

## ON THE WELL-BALANCE PROPERTY OF ROE'S METHOD FOR NONCONSERVATIVE HYPERBOLIC SYSTEMS. APPLICATIONS TO SHALLOW-WATER SYSTEMS

CARLOS PARÉS<sup>1</sup> AND MANUEL CASTRO<sup>1</sup>

**Abstract.** This paper is concerned with the numerical approximations of Cauchy problems for one-dimensional nonconservative hyperbolic systems. The first goal is to introduce a general concept of well-balancing for numerical schemes solving this kind of systems. Once this concept stated, we investigate the well-balance properties of numerical schemes based on the generalized Roe linearizations introduced by [Toumi, *J. Comp. Phys.* **102** (1992) 360–373]. Next, this general theory is applied to obtain well-balanced schemes for solving coupled systems of conservation laws with source terms. Finally, we focus on applications to shallow water systems: the numerical schemes obtained and their properties are compared, in the case of one layer flows, with those introduced by [Bermúdez and Vázquez-Cendón, *Comput. Fluids* **23** (1994) 1049–1071]; in the case of two layer flows, they are compared with the numerical scheme presented by [Castro, Macías and Parés, *ESAIM: M2AN* **35** (2001) 107–127].

**Mathematics Subject Classification.** 65M99, 76B55, 76B70.

Received: November 3, 2003. Revised: June 9, 2004.

### 1. INTRODUCTION

This paper is concerned with the numerical approximations of Cauchy problems for one-dimensional nonconservative hyperbolic systems:

$$\frac{\partial W}{\partial t} + \mathcal{A}(W) \frac{\partial W}{\partial x} = 0, \quad x \in \mathbb{R}, t > 0. \quad (1.1)$$

A first difficulty related with these systems comes from the presence of nonconservative products in (1.1) which makes difficult the definition of weak solutions. Dal Masso, LeFloch, and Murat [11] proposed an interpretation of these products using a family of paths drawn in the phases space.

Some of the usual numerical schemes designed for conservation laws can be adapted to the discretization of the more general system (1.1). This is the case of numerical schemes based on Approximate Riemann Solvers and, in particular, Roe schemes (see [16, 23, 29, 32], . . .): in [35] a general definition of Roe linearizations was introduced, based again on the use of a family of paths.

---

*Keywords and phrases.* Nonconservative hyperbolic systems, well-balanced schemes, Roe method, source terms, shallow-water systems.

<sup>1</sup> Dpto. Análisis Matemático, Facultad de Ciencias, Universidad de Málaga, Campus de Teatinos s/n, 29080-Málaga, Spain.  
e-mail: grupo@anamat.cie.uma.es

Following an idea introduced in [17, 19, 20], a system of conservation laws with source terms or *balance* law

$$\frac{\partial W}{\partial t} + \frac{\partial F}{\partial x}(W, \sigma) = \tilde{S}(W, \sigma) \frac{d\sigma}{dx}, \quad (1.2)$$

$\sigma(x)$  being a known function, can be considered as a particular case of (1.1), if the trivial equation:

$$\frac{\partial \sigma}{\partial t} = 0,$$

is added to the system. In particular, the hyperbolic shallow water system with source terms due to bed elevations or breadth variations can be formulated under this form.

It is known that, in the presence of source terms, standard methods can fail in approximating steady or nearly steady flows (see [25, 30, 33], ...). In the context of shallow water equations Bermúdez and Vázquez-Cendón have shown in [3, 36] that methods based on Approximate Riemann Solvers for the discretization of the flux terms and *upwinding* the source terms suitably solve these difficulties. These authors introduced the condition called *conservation property* or *C-property*: a scheme is said to satisfy this condition if it solves, exactly or up to the second order, the steady state solutions corresponding to water at rest. This idea of constructing numerical schemes that preserve some equilibria, which are called in general *well-balanced* schemes, has been extended in different ways: see [2, 6–9, 12, 14, 17–20, 24, 26, 27, 34, 37].

A more general type of systems that can be considered as a particular case of (1.1) is the class of coupled systems of conservation laws with source terms of the form:

$$\frac{\partial W_k}{\partial t} + \frac{\partial F_k}{\partial x}(W_k, \sigma) = \sum_{l \neq k} \mathcal{B}_{k,l}(W_1, \dots, W_K, \sigma) \cdot \frac{\partial W_l}{\partial x} + \tilde{S}_k(W_1, \dots, W_K, \sigma) \frac{d\sigma}{dx}, \quad k = 1, \dots, K. \quad (1.3)$$

The equations governing the flow of a stratified fluid composed by two superposed shallow layers of immiscible liquids can be formulated under this form (see [4]). In the cited work, some numerical schemes generalizing those introduced in [3, 36] for balance laws, were presented. The well-balance property required to the numerical schemes was the natural extension of the *C-property*: equilibria corresponding to water at rest had to be preserved. Recently, this work has been extended to the more complex system corresponding to 1d bilayer shallow water flows in symmetric channels with irregular geometries in [5]. In [12] a more general definition of the well-balance property for numerical schemes solving (1.3) has been introduced, as well as a family of numerical schemes satisfying this extended property.

Systems with similar characteristics also appear in other flow models such as boiling flows and two-phase flows (see [13]).

The main goal of this work is to provide a theoretical framework for a general definition of *well-balanced* numerical schemes, and a methodology for the construction of well-balanced Roe schemes with special emphasis on the particular case of coupled systems of conservation laws with source terms.

This paper is organized as follows. In Section 2, we recall briefly the notion of a weak solution of (1.1) and the general form of a Roe scheme as proposed in [11] and [35], respectively. In Section 3, we introduce a general definition of well-balanced numerical schemes for solving (1.1) and we give some general results concerning the well-balance properties of Roe schemes. In particular, we pay attention to Roe schemes based on the simplest choice of paths: the family of segments. In Section 4, coupled systems of conservations laws with source terms (1.3) are considered: on the basis of the general results of Section 3, we discuss firstly how to construct a Roe scheme when Roe matrices are known for each particular conservation law involved. Then, we describe how these schemes can be adapted to the more difficult case of *resonant problems* in which the system becomes nonstrictly hyperbolic.

In Section 5 this general methodology is applied to some systems related to shallow water flows, recovering some known well-balanced schemes, or resulting in new schemes. In the case of flows through channels with rectangular cross-section of constant breadth and varying depth, the *Q*-scheme of Roe introduced in [3, 36] is

recovered for one-layer flows, and the generalized  $Q$ -scheme of Roe presented in [4] for two-layer flows. The previously known results concerning the well-balanced character of these schemes are easily deduced from the general theory. In the case of flows through channels with rectangular cross-section of varying breadth and constant depth, we obtain two different schemes depending on the choice of the family of paths. One of them is similar to the numerical scheme presented in [14], but the other has not been yet described, in our knowledge.

In view of these results, the methodology developed seems to be useful for designing Roe-type schemes that are well-balanced for general nonconservative hyperbolic systems, and it can be also a first step for designing higher order well-balanced schemes for this kind of systems.

## 2. ROE METHODS FOR NONCONSERVATIVE SYSTEMS

We consider the problem:

$$\frac{\partial W}{\partial t} + \mathcal{A}(W) \frac{\partial W}{\partial x} = 0, \quad x \in \mathbb{R}, t > 0. \tag{2.4}$$

We suppose that the unknown function  $W(x, t)$  takes its values inside an open convex set  $\Omega$  included in  $\mathbb{R}^N$  and that  $W \in \Omega \rightarrow \mathcal{A}(W)$  is a smooth locally bounded map. We suppose that system (2.4) is strictly hyperbolic, that is, for each  $W \in \Omega$  the matrix  $\mathcal{A}(W)$  has  $N$  real distinct eigenvalues

$$\lambda_1(W) < \dots < \lambda_N(W),$$

and associated eigenvectors  $R_j(W)$ ,  $j = 1, \dots, N$ . Moreover, we suppose that for each integer  $j \in \{1, \dots, N\}$  the  $j$ th field is either genuinely nonlinear:

$$R_j(W) \cdot \nabla \lambda_j(W) \neq 0, \forall W \in \Omega,$$

or linearly degenerate:

$$R_j(W) \cdot \nabla \lambda_j(W) = 0, \forall W \in \Omega.$$

Observe that the system contains a nonconservative product  $\mathcal{A}(W) \cdot W_x$  which, in general, cannot make sense within the framework of the theory of distributions. After the theory developed by Dal Masso, LeFloch, and Murat [11], a rigorous definition of weak solutions can be performed using a family of paths in  $\Omega$  defined as follows:

**Definition 1.** A family of paths in  $\Omega \subset \mathbb{R}^N$  is a locally Lipschitz map  $\Phi : [0, 1] \times \Omega \times \Omega \rightarrow \Omega$  such that:

- $\Phi(0; W_L, W_R) = W_L$ ,  $\Phi(1; W_L, W_R) = W_R$  for any  $W_L, W_R$  in  $\Omega$ ;
- for every bounded set  $\mathcal{O}$  of  $\Omega$  there exists  $k$  such that:

$$\left| \frac{\partial \Phi}{\partial s}(s; W_L, W_R) \right| \leq k |W_L - W_R| \tag{2.5}$$

for any  $W_L, W_R \in \mathcal{O}$  and  $s \in [0, 1]$ ;

- for every bounded set  $\mathcal{O}$  of  $\Omega$  there exists  $K$  such that:

$$\left| \frac{\partial \Phi}{\partial s}(s; W_L^1, W_R^1) - \frac{\partial \Phi}{\partial s}(s; W_L^2, W_R^2) \right| \leq K (|W_L^1 - W_L^2| + |W_R^1 - W_R^2|)$$

for any  $W_L^1, W_R^1, W_L^2, W_R^2 \in \mathcal{O}$  and  $s \in [0, 1]$ .

Once a family of paths chosen, given a function  $W \in (L^\infty(\mathbb{R} \times \mathbb{R}^+) \cup BV(\mathbb{R} \times \mathbb{R}^+))^N$  it is possible to give a sense to the nonconservative product as a Borel measure (see [11] for details), which is denoted  $[\mathcal{A}(W) \cdot W_x]_\Phi$  and weak solutions are the functions satisfying the equality

$$W_t + [\mathcal{A}(W) \cdot W_x]_\Phi = 0.$$

Across a discontinuity, weak solutions satisfy the generalized Rankine-Hugoniot condition:

$$\int_0^1 (\xi \cdot \mathcal{I} - \mathcal{A}(\Phi(s; W^-, W^+))) \frac{\partial \Phi}{\partial s}(s; W^-, W^+) ds = 0, \tag{2.6}$$

where  $\xi$  denotes the speed of the discontinuity,  $\mathcal{I}$  the identity matrix, and  $W^-, W^+$  the limits to the left and to the right of the solution.

Observe that in the particular case of a conservation law, *i.e.* if there exists a smooth function  $F : \Omega \rightarrow \mathbb{R}$  such that, for any  $W \in \Omega$ ,  $\mathcal{A}(W)$  is the Jacobian matrix of  $F$ , condition (2.6) coincides with the usual Rankine-Hugoniot conditions for any choice of the family of paths.

Together with this definition of weak solutions, we assume here the Lax’s concept of *entropic solution*:

**Definition 2.** A weak solution of (2.4) will be said an *entropic solution* if, at each discontinuity  $\Sigma$  there exists  $k \in \{1, \dots, N\}$  such that:

$$\lambda_k(W^+) < \xi < \lambda_{k+1}(W^+); \quad \lambda_{k-1}(W^-) < \xi < \lambda_k(W^-) \tag{2.7}$$

if the characteristic field is genuinely nonlinear or

$$\lambda_k(W^-) = \xi = \lambda_k(W^+) \tag{2.8}$$

if the characteristic field is linearly degenerate.

These definitions of weak and entropic solution allow to extend to systems (2.4) the theory of simple waves of hyperbolic systems of conservation laws and the results concerning the solutions of Riemann problems (see [11]).

The choice of the family of paths is important as it determines the propagation speed of shocks. The simplest choice is given by the family of segments:

$$\Phi(s; W_L, W_R) = W_L + s(W_R - W_L), \tag{2.9}$$

that corresponds to the definition of nonconservative product proposed by Volpert [38]. In general, in practical applications, its selection has to be based on the physical background (see [22, 28] for instance). Nevertheless, in the particular case in which  $W_L$  and  $W_R$  are linked by an integral curve of a linearly degenerate field, the natural choice of  $\Phi(s; W_L, W_R)$  is a parameterization of the arc of the curve delimited by  $W_L$  and  $W_R$ . In effect, this choice ensures that the contact discontinuity:

$$W(x, t) = \begin{cases} W_L & \text{if } x < \xi t, \\ W_R & \text{if } x > \xi t, \end{cases}$$

where  $\xi$  is the (constant) value of the corresponding eigenvalue through the integral curve, verifies (2.6) and thus is a weak solution of (2.4).

In [35] a generalization of Roe methods to systems of the form (2.4) was introduced. These methods are based on the following general definition of a *Roe linearization*:

**Definition 3.** Given a family of paths  $\Psi$ , a matrix function  $\mathcal{A}_\Psi : \Omega \times \Omega \rightarrow \mathcal{M}_N(\mathbb{R})$  is called a Roe linearization if it satisfies:

- for any  $W_L, W_R \in \Omega$ ,  $\mathcal{A}_\Psi(W_L, W_R)$  has  $N$  real distinct eigenvalues;
- $\mathcal{A}_\Psi(W, W) = \mathcal{A}(W)$ , for all  $W \in \Omega$ ;
- for any  $W_L, W_R \in \Omega$ :

$$\mathcal{A}_\Psi(W_L, W_R) \cdot (W_R - W_L) = \int_0^1 \mathcal{A}(\Psi(s; W_L, W_R)) \frac{\partial \Psi}{\partial s}(s; W_L, W_R) ds. \tag{2.10}$$

Once the linearization chosen, to discretize the system a set of computing cells  $I_i = [x_{i-1/2}, x_{i+1/2}]$ ,  $i \in \mathbb{Z}$  is chosen. For the sake of simplicity, we assume that these cells have a constant size  $\Delta x$ , and that

$$x_{i+1/2} = i\Delta x.$$

$x_i = (i - 1/2)\Delta x$  is the center of the cell  $I_i$ . Let  $\Delta t$  be the time step and  $t^n = n\Delta t$ .

