F}^\rho\right) + o(\Delta Z) \\ &= -\rho \Delta Z + o(\Delta Z),\end{aligned}\tag{3.66}$$

which yields (iii). \square

3.5. Comments on the consistency at sonic points. The proof of (iii) in the previous theorem when we are close to a sonic point (ρ, u) involves some problems. Consider the case when $U_l, U_r \rightarrow U$, $\Delta Z \rightarrow 0$, with $u^2 = p'(\rho) \neq 0$. Assume as previously that $\Delta Z \geq 0$. If (ρ_l, u_l) is supersonic, then (3.59) is valid.

Otherwise, if (ρ_l, u_l) is subsonic and $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) - \Delta Z \geq m_s(\rho_l u_l)$, the computations made in the proof of consistency in the subsonic case can be followed, with the result

$$\mathcal{F}_r^{\rho u} - \mathcal{F}_l^{\rho u} = -\rho \Delta Z + O(\Delta Z) + O\left(|\Delta Z|^{1/2}|U_l - U_r|\right). \quad (3.67)$$

Assume now that either (ρ_l, u_l) is sonic, or (ρ_l, u_l) is subsonic with $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) - \Delta Z < m_s(\rho_l u_l)$. In any case we have $\frac{u_l^2}{2} + \left(e + \frac{p}{\rho}\right)(\rho_l) - \Delta Z \leq m_s(\rho_l u_l)$, thus $\rho_l^* = \rho_s(\rho_l u_l) > 0$, and since

$$0 \leq f(\rho_l u_l, \rho_l) - m_s(\rho_l u_l) \leq \Delta Z \quad (3.68)$$

and $\rho_s(\rho_l u_l)$ is the minimum point of the function $f(\rho_l u_l, \cdot)$, we deduce that

$$\rho_s(\rho_l u_l) - \rho_l = O(|\Delta Z|^{1/2}), \quad u_l^* = \frac{\rho_l u_l}{\rho_s(\rho_l u_l)} = u_l + O(|\Delta Z|^{1/2}). \quad (3.69)$$

Thus, assuming that the numerical flux \mathcal{F} is Lipschitz continuous,

$$\frac{\rho_l - \rho_l^*}{\rho_l^*} \left(\mathcal{F}^{\rho u} - p(\rho_l^*) - u_l^* \mathcal{F}^p \right) = O(\Delta Z) + O\left(|\Delta Z|^{1/2}|U_l - U_r|\right). \quad (3.70)$$

Now, from (3.68) we deduce $\left(e + \frac{p}{\rho}\right)(\rho_l^*) = \left(e + \frac{p}{\rho}\right)(\rho_l) + \frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} + O(\Delta Z)$, and using analogous calculations to the ones applied in the proof of the consistency in the subsonic case, one gets

$$\begin{aligned} \rho_l^* &= \rho_l + \frac{\rho_l}{p'(\rho_l)} \left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} \right) + O(\Delta Z), \\ p(\rho_l^*) &= p(\rho_l) + \rho_l \left(\frac{u_l^2}{2} - \frac{(u_l^*)^2}{2} \right) + O(\Delta Z), \end{aligned} \quad (3.71)$$

and

$$\begin{aligned} \mathcal{F}_r^{\rho u} - \mathcal{F}_l^{\rho u} &= p(\rho_l^*) - p(\rho_l) - T_l \\ &= p(\rho_l^*) - p(\rho_l) + (u_l^* - u_l) \mathcal{F}^p + O(\Delta Z) + O\left(|\Delta Z|^{1/2}|U_l - U_r|\right) \\ &= (u_l - u_l^*) \left(\rho_l \frac{u_l + u_l^*}{2} - \mathcal{F}^p \right) + O(\Delta Z) + O\left(|\Delta Z|^{1/2}|U_l - U_r|\right) \\ &= O(\Delta Z) + O\left(|\Delta Z|^{1/2}|U_l - U_r|\right). \end{aligned} \quad (3.72)$$

Thus in any case, (3.67) is valid.