As usual, we denote by  $W_i^n$  the approximation of the cell averages of the exact solution provided by the numerical scheme:

$$W_i^n \cong \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} W(x, t^n) dx.$$

In order to calculate these approximations, we first introduce the *intermediate* matrices

$$\mathcal{A}_{i+1/2} = \mathcal{A}_\Psi(W_i^n, W_{i+1}^n). \tag{2.11}$$

The eigenvalues of this matrix will be denoted by:

$$\lambda_1^{i+1/2} < \lambda_2^{i+1/2} < \dots < \lambda_N^{i+1/2},$$

and a set of associated eigenvectors by  $\{R_i^{i+1/2}\}_{i=1}^N$ . We will denote by  $\mathcal{K}_{i+1/2}$  the  $N \times N$  matrix whose columns are these eigenvectors and by  $\mathcal{L}_{i+1/2}$  the diagonal matrix whose coefficients are the eigenvalues. We also introduce the matrices  $\mathcal{L}_{i+1/2}^+$ ,  $\mathcal{L}_{i+1/2}^-$ ,  $\mathcal{A}_{i+1/2}^+$ ,  $\mathcal{A}_{i+1/2}^-$  as usual

$$\mathcal{L}_{i+1/2}^\pm = \begin{bmatrix} (\lambda_1^{i+1/2})^\pm & & 0 \\ & \ddots & \\ 0 & & (\lambda_N^{i+1/2})^\pm \end{bmatrix}, \quad \mathcal{A}_{i+1/2}^\pm = \mathcal{K}_{i+1/2} \mathcal{L}_{i+1/2}^\pm \mathcal{K}_{i+1/2}^{-1}. \tag{2.12}$$

The numerical scheme progresses in time as follows: once the approximations at time  $t^n$ ,  $W_i^n$ , have been calculated, a Linear Riemann problem is considered at each intercell  $x_{i+1/2}$  whose matrix is  $\mathcal{A}_{i+1/2}$  and whose states to the left and to the right are respectively  $W_i^n$  and  $W_{i+1}^n$ . The approximations at the time  $t^{n+1}$ ,  $W_i^{n+1}$ , are obtained by averaging in the cells the solutions of these linear problems. As in the case of systems of conservation laws, some calculations allow to show that, under the hypothesis:

$$x_{i-1/2} + \lambda_N^{i-1/2} \Delta t \leq x_i \leq x_{i+1/2} + \lambda_1^{i+1/2} \Delta t, \tag{2.13}$$

the approximations at time  $t^{n+1}$  can be obtained by the formula:

$$W_i^{n+1} = W_i^n - \frac{\Delta t}{\Delta x} \left( \mathcal{A}_{i-1/2}^+ \cdot (W_i^n - W_{i-1}^n) + \mathcal{A}_{i+1/2}^- \cdot (W_{i+1}^n - W_i^n) \right), \tag{2.14}$$

which is the general expression of a Roe scheme for (2.4).

The best choice of the family of paths  $\Psi$  appearing in the definition of the Roe matrix is the family  $\Phi$  selected for the definition of weak solutions: this choice assures that, if the states  $W_i^n$  and  $W_{i+1}^n$  can be linked by a discontinuity satisfying (2.6) that propagates at velocity  $\xi$  then:

$$(\mathcal{A}_{i+1/2} - \xi \cdot \mathcal{I}) \cdot (W_{i+1}^n - W_i^n) = 0, \tag{2.15}$$

that is, the difference of the states is an eigenvector of the intermediate matrix associated to  $\xi$ . As a consequence, the solution of the linear Riemann problem agrees with the exact solution.

Nevertheless the construction of a Roe linearization based on the family of paths  $\Phi$  can be difficult in practice, as we shall see in the applications. Therefore, we are mainly interested in the more general case in which  $\Psi \neq \Phi$ . In particular, we will pay attention to the simplest choice:

$$\Psi(s; W_L, W_R) = W_L + s(W_R - W_L), \tag{2.16}$$

which coincides with the family of paths used in Volpert’s definition (2.9).

If  $\Phi \neq \Psi$  and the states  $W_i^n$  and  $W_{i+1}^n$  can be linked by a shock or a contact discontinuity that propagates at velocity  $\xi$ , the identity (2.15) doesn’t hold in general. In that case, the following estimate can be easily obtained from (2.5) and (2.6):

$$\begin{aligned} |(\mathcal{A}_{i+1/2} - \xi \cdot \mathcal{I}) \cdot (W_{i+1}^n - W_i^n)| &\leq K_1 \int_0^1 \left| \frac{\partial \Phi}{\partial s}(s; W_i^n, W_{i+1}^n) - \frac{\partial \Psi}{\partial s}(s; W_i^n, W_{i+1}^n) \right| ds \\ &+ kK_2 |W_{i+1}^n - W_i^n| \int_0^1 |\Phi(s; W_i^n, W_{i+1}^n) - \Psi(s; W_i^n, W_{i+1}^n)| ds, \end{aligned} \tag{2.17}$$

where  $k$  is the constant appearing in (2.5) corresponding to a compact set  $\mathcal{K} \subset \Omega$  containing both paths;  $K_1$  is given by:

$$K_1 = \max_{W \in \mathcal{K}} \|\mathcal{A}(W)\|, \tag{2.18}$$

$\|\cdot\|$  being the matrix norm associated to the Euclidean norm in  $\mathbb{R}^N$ ,  $|\cdot|$ ; and finally  $K_2$  is the best constant such that:

$$\|\mathcal{A}(W_1) - \mathcal{A}(W_2)\| \leq K_2 |W_1 - W_2|, \quad \forall W_1, W_2 \in \mathcal{K}. \tag{2.19}$$

Therefore, the difference of states can be viewed as an *approximate* eigenvector associated to  $\xi$  whose error bound depends on the difference between the paths and its derivatives, the distance between the states, and the matrix  $\mathcal{A}$ . Observe that, if we let  $\Delta x$  and  $\Delta t$  tend to zero assuming that the numerical solutions converge to a function having discontinuities, and  $W_i^n, W_{i+1}^n$  converge respectively to the limits  $W^-, W^+$  at both sides of a discontinuity, the error bound (2.17) does not converge to zero, unless  $\Psi = \Phi$ . Therefore, the convergence of the numerical scheme can fail when the solutions to be approached involve discontinuities linking two states  $W^-, W^+$  such that  $\Phi(\cdot; W^-, W^+) \neq \Psi(\cdot; W^-, W^+)$ .

**Remark 1.** Observe that in the deduction of the schemes a CFL-like requirement (2.13) has been imposed. In practice, the following condition can be used:

$$\max \left\{ \left| \lambda_l^{i+1/2} \right|, 1 \leq l \leq N, i \in \mathbb{Z} \right\} \frac{\Delta t}{\Delta x} \leq \gamma \tag{2.20}$$

with  $0 < \gamma \leq 1$ .

**Remark 2.** As in the case of systems of conservation laws, when sonic rarefaction waves appear it is necessary to modify the approximate Riemann problem solver in order to obtain entropy-satisfying solutions. The Harten-Hyman Entropy Fix technique [21, 23], for instance, can be easily adapted to this problem.

### 3. WELL-BALANCING

Well-balancing is related to the numerical approximation of equilibria, *i.e.*, steady state solutions. Observe that system (2.4) can only have regular nontrivial steady state solutions if it has linearly degenerate fields: if  $W(x)$  is a regular steady state solution it satisfies

$$\mathcal{A}(W(x)) \cdot W'(x) = 0, \quad \forall x \in \mathbb{R},$$

and then 0 is an eigenvalue of  $\mathcal{A}(W(x))$  for all  $x$  and  $W'(x)$  is an associated eigenvector. Therefore, the solution can be interpreted as a parameterization of an integral curve of a linearly degenerate characteristic field whose corresponding eigenvalue takes the value 0 through the curve. In order to define the concept of well-balancing, let us introduce the set  $\Gamma$  of all the integral curves  $\gamma$  of a linearly degenerate field of  $\mathcal{A}(W)$  such that the corresponding eigenvalue vanishes on  $\Gamma$ . Clearly, if  $\Gamma$  is empty, the well-balance property doesn't make sense.

**Definition 4.** Given a curve  $\gamma \in \Gamma$ , a numerical scheme for solving (2.4):

$$W_j^{n+1} = W_j^n + \frac{\Delta t}{\Delta x} H(W_{j-p}^n, \dots, W_{j+q}^n) \quad (3.21)$$

is said to be exactly well-balanced for  $\gamma$  if, given any  $\mathcal{C}^1$  function  $x \in (\alpha, \beta) \subset \mathbb{R} \rightarrow W(x) \in \Omega$  such that

$$W(x) \in \gamma, \quad \forall x \in (\alpha, \beta), \quad (3.22)$$

and  $p+q+1$  points in  $(\alpha, \beta)$   $x_{-p}, \dots, x_q$  such that:

$$x_{-p} < \dots < x_q; \quad x_{i+1} - x_i = \Delta x, \quad i = -p, \dots, q-1, \quad (3.23)$$

then

$$H(W(x_{-p}), \dots, W(x_q)) = 0. \quad (3.24)$$

The scheme is said to be well-balanced with order  $k$  for  $\gamma$  if, given any  $\mathcal{C}^{k+1}$  function  $W$  and any set of points  $\{x_{-q}, \dots, x_p\}$  satisfying (3.22), (3.23), then:

$$|H(W(x_{-p}), \dots, W(x_q))| = O(\Delta x^{k+1}). \quad (3.25)$$

Finally, the scheme is said to be exactly well-balanced or well-balanced with order  $k$  if these properties are satisfied for any curve of  $\Gamma$ .

In the definition above, the constant in the expression  $O(\Delta x^k)$  is allowed to depend continuously on  $x_0$ .

We only consider 1-level schemes and uniform meshes in order to avoid an excess of notation, but the definition can be easily extended to general schemes.

Well-balanced schemes solve correctly steady state solutions in the following sense: given a smooth steady state solution  $x \in \mathbb{R} \rightarrow W(x)$  of (2.4), if we apply a numerical scheme to the initial values  $\{W_i^0\}_{i \in \mathbb{Z}}$  given by

$$W_i^0 = W(x_i), \quad \forall i \in \mathbb{Z}, \quad (3.26)$$

using a time step  $\Delta t = T/N$ , for some constant  $T > 0$ , then:

$$W_i^n = W(x_i), \quad \forall n \in \mathbb{N}, \quad i \in \mathbb{Z},$$

if the scheme is exactly well-balanced, and

$$W_i^n = W(x_i) + O(\Delta x^k), \quad \forall i \in \mathbb{Z}, \quad n = 1, \dots, N,$$

if the scheme satisfies the approximated well-balance property with order  $k$  and the estimate (3.25) is satisfied uniformly for  $W$ .

The well-balance property of Roe schemes is strongly related to its ability for approaching solutions involving contact discontinuities: again, it depends on how close the paths of the families  $\Phi$  and  $\Psi$  are.

**Theorem 1.** Consider a Roe scheme (2.14) associated to a family of paths  $\Psi$  for solving (2.4) and let  $\gamma$  be a curve belonging to  $\Gamma$ . If for any two states  $W_L, W_R \in \gamma$ , the path  $\Psi(s; W_L, W_R)$  is a parameterization of the arc of  $\gamma$  delimited by  $W_L$  and  $W_R$ , then the numerical scheme is exactly well-balance for  $\gamma$ .

*Proof.* Let  $W$  be a  $C^1$  function from  $(\alpha, \beta)$  to  $\Omega$  satisfying (3.22) and  $\{x_{-1}, x_0, x_1\}$  three points in  $(\alpha, \beta)$  satisfying (3.23). As  $\Psi(s; W(x_i), W(x_{i+1}))$  is a parameterization of an arc of  $\gamma$  which is in  $\Gamma$  we have for almost every  $s$ :

$$\mathcal{A}(\Psi(s; W(x_i), W(x_{i+1}))) \cdot \frac{\partial \Psi}{\partial s}(s; W(x_i), W(x_{i+1})) = 0.$$

From this equality and (2.10) it can be easily deduced that:

$$\mathcal{A}_{i+1/2} \cdot (W(x_{i+1}) - W(x_i)) = 0, \quad i = -1, 0.$$

The exact well-balance property is proved by using the expression of the scheme (2.14) and the matrix equalities:

$$\mathcal{A}_{i+1/2}^\pm = \mathcal{P}_{i+1/2}^\pm \mathcal{A}_{i+1/2}, \quad i = -1, 0, \tag{3.27}$$

where

$$\mathcal{P}_{i+1/2}^\pm = \frac{1}{2} \mathcal{K}_{i+1/2} (\mathcal{I} \pm \text{sgn}(\mathcal{L}_{i+1/2})) \mathcal{K}_{i+1/2}^{-1}. \quad \square$$

**Corollary 1.** *If  $\Psi = \Phi$ , the numerical scheme is exactly well-balanced.*

*Proof.* This result is easily deduced from the theorem, just taking into account that the path  $\Phi(s; W_L, W_R)$  connecting two states belonging to a curve of a linearly degenerate field is a parameterization of the arc delimited by these states.  $\square$

**Theorem 2.** *Let  $\gamma$  be a curve belonging to  $\Gamma$ . Let us suppose that there exists  $p \in \mathbb{N}$  such that the following estimate holds*

$$\int_0^1 \left| (b-a)W'(a+s(b-a)) - \frac{\partial \Psi}{\partial s}(s; W(a), W(b)) \right| ds = O((b-a)^{p+1}), \tag{3.28}$$

for any  $C^{p+1}$  function  $x \in (\alpha, \beta) \subset \mathbb{R} \rightarrow W(x) \in \Omega$  satisfying (3.22), and any  $a < b$  in  $(\alpha, \beta)$ , then the scheme is well-balanced with order  $p$  for  $\gamma$ .

*Proof.* Let  $W$  be a  $C^{p+1}$  function satisfying (3.22), and  $\{x_{-1}, x_0, x_1\}$  three points in  $(\alpha, \beta)$  satisfying (3.23). We consider a compact set  $\mathcal{K}$  containing the graph of the function  $W : [x_{-1}, x_1] \rightarrow \Omega$ , as well as the paths of the family  $\Psi$  linking  $W(x_{-1}), W(x_0)$ , and  $W(x_0), W(x_1)$ . Let  $k$  be the constant in (2.5) corresponding to the bounded set  $\mathcal{K}$ , and  $K_1, K_2$  the constants defined by (2.18), (2.19), respectively. The following estimate can be obtained

$$\begin{aligned} |\mathcal{A}_{i+1/2} \cdot (W(x_{i+1}) - W(x_i))| &\leq K_1 \int_0^1 \left| \Delta x W'(x_i + s\Delta x) - \frac{\partial \Psi}{\partial s}(s; W(x_i), W(x_{i+1})) \right| ds \\ &+ kK_2 |W(x_{i+1}) - W(x_i)| \int_0^1 |W(x_i + s\Delta x) - \Psi(s; W(x_i), W(x_{i+1}))| ds, \quad i = -1, 0, \end{aligned} \tag{3.29}$$

in the same manner as (2.17). Because of (3.28) the first term on the right-hand side is  $O(\Delta x^{p+1})$ . Moreover, given  $t \in [0, 1]$  we deduce from (3.28):

$$\int_0^1 |W(x_i + s\Delta x) - \Psi(s; W(x_i), W(x_{i+1}))| ds = O(\Delta x^{p+1}),$$

and then, the second term is at least  $O(\Delta x^{p+2})$ . The proof is concluded by using again (3.27) and the expression of the scheme (2.14).  $\square$



In the applications, we will be concerned with Roe schemes based on the family of paths (2.16): this is the natural choice when Volpert's concept of nonconservative products has been chosen or when it is difficult to construct a Roe scheme based on the family of paths chosen in the definition of weak solutions. Besides of its simplicity, the choice (2.16) produces numerical schemes that are well-balanced with order 2:

**Theorem 3.** *A Roe scheme (2.14) based on the family of paths  $\Psi$  given by (2.16) for solving (2.4) is well-balanced with order 2. Moreover, if  $\gamma \in \Gamma$  is a straight line, the numerical scheme is exactly well-balanced for  $\gamma$ .*

*Proof.* Given a curve  $\gamma \in \Gamma$  and a  $C^3$  function  $W$  from  $(\alpha, \beta)$  to  $\Omega$  satisfying (3.22), the well-balance property with order 1 can be easily shown by applying Theorem 2: using Taylor developments, (3.28) can be easily proved with  $p = 1$ . Nevertheless, this result can be improved as follows: given  $j, l \in \{1, \dots, N\}$  we consider the expression

$$\mathcal{A}_{j,l}(W(a) + s(W(b) - W(a)))(w_l(b) - w_l(a)),$$

where  $w_l(x)$  denotes the  $l$ th component of  $W(x)$ , and we perform Taylor developments of  $\mathcal{A}_{j,l}(\cdot)$  and  $w_l(\cdot)$  centered respectively at  $W(c)$  and  $c$ , with  $c = (a + b)/2$ . We obtain:

$$\begin{aligned} \mathcal{A}_{j,l}(W(a) + s(W(b) - W(a)))(w_l(b) - w_l(a)) &= (b - a)\mathcal{A}_{j,l}(W(c))w'_l(c) \\ &\quad + (s - 1/2)(b - a)^2 \sum_{k=1}^N \partial_k \mathcal{A}_{j,l}(W(c))w'_k(c)w'_l(c) + O((b - a)^3). \end{aligned}$$

Adding up these equalities and integrating, we obtain:

$$\int_0^1 \mathcal{A}(\Psi(s; a, b)) \cdot \frac{\partial \Psi}{\partial s}(s; W(a), W(b)) ds = O((b - a)^3). \tag{3.30}$$

The well-balance property with order 2 is concluded from (3.30) reasoning as in the proof of Theorem 2.

Finally, if  $\gamma \in \Gamma$  is a straight line, the exact well-balance property of the numerical scheme for  $\gamma$  is trivially deduced from Theorem 1. □

**Corollary 2.** *If Volpert's definition is assumed for the nonconservative products, that is, if the family of paths  $\Phi$  given by (2.9) is chosen in the definition of weak solutions, then a Roe scheme based on the family of paths (2.16) is exactly well-balanced.*

The approximate well-balance property with order 2 of the numerical schemes based on the choice (2.16) is inherited by those based on paths that are segments for some choice of coordinates in  $\Omega$ . More precisely, we have the following result:

**Theorem 4.** *Let us suppose that  $\mathcal{T} : \Omega \rightarrow \Omega^* \subset \mathbb{R}^N$  is a  $C^3$  one-to-one function with a differentiable inverse function  $\mathcal{S}$ . We will use the following notations:*

$$W^* = \mathcal{T}(W); \quad \mathcal{A}^*(W^*) = \mathcal{A}(\mathcal{S}(W)).$$

*Let us suppose that  $\Psi$  is the family of paths given by:*

$$\Psi(s; W_L, W_R) = \mathcal{S}(W_L^* + s(W_R^* - W_L^*)). \tag{3.31}$$

*A Roe scheme (2.14) based on the family of paths  $\Psi$  given by (3.31) for solving (2.4) is well-balanced with order 2. Moreover, if  $\gamma \in \Gamma$  is such that  $\mathcal{T}(\gamma)$  is a straight line, the numerical scheme is exactly well-balanced for  $\gamma$ .*

*Proof.* Given a curve  $\gamma \in \Gamma$  and a  $\mathcal{C}^3$  function  $W$  from  $(\alpha, \beta)$  to  $\Gamma$  satisfying (3.22), let us define the function  $W^*$  by:

$$W^*(x) = \mathcal{T}(W(x)), \quad x \in (\alpha, \beta).$$

On the one hand,  $W^*$  satisfies the equation:

$$\mathcal{B}(W^*) \cdot W_x^* = 0,$$

where

$$\mathcal{B}(W^*) = \mathcal{A}^*(W^*) \cdot DS(W^*).$$

On the other hand, given  $a, b \in (\alpha, \beta)$ , from (2.10) and (3.31) it can be easily deduced that:

$$\mathcal{A}_\Psi(W(a), W(b)) \cdot (W(b) - W(a)) = \int_0^1 \mathcal{B}(W^*(a) + s(W^*(b) - W^*(a))) \cdot (W^*(b) - W^*(a)) \, ds.$$

The approximate well-balance property with order 2 of the numerical scheme is deduced exactly as in the proof of Theorem 3. The exact well-balance property for curves  $\gamma$  that are straight lines in the coordinates  $W^*$  is trivially deduced from Theorem 1. □

#### 4. COUPLED SYSTEMS OF CONSERVATION LAWS WITH SOURCE TERMS

##### 4.1. Equations

In [4] the numerical resolution of an abstract problem consisting of several systems of conservation laws with source terms coupled to each other by nonconservative products was considered. In this section, we consider a more general problem that fits into the abstract frame (2.4).

The equations considered are as follows:

$$\frac{\partial W_k}{\partial t} + \frac{\partial F_k}{\partial x}(W_k, \sigma) = \sum_{l \neq k} \mathcal{B}_{k,l}(W_1, \dots, W_K, \sigma) \cdot \frac{\partial W_l}{\partial x} + \tilde{S}_k(W_1, \dots, W_K, \sigma) \frac{d\sigma}{dx}, \quad k = 1, \dots, K, \quad (4.32)$$

where

$$W_k(x, t) = \begin{bmatrix} w_1^k(x, t) \\ w_2^k(x, t) \\ \vdots \\ w_{N_k}^k(x, t) \end{bmatrix} \in \mathbb{R}^{N_k},$$

$\sigma(x)$  is a known function from  $\mathbb{R}$  to  $\mathbb{R}$ ;  $F_k$  is a regular function from  $\Omega_k \times \mathbb{R}$  to  $\mathbb{R}^{N_k}$ ,  $\Omega_k$  being an open convex subset of  $\mathbb{R}^{N_k}$ ;  $\mathcal{B}_{k,l}$  is a regular matrix function from  $\Omega = \Omega_1 \times \dots \times \Omega_K \times \mathbb{R} \subset \mathbb{R}^{N+1}$ , being  $N = N_1 + N_2 + \dots + N_K$ , to  $\mathcal{M}_{N_k \times N_l}$ , and  $\tilde{S}_j$  a function from  $\Omega$  to  $\mathbb{R}^{N_j}$ . We assume without loss of generality that this function has the following form:

$$\tilde{S}_k(W_1, \dots, W_K, \sigma) = S_k(W_1, \dots, W_K, \sigma) + \frac{\partial F_k}{\partial \sigma}(W_k, \sigma), \quad (4.33)$$

for a regular function  $S_k$ .

In the absence of the right-hand side, the  $k$ th equation of (4.32) is a system of conservation laws with flux  $F_k$ . We denote by  $\mathcal{J}_k(W_k, \sigma)$  its Jacobian matrix:

$$\mathcal{J}_k(W_k, \sigma) = \frac{\partial F_k}{\partial W_k}(W_k, \sigma).$$

Notice that, if  $K = 1$ , (4.32) is a system of conservation laws with source terms:

$$\frac{\partial W}{\partial t} + \frac{\partial F}{\partial x}(W, \sigma) = S(W, \sigma) \frac{d\sigma}{dx} + \frac{\partial F}{\partial \sigma}(W, \sigma) \frac{d\sigma}{dx}. \tag{4.34}$$

System (4.32) can be also written in a more compact form as follows:

$$\mathbf{W}_t + \mathbf{F}(\mathbf{W}, \sigma)_x = \mathbf{B}(\mathbf{W}, \sigma) \cdot \mathbf{W}_x + \mathbf{S}(\mathbf{W}, \sigma) \frac{d\sigma}{dx} + \frac{\partial \mathbf{F}}{\partial \sigma}(\mathbf{W}, \sigma) \frac{d\sigma}{dx} \tag{4.35}$$

where

$$\mathbf{W} = \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_K \end{bmatrix}, \quad \mathbf{F}(\mathbf{W}, \sigma) = \begin{bmatrix} F_1(W_1, \sigma) \\ \vdots \\ F_K(W_K, \sigma) \end{bmatrix},$$

$\mathbf{B}$  represents the  $N \times N$  matrix with the following block structure:

$$\mathbf{B}(\mathbf{W}, \sigma) = \begin{bmatrix} 0 & \mathbf{B}_{1,2}(\mathbf{W}, \sigma) & \dots & \mathbf{B}_{1,K}(\mathbf{W}, \sigma) \\ \mathbf{B}_{2,1}(\mathbf{W}, \sigma) & 0 & \dots & \mathbf{B}_{2,K}(\mathbf{W}, \sigma) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{B}_{K,1}(\mathbf{W}, \sigma) & \mathbf{B}_{K,2}(\mathbf{W}, \sigma) & \dots & 0 \end{bmatrix}, \tag{4.36}$$

and, finally,  $\mathbf{S}(\mathbf{W}, \sigma)$  is the vector:

$$\mathbf{S}(\mathbf{W}, \sigma) = \begin{bmatrix} S_1(W_1, \dots, W_K, \sigma) \\ \vdots \\ S_K(W_1, \dots, W_K, \sigma) \end{bmatrix}.$$

Using these notations, the function  $\mathbf{F}$  can be viewed as a global flux function whose Jacobian matrix is:

$$\frac{\partial \mathbf{F}}{\partial \mathbf{W}} = \mathcal{J}(\mathbf{W}, \sigma) = \begin{bmatrix} \mathcal{J}_1(W_1, \sigma) & 0 & \dots & 0 \\ 0 & \mathcal{J}_2(W_2, \sigma) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathcal{J}_K(W_K, \sigma) \end{bmatrix}.$$

Following the idea developed in [17, 18] for conservation laws with source terms, if we add to (4.32) the trivial equation:

$$\frac{\partial \sigma}{\partial t} = 0,$$

the problem can be written under the form (2.4):

$$\widetilde{\mathbf{W}}_t + \widetilde{\mathcal{A}}(\widetilde{\mathbf{W}}) \cdot \widetilde{\mathbf{W}}_x = \mathbf{0}, \tag{4.37}$$

where  $\widetilde{\mathbf{W}}$  is the augmented vector:

$$\widetilde{\mathbf{W}} = \begin{bmatrix} \mathbf{W} \\ \sigma \end{bmatrix},$$

and the block structure of the  $(N + 1) \times (N + 1)$  matrix  $\widetilde{\mathcal{A}}$  is given by

$$\widetilde{\mathcal{A}}(\widetilde{\mathbf{W}}) = \left[ \begin{array}{c|c} \mathcal{A}(\widetilde{\mathbf{W}}) & -\widetilde{\mathcal{S}}(\widetilde{\mathbf{W}}, \sigma) \\ \hline 0 & 0 \end{array} \right]. \tag{4.38}$$

Here  $\mathcal{A}(\widetilde{\mathbf{W}})$  represents the  $N \times N$  matrix:

$$\mathcal{A}(\widetilde{\mathbf{W}}) = \mathcal{J}(\mathbf{W}, \sigma) - \mathcal{B}(\mathbf{W}, \sigma).$$

We assume that the matrix  $\mathcal{A}(\widetilde{\mathbf{W}})$  has  $N$  real distinct eigenvalues

$$\lambda_1(\widetilde{\mathbf{W}}) < \dots < \lambda_N(\widetilde{\mathbf{W}}),$$

and associated eigenvectors  $\mathbf{R}_j(\widetilde{\mathbf{W}})$ ,  $j = 1, \dots, N$ . If these eigenvalues do not vanish, (4.37) is a strictly hyperbolic system:  $\widetilde{\mathcal{A}}(\widetilde{\mathbf{W}})$  has  $N + 1$  distinct real eigenvalues:

$$\lambda_1(\widetilde{\mathbf{W}}), \dots, \lambda_N(\widetilde{\mathbf{W}}), 0,$$

with associated eigenvectors:

$$\widetilde{\mathbf{R}}_1(\widetilde{\mathbf{W}}), \dots, \widetilde{\mathbf{R}}_{N+1}(\widetilde{\mathbf{W}})$$

given by

$$\widetilde{\mathbf{R}}_i(\widetilde{\mathbf{W}}) = \begin{bmatrix} \mathbf{R}_i(\widetilde{\mathbf{W}}) \\ 0 \end{bmatrix}, \quad i = 1, \dots, N; \quad \widetilde{\mathbf{R}}_{N+1}(\widetilde{\mathbf{W}}) = \begin{bmatrix} \mathcal{A}(\widetilde{\mathbf{W}})^{-1} \cdot \mathbf{S}(\widetilde{\mathbf{W}}) \\ 1 \end{bmatrix}.$$

Clearly, the  $(N + 1)$ -th field is linearly degenerate and, for the sake of simplicity, we will suppose that it is the only one. The integral curves of the linearly degenerate field are given by those of the o.d.e. system:

$$\frac{d\widetilde{\mathbf{W}}}{ds} = \widetilde{\mathbf{R}}_{N+1}(\widetilde{\mathbf{W}}). \tag{4.39}$$

When one of the eigenvalues of  $\mathcal{A}(\widetilde{\mathbf{W}})$  vanishes, (4.37) becomes nonstrictly hyperbolic. Problems where this situation arises are called *resonant* (see [15]). The definition and the analysis of weak solutions are much more difficult in that case: in particular, Riemann Problems may have more than one entropic solution [1]. In the following two subsections we only consider the strictly hyperbolic case: we discuss the definition of weak solutions and the construction of Roe schemes, respectively. We will only consider the nonstrictly hyperbolic case in the third subsection, when the adaptation of Roe schemes to resonant problems is discussed.

### 4.2. Weak solutions

In order to define the weak solutions of (4.37), first of all a family of paths  $\widetilde{\Phi}(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R)$  has to be chosen. Given two states  $\widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R$  the following notations will be used:

$$\widetilde{\mathbf{W}}_L = \begin{bmatrix} \mathbf{W}_L \\ \sigma_L \end{bmatrix}, \quad \widetilde{\mathbf{W}}_R = \begin{bmatrix} \mathbf{W}_R \\ \sigma_R \end{bmatrix}$$

$$[\widetilde{\mathbf{W}}] = \widetilde{\mathbf{W}}_R - \widetilde{\mathbf{W}}_L = \begin{bmatrix} [\mathbf{W}] \\ [\sigma] \end{bmatrix}.$$

Let us denote by  $\Phi_j(s; \mathbf{W}_L, \mathbf{W}_R)$  the  $j$ th component of a path  $\widetilde{\Phi}(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R)$ . The following notation will be also used:

$$\widetilde{\Phi}(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) = \begin{bmatrix} \Phi(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) \\ \Phi_{N+1}(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) \end{bmatrix} = \begin{bmatrix} \Phi_1(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) \\ \vdots \\ \Phi_N(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) \\ \Phi_{N+1}(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) \end{bmatrix}.$$

Using these notations, some straightforward calculations show that, once the family of paths is chosen, the generalized Rankine-Hugoniot condition (2.6) can be rewritten as follows:

$$\begin{cases} \xi[\mathbf{W}] = \mathbf{F}(\mathbf{W}_R, \sigma_R) - \mathbf{F}(\mathbf{W}_L, \sigma_L) - \mathbf{B}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) - \tilde{\mathbf{S}}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R); \\ \xi[\sigma] = 0; \end{cases} \tag{4.40}$$

where:

$$\mathbf{B}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) = \int_0^1 \mathbf{B}(\tilde{\Phi}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R)) \cdot \frac{\partial \tilde{\Phi}}{\partial s}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) ds, \tag{4.41}$$

$$\tilde{\mathbf{S}}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) = \int_0^1 \left( \mathbf{S} + \frac{\partial \mathbf{F}}{\partial \sigma} \right) (\tilde{\Phi}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R)) \cdot \frac{\partial \Phi_{N+1}}{\partial s}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) ds. \tag{4.42}$$

Notice that, accordingly to the second condition in (4.40), the discontinuities appearing in weak solutions have to be either stationary or they develop in regions where  $\sigma$  is continuous.