This property (3.67) does not mean consistency in the sense of (2.17), but anyway, one can make the following observation. Assume that at the point considered, one has $dZ/dx = 0$. Then $\Delta Z = o(\Delta x)$, $U_r - U_l = O(\Delta x)$, and therefore (3.67) yields

$$\mathcal{F}_r^{\rho u} - \mathcal{F}_l^{\rho u} = o(\Delta x), \quad (3.73)$$

which means consistency with the vanishing source. Since the condition $dZ/dx = 0$ is generically satisfied at sonic points (see [10]), it justifies the global consistency of the scheme, except maybe close to the point $(\rho, u) = (0, 0)$. We shall see in the numerical computations that even if the numerical fluxes can sometimes take large values, the scheme behaves reasonably well in the presence of data close to $(\rho, u) = (0, 0)$.

4. APPLICATION TO THE EULER-POISSON SYSTEM

Let us consider the Euler-Poisson system

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = -\rho \partial_x \phi, \\ -\partial_{xx}^2 \phi = \rho - \rho_b, \end{cases} \quad (4.1)$$

where $\rho_b \equiv \rho_b(x) \geq 0$ is given. The system (4.1) is set for $t > 0$, $0 < x < l$, with initial and boundary conditions

$$\begin{cases} \rho(t=0, \cdot) = \rho_b, \\ u(t=0, \cdot) = u_0, \\ \phi(l) - \phi(0) = V, \\ \rho u(t, x=0) = q_0 \geq 0, \\ \rho(t, x=l) = \rho_b(l). \end{cases} \quad (4.2)$$

As usual, one has to complete (4.1) by an entropy inequality. In order to describe the steady-states, we subtract u times the first equation in (4.1) to the second, we divide the result by ρ , and get

$$\partial_t u + \partial_x(u^2/2 + e(\rho) + p(\rho)/\rho + \phi) = 0. \quad (4.3)$$

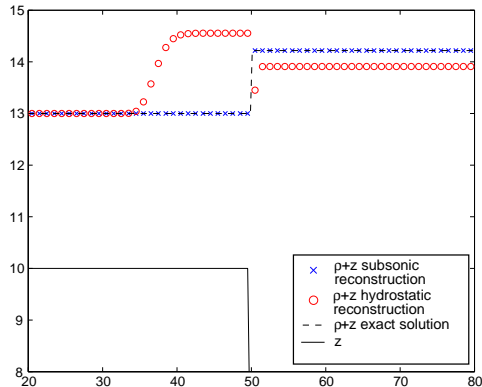
Therefore, the steady-states are determined by the relations

$$\begin{cases} \rho u = Cst, \\ \frac{u^2}{2} + e(\rho) + \frac{p(\rho)}{\rho} + \phi = Cst. \end{cases} \quad (4.4)$$

We can observe that the Euler-Poisson system is of the type (2.1), where the bottom Z is replaced by a function ϕ that is time-dependent.

The subsonic reconstruction scheme can be applied to this system by "freezing" the potential on a time interval, as follows. Given an approximation of ϕ at time t_n , $\phi^n = \phi(t_n, \cdot)$, we solve the system

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = -\rho \partial_x \phi^n \end{cases} \quad (4.5)$$

FIGURE 2. Subsonic steady-state, $\rho + z$ at $t = 5$

in the interval $[t_n, t_{n+1})$ using the subsonic reconstruction scheme, where ϕ^n stands for Z . We obtain approximations ρ^{n+1} , u^{n+1} . Finally, we solve the ODE

$$-\partial_{xx}^2 \phi^{n+1} = \rho^{n+1} - \rho_b \quad (4.6)$$

in order to get the new potential. It is obvious that this algorithm is well-balanced, since the freezing of the potential does not introduce any error in the case of a steady state.

The interest of the subsonic reconstruction scheme in this context is the ability to compute with high accuracy flows which are close to a subsonic steady state with constant discharge $\rho u \neq 0$.

5. NUMERICAL RESULTS

5.1. Saint-Venant system. In order to evaluate our method, we compare the subsonic reconstruction scheme described here to the original hydrostatic reconstruction scheme of [1]. We use first-order resolution, with the CFL 1 condition induced by the nonnegativity of density (see the proof of (i) in Theorem 3.12). A second-order extension can be used as in [1], and in that case no major differences are observed between the two reconstructions. The numerical flux \mathcal{F} chosen here is the one obtained from the Suliciu relaxation system described in [7]. We take $p(\rho) = g\rho^2/2$, $g = 9.81$ and use 100 points in the considered interval in each case.