When the function  $\sigma$  is constant,  $\sigma(x) = \tilde{\sigma}$ , (4.32) reduces to the homogeneous problem:

$$\frac{\partial W_k}{\partial t} + \frac{\partial F_k}{\partial x}(W_k, \tilde{\sigma}) = \sum_{l \neq k} \mathcal{B}_{k,l}(W_1, \dots, W_K, \tilde{\sigma}) \cdot \frac{\partial W_l}{\partial x} \quad k = 1, \dots, K. \tag{4.43}$$

Therefore, the definitions of weak solution of systems (4.43) and (4.32) have to be consistent in the following sense: if  $(W_1, \dots, W_N)$  solves (4.43), then  $(W_1, \dots, W_N, \tilde{\sigma})$  has to be a solution of (4.32). This is achieved if the family of paths are constructed as follows: first, for any given value of  $\sigma$ , say  $\tilde{\sigma}$ , a family of paths  $\tilde{\Phi}^{\tilde{\sigma}}$  is chosen in order to define the weak solutions of the system (4.43). Then, the family of paths  $\tilde{\Phi}$  is constructed, if possible, in such a way that:

$$\tilde{\Phi}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) = \tilde{\Phi}^{\tilde{\sigma}}(s; \mathbf{W}_L, \mathbf{W}_R); \tag{4.44}$$

$$\Phi_{N+1}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) = \tilde{\sigma}; \tag{4.45}$$

when the states  $\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R$  are such that  $\sigma_L = \sigma_R = \tilde{\sigma}$ .

Another natural requirement in choosing the family of paths is that, if  $\tilde{\mathbf{W}}_L$  and  $\tilde{\mathbf{W}}_R$  lie on an integral curve of (4.39), the contact discontinuity connecting both states has to be a weak solution of (4.32). This is achieved by choosing  $\tilde{\Phi}(s; \mathbf{W}_L, \mathbf{W}_R)$  as a parameterization of the arc of this curve linking the states. More precisely, the path has to satisfy:

$$\frac{d\tilde{\Phi}}{ds} = \alpha(s) \tilde{\mathbf{R}}_{N+1}(\tilde{\Phi}(s)), \quad \tilde{\Phi}(0) = \tilde{\mathbf{W}}_L, \quad \tilde{\Phi}(1) = \tilde{\mathbf{W}}_R; \tag{4.46}$$

for some non-vanishing scalar function  $\alpha(s)$ .

In the particular case  $K = 1$  the above requirements completely determine the family of paths: notice firstly that, in this case, the homogeneous system (4.43) is a pure conservation law and thus the definition of weak solutions is independent of the family of paths: we can choose, for instance the family of segments:

$$\tilde{\Phi}^{\tilde{\sigma}}(s; \mathbf{W}_L, \mathbf{W}_R) = \mathbf{W}_L + s(\mathbf{W}_R - \mathbf{W}_L).$$

This choice, together with (4.44) and (4.45), determines the path connecting two states such that  $\sigma_R = \sigma_L = \tilde{\sigma}$ , which is again the segment:

$$\tilde{\Phi}(s; \tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R) = \tilde{\mathbf{W}}_L + s(\tilde{\mathbf{W}}_R - \tilde{\mathbf{W}}_L). \tag{4.47}$$

On the other hand, the path connecting two states that lie on the same integral curve of the  $(N+1)$ -characteristic field is given by (4.46). In general, the path connecting two arbitrary states has to be composed of segments (4.47) lying on hyperplanes  $\sigma = cte$  and arcs of integral curves of (4.39). The difficulty comes from the calculation of

the intermediate points linking these segments and arcs. In fact, in order to calculate these points, the exact Riemann problem associated to equation (4.37) with initial conditions given by  $\widetilde{\mathbf{W}}_L$  and  $\widetilde{\mathbf{W}}_R$ , has to be solved, and this is in general a difficult task.

In the case  $K > 1$ , some extra information is needed in order to define the family of paths, that may come from physical consideration (see [10]) or by taking the limit of viscous profiles (see [28]). Nevertheless, in the particular case where the matrices  $\mathcal{B}_{k,l}$  satisfy:

$$\mathcal{B}_{k,l} = \mathcal{B}_{k,l}(W_k, \sigma) \text{ and } \mathcal{B}_{k,l}(0, \sigma) = 0, \tag{4.48}$$

another natural requirement can be imposed to weak solutions of (4.32). In this case, it can be easily verified that, if  $W_k$  is a classical solution of the conservation law:

$$\frac{\partial W_k}{\partial t} + \frac{\partial F_k}{\partial x}(W_k, \tilde{\sigma}) = 0, \tag{4.49}$$

for a given  $\tilde{\sigma}$ , then  $(0, \dots, W_k, \dots, 0, \tilde{\sigma})$  is a classical solution of (4.32). Therefore, it is natural to define the weak solutions so that this relationship remains valid. This is achieved by choosing a family of paths satisfying also:

$$\Phi_j(s; \widetilde{\mathbf{W}}_L, \widetilde{\mathbf{W}}_R) = 0, \quad \forall j \neq i, \tag{4.50}$$

when the states to link are of the form:

$$\widetilde{\mathbf{W}}_L = \begin{bmatrix} 0 \\ \vdots \\ W_{i,L} \\ \vdots \\ 0 \\ \tilde{\sigma} \end{bmatrix}, \quad \widetilde{\mathbf{W}}_R = \begin{bmatrix} 0 \\ \vdots \\ W_{i,R} \\ \vdots \\ 0 \\ \tilde{\sigma} \end{bmatrix}. \tag{4.51}$$

For instance, the segment linking the states can be chosen in this case.

Once  $\Phi^{\tilde{\sigma}}$  chosen, the family of paths  $\tilde{\Phi}$  is completely determined. On the one hand, the paths linking two states satisfying  $\sigma_L = \sigma_R$  is given by (4.44), (4.45), and if they lie on the same integral curve, by (4.46). As in the case  $K = 1$ , the path linking two arbitrary states has to be composed by pieces of these types the curves, and the calculation of the intermediate points requires in general the resolution of a Riemann problem.

### 4.3. Roe schemes

In the description of the numerical schemes we will use the following notations: the approximation of the solution in the cell  $I_i$  at time  $n\Delta t$  will be represented by

$$\widetilde{\mathbf{W}}_i^n = \begin{bmatrix} \mathbf{W}_i^n \\ \sigma_i \end{bmatrix} = \begin{bmatrix} W_{i,1}^n \\ \vdots \\ W_{i,K}^n \\ \sigma_i \end{bmatrix}.$$

On the other hand, we will denote by  $\tilde{\Psi}$  the family of paths chosen in order to construct the Roe linearization. We will use the notation:

$$\tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) = \begin{bmatrix} \Psi(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \\ \Psi_N(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \\ \Psi_{N+1}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \end{bmatrix} = \begin{bmatrix} \Psi_1(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \\ \vdots \\ \Psi_N(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \\ \Psi_{N+1}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \end{bmatrix}.$$

It is natural to look for Roe matrices with the same structure (4.38) of  $\tilde{\mathcal{A}}(\tilde{\mathbf{W}})$ :

$$\tilde{\mathcal{A}}_{i+1/2} = \begin{bmatrix} \mathcal{A}_{i+1/2} & -\mathcal{S}_{i+1/2} \\ 0 & 0 \end{bmatrix}. \tag{4.52}$$

For Roe matrices with this structure (2.10) can be rewritten as follows:

$$\mathcal{A}_{i+1/2} \cdot (\mathbf{W}_{i+1}^n - \mathbf{W}_i^n) - \mathcal{S}_{i+1/2}(\sigma_{i+1} - \sigma_i) = \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1}) - \mathbf{F}(\mathbf{W}_i^n, \sigma_i) - \mathbf{B}_{\tilde{\Psi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) - \tilde{\mathcal{S}}_{\tilde{\Psi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n); \tag{4.53}$$

where  $\mathbf{B}_{\tilde{\Psi}}$  and  $\mathcal{S}_{\tilde{\Psi}}$  are defined by (4.41), (4.42) substituting  $\tilde{\Phi}$  by  $\tilde{\Psi}$  and  $(\tilde{\mathbf{W}}_L, \tilde{\mathbf{W}}_R)$  by  $(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n)$ .

Let us suppose that, for any fixed value of  $\sigma$ , Roe matrices can be calculated for each pure system of conservation laws:

$$\frac{\partial W_k}{\partial t} + \frac{\partial F_k}{\partial x}(W_k, \sigma) = 0. \tag{4.54}$$

That is, we suppose that, given  $\mathbf{W}_i^n, \mathbf{W}_{i+1}^n, \sigma$ , we can calculate for each  $k$  a  $N_k \times N_k$  matrix  $\mathcal{J}_{i+1/2,k}^\sigma$  such that:

$$\mathcal{J}_{i+1/2,k}^\sigma \cdot (W_{i+1,k}^n - W_{i,k}^n) = F_k(W_{i+1,k}^n, \sigma) - F_k(W_{i,k}^n, \sigma), \quad 1 \leq k \leq K.$$

Let us suppose also that it is possible to calculate a value of  $\sigma, \sigma_{i+1/2}$ , a  $N \times N$  matrix  $\mathcal{B}_{i+1/2}$ , and a vector  $\mathcal{S}_{i+1/2}$  such that the following identities hold:

$$\begin{aligned} & \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1}) - \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1/2}) + \mathbf{F}(\mathbf{W}_i^n, \sigma_{i+1/2}) - \mathbf{F}(\mathbf{W}_i^n, \sigma_i) \\ &= \int_0^1 \frac{\partial \mathbf{F}}{\partial \sigma} \left( \tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \right) \cdot \frac{\partial \Psi_{N+1}}{\partial s} \left( s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n \right) ds; \end{aligned} \tag{4.55}$$

$$\mathcal{B}_{i+1/2} \cdot (\mathbf{W}_{i+1}^n - \mathbf{W}_i^n) = \int_0^1 \mathcal{B} \left( \tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \right) \cdot \frac{\partial \Psi}{\partial s} \left( s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n \right) ds; \tag{4.56}$$

$$\mathcal{S}_{i+1/2}(\sigma_{i+1} - \sigma_i) = \int_0^1 \mathcal{S} \left( \tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \right) \cdot \frac{\partial \Psi_{N+1}}{\partial s} \left( s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n \right) ds. \tag{4.57}$$

Then, if we define the matrices:

$$\mathcal{J}_{i+1/2} = \begin{bmatrix} \mathcal{J}_{i+1/2,1}^{\sigma_{i+1/2}} & 0 & \dots & 0 \\ 0 & \mathcal{J}_{i+1/2,2}^{\sigma_{i+1/2}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathcal{J}_{i+1/2,K}^{\sigma_{i+1/2}} \end{bmatrix}, \tag{4.58}$$

$$\mathcal{A}_{i+1/2} = \mathcal{J}_{i+1/2} - \mathcal{B}_{i+1/2}, \tag{4.59}$$

an easy calculation shows that the corresponding matrix (4.52) satisfies (4.53). Therefore, if this matrix has  $N + 1$  real distinct eigenvalues, it provides a Roe linearization for (4.37).

Observe that the condition concerning the eigenvalues is satisfied if  $\mathcal{A}_{i+1/2}$  has  $N$  real distinct eigenvalues

$$\lambda_{i+1/2,1} < \dots < \lambda_{i+1/2,N},$$

and associated eigenvectors  $\mathbf{R}_{i+1/2,j}$ ,  $j = 1, \dots, N$ : in this case,  $\tilde{\mathcal{A}}_{i+1/2}$  has  $N + 1$  distinct real eigenvalues:

$$\lambda_{i+1/2,1}, \dots, \lambda_{i+1/2,N}, 0,$$

with associated eigenvectors:

$$\tilde{\mathbf{R}}_{i+1/2,1}, \dots, \tilde{\mathbf{R}}_{i+1/2,N+1}$$

given by

$$\tilde{\mathbf{R}}_{i+1/2,j} = \begin{bmatrix} \mathbf{R}_{i+1/2,j} \\ 0 \end{bmatrix}, \quad j = 1, \dots, N; \quad \tilde{\mathbf{R}}_{i+1/2,N+1} = \begin{bmatrix} \mathcal{A}_{i+1/2}^{-1} \cdot \mathbf{S}_{i+1/2} \\ 1 \end{bmatrix}. \tag{4.60}$$

In this case, some algebraic manipulations allow to show the following identities:

$$\tilde{\mathcal{A}}_{i+1/2}^{\pm} = \left[ \begin{array}{c|c} \mathcal{A}_{i+1/2}^{\pm} & -\mathcal{A}_{i+1/2}^{\pm} \mathcal{A}_{i+1/2}^{-1} \mathbf{S}_{i+1/2} \\ \hline 0 & 0 \end{array} \right], \tag{4.61}$$

$$|\tilde{\mathcal{A}}_{i+1/2}| = \left[ \begin{array}{c|c} |\mathcal{A}_{i+1/2}| & -|\mathcal{A}_{i+1/2}| \mathcal{A}_{i+1/2}^{-1} \mathbf{S}_{i+1/2} \\ \hline 0 & 0 \end{array} \right]. \tag{4.62}$$

Using these equalities, the numerical scheme (2.14) can be rewritten under the equivalent form:

$$\begin{aligned} \mathbf{W}_i^{n+1} &= \mathbf{W}_i^n + \frac{\Delta t}{\Delta x} (\mathbf{F}_{i-1/2} - \mathbf{F}_{i+1/2}) \\ &+ \frac{\Delta t}{2\Delta x} (\mathcal{B}_{i-1/2} \cdot (\mathbf{W}_i^n - \mathbf{W}_{i-1}^n) + \mathcal{B}_{i+1/2} \cdot (\mathbf{W}_{i+1}^n - \mathbf{W}_i^n)) \\ &+ \frac{\Delta t}{\Delta x} (\mathcal{P}_{i-1/2}^+ \mathbf{S}_{i-1/2} (\sigma_i - \sigma_{i-1}) + \mathcal{P}_{i+1/2}^- \mathbf{S}_{i+1/2} (\sigma_{i+1} - \sigma_i)) \\ &+ \frac{\Delta t}{2\Delta x} (\mathbf{V}_{i-1/2} + \mathbf{V}_{i+1/2}) \end{aligned} \tag{4.63}$$

where

$$\mathbf{F}_{i+1/2} = \frac{1}{2} (\mathbf{F}(\mathbf{W}_i^n, \sigma_i) + \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1})) - \frac{1}{2} |\mathcal{A}_{i+1/2}| \cdot (\mathbf{W}_{i+1}^n - \mathbf{W}_i^n), \tag{4.64}$$

$$\mathcal{P}_{i+1/2}^{\pm} = \frac{1}{2} (\mathcal{I} \pm |\mathcal{A}_{i+1/2}| \mathcal{A}_{i+1/2}^{-1}). \tag{4.65}$$



These latter matrices can be also be written under the form:

$$\mathcal{P}_{i+1/2}^\pm = \frac{1}{2} \mathcal{K}_{i+1/2} (\mathcal{I} \pm \text{sgn}(\mathcal{L})_{i+1/2}) \mathcal{K}_{i+1/2}^{-1}, \tag{4.66}$$

where  $\mathcal{K}_{i+1/2}$  is the  $N \times N$  matrix whose columns are the eigenvectors  $\mathbf{R}_{i+1/2,1}, \dots, \mathbf{R}_{i+1/2,N}$  and  $\text{sgn}(\mathcal{L})_{i+1/2}$  is the diagonal matrix whose coefficients are the signs of the eigenvalues  $\lambda_{i+1/2,1}, \dots, \lambda_{i+1/2,N}$ . Finally

$$\mathbf{V}_{i+1/2} = \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1}) - \mathbf{F}(\mathbf{W}_{i+1}^n, \sigma_{i+1/2}) + \mathbf{F}(\mathbf{W}_i^n, \sigma_{i+1/2}) - \mathbf{F}(\mathbf{W}_i^n, \sigma_i),$$

or, equivalently (see (4.55)):

$$\mathbf{V}_{i+1/2} = \int_0^1 \frac{\partial \mathbf{F}}{\partial \sigma} \left( \tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) \right) \cdot \frac{\partial \Psi_{N+1}}{\partial s} \left( s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n \right) ds.$$

Under the formulation (4.63) the numerical scheme has the form of a  $Q$ -scheme of Roe upwinding the source term in the sense defined in [3] for balance laws or in [4] for coupled systems of conservation laws with source term. Nevertheless, there is a slight difference: in the cited articles, it was assumed that the Roe matrices could be calculated by evaluating the Jacobians at an intermediate state, and  $\mathbf{B}_{i+1/2}$ ,  $\mathbf{S}_{i+1/2}$  were obtained by evaluating  $\mathbf{B}$  and  $\mathbf{S}$  at these intermediate states. Nevertheless, in the particular case of the one-layer or the two-layer shallow water system studied below, these choices of  $\mathbf{B}_{i+1/2}$ ,  $\mathbf{S}_{i+1/2}$  satisfy (4.56), (4.57) and, therefore, the numerical schemes presented in the cited works are particular cases of the Roe schemes studied here.

Concerning the family of paths  $\tilde{\Psi}$ , the optimal choice would be  $\tilde{\Psi} = \tilde{\Phi}$  as it gives exactly well-balanced schemes that can handle correctly with shocks and contact discontinuities, but the construction of Roe linearizations based on  $\tilde{\Psi}$  can be in practice a very difficult and expensive task. Therefore, in the applications below, we will consider Roe linearizations based on families of paths that are segments for a convenient choice of coordinates. Theorems 3 and 4 assure that the numerical schemes obtained are well-balanced with order 2.

Due to this choice of paths, the resulting numerical schemes cannot be expected to give good numerical solutions when they are applied to problems where  $\sigma$  present discontinuities. Nevertheless, they are expected to handle correctly with shocks, in the sense discussed above: if  $\tilde{\mathbf{W}}_i^n$  and  $\tilde{\mathbf{W}}_{i+1}^n$  can be linked by a shock whose speed is  $\xi$ , the matrix  $\tilde{\mathcal{A}}_{i+1/2}$  is expected to have  $\xi$  as an eigenvalue and the difference of the states as an associated eigenvector. Clearly, this is achieved if the paths  $\tilde{\Phi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n)$  and  $\tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n)$  coincide, but this is not a necessary condition: if (4.40) and (4.53) are compared, it can be easily shown that this property is satisfied if the following equality holds:

$$\mathbf{B}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) + \tilde{\mathbf{S}}_{\tilde{\Phi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) = \mathbf{B}_{\tilde{\Psi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) + \tilde{\mathbf{S}}_{\tilde{\Psi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n). \tag{4.67}$$

In the particular case  $K = 1$ , when two states  $\tilde{\mathbf{W}}_i^n$  and  $\tilde{\mathbf{W}}_{i+1}^n$  can be linked by a shock, then  $\sigma_i = \sigma_{i+1}$  and (4.67) reduces to:

$$\tilde{\mathbf{S}}_{\tilde{\Psi}}(\tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) = 0, \tag{4.68}$$

which is therefore a sufficient condition to assure that the shock is correctly treated. This equality is trivially satisfied if the path  $\tilde{\Psi}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n)$  is such that:

$$\Psi_{N+1}(s; \tilde{\mathbf{W}}_i^n, \tilde{\mathbf{W}}_{i+1}^n) = \sigma_i = \sigma_{i+1}, \quad \forall s. \tag{4.69}$$

Again, this is not a necessary condition: in Section 5.2 a numerical scheme satisfying (4.68) but not (4.69) will be described.

**Remark 3.** Let us suppose that there is a discontinuity of  $\sigma$  placed at the intercell  $x_{i+1/2}$  and that  $\widetilde{\mathbf{W}}_i^n, \widetilde{\mathbf{W}}_{i+1}^n$  can be linked by a contact discontinuity. If there exists  $\bar{s} \in [0, 1]$  such that:

$$\widetilde{\mathbf{W}}_{i+1}^n - \widetilde{\mathbf{W}}_i^n = \widetilde{\Phi}(1; \widetilde{\mathbf{W}}_i^n, \widetilde{\mathbf{W}}_{i+1}^n) - \widetilde{\Phi}(0; \widetilde{\mathbf{W}}_i^n, \widetilde{\mathbf{W}}_{i+1}^n) = \frac{\partial \widetilde{\Phi}}{\partial s}(\bar{s}; \widetilde{\mathbf{W}}_i^n, \widetilde{\mathbf{W}}_{i+1}^n),$$

then

$$\widetilde{\mathcal{A}}_{i+1/2} = \widetilde{\mathcal{A}}(\widetilde{\Phi}(\bar{s}; \widetilde{\mathbf{W}}_i^n, \widetilde{\mathbf{W}}_{i+1}^n))$$

would be a good choice for the intermediate matrix, as  $\widetilde{\mathbf{W}}_{i+1}^n - \widetilde{\mathbf{W}}_i^n$  is an eigenvector associated to the eigenvalue 0. As a consequence, the solution of the Linear Riemann Problem is a stationary shock placed at the intercell. Unfortunately, even if such a value of  $\bar{s}$  can be calculated, it is not easy to construct a suitable numerical scheme combining this choice of intermediate matrix with a Roe linearization, as we shall see in the applications.

#### 4.4. The nonstrictly hyperbolic case

In [12, 36] the numerical treatment of resonant problems was performed by introducing some corrections in numerical schemes of the form (4.63): the projection matrices (4.66) were defined so that for every vanishing eigenvalue of the Jacobian matrix, no unwinding of the source term was performed in the corresponding characteristic direction. The goal of this subsection is to give a mathematical deduction of these corrected schemes by considering some resonant Linear Riemann Problems at intercells where this situation arises.

Let us suppose that we have chosen Roe matrices  $\widetilde{\mathcal{A}}_{i+1/2}$  of the form (4.52) and that  $\lambda_{i+1/2,j} = 0$  for some  $i \in \mathbb{Z}, j \in \{1, \dots, N\}$ . In this case, 0 is a double eigenvalue for  $\widetilde{\mathcal{A}}_{i+1/2}$  and the Linear Riemann Problem to solve at the intercell  $x_{i+1/2}$  is nonstrictly hyperbolic.

Two different situations may arise: let us suppose firstly that there exists a vector  $\mathbf{T}_{i+1/2} \in \mathbb{R}^N$  such that

$$\mathcal{A}_{i+1/2} \cdot \mathbf{T}_{i+1/2} = \mathbf{S}_{i+1/2}.$$

In this case  $\widetilde{\mathcal{A}}_{i+1/2}$  still has a complete set of eigenvectors

$$\widetilde{\mathbf{R}}_{i+1/2,1}, \dots, \widetilde{\mathbf{R}}_{i+1/2,N+1}$$

whose first  $N$  components are given by (4.60) and

$$\widetilde{\mathbf{R}}_{i+1/2,N+1} = \begin{bmatrix} \mathbf{T}_{i+1/2} \\ 1 \end{bmatrix}. \tag{4.70}$$

The solution of the Linear Riemann Problem is self-similar,  $\widetilde{U}(\frac{x-x_{i+1/2}}{t})$ , and, setting

$$\widetilde{\mathbf{W}}_i^n = \sum_{k=1}^{N+1} \alpha_k^i \widetilde{\mathbf{R}}_k; \quad \widetilde{\mathbf{W}}_{i+1}^n = \sum_{k=1}^{N+1} \alpha_k^{i+1} \widetilde{\mathbf{R}}_k,$$

it satisfies:

$$\widetilde{U}(0)^- = \sum_{k=1}^{j-1} \alpha_k^{i+1} \widetilde{\mathbf{R}}_k + \sum_{k=j}^{N+1} \alpha_k^i \widetilde{\mathbf{R}}_k, \tag{4.71}$$

$$\widetilde{U}(0)^+ = \sum_{k=1}^j \alpha_k^{i+1} \widetilde{\mathbf{R}}_k + \sum_{k=j+1}^N \alpha_k^i \widetilde{\mathbf{R}}_k + \alpha_{N+1}^{i+1} \widetilde{\mathbf{R}}_{N+1}. \tag{4.72}$$

Using these expressions, the following equalities can be easily verified:

$$\tilde{\mathcal{A}}_{i+1/2} \cdot (\tilde{U}(0)^+ - \tilde{U}(0)^-) = 0; \tag{4.73}$$

$$\tilde{\mathcal{A}}_{i+1/2} \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{U}(0)^+) = \tilde{\mathcal{A}}_{i+1/2}^+ \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{\mathbf{W}}_i^n), \tag{4.74}$$

$$\tilde{\mathcal{A}}_{i+1/2} \cdot (\tilde{U}(0)^- - \tilde{\mathbf{W}}_i^n) = \tilde{\mathcal{A}}_{i+1/2}^- \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{\mathbf{W}}_i^n). \tag{4.75}$$

As a consequence, the expression of the numerical scheme coincides with the corresponding to the strictly hyperbolic case. It can be easily shown that the matrices  $\tilde{\mathcal{A}}_{i+1/2}^\pm$  and  $|\tilde{\mathcal{A}}_{i+1/2}|$  are now as follows:

$$\tilde{\mathcal{A}}_{i+1/2}^\pm = \left[ \begin{array}{c|c} \mathcal{A}_{i+1/2}^\pm & -\mathcal{A}_{i+1/2}^\pm \mathbf{T}_{i+1/2} \\ \hline 0 & 0 \end{array} \right], \tag{4.76}$$

$$|\tilde{\mathcal{A}}_{i+1/2}| = \left[ \begin{array}{c|c} |\mathcal{A}_{i+1/2}| & -|\mathcal{A}_{i+1/2}| \mathbf{T}_{i+1/2} \\ \hline 0 & 0 \end{array} \right]. \tag{4.77}$$

Moreover, if  $\mathbf{T}_{i+1/2}$  is chosen as follows:

$$\mathbf{T}_{i+1/2} = \mathcal{K}_{i+1/2} \cdot \mathcal{D}_{i+1/2} \cdot \mathcal{K}_{i+1/2}^{-1} \cdot \mathbf{S}_{i+1/2} \tag{4.78}$$

where  $\mathcal{D}_{i+1/2}$  is the  $N \times N$  diagonal matrix with coefficients

$$\lambda_{i+1/2,1}^{-1}, \dots, \lambda_{i+1/2,j-1}^{-1}, 0, \lambda_{i+1/2,j+1}^{-1}, \dots, \lambda_{i+1/2,N}^{-1},$$

we obtain:

$$-|\mathcal{A}_{i+1/2}| \mathbf{T}_{i+1/2} = \mathcal{K}_{i+1/2} \cdot \text{sgn}(\mathcal{L})_{i+1/2} \mathcal{K}_{i+1/2}^{-1} \cdot \mathbf{S}_{i+1/2}, \tag{4.79}$$

where  $\text{sgn}(\mathcal{L})_{i+1/2}$  is defined as in the previous subsection, assuming that the sign of 0 is 0. Therefore, the scheme can also be written under the form (4.63) using (4.66) for upwinding the source terms.

Let us consider now the case in which  $\mathbf{S}_{i+1/2}$  does not belong to the image of  $\mathcal{A}_{i+1/2}$ . A vector  $\mathbf{T}_{i+1/2}$  can be chosen such that:

$$\mathcal{A}_{i+1/2} \mathbf{T}_{i+1/2} - \mathbf{S}_{i+1/2} = \alpha \mathbf{R}_{i+1/2,j}, \tag{4.80}$$

for some  $\alpha \in \mathbb{R}$ . If we consider the basis

$$\tilde{\mathbf{R}}_{i+1/2,1}, \dots, \tilde{\mathbf{R}}_{i+1/2,N+1},$$

whose first  $N$  components are given again by (4.60) and  $\tilde{\mathbf{R}}_{i+1/2,N+1}$  by (4.70), the expression of  $\tilde{\mathcal{A}}_{i+1/2}$  in this basis has the Jordan structure:

$$\left[ \begin{array}{c|c} \mathcal{L}_{i+1/2} & \alpha \mathbf{e}_j \\ \hline 0 & 0 \end{array} \right],$$

where  $\mathcal{L}_{i+1/2}$  is the diagonal matrix whose coefficients are  $\lambda_{i+1/2,1}, \dots, \lambda_{i+1/2,N}$  and  $\mathbf{e}_j$  represents the  $j$ th element of the canonical basis of  $\mathbb{R}^N$ .