We consider first a subsonic steady-state in the interval $(0, 100)$. The initial data are given by

$$\rho_0(x) = \begin{cases} 3 & \text{if } x \leq 50, \\ 14.2175 & \text{if } x > 50, \end{cases} \quad u_0(x) = \begin{cases} 5 & \text{if } x \leq 50, \\ 1.055 & \text{if } x > 50, \end{cases} \quad (5.1)$$

$$z(x) = \begin{cases} 10 & \text{if } x \leq 50, \\ 0 & \text{if } x > 50. \end{cases} \quad (5.2)$$

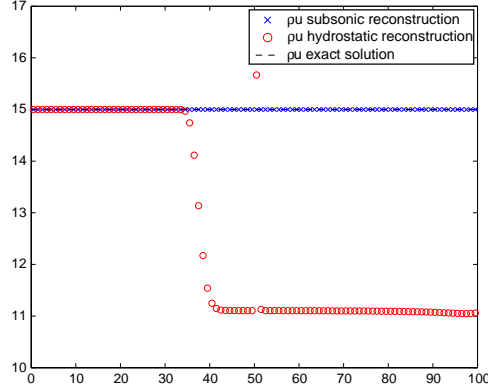


FIGURE 3. Subsonic steady-state, ρu at $t = 5$

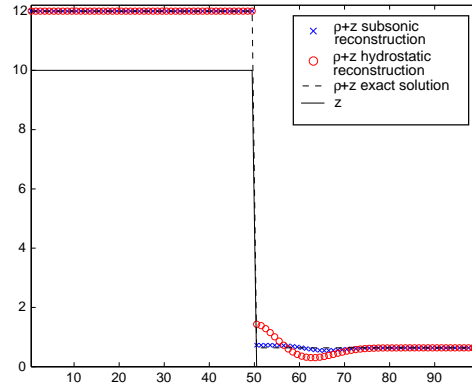


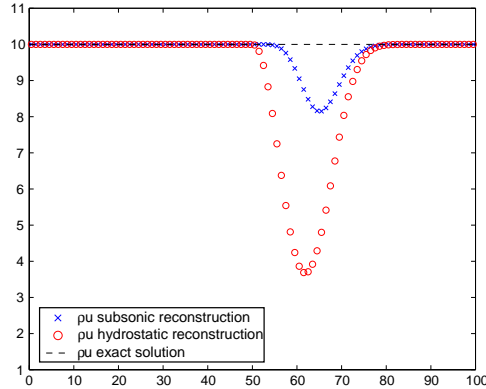
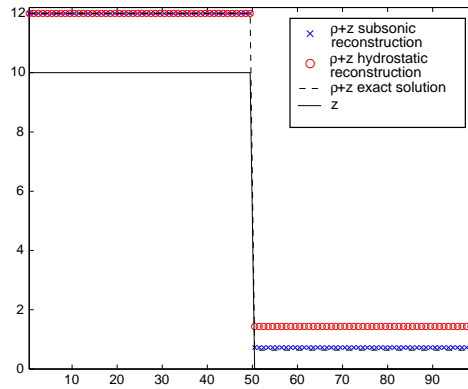
FIGURE 4. Supersonic steady-state, $\rho + z$ at $t = 1$

The results are shown on Figures 2 and 3. As we see, the subsonic reconstruction scheme maintains the subsonic steady-state, as we have proved. This is an improvement with respect to the hydrostatic reconstruction scheme which does not.

Then, we consider a supersonic steady-state, still in the interval $(0, 100)$,

$$\rho_0(x) = \begin{cases} 2 & \text{if } x \leq 50, \\ 0.635 & \text{if } x > 50, \end{cases} \quad u_0(x) = \begin{cases} 5 & \text{if } x \leq 50, \\ 15.7474 & \text{if } x > 50, \end{cases} \quad (5.3)$$

$$z(x) = \begin{cases} 10 & \text{if } x \leq 50, \\ 0 & \text{if } x > 50. \end{cases} \quad (5.4)$$

FIGURE 5. Supersonic steady-state, ρu at $t = 1$ FIGURE 6. Supersonic steady-state, $\rho + z$ at $t = 25$

We see in Figures 4, 5, 6, 7, that neither of the two schemes gives the right solution, but the subsonic reconstruction scheme is more accurate.