The Linear Riemann Problem is not hyperbolic, but it can be easily solved in the characteristic coordinates. The solution is again self-similar and it satisfies (4.71), (4.72). Nevertheless, the equality (4.73) does not hold in this case. In order to define the numerical scheme, the solution of the Linear Riemann problem at  $x = 0$  can be approached by:

$$\tilde{U}(0) \cong \tilde{U}_0 = \frac{1}{2} (\tilde{U}(0)^+ + \tilde{U}(0)^-).$$

With this approximation, the following equalities similar to (4.74), (4.75) hold:

$$\begin{aligned} \tilde{\mathcal{A}}_{i+1/2} \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{\mathbf{U}}_0) &= \tilde{\mathcal{A}}_{i+1/2}^+ \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{\mathbf{W}}_i^n), \\ \tilde{\mathcal{A}}_{i+1/2} \cdot (\tilde{\mathbf{U}}_0 - \tilde{\mathbf{W}}_i^n) &= \tilde{\mathcal{A}}_{i+1/2}^- \cdot (\tilde{\mathbf{W}}_{i+1}^n - \tilde{\mathbf{W}}_i^n), \end{aligned}$$

if the matrices  $\tilde{\mathcal{A}}_{i+1/2}^\pm$  are defined as follows:

$$\tilde{\mathcal{A}}_{i+1/2}^\pm = \tilde{\mathcal{K}}_{i+1/2} \cdot \left[ \frac{\mathcal{L}_{i+1/2}^\pm}{0} \middle| \frac{\frac{\alpha}{2} \mathbf{e}_j}{0} \right] \cdot \tilde{\mathcal{K}}_{i+1/2}^{-1}, \tag{4.81}$$

being  $\tilde{\mathcal{K}}_{i+1/2}$  the  $(N + 1) \times (N + 1)$  matrix whose columns are the eigenvectors  $\tilde{\mathbf{R}}_k$ ,  $k = 1, \dots, N + 1$ . The block structure of these matrices are:

$$\tilde{\mathcal{A}}_{i+1/2}^\pm = \left[ \frac{\mathcal{A}_{i+1/2}^\pm}{0} \middle| \frac{-\mathcal{A}_{i+1/2}^\pm \mathbf{T}_{i+1/2} + \frac{\alpha}{2} \mathbf{R}_j}{0} \right]. \tag{4.82}$$

With these definitions, the equality

$$\tilde{\mathcal{A}}_{i+1/2}^+ + \tilde{\mathcal{A}}_{i+1/2}^- = \tilde{\mathcal{A}}_{i+1/2}$$

still holds, and we can define:

$$\left| \tilde{\mathcal{A}}_{i+1/2} \right| = \tilde{\mathcal{A}}_{i+1/2}^+ - \tilde{\mathcal{A}}_{i+1/2}^-,$$

whose expression is given again by (4.77).

Choosing again the vector  $\mathbf{T}_{i+1/2}$  given by (4.78), (4.80) is satisfied for a convenient value of  $\alpha$ . Therefore, this vector satisfies (4.79) and the scheme can also be written under the form (4.63) using (4.66) for upwinding the source terms.

## 5. APPLICATIONS

### 5.1. Shallow water equations with depth variations

We consider the one-dimensional shallow water system:

$$\begin{cases} \frac{\partial h}{\partial t} + \frac{\partial q}{\partial x} = 0, \\ \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} \left( \frac{q^2}{h} + \frac{g}{2} h^2 \right) = gh \frac{dH}{dx}, \end{cases} \tag{5.83}$$

which are the equations governing the flow of a shallow layer of fluid through a straight channel with a constant rectangular cross-section. The coordinate  $x$  refers to the axis of the channel and  $t$  is time;  $q(x, t)$  and  $h(x, t)$  represent the mass-flow and the thickness;  $g$  is gravity and  $H(x)$  the depth function measured from a fixed level of reference. The fluid is supposed to be homogeneous and inviscid.

System (5.83) is a problem of the form (4.34) with  $N = N_1 = 2$ ,  $W = [h, q]^T$ ,  $\sigma = H$ , and

$$\begin{aligned} F(W) &= \begin{bmatrix} q \\ \frac{q^2}{h} + \frac{g}{2} h^2 \end{bmatrix}, \\ S(W) &= \begin{bmatrix} 0 \\ gh \end{bmatrix}. \end{aligned} \tag{5.84}$$

The system can be written under the form (2.4) with

$$\widetilde{W} = \begin{bmatrix} h \\ q \\ H \end{bmatrix},$$

and

$$\widetilde{\mathcal{A}}(\widetilde{W}) = \begin{bmatrix} 0 & 1 & 0 \\ -u^2 + c^2 & 2u & -c^2 \\ 0 & 0 & 0 \end{bmatrix}, \quad (5.85)$$

where  $u = q/h$  represents the averaged velocity and  $c = \sqrt{gh}$ .

The eigenvalues of this matrix are:

$$\lambda_1 = u - c, \quad \lambda_2 = u + c, \quad 0.$$

If  $\lambda_i \neq 0$ ,  $i = 1, 2$ , the system is strictly hyperbolic and a complete set of eigenvectors is given by:

$$\widetilde{R}_i(W) = \begin{bmatrix} 1 \\ \lambda_i \\ 0 \end{bmatrix}, \quad i = 1, 2; \quad \widetilde{R}_3(W) = \begin{bmatrix} 1 \\ 0 \\ 1 - Fr^2 \end{bmatrix};$$

where

$$Fr = \frac{u}{c} \quad (5.86)$$

is the *Froude number*. This problem can also lose its strictly hyperbolic character when  $\lambda_1 = \lambda_2$ . This situation arises only when the thickness of the layer  $h$  vanishes, *i.e.* in wet/dry fronts. These situations will not be considered here.

The integral curves of the 3rd-field are given by:

$$q = ct., \quad h + \frac{q^2}{2gh^2} - H = ct. \quad (5.87)$$

Observe that in the particular case  $q = 0$ , these curves are straight lines in the  $q, h, H$ -space:

$$q = 0, \quad h - H = ct.$$

Accordingly to Section 4.2, two kind of discontinuities may appear in the weak solutions of the system: on the one hand, 1 or 2-shocks satisfying the usual Rankine-Hugoniot conditions and linking two states  $\widetilde{W}_L$  and  $\widetilde{W}_R$  such that  $H_L = H_R$ . On the other hand, contact discontinuities standing at points where the bottom presents a jump. Two states  $\widetilde{W}_L$  and  $\widetilde{W}_R$  can be connected by such a discontinuity if they lie on a curve of the family (5.87). If we define:

$$\eta_L = h_L - H_L, \quad \eta_R = h_R - H_R,$$

that can be interpreted as the elevation of the free surface over the level of reference at both sides of the discontinuity,  $\widetilde{W}_L$  and  $\widetilde{W}_R$  can be connected by a contact discontinuity if they verify:

$$q_L = q_R, \quad g\eta_L + \frac{u_L^2}{2} = g\eta_R + \frac{u_R^2}{2}.$$

The second equality can be interpreted in terms of energy: the total mechanical energy is preserved through a contact discontinuity, while there is always a loss of energy across entropic 1 or 2-shock waves (see [31]).

Given a bottom function  $H(x)$ , steady states solutions of (5.83) can be calculated by parameterizing the curves (5.87) with  $x$ . In particular, if  $q = 0$  we obtain the solutions:

$$h(x) = H(x) + ct., \quad q = 0, \quad (5.88)$$

representing water at rest. Another steady state solution is given by:

$$h = 0, \quad q = 0. \tag{5.89}$$

which corresponds to *vacuum*.

It is easy to construct a Roe scheme based on the family of paths (2.16) for (5.83), following the indications of Section 4.3. First we consider the Roe matrices for the homogeneous shallow water system:

$$\mathcal{J}_{i+1/2} = \begin{bmatrix} 0 & 1 \\ -(u_{i+1/2}^n)^2 + (c_{i+1/2}^n)^2 & 2u_{i+1/2}^n \end{bmatrix} \tag{5.90}$$

where

$$u_{i+1/2}^n = \frac{\sqrt{h_i^n}u_i^n + \sqrt{h_{i+1}^n}u_{i+1}^n}{\sqrt{h_i^n} + \sqrt{h_{i+1}^n}}, \quad c_{i+1/2}^n = \sqrt{g \frac{h_i^n + h_{i+1}^n}{2}}. \tag{5.91}$$

Then, we calculate:

$$S_{i+1/2} = \int_0^1 S(\widetilde{W}_i^n + s(\widetilde{W}_{i+1}^n - \widetilde{W}_i^n)) ds = \begin{bmatrix} 0 \\ g \frac{h_i^n + h_{i+1}^n}{2} \end{bmatrix}.$$

The Roe matrices are thus as follows:

$$\mathcal{A}_{i+1/2} = \begin{bmatrix} 1 & 0 & 0 \\ -(u_{i+1/2}^n)^2 + (c_{i+1/2}^n)^2 & 2u_{i+1/2}^n & -(c_{i+1/2}^n)^2 \\ 0 & 0 & 0 \end{bmatrix}. \tag{5.92}$$

As it was pointed out in [18], the formulation of this scheme under the form (4.63) coincides with the *Q*-scheme of Roe scheme upwinding the source term presented in [36].

Accordingly to Theorem 2, this numerical scheme is well-balanced with order 2. Moreover, it is exactly well-balanced for the integral curves of the linearly degenerate field of equations

$$q = 0, \quad h - H = cte,$$

or

$$q = 0, \quad h = 0.$$

In other terms: the scheme approximates the steady state solutions with order 2, and it solves exactly those corresponding to water at rest or vacuum situations. In fact, these well-balance properties of the scheme under its formulation (4.63) were known yet: the exactly well-balance property for solutions corresponding to water at rest were proved in [36] (where this property was called *C*-property), and the well-balance property with order 2 for general steady state solutions in [12].

With this scheme, contact discontinuities related to bottom jumps can not be expected to be well captured. To illustrate this, let us consider the following example: we consider a straight channel whose depth function is given by:

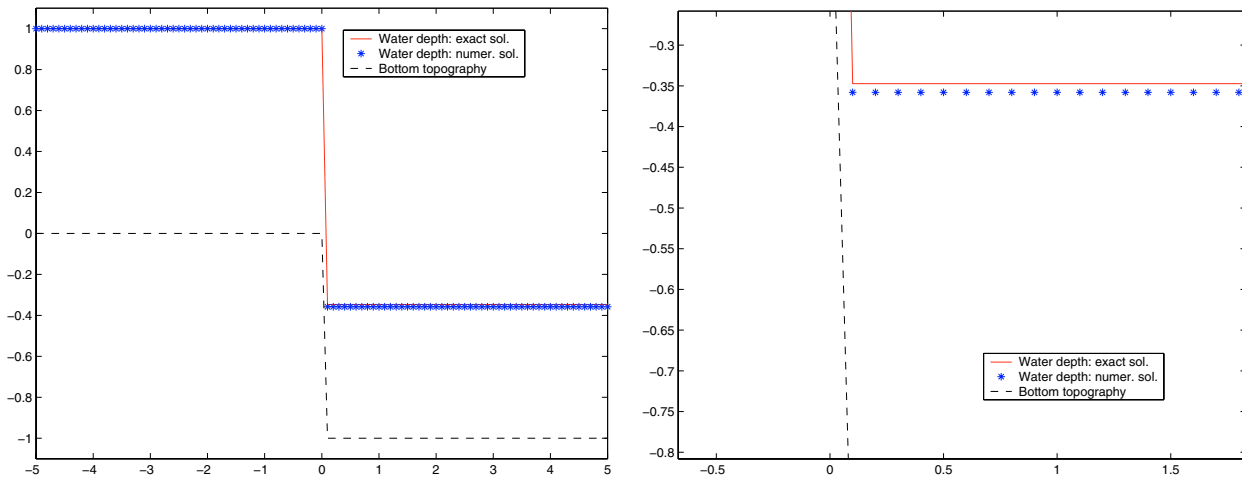
$$H(x) = \begin{cases} 0 & \text{if } x < 0; \\ 1 & \text{if } x > 0; \end{cases}$$

and we take initial conditions:

$$W(x, 0) = \begin{cases} W_L & \text{if } x < 0; \\ W_R & \text{if } x > 0; \end{cases}$$

where

$$W_L = \left[ \frac{1.0}{\sqrt{2g}} \right], \quad g = 9.81,$$



(a) Bottom topography and the computed water depth at the stationary state. Comparison with the exact solution.

(b) Zoom of Figure 1a.

FIGURE 1. A solution with a contact discontinuity related with a bottom jump: comparison between the stationary computed solution and the exact one.

and  $W_R$  is a state that can be connected to  $W_L$  by an entropic contact discontinuity. This state is calculated by using (5.87): the first equation gives  $q_R = \sqrt{2g}$  and the second one gives a cubic equation for  $h_R$  that is solved numerically. The only positive root giving a state that can be connected to  $W_L$  by an entropic contact discontinuity is selected:

$$W_R = \begin{bmatrix} 0.6527036446614 \\ \sqrt{2g} \end{bmatrix}.$$

We have applied the numerical scheme to this Riemann problem. As boundary conditions, the state  $W_L$  is imposed upstream and free flow condition is imposed downstream. These are the natural boundary conditions for this problem as  $\lambda_1$  and  $\lambda_2$  are both positives in all the domain. As expected, the steady state corresponding to the contact discontinuity linking the states is not well captured: small waves develop downwards and after some iterations, a steady state solution is obtained. Figure 1 depicts the water depth at the steady state solution obtained with  $\Delta x = 0.1$  and  $CFL = 0.99$ . While mass-fluxes and the thickness to the left of  $x = 0$  are computed exactly, there are some small differences in the thickness at the right (see Fig. 1a). Nevertheless, these small differences are not due to the discretization error but to a consistence error: they remain as  $\Delta x$  tends to 0. To observe this, in Table 1 we show the distance measured in the  $L^\infty$  norm between the thickness corresponding to the exact steady state solution and the numerical approximation obtained with  $CFL = 0.99$  and decreasing values of  $\Delta x$ .

As a consequence, the scheme can not be used to accurately simulate shallow flows over a bottom with jumps.

**Remark 4.** When the states  $\tilde{W}_i^n, \tilde{W}_{i+1}^n$  can be linked by a contact discontinuity, a good choice for the intermediate matrix can be calculated following the indications proposed in Remark 3:

$$\tilde{\mathcal{A}}_{i+1/2} = \begin{bmatrix} 1 & 0 & 0 \\ -(u_{i+1/2}^n)^2 + (c_{i+1/2}^n)^2 & 2u_{i+1/2}^n & -(c_{i+1/2}^n)^2 \\ 0 & 0 & 0 \end{bmatrix},$$

where:

$$u_{i+1/2}^n = \frac{q_{i+1/2}^n}{h_{i+1/2}^n}; \quad c_{i+1/2}^n = \sqrt{gh_{i+1/2}^n},$$

TABLE 1. Approximation of a contact discontinuity:  $L^\infty$  distance between the thickness corresponding to the stationary solution and its numerical approximation.

| $\Delta x$ | Error: $L^\infty$ norm |
|------------|------------------------|
| $10^{-1}$  | 0.01063                |
| $10^{-2}$  | 0.01063                |
| $10^{-3}$  | 0.01063                |
| $10^{-4}$  | 0.01063                |

with

$$q_{i+1/2}^n = q_i^n = q_{i+1}^n, \quad h_{i+1/2}^n = \left( \frac{(q_{i+1/2}^n)^2 [h]}{[h] - [H]} \right)^{\frac{1}{3}},$$

$[h], [H]$  being the jumps:

$$[h] = h_{i+1}^n - h_i^n, \quad [H] = H_{i+1} - H_i.$$

With this choice of intermediate matrix, the numerical scheme solves exactly every contact discontinuity and every regular stationary solution. Nevertheless, as it was said in Remark 3, it is difficult to combine this choice with the Roe linearization given above: in practice, in order to choose one or another type of intermediate matrix, a test has to be implemented to decide whether or not the states at two neighbor cells lie on the same integral curve of the linearly degenerate field, that is, if the total energy is preserved. Nevertheless, due to the discretization and round-off errors, this test has to be done with a finite precision bounded by a chosen parameter  $\varepsilon$ . We have observed in the numerical experiments that, depending of the value of this parameter, the resulting scheme can capture accurately contact discontinuities or stationary shocks, but not both of them at the same time.

### 5.2. Shallow water equations with breadth variations

Let us consider now the one-dimensional shallow water system:

$$\begin{cases} \frac{\partial A}{\partial t} + \frac{\partial q}{\partial x} = 0, \\ \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} \left( \frac{q^2}{A} + \frac{g}{2} \frac{A^2}{\sigma} \right) = g \frac{A^2}{\sigma^2} \frac{d\sigma}{dx} - \frac{g}{2} \frac{A^2}{\sigma^2} \frac{d\sigma}{dx}, \end{cases} \tag{5.93}$$

which are the equations governing the flow of a shallow layer of homogeneous inviscid fluid through a symmetric channel with a straight axis, a flat bottom and rectangular cross-sections of varying breadth. Again, the coordinate  $x$  refers to the axis of the channel,  $t$  is time, and  $g$  is gravity;  $q(x, t)$  and  $A(x, t)$  represent the mass-flow and the wetted cross-section respectively; and  $\sigma(x)$  is the breadth function.  $A(x, t)$  is related to the thickness  $h(x, t)$  by the formula:

$$A(x, t) = h(x, t)\sigma(x, t),$$

and  $q(x, t)$  to the velocity  $u(x, t)$  by:

$$q(x, t) = u(x, t)A(x, t).$$



System (5.93) can be written under the form (4.34) with  $N = N_1 = 2$ ,  $W = [A, q]^T$ , and

$$\begin{aligned}
 F(W, \sigma) &= \begin{bmatrix} q \\ \frac{q^2}{A} + \frac{g}{2} \frac{A^2}{\sigma} \end{bmatrix}, \\
 S(W, \sigma) &= \begin{bmatrix} 0 \\ g \frac{A^2}{\sigma^2} \end{bmatrix}.
 \end{aligned}
 \tag{5.94}$$

Clearly:

$$\frac{\partial F}{\partial \sigma}(W, \sigma) = \begin{bmatrix} 0 \\ -\frac{g}{2} \frac{A^2}{\sigma^2} \end{bmatrix}.$$

The system can be written under the form (2.4) with

$$\widetilde{W} = \begin{bmatrix} A \\ q \\ \sigma \end{bmatrix},$$

and

$$\widetilde{\mathcal{A}}(\widetilde{W}) = \begin{bmatrix} 0 & 1 & 0 \\ -u^2 + c^2 & 2u & -hc^2 \\ 0 & 0 & 0 \end{bmatrix},
 \tag{5.95}$$

where  $c = \sqrt{gh}$ .

The eigenvalues of this matrix are again:

$$\lambda_1 = u - c, \quad \lambda_2 = u + c, \quad 0.$$

If  $c \neq 0$ ,  $\lambda_i \neq 0$ ,  $i = 1, 2$ , the system is strictly hyperbolic and a complete set of eigenvectors is given by:

$$\widetilde{R}_i(W) = \begin{bmatrix} 1 \\ \lambda_i \\ 0 \end{bmatrix}, \quad i = 1, 2; \quad \widetilde{R}_3(W) = \begin{bmatrix} 1 \\ 0 \\ (1 - Fr^2)/h \end{bmatrix};$$

where  $Fr$  is given by (5.86). The integral curves of the 3rd-field are given by:

$$q = ct., \quad \frac{A}{\sigma} + \frac{q^2}{2gA^2} = ct.
 \tag{5.96}$$

In the particular case  $q = 0$  these curves are straight lines.

As in the previous case, two kind of discontinuities may appear in the weak solutions: 1 or 2-shocks satisfying the usual Rankine-Hugoniot conditions and evolving in zones when the breadth function is regular, and stationary contact discontinuities standing at points where the breadth function presents a jump. Again, two states  $\widetilde{W}_L$  and  $\widetilde{W}_R$  can be connected by a contact discontinuity when the total mechanical energy is preserved.

Steady state solutions corresponding to water at rest are given by:

$$\frac{A}{\sigma} = ct., \quad q = 0,
 \tag{5.97}$$

and those corresponding to vacuum by:

$$A = 0, \quad q = 0.
 \tag{5.98}$$

In this case the construction of a Roe scheme based on the family of paths (2.16) is not as easy as in the previous case. First the Roe matrices  $\mathcal{J}_{i+1/2}$  given by (5.90), (5.91) are again considered. Then, the numerical source term has to satisfy:

$$S_{i+1/2} = \int_0^1 S(\widetilde{W}_i^n + s(\widetilde{W}_{i+1}^n - \widetilde{W}_i^n)) ds.$$

After some calculations, the following expression is obtained:

$$S_{i+1/2} = \begin{bmatrix} 0 \\ s_{i+1/2} \end{bmatrix};$$

where  $s_{i+1/2}$  is equal to 0 if  $\sigma_i = \sigma_{i+1}$ , and it is given by:

$$\frac{g}{(\sigma_i - \sigma_{i+1})^2} \left( 2 \frac{A_{i+1} - A_i}{\sigma_{i+1} - \sigma_i} (A_i \sigma_{i+1} - A_{i+1} \sigma_i) \log \left( \frac{\sigma_{i+1}}{\sigma_i} \right) - 4A_i A_{i+1} + \left( \frac{A_{i+1}^2}{\sigma_{i+1}} + \frac{A_i^2}{\sigma_i} \right) (\sigma_i + \sigma_{i+1}) \right),$$

otherwise.

Finally an intermediate value of  $\sigma$  satisfying (4.55) has to be calculated. If  $\sigma_i \neq \sigma_{i+1}$  this is done by solving:

$$\frac{A_{i+1}^2 - A_i^2}{\sigma_{i+1/2}} = \frac{A_{i+1}^2}{\sigma_{i+1}} - \frac{A_i^2}{\sigma_i} + s_{i+1/2} \frac{\sigma_{i+1} - \sigma_i}{g}.$$

We deduce from Theorem 3 that the resulting scheme is well-balanced with order 2 and exactly well-balanced for steady state solutions corresponding to water at rest or vacuum. Moreover, the numerical scheme can handle correctly with the shocks, as (4.69) is satisfied when  $\sigma_i = \sigma_{i+1}$ . In our knowledge, this scheme has not been yet described.

A numerical scheme with the same properties and a more suitable expression can be constructed by choosing a family of paths which are segments in the variables:

$$\widetilde{W}^* = \begin{bmatrix} A \\ q \\ h \end{bmatrix}.$$

More precisely, given  $\widetilde{W}_L = [A_L, q_L, \sigma_L]^T$ ;  $\widetilde{W}_R = [A_R, q_R, \sigma_R]^T$  and  $s \in [0, 1]$ , we define:

$$\widetilde{\Psi}(s; W_L, W_R) = \begin{bmatrix} A_L + s(A_R - A_L) \\ q_L + s(q_R - q_L) \\ \frac{A_L + s(A_R - A_L)}{h_L + s(h_R - h_L)} \end{bmatrix} \tag{5.99}$$

where  $h_L = A_L/\sigma_L$ ,  $h_R = A_R/\sigma_R$ .

With this choice of paths, some easy calculations allow to show that:

$$\begin{aligned} S_{i+1/2}(\sigma_{i+1} - \sigma_i) &= \int_0^1 S(\widetilde{\Psi}(s; \widetilde{W}_i, \widetilde{W}_{i+1})) \frac{d}{ds} \left( \frac{A_L + s(A_R - A_L)}{h_L + s(h_R - h_L)} \right) ds \\ &= \begin{bmatrix} 0 \\ -A_{i+1/2}(h_{i+1} - h_i) + h_{i+1/2}(A_{i+1} - A_i) \end{bmatrix}, \end{aligned}$$

where

$$A_{i+1/2} = \frac{A_i + A_{i+1}}{2}, \quad h_{i+1/2} = \frac{h_i + h_{i+1}}{2}.$$

Finally, it can be easily shown that

$$\sigma_{i+1/2} = \frac{A_{i+1/2}}{h_{i+1/2}}$$

satisfies (4.55).

From Theorem 4 we deduce that this numerical scheme is also well-balanced with order 2. Moreover it is exactly well-balanced for steady state solutions corresponding to water at rest or vacuum, as the corresponding curves in the  $(A, q, h)$  variables are the straight lines of equation  $h = ct$ .

Observe that, if  $\sigma_i = \sigma_{i+1}$  the identity (4.69) is not satisfied in general. Nevertheless, in this case it can be easily proved that:

$$-A_{i+1/2}(h_{i+1} - h_i) + h_{i+1/2}(A_{i+1} - A_i) = 0,$$

and thus:

$$\int_0^1 S\left(\tilde{\Psi}(s; \tilde{W}_i, \tilde{W}_{i+1})\right) \frac{d}{ds} \left( \frac{A_L + s(A_R - A_L)}{h_L + s(h_R - h_L)} \right) ds = 0.$$

As  $\frac{\partial F}{\partial \sigma}(W, \sigma) = -\frac{1}{2}S(W, \sigma)$  the following identity also holds:

$$\int_0^1 \frac{\partial F}{\partial \sigma}(W, \sigma) \left( \tilde{\Psi}(s; \tilde{W}_i, \tilde{W}_{i+1}) \right) \frac{d}{ds} \left( \frac{A_L + s(A_R - A_L)}{h_L + s(h_R - h_L)} \right) ds = 0.$$

Therefore, the numerical scheme satisfies (4.68) when  $\sigma_i = \sigma_{i+1}$  and, as a consequence, the shocks are correctly taken into account.

This last numerical scheme is similar to the one introduced in [14] for a more general problem corresponding to a symmetric channel with an arbitrary cross-section.

Systematic comparison between these two schemes are in progress.

### 5.3. The bilayer shallow water system

We consider in this paragraph the system of partial differential equations governing the one-dimensional flow of two superposed immiscible layers of shallow water fluids studied in [4]:

$$\left\{ \begin{array}{l} \frac{\partial h_1}{\partial t} + \frac{\partial q_1}{\partial x} = 0, \\ \frac{\partial q_1}{\partial t} + \frac{\partial}{\partial x} \left( \frac{q_1^2}{h_1} + \frac{g}{2} h_1^2 \right) = -gh_1 \frac{\partial h_2}{\partial x} + gh_1 \frac{dH}{dx}, \\ \frac{\partial h_2}{\partial t} + \frac{\partial q_2}{\partial x} = 0, \\ \frac{\partial q_2}{\partial t} + \frac{\partial}{\partial x} \left( \frac{q_2^2}{h_2} + \frac{g}{2} h_2^2 \right) = -\frac{\rho_1}{\rho_2} gh_2 \frac{\partial h_1}{\partial x} + gh_2 \frac{dH}{dx}. \end{array} \right. \tag{5.100}$$

In these equations, index 1 makes reference to the upper layer and index 2 to the lower one. The fluid is assumed to occupy a straight channel with constant rectangular cross-section and constant width. The coordinate  $x$  refers to the axis of the channel,  $t$  is the time, and  $g$  is gravity.  $H(x)$  represents the depth function measured from a fixed level of reference. Each layer is assumed to have a constant density,  $\rho_i$ ,  $i = 1, 2$  ( $\rho_1 < \rho_2$ ). The unknowns  $q_i(x, t)$  and  $h_i(x, t)$  represent respectively the mass-flow and the thickness of the  $i$ th layer at the section of coordinate  $x$  at time  $t$ .

System (5.100) can be written under the form (4.32) with  $K = 2, N_1 = 2, N_2 = 2, N = 4, \sigma = H,$

$$W_j(x, t) = \begin{bmatrix} h_j(x, t) \\ q_j(x, t) \end{bmatrix}, \quad j = 1, 2. \tag{5.101}$$

Here,  $F_1 = F_2 = F,$  where  $F$  is given by (5.84), and

$$S_1(W_1, W_2) = \begin{bmatrix} 0 \\ gh_1 \end{bmatrix}, \quad S_2(W_1, W_2) = \begin{bmatrix} 0 \\ gh_2 \end{bmatrix},$$

$$\mathcal{B}_{1,2}(W_1, W_2) = \begin{bmatrix} 0 & 0 \\ -gh_1 & 0 \end{bmatrix}, \quad \mathcal{B}_{2,1}(W_1, W_2) = \begin{bmatrix} 0 & 0 \\ -grh_2 & 0 \end{bmatrix},$$

with

$$r = \frac{\rho_1}{\rho_2}.$$

The system can also be written under the forms (4.35) or (4.37). In this last formulation:

$$\widetilde{\mathbf{W}} = \begin{bmatrix} W_1 \\ W_2 \\ H \end{bmatrix}, \tag{5.102}$$

$$\widetilde{\mathcal{A}}(\widetilde{\mathbf{W}}) = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ -u_1^2 + c_1^2 & 2u_1 & c_1^2 & 0 & -c_1^2 \\ 0 & 0 & 0 & 1 & 0 \\ rc_2^2 & 0 & -u_2^2 + c_2^2 & 2u_2 & -c_2^2 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \tag{5.103}$$

where  $u_i = q_i/h_i$  represents the averaged velocity of the  $i$ th layer and  $c_i = \sqrt{gh_i}, i = 1, 2.$

The eigenvalues of this matrix may become complex, corresponding to the development of shear instabilities. In this work only the case where the matrix  $\mathcal{A}$  has real eigenvalues is considered, *i.e.*, the flow is supposed to be stable and the system hyperbolic.

The eigenvector associated to the eigenvalue 0 is:

$$\widetilde{\mathbf{R}}_5 = \begin{bmatrix} -Fr_2^2 \\ 0 \\ 1 - Fr_1^2 \\ 0 \\ 1 - Fr_c^2 \end{bmatrix},$$

where:

$$Fr_i = \frac{u_i}{\sqrt{g_r h_i}}, \quad i = 1, 2,$$

with

$$g_r = (1 - r)g,$$

and

$$Fr_c^2 = Fr_1^2 + Fr_2^2 - (1 - r)Fr_1^2 Fr_2^2.$$

$Fr_i, i = 1, 2,$  are the *internal Froude numbers*;  $Fr_c,$  the *composite Froude number*; and  $g_r$  the *reduced gravity*.

The integral curves of the 5th characteristic field are calculated by integrating (4.39):

$$\begin{cases} q_1 = ct. \\ \frac{u_1^2}{2} - \frac{u_2^2}{2} + g_r h_1 = ct. \\ q_2 = ct. \\ \frac{u_1^2}{2} + g(h_1 + h_2 - H) = ct. \end{cases} \tag{5.104}$$

In this case it is not easy to verify the genuinely nonlinear character of the 4 other characteristic fields, as the eigenvalues and eigenvectors can not be explicitly written in a simple manner. Nevertheless, this fact can be easily proved if  $r = 0$  as, in this case, the system reduces to two decoupled shallow water systems. As a consequence, using a continuity argument, this is also true at least for small values of  $r$ . We assume here that this hypothesis is satisfied.