We consider now a classical transcritical shock test. The space domain is $(0, 25)$, the initial data are $\rho_0(x) = 0.33$, $u_0(x) = 0.18/0.33$ and the topography is

$$z(x) = \begin{cases} 0.2 - 0.05(x - 10)^2 & \text{if } 8 < x < 12, \\ 0 & \text{otherwise.} \end{cases} \quad (5.5)$$

The boundary conditions are taken $\rho u(x = 0) = 0.18$ and $\rho(x = 25) = 0.33$. The results are shown on Figures 8 and 9. We see that the subsonic reconstruction scheme gives a solution which is sharper on the left part of the discontinuity than the hydrostatic reconstruction.

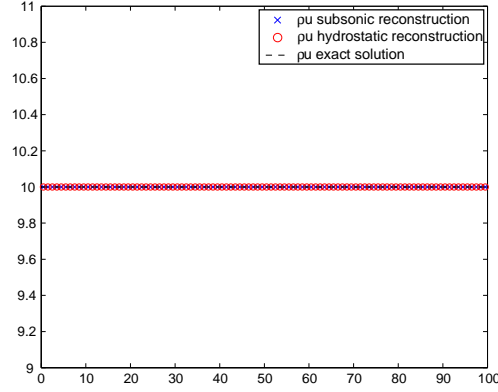


FIGURE 7. Supersonic steady-state, ρu at $t = 25$

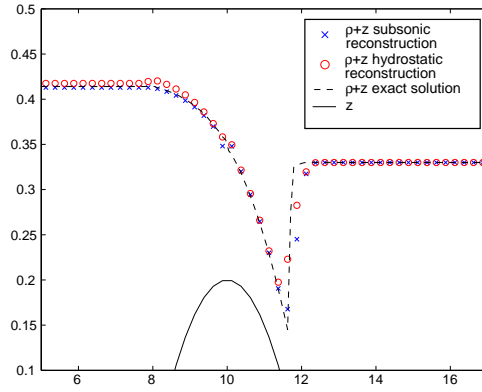


FIGURE 8. Transcritical flow with shock, $\rho + z$ at $t = 200$

5.2. Euler-Poisson System. We solve the Euler-Poisson system (4.1), (4.2) with $p(\rho) = \kappa\rho^\gamma$,

$$\kappa = 1, \quad \gamma = 1.1, \tag{5.6}$$

for $x \in (0, 0.6)$ and with initial and boundary conditions

$$\rho_b(x) = \begin{cases} 1 & \text{if } x \in (0.1, 0.5), \\ 100 & \text{otherwise,} \end{cases} \tag{5.7}$$

$$u_0 = 0, \quad V = -1.$$

In Figure 10 we show the result with the boundary discharge $q_0 = 10$. We use 100 points in space and the final time is $t = 100$. A second-order reconstruction is used. As we see, we have reached a subsonic equilibrium where $u^2/2 + e + p/\rho + \phi$ is constant and q is almost constant.

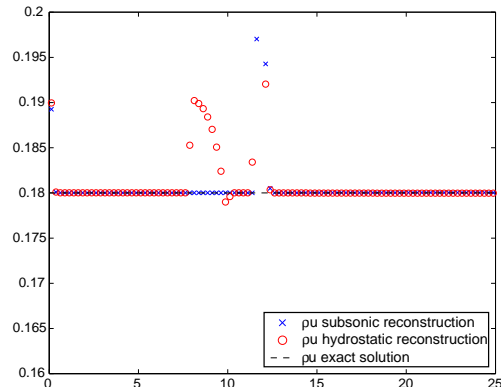


FIGURE 9. Transcritical flow with shock, ρu at $t = 200$

Finally the same test with boundary discharge $q_0 = 40$ is shown on Figure 11. Even if we have reached a steady-state, q and $u^2/2 + e + p/\rho + \phi$ are not constant. We observe a jump at the point where there is a change from a supersonic regime to a subsonic one. This is not surprising since the scheme is not exact for supersonic states.