Concerning the definition of weak solutions, accordingly to Section 4.2 a family of paths has to be chosen firstly for the homogeneous system (4.43), that is, the system corresponding to a channel with a constant depth  $H$ . As this homogeneous problem is again independent of  $H$ , only a family of paths  $\Phi$  has to be chosen. At this stage, we have considered the family of segments:

$$\Phi(s; \mathbf{W}_L, \mathbf{W}_R) = \mathbf{W}_L + s(\mathbf{W}_R - \mathbf{W}_L),$$

that is, we assume Volpert's definition for the nonconservative products. Nevertheless, further investigations will be performed in order to look for a family of paths closer to the physical background of the problem.

Once  $\Phi$  chosen, the family of paths  $\tilde{\Phi}$  corresponding to the nonhomogeneous case (5.100) is determined (see Sect. 4.2). Again, the calculation of the path linking two arbitrary states requires to solve a Riemann problem, which is a very difficult task in this case.

The steady states solutions of (5.83) can be calculated by parameterizing the curves (5.104) with  $x$ . In particular, if  $q_1 = 0, q_2 = 0$ , we obtain the solutions:

$$q_1 = 0, \quad h_1(x) = ct., \quad q_2 = 0, \quad h_2(x) - H(x) = ct., \tag{5.105}$$

representing water at rest. Another steady state solution corresponding to *vacuum* is given by:

$$h_i = 0, \quad q_i = 0, \quad i = 1, 2. \tag{5.106}$$

We consider again the family of segments to construct the Roe linearization. Following Section 4.3, we first consider the Roe matrices for each homogeneous shallow water system:

$$\mathcal{J}_{i+1/2,k} = \begin{bmatrix} 1 & 0 \\ -(u_{i+1/2,k}^n)^2 + (c_{i+1/2,k}^n)^2 & 2u_{i+1/2,k}^n \end{bmatrix}, \quad k = 1, 2,$$

where

$$u_{i+1/2,k}^n = \frac{\sqrt{h_{i,k}^n} u_{i,k}^n + \sqrt{h_{i+1,k}^n} u_{i+1,k}^n}{\sqrt{h_{i,k}^n} + \sqrt{h_{i+1,k}^n}}, \quad c_{i+1/2,k}^n = \sqrt{g \frac{h_{i,k}^n + h_{i+1,k}^n}{2}}, \tag{5.107}$$

and then, we calculate  $\mathbf{B}_{i+1/2}$  and  $\mathbf{S}_{i+1/2}$  by (4.56), (4.57). The Roe matrices  $\tilde{\mathbf{A}}_{i+1/2}$  obtained are as follows:

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ -(u_{i+1/2,1}^n)^2 + (c_{i+1/2,1}^n)^2 & 2u_{i+1/2,1}^n & (c_{i+1/2,1}^n)^2 & 0 & -(c_{i+1/2,1}^n)^2 \\ 0 & 0 & 0 & 1 & 0 \\ r(c_{i+1/2,2}^n)^2 & 0 & -(u_{i+1/2,2}^n)^2 + (c_{i+1/2,2}^n)^2 & 2u_{i+1/2,2}^n & -(c_{i+1/2,2}^n)^2 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}. \tag{5.108}$$

If the numerical scheme is expressed under the formulation (4.63), it coincides exactly with the generalized  $Q$ -scheme of Roe upwinding the source term introduced in [4].

Accordingly to Theorem 3, the numerical scheme is well-balanced with order 2. Moreover, it is exactly well-balanced for the integral curves of the linearly degenerate field:

$$q_1 = 0, \quad h_1 = ct., \quad q_2 = 0, \quad h_2 - H = ct.,$$

or

$$q_1 = 0, \quad h_1 = 0, \quad q_2 = 0, \quad h_2 = 0,$$

which are straight lines. In other terms: the scheme approaches the steady state solutions with order 2, and it solves exactly those corresponding to water at rest or vacuum. The exact well-balance property for water at rest was proved, for the formulation (4.63), in [5].

In [4, 5], some numerical solutions obtained with this scheme (or with its natural extension to more general cases) have been compared successfully with experimental measurements or with analytical steady state solutions. Nevertheless, as in the case of the shallow water system, contact discontinuities related to bottom jumps can not be expected to be well captured.

As Volpert’s definition is assumed for the nonconservative products in the homogenous case, the paths used in the definitions of weak solutions and Roe matrices coincide. Therefore,  $j$ -shocks,  $1 \leq j \leq 4$ , are expected to be correctly approached. Let us verify this with a numerical example. To begin with, we calculate a steady state solution of the homogeneous system ( $H$  is constant) representing an *internal hydraulic jump* as follows: we fix the values of the constants  $g = 10$ ,  $r = 0.02$ , we put  $\xi = 0$  and we choose the state to the left of the shock as:

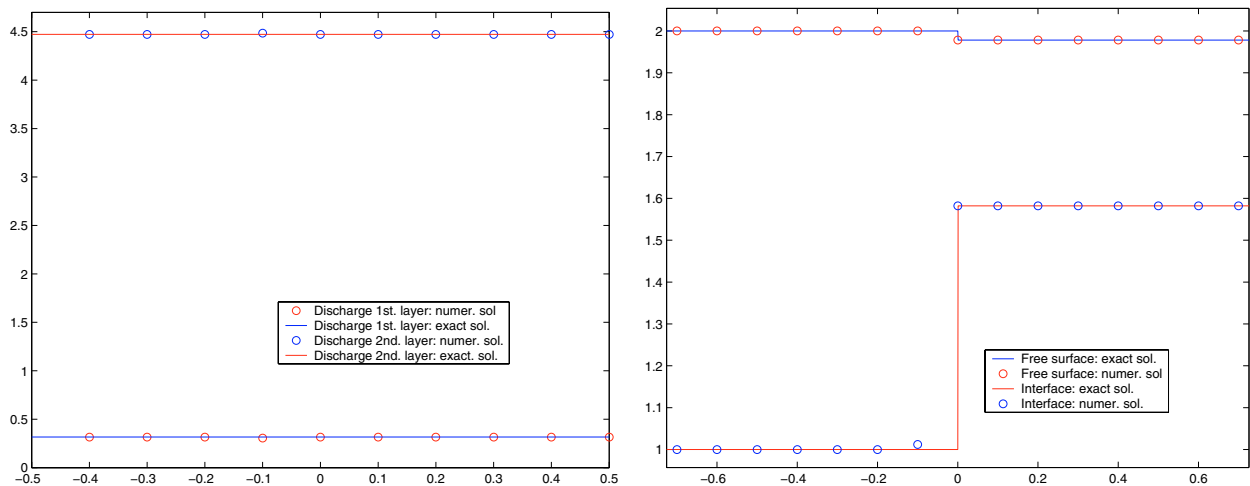
$$\mathbf{W}_L = \begin{bmatrix} 1 \\ \sqrt{0.1} \\ 1 \\ \sqrt{20} \end{bmatrix}.$$

Then, the generalized Rankine-Hugoniot conditions give a nonlinear system for the state to the right that is solved numerically. Among the set of solutions obtained, we select the only one satisfying Lax’s entropy condition, which is:

$$\mathbf{W}_R = \begin{bmatrix} 0.396156 \\ \sqrt{0.1} \\ 1.5820186 \\ \sqrt{20} \end{bmatrix}.$$

The generalized Rankine-Hugoniot conditions are satisfied by  $\mathbf{W}_L, \mathbf{W}_R, \xi = 0$  with an error of about  $10^{-5}$ .

We have applied the numerical scheme to the Riemann problem with the initial condition given by these two states. As the jump condition is not satisfied in an exact manner, initially some small waves develop. After some iterations, a steady state solution is obtained. Figure 2 depicts the steady state solution obtained with  $\Delta x = 0.066666666$  and  $CFL = 0.99$ . The  $L^\infty$  distance between the exact and the numerical solution is of order  $10^{-3}$  and it can not be appreciated in Figure 2.



(a) Zoom at  $x = 0$  of the computed discharges at the stationary state. Comparison with the exact solution.

(b) Zoom at  $x = 0$  of the computed free surface and interface at stationary state. Comparison with the exact solution.

FIGURE 2. Solution with a stationary shock of the two-layer shallow-water system: comparison between the numerical and the exact solution.

*Acknowledgements.* This research has been partially supported by the Spanish Government Research projects REN2000-1168-C02-01 and BFM2003-07530-C02-02.

## REFERENCES

- [1] N. Andronov and G. Warnecke, On the solution to the Riemann problem for the compressible duct flow. *SIAM J. Appl. Math.* **64** (2004) 878–901.
- [2] F. Bouchut, An introduction to finite volume methods for hyperbolic systems of conservation laws with source, in *Free surface geophysical flows. Tutorial Notes*. INRIA, Rocquencourt (2002).
- [3] A. Bermúdez and M.E. Vázquez, Upwind methods for hyperbolic conservation laws with source terms. *Comput. Fluids* **23** (1994) 1049–1071.
- [4] M.J. 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. *ESAIM: M2AN* **35** (2001) 107–127.
- [5] M.J. Castro, J.A. García-Rodríguez, J.M. González-Vida, J. Macías, C. Parés and M.E. Vázquez-Cendón, Numerical simulation of two-layer Shallow Water flows through channels with irregular geometry. *J. Comp. Phys.* **195** (2004) 202–235.
- [6] T. Chacón, A. Domínguez and E.D. Fernández, A family of stable numerical solvers for Shallow Water equations with source terms. *Comp. Meth. Appl. Mech. Eng.* **192** (2003) 203–225.
- [7] T. Chacón, A. Domínguez and E.D. Fernández, An entropy-correction free solver for non-homogeneous shallow water equations. *ESAIM: M2AN* **37** (2003) 755–772.
- [8] T. Chacón, E.D. Fernández and M. Gómez Mármol, A flux-splitting solver for shallow water equations with source terms. *Int. Jour. Num. Meth. Fluids* **42** (2003) 23–55.
- [9] T. Chacón, A. Domínguez and E.D. Fernández, Asymptotically balanced schemes for non-homogeneous hyperbolic systems – application to the Shallow Water equations. *C.R. Acad. Sci. Paris, Ser. I* **338** (2004) 85–90.
- [10] J.F. Colombeau, A.Y. Le Roux, A. Noussair and B. Perrot, Microscopic profiles of shock waves and ambiguities in multiplications of distributions. *SIAM J. Num. Anal.* **26** (1989) 871–883.
- [11] G. Dal Masso, P.G. LeFloch and F. Murat, Definition and weak stability of nonconservative products. *J. Math. Pures Appl.* **74** (1995) 483–548.
- [12] E.D. Fernández Nieto, *Aproximación Numérica de Leyes de Conservación Hiperbólicas No Homogéneas. Aplicación a las Ecuaciones de Aguas Someras*. Ph.D. Thesis, Universidad de Sevilla (2003).
- [13] A.C. Fowler, *Mathematical Model in the Applied Sciences*. Cambridge (1997).
- [14] P. García-Navarro and M.E. Vázquez-Cendón, On numerical treatment of the source terms in the shallow water equations. *Comput. Fluids* **29** (2000) 17–45.
- [15] P. Goatin and P.G. LeFloch, The Riemann problem for a class of resonant hyperbolic systems of balance laws, preprint (2003).

- [16] E. Godlewski and P.A. Raviart, *Numerical Approximation of Hyperbolic Systems of Conservation Laws*. Springer-Verlag, New York (1996).
- [17] L. Gosse, A well-balanced flux-vector splitting scheme designed for hyperbolic systems of conservation laws with source terms. *Comp. Math. Appl.* **39** (2000) 135–159.
- [18] L. Gosse, A well-balanced scheme using non-conservative products designed for hyperbolic system of conservation laws with source terms. *Mat. Mod. Meth. Appl. Sc.* **11** (2001) 339–365.
- [19] J.M. Greenberg and A.Y. LeRoux, A well balanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM J. Numer. Anal.* **33** (1996) 1–16.
- [20] J.M. Greenberg, A.Y. LeRoux, R. Baraille and A. Noussair, Analysis and approximation of conservation laws with source terms. *SIAM J. Numer. Anal.* **34** (1997) 1980–2007.
- [21] A. Harten and J.M. Hyman, Self-adjusting grid methods for one-dimensional hyperbolic conservation laws. *J. Comp. Phys.* **50** (1983) 235–269.
- [22] P.G. LeFloch, Propagating phase boundaries; formulation of the problem and existence *via* Glimm scheme. *Arch. Rat. Mech. Anal.* **123** (1993) 153–197.
- [23] R. LeVeque, *Numerical Methods for Conservation Laws*. Birkhäuser (1990).
- [24] R. LeVeque, Balancing source terms and flux gradients in high-resolution Godunov methods: the quasi-steady wave-propagation algorithm. *J. Comp. Phys.* **146** (1998) 346–365.
- [25] R. LeVeque, *Finite Volume Methods for Hyperbolic Problems*. Cambridge University Press (2002).
- [26] B. Perthame and C. Simeoni, A kinetic scheme for the Saint-Venant system with a source term. *Calcolo* **38** (2001) 201–231.
- [27] B. Perthame and C. Simeoni, *Convergence of the upwind interface source method for hyperbolic conservation laws*, in *Proc. of Hyp 2002*, Thou and Tadmor Eds., Springer (2003).
- [28] P.A. Raviart and L. Sainsaulieu, A nonconservative hyperbolic system modeling spray dynamics. I. Solution of the Riemann problem. *Math. Mod. Meth. Appl. Sci.* **5** (1995) 297–333.
- [29] P.L. Roe, Approximate Riemann solvers, parameter vectors and difference schemes. *J. Comp. Phys.* **43** (1981) 357–371.
- [30] P.L. Roe, *Upwinding difference schemes for hyperbolic conservation laws with source terms*, in *Proc. of the Conference on Hyperbolic Problems*, Carasso, Raviart and Serre Eds., Springer (1986) 41–51.
- [31] J.J. Stoker, *Water Waves*. Interscience, New York (1957).
- [32] E.F. Toro, *Riemann Solvers and Numerical Methods for Fluid Dynamics. A Practical Introduction*. Springer-Verlag (1997).
- [33] E.F. Toro, *Shock-Capturing Methods for Free-Surface Shallow Flows*. Wiley (2001).
- [34] E.F. Toro and M.E. Vázquez-Cendón, *Model hyperbolic systems with source terms: exact and numerical solutions*, in *Proc. of Godunov methods: Theory and Applications* (2000).
- [35] I. Toumi, A weak formulation of Roe’s approximate Riemann Solver. *J. Comp. Phys.* **102** (1992) 360–373.
- [36] M.E. Vázquez-Cendón, *Estudio de Esquemas Descentrados para su Aplicación a las Leyes de Conservación Hiperbólicas con Términos Fuente*. Ph.D. Thesis, Universidad de Santiago de Compostela (1994).
- [37] M.E. Vázquez-Cendón, Improved treatment of source terms in upwind schemes for the shallow water equations in channels with irregular geometry. *J. Comp. Phys.* **148** (1999) 497–526.
- [38] A.I. Volpert, The space BV and quasilinear equations. *Math. USSR Sbornik* **73** (1967) 225–267.