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 (2004), 2050–2065.
- [2] E. Audusse, M.-O. Bristeau, and B. Perthame, Kinetic schemes for Saint-Venant equations with source terms on unstructured grids, *INRIA Report*, RR-3989, 2000.
- [3] D. S. Bale, R. J. Leveque, S. Mitran, and J. A. Rossmannith, A wave propagation method for conservation laws and balance laws with spatially varying flux functions, *SIAM J. Sci. Comput.*, 24 (2002), 955–978.
- [4] A. Bermúdez, and M.E. Vázquez, Upwind methods for hyperbolic conservation laws with source terms, *Comput. Fluids*, 23 (1994), 1049–1071.
- [5] R. Botchorishvili, B. Perthame, and A. Vasseur, Equilibrium schemes for scalar conservation laws with stiff sources, *Math. Comp.*, 72 (2003), 131–157.
- [6] N. Botta, R. Klein, S. Langenberg, and S. Lützenkirchen, Well balanced finite volume methods for nearly hydrostatic flows, *J. Comput. Phys.*, 196 (2004), 539–565.
- [7] F. Bouchut, *Nonlinear stability of finite volume methods for hyperbolic conservation laws and well-balanced schemes for sources*, Frontiers in Mathematics, Birkhäuser Verlag, Basel, 2004.
- [8] M. Castro, J. Macías, and C. Parés, A Q -scheme for a class of systems of coupled conservation laws with source term, application to a two-layer 1-D shallow water system, *M2AN Math. Model. Numer. Anal.*, 35 (2001), 107–127.
- [9] M. Castro, A. Pardo Milanés, C. Parés, Well-balanced numerical schemes based on a generalized hydrostatic reconstruction technique, *Math. Models Methods Appl. Sci.*, 17 (2007), 2055–2113.
- [10] A. Chinnayya, A.-Y. LeRoux, and N. Seguin, A well-balanced numerical scheme for the approximation of the shallow-water equations with topography: the resonance phenomenon, *Int. J. Finite Volume*, 1 (2004), 1–33.

- [11] T. Gallouët, J.-M. Hérard, and N. Seguin, Some approximate Godunov schemes to compute shallow-water equations with topography, *Comput. & Fluids*, 32 (2003), 479–513.
- [12] L. Gosse, A well-balanced flux-vector splitting scheme designed for hyperbolic systems of conservation laws with source terms, *Comput. Math. Appl.*, 39 (2000), 135–159.
- [13] L. Gosse and A.-Y. Leroux, Un schéma-équilibre adapté aux lois de conservation scalaires non-homogènes, *C. R. Acad. Sci. Paris Sér. I Math.*, 323 (1996), 543–546.
- [14] J. M. Greenberg and A. Y. Leroux, A well-balanced scheme for the numerical processing of source terms in hyperbolic equations, *SIAM J. Numer. Anal.*, 33 (1996), 1–16.
- [15] S. Jin, A steady-state capturing method for hyperbolic systems with geometrical source terms, *M2AN Math. Model. Numer. Anal.*, 35 (2001), 631–645.
- [16] A. Kurganov, and D. Levy, Central-upwind schemes for the Saint-Venant system, *M2AN Math. Model. Numer. Anal.*, 36 (2002), 397–425.
- [17] B. Perthame and C. Simeoni, A kinetic scheme for the Saint-Venant system with a source term, *Calcolo*, 38 (2001), 201–231.
- [18] M. E. Vázquez-Cendón, Improved treatment of source terms in upwind schemes for the shallow water equations in channels with irregular geometry, *J. Comput. Phys.*, 148 (1999), 497–526.

FRANÇOIS BOUCHUT, DMA, CNRS & ÉCOLE NORMALE SUPÉRIEURE, 45 RUE D'ULM, F-75230 PARIS CEDEX 05, FRANCE

E-mail address: Francois.Bouchut@ens.fr

TOMÁS MORALES, DEPARTAMENTO DE MATEMÁTICAS, EDIFICIO ALBERT EINSTEIN (C2), CAMPUS DE RABANALES, UNIVERSIDAD DE CÓRDOBA, 14071-CÓRDOBA, SPAIN

E-mail address: ma1molut@uco.es

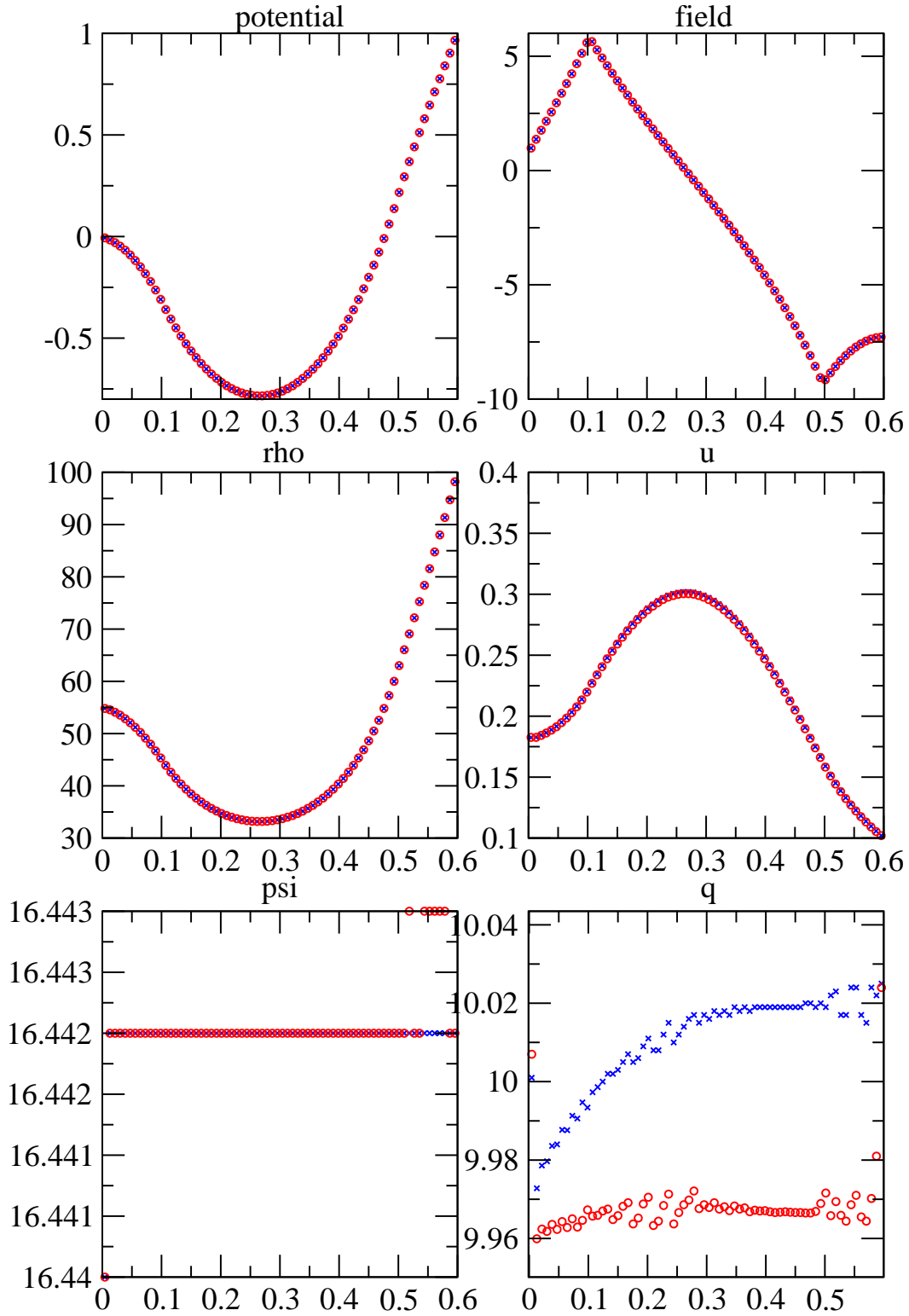


FIGURE 10. Crosses: subsonic reconstruction, circles: hydrostatic reconstruction. The quantity ψ represents $u^2/2 + e + p/\rho + \phi$, and $q = \rho u$.

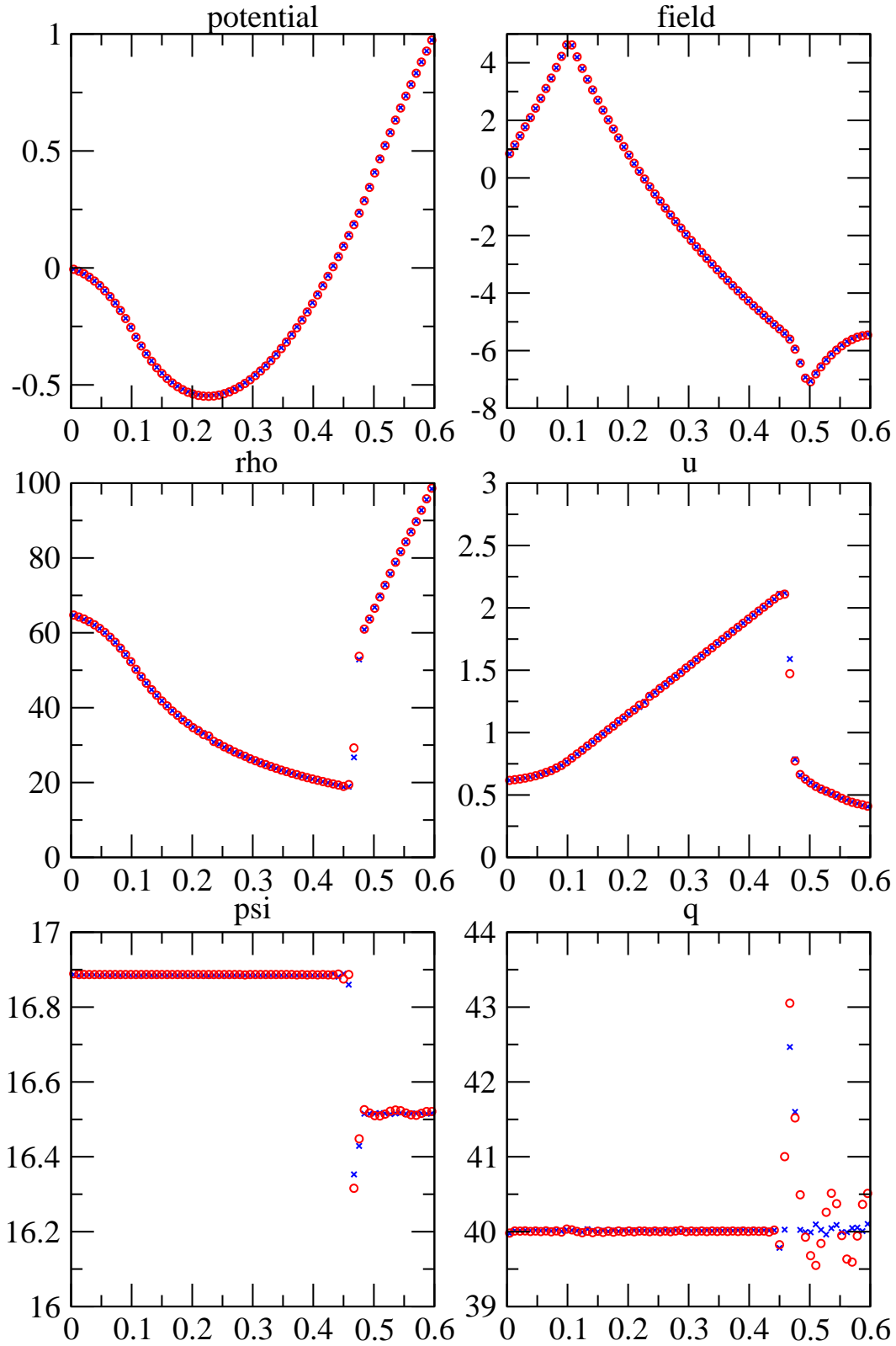


FIGURE 11. Crosses: subsonic reconstruction, circles: hydrostatic reconstruction. The quantity ψ represents $u^2/2 + e + p/\rho + \phi$, and $q = \rho u$.